



THE  
POWER  
TO KNOW.

# **SAS/ETS<sup>®</sup> 13.1**

## **User's Guide**

### **High-Performance Procedures**

The correct bibliographic citation for this manual is as follows: SAS Institute Inc. 2013. *SAS/ETS® 13.1 User's Guide: High-Performance Procedures*. Cary, NC: SAS Institute Inc.

**SAS/ETS® 13.1 User's Guide: High-Performance Procedures**

Copyright © 2013, SAS Institute Inc., Cary, NC, USA

All rights reserved. Produced in the United States of America.

**For a hard-copy book:** No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior written permission of the publisher, SAS Institute Inc.

**For a web download or e-book:** Your use of this publication shall be governed by the terms established by the vendor at the time you acquire this publication.

The scanning, uploading, and distribution of this book via the Internet or any other means without the permission of the publisher is illegal and punishable by law. Please purchase only authorized electronic editions and do not participate in or encourage electronic piracy of copyrighted materials. Your support of others' rights is appreciated.

**U.S. Government License Rights; Restricted Rights:** The Software and its documentation is commercial computer software developed at private expense and is provided with RESTRICTED RIGHTS to the United States Government. Use, duplication or disclosure of the Software by the United States Government is subject to the license terms of this Agreement pursuant to, as applicable, FAR 12.212, DFAR 227.7202-1(a), DFAR 227.7202-3(a) and DFAR 227.7202-4 and, to the extent required under U.S. federal law, the minimum restricted rights as set out in FAR 52.227-19 (DEC 2007). If FAR 52.227-19 is applicable, this provision serves as notice under clause (c) thereof and no other notice is required to be affixed to the Software or documentation. The Government's rights in Software and documentation shall be only those set forth in this Agreement.

SAS Institute Inc., SAS Campus Drive, Cary, North Carolina 27513-2414.

December 2013

SAS provides a complete selection of books and electronic products to help customers use SAS® software to its fullest potential. For more information about our offerings, visit [support.sas.com/bookstore](http://support.sas.com/bookstore) or call 1-800-727-3228.

SAS® and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. ® indicates USA registration.

Other brand and product names are trademarks of their respective companies.



# Gain Greater Insight into Your SAS<sup>®</sup> Software with SAS Books.

Discover all that you need on your journey to knowledge and empowerment.



[support.sas.com/bookstore](http://support.sas.com/bookstore)  
for additional books and resources.





# Contents

---

Chapter 1.	What's New in SAS/ETS 13.1 High-Performance Procedures . . . . .	1
Chapter 2.	Introduction . . . . .	7
Chapter 3.	Shared Concepts and Topics . . . . .	11
Chapter 4.	The HPCDM Procedure (Experimental) . . . . .	43
Chapter 5.	The HPCOPULA Procedure (Experimental) . . . . .	113
Chapter 6.	The HPCOUNTREG Procedure . . . . .	127
Chapter 7.	The HPPANEL Procedure (Experimental) . . . . .	159
Chapter 8.	The HPQLIM Procedure . . . . .	185
Chapter 9.	The HPSEVERITY Procedure . . . . .	237

<b>Subject Index</b>	<b>387</b>
----------------------	------------

<b>Syntax Index</b>	<b>390</b>
---------------------	------------



# Credits and Acknowledgments

---

## Credits

---

### Documentation

Editing	Anne Baxter, Ed Huddleston
Documentation Support	Tim Arnold

---

### Software

The procedures in this book were implemented by the following members of the development staff. Program development includes design, programming, debugging, support, and documentation. In the following list, the names of the developers who currently provide primary support are listed first; other developers and previous developers are also listed.

HPCDM	Mahesh Joshi, Richard Potter, Jan Chvosta
HPCOPULA	Hao Chen, Richard Potter, Jan Chvosta
HPCOUNTREG	Richard Potter, Jan Chvosta
HPPANEL	Linxia Ren, Richard Potter
HPQLIM	Christian Macaro, Richard Potter, Jan Chvosta
HPSEVERITY	Mahesh Joshi
High-performance computing foundation	Steve E. Krueger
High-performance analytics foundation	Robert Cohen, Georges H. Guirguis, Trevor Kearney, Richard Knight, Gang Meng, Oliver Schabenberger, Charles Shorb, Tom P. Weber
Numerical routines	Georges H. Guirguis

The following people contribute with their leadership and support: Chris Bailey, Tanya Balan, David Pope, Oliver Schabenberger, Renee Sciortino.

---

### Testing

Ming Chun-Chang, Bruce Elsheimer, Sanggohn Han, Oleksiy Tokovenko, Jim McKenzie, Dright Ho, Girija Gavankar, Tim Carter, Bengt Pederson, Cheryl LeSaint, Jim Metcalf

---

## **Internationalization Testing**

Alex Chai, Mi-Na Chu, Jacky Dong, Feng Gao, Masayuki Iizuka, David Li, Lan Luan, Haiyong Rong, Bin Sun, Frank Wang, Lina Xu, Catherine Yang

---

## **Technical Support**

Phil Gibbs, Wen Bardsley, David Schlotzhauer, Donna Woodward

---

## **Acknowledgments**

Many people make significant and continuing contributions to the development of SAS software products.

The final responsibility for the SAS System lies with SAS alone. We hope that you will always let us know your opinions about the SAS System and its documentation. It is through your participation that SAS software is continuously improved.



# Chapter 1

## What's New in SAS/ETS 13.1

### High-Performance Procedures

#### Contents

Overview . . . . .	1
Highlights of Changes and Enhancements . . . . .	1
HPCDM Procedure (Experimental) . . . . .	2
HPCOPULA Procedure (Experimental) . . . . .	3
HPCOUNTREG Procedure . . . . .	3
HPPANEL Procedure (Experimental) . . . . .	3
HPQLIM Procedure . . . . .	4
HPSEVERITY Procedure . . . . .	4
References . . . . .	5

## Overview

This chapter summarizes the new features available in SAS/ETS 13.1 high-performance procedures.

If you have used SAS/ETS high-performance procedures in the past, you can review this chapter to learn about the new features that have been added. When you see a new feature that might be useful for your work, turn to the appropriate chapter to read about the feature in detail.

## Highlights of Changes and Enhancements

The following procedures have been added to SAS/ETS high-performance procedures:

- HPCDM procedure
- HPCOPULA procedure
- HPPANEL procedure

New features have been added to the following SAS/ETS high-performance procedures:

- HPCOUNTREG procedure

- HPQLIM procedure
- HPSEVERITY procedure

---

## HPCDM Procedure (Experimental)

The **HPCDM procedure** is a new, experimental high-performance procedure in SAS/ETS. The HPCDM procedure estimates a compound distribution model (CDM), which is the distribution of an aggregate loss that you expect to see in a given period of time. The aggregate loss depends on the number of loss events that occur in a given period of time and the severity (magnitude) of each loss event.

If you have estimated the probability distribution model of the frequency of loss events by using the COUNTREG procedure in SAS/ETS, and if you have estimated the probability distribution model of the severity of each loss by using the SEVERITY procedure in SAS/ETS, then PROC HPCDM combines the frequency and severity models to estimate the distribution of an aggregate loss. The probability distribution model of the aggregate loss is referred to as the compound distribution model (CDM).

At its core, PROC HPCDM uses the Monte Carlo simulation method to generate a random sample of the aggregate loss. The random sample is used to compute empirical estimates of various summary statistics and percentiles of the compound distribution. The HPCDM procedure also has the following features:

- PROC HPCDM performs a scenario or what-if analysis when the distributions of severity and count depend on a set of exogenous variables. You can simulate a scenario for one entity that is subject to loss events. For example, you can estimate the distribution of the aggregate loss that one insurance policyholder might incur by using the policyholder's characteristics. You can also simulate a scenario that consists of multiple entities that are subject to the loss events. For example, you can estimate the distribution of the aggregate loss that is incurred together by a portfolio of financial instruments by using the characteristics of each instrument as well as the macroeconomic factors that affect each of them.
- PROC HPCDM can perform a parameter perturbation analysis to assess the effect of the uncertainty in the estimates of severity and frequency model parameters on the uncertainty of various statistics of the compound distribution. For example, if you are interested in computing the value-at-risk (VaR), which is usually the 97.5th or the 99.5th percentile of the aggregate loss distribution, then by using perturbation analysis of PROC HPCDM, you can get an estimate of not only the mean expected VaR but also the standard error that is associated with VaR. Both estimates can help you maintain an adequate amount of capital for regulatory compliance as well as a healthy, solvent business.
- PROC HPCDM estimates the compound distribution of an aggregate *adjusted* loss when you specify SAS programming statements to adjust the magnitude of each simulated loss. For example, this is helpful in estimating the distribution of the total amount paid by an insurance company to its policyholders in a given period of time. The amount paid for each loss depends on the magnitude of the loss as well as the various provisions in the insurance policy, such as the deductible and policy limit.

---

## HPCOPULA Procedure (Experimental)

The new [HPCOPULA procedure](#) is a high-performance version of the COPULA procedure in SAS/ETS software, which enables you to simulate realizations of multivariate distributions by using the copula approach.

The HPCOPULA procedure supports simulation of multivariate distributions for the following types of copulas:

- normal copula
- $t$  copula
- Clayton copula
- Gumbel copula
- Frank copula

---

## HPCOUNTREG Procedure

The following features have been added to the [HPCOUNTREG procedure](#):

- support for the BY statement

---

## HPPANEL Procedure (Experimental)

The new [HPPANEL procedure](#) is a high-performance version of the PANEL procedure in SAS/ETS software, which analyzes a class of linear econometric panel data models.

The HPPANEL procedure provides the following four models:

- one-way fixed-effects model
- two-way fixed-effects model
- one-way random-effects model
- two-way random-effects model

The following features have been added:

- The TEST statement performs Wald, Lagrange multiplier, and likelihood ratio tests of linear hypotheses about the regression parameters.
- You can request the restricted estimator by using the RESTRICT statement.
- You can request the Hausman (1978) specification test for random effects.
- For random-effects models, variance components are calculated by using methods described by Fuller and Battese (1974); Wansbeek and Kapteyn (1989); Wallace and Hussain (1969); Nerlove (1971).

---

## HPQLIM Procedure

The following features have been added to the HPQLIM procedure:

- Discrete choice models. Most of the discrete choice models available in the QLIM procedure are now available in the HPQLIM procedure. The HPQLIM procedure also support the BY statement. The main features are as follows:
  - support for binary models: probit and logit models
  - support for ordinal models: ordered probit and ordered logit models
  - support for the BY statement

---

## HPSEVERITY Procedure

The following features have been added to the HPSEVERITY procedure:

- The HPSEVERITY procedure now supports the BY statement in the single-machine mode of execution. For the distributed mode, the BY statement is supported only in the client-data (or local-data) mode of the distributed computing model.
- The HPSEVERITY procedure now supports a new OFFSET= option in the SCALEMODEL statement to specify an offset variable in the scale regression model. An offset variable is a regressor that has a fixed regression coefficient of 1. An offset variable is useful for modeling the scale parameter per unit of some measure of exposure and for correcting the omitted variable bias.
- The HPSEVERITY procedure now supports a new OUTSCORELIB statement to create scoring functions. A scoring function is derived from a distribution function such that the parameter estimates of a fitted severity model are encoded inside the scoring function. You can use a scoring function to evaluate the corresponding distribution function without having to write a potentially complex SAS program to extract the parameter estimates. If you have specified a scale regression model, then you can use the scoring functions to score new observations by easily computing various statistics of interest, such as the mean and quantiles that depend on the regressor data. **NOTE:** This is an experimental feature.

The specification of custom objective function was an experimental feature in previous releases. It is now production.

---

## References

- Fuller, W. A. and Battese, G. E. (1974), “Estimation of Linear Models with Crossed-Error Structure,” *Journal of Econometrics*, 2, 67–78.
- Hausman, J. A. (1978), “Specification Tests in Econometrics,” *Econometrica*, 46, 1251–1271.
- Nerlove, M. (1971), “Further Evidence on the Estimation of Dynamic Relations from a Time Series of Cross Sections,” *Econometrica*, 39, 359–382.
- Wallace, T. and Hussain, A. (1969), “The Use of Error Components Model in Combining Cross Section with Time Series Data,” *Econometrica*, 37, 55–72.
- Wansbeek, T. and Kapteyn, A. (1989), “Estimation of the Error-Components Model with Incomplete Panels,” *Journal of Econometrics*, 41, 341–361.



# Chapter 2

## Introduction

### Contents

Overview of SAS/ETS High-Performance Procedures . . . . .	7
About This Book . . . . .	7
Chapter Organization . . . . .	7
Typographical Conventions . . . . .	8
Options Used in Examples . . . . .	9
Online Documentation . . . . .	9
SAS Technical Support Services . . . . .	9

---

## Overview of SAS/ETS High-Performance Procedures

SAS/ETS high-performance procedures provide econometric modeling tools that have been specially developed to take advantage of parallel processing in both multithreaded single-machine mode and distributed multiple-machine mode. Econometric modeling methods include regression for count data, models for the severity of losses or other events, and regression models for qualitative and limited dependent variables.

In addition to the high-performance econometric procedures described in this book, SAS/ETS includes high-performance utility procedures, which are described in *Base SAS Procedures Guide: High-Performance Procedures*. You can run all these procedures in single-machine mode without licensing SAS High-Performance Econometrics. However, to run these procedures in distributed mode, you must license SAS High-Performance Econometrics.

---

## About This Book

This book assumes that you are familiar with Base SAS software and with the books *SAS Language Reference: Concepts* and *Base SAS Procedures Guide*. It also assumes that you are familiar with basic SAS System concepts, such as using the DATA step to create SAS data sets and using Base SAS procedures (such as, the PRINT and SORT procedures) to manipulate SAS data sets.

---

## Chapter Organization

This book is organized as follows:

Chapter 2, this chapter, provides an overview of SAS/ETS high-performance procedures.

Chapter 3, “[Shared Concepts and Topics](#),” describes the modes in which SAS/ETS high-performance procedures can execute.

Subsequent chapters describe the individual procedures. These chapters appear in alphabetical order by procedure name. Each chapter is organized as follows:

- The “Overview” section provides a brief description of the analysis provided by the procedure.
- The “Getting Started” section provides a quick introduction to the procedure through a simple example.
- The “Syntax” section describes the SAS statements and options that control the procedure.
- The “Details” section discusses methodology and other topics, such as ODS tables.
- The “Examples” section contains examples that use the procedure.
- The “References” section contains references for the methodology.

---

## Typographical Conventions

This book uses several type styles for presenting information. The following list explains the meaning of the typographical conventions used in this book:

roman	is the standard type style used for most text.
UPPERCASE ROMAN	is used for SAS statements, options, and other SAS language elements when they appear in the text. However, you can enter these elements in your own SAS programs in lowercase, uppercase, or a mixture of the two.
<b>UPPERCASE BOLD</b>	is used in the “Syntax” sections’ initial lists of SAS statements and options.
<i>oblique</i>	is used in the syntax definitions and in text to represent arguments for which you supply a value.
VariableName	is used for the names of variables and data sets when they appear in the text.
<b>bold</b>	is used to for matrices and vectors.
<i>italic</i>	is used for terms that are defined in the text, for emphasis, and for references to publications.
<b>monospace</b>	is used for example code. In most cases, this book uses lowercase type for SAS code.



---

## Options Used in Examples

Most of the output shown in this book is produced with the following SAS System options:

```
options linesize=80 pagesize=500 nonumber nodate;
```

The HTMLBLUE style is used to create the HTML output and graphs that appear in the online documentation. A style template controls stylistic elements such as colors, fonts, and presentation attributes. The style template is specified in the ODS HTML statement as follows:

```
ods html style=HTMLBlue;
```

If you run the examples, your output might be slightly different, because of the SAS System options you use and the precision that your computer uses for floating-point calculations.

---

## Online Documentation

This documentation is available online with the SAS System. To access documentation for the SAS/ETS high-performance procedures from the SAS windowing environment, select **Help** from the main menu and then select **SAS Help and Documentation**. On the **Contents** tab, expand the **SAS Products, SAS/ETS, and SAS/ETS User's Guide: High-Performance Procedures** items. Then expand chapters and click on sections. You can search the documentation by using the **Search** tab.

You can also access the documentation by going to <http://support.sas.com/documentation>.

---

## SAS Technical Support Services

The SAS Technical Support staff is available to respond to problems and answer technical questions regarding the use of high-performance procedures. Go to <http://support.sas.com/techsup> for more information.



# Chapter 3

## Shared Concepts and Topics

### Contents

Overview . . . . .	11
Processing Modes . . . . .	12
Single-Machine Mode . . . . .	12
Distributed Mode . . . . .	12
Symmetric and Asymmetric Distributed Modes . . . . .	13
Controlling the Execution Mode with Environment Variables and Performance State- ment Options . . . . .	13
Determining Single-Machine Mode or Distributed Mode . . . . .	15
Alongside-the-Database Execution . . . . .	19
Alongside-LASR Distributed Execution . . . . .	22
Running High-Performance Analytical Procedures Alongside a SAS LASR Analytic Server in Distributed Mode . . . . .	22
Starting a SAS LASR Analytic Server Instance . . . . .	22
Associating a SAS Libref with the SAS LASR Analytic Server Instance . . . . .	23
Running a High-Performance Analytical Procedure Alongside the SAS LASR Ana- lytic Server Instance . . . . .	23
Terminating a SAS LASR Analytic Server Instance . . . . .	24
Alongside-LASR Distributed Execution on a Subset of the Appliance Nodes . . . . .	24
Running High-Performance Analytical Procedures in Asymmetric Mode . . . . .	25
Running in Symmetric Mode . . . . .	25
Running in Asymmetric Mode on One Appliance . . . . .	27
Running in Asymmetric Mode on Distinct Appliances . . . . .	27
Alongside-HDFS Execution . . . . .	30
Alongside-HDFS Execution by Using the SASHDAT Engine . . . . .	31
Alongside-HDFS Execution by Using the Hadoop Engine . . . . .	32
Output Data Sets . . . . .	36
Working with Formats . . . . .	37
PERFORMANCE Statement . . . . .	39

### Overview

This chapter describes the modes of execution in which SAS high-performance analytical procedures can execute. If you have SAS/ETS installed, you can run any procedure in this book on a single machine.

However, to run procedures in this book in distributed mode, you must also have SAS High-Performance Econometrics software installed. For more information about these modes, see the next section.

This chapter provides details of how you can control the modes of execution and includes the syntax for the PERFORMANCE statement, which is common to all high-performance analytical procedures.

---

## Processing Modes

---

### Single-Machine Mode

Single-machine mode is a computing model in which multiple processors or multiple cores are controlled by a single operating system and can access shared resources, such as disks and memory. In this book, single-machine mode refers to an application running multiple concurrent threads on a multicore machine in order to take advantage of parallel execution on multiple processing units. More simply, single-machine mode for high-performance analytical procedures means multithreading on the client machine.

All high-performance analytical procedures are capable of running in single-machine mode, and this is the default mode when a procedure runs on the client machine. The procedure uses the number of CPUs (cores) on the machine to determine the number of concurrent threads. High-performance analytical procedures use different methods to map core count to the number of concurrent threads, depending on the analytic task. Using one thread per core is not uncommon for the procedures that implement data-parallel algorithms.

---

### Distributed Mode

Distributed mode is a computing model in which several nodes in a distributed computing environment participate in the calculations. In this book, the distributed mode of a high-performance analytical procedure refers to the procedure performing the analytics on an appliance that consists of a cluster of nodes. This appliance can be one of the following:

- a database management system (DBMS) appliance on which the SAS High-Performance Analytics infrastructure is also installed
- a cluster of nodes that have the SAS High-Performance Analytics infrastructure installed but no DBMS software installed

Distributed mode has several variations:

- Client-data (or local-data) mode: The input data for the analytic task are not stored on the appliance or cluster but are distributed to the distributed computing environment by the SAS High-Performance Analytics infrastructure when the procedure runs.
- Alongside-the-database mode: The data are stored in the distributed database and are read from the DBMS in parallel into a high-performance analytical procedure that runs on the database appliance.

- Alongside-HDFS mode: The data are stored in the Hadoop Distributed File System (HDFS) and are read in parallel from the HDFS. This mode is available if you install the SAS High-Performance Deployment of Hadoop on the appliance or when you configure a Cloudera 4 Hadoop deployment on the appliance to operate with the SAS High-Performance Analytics infrastructure. For more information about installing the SAS High-Performance Deployment of Hadoop, see the *SAS High-Performance Analytics Infrastructure: Installation and Configuration Guide*.
- Alongside-LASR mode: The data are loaded from a SAS LASR Analytic Server that runs on the appliance.

---

## Symmetric and Asymmetric Distributed Modes

SAS high-performance analytical procedures can run alongside the database or alongside HDFS in asymmetric mode. The primary reason for providing the asymmetric mode is to enable you to manage and house data on one appliance (the data appliance) and to run the high-performance analytical procedure on a second appliance (the computing appliance). You can also run in asymmetric mode on a single appliance that functions as both the data appliance and the computing appliance. This enables you to run alongside the database or alongside HDFS, where computations are done on a different set of nodes from the nodes that contain the data. The following subsections provide more details.

### **Symmetric Mode**

When SAS high-performance analytical procedures run in symmetric distributed mode, the data appliance and the computing appliance must be the same appliance. Both the SAS Embedded Process and the high-performance analytical procedures execute in a SAS process that runs on the same hardware where the DBMS process executes. This is called symmetric mode because the number of nodes on which the DBMS executes is the same as the number of nodes on which the high-performance analytical procedures execute. The initial data movement from the DBMS to the high-performance analytical procedure does not cross node boundaries.

### **Asymmetric Mode**

When SAS high-performance analytical procedures run in asymmetric distributed mode, the data appliance and computing appliance are usually distinct appliances. The high-performance analytical procedures execute in a SAS process that runs on the computing appliance. The DBMS and a SAS Embedded Process run on the data appliance. Data are requested by a SAS data feeder that runs on the computing appliance and communicates with the SAS Embedded Process on the data appliance. The SAS Embedded Process transfers the data in parallel to the SAS data feeder that runs on each of the nodes of the computing appliance. This is called asymmetric mode because the number of nodes on the data appliance does not need to be the same as the number of nodes on the computing appliance.

---

## Controlling the Execution Mode with Environment Variables and Performance Statement Options

You control the execution mode by using environment variables or by specifying options in the **PERFORMANCE** statement in high-performance analytical procedures, or by a combination of these methods.

The important environment variables follow:

- *grid host* identifies the domain name system (DNS) or IP address of the appliance node to which the SAS High-Performance Econometrics software connects to run in distributed mode.
- *installation location* identifies the directory where the SAS High-Performance Econometrics software is installed on the appliance.
- *data server* identifies the database server on Teradata appliances as defined in the *hosts* file on the client. This data server is the same entry that you usually specify in the `SERVER=` entry of a `LIBNAME` statement for Teradata. For more information about specifying `LIBNAME` statements for Teradata and other engines, see the DBMS-specific section of *SAS/ACCESS for Relational Databases: Reference* for your engine.
- *grid mode* specifies whether the high-performance analytical procedures execute in symmetric or asymmetric mode. Valid values for this variable are `'sym'` for symmetric mode and `'asym'` for asymmetric mode. The default is symmetric mode.

You can set an environment variable directly from the SAS program by using the `OPTION SET=` command. For example, the following statements define three variables for a Teradata appliance (the grid mode is the default symmetric mode):

```
option set=GRIDHOST      ="hpa.sas.com";
option set=GRIDINSTALLLOC="/opt/TKGrid";
option set=GRIDDATASERVER="myserver";
```

Alternatively, you can set the parameters in the `PERFORMANCE` statement in high-performance analytical procedures. For example:

```
performance host      ="hpa.sas.com"
               install  ="/opt/TKGrid"
               dataserver="myserver";
```

The following statements define three variables that are needed to run asymmetrically on a computing appliance.

```
option set=GRIDHOST      ="compute_appliance.sas.com";
option set=GRIDINSTALLLOC="/opt/TKGrid";
option set=GRIDMODE      ="asym";
```

Alternatively, you can set the parameters in the `PERFORMANCE` statement in high-performance analytical procedures. For example:

```
performance host      ="compute_appliance.sas.com"
               install  ="/opt/TKGrid"
               gridmode  ="asym"
```

A specification in the `PERFORMANCE` statement overrides a specification of an environment variable without resetting its value. An environment variable that you set in the SAS session by using an `OPTION SET=` command remains in effect until it is modified or until the SAS session terminates.

Specifying a data server is necessary only on Teradata systems when you do not explicitly set the GRIDMODE environment variable or specify the GRIDMODE= option in the **PERFORMANCE** statement. The data server specification depends on the entries in the (client) *hosts* file. The file specifies the server (suffixed by *cop* and a number) and an IP address. For example:

```
myservercop1 33.44.55.66
```

The key variable that determines whether a high-performance analytical procedure executes in single-machine or distributed mode is the *grid host*. The installation location and data server are needed to ensure that a connection to the grid host can be made, given that a host is specified. This book assumes that the installation location and data server (if necessary) have been set by your system administrator.

The following sets of SAS statements are functionally equivalent:

```
proc hpreduce;
  reduce unsupervised x;;
  performance host="hpa.sas.com";
run;

option set=GRIDHOST="hpa.sas.com";
proc hpreduce;
  reduce unsupervised x;;
run;
```

---

## Determining Single-Machine Mode or Distributed Mode

High-performance analytical procedures use the following rules to determine whether they run in single-machine mode or distributed mode:

- If a grid host is not specified, the analysis is carried out in single-machine mode on the client machine that runs the SAS session.
- If a grid host is specified, the behavior depends on whether the execution is alongside the database or alongside HDFS. If the data are local to the client (that is, not stored in the distributed database or HDFS on the appliance), you need to use the **NODES=** option in the **PERFORMANCE** statement to specify the number of nodes on the appliance or cluster that you want to engage in the analysis. If the procedure executes alongside the database or alongside HDFS, you do not need to specify the **NODES=** option.

The following example shows single-machine and client-data distributed configurations for a data set of 100,000 observations that are simulated from a logistic regression model. The following DATA step generates the data:

```
data simData;
  array _a{8} _temporary_ (0,0,0,1,0,1,1,1);
  array _b{8} _temporary_ (0,0,1,0,1,0,1,1);
  array _c{8} _temporary_ (0,1,0,0,1,1,0,1);
```

```

do obsno=1 to 100000;
  x  = rantbl(1,0.28,0.18,0.14,0.14,0.03,0.09,0.08,0.06);
  a  = _a{x};
  b  = _b{x};
  c  = _c{x};
  x1 = int(ranuni(1)*400);
  x2 = 52 + ranuni(1)*38;
  x3 = ranuni(1)*12;
  lp = 6. -0.015*(1-a) + 0.7*(1-b) + 0.6*(1-c) + 0.02*x1 -0.05*x2 - 0.1*x3;
  y  = ranbin(1,1,(1/(1+exp(lp))));
  output;
end;
drop x lp;
run;

```

The following statements run PROC HPLOGISTIC to fit a logistic regression model:

```

proc hplogistic data=simData;
  class a b c;
  model y = a b c x1 x2 x3;
run;

```

Figure 3.1 shows the results from the analysis.

**Figure 3.1** Results from Logistic Regression in Single-Machine Mode

The HPLOGISTIC Procedure	
Performance Information	
Execution Mode	Single-Machine
Number of Threads	4
Model Information	
Data Source	WORK.SIMDATA
Response Variable	y
Class Parameterization	GLM
Distribution	Binary
Link Function	Logit
Optimization Technique	Newton-Raphson with Ridging



**Figure 3.1** *continued*

Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

The entries in the “Performance Information” table show that the HPLOGISTIC procedure runs in single-machine mode and uses four threads, which are chosen according to the number of CPUs on the client machine. You can force a certain number of threads on any machine that is involved in the computations by specifying the **NTHREADS** option in the **PERFORMANCE** statement. Another indication of execution on the client is the following message, which is issued in the SAS log by all high-performance analytical procedures:

**NOTE: The HPLOGISTIC procedure is executing on the client.**

The following statements use 10 nodes (in distributed mode) to analyze the data on the appliance; results appear in [Figure 3.2](#):

```
proc hplogistic data=simData;
  class a b c;
  model y = a b c x1 x2 x3;
  performance host="hpa.sas.com" nodes=10;
run;
```

**Figure 3.2** Results from Logistic Regression in Distributed Mode

The HPLOGISTIC Procedure	
Performance Information	
Host Node	hpa.sas.com
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	10
Number of Threads per Node	24

Figure 3.2 continued

Model Information					
Data Source	WORK.SIMDATA				
Response Variable	y				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				

Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

The specification of a host causes the “Performance Information” table to display the name of the host node of the appliance. The “Performance Information” table also indicates that the calculations were performed in a distributed environment on the appliance. Twenty-four threads on each of 10 nodes were used to perform the calculations—for a total of 240 threads.

Another indication of distributed execution on the appliance is the following message, which is issued in the SAS log by all high-performance analytical procedures:

**NOTE: The HPLOGISTIC procedure is executing in the distributed computing environment with 10 worker nodes.**

You can override the presence of a grid host and force the computations into single-machine mode by specifying the `NODES=0` option in the `PERFORMANCE` statement:

```
proc hplogistic data=simData;
  class a b c;
  model y = a b c x1 x2 x3;
  performance host="hpa.sas.com" nodes=0;
run;
```

Figure 3.3 shows the “Performance Information” table. The numeric results are not reproduced here, but they agree with the previous analyses, which are shown in Figure 3.1 and Figure 3.2.

**Figure 3.3** Single-Machine Mode Despite Host Specification

The HPLOGISTIC Procedure	
Performance Information	
Execution Mode	Single-Machine
Number of Threads	4

The “Performance Information” table indicates that the HPLOGISTIC procedure executes in single-machine mode on the client. This information is also reported in the following message, which is issued in the SAS log:

**NOTE: The HPLOGISTIC procedure is executing on the client.**

In the analysis shown previously in [Figure 3.2](#), the data set `Work.simData` is local to the client, and the HPLOGISTIC procedure distributed the data to 10 nodes on the appliance. The High-Performance Analytics infrastructure does not keep these data on the appliance. When the procedure terminates, the in-memory representation of the input data on the appliance is freed.

When the input data set is large, the time that is spent sending client-side data to the appliance might dominate the execution time. In practice, transfer speeds are usually lower than the theoretical limits of the network connection or disk I/O rates. At a transfer rate of 40 megabytes per second, sending a 10-gigabyte data set to the appliance requires more than four minutes. If analytic execution time is in the range of seconds, the “performance” of the process is dominated by data movement.

The alongside-the-database execution model, unique to high-performance analytical procedures, enables you to read and write data in distributed form from the database that is installed on the appliance.

---

## Alongside-the-Database Execution

High-performance analytical procedures interface with the distributed database management system (DBMS) on the appliance in a unique way. If the input data are stored in the DBMS and the grid host is the appliance that houses the data, high-performance analytical procedures create a distributed computing environment in which an analytic process is co-located with the nodes of the DBMS. Data then pass from the DBMS to the analytic process on each node. Instead of moving across the network and possibly back to the client machine, the data pass locally between the processes on each node of the appliance.

Because the analytic processes on the appliance are separate from the database processes, the technique is referred to as alongside-the-database execution in contrast to in-database execution, where the analytic code executes in the database process.

In general, when you have a large amount of input data, you can achieve the best performance from high-performance analytical procedures if execution is alongside the database.

Before you can run alongside the database, you must distribute the data to the appliance. The following statements use the HPDS2 procedure to distribute the data set `Work.simData` into the `mydb` database on the `hpa.sas.com` appliance. In this example, the appliance houses a Greenplum database.

```
option set=GRIDHOST="hpa.sas.com";
libname applianc greenplm
      server  ="hpa.sas.com"
      user    =XXXXXX
      password=YYYYY
      database=mydb;

proc datasets lib=applianc nolist; delete simData;
proc hpds2 data=simData
      out =applianc.simData(distributed_by='distributed randomly');
  performance commit=10000 nodes=all;
  data DS2GTF.out;
    method run();
    set DS2GTF.in;
  end;
enddata;
run;
```

If the output table `applianc.simData` exists, the DATASETS procedure removes the table from the Greenplum database because a DBMS does not usually support replacement operations on tables.

Note that the libref for the output table points to the appliance. The data set option informs the HPDS2 procedure to distribute the records randomly among the data segments of the appliance. The statements that follow the **PERFORMANCE** statement are the DS2 program that copies the input data to the output data without further transformations.

Because you loaded the data into a database on the appliance, you can use the following HPLOGISTIC statements to perform the analysis on the appliance in the alongside-the-database mode. These statements are almost identical to the first PROC HPLOGISTIC example in the previous section, which executed in single-machine mode.

```
proc hplogistic data=applianc.simData;
  class a b c;
  model y = a b c x1 x2 x3;
run;
```

The subtle differences are as follows:

- The grid host environment variable that you specified in an `OPTION SET=` command is still in effect.
- The `DATA=` option in the high-performance analytical procedure uses a libref that identifies the data source as being housed on the appliance. This libref was specified in a prior `LIBNAME` statement.

Figure 3.4 shows the results from this analysis. The “Performance Information” table shows that the execution was in distributed mode. In this case the execution was alongside the Greenplum database. The numeric results agree with the previous analyses, which are shown in Figure 3.1 and Figure 3.2.

**Figure 3.4** Alongside-the-Database Execution on Greenplum

The HPLOGISTIC Procedure					
Performance Information					
Host Node	hpa.sas.com				
Execution Mode	Distributed				
Grid Mode	Symmetric				
Number of Compute Nodes	8				
Number of Threads per Node	24				
Model Information					
Data Source	SIMDATA				
Response Variable	y				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				
Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

When high-performance analytical procedures execute symmetrically alongside the database, any nonzero specification of the **NODES=** option in the **PERFORMANCE** statement is ignored. If the data are read alongside the database, the number of compute nodes is determined by the layout of the database and cannot be modified. In this example, the appliance contains 16 nodes. (See the “Performance Information” table.)

However, when high-performance analytical procedures execute asymmetrically alongside the database, the number of compute nodes that you specify in the **PERFORMANCE** statement can differ from the number of nodes across which the data are partitioned. For an example, see the section “[Running High-Performance Analytical Procedures in Asymmetric Mode](#)” on page 25.

---

## Alongside-LASR Distributed Execution

You can execute high-performance analytical procedures in distributed mode alongside a SAS LASR Analytic Server. When high-performance analytical procedures execute in this mode, the data are preloaded in distributed form in memory that is managed by a LASR Analytic Server. The data on the nodes of the appliance are accessed in parallel in the process that runs the LASR Analytic Server, and they are transferred to the process where the high-performance analytical procedure runs. In general, each high-performance analytical procedure copies the data to memory that persists only while that procedure executes. Hence, when a high-performance analytical procedure runs alongside a LASR Analytic Server, both the high-performance analytical procedure and the LASR Analytic Server have a copy of the subset of the data that is used by the high-performance analytical procedure. The advantage of running high-performance analytical procedures alongside a LASR Analytic Server (as opposed to running alongside a DBMS table or alongside HDFS) is that the initial transfer of data from the LASR Analytic Server to the high-performance analytical procedure is a memory-to-memory operation that is faster than the disk-to-memory operation when the procedure runs alongside a DBMS or HDFS. When the cost of preloading a table into a LASR Analytic Server is amortized by multiple uses of these data in separate runs of high-performance analytical procedures, using the LASR Analytic Server can result in improved performance.

---

## Running High-Performance Analytical Procedures Alongside a SAS LASR Analytic Server in Distributed Mode

This section provides an example of steps that you can use to start and load data into a SAS LASR Analytic Server instance and then run high-performance analytical procedures alongside this LASR Analytic Server instance.

---

### Starting a SAS LASR Analytic Server Instance

The following statements create a SAS LASR Analytic Server instance and load it with the `simData` data set that is used in the preceding examples. The data that are loaded into the LASR Analytic Server persist in memory across procedure boundaries until these data are explicitly deleted or until the server instance is terminated.

```
proc lasr port=12345
      data=simData
      path="/tmp/";
      performance host="hpa.sas.com" nodes=ALL;
run;
```

The `PORT=` option specifies a network port number to use. The `PATH=` option specifies the directory in which the server and table signature files are to be stored. The specified directory must exist on each machine in the cluster. The `DATA=` option specifies the name of a data set that is loaded into this LASR Analytic Server instance. (You do not need to specify the `DATA=` option at this time because you can add tables to

the LASR Analytic Server instance at any stage of its life.) For more information about starting and using a LASR Analytic Server, see the *SAS LASR Analytic Server: Administration Guide*.

The `NODES=ALL` option in the **PERFORMANCE** statement specifies that the LASR Analytic Server run on all the nodes on the appliance. You can start a LASR Analytic Server on a subset of the nodes on an appliance, but this might affect whether high-performance analytical procedures can run alongside the LASR Analytic Server. For more information, see the section “[Alongside-LASR Distributed Execution on a Subset of the Appliance Nodes](#)” on page 24.

Figure 3.5 shows the “Performance Information” table, which shows that the LASR procedure executes in distributed mode on 16 nodes.

**Figure 3.5** Performance Information

The LASR Procedure	
Performance Information	
Host Node	hpa.sas.com
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	8

---

## Associating a SAS Libref with the SAS LASR Analytic Server Instance

The following statements use a `LIBNAME` statement that associates a SAS libref (named `MyLasr`) with tables on the server instance as follows:

```
libname MyLasr sasiola port=12345;
```

The `SASIOLA` option requests that the `MyLasr` libref use the `SASIOLA` engine, and the `PORT=` value associates this libref with the appropriate server instance. For more information about creating a libref that uses the `SASIOLA` engine, see the *SAS LASR Analytic Server: Administration Guide*.

---

## Running a High-Performance Analytical Procedure Alongside the SAS LASR Analytic Server Instance

You can use the `MyLasr` libref to specify the input data for high-performance analytical procedures. You can also create output data sets in the SAS LASR Analytic Server instance by using this libref to request that the output data set be held in memory by the server instance as follows:

```
proc hplogistic data=MyLasr.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  output out=MyLasr.simulateScores pred=PredictedProbability;
run;
```

Because you previously specified the GRIDHOST= environment variable and the input data are held in distributed form in the associated server instance, this PROC HPLOGISTIC step runs in distributed mode alongside the LASR Analytic Server, as indicated in the “Performance Information” table shown in [Figure 3.6](#).

**Figure 3.6** Performance Information

Performance Information	
Host Node	hpa.sas.com
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	8
Number of Threads per Node	24

The preceding OUTPUT statement creates an output table that is added to the LASR Analytic Server instance. Output data sets do not have to be created in the same server instance that holds the input data. You can use a different LASR Analytic Server instance to hold the output data set. However, in order for the output data to be created in alongside mode, all the nodes that are used by the server instance that holds the input data must also be used by the server instance that holds the output data.

---

## Terminating a SAS LASR Analytic Server Instance

You can continue to run high-performance analytical procedures and add and delete tables from the SAS LASR Analytic Server instance until you terminate the server instance as follows:

```
proc lasr term port=12345;
run;
```

---

## Alongside-LASR Distributed Execution on a Subset of the Appliance Nodes

When you run PROC LASR to start a SAS LASR Analytic Server, you can specify the NODES= option in a **PERFORMANCE** statement to control how many nodes the LASR Analytic Server executes on. Similarly, a high-performance analytical procedure can execute on a subset of the nodes either because you specify the NODES= option in a **PERFORMANCE** statement or because you run alongside a DBMS or HDFS with an input data set that is distributed on a subset of the nodes on an appliance. In such situations, if a high-performance analytical procedure uses nodes on which the LASR Analytic Server is not running, then running alongside LASR is not supported. You can avoid this issue by specifying the NODES=ALL in the **PERFORMANCE** statement when you use PROC LASR to start the LASR Analytic Server.



## Running High-Performance Analytical Procedures in Asymmetric Mode

This section provides examples of how you can run high-performance analytical procedures in asymmetric mode. It also includes examples that run in symmetric mode to highlight differences between the modes. For a description of asymmetric mode, see the section “Symmetric and Asymmetric Distributed Modes” on page 13.

Asymmetric mode is commonly used when the data appliance and the computing appliance are distinct appliances. In order to be able to use an appliance as a data provider for high-performance analytical procedures that run in asymmetric mode on another appliance, it is not necessary that SAS High-Performance Econometrics be installed on the data appliance. However, it is essential that a SAS Embedded Process be installed on the data appliance and that SAS High-Performance Econometrics be installed on the computing appliance.

The following examples use a 24-node data appliance named “data\_appliance.sas.com,” which houses a Teradata DBMS and has a SAS Embedded Process installed. Because SAS High-Performance Econometrics is also installed on this appliance, it can be used to run high-performance analytical procedures in both symmetric and asymmetric modes.

The following statements load the simData data set of the preceding sections onto the data appliance:

```
libname dataLib teradata
    server    ="tera2650"
    user      =XXXXXX
    password=YYYYY
    database=mydb;

data dataLib.simData;
    set simData;
run;
```

**NOTE:** You can provision the appliance with data even if SAS High-Performance Econometrics software is not installed on the appliance.

The following subsections show how you can run the HPLOGISTIC procedure symmetrically and asymmetrically on a single data appliance and asymmetrically on distinct data and computing appliances.

---

## Running in Symmetric Mode

The following statements run the HPLOGISTIC procedure in symmetric mode on the data appliance:

```
proc hplogistic data=dataLib.simData;
    class a b c;
    model y = a b c x1 x2 x3;
    performance host      = "data_appliance.sas.com"
                nodes      = 10
                gridmode    = sym;
```

```
run;
```

Because you explicitly specified the `GRIDMODE=` option, you do not need to also specify the `DATASERVER=` option in the `PERFORMANCE` statement. Figure 3.7 shows the results of this analysis.

**Figure 3.7** Alongside-the-Database Execution in Symmetric Mode on Teradata

The HPLOGISTIC Procedure					
Performance Information					
Host Node	data_appliance.sas.com				
Execution Mode	Distributed				
Grid Mode	Symmetric				
Number of Compute Nodes	24				
Number of Threads per Node	24				
Model Information					
Data Source	simData				
Response Variable	y				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				
Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

The “Performance Information” table shows that the execution occurs in symmetric mode on the 24 nodes of the data appliance. In this case, the `NODES=10` option in the `PERFORMANCE` statement is ignored because the number of nodes that are used is determined by the number of nodes across which the data are distributed, as indicated in the following warning message in the SAS log:

```
WARNING: The NODES=10 option in the PERFORMANCE statement is ignored because
you are running alongside the distributed data source
DATALIB.simData.DATA. The number of compute nodes is determined by the
configuration of the distributed DBMS.
```

## Running in Asymmetric Mode on One Appliance

You can switch to running the HPLOGISTIC procedure in asymmetric mode by specifying the GRIDMODE=ASYM option in the **PERFORMANCE** statement as follows:

```
proc hplogistic data=dataLib.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  performance host      = "data_appliance.sas.com"
              nodes      = 10
              gridmode    = asym;
run;
```

Figure 3.8 shows the “Performance Information” table.

**Figure 3.8** Alongside Teradata Execution in Asymmetric Mode

The HPLOGISTIC Procedure	
Performance Information	
Host Node	data_appliance.sas.com
Execution Mode	Distributed
Grid Mode	Asymmetric
Number of Compute Nodes	10
Number of Threads per Node	24

You can see that now the grid mode is asymmetric. Furthermore, the NODES=10 option that you specified in the **PERFORMANCE** statement is honored. The data are moved in parallel from the 24 nodes on which the data are stored to the 10 nodes on which the execution occurs. The numeric results are not reproduced here, but they agree with the previous analyses.

## Running in Asymmetric Mode on Distinct Appliances

Usually, there is no advantage to executing high-performance analytical procedures in asymmetric mode on one appliance, because data might have to be unnecessarily moved between nodes. The following example demonstrates the more typical use of asymmetric mode. In this example, the specified grid host “compute\_appliance.sas.com” is a computing appliance that has 15 compute nodes, and it is a different appliance from the 24-node data appliance “data\_appliance.sas.com,” which houses the Teradata DBMS where the data reside.

The advantage of using different computing and data appliances is that the data appliance is not affected by the execution of high-performance analytical procedures except during the initial parallel data transfer. A potential disadvantage of this asymmetric mode of execution is that the performance can be limited by the bandwidth with which data can be moved between the appliances. However, because this data movement takes place in parallel from the nodes of the data appliance to the nodes of the computing appliance, this

potential performance bottleneck can be overcome with appropriately provisioned hardware. The following statements show how this is done:

```
proc hplogistic data=dataLib.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  performance host      = "compute_appliance.sas.com"
               gridmode = asym;
run;
```

Figure 3.9 shows the “Performance Information” table.

**Figure 3.9** Asymmetric Mode with Distinct Data and Computing Appliances

The HPLOGISTIC Procedure	
Performance Information	
Host Node	compute_appliance.sas.com
Execution Mode	Distributed
Grid Mode	Asymmetric
Number of Compute Nodes	15
Number of Threads per Node	24

PROC HPLOGISTIC ran on the 15 nodes of the computing appliance, even though the data are partitioned across the 24 nodes of the data appliance. The numeric results are not reproduced here, but they agree with the previous analyses shown in Figure 3.1 and Figure 3.2.

Every time you run a high-performance analytical procedure in asymmetric mode that uses different computing and data appliances, data are transferred between these appliances. If you plan to make repeated use of the same data, then it might be advantageous to temporarily persist the data that you need on the computing appliance. One way to persist the data is to store them as a table in a SAS LASR Analytic Server that runs on the computing appliance. By running PROC LASR in asymmetric mode, you can load the data in parallel from the data appliance nodes to the nodes on which the LASR Analytic Server runs on the computing appliance. You can then use a LIBNAME statement that associates a SAS libref with tables on the LASR Analytic Server. The following statements show how you do this:

```
proc lasr port=54321
  data=dataLib.simData
  path="/tmp/";
  performance host      = "compute_appliance.sas.com"
               gridmode = asym;
run;

libname MyLasr sasiola tag="dataLib" port=54321 host="compute_appliance.sas.com" ;
```

Figure 3.10 show the “Performance Information” table.

**Figure 3.10** PROC LASR Running in Asymmetric Mode

The LASR Procedure	
Performance Information	
Host Node	compute_appliance.sas.com
Execution Mode	Distributed
Grid Mode	Asymmetric
Number of Compute Nodes	15

PROC LASR ran in asymmetric mode on the computing appliance, which has 15 compute nodes. In this mode, the data are loaded in parallel from the 24 data appliance nodes to the 15 compute nodes on the computing appliance. By default, all the nodes on the computing appliance are used. You can use the `NODES=` option in the `PERFORMANCE` statement to run the LASR Analytic Server on a subset of the nodes on the computing appliance. If you omit the `GRIDMODE=ASYM` option from the `PERFORMANCE` statement, PROC LASR still runs successfully but much less efficiently. The Teradata access engine transfers the `simData` data set to a temporary table on the client, and the High-Performance Analytics infrastructure then transfers these data from the temporary table on the client to the grid nodes on the computing appliance.

After the data are loaded into a LASR Analytic Server that runs on the computing appliance, you can run high-performance analytical procedures alongside this LASR Analytic Server. Because these procedures run on the same computing appliance where the LASR Analytic Server is running, it is best to run these procedures in symmetric mode, which is the default or can be explicitly specified in the `GRIDMODE=SYM` option in the `PERFORMANCE` statement. The following statements provide an example. The `OUTPUT` statement creates an output data set that is held in memory by the LASR Analytic Server. The data appliance has no role in executing these statements.

```
proc hplogistic data=MyLasr.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  output out=MyLasr.myOutputData pred=myPred;
  performance host = "compute_appliance.sas.com";
run;
```

The following note, which appears in the SAS log, confirms that the output data set is created successfully:

**NOTE:** The table `DATALIB.MYOUTPUTDATA` has been added to the LASR Analytic Server with port 54321. The Libname is `MYLASR`.

You can use the `dataLib` libref that you used to load the data onto the data appliance to create an output data set on the data appliance. In order for this output to be directly written in parallel from the nodes of the computing appliance to the nodes of the data appliance, you need to run the `HPLOGISTIC` procedure in asymmetric mode by specifying the `GRIDMODE=ASYM` option in the `PERFORMANCE` statement as follows:

```
proc hplogistic data=MyLasr.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  output out=dataLib.myOutputData pred=myPred;
  performance host      = "compute_appliance.sas.com"
              gridmode   = asym;
run;
```

The following note, which appears in the SAS log, confirms that the output data set is created successfully on the data appliance:

**NOTE: The data set DATALIB.myOutputData has 100000 observations and 1 variables.**

When you run a high-performance analytical procedure on a computing appliance and either read data from or write data to a different data appliance, it is important to run the high-performance analytical procedures in asymmetric mode so that the Read and Write operations take place in parallel without any movement of data to and from the SAS client. If you omit running the preceding PROC HPLOGISTIC step in asymmetric mode, then the output data set would be created much less efficiently: the output data would be moved sequentially to a temporary table on the client, after which the Teradata access engine sequentially would write this table to the data appliance.

When you no longer need the data in the SAS LASR Analytic Server, you should terminate the server instance as follows:

```
proc lasr term port=54321;
  performance host="compute_appliance.sas.com";
run;
```

If you configured Hadoop on the computing appliance, then you can create output data tables that are stored in the HDFS on the computing appliance. You can do this by using the SASHDAT engine as described in the section “[Alongside-HDFS Execution](#)” on page 30.

---

## Alongside-HDFS Execution

Running high-performance analytical procedures alongside HDFS shares many features with running alongside the database. You can execute high-performance analytical procedures alongside HDFS by using either the SASHDAT engine or the Hadoop engine.

You use the SASHDAT engine to read and write data that are stored in HDFS in a proprietary SASHDAT format. In SASHDAT format, metadata that describe the data in the Hadoop files are included with the data. This enables you to access files in SASHDAT format without supplying any additional metadata. Additionally, you can also use the SASHDAT engine to read data in CSV (comma-separated value) format, but you need supply metadata that describe the contents of the CSV data. The SASHDAT engine provides highly optimized access to data in HDFS that are stored in SASHDAT format.

The Hadoop engine reads data that are stored in various formats from HDFS and writes data to HDFS in CSV format. This engine can use metadata that are stored in Hive, which is a data warehouse that supplies

metadata about data that are stored in Hadoop files. In addition, this engine can use metadata that you create by using the HDMD procedure.

The following subsections provide details about using the SASHDAT and Hadoop engines to execute high-performance analytical procedures alongside HDFS.

---

## Alongside-HDFS Execution by Using the SASHDAT Engine

If the grid host is a cluster that houses data that have been distributed by using the SASHDAT engine, then high-performance analytical procedures can analyze those data in the alongside-HDFS mode. The procedures use the distributed computing environment in which an analytic process is co-located with the nodes of the cluster. Data then pass from HDFS to the analytic process on each node of the cluster.

Before you can run a procedure alongside HDFS, you must distribute the data to the cluster. The following statements use the SASHDAT engine to distribute to HDFS the `simData` data set that was used in the previous two sections:

```
option set=GRIDHOST="hpa.sas.com";

libname hdatLib sashdat
      path="/hps";

data hdatLib.simData (replace = yes) ;
  set simData;
run;
```

In this example, the `GRIDHOST` is a cluster where the SAS Data in HDFS Engine is installed. If a data set that is named `simData` already exists in the `hps` directory in HDFS, it is overwritten because the `REPLACE=YES` data set option is specified. For more information about using this `LIBNAME` statement, see the section “`LIBNAME` Statement for the SAS Data in HDFS Engine” in the *SAS LASR Analytic Server: Administration Guide*.

The following `HPLOGISTIC` procedure statements perform the analysis in alongside-HDFS mode. These statements are almost identical to the `PROC HPLOGISTIC` example in the previous two sections, which executed in single-machine mode and alongside-the-database distributed mode, respectively.

```
proc hplogistic data=hdatLib.simData;
  class a b c;
  model y = a b c x1 x2 x3;

run;
```

Figure 3.11 shows the “Performance Information” table. You see that the procedure ran in distributed mode. The numeric results shown in Figure 3.12 agree with the previous analyses shown in Figure 3.1, Figure 3.2, and Figure 3.4.

**Figure 3.11** Alongside-HDFS Execution Performance Information

Performance Information	
Host Node	hpa.sas.com
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	206
Number of Threads per Node	8

**Figure 3.12** Alongside-HDFS Execution Model Information

Model Information					
Data Source	HDATLIB.SIMDATA				
Response Variable	y				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				

Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

## Alongside-HDFS Execution by Using the Hadoop Engine

The following LIBNAME statement sets up a libref that you can use to access data that are stored in HDFS and have metadata in Hive:

```
libname hdoopLib hadoop
server      = "hpa.sas.com"
user        = XXXXX
password    = YYYYY
database    = myDB
config      = "demo.xml" ;
```



For more information about LIBNAME options available for the Hadoop engine, see the LIBNAME topic in the Hadoop section of *SAS/ACCESS for Relational Databases: Reference*. The configuration file that you specify in the CONFIG= option contains information that is needed to access the Hive server. It also contains information that enables this configuration file to be used to access data in HDFS without using the Hive server. This information can also be used to specify replication factors and block sizes that are used when the engine writes data to HDFS. The following XML shows the contents of the file demo.xml that is used in this example:

```
<configuration>
  <property>
    <name>fs.default.name</name>
    <value>hdfs://hpa.sas.com:8020</value>
  </property>
  <property>
    <name>mapred.job.tracker</name>
    <value>hpa.sas.com:8021</value>
  </property>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>
  <property>
    <name>dfs.block.size</name>
    <value>33554432</value>
  </property>
</configuration>
```

The following DATA step uses the Hadoop engine to distribute to HDFS the simData data set that was used in the previous sections. The engine creates metadata for the data set in Hive.

```
data hdoopLib.simData;
  set simData;
run;
```

After you have loaded data or if you are accessing preexisting data in HDFS that have metadata in Hive, you can access this data alongside HDFS by using high-performance analytics procedures. The following HPLOGISTIC procedure statements perform the analysis in alongside-HDFS mode. These statements are similar to the PROC HPLOGISTIC example in the previous sections. However, whenever you use the Hadoop engine, you must execute the analysis in asymmetric mode to cause the execution to occur alongside HDFS.

```
proc hplogistic data=hdoopLib.simData;
  class a b c;
  model y = a b c x1 x2 x3;
  performance host      = "compute_appliance.sas.com"
                gridmode = asym;
run;
```

Figure 3.13 shows the “Performance Information” table. You see that the procedure ran asymmetrically in distributed mode. The numeric results shown in Figure 3.14 agree with the previous analyses.

**Figure 3.13** Alongside-HDFS Execution by Using the Hadoop Engine

The HPLOGISTIC Procedure	
Performance Information	
Host Node	compute_appliance.sas.com
Execution Mode	Distributed
Grid Mode	Asymmetric
Number of Compute Nodes	15
Number of Threads per Node	24

**Figure 3.14** Alongside-HDFS Execution by Using the Hadoop Engine

Model Information					
Data Source	HDOOPLIB.SIMDATA				
Response Variable	y				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				
Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	5.7011	0.2539	Infty	22.45	<.0001
a 0	-0.01020	0.06627	Infty	-0.15	0.8777
a 1	0	.	.	.	.
b 0	0.7124	0.06558	Infty	10.86	<.0001
b 1	0	.	.	.	.
c 0	0.8036	0.06456	Infty	12.45	<.0001
c 1	0	.	.	.	.
x1	0.01975	0.000614	Infty	32.15	<.0001
x2	-0.04728	0.003098	Infty	-15.26	<.0001
x3	-0.1017	0.009470	Infty	-10.74	<.0001

The Hadoop engine also enables you to access tables in HDFS that are stored in various formats and that are not registered in Hive. You can use the HDMD procedure to generate metadata for tables that are stored in the following file formats:

- delimited text
- fixed-record length binary
- sequence files
- XML text

To read any other kind of file in Hadoop, you can write a custom file reader plug-in in Java for use with PROC HDMD. For more information about LIBNAME options available for the Hadoop engine, see the LIBNAME topic in the Hadoop section of *SAS/ACCESS for Relational Databases: Reference*.

The following example shows how you can use PROC HDMD to register metadata for CSV data independently from Hive and then analyze these data by using high-performance analytics procedures. The CSV data in the table csvExample.csv is stored in HDFS in the directory /user/demo/data. Each record in this table consists of the following fields, in the order shown and separated by commas.

1. a string of at most six characters
2. a numeric field with values of 0 or 1
3. a numeric field with real numbers

Suppose you want to fit a logistic regression model to these data, where the second field represents a target variable named Success, the third field represents a regressor named Dose, and the first field represents a classification variable named Group.

The first step is to use PROC HDMD to create metadata that are needed to interpret the table, as in the following statements:

```
libname hdoopLib hadoop
    server      = "hpa.sas.com"
    user        = XXXXX
    password    = YYYY
    HDFS_PERMDIR = "/user/demo/data"
    HDFS_METADIR = "/user/demo/meta"
    config      = "demo.xml"
    DBCREATE_TABLE_EXTERNAL=YES;

proc hdmd name=hdoopLib.csvExample data_file='csvExample.csv'
    format=delimited encoding=utf8 sep = ',';

    column Group    char(6);
    column Success  double;
    column Dose     double;
run;
```

The metadata that are created by PROC HDMD for this table are stored in the directory /user/demo/meta that you specified in the HDFS\_METADIR = option in the preceding LIBNAME statement. After you create the metadata, you can execute high-performance analytics procedures with these data by using the hdoopLib libref. For example, the following statements fit a logistic regression model to the CSV data that are stored in csvExample.csv table.

```
proc hplogistic data=hdoopLib.csvExample;
    class Group;
    model Success = Dose;
    performance host      = "compute_appliance.sas.com"
                  gridmode = asym;
run;
```

Figure 3.15 shows the results of this analysis. You see that the procedure ran asymmetrically in distributed mode. The metadata that you created by using the HDMD procedure have been used successfully in executing this analysis.

**Figure 3.15** Alongside-HDFS Execution with CSV Data

The HPLOGISTIC Procedure					
Performance Information					
Host Node	compute_appliance.sas.com				
Execution Mode	Distributed				
Grid Mode	Asymmetric				
Number of Compute Nodes	15				
Number of Threads per Node	24				
Model Information					
Data Source	GRIDLIB.CSVEXAMPLE				
Response Variable	Success				
Class Parameterization	GLM				
Distribution	Binary				
Link Function	Logit				
Optimization Technique	Newton-Raphson with Ridging				
Class Level Information					
Class	Levels	Values			
Group	3	group1 group2 group3			
Number of Observations Read		1000			
Number of Observations Used		1000			
Parameter Estimates					
Parameter	Estimate	Standard Error	DF	t Value	Pr >  t
Intercept	0.1243	0.1295	Infty	0.96	0.3371
Dose	-0.2674	0.2216	Infty	-1.21	0.2277

## Output Data Sets

In the alongside-the-database mode, the data are read in distributed form, minimizing data movement for best performance. Similarly, when you write output data sets and a high-performance analytical procedure executes in distributed mode, the data can be written in parallel into the database.

For example, in the following statements, the HPLOGISTIC procedure executes in distributed mode by using eight nodes on the appliance to perform the logistic regression on `work.simData`:

```
proc hplogistic data=simData;
  class a b c;
  model y = a b c x1 x2 x3;
  id a;
  output out=applianc.simData_out pred=p;
  performance host="hpa.sas.com" nodes=8;
run;
```

The output data set `applianc.simData_out` is written in parallel into the database. Although the data are fed on eight nodes, the database might distribute the data on more nodes.

When a high-performance analytical procedure executes in single-machine mode, all output objects are created on the client. If the `libref` of the output data sets points to the appliance, the data are transferred to the database on the appliance. This can lead to considerable performance degradation compared to execution in distributed mode.

Many procedures in SAS software add the variables from the input data set when an observationwise output data set is created. The assumption of high-performance analytical procedures is that the input data sets can be large and contain many variables. For performance reasons, the output data set contains the following:

- variables that are explicitly created by the statement
- variables that are listed in the `ID` statement
- distribution keys or hash keys that are transferred from the input data set

Including this information enables you to add to the output data set information necessary for subsequent SQL joins without copying the entire input data set to the output data set.

---

## Working with Formats

You can use SAS formats and user-defined formats with high-performance analytical procedures as you can with other procedures in the SAS System. However, because the analytic work is carried out in a distributed environment and might depend on the formatted values of variables, some special handling can improve the efficiency of work with formats.

High-performance analytical procedures examine the variables that are used in an analysis for association with user-defined formats. Any user-defined formats that are found by a procedure are transmitted automatically to the appliance. If you are running multiple high-performance analytical procedures in a SAS session and the analysis variables depend on user-defined formats, you can preprocess the formats. This step involves generating an XML stream (a file) of the formats and passing the stream to the high-performance analytical procedures.

Suppose that the following formats are defined in your SAS program:

```
proc format;
  value YesNo      1='Yes'      0='No';
  value checkThis  1='ThisisOne' 2='ThisisTwo';
  value $cityChar  1='Portage'   2='Kinston';
run;
```

The next group of SAS statements create the XML stream for the formats in the file *Myfmt.xml*, associate that file with the file reference *myxml*, and pass the file reference with the *FMTLIBXML=* option in the PROC HPLOGISTIC statement:

```
filename myxml 'Myfmt.xml';
libname myxml XML92 xmltype=sasfmt tagset=tagsets.XMLsuv;
proc format cntlout=myxml.allfmts;
run;

proc hplogistic data=six fmtlibxml=myxml;
  class wheeze cit age;
  format wheeze best4. cit $cityChar.;
  model wheeze = cit age;
run;
```

Generation and destruction of the stream can be wrapped in convenience macros:

```
%macro Make_XMLStream(name=tempxml);
  filename &name 'fmt.xml';
  libname &name XML92 xmltype=sasfmt tagset=tagsets.XMLsuv;
  proc format cntlout=&name..allfmts;
  run;
%mend;

%macro Delete_XMLStream(fref);
  %let rc=%sysfunc(fdelete(&fref));
%mend;
```

If you do not pass an XML stream to a high-performance analytical procedure that supports the *FMTLIBXML=* option, the procedure generates an XML stream as needed when it is invoked.

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement defines performance parameters for multithreaded and distributed computing, passes variables that describe the distributed computing environment, and requests detailed results about the performance characteristics of a high-performance analytical procedure.

You can also use the PERFORMANCE statement to control whether a high-performance analytical procedure executes in single-machine or distributed mode.

You can specify the following *performance-options* in the PERFORMANCE statement:

### **COMMIT=*n***

requests that the high-performance analytical procedure write periodic updates to the SAS log when observations are sent from the client to the appliance for distributed processing.

High-performance analytical procedures do not have to use input data that are stored on the appliance. You can perform distributed computations regardless of the origin or format of the input data, provided that the data are in a format that can be read by the SAS System (for example, because a SAS/ACCESS engine is available).

In the following example, the HPREG procedure performs LASSO variable selection where the input data set is stored on the client:

```
proc hpreg data=work.one;
  model y = x1-x500;
  selection method=lasso;
  performance nodes=10 host='mydca' commit=10000;
run;
```

In order to perform the work as requested using 10 nodes on the appliance, the data set Work.One needs to be distributed to the appliance.

High-performance analytical procedures send the data in blocks to the appliance. Whenever the number of observations sent exceeds an integer multiple of the COMMIT= size, a SAS log message is produced. The message indicates the actual number of observations distributed, and not an integer multiple of the COMMIT= size.

### **DATASERVER="name"**

specifies the name of the server on Teradata systems as defined through the *hosts* file and as used in the LIBNAME statement for Teradata. For example, assume that the *hosts* file defines the server for Teradata as follows:

```
myservercop1 33.44.55.66
```

Then a LIBNAME specification would be as follows:

```
libname TDLib teradata server=myserver user= password= database= ;
```

A PERFORMANCE statement to induce running alongside the Teradata server would specify the following:

```
performance dataserver="myserver";
```

The DATASERVER= option is not required if you specify the GRIDMODE=option in the PERFORMANCE statement or if you set the GRIDMODE environment variable.

Specifying the DATASERVER= option overrides the GRIDDATASERVER environment variable.

### DETAILS

requests a table that shows a timing breakdown of the procedure steps.

**GRIDHOST=***"name"*

**HOST=***"name"*

specifies the name of the appliance host in single or double quotation marks. If this option is specified, it overrides the value of the GRIDHOST environment variable.

**GRIDMODE=SYM | ASYM**

**MODE=SYM | ASYM**

specifies whether the high-performance analytical procedure runs in symmetric (SYM) mode or asymmetric (ASYM) mode. The default is GRIDMODE=SYM. For more information about these modes, see the section [“Symmetric and Asymmetric Distributed Modes”](#) on page 13.

If this option is specified, it overrides the GRIDMODE environment variable.

**GRIDTIMEOUT=s**

**TIMEOUT=s**

specifies the time-out in seconds for a high-performance analytical procedure to wait for a connection to the appliance and establish a connection back to the client. The default is 120 seconds. If jobs are submitted to the appliance through workload management tools that might suspend access to the appliance for a longer period, you might want to increase the time-out value.

**INSTALL=***"name"*

**INSTALLLOC=***"name"*

specifies the directory in which the shared libraries for the high-performance analytical procedure are installed on the appliance. Specifying the INSTALL= option overrides the GRIDINSTALLLOC environment variable.

**LASRSERVER=***"path"*

**LASR=***"path"*

specifies the fully qualified path to the description file of a SAS LASR Analytic Server instance. If the input data set is held in memory by this LASR Analytic Server instance, then the procedure runs alongside LASR. This option is not needed to run alongside LASR if the DATA= specification of the input data uses a libref that is associated with a LASR Analytic Server instance. For more information, see the section [“Alongside-LASR Distributed Execution”](#) on page 22 and the *SAS LASR Analytic Server: Administration Guide*.



**NODES=ALL** | *n*

**NNODES=ALL** | *n*

specifies the number of nodes in the distributed computing environment, provided that the data are not processed alongside the database.

Specifying NODES=0 indicates that you want to process the data in single-machine mode on the client machine. If the input data are not alongside the database, this is the default. The high-performance analytical procedures then perform the analysis on the client. For example, the following sets of statements are equivalent:

```
proc hplogistic data=one;
  model y = x;
run;
```

```
proc hplogistic data=one;
  model y = x;
  performance nodes=0;
run;
```

If the data are not read alongside the database, the NODES= option specifies the number of nodes on the appliance that are involved in the analysis. For example, the following statements perform the analysis in distributed mode by using 10 units of work on the appliance that is identified in the [HOST=](#) option:

```
proc hplogistic data=one;
  model y = x;
  performance nodes=10 host="hpa.sas.com";
run;
```

If the number of nodes can be modified by the application, you can specify a NODES=*n* option, where *n* exceeds the number of physical nodes on the appliance. The SAS High-Performance Econometrics software then *oversubscribes* the nodes and associates nodes with multiple units of work. For example, on a system that has 16 appliance nodes, the following statements oversubscribe the system by a factor of 3:

```
proc hplogistic data=one;
  model y = x;
  performance nodes=48 host="hpa.sas.com";
run;
```

Usually, it is not advisable to oversubscribe the system because the analytic code is optimized for a certain level of multithreading on the nodes that depends on the CPU count. You can specify `NODES=ALL` if you want to use all available nodes on the appliance without oversubscribing the system.

If the data are read alongside the distributed database on the appliance, specifying a nonzero value for the `NODES=` option has no effect. The number of units of work in the distributed computing environment is then determined by the distribution of the data and cannot be altered. For example, if you are running alongside an appliance with 24 nodes, the `NODES=` option in the following statements is ignored:

```
libname GPLib greenplm server=gpdca user=XXX password=YYY
      database=ZZZ;
proc hplogistic data=gplib.one;
  model y = x;
  performance nodes=10 host="hpa.sas.com";
run;
```

#### **NTHREADS=*n***

#### **THREADS=*n***

specifies the number of threads for analytic computations and overrides the SAS system option `THREADS` | `NOTHREADS`. If you do not specify the `NTHREADS=` option, the number of threads is determined based on the number of CPUs on the host on which the analytic computations execute. The algorithm by which a CPU count is converted to a thread count is specific to the high-performance analytical procedure. Most procedures create one thread per CPU for the analytic computations.

By default, high-performance analytical procedures execute in multiple concurrent threads unless multithreading has been turned off by the `NOTHREADS` system option or you force single-threaded execution by specifying `NTHREADS=1`. The largest number that can be specified for *n* is 256. Individual high-performance analytical procedures can impose more stringent limits if called for by algorithmic considerations.

**NOTE:** The SAS system options `THREADS` | `NOTHREADS` apply to the client machine on which the SAS high-performance analytical procedures execute. They do not apply to the compute nodes in a distributed environment.

## Chapter 4

# The HPCDM Procedure (Experimental)

### Contents

---

Overview: HPCDM Procedure . . . . .	44
Getting Started: HPCDM Procedure . . . . .	46
Estimating a Simple Compound Distribution Model . . . . .	46
Analyzing the Effect of Parameter Uncertainty on the Compound Distribution . . . . .	50
Scenario Analysis . . . . .	52
Syntax: HPCDM Procedure . . . . .	61
Functional Summary . . . . .	62
PROC HPCDM Statement . . . . .	63
BY Statement . . . . .	68
DISTBY Statement . . . . .	68
EXTERNALCOUNTS Statement . . . . .	69
OUTPUT Statement . . . . .	69
OUTSUM Statement . . . . .	70
PERFORMANCE Statement . . . . .	73
SEVERITYMODEL Statement . . . . .	73
Programming Statements . . . . .	74
Details: HPCDM Procedure . . . . .	74
Specifying Scenario Data in the DATA= Data Set . . . . .	74
Simulation Procedure . . . . .	75
Simulation of Adjusted Compound Distribution Sample . . . . .	82
Parameter Perturbation Analysis . . . . .	89
Descriptive Statistics . . . . .	90
Input Specification . . . . .	92
Output Data Sets . . . . .	93
Displayed Output . . . . .	95
ODS Graphics . . . . .	96
Examples: HPCDM Procedure . . . . .	98
Example 4.1: Estimating the Probability Distribution of Insurance Payments . . . . .	98
Example 4.2: Using Externally Simulated Count Data . . . . .	103
References . . . . .	111

---

## Overview: HPCDM Procedure

In many loss modeling applications, the loss events are analyzed by modeling the severity (magnitude) of loss and the frequency (count) of loss separately. The primary goal of preparing these models is to estimate the aggregate loss—that is, the total loss that occurs over a period of time for which the frequency model is applicable. For example, an insurance company might want to assess the expected and worst-case losses for a particular business line, such as automobile insurance, over an entire year given the models for the number of losses in a year and the severity of each loss. A bank might want to assess the value-at-risk (VaR), a measure of the worst-case loss, for a portfolio of assets given the frequency and severity models for each asset type.

Loss severity and loss frequency are random variables, so the aggregate loss is also a random variable. Instead of preparing a point estimate of the expected aggregate loss, it is more desirable to estimate its probability distribution, because this enables you to infer various aspects of the aggregate loss such as measures of location, scale (variability), and shape in addition to percentiles. For example, the value-at-risk that banks or insurance companies use to compute regulatory capital requirements is usually the estimate of the 97.5th or 99th percentile from the aggregate loss distribution.

Let  $N$  represent the frequency random variable for the number of loss events that occur in the time period of interest. Let  $X$  represent the severity random variable for the magnitude of one loss event. Then, the aggregate loss  $S$  is defined as

$$S = \sum_{j=1}^N X_j$$

The goal is to estimate the probability distribution of  $S$ . Let  $F_X(x)$  denote the cumulative distribution function (CDF) of  $X$ ,  $F_X^{*n}(x)$  denote the  $n$ -fold convolution of the CDF of  $X$ , and  $\Pr(N = n)$  denote the probability of seeing  $n$  losses as per the frequency distribution. The CDF of  $S$  is theoretically computable as

$$F_S(s) = \sum_{n=0}^{\infty} \Pr(N = n) \cdot F_X^{*n}(x)$$

This probability distribution model of  $S$ , characterized by the CDF  $F_S(s)$ , is referred to as a *compound distribution model* (CDM). The HPCDM procedure computes an estimate of the CDM, given the distribution models of  $X$  and  $N$ .

PROC HPCDM accepts the severity model of  $X$  as estimated by the SEVERITY procedure. It accepts the frequency model of  $N$  as estimated by the COUNTREG procedure. Both the SEVERITY and COUNTREG procedures are part of SAS/ETS software. Both procedures allow models of  $X$  and  $N$  to be conditional on external factors (regressors). In particular, you can model the severity distribution such that its scale parameter depends on severity regressors, and you can model the frequency distribution such that its mean depends on frequency regressors. The frequency model can also be a zero-inflated model. PROC HPCDM uses the estimates of model parameters and the values of severity and frequency regressors to estimate the compound distribution model.

Direct computation of  $F_S$  is usually a difficult task because of the need to compute the  $n$ -fold convolution. Klugman, Panjer, and Willmot (1998, Ch. 4) suggest some relatively efficient recursion and inversion methods for certain combinations of severity and frequency distributions. However, those methods assume that distributions of  $N$  and  $X$  are fixed and all  $X$ s are identically distributed. When the distributions of  $X$

and  $N$  are conditional on regressors, each set of regressor values results in a different distribution. So you must repeat the recursion and inversion methods for each combination of regressor values, and this repetition makes these methods prohibitively expensive. PROC HPCDM instead estimates the compound distribution by using a Monte Carlo simulation method, which can use all available computational resources to generate a sufficiently large, representative sample of the compound distribution while accommodating the dependence of distributions of  $X$  and  $N$  on external factors. Conceptually, the simulation method works as follows:

1. Use the specified frequency model to draw a value  $N$ , which represents the number of loss events.
2. Use the specified severity model to draw  $N$  values, each of which represents the magnitude of loss for each of the  $N$  loss events.
3. Add the  $N$  severity values from step 2 to compute aggregate loss  $S$  as

$$S = \sum_{j=1}^N X_j$$

This forms one sample point of the CDM.

Steps 1 through 3 are repeated  $M$  number of times, where  $M$  is specified by you, to obtain the representative sample of the CDM. PROC HPCDM analyzes this sample to compute empirical estimates of various summary statistics of the compound distribution such as the mean, variance, skewness, and kurtosis in addition to percentiles such as the median, the 95th percentile, the 99th percentile, and so on. You can also use PROC HPCDM to write the entire simulated sample to an output data set and to produce the plot of the empirical distribution function (EDF), which serves as a nonparametric estimate of  $F_S$ .

The simulation process gets more complicated when the frequency and severity models contain regression effects. The CDM is then conditional on the given values of regressors. The simulation process essentially becomes a scenario analysis, because you need to specify the expected values of the regressors that together represent the scenario for which you want to estimate the CDM. PROC HPCDM enables you to specify an input data set that contains the scenario. If you are modeling a group of entities together (such as a portfolio of multiple assets or a group of insurance policies), each with a different set of characteristics, then the scenario consists of more than one observation, and each observation corresponds to a different entity. PROC HPCDM enables you to specify such a group scenario in the input data set and performs a realistic simulation of loss events that each entity can generate.

PROC HPCDM also enables you to specify externally simulated counts. This is useful if you have an empirical frequency model or if you estimate the frequency model by using a method other than PROC COUNTREG and simulate counts by using such a model. You can specify  $M$  replications of externally simulated counts. For each of the replications, in step 1 of the simulation, instead of using the frequency model, PROC HPCDM uses the count  $N$  that you specify. If the severity model contains regression effects, then you can specify the scenario to simulate for each of the  $M$  replications.

If the parameters of your severity and frequency models have uncertainty associated with them, and they usually do, then you can use PROC HPCDM to conduct parameter perturbation analysis to assess the effect of parameter uncertainty on the estimates of CDM. If you specify that  $P$  perturbed samples be generated, then the parameter set is perturbed  $P$  times, and each time PROC HPCDM makes a random draw from either the univariate normal distribution of each parameter or the multivariate normal distribution over all parameters. For each of the  $P$  perturbed parameter sets, a full compound distribution sample is simulated and summarized. This process yields  $P$  number of estimates for each summary statistic and percentile, which are then used to provide you with estimates of the location and variability of each summary statistic and percentile.

You can also use PROC HPCDM to compute the distribution of an aggregate *adjusted* loss. For example, in insurance applications, you might want to compute the distribution of the *amount paid* in a given time period after applying adjustments such as deductible and policy limit to each individual loss. PROC HPCDM enables you to specify SAS programming statements to adjust each severity value. If  $X_j^a$  represents the adjusted severity value, then PROC HPCDM computes  $S^a$ , an aggregate adjusted loss, as

$$S^a = \sum_{j=1}^N X_j^a$$

All the analyses that PROC HPCDM conducts for the aggregate unadjusted loss, including scenario analysis and parameter perturbation analysis, are also conducted for the aggregate adjusted loss, thereby giving you a comprehensive picture of the adjusted compound distribution model.

---

## Getting Started: HPCDM Procedure

This section outlines the use of the HPCDM procedure to fit compound distribution models. The examples are intended as a gentle introduction to some of the features of the procedure.

---

### Estimating a Simple Compound Distribution Model

This example illustrates the simplest use of PROC HPCDM. Assume that you are an insurance company that has used the historical data about the number of losses per year and the severity of each loss to determine that the Poisson distribution is the best distribution for the loss frequency and that the gamma distribution is the best distribution for the severity of each loss. Now, you want to estimate the distribution of an aggregate loss to determine the worst-case loss that can be incurred by your policyholders in a year. In other words, you want to estimate the compound distribution of  $S = \sum_{i=1}^N X_i$ , where the loss frequency,  $N$ , follows the fitted Poisson distribution and the severity of each loss event,  $X_i$ , follows the fitted gamma distribution.

If your historical count and severity data are stored in the data sets `Work.ClaimCount` and `Work.ClaimSev`, respectively, then you need to ensure that you use the following PROC COUNTREG and PROC SEVERITY steps to fit and store the parameter estimates of the frequency and severity models:

```
/* Fit an intercept-only Poisson count model and
   write estimates to an item store */
proc countreg data=claimcount;
  model numLosses= / dist=poisson;
  store countStorePoisson;
run;

/* Fit severity models and write estimates to a data set */
proc severity data=claimsev criterion=aicc outest=sevest covout plots=none;
  loss lossValue;
  dist _predefined_;
run;
```

The STORE statement in the PROC COUNTREG step saves the count model information, including the parameter estimates, in the Work.CountStorePoisson item store. An item store contains the model information in a binary format that cannot be modified after it is created. You can examine the contents of an item store that is created by a PROC COUNTREG step by specifying a combination of the RESTORE= option and the SHOW statement in another PROC COUNTREG step. For more information, see Chapter 11, “The COUNTREG Procedure” (*SAS/ETS User’s Guide*),.

The OUTEST= option in the PROC SEVERITY statement stores the estimates of all the fitted severity models in the Work.SevEst data set. Let the best severity model that the PROC SEVERITY step chooses be the gamma distribution model.

You can now submit the following PROC HPCDM step to simulate an aggregate loss sample of size 10,000 by specifying the count model’s item store in the COUNTSTORE= option and the severity model’s data set of estimates in the SEVERITYEST= option:

```
/* Simulate and estimate Poisson-gamma compound distribution model */
proc hpcdm countstore=countStorePoisson severityest=sevest
    seed=13579 nreplicates=10000 plots=(edf(alpha=0.05) density)
    print=(summarystatistics percentiles);
    severitymodel gamma;
    output out=aggregateLossSample samplevar=aggloss;
    outsum out=aggregateLossSummary mean stddev skewness kurtosis
        p01 p05 p95 p995=var pctlpts=90 97.5;
run;
```

The SEVERITYMODEL statement requests that an aggregate sample be generated by compounding only the gamma distribution and the frequency distribution. Specifying the SEED= value helps you get an identical sample each time you execute this step, provided that you use the same execution environment. In the single-machine mode of execution, the execution environment is the combination of the operating environment and the number of threads that are used for execution. In the distributed computing mode, the execution environment is the combination of the operating environment, the number of nodes, and the number of threads that are used for execution on each node.

Upon completion, PROC HPCDM creates the two output data sets that you specify in the OUT= options of the OUTPUT and OUTSUM statements. The Work.AggregateLossSample data set contains 10,000 observations such that the value of the AggLoss variable in each observation represents one possible aggregate loss value that you can expect to see in one year. Together, the set of the 10,000 values of the AggLoss variable represents one sample of compound distribution. PROC HPCDM uses this sample to compute the empirical estimates of various summary statistics and percentiles of the compound distribution. The Work.AggregateLossSummary data set contains the estimates of mean, standard deviation, skewness, and kurtosis that you specify in the OUTSUM statement. It also contains the estimates of the 1st, 5th, 90th, 95th, 97.5th, and 99.5th percentiles that you specify in the OUTSUM statement. The value-at-risk (VaR) is an aggregate loss value such that there is a very low probability that an observed aggregate loss value exceeds the VaR. One of the commonly used probability levels to define VaR is 0.005, which makes the 99.5th percentile an empirical estimate of the VaR. Hence, the OUTSUM statement of this example stores the 99.5th percentile in a variable named VaR. VaR is one of the widely used measures of worst-case risk.

Some of the default output and some of the output that you have requested by specifying the PRINT= option are shown in [Figure 4.1](#).

**Figure 4.1** Information, Summary Statistics, and Percentiles of the Poisson-Gamma Compound Distribution

The HPCDM Procedure			
Severity Model: Gamma			
Count Model: Poisson			
Compound Distribution Information			
Severity Model	Gamma Distribution		
Count Model	Poisson Model in Item Store .COUNTSTOREPOISSON		
Sample Summary Statistics			
Mean	4062.8	Median	3349.7
Standard Deviation	3429.6	Interquartile Range	4456.4
Variance	11761948.0	Minimum	0
Skewness	1.14604	Maximum	26077.4
Kurtosis	1.76466	Sample Size	10000
Sample Percentiles			
Percentile		Value	
1		0	
5		0	
25		1449.1	
50		3349.7	
75		5905.5	
90		8792.6	
95		10672.5	
97.5		12391.7	
99		14512.5	
99.5		15877.9	
Percentile Method = 5			

The “Sample Summary Statistics” table indicates that for the given parameter estimates of the Poisson frequency and gamma severity models, you can expect to see a mean aggregate loss of 4,062.8 and a median aggregate loss of 3,349.7 in a year. The “Sample Percentiles” table indicates that there is a 0.5% chance that the aggregate loss exceeds 15,877.9, which is the VaR estimate, and a 2.5% chance that the aggregate loss exceeds 12,391.7. These summary statistic and percentile estimates provide a quantitative picture of the compound distribution. You can also visually analyze the compound distribution by examining the plots that PROC HPCDM prepares. The first plot in [Figure 4.2](#) shows the empirical distribution function (EDF), which is a nonparametric estimate of the cumulative distribution function (CDF). The second plot shows the histogram and the kernel density estimate, which are nonparametric estimates of the probability density function (PDF).



**Figure 4.2** Nonparametric CDF and PDF Plots of the Poisson-Gamma Compound Distribution

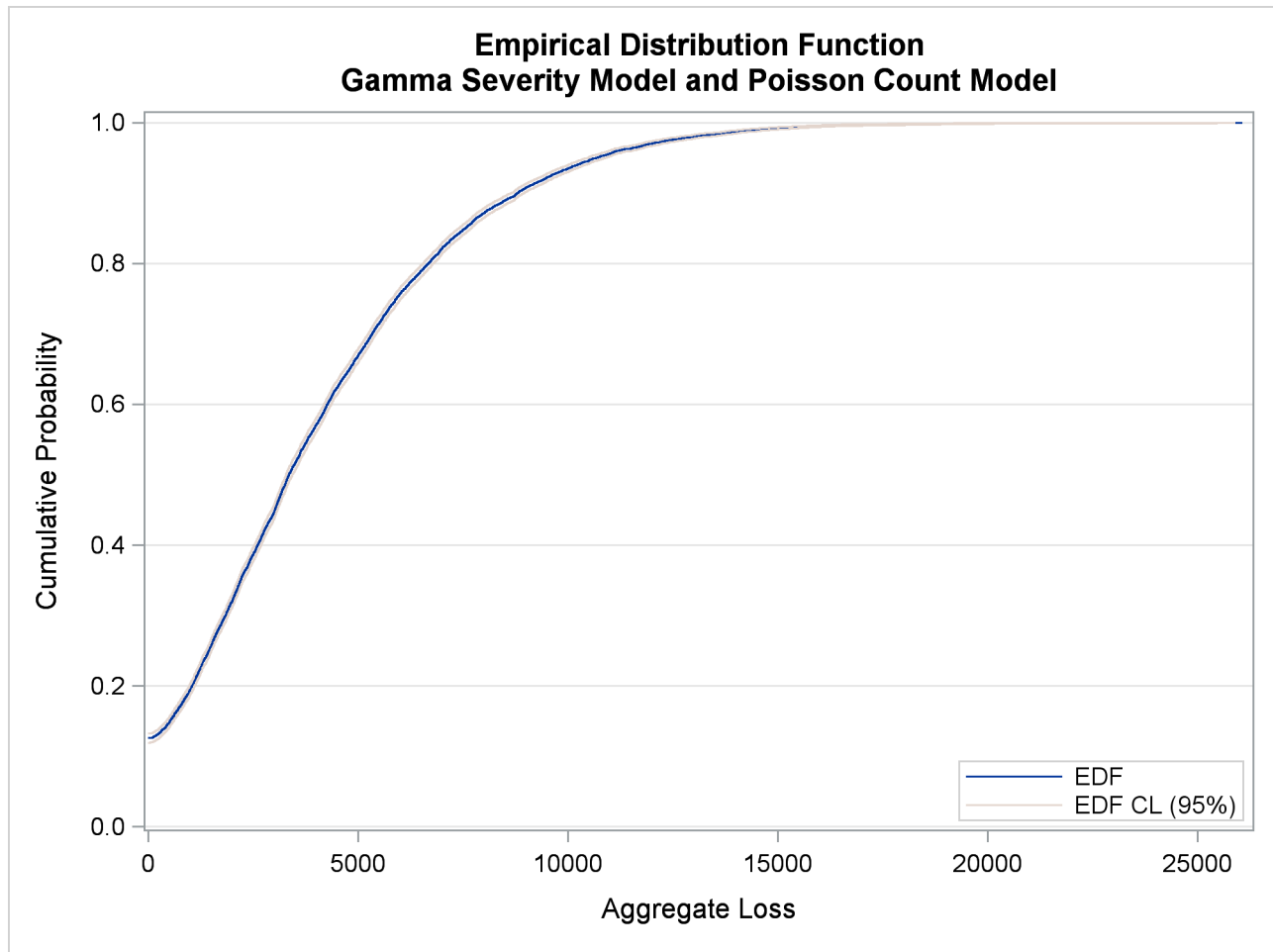
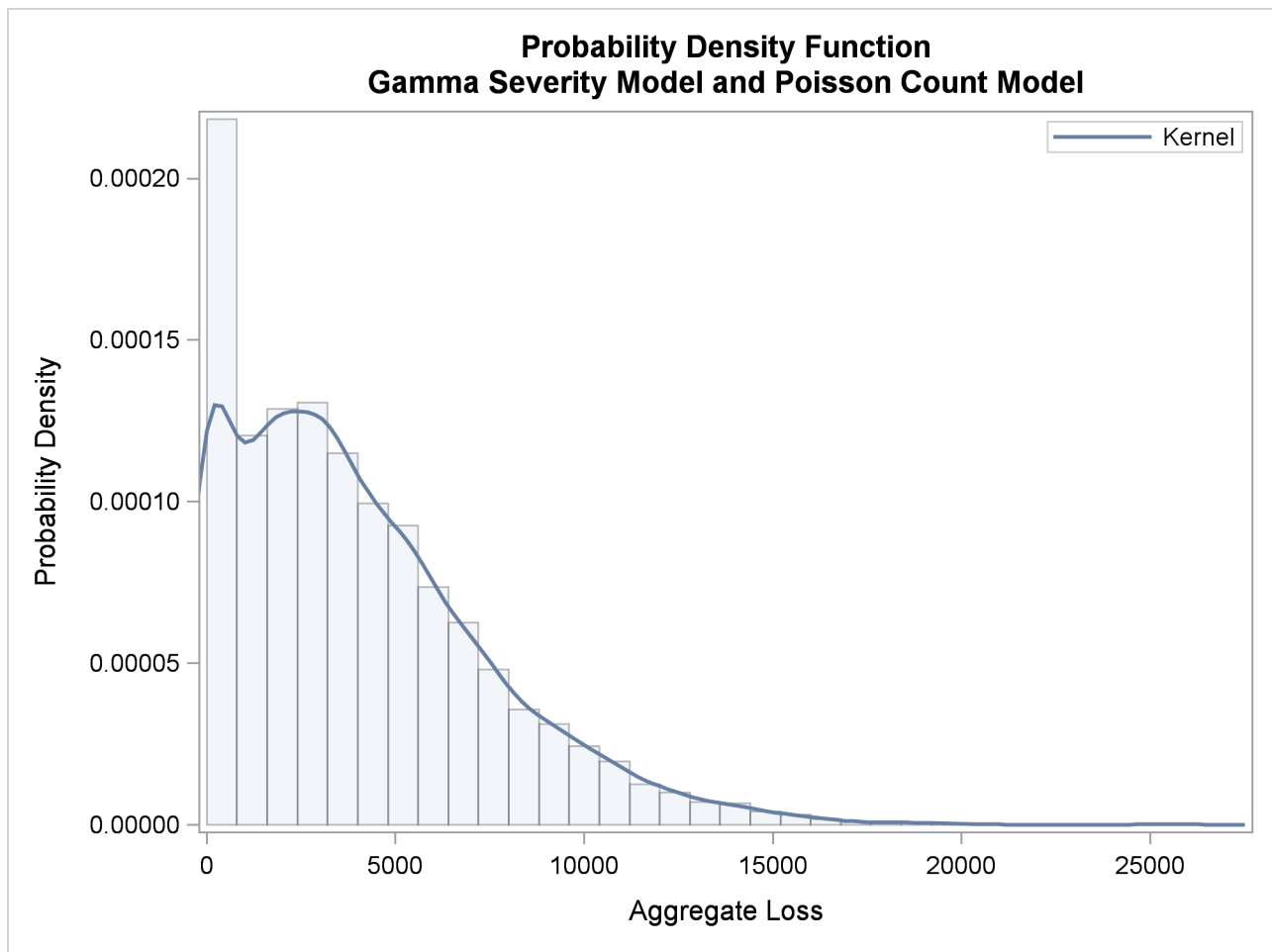


Figure 4.2 continued



The plots confirm the right skew that is indicated by the estimate of skewness in Figure 4.1 and a relatively fat tail, which is indicated by comparing the maximum and the 99.5th percentiles in Figure 4.1.

## Analyzing the Effect of Parameter Uncertainty on the Compound Distribution

Continuing with the previous example, note that you have fitted the frequency and severity models by using the historical data. Even if you choose the best-fitting models, the true underlying models are not known exactly. This fact is reflected in the uncertainty that is associated with the parameters of your models. Any compound distribution estimate that is computed by using these uncertain parameter estimates is inherently uncertain. You can request that PROC HPCDM conduct parameter perturbation analysis, which assesses the effect of the parameter uncertainty on the estimates of the compound distribution by simulating multiple samples, each with parameters that are randomly perturbed from their mean estimates.

The following PROC HPCDM step adds the NPerturbedSamples= option to the PROC HPCDM statement to request that perturbation analysis be conducted and the PRINT=PERTURBSUMMARY option to request that a summary of the perturbation analysis be displayed:

```

/* Conduct parameter perturbation analysis of
the Poisson-gamma compound distribution model */
proc hpcdm countstore=countStorePoisson severityest=sevest
    seed=13579 nreplicates=10000 nperturbedsamples=30
    print(only)=(perturbsummary) plots=none;
severitymodel gamma;
output out=aggregateLossSample samplevar=aggloss;
outsum out=aggregateLossSummary mean stddev skewness kurtosis
    p01 p05 p95 p995=var pctlpts=90 97.5;
run;

```

The Work.AggregateLossSummary data set contains the specified summary statistics and percentiles for all 30 perturbed samples. You can identify a perturbed sample by the value of the `_DRAWID_` variable. The first few observations of the Work.AggregateLossSummary data set are shown in Figure 4.3. For the first observation, the value of the `_DRAWID_` variable is 0, which represents an unperturbed sample—that is, the aggregate sample that is simulated without perturbing the parameters from their means.

**Figure 4.3** Summary Statistics and Percentiles of the Perturbed Samples

<u>_SEVERITYMODEL_</u>	<u>_COUNTMODEL_</u>	<u>_DRAWID_</u>	<u>_SAMPLEVAR_</u>	<u>N</u>	<u>MEAN</u>	<u>STDDEV</u>				
Gamma	Poisson	0	aggloss	10000	4062.76	3429.57				
Gamma	Poisson	1	aggloss	10000	4008.04	3406.22				
Gamma	Poisson	2	aggloss	10000	4426.67	3719.94				
Gamma	Poisson	3	aggloss	10000	3991.87	3480.10				
Gamma	Poisson	4	aggloss	10000	3807.58	3303.61				
Gamma	Poisson	5	aggloss	10000	4083.70	3429.83				
Gamma	Poisson	6	aggloss	10000	4185.82	3525.20				
Gamma	Poisson	7	aggloss	10000	3882.99	3372.81				
Gamma	Poisson	8	aggloss	10000	4092.94	3483.60				
Gamma	Poisson	9	aggloss	10000	4039.82	3454.69				
Gamma	Poisson	10	aggloss	10000	3851.17	3287.52				
<u>SKEWNESS</u>	<u>KURTOSIS</u>	<u>P01</u>	<u>P05</u>	<u>P90</u>	<u>P95</u>	<u>P97_5</u>	<u>var</u>			
1.14604	1.76466	0	0	8792.64	10672.49	12391.70	15877.89			
1.10747	1.43304	0	0	8658.62	10521.82	12279.33	16152.05			
1.14337	1.66525	0	0	9484.05	11522.70	13523.54	17575.20			
1.23233	2.07634	0	0	8672.80	10568.25	12472.90	16969.77			
1.08965	1.15633	0	0	8375.09	10319.59	11884.11	15255.16			
1.08043	1.31018	0	0	8836.78	10707.19	12399.09	16236.24			
1.12642	1.49282	0	0	9095.46	11056.46	12752.18	16519.99			
1.22931	1.95615	0	0	8515.35	10371.84	12245.23	16153.91			
1.10040	1.47077	0	0	8923.13	10757.13	12522.34	16275.95			
1.17185	1.84608	0	0	8696.09	10679.34	12611.43	16350.84			
1.12302	1.60240	0	0	8383.29	10129.41	11725.89	15303.35			

The `PRINT=PERTURBSUMMARY` option in the preceding PROC HPCDM step produces the “Sample Perturbation Analysis” and “Sample Percentile Perturbation Analysis” tables that are shown in Figure 4.4. The tables show that you can expect a mean aggregate loss of about 4,049.1 and the standard error of the mean is 193.6. If you want to use the VaR estimate to determine the amount of reserves that you need to maintain to cover the worst-case loss, then you should consider not only the mean estimate of the 99.5th

percentile, which is about 16,339.1, but also the standard error of 692.8 to account for the effect of uncertainty in your frequency and severity parameter estimates.

**Figure 4.4** Summary of Perturbation Analysis of the Poisson-Gamma Compound Distribution

<p>The HPCDM Procedure Severity Model: Gamma Count Model: Poisson</p>		
Sample Perturbation Analysis		
Statistic	Estimate	Standard Error
Mean	4049.1	193.55480
Standard Deviation	3448.5	132.43375
Variance	11909479	919586.4
Skewness	1.14075	0.04610
Kurtosis	1.64953	0.27146
<p>Number of Perturbed Samples = 30 Size of Each Sample = 10000</p>		
Sample Percentile Perturbation Analysis		
Percentile	Estimate	Standard Error
0	0	0
1	0	0
5	0	0
25	1386.8	114.41389
50	3368.2	185.13314
75	5944.8	265.53061
90	8756.0	365.86765
95	10663.6	441.16381
97.5	12454.8	519.67311
99	14685.6	620.49261
99.5	16339.1	692.79352
<p>Number of Perturbed Samples = 30 Size of Each Sample = 10000</p>		

## Scenario Analysis

The distributions of loss frequency and loss severity often depend on exogenous variables (regressors). For example, the number of losses and the severity of each loss that an automobile insurance policyholder incurs might depend on the characteristics of the policyholder and the characteristics of the vehicle. When you fit frequency and severity models, you need to account for the effects of such regressors on the probability distributions of the counts and severity. The COUNTREG procedure enables you to model regression effects

on the mean of the count distribution, and the SEVERITY procedure enables you to model regression effects on the scale parameter of the severity distribution. When you use these models to estimate the compound distribution model of the aggregate loss, you need to specify a set of values for all the regressors, which represents the state of the world for which the simulation is conducted. This is referred to as the what-if or scenario analysis.

Consider that you, as an automobile insurance company, have postulated that the distribution of the loss event frequency depends on five regressors (external factors): age of the policyholder, gender, type of car, annual miles driven, and policyholder's education level. Further, the distribution of the severity of each loss depends on three regressors: type of car, safety rating of the car, and annual household income of the policyholder (which can be thought of as a proxy for the luxury level of the car). Note that the frequency model regressors and severity model regressors can be different, as illustrated in this example.

Let these regressors be recorded in the variables Age (scaled by a factor of 1/50), Gender (1: female, 2: male), CarType (1: sedan, 2: sport utility vehicle), AnnualMiles (scaled by a factor of 1/5,000), Education (1: high school graduate, 2: college graduate, 3: advanced degree holder), CarSafety (scaled to be between 0 and 1, the safest being 1), and Income (scaled by a factor of 1/100,000), respectively. Let the historical data about the number of losses that various policyholders incur in a year be recorded in the NumLoss variable of the Work.LossCounts data set, and let the severity of each loss be recorded in the LossAmount variable of the Work.Losses data set.

The following PROC COUNTREG step fits the count regression model and stores the fitted model information in the Work.CountregModel item store:

```
/* Fit negative binomial frequency model for the number of losses */
proc countreg data=losscounts;
  model numloss = age gender carType annualMiles education / dist=negbin;
  store work.countregmodel;
run;
```

You can examine the parameter estimates of the count model that are stored in the Work.CountregModel item store by submitting the following statements:

```
/* Examine the parameter estimates for the model in the item store */
proc countreg restore=work.countregmodel;
  show parameters;
run;
```

The “Parameter Estimates” table that is displayed by the SHOW statement is shown in [Figure 4.5](#).

**Figure 4.5** Parameter Estimates of the Count Regression Model

The COUNTREG Procedure					
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Approx Pr >  t
Intercept	1	0.910479	0.090515	10.06	<.0001
age	1	-0.626803	0.058547	-10.71	<.0001
gender	1	1.025034	0.032099	31.93	<.0001
carType	1	0.615165	0.031153	19.75	<.0001
annualMiles	1	-1.010276	0.017512	-57.69	<.0001
education	1	-0.280246	0.021677	-12.93	<.0001
_Alpha	1	0.318403	0.020090	15.85	<.0001

The following PROC SEVERITY step fits the severity scale regression models for all the common distributions that are predefined in PROC SEVERITY:

```
/* Fit severity models for the magnitude of losses */
proc severity data=losses plots=none outest=work.sevregest print=all;
  loss lossamount;
  scalemodel carType carSafety income;
  dist _predef_;
  nloptions maxiter=100;
run;
```

The comparison of fit statistics of various scale regression models is shown in Figure 4.6. The scale regression model that is based on the lognormal distribution is deemed the best-fitting model according to the likelihood-based statistics, whereas the scale regression model that is based on the generalized Pareto distribution (GPD) is deemed the best-fitting model according to the EDF-based statistics.

**Figure 4.6** Severity Model Comparison

The SEVERITY Procedure					
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Burr	127231	127243	127243	127286	7.75407
Exp	128431	128439	128439	128467	6.13537
Gamma	128324	128334	128334	128370	7.54562
Igauss	127434	127444	127444	127480	6.15855
Logn	127062*	127072*	127072*	127107*	6.77687
Pareto	128166	128176	128176	128211	5.37453
Gpd	128166	128176	128176	128211	5.37453*
Weibull	128429	128439	128439	128475	6.21268
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
Burr	224.47578		27.41346		
Exp	181.83094		12.33919		
Gamma	276.13156		24.59515		
Igauss	211.51908		17.70942		
Logn	212.70400		21.47945		
Pareto	110.53673		7.07119		
Gpd	110.53660*		7.07116*		
Weibull	190.81178		13.45425		
Note: The asterisk (*) marks the best model according to each column's criterion.					

Now, you are ready to analyze the distribution of the aggregate loss that can be expected from a specific policyholder—for example, a 59-year-old male policyholder with an advanced degree who earns 159,870 and drives a sedan that has a very high safety rating about 11,474 miles annually. First, you need to encode and scale this information into the appropriate regressor variables of a data set. Let that data set be named `Work.SinglePolicy`, with an observation as shown in [Figure 4.7](#).

**Figure 4.7** Scenario Analysis Data for One Policyholder

age	gender	car Type	annual Miles	education	car Safety	income
1.18	2	1	2.2948	3	0.99532	1.5987

Now, you can submit the following PROC HPCDM step to analyze the compound distribution of the aggregate loss that is incurred by the policyholder in the `Work.SinglePolicy` data set in a given year by using the frequency model from the `Work.CountregModel` item store and the two best severity models, lognormal and GPD, from the `Work.SevRegEst` data set:

```
/* Simulate the aggregate loss distribution for the scenario
   with single policyholder */
proc hpcdm data=singlePolicy nreplicates=10000 seed=13579 print=all
    countstore=work.countregmodel severityest=work.sevregest;
    severitymodel logn gpd;
    outsum out=onepolicysum mean stddev skew kurtosis median
    pctlpts=97.5 to 99.5 by 1;
run;
```

The displayed results from the preceding PROC HPCDM step are shown in [Figure 4.8](#).

When you use a severity scale regression model, it is recommended that you verify the severity scale regressors that are used by PROC HPCDM by examining the Scale Model Regressors row of the “Compound Distribution Information” table. PROC HPCDM detects the severity regressors automatically by examining the variables in the `SEVERITYEST=` and `DATA=` data sets. If those data sets contain variables that you did not include in the `SCALEMODEL` statement in PROC SEVERITY, then such variables can be treated as severity regressors. One common mistake that can lead to this situation is to fit a severity model by using the `BY` statement and forget to specify the identical `BY` statement in the PROC HPCDM step; this can cause PROC HPCDM to treat `BY` variables as scale model regressors. In this example, [Figure 4.8](#) confirms that the correct set of scale model regressors is detected.

**Figure 4.8** Scenario Analysis Results for One Policyholder with Lognormal Severity Model

<p>The HPCDM Procedure  Severity Model: Logn  Count Model: NegBin(p=2)</p>	
<p>Compound Distribution Information</p>	
Severity Model	Lognormal Distribution
Scale Model Regressors	carType carSafety income
Count Model	NegBin(p=2) Model in Item Store WORK.COUNTREGMODEL

Figure 4.8 continued

Sample Summary Statistics			
Mean	209.76287	Median	0
Standard Deviation	559.78686	Interquartile Range	126.67003
Variance	313361.3	Minimum	0
Skewness	5.04884	Maximum	9998.2
Kurtosis	39.24667	Sample Size	10000

Sample Percentiles	
Percentile	Value
0	0
1	0
5	0
25	0
50	0
75	126.67003
95	1215.9
97.5	1863.1
98.5	2288.9
99	2737.3
99.5	3529.0

Percentile Method = 5

The “Sample Summary Statistics” and “Sample Percentiles” tables in Figure 4.8 show estimates of the aggregate loss distribution for the lognormal severity model. The average expected loss is about 210, and the worst-case loss, if approximated by the 97.5th percentile, is about 1,863. The percentiles table shows that the distribution is highly skewed to the right; this is also confirmed by the skewness estimate. The median estimate of 0 can be interpreted in two ways. One way is to conclude that the policyholder will not incur any loss in 50% of the years during which he or she is insured. The other way is to conclude that 50% of policyholders who have the characteristics of this policyholder will not incur any loss in a given year. However, there is a 2.5% chance that the policyholder will incur a loss that exceeds 1,863 in any given year and a 0.5% chance that the policyholder will incur a loss that exceeds 3,529 in any given year.

If the aggregate loss sample is simulated by using the GPD severity model, then the results are as shown in Figure 4.9. The average and worst-case losses are 211 and 1,856, respectively. These estimates are very close to the values that are predicted by the lognormal severity model.



**Figure 4.9** Scenario Analysis Results for One Policyholder with GPD Severity Model

The HPCDM Procedure			
Severity Model: Gpd			
Count Model: NegBin(p=2)			
Compound Distribution Information			
Severity Model	Generalized Pareto Distribution		
Scale Model Regressors	carType carSafety income		
Count Model	NegBin(p=2) Model in Item Store WORK.COUNTREGMODEL		
Sample Summary Statistics			
Mean	211.16729	Median	0
Standard Deviation	539.58331	Interquartile Range	121.18808
Variance	291150.1	Minimum	0
Skewness	4.44116	Maximum	7349.2
Kurtosis	29.03404	Sample Size	10000
Sample Percentiles			
Percentile		Value	
0		0	
1		0	
5		0	
25		0	
50		0	
75		121.18808	
95		1259.5	
97.5		1855.7	
98.5		2288.5	
99		2577.5	
99.5		3294.9	
Percentile Method = 5			

The scenario that you just analyzed contains only one policyholder. You can extend the scenario to include multiple policyholders. Let the Work.GroupOfPolicies data set record information about five different policyholders, as shown in [Figure 4.10](#).

**Figure 4.10** Scenario Analysis Data for Multiple Policyholders

policyholder Id	age	gender	car Type	annual Miles	education	car Safety	income
1	1.18	2	1	2.2948	3	0.99532	1.59870
2	0.66	2	1	2.6718	2	0.86412	0.84459
3	0.64	2	2	1.9528	1	0.86478	0.50177
4	0.46	1	2	2.6402	2	0.27062	1.18870
5	0.62	1	1	1.7294	1	0.32830	0.37694

The following PROC HPCDM step conducts a scenario analysis for the aggregate loss that is incurred by all five policyholders in the Work.GroupOfPolicies data set together in one year:

```
/* Simulate the aggregate loss distribution for the scenario
   with multiple policyholders */
proc hpcdm data=groupOfPolicies nreplicates=10000 seed=13579 print=all
      countstore=work.countregmodel severityest=work.sevregest
      plots=(conditionaldensity(rightq=0.95)) nperturbedSamples=50;
      severitymodel logn gpd;
      outsum out=multipolicysum mean stddev skew kurtosis median
      pctlpts=97.5 to 99.5 by 1;
run;
```

The preceding PROC HPCDM step conducts perturbation analysis by simulating 50 perturbed samples. The perturbation summary results for the lognormal severity model are shown in [Figure 4.11](#), and the results for the GPD severity model are shown in [Figure 4.12](#). If the severity of each loss follows the fitted lognormal distribution, then you can expect that the group of policyholders together incurs an average loss of  $5,333 \pm 547$  and a worst-case loss of  $26,416 \pm 2,681$  when you define the worst-case loss as the 97.5th percentile.

**Figure 4.11** Perturbation Analysis of Losses from Multiple Policyholders with Lognormal Severity Model

The HPCDM Procedure		
Severity Model: Logn		
Count Model: NegBin(p=2)		
Compound Distribution Information		
Severity Model	Lognormal Distribution	
Scale Model Regressors	carType carSafety income	
Count Model	NegBin(p=2) Model in Item Store WORK.COUNTREGMODEL	
Sample Perturbation Analysis		
Statistic	Estimate	Standard Error
Mean	5333.4	547.02353
Standard Deviation	7428.7	729.93126
Variance	55718988	11240603
Skewness	2.99560	0.18583
Kurtosis	14.12580	2.83814
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		

Figure 4.11 *continued*

Sample Percentile Perturbation Analysis		
Percentile	Estimate	Standard Error
1	0	0
5	0	0
25	727.92534	113.89462
50	2589.4	302.39245
75	6919.1	718.06509
95	20059.4	1971.8
97.5	26416.3	2681.2
98.5	31256.8	3136.3
99	35166.1	3403.8
99.5	42119.0	4099.4
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		

If the severity of each loss follows the fitted GPD distribution, then you can expect an average loss of  $5,303 \pm 553$  and a worst-case loss of  $25,885 \pm 2,936$ .

If you decide to use the 99.5th percentile to define the worst-case loss, then the worst-case loss is  $42,119 \pm 4,099$  for the lognormal severity model and  $40,626 \pm 5,022$  for the GPD severity model. The numbers for lognormal and GPD are well within one standard error of each other, which indicates that the aggregate loss distribution is less sensitive to the choice of these two severity distributions in this particular example; you can use the results from either of them.

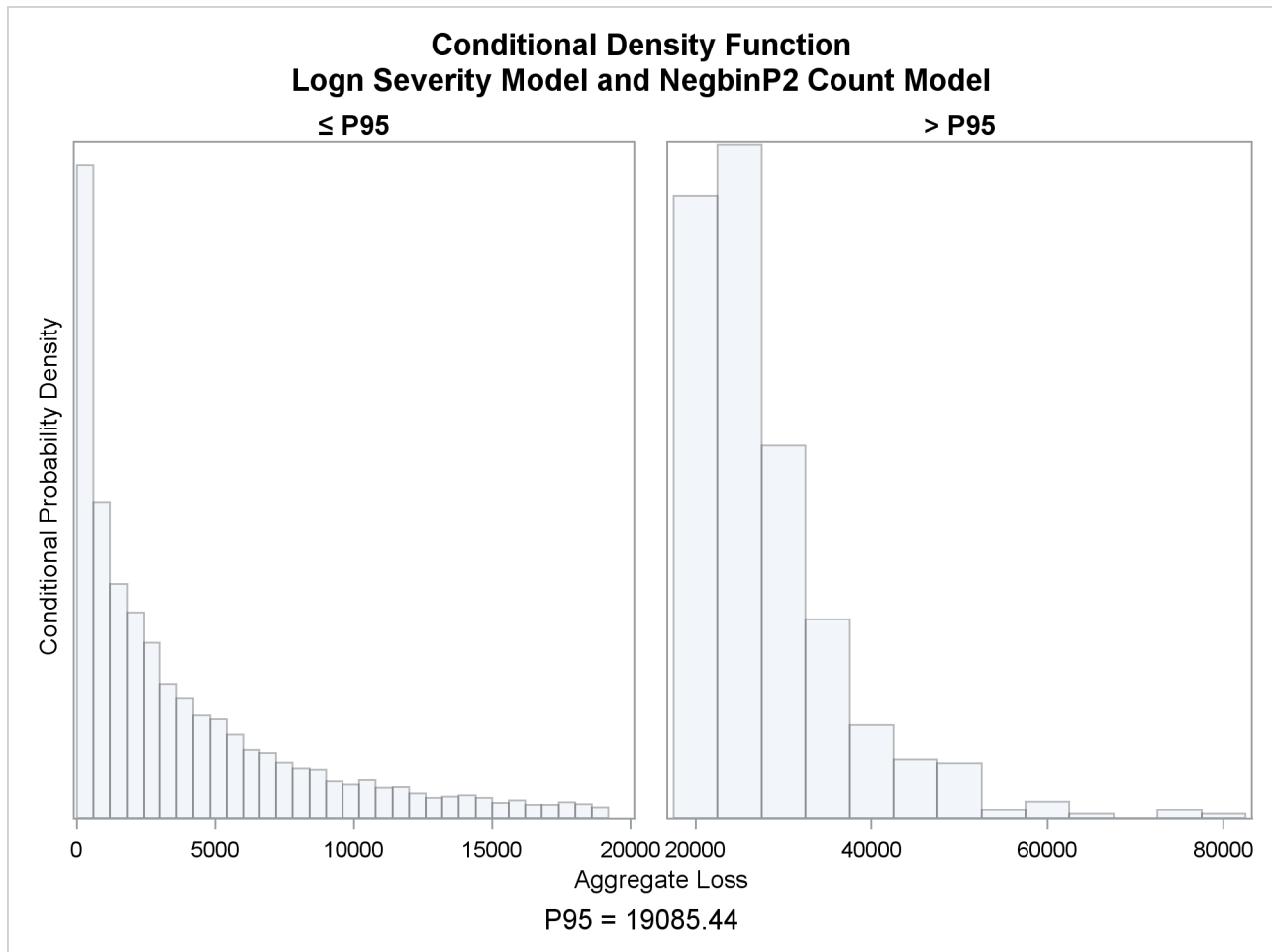
Figure 4.12 Perturbation Analysis of Losses from Multiple Policyholders with GPD Severity Model

The HPCDM Procedure	
Severity Model: Gpd	
Count Model: NegBin(p=2)	
Compound Distribution Information	
Severity Model	Generalized Pareto Distribution
Scale Model Regressors	carType carSafety income
Count Model	NegBin(p=2) Model in Item Store WORK.COUNTREGMODEL

Figure 4.12 continued

Sample Perturbation Analysis		
Statistic	Estimate	Standard Error
Mean	5302.7	552.73848
Standard Deviation	7280.4	836.84441
Variance	53704766	12319537
Skewness	2.90094	0.19750
Kurtosis	13.29785	2.75709
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		
Sample Percentile Perturbation Analysis		
Percentile	Estimate	Standard Error
1	0	0
5	0	0
25	708.74349	102.91404
50	2616.7	282.24080
75	6974.4	697.83278
95	19747.8	2171.9
97.5	25885.1	2936.3
98.5	30424.8	3525.2
99	34213.3	4066.3
99.5	40626.2	5022.4
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		

The PLOTS=CONDITIONALDENSITY option that is used in the preceding PROC HPCDM step prepares the conditional density plots for the body and right-tail regions of the density function of the aggregate loss. The plots for the aggregate loss sample that is generated by using the lognormal severity model are shown in Figure 4.13. The plot on the left side is the plot of  $\Pr(Y|Y \leq 19,085)$ , where the limit 19,085 is the 95th percentile as specified by the RIGHTQ=0.95 option. The plot on the right side is the plot of  $\Pr(Y|Y > 19,085)$ , which helps you visualize the right-tail region of the density function. You can also request the plot of the left tail by specifying the LEFTQ= suboption of the CONDITIONALDENSITY option if you want to explore the details of the left tail region. Note that the conditional density plots are always produced by using the unperturbed sample.

**Figure 4.13** Conditional Density Plots for the Aggregate Loss of Multiple Policyholders

## Syntax: HPCDM Procedure

The following statements are used with the HPCDM procedure:

```

PROC HPCDM options ;
  BY variable-list ;
  DISTBY replication-id-variable ;
  SEVERITYMODEL severity-model-list ;
  EXTERNALCOUNTS COUNT=frequency-variable < ID=replication-id-variable > ;
  OUTPUT OUT=SAS-data-set < variable-name-options > < / out-option > ;
  OUTSUM OUT=SAS-data-set statistic-keyword< =variable-name > < ... statistic-keyword< =variable-
    name>> < outsum-options > ;
  PERFORMANCE options ;
  Programming statements ;

```

## Functional Summary

Table 4.1 summarizes the statements and options available in the HPCDM procedure.

**Table 4.1** HPCDM Functional Summary

Description	Statement	Option
<b>Statements</b>		
Specifies the names of severity distribution models	SEVERITYMODEL	
Specifies externally simulated count data	EXTERNALCOUNTS	
Specifies where and how the full simulated samples are written	OUTPUT	
Specifies where and how the summary statistics of simulated samples are written	OUTSUM	
Specifies performance options	PERFORMANCE	
Specifies programming statements that define an objective function	Programming statements	
<b>Data Set Options</b>		
Specifies the input data set	PROC HPCDM	DATA=
Specifies the output data set for the full simulated samples	OUTPUT	OUT=
Specifies the output data set for the summary statistics	OUTSUM	OUT=
<b>Model Input Options</b>		
Specifies the variable that contains externally simulated counts	EXTERNALCOUNTS	COUNT=
Specifies the item store that contains the frequency (count) model	PROC HPCDM	COUNTSTORE=
Specifies the replicate identifier variable for external counts	EXTERNALCOUNTS	ID=
Specifies the input data set for parameter estimates of the severity models	PROC HPCDM	SEVERITYEST=
<b>Simulation Options</b>		
Specifies the adjusted severity symbol in the programming statements	PROC HPCDM	ADJUSTEDSEVERITY=
Specifies the number of parameter-perturbed samples to be simulated	PROC HPCDM	NPERTURBEDSAMPLES=
Specifies a number that controls the size of the simulated sample	PROC HPCDM	NREPLICATES=
Specifies a seed for the internal pseudo-random number generator	PROC HPCDM	SEED=

Table 4.1 *continued*

Description	Statement	Option
<b>Output Preparation Options</b>		
Specifies the variable for the aggregate adjusted loss sample	OUTPUT	ADJSAMPLEVAR=
Specifies the names of the variables for percentiles	OUTSUM	PCTLNAME=
Specifies the decimal precision to form default percentile variable names	OUTSUM	PCTLNDEC=
Specifies percentiles to compute and report	OUTSUM	PCTLPTS=
Specifies the method to compute the percentiles	PROC HPCDM	PCTLDEF=
Specifies that all perturbed samples be written to the OUT= data set	OUTPUT	PERTURBOUT
Specifies the variable for the aggregate loss sample	OUTPUT	SAMPLEVAR=
Specifies the denominator for computing second- and higher-order moments	PROC HPCDM	VARDEF=
<b>Displayed Output and Plotting Options</b>		
Suppresses all displayed and graphical output	PROC HPCDM	NOPRINT
Specifies which displayed output to prepare	PROC HPCDM	PRINT=
Specifies which graphical output to prepare	PROC HPCDM	PLOTS=

## PROC HPCDM Statement

### PROC HPCDM *options* ;

The PROC HPCDM statement invokes the procedure. You can specify the following *options*, which are listed in alphabetical order.

**ADJUSTEDSEVERITY=***symbol-name*

**ADJSEV=***symbol-name*

names the symbol that represents the adjusted severity value in the SAS programming statements that you specify. The *symbol-name* is a SAS name that conforms to the naming conventions of a SAS variable. For more information, see the section “[Programming Statements](#)” on page 74.

**COUNTSTORE=***SAS-item-store*

names the item store that contains all the information about the frequency (count) model. The COUNTREG procedure generates this item store when you use the STORE statement.

The exogenous variables in the frequency model, if any, are deduced from this item store. The DATA= data set must contain all those variables.

You must specify this option if you do not specify the EXTERNALCOUNTS statement. This option is ignored if you specify the EXTERNALCOUNTS statement, because PROC HPCDM does not need to simulate frequency counts internally when you specify externally simulated counts.

In SAS/ETS 13.1, PROC COUNTREG can create an item store only when you fit one count model and when you do not specify the BY statement in PROC COUNTREG. So if you specify the COUNTSTORE= option, then you cannot specify the BY statement in PROC HPCDM, and vice versa.

**DATA=SAS-data-set**

names the input data set that contains the values of regression variables in frequency or severity models and severity adjustment variables that you use in the programming statements.

The DATA= data set specifies information about the scenario for which you want to estimate the aggregate loss distribution. The contents of the data set are interpreted differently based on whether you specify the EXTERNALCOUNTS statement. For more information, see the section “[Specifying Scenario Data in the DATA= Data Set](#)” on page 74.

**NOPRINT**

turns off all displayed and graphical output. If you specify this option, then PROC HPCDM ignores any value that you specify for the PRINT= or PLOTS= option.

**NPERTURBEDSAMPLES=number**

**NPERTURB=number**

requests that parameter perturbation analysis be conducted. The model parameters are perturbed the specified *number* of times and a separate full sample is simulated for each set of perturbed parameter values. The summary statistics and percentiles are computed for each such perturbed sample, and their values are aggregated across the samples to compute the mean and standard deviation of each summary statistic and percentile.

The parameter perturbation procedure makes random draws of parameter values from a multivariate normal distribution if the covariance estimates of the parameters are available in the SEVERITYEST= data set for the severity model and in the COUNTSTORE= store for the count model. If covariance estimates are not available, then for each parameter, a random draw is made from the univariate normal distribution that has mean and standard deviation equal to the point estimate and the standard error, respectively, of that parameter. If neither covariance nor standard error estimates are available, then perturbation analysis is not conducted.

If you specify the PRINT=ALL or PRINT=PERTURBSUMMARY option, then the summary of perturbation analysis is printed for the core summary statistics and the percentiles of the aggregate loss distribution. If you specify the OUTSUM statement, then the requested summary statistics are written to the OUTSUM= data set for each perturbed sample. You can also optionally request that each perturbed sample be written in its entirety to the OUT= data set by specifying the PERTURBOUT option in the OUTPUT statement.

For more information on the parameter perturbation analysis, see the section “[Parameter Perturbation Analysis](#)” on page 89.

**NREPLICATES=number**

**NREP=number**

specifies a *number* that controls the size of the compound distribution sample that PROC HPCDM simulates. The *number* is interpreted differently based on whether you specify the EXTERNALCOUNTS statement.

If you do not specify the EXTERNALCOUNTS statement, then the sample size is equal to the *number* that you specify for this option. If you do not specify this option, then a default value of 100,000 is used.



If you specify the EXTERNALCOUNTS statement, then the number of replicates that you specify in the DATA= data set is multiplied by the *number* that you specify for this option to get the total size of the compound distribution sample. If you do not specify this option, then a default value of 1 is used.

**PCTLDEF=percentile-method**

specifies the method to compute the percentiles of the compound distribution. The *percentile-method* can be 1, 2, 3, 4, or 5. The default method is 5. For more information, see the description of the PCTLDEF= option in the UNIVARIATE procedure in the *Base SAS Procedures Guide: Statistical Procedures*.

**PLOTS <(global-plot-options)> =plot-request-option**

**PLOTS <(global-plot-options)> =(plot-request-option . . . plot-request-option)**

specifies the desired graphical output.

By default, the HPCDM procedure produces no graphical output.

You can specify the following *global-plot-option*:

**ONLY**

turns off the default graphical output and prepares only the requested plots.

If you specify more than one *plot-request-option*, then separate them with spaces and enclose them in parentheses. The following *plot-request-options* are available:

**ALL**

displays all the graphical output.

**CONDITIONALDENSITY (conditional-density-plot-options)**

**CONDPDF (conditional-density-plot-options)**

prepares a group of plots of the conditional density functions estimates. The group contains at most three plots, each conditional on the value of the aggregate loss being in one of the three regions that are defined by the quantiles that you specify in the following *conditional-density-plot-options*:

**LEFTQ=number**

specifies the quantile in the range (0,1) that marks the end of the left-tail region. If you specify a value of *l* for *number*, then the left-tail region is defined as the set of values that are less than or equal to  $q_l$ , where  $q_l$  is the *l*th quantile. For the left-tail region, nonparametric estimates of the conditional probability density function  $f_S^l(s) = \Pr[S = s | S \leq q_l]$  are plotted. The value of  $q_l$  is estimated by the 100/*l*th percentile of the simulated compound distribution sample.

If you do not specify this option or you specify a missing value for this option, then the left-tail region is not plotted.

**RIGHTQ=number**

specifies the quantile in the range (0,1) that marks the beginning of the right-tail region. If you specify a value of *r* for *number*, then the right-tail region is defined as the set of values that are greater than  $q_r$ , where  $q_r$  is the *r*th quantile. For the right-tail region, nonparametric estimates of the conditional probability density function  $f_S^r(s) = \Pr[S = s | S > q_r]$  are plotted. The value of  $q_r$  is estimated by the 100/*r*th percentile of the simulated compound distribution sample.

If you do not specify this option or you specify a missing value for this option, then the right-tail region is not plotted.

You must specify nonmissing value for at least one of the preceding two options. For the region between the LEFTQ= and RIGHTQ= quantiles, which is referred to as the central or body region, nonparametric estimates of the conditional probability density function  $f_S^c(s) = \Pr[S = s | q_l < S \leq q_r]$  are plotted. If you do not specify a LEFTQ= value, then  $q_l$  is assumed to be 0. If you do not specify a RIGHTQ= value, then  $q_r$  is assumed to be  $\infty$ .

### DENSITY

prepares a plot of the nonparametric estimates of the probability density function (in particular, histogram and kernel density estimates) of the compound distribution.

### EDF <(edf-plot-option)>

prepares a plot of the nonparametric estimates of the cumulative distribution function of the compound distribution.

You can request that the confidence interval be plotted by specifying the following *edf-plot-option*:

#### ALPHA=*number*

specifies the confidence level in the (0,1) range that is used for computing the confidence intervals for the EDF estimates. If you specify a value of  $\alpha$  for *number*, then the upper and lower confidence limits for the confidence level of  $100(1 - \alpha)$  are plotted.

### NONE

displays none of the graphical output. If you specify this option, then it overrides all other plot request options. The default graphical output is also suppressed.

Note that if the simulated sample size is large, then it can take a significant amount of time and memory to prepare the plots.

### PRINT <(global-display-option)> =*display-option*

### PRINT <(global-display-option)> =(*display-option* . . . *display-option*)

specifies the desired displayed output. If you specify more than one *display-option*, then separate them with spaces and enclose them in parentheses.

You can specify the following *global-display-option*:

### ONLY

turns off the default displayed output and displays only the requested output.

You can specify the following *display-options*:

### ALL

displays all the output.

### NONE

displays none of the output. If you specify this option, then it overrides all other display options. The default displayed output is also suppressed.

**PERCENTILES**

displays the percentiles of the compound distribution sample. This includes all the predefined percentiles, percentiles that you request in the OUTSUM statement, and percentiles that you specify for preparing conditional density plots.

**PERTURBSUMMARY**

displays the mean and standard deviation of the summary statistics and percentiles that are taken across all the samples produced by perturbing the model parameters. This option is valid only if you specify the NPerturbedSamples= option in the PROC HPCDM statement.

**SUMMARYSTATISTICS | SUMSTAT**

displays the summary statistics of the compound distribution sample.

If you do not specify the PRINT= option or the ONLY *global-display-option*, then the default displayed output is equivalent to specifying PRINT=(SUMMARYSTATISTICS).

**SEED=number**

specifies the integer to use as the seed in generating the pseudo-random numbers that are used for simulating severity and frequency values. If you do not specify the seed or if *number* is negative or 0, then the time of day from the computer's clock is used as the seed.

**SEVERITYEST=SAS-data-set**

names the input data set that contains the parameter estimates for the severity model. The format of this data set must be the same as the OUTEST= data set that is produced by the SEVERITY procedure.

The names of the regression variables in the scale regression model, if any, are deduced from this data set. In particular, PROC HPCDM assumes that all the variables in the SEVERITYEST= data set that do not appear in the following list are scale regression variables:

- BY variables
- \_MODEL\_, \_TYPE\_, \_NAME\_, and \_STATUS\_ variables
- variables that represent distribution parameters

The DATA= data set must contain all the regressors in the scale regression model.

To ensure that PROC HPCDM correctly matches the values of regressors and the values of regression parameter estimates, you might need to rename the regressors in the DATA= data set so that their names match the names of the regressors that you specify in the SCALEMODEL statement of the PROC SEVERITY step that fits the severity model.

**VARDEF=divisor**

specifies the divisor to use in the calculation of variance, standard deviation, kurtosis, and skewness of the compound distribution sample. If the sample size is  $N$ , then you can specify one of the following values for the *divisor*:

**DF**

sets the divisor for variance to  $N - 1$ . This is the default. This also changes the definitions of skewness and kurtosis.

**N**sets the divisor to  $N$ .

For more information, see the section “Descriptive Statistics” on page 90.

---

## BY Statement

**BY** *variable-list* ;

You can use the BY statement in the HPCDM procedure to process the input data set in groups of observations defined by the BY variables.

If you specify the BY statement, then PROC HPCDM expects the input data set to be sorted in the order of the BY variables unless you specify the NOTSORTED option.

The BY statement is always supported in the single-machine mode of execution. For the distributed mode, it is supported only when the DATA= data set resides on the client machine. In other words, the BY statement is supported only in the client-data (or local-data) mode of the distributed computing model and not for any of the alongside modes, such as the alongside-the-database or alongside HDFS mode.

If you specify the COUNTSTORE= option, then the BY group processing is not supported.

---

## DISTBY Statement

**DISTBY** *replication-id-variable* ;

A DISTBY statement is necessary if and only if you specify an ID= variable in the EXTERNALCOUNTS statement. In fact, the *replication-id-variable* must be the same as the ID= variable. This is especially important in the distributed mode of execution, because when the observations in the DATA= data set are distributed to the grid nodes, by specifying the *replication-id-variable* as a DISTBY variable, you are instructing PROC HPCDM to make sure that the observations that have the same value for the *replication-id-variable* are always kept together on one grid node. This is required for correct simulation of the CDM in the presence of the ID= variable.

Contrast this to the BY variables that you specify in the BY statement. The observations of a BY group might be split across all the nodes of the grid, but the observations of a DISTBY group, which is nested within a BY group, are never split across the nodes of the grid.

The *replication-id-variable* must not appear in the BY statement. However, the DATA= data set must be sorted as if the *replication-id-variable* were listed after the BY variables in the BY statement.

Even though the DISTBY statement is important primarily in distributed mode, you must also specify it in single-machine mode.

---

## EXTERNALCOUNTS Statement

**EXTERNALCOUNTS** *COUNT=frequency-variable* < *ID=replication-id-variable* > ;

The EXTERNALCOUNTS statement enables you to specify externally simulated frequency counts. By default, PROC HPCDM internally simulates the number of loss events by using the frequency model input (COUNTSTORE= item store). However, if you specify the EXTERNALCOUNTS statement, then PROC HPCDM uses the counts that you specify in the DATA= data set and simulates only the severity values internally.

If you specify more than one EXTERNALCOUNTS statement, only the first one is used.

You must specify the following option in the EXTERNALCOUNTS statement:

**COUNT**=*count-variable*

specifies the variable that contains the simulated counts. This variable must be present in the DATA= data set.

You can also specify the following option in the EXTERNALCOUNTS statement:

**ID**=*replication-id-variable*

specifies the variable that contains the replicate identifier. This variable must be present in the DATA= data set. Furthermore, you must specify the DISTBY statement with *replication-id-variable* as the only DISTBY variable to ensure correct simulation.

The observations of DATA= data set must be arranged such that the values of the ID= variable are in increasing order in each BY group or in the entire data set if you do not specify the BY statement.

If you do not specify the ID= option, then PROC HPCDM assumes that each observation represents one replication. In other words, the observation number serves as the default replication identifier.

The simulation process of using the external counts to generate the compound distribution sample is described in the section “[Simulation with External Counts](#)” on page 78.

---

## OUTPUT Statement

**OUTPUT** *OUT=SAS-data-set* < *variable-name-options* > < / *out-option* > ;

The OUTPUT statement enables you to specify the data set to output the generated compound distribution sample.

If you specify more than one OUTPUT statement, only the first one is used.

You must specify the output data set by using the following option:

**OUT**=*SAS-data-set*

**OUTSAMPLE**=*SAS-data-set*

specifies the output data set to contain the simulated compound distribution sample. If you specify programming statements to adjust individual severity values, then this data set contains both unadjusted and adjusted samples.

You can specify the following *variable-name-options* to control the names of the variables created in the OUT= data set:

**ADJSAMPLEVAR=***variable-name*

specifies the name of the variable to contain the adjusted compound distribution sample in the OUT= data set. If you do not specify ADJSAMPLEVAR= option, then “\_AGGADJSEV\_” is used by default.

This option is ignored if you do not specify the **ADJUSTEDSEVERITY=** option and the programming statements to adjust the simulated severity values.

**SAMPLEVAR=***variable-name*

specifies the name of the variable to contain the simulated sample in the OUT= data set. If you do not specify SAMPLEVAR= option, then “\_AGGSEV\_” is used by default.

Further, you can request that the perturbed samples be written to the OUT= data set by specifying the following *out-option*:

**PERTURBOUT**

specifies that all the perturbed samples be written to the OUT= data set. Each perturbed sample is identified by the \_DRAWID\_ variable in the OUT= data set. A value of 0 for the \_DRAWID\_ variable indicates an unperturbed sample.

Separate compound distribution samples are generated for each combination of specified severity and frequency models. The \_SEVERITYMODEL\_ and \_COUNTMODEL\_ columns in the OUT= data set identify the severity and frequency models, respectively, that are used to generate the sample in the SAMPLEVAR= and ADJSAMPLEVAR= variables.

---

## OUTSUM Statement

**OUTSUM** *OUT=SAS-data-set statistic-keyword* <=*variable-name*> <... *statistic-keyword* <=*variable-name*>> <*outsum-options*> ;

The OUTSUM statement enables you to specify the data set in which PROC HPCDM writes the summary statistics of the compound distribution samples.

If you specify more than one OUTSUM statement, only the first one is used.

You must specify the output data set by using the following option:

**OUT=***SAS-data-set*

**OUTSUM=***SAS-data-set*

specifies the output data set that contains the summary statistics of each of the simulated compound distribution samples. You can control the summary statistics that appear in this data set by specifying different *statistic-keywords* and *outsum-options*.

You can request that one or more predefined statistics of the compound distribution sample be written to the OUTSUM= data set. For each specification of the form *statistic-keyword* <=*variable-name*>, the statistic that is specified by the *statistic-keyword* is written to a variable named *variable-name*. If you do not specify the *variable-name*, then the statistic is written to a variable named *statistic-keyword*. You can specify the following *statistic-keywords*:

**KURTOSIS****KURT**

specifies the kurtosis of the compound distribution sample.

**MEAN**

specifies the mean of the compound distribution sample.

**MEDIAN****Q2****P50**

specifies the median (the 50th percentile) of the compound distribution sample.

**P01**

specifies the 1st percentile of the compound distribution sample.

**P05**

specifies the 5th percentile of the compound distribution sample.

**P95**

specifies the 95th percentile of the compound distribution sample.

**P99**

specifies the 99th percentile of the compound distribution sample.

**P99\_5****P995**

specifies the 99.5th percentile of the compound distribution sample.

**Q1****P25**

specifies the lower or 1st quartile (the 25th percentile) of the compound distribution sample.

**Q3****P75**

specifies the upper or 3rd quartile (the 75th percentile) of the compound distribution sample.

**QRANGE**

specifies the interquartile range ( $Q3 - Q1$ ) of the compound distribution sample.

**SKEWNESS****SKEW**

specifies the skewness of the compound distribution sample.

**STDDEV****STD**

specifies the standard deviation of the compound distribution sample.

All percentiles are computed by using the method that you specify for the **PCTLDEF=** option in the PROC HPCDM statement. You can also request additional percentiles to be reported in the OUTSUM= data set by specifying the following *outsum-options*:

**PCTLPTS=percentile-list**

specifies one or more percentiles that you want to be computed and written to the OUTSUM= data set. This option is useful if you need to request percentiles that are not available in the preceding list of *statistic-keyword* values. Each percentile value must belong to the (0,100) open interval. The *percentile-list* is a comma-separated list of numbers. You can also use a list notation of the form “<number1> to <number2> by <increment>”. For example, the following two options are equivalent:

```
pctlpts=10, 20, 99.6, 99.7, 99.8, 99.9
pctlpts=10, 20, 99.6 to 99.9 by 0.1
```

The name of the variable for a given percentile value is decided by the PCTLNAME= option.

**PCTLNAME=percentile-variable-name-list**

specifies the names of the variables that contain the estimates of the percentiles that you request by using the PCTLPTS= option.

If you do not specify the PCTLNAME= option, then each percentile value  $t$  in the list of values in the PCTLPTS= option is written to the variable named “Pt,” where the decimal point in  $t$ , if any, is replaced by an underscore.

The *percentile-variable-name-list* is a space-separated list of names. You can also use a shortcut notation of <prefix> $m$ –<prefix> $n$  for two integers  $m$  and  $n$  ( $m < n$ ) to generate the following list of names: <prefix> $m$ , <prefix> $m + 1$ , ..., and <prefix> $n$ . For example, the following two options are equivalent:

```
pctlname=p1 p2 pc5 pc6 pc7 pc8 pc9 pc10
pctlname=p1 p2 pc5-pc10
```

The name in  $j$ th position of the expanded name list of the PCTLNAME= option is used to create a variable for a percentile value in the  $j$ th position of the expanded value list of the PCTLPTS= option. If you specify  $k_n$  names in the PCTLNAME= option and  $k_v$  percentile values in the PCTLPTS= option, and if  $k_n < k_v$ , then the first  $k_n$  percentiles are written to the variables that you specify and the remaining  $k_v - k_n$  percentiles are written to the variables that have the name of the form  $Pt$ , where  $t$  is the text representation of the percentile value that is formed by retaining at most PCTLNDEC= digits after the decimal point and replacing the decimal point with an underscore (‘\_’). For example, assume you specify the options

```
pctlpts=10, 20, 99.3 to 99.5 by 0.1, 99.995
pctlname=pten ptwenty ninenine3-ninenine5
```

Then PROC HPCDM writes the 10th and 20th percentiles to pten and ptwenty variables, respectively; the 99.3rd through 99.5th percentiles to ninenine3, ninenine4, and ninenine5 variables, respectively; and the remaining 99.995th percentile to the P99\_995 variable.



If a percentile value in the PCTLPTS= option matches a percentile value implied by one of the predefined percentile statistics and you specify the corresponding *statistic-keyword*, then the variable name that is implied by the *statistic-keyword*<=*variable-name*> specification takes precedence over the name that you specify in the PCTLNAME= option. For example, assume you specify the predefined percentile statistic of P95 as in the OUTSUM statement

```
outsum out=mypctls p95=ninetyfifth
      pctlpts=95 to 99 by 1 pctlname=pct95-pct99;
```

Then the 95th percentile is written to the ninetyfifth variable instead of the pct95 variable that the PCTLNAME= option implies.

#### **PCTLNDEC=integer-value**

specifies the maximum number of decimal places to use while creating the names of the variables for the percentile values in the PCTLPTS= option. The default value is 3. For example, for a percentile value of 99.9995, PROC HPCDM creates a variable named P99\_999 by default, but if you specify PCTLNDEC=4, then the variable is named P99\_9995.

The PCTLNDEC= option is used only for percentile values for which you do not specify a name in the PCTLNAME= option.

Note that all variable names in the OUTSUM= data set have a limit of 32 characters. If a name exceeds that limit, then it is truncated to contain only the first 32 characters. For more information about the variables in the OUTSUM= data set, see the section “[Output Data Sets](#)” on page 93.

---

## **PERFORMANCE Statement**

### **PERFORMANCE** *options* ;

The PERFORMANCE statement defines performance parameters for distributed and multithreaded computing, passes variables that describe the distributed computing environment, and requests detailed results about the performance characteristics of PROC HPCDM.

You can also use the PERFORMANCE statement to control whether a high-performance analytical procedure executes in single-machine or distributed mode.

For more information about the PERFORMANCE statement, see the section “[PERFORMANCE Statement](#)” on page 39 of Chapter 3, “[Shared Concepts and Topics](#).”

---

## **SEVERITYMODEL Statement**

### **SEVERITYMODEL** *severity-model-list* ;

The SEVERITYMODEL statement specifies one or more severity distribution models that you want to use in simulating a compound distribution sample. The *severity-model-list* is a space-separated list of names of severity models that you would use with PROC SEVERITY’s DIST statement. The [SEVERITYEST=](#) data set must contain all the severity models in the list. If you specify a name that does not appear in the `_MODEL_` column of the SEVERITYEST= data set, then that name is ignored.

If you specify more than one SEVERITYMODEL statement, only the first one is used.

If you do not specify a SEVERITYMODEL statement, then this is equivalent to specifying all the severity models that appear in the SEVERITYEST= data set.

A compound distribution sample is generated for each of the severity models by compounding that severity model with the frequency model that you specify in the COUNTSTORE= item store or the external frequency model that is encoded by the COUNT= variable that you specify in the EXTERNALCOUNTS statement.

---

## Programming Statements

In PROC HPCDM, you can use a series of programming statements that use variables in the DATA= data set to adjust an individual severity value. The adjusted severity values are aggregated to form a separate adjusted compound distribution sample.

The programming statements are executed for each simulated individual severity value. The observation of the input data set that is used to evaluate the programming statements is determined by the simulation procedure that is described in the section “Simulation Procedure” on page 75.

For more information, see the section “Simulation of Adjusted Compound Distribution Sample” on page 82.

---

## Details: HPCDM Procedure

---

### Specifying Scenario Data in the DATA= Data Set

A scenario represents a state of the world for which you want to estimate the distribution of aggregate losses. The state consists of one or more entities that generate the loss events. For example, an entity might be an individual who has an insurance policy or an organization that has a workers’ compensation policy. Each entity has some characteristics of its own and some external factors that affect the frequency with which it generates the losses and the severity of each loss. For example, characteristics of an individual with an automobile insurance policy can include various demographics of the individual and various features of the automobile. Characteristics of an organization with a workers’ compensation policy can be the number of employees, revenue, ratio of temporary to permanent employees, and so on. The organization can also be affected by external macroeconomic factors such as GDP and unemployment of the country where the organization operates and factors that affect its industry. You need to quantify and specify all these characteristics as external factors (regressors) when you fit severity and frequency models.

You should specify all the information about a scenario in the DATA= data set that you specify in the PROC HPCDM statement. Each observation in the DATA= data set encodes the characteristics of an entity. For proper simulation of counts and severities, you must specify in the DATA= data set all the characteristics that you use as regressors in the frequency and severity models. All the regressors are expected to have nonmissing values. If any of the regressors have a missing value in an observation, then that observation is ignored.

The information in the DATA= data set is interpreted as follows, based on whether you specify the EXTERNALCOUNTS statement:

- If you do not specify the EXTERNALCOUNTS statement, then all the observations in the data set, or in one BY group if you specify the BY statement, form a scenario. The observations are used together to compute one random draw from the compound distribution. The total number of draws is equal to the value that you specify in the NREPLICATES= option. The simulation process is described in the section “[Simulation with Regressors and No External Counts](#)” on page 76 and illustrated using an example in the section “[Illustration of Aggregate Loss Simulation Process](#)” on page 77.
- If you specify the EXTERNALCOUNTS statement, then the DATA= data set is expected to contain multiple replications (draws) of the frequency counts that you simulate externally for a scenario. The DATA= data set must contain the COUNT= variable that you specify in the EXTERNALCOUNTS statement. The replications are identified by the observation number or the ID= variable that you specify in the EXTERNALCOUNTS statement. For each observation in a given replication, the COUNT= variable is expected to contain the count of losses that are generated by the entity associated with that observation. All the observations of a given replication are used together to compute one random draw from the compound distribution. The size of the compound distribution sample is equal to the number of distinct replications that you specify in the DATA= data set, multiplied by the value that you specify in the NREPLICATES= option. The simulation process is described in the section “[Simulation with External Counts](#)” on page 78 and illustrated using an example in the section “[Illustration of the Simulation Process with External Counts](#)” on page 79.

In both cases, an observation can also contain severity adjustment variables that you can use to adjust the severity of the losses generated by that entity, based on some policy rules. For more information about simulating the adjusted compound distribution sample, see the section “[Simulation of Adjusted Compound Distribution Sample](#)” on page 82.

---

## Simulation Procedure

PROC HPCDM selects a simulation procedure based on whether you specify external counts or you request that PROC HPCDM simulate the counts, and whether the severity or frequency models contain regression effects. The following sections describe the process for the different scenarios.

### Simulation with No Regressors and No External Counts

If you specify severity and frequency models that have no regression effects in them, and if you do not specify externally simulated counts in the EXTERNALCOUNTS statement, then PROC HPCDM uses the following simulation procedure.

The process is described for one severity distribution, *dist*. If you specify multiple severity distributions in the SEVERITYMODEL statement, then the process is repeated for each specified distribution.

The following steps are repeated  $M$  times to generate a compound distribution sample of size  $M$ , where  $M$  is the value that you specify in the NREPLICATES= option or the default value of that option:

1. Use the frequency model that you specify in the COUNTSTORE= option to draw a value  $N$  from the count distribution.  $N$  is the number of loss events that are expected to occur in the time period that is being simulated.

2. Draw  $N$  values,  $X_j$  ( $j = 1, \dots, N$ ), from the severity distribution *dist* with parameters that you specify in the SEVERITYEST= data set.
3. Add the  $N$  severity values that are drawn in step 2 to compute one point  $S$  from the compound distribution as

$$S = \sum_{j=1}^N X_j$$

Note that although it is more common to fit the frequency model with regressors, PROC COUNTREG enables you to fit a frequency model without regressors. If you do not specify any regressors in the MODEL statement of the COUNTREG procedure, then it fits a model that contains only an intercept.

### Simulation with Regressors and No External Counts

If the severity or frequency models have regression effects and if you do not specify externally simulated counts in the EXTERNALCOUNTS statement, then you must specify a DATA= data set to provide values of the regression variables, which together represent a scenario for which you want to simulate the CDM. In this case, PROC HPCDM uses the following simulation procedure.

The process is described for one severity distribution. If you specify multiple severity distributions in the SEVERITYMODEL statement, then the process is repeated for each specified distribution.

Note that you are doing scenario analysis when regression effects are present. Let  $K$  denote the number of observations that form the scenario. This is the number of observations either in the current BY group or in the entire DATA= data set if you do not specify the BY statement. If  $K > 1$ , then you are modeling the scenario for a group of entities. If  $K = 1$ , then you are modeling the scenario for one entity.

The following steps are repeated  $M$  times to generate a compound distribution sample of size  $M$ , where  $M$  is the value that you specify in the NREPLICATES= option or the default value of that option:

1. For each observation  $k$  ( $k = 1, \dots, K$ ), a count  $N_k$  is drawn from the frequency model that you specify in the COUNTSTORE= option. The parameters of this model are determined by the frequency regressors in observation  $k$ .  $N_k$  represents the number of loss events that are expected to be generated by entity  $k$  in the time period that is being simulated.
2. Counts from all observations are added to compute  $N = \sum_{k=1}^K N_k$ .  $N$  is the total number of loss events that are expected to occur in the time period that is being simulated.
3.  $N$  number of random draws are made from the severity distribution, and they are added to generate one point of the compound distribution sample. Each of the  $N$  draws uses one of the  $K$  observations. If you specify a scale regression model for the severity distribution, then the scale parameter of the severity distribution is determined by the values of the severity regressors in the observation that is chosen for that draw.

If you specify the BY statement, then a separate sample of size  $M$  is created for each BY group in the DATA= data set.

### Illustration of Aggregate Loss Simulation Process

As an illustration of the simulation process, consider a very simple example of analyzing the distribution of an aggregate loss that is incurred by a set of policyholders of an automobile insurance company in a period of one year. It is postulated that the frequency and severity distributions depend on three variables: Age, Gender (1: female, 2: male), and CarType (1: sedan, 2: sport utility vehicle). So these variables are used as regressors while you fit the count model and severity scale regression model by using the COUNTREG and SEVERITY procedures, respectively. Now, consider that you want to use the fitted frequency and severity models to estimate the distribution of the aggregate loss that is incurred by a set of five policyholders together. Let the characteristics of the five policyholders be encoded in a SAS data set named Work.Scenario that has the following contents:

Obs	age	gender	carType
1	30	2	1
2	25	1	2
3	45	2	2
4	33	1	1
5	50	1	1

The column Obs contains the observation number. It is shown only for the purpose of illustration. It need not be present in the data set. The following PROC HPCDM step simulates the scenario in the Work.Scenario data set:

```
proc hpcdm data=scenario
    severityest=<severity parameter estimates data set>
    countstore=<count model store> nreplicates=<sample size>;
    severitymodel <severity distribution name(s)>;
run;
```

The following process generates a sample from the aggregate loss distribution for the scenario in the Work.Scenario data set:

1. Use the values Age=30, Gender=2, and CarType=1 in the first observation to draw a count from the count distribution. Let that count be 2. Repeat the process for the remaining four observations. Let the counts be as shown in the Count column in the following table:

Obs	age	gender	carType	count
1	30	2	1	2
2	25	1	2	1
3	45	2	2	2
4	33	1	1	3
5	50	1	1	0

Note that the Count column is shown for illustration only; it is not added as a variable to the DATA= data set.

2. The simulated counts from all the observations are added to get a value of  $N = 8$ . This means that for this particular sample point, you expect a total of eight loss events in a year from these five policyholders.

- For the first observation, the scale parameter of the severity distribution is computed by using the values Age=30, Gender=2, and CarType=1. That value of the scale parameter is used together with estimates of the other parameters from the SEVERITYEST= data set to make two draws from the severity distribution. Each of the draws simulates the magnitude of the loss that is expected from the first policyholder. The process is repeated for the remaining four policyholders. The fifth policyholder does not generate any loss event for this particular sample point, so no severity draws are made by using the fifth observation. Let the severity draws, rounded to integers for convenience, be as shown in the `_SEV_` column in the following table:

Obs	age	gender	carType	count	<code>_sev_</code>		
1	30	2	1	2	350	2100	
2	25	1	2	1	4500		
3	45	2	2	2	700	4300	
4	33	1	1	3	600	1500	950
5	50	1	1	0			

Note that the `_SEV_` column is shown for illustration only; it is not added as a variable to the DATA= data set.

PROC HPCDM adds the severity values of the eight draws to compute an aggregate loss value of 15,000. After recording this amount in the sample, the process returns to step 1 to compute the next point in the aggregate loss sample. For example, in the second iteration, the count distribution of each policyholder might generate one loss event for a total of five loss events, and the five severity draws from the severity distributions that govern each of the policyholders might add up to 5,000. Then, the value of 5,000 is recorded as the second point in the aggregate loss sample. The process continues until  $M$  aggregate loss sample points are simulated, where the  $M$  is the value that you specify in the `NREPLICATES=` option.

## Simulation with External Counts

If you specify externally simulated counts by using the `EXTERNALCOUNTS` statement, then each replication in the input data set represents the loss events generated by an entity. An entity can be an individual or organization for which you want to estimate the compound distribution. If an entity has any characteristics that are used as external factors (regressors) in developing the severity scale regression model, then you must specify the values of those factors in the DATA= data set. If you specify the `ID=` variable, then multiple observations for the same replication ID represent different entities in a group for which you are simulating the CDM.

PROC HPCDM uses the following simulation procedure in the presence of externally simulated counts.

The process is described for one severity distribution. If you specify multiple severity distributions in the `SEVERITYMODEL` statement, then the process is repeated for each specified distribution.

Let there be  $M$  distinct replications in the current BY group of the DATA= data set or in the entire DATA= data set if you do not specify the BY statement. A replication is identified by either the observation number or the value of the `ID=` variable that you specify in the `EXTERNALCOUNTS` statement.

For each of the  $M$  values of the replication identifier, the following steps are executed  $R$  times, where  $R$  is the value of the `NREPLICATES=` option or the default value of that option:

1. Compute the total number of losses,  $N$ . If there are  $K$  ( $K \geq 1$ ) observations for the current value of the replication identifier, then  $N = \sum_{k=1}^K N_k$ , where  $N_k$  is the value of the COUNT= variable for observation  $k$ .
2.  $N$  number of random draws are made from the severity distribution, and they are added to generate one point of the compound distribution sample.

This process generates a compound distribution sample of size  $M \times R$ . If you specify the BY statement, then a separate sample of size  $M \times R$  is created for each BY group in the DATA= data set.

### **Illustration of the Simulation Process with External Counts**

In order to illustrate the simulation process, consider the following simple example. In this example, your severity model does not contain any regressors. An example that uses a severity scale regression model is illustrated later. Assume that you have made 10 random draws from an external count model and recorded them in the ExtCount variable of a SAS data set named Work.Counts1 as follows:

Obs	extCount
1	3
2	2
3	0
4	1
5	3
6	4
7	1
8	2
9	0
10	5

Because the data set does not contain an ID= variable, the observation number that is shown in the Obs column acts as the replicate identifier. The following PROC HPCDM step simulates an aggregate loss sample by using the Work.Counts1 data set:

```
proc hpcdm data=work.counts1 nreplicates=5
    severityest=<severity parameter estimates data set>;
    severitymodel <severity distribution name(s)>;
    externalcounts count=extCount;
run;
```

The simulation process works as follows:

1. For the first replication, which is associated with the first observation, three severity values are drawn from the severity distribution by using the parameter estimates that you specify in the SEVERITYEST= data set. If the severity values are 150, 500, and 320, then their sum of 970 is recorded as the first point of the aggregate loss sample. Because the value of the NREPLICATES= option is 5, this process of drawing three severity values and adding them to form a point of the aggregate loss sample is repeated four more times to generate a total of five sample points that correspond to the first observation.
2. For the second replication, two severity values are drawn from the severity distribution. If the severity values are 450 and 100, then their sum of 550 is recorded as a point of the aggregate loss sample. This process of drawing two severity values and adding them to form a point of the aggregate loss sample is repeated four more times to generate a total of five sample points that correspond to the second observation.



3. The process continues until all the replications, which are observations in this case, are exhausted. PROC HPCDM ignores any replications that generate zero number of losses. In this example, replication 3 is ignored.

The process results in an aggregate loss sample of size 50, which is equal to the number of replications in the data set (10) multiplied by the value of the NREPLICATES= option (5).

Now, consider an example in which the severity models in the SEVERITYEST= data set are scale regression models. In this case, the severity distribution that is used for drawing the severity value is decided by the values of regressors in the observation that is being processed. Consider that you want to simulate the aggregate loss that is incurred by one policyholder and you have recorded, in the ExtCount variable, the results of 10 random draws from an external count model. The DATA= data set has the following contents:

Obs	age	gender	carType	extCount
1	30	2	1	5
2	30	2	1	2
3	30	2	1	0
4	30	2	1	1
5	30	2	1	3
6	30	2	1	4
7	30	2	1	1
8	30	2	1	2
9	30	2	1	0
10	30	2	1	5

The simulation process in this case is the same as the process in the previous case of no regressors, except that the severity distribution that is used for drawing the severity values has a scale parameter that is determined by the values of the regressors Age, Gender, and CarType in the observation that is being processed. In this particular example, all observations have the same value for all regressors, indicating that you are modeling a scenario in which the characteristics of the policyholder do not change during the time for which you have simulated the number of events. You can also model a scenario in which the characteristics of the policyholder change by recording those changes in the values of the appropriate regressors.

Extending this example further, consider that you want to analyze the distribution of the aggregate loss that is incurred by a group of policyholders, as in the example in the section “[Illustration of Aggregate Loss Simulation Process](#)” on page 77. Let the Work.Counts2 data set record multiple replications of the number of losses that might be generated by each policyholder. The contents of the Work.Counts2 data set are as follows:

Obs	replicateId	age	gender	carType	extCount
1	1	30	2	1	2
2	1	25	1	2	1
3	1	45	2	2	3
4	1	33	1	1	5
5	1	50	1	1	1
6	2	30	2	1	3
7	2	25	1	2	2
8	2	33	1	1	4
9	2	50	1	1	1

The ReplicateId variable records the identifier for the replication. Each replication contains multiple observations, such that each observation represents one of the policyholders that you are analyzing. For



simplicity, only the first two replications are shown here. Note that you do not need to record an observation if the number of loss events that are generated by a policyholder is 0, as shown for the second replication, which does not have any data for the 45-year-old policyholder.

The following PROC HPCDM step simulates an aggregate loss sample by using the Work.Counts2 data set:

```
proc hpcdm data=work.counts2 nreplicates=3
    severityest=<severity parameter estimates data set>;
    severitymodel <severity distribution name(s)>;
    distby replicateId;
    externalcounts count=extCount id=replicateId;
    output out=aggloss samplevar=totalLoss;
run;
```

When you specify an ID= variable in the EXTERNALCOUNTS statement, you must specify the same ID= variable in the DISTBY statement in order for the procedure to work correctly in a distributed computing environment. Further, the DATA= set must be sorted in ascending order of the ID= variable values.

The simulation process works as follows:

1. First, the five observations of the first replication (ReplicateId=1) are analyzed. For the first observation (Obs=1), the scale parameter of the severity distribution is computed by using the values Age=30, Gender=2, and CarType=1. That value of the scale parameter is used together with estimates of the other parameters from the SEVERITYEST= data set to make two draws from the severity distribution. Next, the regressor values of the second observation are used to compute the scale parameter of the severity distribution, which is used to make one severity draw. The process continues such that the regressor values in the third, fourth, and fifth observations are used to decide the severity distribution to make three, five, and one draws from, respectively. Let the severity values that are drawn from the observations of this replication be as shown in the `_SEV_` column in the following table, where the `_SEV_` column is shown for illustration only; it is not added as a variable to the DATA= data set:

Obs	replicateId	age	gender	carType	extCount	_sev_
1	1	30	2	1	2	700 500
2	1	25	1	2	1	5000
3	1	45	2	2	3	900 1400 300
4	1	33	1	1	5	350 2000 150 800 600
5	1	50	1	1	1	250

The values of all 12 severity draws are added to compute and record the value of 12,950 as the first point of the aggregate loss sample. Because you specify NREPLICATES=3 in the PROC HPCDM step, this process of making 12 severity draws from the respective observations is repeated two more times to generate a total of three sample points for the first replication.

2. The four observations of the second replication (ReplicateId=2) are analyzed next to draw three, two, four, and one severity values from the severity distributions, with scale parameters that are decided by the regressor values in the sixth, seventh, eighth, and ninth observations, respectively. The 10 severity values are added to form a point of the aggregate loss sample. This process of making 10 severity draws from the respective observations is repeated two more times to generate a total of three sample points for the second replication.

If your Work.Counts2 data set contains 10,000 distinct values of ReplicateId, then 30,000 observations are written to the Work.AgglLoss data set that you specify in the OUTPUT statement of the preceding PROC HPCDM step. Because you specify SAMPLEVAR=TotalLoss in the OUTPUT statement, the aggregate loss sample is available in the TotalLoss column of the Work.AgglLoss data set.

## Simulation of Adjusted Compound Distribution Sample

If you specify programming statements that adjust the severity value, then a separate adjusted compound distribution sample is also generated.

Your programming statements are expected to implement an adjustment function  $f$  that uses the unadjusted severity value,  $X_j$ , to compute and return an adjusted severity value,  $X_j^a$ . To compute  $X_j^a$ , you might also use the sum of unadjusted severity values and the sum of adjusted severity values.

Formally, if  $N$  denotes the number of loss events that are to be simulated for the current replication of the simulation process, then for the severity draw,  $X_j$ , of the  $j$ th loss event ( $j = 1, \dots, N$ ), the adjusted severity value is

$$X_j^a = f(X_j, S_{j-1}, S_{j-1}^a)$$

where  $S_{j-1} = \sum_{l=1}^{j-1} X_l$  is the aggregate unadjusted loss before  $X_j$  is generated and  $S_{j-1}^a = \sum_{l=1}^{j-1} X_l^a$  is the aggregate adjusted loss before  $X_j$  is generated. The initial values of both types of aggregate losses are set to 0. In other words,  $S_0 = 0$  and  $S_0^a = 0$ .

The aggregate adjusted loss for the replication is  $S_N^a$ , which is denoted by  $S^a$  for simplicity, and is defined as

$$S^a = \sum_{j=1}^N X_j^a$$

In your programming statements that implement  $f$ , you can use the following keywords as placeholders for the input arguments of the function  $f$ :

### \_SEV\_

indicates the placeholder for  $X_j$ , the unadjusted severity value. PROC HPCDM generates this value as described in the section “Simulation with No Regressors and No External Counts” on page 75 (step 2) or the section “Simulation with Regressors and No External Counts” on page 76 (step 3). PROC HPCDM supplies this value to your program.

### \_CUMSEV\_

indicates the placeholder for  $S_{j-1}$ , the sum of unadjusted severity values that PROC HPCDM generates before  $X_j$  is generated. PROC HPCDM supplies this value to your program.

### \_CUMADJSEV\_

indicates the placeholder for  $S_{j-1}^a$ , the sum of adjusted severity values that are computed by your programming statements before  $X_j$  is generated and adjusted. PROC HPCDM supplies this value to your program.

In your programming statements, you must assign the value of  $X_j^a$ , the output of function  $f$ , to a symbol that you specify in the ADJUSTEDSEVERITY= option in the PROC HPCDM statement. PROC HPCDM uses the final assigned value of this symbol as the value of  $X_j^a$ .

You can use most DATA step statements and functions in your program. The DATA step file and the data set I/O statements (for example, INPUT, FILE, SET, and MERGE) are not available. However, some functionality of the PUT statement is supported. For more information, see the section “PROC FCMP and DATA Step Differences” in *Base SAS Procedures Guide*.

The simulation process that generates the aggregate adjusted loss sample is identical to the process that is described in the section “[Simulation with Regressors and No External Counts](#)” on page 76 or the section “[Simulation with External Counts](#)” on page 78, except that after making each of the  $N$  severity draws, PROC HPCDM executes your severity adjustment programming statements to compute the adjusted severity ( $X_j^a$ ). All the  $N$  adjusted severity values are added to compute  $S^a$ , which forms a point of the aggregate adjusted loss sample. The process is illustrated using an example in the section “[Illustration of Aggregate Adjusted Loss Simulation Process](#)” on page 85.

## Using Severity Adjustment Variables

If you do not specify the DATA= data set, then your ability to adjust the severity value is limited, because you can use only the current severity draw, sums of unadjusted and adjusted severity draws that are made before the current draw, and some constant numbers to encode your adjustment policy. That is sufficient if you want to estimate the distribution of aggregate adjusted loss for only one entity. However, if you are simulating a scenario that contains more than one entity, then it might be more useful if the adjustment policy depends on factors that are specific to each entity that you are simulating. To do that, you must specify the DATA= data set and encode such factors as *adjustment variables* in the DATA= data set. Let  $A$  denote the set of values of the adjustment variables. Then, the form of the adjustment function  $f$  that computes the adjusted severity value becomes

$$X_j^a = f(X_j, S_{j-1}, S_{j-1}^a, A)$$

PROC HPCDM reads the values of adjustment variables from the DATA= data set and supplies the set of those values ( $A$ ) to your severity adjustment program. For an invocation of  $f$  with an unadjusted severity value of  $X_j$ , the values in set  $A$  are read from the same observation that is used to simulate  $X_j$ .

All adjustment variables that you use in your program must be present in the DATA= data set. You must not use any keyword for a placeholder symbol as a name of any variable in the DATA= data set, whether the variable is a severity adjustment variable or a regressor in the frequency or severity model. Further, the following restrictions apply to the adjustment variables:

- You can use only numeric-valued variables in PROC HPCDM programming statements. This restriction also implies that you cannot use SAS functions or call routines that require character-valued arguments, unless you pass those arguments as constant (literal) strings or characters.
- You cannot use functions that create lagged versions of a variable in PROC HPCDM programming statements. If you need lagged versions, then you can use a DATA step before the PROC HPCDM step to add those versions to the input data set.

The use of adjustment variables is illustrated using an example in the section “[Illustration of Aggregate Adjusted Loss Simulation Process](#)” on page 85.

### Aggregate Adjusted Loss Simulation for a Multi-entity Scenario

If you are simulating a scenario that consists of multiple entities, then you can use some additional pieces of information in your severity adjustment program. Let the scenario consist of  $K$  entities and let  $N_k$  denote the number of loss events that are incurred by  $k$ th entity ( $k = 1, \dots, K$ ) in the current iteration of the simulation process. The total number of severity draws that need to be made is  $N = \sum_{k=1}^K N_k$ . The aggregate adjusted loss is now defined as

$$S^a = \sum_{k=1}^K \sum_{j=1}^{N_k} X_{k,j}^a$$

where  $X_{k,j}^a$  is an adjusted severity value of the  $j$ th draw ( $j = 1, \dots, N_k$ ) for the  $k$ th entity, and the form of the adjustment function  $f$  that computes  $X_{k,j}^a$  is

$$X_{k,j}^a = f(X_{k,j}, S_{k,j-1}, S_{k,j-1}^a, S_{n-1}, S_{n-1}^a, A)$$

where  $X_{k,j}$  is the value of the  $j$ th draw of unadjusted severity for the  $k$ th entity.  $S_{k,j-1} = \sum_{l=1}^{j-1} X_{k,l}$  and  $S_{k,j-1}^a = \sum_{l=1}^{j-1} X_{k,l}^a$  are the aggregate unadjusted loss and the aggregate adjusted loss, respectively, for the  $k$ th entity before  $X_{k,j}$  is generated. The index  $n$  ( $n = 1, \dots, N$ ) keeps track of the total number of severity draws, across all entities, that are made before  $X_{k,j}$  is generated. So  $S_{n-1} = \sum_{l=1}^{n-1} X_l$  and  $S_{n-1}^a = \sum_{l=1}^{n-1} X_l^a$  are the aggregate unadjusted loss and aggregate adjusted loss, respectively, for all the entities that are processed before  $X_{k,j}$  is generated. Note that  $S_{n-1}$  and  $S_{n-1}^a$  include the  $j - 1$  draws that are made for the  $k$ th entity before  $X_{k,j}$  is generated.

The initial values of all types of aggregate losses are set to 0. In other words,  $S_0 = 0$ ,  $S_0^a = 0$ , and for all values of  $k$ ,  $S_{k,0} = 0$  and  $S_{k,0}^a = 0$ .

PROC HPCDM uses the final value that you assign to the **ADJUSTEDSEVERITY=** symbol in your programming statements as the value of  $X_{k,j}^a$ .

In your severity adjustment program, you can use the following two additional placeholder keywords:

#### **\_CUMSEVFOROBS\_**

indicates the placeholder for  $S_{k,j-1}$ , which is the total loss that is incurred by the  $k$ th entity before the current loss event. PROC HPCDM supplies this value to your program.

#### **\_CUMADJSEVFOROBS\_**

indicates the placeholder for  $S_{k,j-1}^a$ , which is the total adjusted loss that is incurred by the  $k$ th entity before the current loss event. PROC HPCDM supplies this value to your program.

The previously described placeholder symbols **\_CUMSEV\_** and **\_CUMADJSEV\_** represent  $S_{n-1}$  and  $S_{n-1}^a$ , respectively. If you have only one entity in the scenario ( $K = 1$ ), then the values of **\_CUMSEVFOROBS\_** and **\_CUMADJSEVFOROBS\_** are identical to the values of **\_CUMSEV\_** and **\_CUMADJSEV\_**, respectively.

There is one caveat when a scenario consists of more than one entity ( $K > 1$ ) and when you use any of the symbols for cumulative severity values (**\_CUMSEV\_**, **\_CUMADJSEV\_**, **\_CUMSEVFOROBS\_**, or **\_CUMADJSEVFOROBS\_**) in your programming statements. In this case, to make the simulation realistic, it is important to randomize the order of  $N$  severity draws across  $K$  entities. For more information, see the section “Randomizing the Order of Severity Draws across Observations of a Scenario” on page 87.

## Illustration of Aggregate Adjusted Loss Simulation Process

This section continues the example in the section “Simulation with Regressors and No External Counts” on page 76 to illustrate the simulation of aggregate adjusted loss.

Recall that the earlier example simulates a scenario that consists of five policyholders. Assume that you want to compute the distribution of the aggregate amount paid to all the policyholders in a year, where the payment for each loss is decided by a deductible and a per-payment limit. To begin with, you must record the deductible and limit information in the input DATA= data set. The following table shows the DATA= data set from the earlier example, extended to include two variables, Deductible and Limit:

Obs	age	gender	carType	deductible	limit
1	30	2	1	250	5000
2	25	1	2	500	3000
3	45	2	2	100	2000
4	33	1	1	500	5000
5	50	1	1	200	2000

The variables Deductible and Limit are referred to as severity adjustment variables, because you need to use them to compute the adjusted severity. Let AmountPaid represent the value of adjusted severity that you are interested in. Further, let the following SAS programming statements encode your logic of computing the value of AmountPaid:

```
amountPaid = MAX(_sev_ - deductible, 0);
amountPaid = MIN(amountPaid, MAX(limit - _cumadjsevforobs_, 0));
```

PROC HPCDM supplies your program with values of the placeholder symbols `_SEV_` and `_CUMADJSEVFOROBS_`, which represent the value of the current unadjusted severity draw and the sum of adjusted severity values from the previous draws, respectively, for the observation that is being processed. The use of `_CUMADJSEVFOROBS_` helps you ensure that the payment that is made to a given policyholder in a year does not exceed the limit that is recorded in the Limit variable.

In order to simulate a sample for the aggregate of AmountPaid, you need to submit a PROC HPCDM step whose structure is like the following:

```
proc hpcdm data=<data set name> adjustedseverity=amountPaid
    severityest=<severity parameter estimates data set>
    countstore=<count model store>;
    severitymodel <severity distribution name(s)>;

    amountPaid = MAX(_sev_ - deductible, 0);
    amountPaid = MIN(amountPaid, MAX(limit - _cumadjsevforobs_, 0));
run;
```

The simulation process of one replication that generates one point of the aggregate loss sample and the corresponding point of the aggregate adjusted loss sample is as follows:

1. Use the values Age=30, Gender=2, and CarType=1 in the first observation to draw a count from the count distribution. Let that count be 3. Repeat the process for the remaining four observations. Let the counts be as shown in the Count column in the following table:

Obs	age	gender	carType	deductible	limit	count
1	30	2	1	250	5000	2
2	25	1	2	500	3000	1
3	45	2	2	100	2000	2
4	33	1	1	500	5000	3
5	50	1	1	200	2000	0

Note that the Count column is shown for illustration only; it is not added as a variable to the DATA= data set.

- The simulated counts from all the observations are added to get a value of  $N = 8$ . This means that for this particular replication, you expect a total of eight loss events in a year from these five policyholders.
- For the first observation, the scale parameter of the severity distribution is computed by using the values Age=30, Gender=2, and CarType=1. That value of the scale parameter is used together with estimates of the other parameters from the SEVERITYEST= data set to make two draws from the severity distribution. The process is repeated for the remaining four policyholders. The fifth policyholder does not generate any loss event for this particular replication, so no severity draws are made by using the fifth observation. Let the severity draws, rounded to integers for convenience, be as shown in the `_SEV_` column in the following table, where the `_SEV_` column is shown for illustration only; it is not added as a variable to the DATA= data set:

Obs	age	gender	carType	deductible	limit	count	<code>_sev_</code>	
1	30	2	1	250	5000	2	350	2100
2	25	1	2	500	3000	1	4500	
3	45	2	2	100	2000	2	700	4300
4	33	1	1	200	5000	3	600	1500 950
5	50	1	1	200	2000	0		

The sample point for the aggregate unadjusted loss is computed by adding the severity values of eight draws, which gives an aggregate loss value of 15,000. The unadjusted aggregate loss is also referred to as the ground-up loss.

For each of the severity draws, your severity adjustment programming statements are executed to compute the adjusted severity, which is the value of AmountPaid in this case. For the draws in the preceding table, the values of AmountPaid are as follows:

Obs	deductible	limit	<code>_sev_</code>	<code>_cumadjsevforobs_</code>	amountPaid
1	250	5000	350	0	100
1	250	5000	2100	100	1850
2	500	3000	4500	0	3000
3	100	2000	700	0	600
3	100	2000	4300	600	1400
4	200	5000	600	0	400
4	200	5000	1500	400	1300
4	200	5000	950	1700	750

The adjusted severity values are added to compute the cumulative payment value of 9,400, which forms the first sample point for the aggregate adjusted loss.

After recording the aggregate unadjusted and aggregate adjusted loss values in their respective samples, the process returns to step 1 to compute the next sample point unless the specified number of sample points have been simulated.

In this particular example, you can verify that the order in which the 8 loss events are simulated does not affect the aggregate adjusted loss. As a simple example, consider the following order of draws that is different from the consecutive order that was used in the preceding table:

Obs	deductible	limit	_sev_	_cumadjsevforobs_	amountPaid
4	200	5000	600	0	400
3	100	2000	4300	0	2000
1	250	5000	350	0	100
3	100	2000	700	2000	0
4	200	5000	950	400	750
1	250	5000	2100	100	1850
2	500	3000	4500	0	3000
4	200	5000	1500	1150	1300

Although the payments that are made for individual loss events differ, the aggregate adjusted loss is still 9,400.

However, in general, when you use a cumulative severity value such as `_CUMADJSEVFOROBS_` in your program, the order in which the draws are processed affects the final value of aggregate adjusted loss. For more information, see the sections “[Randomizing the Order of Severity Draws across Observations of a Scenario](#)” on page 87 and “[Illustration of the Need to Randomize the Order of Severity Draws](#)” on page 88.

## Randomizing the Order of Severity Draws across Observations of a Scenario

If you specify a scenario that consists of a group of more than one entity, then it is assumed that each entity generates its loss events independently from other entities. In other words, the time at which the loss event of one entity is generated or recorded is independent of the time at which the loss event of another entity is generated or recorded. If entity  $k$  generates  $N_k$  loss events, then the total number of loss events for a group of  $K$  entities is  $N = \sum_{k=1}^K N_k$ . To simulate the aggregate loss for this group,  $N$  severity draws are made and aggregated to compute one point of the compound distribution sample. However, to honor the assumption of independence among entities, the order of those  $N$  severity draws must be randomized across  $K$  entities such that no entity is preferred over another.

The  $K$  entities are represented by  $K$  observations of the scenario in the `DATA=` data set. If you specify external counts, the  $K$  observations correspond to the observations that have the same replication identifier value. If you do not specify the external counts, then the  $K$  observations correspond to all the observations in the `BY` group or in the entire `DATA=` set if you do not specify the `BY` statement.

The randomization process over  $K$  observations is implemented as follows. First, one of the  $K$  observations is chosen at random and one severity value is drawn from the severity distribution implied by that observation, then another observation is chosen at random and one severity value is drawn from its implied severity distribution, and so on. In each step, the total number of events that are simulated for the selected observation  $k$  is incremented by 1. When all  $N_k$  events for an observation  $k$  are simulated, observation  $k$  is retired and the process continues with the remaining observations until a total of  $N$  severity draws are made. Let  $k(j)$  denote a function that implements this randomization by returning an observation  $k$  ( $k = 1, \dots, K$ ) for the  $j$ th draw ( $j = 1, \dots, N$ ). The aggregate loss computation can then be formally written as

$$S = \sum_{j=1}^N X_{k(j)}$$

where  $X_{k(j)}$  denotes the severity value that is drawn by using observation  $k(j)$ .



If you do not specify a scale regression model for severity, then all severity values are drawn from the same severity distribution. However, if you specify a scale regression model for severity, then the severity draw is made from the severity distribution that is determined by the values of regressors in observation  $k$ . In particular, the scale parameter of the distribution depends on the values of regressors in observation  $k$ . If  $R(l)$  denotes the scale regression model for observation  $l$  and  $X_{R(l)}$  denotes the severity value drawn from scale regression model  $R(l)$ , then the aggregate loss computation can be formally written as

$$S = \sum_{j=1}^N X_{R(k(j))}$$

This randomization process is especially important in the context of simulating an adjusted compound distribution sample when your severity adjustment program uses the aggregate adjusted severity observed so far to adjust the next severity value. For an illustration of the need to randomize in such cases, see the next section.

### **Illustration of the Need to Randomize the Order of Severity Draws**

This section uses the example of the section “Illustration of Aggregate Adjusted Loss Simulation Process” on page 85, but with the following PROC HPCDM step:

```
proc hpcdm data=<data set name> adjustedseverity=amountPaid
    severityest=<severity parameter estimates data set>
    countstore=<count model store>;
    severitymodel <severity distribution name(s)>;

    if (_cumadjsev_ > 15000) then
        amountPaid = 0;
    else do;
        penaltyFactor = MIN(3, 15000/(15000 - _cumadjsev_));
        amountPaid = MAX(0, _sev_ - deductible * penaltyFactor);
    end;
run;
```

The severity adjustment statements in the preceding steps compute the value of AmountPaid by using the following provisions in the insurance policy:

- There is a limit of 15,000 on the total amount that can be paid in a year to the group of policyholders that is being simulated. The amount of payment for each loss event depends on the total amount of payments before that loss event.
- The penalty for incurring more losses is imposed in the form of an increased deductible. In particular, the deductible is increased by the ratio of the maximum cumulative payment (15,000) to the amount that remains available to pay for future losses in the year. The factor by which the deductible can be raised has a limit of three.

This example illustrates only step 3 of the simulation process, where randomization is done. It assumes that step 2 of the simulation process is identical to the step 2 in the example in the section “Illustration of Aggregate Adjusted Loss Simulation Process” on page 85. At the beginning of step 3, let the severity draws from all the observations be as shown in the `_SEV_` column in the following table:



Obs	age	gender	carType	deductible	count	_sev_	
1	30	2	1	250	2	350	2100
2	25	1	2	500	1	4500	
3	45	2	2	100	2	700	4300
4	33	1	1	200	3	600	1500 950
5	50	1	1	200	0		

If the order of these eight draws is not randomized, then all the severity draws for the first observation are adjusted before all the severity draws of the second observation, and so on. The execution of the severity adjustment program leads to the following sequence of values for AmountPaid:

Obs	deductible	_sev_	_cumadjsev_	penaltyFactor	amountPaid
1	250	350	0	1	100
1	250	2100	100	1.0067	1848.32
2	500	4500	1948.32	1.1493	3925.36
3	100	700	5873.68	1.6436	535.64
3	100	4300	6409.32	1.7461	4125.39
4	200	600	10534.72	3	0
4	200	1500	10534.72	3	900
4	200	950	11434.72	3	350

The preceding sequence of simulating loss events results in a cumulative payment of 11,784.72.

If the sequence of draws is randomized over observations, then the computation of the cumulative payment might proceed as follows for one instance of randomization:

Obs	deductible	_sev_	_cumadjsev_	penaltyFactor	amountPaid
2	500	4500	0	1	4000
1	250	350	4000	1.3636	9.09
3	100	700	4009.09	1.3648	563.52
4	200	950	4572.61	1.4385	662.30
4	200	1500	5234.91	1.5361	1192.78
1	250	2100	6427.69	1.7498	1662.54
4	200	600	8090.24	2.1708	165.83
3	100	4300	8256.07	2.2242	4077.58

In this example, a policyholder is identified by the value in the Obs column. As the table indicates, PROC HPCDM randomizes the order of loss events not only across policyholders but also across the loss events that a given policyholder incurs. The particular sequence of loss events that is shown in the table results in a cumulative payment of 12,333.65. This differs from the cumulative payment that results from the previously considered nonrandomized sequence of loss events, which tends to penalize the fourth policyholder by always processing her payments after all other payments, with a possibility of underestimating the total paid amount. This comparison not only illustrates that the order of randomization affects the aggregate adjusted loss sample but also corroborates the arguments about the importance of order randomization that are made at the beginning of the section “[Randomizing the Order of Severity Draws across Observations of a Scenario](#)” on page 87.

## Parameter Perturbation Analysis

It is important to realize that most of the parameters of the frequency and severity models are estimated and there is uncertainty associated with the parameter estimates. Any compound distribution estimate that is computed by using these uncertain parameter estimates is inherently uncertain. The aggregate loss

sample that is simulated by using the mean estimates of the parameters is just one possible sample from the compound distribution. If information about parameter uncertainty is available, then it is recommended that you conduct parameter perturbation analysis that generates multiple samples of the compound distribution, in which each sample is simulated by using a set of perturbed parameter estimates. You can use the `NPURTURBEDSAMPLES=` option in the PROC HPCDM statement to specify the number of perturbed samples to be generated. The set of perturbed parameter estimates is created by making a random draw of the parameter values from their joint probability distribution. If you specify `NPURTURBEDSAMPLES=P`, then PROC HPCDM creates  $P$  sets of perturbed parameters and each set is used to simulate a full aggregate sample. The summary analysis of  $P$  such aggregate loss samples results in a set of  $P$  estimates for each summary statistic and percentile of the compound distribution. The mean and standard deviation of this set of  $P$  estimates quantify the uncertainty that is associated with the compound distribution.

The parameter uncertainty information is available in the form of either the variance-covariance matrix of the parameter estimates or standard errors of the parameters estimates. If the variance-covariance matrix is available and is positive definite, then PROC HPCDM assumes that the joint probability distribution of the parameter estimates is a multivariate normal distribution,  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , where the mean vector  $\boldsymbol{\mu}$  is the set of point parameter estimates and  $\boldsymbol{\Sigma}$  is the variance-covariance matrix. If the variance-covariance matrix is not available or is not positive definite, then PROC HPCDM assumes that each parameter has a univariate normal distribution,  $\mathcal{N}(\mu, \sigma^2)$ , where  $\mu$  is the point estimate of the parameter and  $\sigma$  is the standard error of the parameter estimate.

For severity models, the point parameter estimates are expected to be available in the `SEVERITYEST=` data set in observations for which `_TYPE_='EST'`, the standard errors are expected to be available in the `SEVERITYEST=` data set in observations for which `_TYPE_='STDERR'`, and the variance-covariance matrix is expected to be available in the `SEVERITYEST=` data set in observations for which `_TYPE_='COV'`. If you use the SEVERITY procedure to create the `SEVERITYEST=` data set, then you need to specify the `COVOUT` option in the PROC SEVERITY statement to make the variance-covariance estimates available in the `SEVERITYEST=` data set.

For the frequency model, you must use the COUNTREG procedure to create the `COUNTSTORE=` item store, which always contains the point estimates, standard errors, and variance-covariance matrix of the parameters.

If you specify the `ADJUSTEDSEVERITY=` option in the PROC HPCDM statement, then a separate perturbation analysis is conducted for the distribution of the aggregate adjusted loss.

---

## Descriptive Statistics

This section provides computational details for the descriptive statistics that are computed for each aggregate loss sample. You can also save these statistics in an `OUTSUM=` data set by specifying appropriate keywords in the `OUTSUM` statement.

This section gives specific details about the moment statistics. For more information about the methods of computing percentile statistics, see the description of the `PCTLDEF=` option in the UNIVARIATE procedure in the *Base SAS Procedures Guide: Statistical Procedures*.

Standard algorithms (Fisher 1973) are used to compute the moment statistics. The computational methods that the HPCDM procedure uses are consistent with those that other SAS procedures use for calculating descriptive statistics.

## Mean

The sample mean is calculated as

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

where  $n$  is the size of the generated aggregate loss sample and  $y_i$  is the  $i$ th value of the aggregate loss.

## Standard Deviation

The standard deviation is calculated as

$$s = \sqrt{\frac{1}{d} \sum_{i=1}^n (y_i - \bar{y})^2}$$

where  $n$  is the size of the generated aggregate loss sample,  $y_i$  is the  $i$ th value of the aggregate loss,  $\bar{y}$  is the sample mean, and  $d$  is the divisor controlled by the **VARDEF=** option in the PROC HPCDM statement:

$$d = \begin{cases} n - 1 & \text{if VARDEF=DF (default)} \\ n & \text{if VARDEF=N} \end{cases}$$

## Skewness

The sample skewness, which measures the tendency of the deviations to be larger in one direction than in the other, is calculated as

$$\frac{1}{d_s} \sum_{i=1}^n \left( \frac{y_i - \bar{y}}{s} \right)^3$$

where  $n$  is the size of the generated aggregate loss sample,  $y_i$  is the  $i$ th value of the aggregate loss,  $\bar{y}$  is the sample mean,  $s$  is the sample standard deviation, and  $d_s$  is the divisor controlled by the **VARDEF=** option in the PROC HPCDM statement:

$$d_s = \begin{cases} \frac{(n-1)(n-2)}{n} & \text{if VARDEF=DF (default)} \\ n & \text{if VARDEF=N} \end{cases}$$

If VARDEF=DF, then  $n$  must be greater than 2.

The sample skewness can be positive or negative; it measures the asymmetry of the data distribution and estimates the theoretical skewness  $\sqrt{\beta_1} = \mu_3 \mu_2^{-\frac{3}{2}}$ , where  $\mu_2$  and  $\mu_3$  are the second and third central moments. Observations that are normally distributed should have a skewness near zero.

## Kurtosis

The sample kurtosis, which measures the heaviness of tails, is calculated as in [Table 4.2](#) depending on the value that you specify in the **VARDEF=** option.

**Table 4.2** Formulas for Kurtosis

VARDEF Value	Formula
DF (default)	$\frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum_{i=1}^n \left( \frac{y_i - \bar{y}}{s} \right)^4 - \frac{3(n-1)^2}{(n-2)(n-3)}$
N	$\frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \bar{y}}{s} \right)^4 - 3$

In these formulas,  $n$  is the size of the generated aggregate loss sample,  $y_i$  is the  $i$ th value of the aggregate loss,  $\bar{y}$  is the sample mean, and  $s$  is the sample standard deviation. If VARDEF=DF, then  $n$  must be greater than 3.

The sample kurtosis measures the heaviness of the tails of the data distribution. It estimates the adjusted theoretical kurtosis denoted as  $\beta_2 - 3$ , where  $\beta_2 = \frac{\mu_4}{\mu_2^2}$  and  $\mu_4$  is the fourth central moment. Observations that are normally distributed should have a kurtosis near zero.

## Input Specification

PROC HPCDM accepts the DATA= and SEVERITYEST= data sets and the COUNTSTORE= item store as input. This section details the information that they are expected to contain.

### DATA= Data Set

If you specify the BY statement, then the DATA= data set must contain all the BY variables that you specify in the BY statement and the data set must be sorted by the BY variables unless the BY statement includes the NOTSORTED option.

If the severity models in the SEVERITYEST= data set contain any scale regressors, then all those regressors must be present in the DATA= data set.

If you specify the programming statements to compute an aggregate adjusted loss, and if your specified ADJUSTEDSEVERITY= symbol depends on severity adjustment variables, then the DATA= data set must contain all such variables.

The rest of the contents of the DATA= data set depends on whether you specify the EXTERNALCOUNTS statement. If you specify the EXTERNALCOUNTS statement, then the DATA= data set is expected to contain the COUNT= and ID= variables that you specify in the EXTERNALCOUNTS statement. If you do not specify the EXTERNALCOUNTS statement, then the DATA= data set must contain all the regressors, including zero model regressors, that are present in the count model that the COUNTSTORE= item store contains.

You do not need to specify the DATA= data set if all the following conditions are true:

- You do not specify the BY statement.
- You specify a SEVERITYEST= data set such that none of the severity models are scale regression models.
- You do not specify the EXTERNALCOUNTS statement.
- You specify a COUNTSTORE= item store such that the count model contains no count regressors.
- Your severity adjustment programming statements, if you specify any, do not use any external input.

### SEVERITYEST= Data Set

The SEVERITYEST= data set is expected to contain the parameter estimates of the severity models. This is a required data set; you must specify it whenever you use PROC HPCDM.

The SEVERITYEST= data set must have the same format as the OUTEST= data set that is created by the SEVERITY procedure. For more information, see the description of the OUTEST= data set in the SEVERITY procedure in the *SAS/ETS User's Guide*.

If you specify the BY statement, then the SEVERITYEST= data set must contain all the BY variables that you specify in the BY statement. If you do not specify the NOTSORTED option in the BY statement, then the SEVERITYEST= data set must be sorted by the BY variables.

### COUNTSTORE= Item Store

The COUNTSTORE= item store is expected to be created by using the STORE statement in the COUNTREG procedure. You must specify the COUNTSTORE= item store when you do not specify the EXTERNALCOUNTS statement. For more information, see the description of the STORE statement in the COUNTREG procedure in the *SAS/ETS User's Guide*.

---

## Output Data Sets

PROC HPCDM writes the output data sets that you specify in the OUT= option of the **OUTPUT** and **OUTSUM** statements. The contents of these output data sets are described in the sections “**OUTSAMPLE= Data Set**” on page 93 and “**OUTSUM= Data Set**” on page 94, respectively.

### OUTSAMPLE= Data Set

The OUTSAMPLE= data set records the full sample of the aggregate loss and aggregate adjusted loss.

If you specify the BY statement, then the data are organized in BY groups and the data set contains variables that you specify in the BY statement. In addition, the OUTSAMPLE= data set contains the following variables:

`_SEVERITYMODEL_`

indicates the name of the severity distribution model.

**\_COUNTMODEL\_**

indicates the name of the count model. If you specify the EXTERNALCOUNTS statement, then the value of this variable is “\_EXTERNAL\_”. If you specify the COUNTSTORE= option, then the value of this variable is “\_COUNTSTORE\_”.

**<unadjusted sample variable>**

indicates the value of the unadjusted aggregate loss. The name of this variable is the value of the **SAMPLEVAR=** option in the OUTPUT statement. If you do not specify the SAMPLEVAR= option, then the variable is named **\_AGGSEV\_**.

**<adjusted sample variable>**

indicates the value of the adjusted aggregate loss. This variable is created only when you specify the programming statements and the **ADJUSTEDSEVERITY=** option in the PROC HPCDM statement. The name of this variable is the value of the **ADJSAMPLEVAR=** option in the OUTPUT statement. If you do not specify the ADJSAMPLEVAR= option, then the variable is named **\_AGGADJSEV\_**.

**\_DRAWID\_**

indicates the identifier for the perturbed sample. This variable is created only when you specify the **NPERTURBEDSAMPLES=** option in the PROC HPCDM statement. The value of this variable identifies the perturbed sample. A value of 0 for the **\_DRAWID\_** variable indicates an unperturbed sample.

**OUTSUM= Data Set**

The OUTSUM= data set records the summary statistics and percentiles of the compound distributions of aggregate loss and aggregate adjusted loss. Only the estimates that you request in the OUTSUM statement are written to the OUTSUM= data set. For more information about the method of naming the variables that correspond to the summary statistics or percentiles, see the description of the **OUTSUM** statement.

If you specify the BY statement, then the data are organized in BY groups and the data set contains variables that you specify in the BY statement. In addition, the OUTSUM= data set contains the following variables:

**\_SEVERITYMODEL\_**

indicates the name of the severity distribution model.

**\_COUNTMODEL\_**

indicates the name of the count model. If you specify the EXTERNALCOUNTS statement, then the value of this variable is “\_EXTERNAL\_”. If you specify the COUNTSTORE= option, then the value of this variable is “\_COUNTSTORE\_”.

**\_SAMPLEVAR\_**

indicates the name of the aggregate loss sample. For an unadjusted sample, the value of the variable is the value of the **SAMPLEVAR=** option that you specify in the OUTPUT statement or the default value of “\_AGGSEV\_”. For an adjusted sample, the value of the variable is the value of the **ADJSAMPLEVAR=** option that you specify in the OUTPUT statement or the default value of “\_AGGADJSEV\_”.

**\_DRAWID\_**

indicates the identifier for the perturbed sample. This variable is created only when you specify the **NPERTURBEDSAMPLES=** option in the PROC HPCDM statement. The value of this variable identifies the perturbed sample. A value of 0 for **\_DRAWID\_** indicates an unperturbed sample.

## Displayed Output

The HPCDM procedure optionally produces displayed output by using the Output Delivery System (ODS). All output is controlled by the PRINT= option in the PROC HPCDM statement. Table 4.3 relates the PRINT= options to ODS tables.

**Table 4.3** ODS Tables Produced in PROC HPCDM

ODS Table Name	Description	Option
CompoundInfo	Compound distribution information	Default
DataSummary	Input data summary	Default
Percentiles	Percentiles of the aggregate loss sample	PRINT=PERCENTILES
PerformanceInfo	Execution environment information that pertains to the computational performance	Default
PerturbedPctlSummary	Perturbation analysis of percentiles	PRINT=PERTURBSUMMARY and NPerturbedSamples > 0
PerturbedSummary	Perturbation analysis of summary statistics	PRINT=PERTURBSUMMARY and NPerturbedSamples > 0
SummaryStatistics	Summary statistics of the aggregate loss sample	PRINT=SUMMARYSTATISTICS
Timing	Timing information for various computational stages of the procedure	DETAILS (PERFORMANCE statement)

### PRINT= Option

This section provides detailed descriptions of the tables that are displayed by using different PRINT= options.

- If you do not specify the PRINT= option and if you do not specify the NOPRINT or PRINT=NONE options, then by default PROC HPCDM produces the CompoundInfo, DataSummary, and SummaryStatistics ODS tables.

The “Compound Distribution Information” table (ODS name: CompoundInfo) displays the information about the severity and count models.

The “Input Data Summary” table (ODS name: DataSummary) is displayed when you specify the DATA= data set. The table displays the total number of observations and the valid number of observations in the data set. If you specify the EXTERNALCOUNTS statement, then the table also displays the number of replications and total number of loss events across all replications.

- If you specify PRINT=PERCENTILES, the “Percentiles” table (ODS name: Percentiles) is displayed for the distribution of the aggregate loss. The table contains estimates of all the predefined percentiles in addition to the percentiles that you request in the OUTSUM statement.

If you specify the programming statements and the ADJUSTEDSEVERITY= symbol, then an additional table is displayed for the distribution of the aggregate adjusted loss. This table also contains estimates of all the predefined percentiles in addition to the percentiles that you request in the OUTSUM statement.



- If you specify `PRINT=PERTURBSUMMARY`, two tables are displayed for the distribution of the aggregate loss. The “Perturbed Summary Statistics” table (ODS name: `PerturbedSummary`) displays the summary of the effect of perturbing model parameters on the following five summary statistics of the distribution: mean, standard deviation, variance, skewness, and kurtosis. The “Perturbed Percentiles” table (ODS name: `PerturbedPctlSummary`) displays the perturbation summary for all the predefined percentiles in addition to the percentiles that you request in the `OUTSUM` statement.

The tables are displayed only if you specify a value greater than 0 for the `NPERTURBEDSAMPLES=` option.

If you specify a value of  $P$  for the `NPERTURBEDSAMPLES=` option, then for each summary statistic and percentile, an average and standard error of the set of  $P$  values of that summary statistic or percentile are displayed in the respective perturbation summary tables.

If you specify the programming statements and the `ADJUSTEDSEVERITY=` symbol, then additional perturbation summary tables are displayed for the distribution of the aggregate adjusted loss.

- If you specify `PRINT=SUMMARYSTATISTICS`, the “Summary Statistics” table (ODS name: `SummaryStatistics`) is displayed for the distribution of the aggregate loss. The table contains estimates of the following summary statistics: the number of observations in the sample, maximum value in the sample, minimum value in the sample, mean, median, standard deviation, interquartile range, variance, skewness, and kurtosis.

If you specify the programming statements and the `ADJUSTEDSEVERITY=` symbol, then an additional table of summary statistics is displayed for the distribution of the aggregate adjusted loss.

## Performance Information

The “Performance Information” table (ODS name: `PerformanceInfo`) is produced by default. It displays information about the execution mode. For single-machine mode, the table displays the number of threads that are used. For distributed mode, the table displays the grid mode (symmetric or asymmetric), the number of compute nodes, and the number of threads per node.

If you specify the `DETAILS` option in the `PERFORMANCE` statement, `PROC HPCDM` also produces a “Timing” table (ODS name: `Timing`) that displays elapsed times (absolute and relative) for the main tasks of the procedure.

---

## ODS Graphics

Statistical procedures use ODS Graphics to create graphs as part of their output. ODS Graphics is described in detail in Chapter 21, “Statistical Graphics Using ODS” (*SAS/STAT User’s Guide*).

Before you create graphs, ODS Graphics must be enabled (for example, with the `ODS GRAPHICS ON` statement). For more information about enabling and disabling ODS Graphics, see the section “Enabling and Disabling ODS Graphics” in that chapter.

The overall appearance of graphs is controlled by ODS styles. Styles and other aspects of using ODS Graphics are discussed in the section “A Primer on ODS Statistical Graphics” in that chapter.

This section describes the use of ODS for creating graphics with the `HPCDM` procedure.



**NOTE:** If you request simulation of an aggregate loss sample of large size, either by specifying a large value for the NREPLICATES= option or by including a large number of replicates in the DATA= data set that you specify in conjunction with the EXTERNALCOUNTS statement, then it is recommended that you not request any plots, because creating plots that have large numbers of points can require a very large amount of hardware resources and can take a very long time. You can disable the generation of plots either by submitting the ODS GRAPHICS OFF statement before submitting the PROC HPCDM step or by specifying the PLOTS=NONE option in the PROC HPCDM statement. It is recommended that you request plots only when the sample size is less than 100,000.

## ODS Graph Names

PROC HPCDM assigns a name to each graph that it creates by using ODS. You can use these names to selectively refer to the graphs. The names are listed in [Table 4.4](#).

**Table 4.4** ODS Graphics Produced by PROC HPCDM

ODS Graph Name	Plot Description	PLOTS= Option
ConditionalDensityPlot	Conditional density plot	CONDITIONALDENSITY
DensityPlot	Probability density function plot	DENSITY
EDFPlot	Empirical distribution function plot	EDF

## Conditional Density Plot

The conditional density plot helps you visually analyze two or three regions of the compound distribution by displaying a density function estimate that is conditional on the values of the aggregate loss that fall in those regions. You can specify the region boundaries in terms of quantiles by using the [LEFTQ=](#) and [RIGHTQ=](#) suboptions of the PLOTS=CONDITIONALDENSITY option. This is especially useful if you want to see the distribution of aggregate loss values in the right- and left-tail regions.

If you specify the programming statements and the ADJUSTEDSEVERITY= symbol, then a separate set of conditional density plots are displayed for the aggregate adjusted loss.

## Probability Density Function Plot

The probability density function (PDF) plot shows the nonparametric estimates of the PDF of the aggregate loss distribution. This plot includes histogram and kernel density estimates.

If you specify the programming statements and the ADJUSTEDSEVERITY= symbol, then a separate density plot is displayed for the aggregate adjusted loss.

## Empirical Distribution Function Plot

The empirical density function (EDF) plot shows the nonparametric estimate of the cumulative distribution function of the aggregate loss distribution. You can specify the [ALPHA=](#) suboption of the PLOTS=EDF option to request that the upper and lower confidence limits be plotted for each EDF estimate. By default, the confidence interval is not plotted.

If you specify the programming statements and the ADJUSTEDSEVERITY= symbol, then a separate EDF plot is displayed for the aggregate adjusted loss.

## Examples: HPCDM Procedure

### Example 4.1: Estimating the Probability Distribution of Insurance Payments

The primary outcome of running PROC HPCDM is the estimate of the compound distribution of aggregate loss, given the distributions of frequency and severity of the individual losses. This aggregate loss is often referred to as the ground-up loss. If you are an insurance company or a bank, you are also interested in acting on the ground-up loss by computing an entity that is derived from the ground-up loss. For example, you might want to estimate the distribution of the amount that you are expected to pay for the losses or the distribution of the amount that you can offload onto another organization, such as a reinsurance company. PROC HPCDM enables you to specify a severity adjustment program, which is a sequence of SAS programming statements that adjust the severity of the individual loss event to compute the entity of interest. Your severity adjustment program can use external information that is recorded as variables in the observations of the DATA= data set in addition to placeholder symbols for information that PROC HPCDM generates internally, such as the severity of the current loss event (`_SEV_`) and the sum of the adjusted severity values of the events that have been simulated thus far for the current sample point (`_CUMADJSEV_`). If you are doing a scenario analysis such that a scenario contains more than one observation, then you can also access the cumulative severity and cumulative adjusted severity for the current observation by using the `_CUMSEVFOROBS_` and `_CUMADJSEVFOROBS_` symbols.

This example continues the example of the section “[Scenario Analysis](#)” on page 52 to illustrate how you can estimate the distribution of the aggregate amount that is paid to a group of policyholders. Let the amount that is paid to an individual policyholder be computed by using what is usually referred to as a *disappearing deductible* (Klugman, Panjer, and Willmot 1998, Ch. 2). If  $X$  denotes the ground-up loss that a policyholder incurs,  $d$  denotes the lower limit on the deductible,  $d'$  denotes the upper limit on the deductible, and  $u$  denotes the limit on the total payments that are made to a policyholder in a year, then  $Y$ , the amount that is paid to the policyholder for each loss event, is defined as follows:

$$Y = \begin{cases} 0 & X \leq d \\ d' \frac{X-d}{d'-d} & d < X \leq d' \\ X & d' < X \leq u \\ u & X > u \end{cases}$$

You can encode this logic by using a set of SAS programming statements.

Extend the `Work.GroupOfPolicies` data set in the example in the section “[Scenario Analysis](#)” on page 52 to include the following three additional variables for each policyholder: `LowDeductible` to record  $d$ , `HighDeductible` to record  $d'$ , and `Limit` to record  $u$ . The data set contains the observations as shown in [Output 4.1.1](#).

**Output 4.1.1** Scenario Analysis Data for Multiple Policyholders with Policy Provisions

P o l i c y h o l d e r I d	a g e	g e n d e r	c a t e g o r y	a n n u a l M i l l i o n s	e x p o s e d u n i t s	c l a s s i f i c a t i o n	i n c u r r e n t e x p o s e	l o s s c o s t l i m i t	h i g h D e d u c t i b l e	l o w D e d u c t i b l e	a m o u n t p a i d
1	1.18	2	1	2.2948	3	0.99532	1.59870	400	1400	7500	10000
2	0.66	2	2	2.8148	1	0.05625	0.67539	300	1300	2500	20000
3	0.82	1	2	1.6130	2	0.84146	1.05940	100	1100	5000	10000
4	0.44	1	1	1.2280	3	0.14324	0.24110	300	800	5000	20000
5	0.44	1	1	0.9670	2	0.08656	0.65979	100	1100	5000	20000

The following PROC HPCDM step estimates the compound distributions of the aggregate loss and the aggregate amount that is paid to the group of policyholders in the Work.GroupOfPolicies data set by using the count model that is stored in the Work.CountregModel item store and the lognormal severity model that is stored in the Work.SevRegEst data set:

```

/* Simulate the aggregate loss distribution and aggregate adjusted
   loss distribution for the scenario with multiple policyholders */
proc hpcdm data=groupOfPolicies nreplicates=10000 seed=13579 print=all
    countstore=work.countregmodel severityest=work.sevregest
    plots=(edf pdf) nperturbedSamples=50
    adjustedseverity=amountPaid;
severitymodel logn;

if (_sev_ <= lowDeductible) then
    amountPaid = 0;
else do;
    if (_sev_ <= highDeductible) then
        amountPaid = highDeductible *
            (_sev_-lowDeductible)/(highDeductible-lowDeductible);
    else
        amountPaid = MIN(_sev_, limit); /* imposes per-loss payment limit */
    end;
run;

```

The preceding step uses a severity adjustment program to compute the value of the symbol AmountPaid and specifies that symbol in the ADJUSTEDSEVERITY= option in the PROC HPCDM step. The program is executed for each simulated loss event. The PROC HPCDM supplies your program with the value of the severity in the \_SEV\_ placeholder symbol.

The “Sample Summary Statistics” table in [Output 4.1.2](#) shows the summary statistics of the compound distribution of the aggregate ground-up loss. The “Adjusted Sample Summary Statistics” table shows the summary statistics of the compound distribution of the aggregate AmountPaid. The average aggregate payment is about 4,321, as compared to the average aggregate ground-up loss of 5,883.

**Output 4.1.2** Summary Statistics of Compound Distributions of the Total Loss and Total Amount Paid

The HPCDM Procedure			
Severity Model: Logn			
Count Model: NegBin(p=2)			
Compound Distribution Information			
Severity Model	Lognormal Distribution		
Scale Model Regressors	carType carSafety income		
Count Model	NegBin(p=2) Model in Item Store WORK.COUNTREGMODEL		
Sample Summary Statistics			
Mean	5883.4	Median	3281.0
Standard Deviation	7576.5	Interquartile Range	6878.1
Variance	57402878.5	Minimum	0
Skewness	3.05639	Maximum	107486.7
Kurtosis	15.80583	Sample Size	10000
Adjusted Sample Summary Statistics			
Mean	4320.8	Median	2547.9
Standard Deviation	5321.9	Interquartile Range	5461.3
Variance	28322194.1	Minimum	0
Skewness	2.55004	Maximum	64697.4
Kurtosis	11.26066	Sample Size	10000

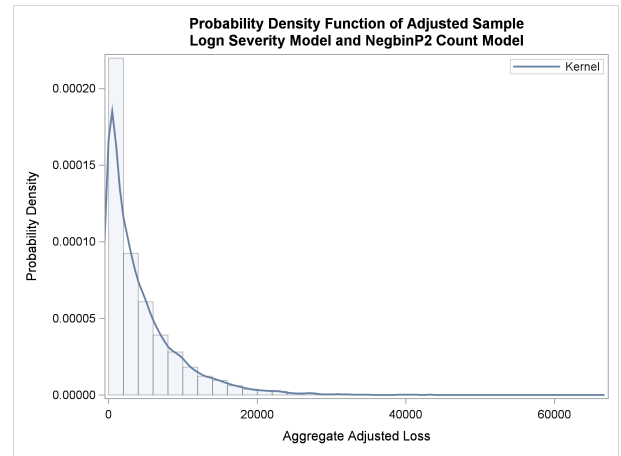
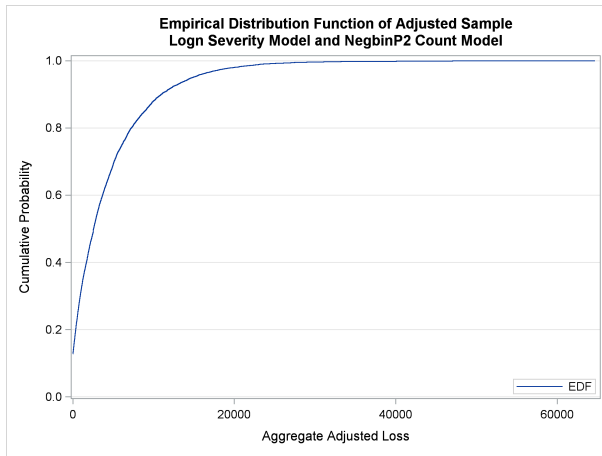
The perturbation summary of the distribution of AmountPaid is shown in [Output 4.1.3](#). It shows that you can expect to pay a median of  $2,568 \pm 270$  to this group of five policyholders in a year. Also, if the 99.5th percentile defines the worst case, then you can expect to pay  $29,098 \pm 2,881$  in the worst-case.

**Output 4.1.3** Perturbation Summary of the Total Amount Paid

Adjusted Sample Percentile Perturbation Analysis		
Percentile	Estimate	Standard Error
1	0	0
5	0	0
25	604.60948	126.17884
50	2568.4	270.29467
75	6126.0	579.35289
95	15245.7	1455.8
99	24920.6	2434.2
99.5	29097.7	2880.8
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		

The empirical distribution function (EDF) and probability density function plots of the aggregate adjusted loss are shown in [Output 4.1.4](#). Both plots indicate a heavy-tailed distribution of the total amount paid.

**Output 4.1.4** PDF and EDF Plots of the Compound Distribution of the Total Amount Paid



Now consider that, in the future, you want to modify the policy provisions to add a limit on the total amount of payment that is made to an individual policyholder in one year and to impose a group limit of 50,000 on the total amount of payments that are made to the group as a whole in one year. You can analyze the effects of these modified policy provisions on the distribution of the aggregate paid amount by recording the individual policyholder's annual limit in the AnnualLimit variable of the input data set and then modifying your severity adjustment program by using the placeholder symbols `_CUMADJSEVFOROBS_` and `_CUMADJSEV_` as shown in the following PROC HPCDM step:

```
/* Simulate the aggregate loss distribution and aggregate adjusted
   loss distribution for the modified set of policy provisions */
proc hpcdm data=groupOfPolicies nreplicates=10000 seed=13579 print=all
    countstore=work.countregmodel severityest=work.sevregist
    plots=none nperturbedSamples=50
    adjustedseverity=amountPaid;
severitymodel logn;

if (_sev_ <= lowDeductible) then
    amountPaid = 0;
else do;
    if (_sev_ <= highDeductible) then
        amountPaid = highDeductible *
            (_sev_-lowDeductible)/(highDeductible-lowDeductible);
    else
        amountPaid = MIN(_sev_, limit); /* imposes per-loss payment limit */

/* impose policyholder's annual limit */
amountPaid = MIN(amountPaid, MAX(0,annualLimit - _cumadjsevforobs_));

/* impose group's annual limit */
amountPaid = MIN(amountPaid, MAX(0,50000 - _cumadjsev_));
end;
run;
```

The results of the perturbation analysis for these modified policy provisions are shown in [Output 4.1.5](#). When compared to the results of [Output 4.1.3](#), the additional policy provisions of restricting the total payment to the policyholder and the group have kept the median payment unchanged, but the provisions have reduced the worst-case payment (99.5th percentile) to  $22,881 \pm 1,119$  from  $29,098 \pm 2,881$ .

**Output 4.1.5** Perturbation Summary of the Total Amount Paid for Modified Policy Provisions

The HPCDM Procedure		
Severity Model: Logn		
Count Model: NegBin(p=2)		
Adjusted Sample Percentile		
Perturbation Analysis		
Percentile	Estimate	Standard Error
0	0	0
1	0	0
5	0	0
25	604.60948	126.17884
50	2568.4	270.29467
75	6126.0	579.35289
95	13865.7	1135.7
99	20869.0	896.03211
99.5	22880.5	1118.9
Number of Perturbed Samples = 50		
Size of Each Sample = 10000		

**Example 4.2: Using Externally Simulated Count Data**

The COUNTREG procedure enables you to estimate count regression models that are based on the most commonly used discrete distributions, such as the Poisson, negative binomial (both  $p = 1$  and  $p = 2$ ), and Conway-Maxwell-Poisson distributions. PROC COUNTREG also enables you to fit zero-inflated models that are based on Poisson, negative binomial ( $p=2$ ), and Conway-Maxwell-Poisson distributions. However, there might be situations in which you want to use some other method of fitting count regression models. For example, if you are modeling the number of loss events that are incurred by two financial instruments such that there is some dependency between the two, then you might use some multivariate frequency modeling methods and simulate the counts for each instrument by using the dependency structure between the count model parameters of the two instruments. As another example, you might want to use different types of count models for different BY groups in your data; this is not possible in PROC COUNTREG in SAS/ETS 13.1 and earlier. So you need to simulate the counts for such BY groups externally. PROC HPCDM enables you to supply externally simulated counts by using the EXTERNALCOUNTS statement. PROC HPCDM then does not need to simulate the counts internally; it simulates only the severity of each loss event by using the severity model estimates in the SEVERITYEST= data set. The process is described and illustrated in the section “[Simulation with External Counts](#)” on page 78.

Consider that you are a bank, and as part of quantifying your operational risk, you want to estimate the aggregate loss distributions for two lines of business, retail banking and commercial banking, by using some key risk indicators (KRIs). Assume that your model fitting and model selection process has determined that the Poisson regression model and negative binomial regression model are the best-fitting count models for number of loss events that are incurred in the retail banking and commercial banking businesses, respectively. Let CorpKRI1, CorpKRI2, CbKRI1, CbKRI2, and CbKRI3 be the KRIs that are used in the count regression model of the commercial banking business, and let CorpKRI1, RbKRI1, and RbKRI2 be the KRIs that are used in the count regression model of the retail banking business. Some examples of corporate-level KRIs (CorpKRI1 and CorpKRI2 in this example) are the ratio of temporary to permanent employees and the

number of security breaches that are reported during a year. Some examples of KRIs that are specific to the commercial banking business (CbKRI1, CbKRI2, and CbKRI3 in this example) are number of credit defaults, proportion of financed assets that are movable, and penalty claims against your bank because of processing delays. Some examples of KRIs that are specific to the retail banking business (RbKRI1 and RbKRI2 in this example) are number of credit cards that are reported stolen, fraction of employees who have not undergone fraud detection training, and number of forged drafts and checks that are presented in a year.

Let the severity of each loss event in the commercial banking business be dependent on two KRIs, CorpKRI1 and CbKRI2. Let the severity of each loss event in the retail banking business be dependent on three KRIs, CorpKRI2, RbKRI1, and RbKRI3. Note that for each line of business, the set of KRIs that are used for the severity model is different from the set of KRIs that are used for the count model, although there is some overlap between the two sets. Further, the severity model for retail banking includes a new regressor (RbKRI3) that is not used for any of the count models. Such use of different sets of KRIs for count and severity models is typical of real-world applications.

Let the parameter estimates of the negative binomial and Poisson regression models, as determined by PROC COUNTREG, be available in the Work.CountEstEx2NB2 and Work.CountEstEx2Poisson data sets, respectively. These data sets are produced by using the OUTEST= option in the respective PROC COUNTREG statements. Let the parameter estimates of the best-fitting severity models, as determined by PROC SEVERITY, be available in the Work.SevEstEx2Best data set. You can find the code to prepare these data sets in the PROC HPCDM sample program *hcdmex02.sas*.

Now, consider that you want to estimate the distribution of the aggregate loss for a scenario, which is represented by a specific set of KRI values. The following DATA step illustrates one such scenario:

```
/* Generate a scenario data set for a single operating condition */
data singleScenario (keep=corpKRI1 corpKRI2 cbKRI1 cbKRI2 cbKRI3
                    rbKRI1 rbKRI2 rbKRI3);
  array x{8} corpKRI1 corpKRI2 cbKRI1 cbKRI2 cbKRI3 rbKRI1 rbKRI2 rbKRI3;
  call streaminit(5151);
  do i=1 to dim(x);
    x(i) = rand('NORMAL');
  end;
  output;
run;
```

The Work.SingleScenario data set contains all the KRIs that are included in the count and severity models of both business lines. Note that if you standardize or scale the KRIs while fitting the count and severity models, then you must apply the same standardization or scaling method to the values of the KRIs that you specify in the scenario. In this particular example, all KRIs are assumed to be standardized.

The following DATA step uses the scenario in the Work.SingleScenario data set to simulate 10,000 replications of the number of loss events that you might observe for each business line and writes the simulated counts to the NumLoss variable of the Work.LossCounts1 data set:

```
/* Simulate multiple replications of the number of loss events that
   you can expect in the scenario being analyzed */
data lossCounts1 (keep=line corpKRI1 corpKRI2 cbKRI2 rbKRI1 rbKRI3 numloss);
  array cxR{3} corpKRI1 rbKRI1 rbKRI2;
  array cbetaR{4} _TEMPORARY_;
  array cxC{5} corpKRI1 corpKRI2 cbKRI1 cbKRI2 cbKRI3;
  array cbetaC{6} _TEMPORARY_;

  retain theta;
```



```

if _n_ = 1 then do;
  call streaminit(5151);
  * read count model estimates *;
  set countEstEx2NB2(where=(line='CommercialBanking' and _type_='PARM'));
  cbetaC(1) = Intercept;
  do i=1 to dim(cxC);
    cbetaC(i+1) = cxC(i);
  end;
  alpha = _Alpha;
  theta = 1/alpha;

  set countEstEx2Poisson(where=(line='RetailBanking' and _type_='PARM'));
  cbetaR(1) = Intercept;
  do i=1 to dim(cxR);
    cbetaR(i+1) = cxR(i);
  end;
end;

set singleScenario;
do iline=1 to 2;
  if (iline=1) then line = 'CommercialBanking';
  else line = 'RetailBanking';
  do repid=1 to 10000;
    nnz = 1;
    maxtries = 5*nnz;
    nc = 0;
    ntries = 0;
    do while (nc < nnz and ntries < maxtries);
      * draw from count distribution *;
      if (iline=1) then do;
        xbeta = cbetaC(1);
        do i=1 to dim(cxC);
          xbeta = xbeta + cxC(i) * cbetaC(i+1);
        end;
        Mu = exp(xbeta);
        p = theta/(Mu+theta);
        numloss = rand('NEGB',p,theta);
      end;
      else do;
        xbeta = cbetaR(1);
        do i=1 to dim(cxR);
          xbeta = xbeta + cxR(i) * cbetaR(i+1);
        end;
        numloss = rand('POISSON', exp(xbeta));
      end;
      if (numloss > 0) then do;
        output;
        nc = nc + 1;
      end;
      ntries = ntries + 1;
    end;
  end;
end;
run;

```

The Work.LossCounts1 data set contains the NumLoss variable in addition to the KRIs that are used by the severity regression model, which are needed by PROC HPCDM to simulate the aggregate loss.

By default, PROC HPCDM computes an aggregate loss distribution by using each of the severity models that you specify in the SEVERITYMODEL statement. However, you can restrict PROC HPCDM to use only a subset of the severity models for a given BY group by modifying the SEVERITYEST= data set to include only the estimates of the desired severity models in each BY group, as illustrated in the following DATA step:

```
/* Keep only the best severity model for each business line
   and set coefficients of unused regressors in each model to 0 */
data sevestEx2Best;
  set sevestEx2;
  if ((line = 'CommercialBanking' and _model_ = 'Logn')) then do;
    corpKRI2 = 0; rbKRI1 = 0; rbKRI3 = 0;
    output;
  end;
  else if ((line = 'RetailBanking' and _model_ = 'Gamma')) then do;
    corpKRI1 = 0; cbKRI2 = 0;
    output;
  end;
run;
```

Note that the preceding DATA step also sets the coefficients of the unused regressors in each model to 0. This is important because PROC HPCDM uses all the regressors that it detects from the SEVERITYEST= data set for each severity model.

Now, you are ready to estimate the aggregate loss distribution for each line of business by submitting the following PROC HPCDM step, in which you specify the EXTERNALCOUNTS statement to request that external counts in the NumLoss variable of the DATA= data set be used for simulation of the aggregate loss:

```
/* Estimate the distribution of the aggregate loss for both
   lines of business by using the externally simulated counts */
proc hpcdm data=lossCounts1 seed=13579 print=all
  severityest=sevestEx2Best;
  by line;
  externalcounts count=numloss;
  severitymodel logn gamma;
run;
```

Each observation in the Work.LossCounts1 data set represents one replication of the external counts simulation process. For each such replication, the preceding PROC HPCDM step makes as many severity draws from the severity distribution as the value of the NumLoss variable and adds the severity values from those draws to compute one sample point of the aggregate loss. The severity distribution that is used for making the severity draws has a scale parameter value that is decided by the KRI values in the given observation and the regression parameter values that are read from the Work.SevEstEx2Best data set.

The summary statistics and percentiles of the aggregate loss distribution for the commercial banking business, which uses the lognormal severity model, are shown in [Output 4.2.1](#). The “Input Data Summary” table indicates that each of the 9,954 observations in the BY group is treated as one replication and that there are a total of 19,241 loss events produced by all the replications together. For the scenario in the Work.SingleScenario data set, you can expect the commercial banking business to incur an average aggregate loss of 651 units, as shown in the “Sample Summary Statistics” table, and the chance that the loss will exceed 4,337 units is 0.5%, as shown in the “Sample Percentiles” table.

**Output 4.2.1** Aggregate Loss Summary for Commercial Banking Business

```

----- line=CommercialBanking -----

      The HPCDM Procedure

      Input Data Summary

      Name                      WORK.LOSSCOUNTS1
      Observations                9954
      Valid Observations          9954
      Replications                 9954
      Total Count                  19241

----- line=CommercialBanking -----

      Sample Summary Statistics

      Mean          651.00065      Median          418.36937
      Standard Deviation 726.61443      Interquartile Range 653.67139
      Variance        527968.5      Minimum          8.00493
      Skewness        3.15153      Maximum         12726.4
      Kurtosis        19.08843      Sample Size       9954

----- line=CommercialBanking -----

      Sample Percentiles

      Percentile      Value
      0                8.00493
      1               29.52879
      5               59.63848
      25             188.00888
      50             418.36937
      75             841.68028
      95             2037.3
      99             3472.2
      99.5           4337.0

      Percentile Method = 5

```

For the retail banking business, which uses the gamma severity model, the “Sample Percentiles” table in [Output 4.2.2](#) indicates that the median operational loss of that business is about 85 units and the chance that the loss will exceed 344 units is about 1%.

**Output 4.2.2** Aggregate Loss Percentiles for Retail Banking Business

----- line=RetailBanking -----	
Sample Percentiles	
Percentile	Value
0	1.19575
1	9.88436
5	19.76335
25	48.97570
50	84.78094
75	141.38838
95	250.92488
99	343.85721
99.5	381.70522
Percentile Method = 5	

When you conduct the simulation and estimation for a scenario that contains only one observation, you assume that the operating environment does not change over the period of time that is being analyzed. That assumption might be valid for shorter durations and stable business environments, but often the operating environments change, especially if you are estimating the aggregate loss over a longer period of time. So you might want to include in your scenario all the possible operating environments that you expect to see during the analysis time period. Each environment is characterized by its own set of KRI values. For example, the operating conditions might change from quarter to quarter, and you might want to estimate the aggregate loss distribution for the entire year. You start the estimation process for such scenarios by creating a scenario data set. The following DATA step creates the Work.MultiConditionScenario data set, which consists of four operating environments, one for each quarter:

```

/* Generate a scenario data set for multiple operating conditions */
data multiConditionScenario (keep=opEnvId corpKRI1 corpKRI2
    cbKRI1 cbKRI2 cbKRI3 rbKRI1 rbKRI2 rbKRI3);
    array x{8} corpKRI1 corpKRI2 cbKRI1 cbKRI2 cbKRI3 rbKRI1 rbKRI2 rbKRI3;
    call streaminit(5151);
    do opEnvId=1 to 4;
        do i=1 to dim(x);
            x(i) = rand('NORMAL');
        end;
        output;
    end;
run;

```

All four observations of the Work.MultiConditionScenario data set together form one scenario. When simulating the external counts for such multi-entity scenarios, one replication consists of the possible number of loss events that can occur as a result of each of the four operating environments. In any given replication, some operating environments might not produce any loss event or all four operating environments might produce some loss events. Assume that you use a DATA step to create the Work.LossCounts2 data set that contains, for each business line, 10,000 replications of the loss counts and that you identify each replication by using the Repld variable. You can find the DATA step code to prepare the Work.LossCounts2 data set in the PROC HPCDM sample program *hcdmex02.sas*.

Output 4.2.3 shows some observations of the Work.LossCounts2 data set for each business line. For the first replication (Repld=1) of the commercial banking business, only operating environment 3 incurs two loss events, whereas the other environments incur no loss events. For the second replication (Repld=2), all operating environments incur at least one loss event. For the first replication (Repld=1) of the retail banking business, operating environments 2, 3, and 4 incur four, one, and four loss events, respectively.

**Output 4.2.3** Snapshot of the External Counts Data with Replication Identifier

		c	c				n
	o	o	o				u
	p	r	r	c	r	r	r
	E	p	p	b	b	b	m
l	n	K	K	K	K	K	e
i	v	R	R	R	R	R	p
n	I	I	I	I	I	I	i
e	d	1	2	2	1	3	s
CommercialBanking	3	-0.29120	-0.45239	0.98855	-0.37208	-1.51534	1 2
CommercialBanking	1	0.45224	0.40661	-0.33680	-1.08692	-2.20557	2 1
CommercialBanking	2	-0.03799	0.98670	-0.03752	1.94589	1.22456	2 3
CommercialBanking	3	-0.29120	-0.45239	0.98855	-0.37208	-1.51534	2 9
CommercialBanking	4	0.87499	-0.67812	-0.04839	-1.44881	0.78221	2 8
CommercialBanking	1	0.45224	0.40661	-0.33680	-1.08692	-2.20557	3 5
CommercialBanking	2	-0.03799	0.98670	-0.03752	1.94589	1.22456	3 1
CommercialBanking	3	-0.29120	-0.45239	0.98855	-0.37208	-1.51534	3 2
RetailBanking	2	-0.03799	0.98670	-0.03752	1.94589	1.22456	1 4
RetailBanking	3	-0.29120	-0.45239	0.98855	-0.37208	-1.51534	1 1
RetailBanking	4	0.87499	-0.67812	-0.04839	-1.44881	0.78221	1 4
RetailBanking	1	0.45224	0.40661	-0.33680	-1.08692	-2.20557	2 2
RetailBanking	2	-0.03799	0.98670	-0.03752	1.94589	1.22456	2 5
RetailBanking	4	0.87499	-0.67812	-0.04839	-1.44881	0.78221	2 3
RetailBanking	1	0.45224	0.40661	-0.33680	-1.08692	-2.20557	3 2
RetailBanking	2	-0.03799	0.98670	-0.03752	1.94589	1.22456	3 3

You can now use this simulated count data to estimate the distribution of the aggregate loss that is incurred in all four operating environments by submitting the following PROC HPCDM step, in which you specify the replication identifier variable Repld in the ID= option of the EXTERNALCOUNTS statement:

```

/* Estimate the distribution of the aggregate loss for both
   lines of business by using the externally simulated counts
   for the multiple operating environments */
proc hpcdm data=lossCounts2 seed=13579 print=all
    severityest=sevestEx2Best plots=density;
    by line;
    distby repid;
    externalcounts count=numloss id=repid;
    severitymodel logn gamma;
run;

```

Note that when you specify the ID= variable in the EXTERNALCOUNTS statement, you must also specify that variable in the DISTBY statement. Within each BY group, for each value of the Repld variable, one point of the aggregate loss sample is simulated by using the process that is described in the section “[Simulation with External Counts](#)” on page 78.

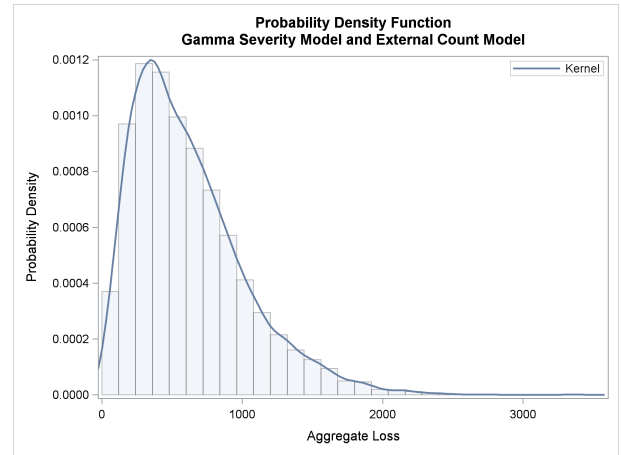
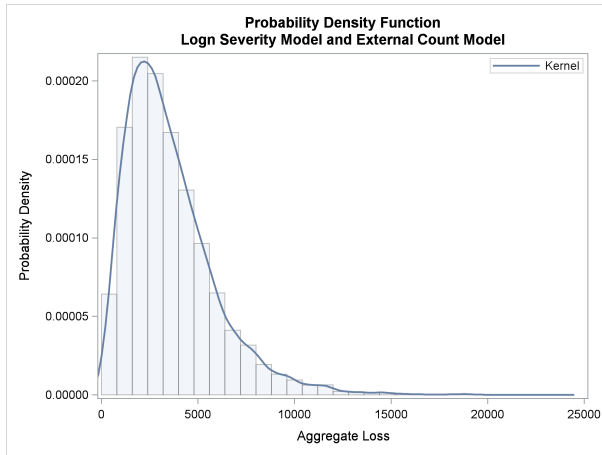
The summary statistics and percentiles of the distribution of the aggregate loss, which is the aggregate of the losses across all four operating environments, are shown in [Output 4.2.4](#) for the commercial banking business. The “Input Data Summary” table indicates that there are 10,000 replications in the BY group and that a total of 98,480 loss events are generated across all replications. The “Sample Percentiles” table indicates that you can expect a median aggregate loss of 3,075 units and a worst-case loss, as defined by the 99.5th percentile, of 13,150 units from the commercial banking business when you combine losses that result from all four operating environments.

**Output 4.2.4** Aggregate Loss Summary for the Commercial Banking Business in Multiple Operating Environments

----- line=CommercialBanking -----	
The HPCDM Procedure	
Input Data Summary	
Name	WORK.LOSSCOUNTS2
Observations	32526
Valid Observations	32526
Replications	10000
Total Count	98480
----- line=CommercialBanking -----	
Sample Percentiles	
Percentile	Value
1	342.53328
5	792.48797
25	1876.8
50	3075.2
75	4694.6
95	8058.6
99	11575.4
99.5	13149.7
Percentile Method = 5	

The probability density functions of the aggregate loss for the commercial and retail banking businesses are shown in [Output 4.2.5](#). In addition to the difference in scales of the losses in the two businesses, you can see that the aggregate loss that is incurred in the commercial banking business has a heavier right tail than the aggregate loss that is incurred in the retail banking business.

**Output 4.2.5** Density Plots of the Aggregate Losses for Commercial Banking (left) and Retail Banking (right) Businesses



## References

- Fisher, R. A. (1973), *Statistical Methods for Research Workers*, 14th Edition, New York: Hafner Publishing.
- Klugman, S. A., Panjer, H. H., and Willmot, G. E. (1998), *Loss Models: From Data to Decisions*, New York: John Wiley & Sons.





## Chapter 5

# The HPCOPULA Procedure (Experimental)

### Contents

---

Overview: HPCOPULA Procedure . . . . .	<b>113</b>
PROC HPCOPULA Features . . . . .	114
Getting Started: HPCOPULA Procedure . . . . .	<b>114</b>
Syntax: HPCOPULA Procedure . . . . .	<b>115</b>
Functional Summary . . . . .	115
PROC HPCOPULA Statement . . . . .	115
DEFINE Statement . . . . .	115
SIMULATE Statement . . . . .	116
PERFORMANCE Statement . . . . .	117
VAR Statement . . . . .	117
Details: HPCOPULA Procedure . . . . .	<b>118</b>
Sklar's Theorem . . . . .	118
Dependence Measures . . . . .	118
Normal Copula . . . . .	119
Student's <i>t</i> copula . . . . .	120
Archimedean Copulas . . . . .	120
OUTUNIFORM= Data Sets . . . . .	122
Examples: HPCOPULA Procedure . . . . .	<b>122</b>
Example 5.1: Simulating Default Times . . . . .	122
References . . . . .	<b>126</b>

---

---

## Overview: HPCOPULA Procedure

The HPCOPULA procedure is a high-performance version of the SAS/ETS COPULA procedure, which simulates data from a specified copula. Unlike the COPULA procedure, which can be run only on an individual workstation, the HPCOPULA procedure takes advantage of a computing environment in which the optimization task can be distributed to one or more nodes. In addition, each node can use one or more threads to perform the optimization on its subset of the data. When several nodes are used and each node uses several threads to carry out its part of the work, the result is a highly parallel computation that provides a dramatic gain in performance.

You can use the HPCOPULA procedure to read and write data in distributed form and perform analyses either in single-machine mode or in distributed mode. For more information about the execution mode of SAS High-Performance Analytics procedures, see the section “[Processing Modes](#)” on page 12 in Chapter 3, “[Shared Concepts and Topics](#).”

The HPCOPULA procedure is specifically designed to operate in the high-performance distributed environment. By default, PROC HPCOPULA performs computations in multiple threads.

---

## PROC HPCOPULA Features

The HPCOPULA procedure enables you to simulate a specified copula, and it supports the following types of copulas:

- normal copula
- $t$  copula
- Archimedean copulas:
  - Clayton copula
  - Frank copula
  - Gumbel copula

---

## Getting Started: HPCOPULA Procedure

This example illustrates the use of PROC HPCOPULA. The data are daily returns on several major stocks. The main purpose of this example is to simulate from the joint distribution of stock returns a new sample of a specified size, provided that the parameter estimates of the copula model that is used are available.

In the following statements, the DEFINE statement specifies a normal copula named COP, and the CORR= option specifies that the data set Estimates be used as the source for the model parameters. The NDRAWS=1000000 option in the SIMULATE statement generates one million observations from the normal copula. The OUTUNIFORM= option specifies the name of the SAS data set to contain the simulated sample that has uniform marginal distributions. The PERFORMANCE statement requests that the analytic computations use two nodes in the distributed computing environment and two threads in each node. Note that this syntax does not require the DATA= option.

```
/* Copula simulation of uniforms */
proc hpcopula;
  var ret_ibm ret_msft ret_bp ret_ko ret_duk;
  define cop normal (corr = estimates);
  simulate cop / ndraws      = 1000000
                outuniform = simulated_uniforms;
  PERFORMANCE nodes=2 nthreads=2 details;
run;
```

The simulated data are contained in the new SAS data set, Simulated\_Uniforms.

## Syntax: HPCOPULA Procedure

The following statements are available in the HPCOPULA procedure:

```
PROC HPCOPULA options ;
  VAR variables ;
  DEFINE name copula-type < ( parameter-value-options ... ) > ;
  SIMULATE < copula-name-list > / options ;
```

## Functional Summary

Table 5.1 summarizes the statements and options that the HPCOPULA procedure uses.

**Table 5.1** PROC HPCOPULA Functional Summary

Description	Statement	Option
<b>Data Set Options</b>		
Specifies the input data set that contains the correlation matrix for elliptical copulas	DEFINE	CORR=
<b>Declaring the Role of Variables</b>		
Specifies the names of the variables to use in copula fitting or in simulation	VAR	
<b>Copula Simulation Options</b>		
Specifies the random sample size	SIMULATE	NDRAWS=
Specifies the random number generator seed	SIMULATE	SEED=
<b>Output Control Options</b>		
Specifies the output data set to contain the random samples from the simulation with uniform marginal distribution	SIMULATE	OUTUNIFORM=

## PROC HPCOPULA Statement

```
PROC HPCOPULA ;
```

The PROC HPCOPULA statement invokes the HPCOPULA procedure.

## DEFINE Statement

```
DEFINE name copula-type < ( parameter-value-options ... ) > ;
```

The DEFINE statement specifies the relevant information about the copula that is used for the simulation. You can specify the following arguments:

<i>name</i>	specifies the name of the copula definition. You can use this <i>name</i> later in the SIMULATE statement.
<i>copula-type</i>	specifies the type of copula. You must specify one of the following copula types, which are described in the section “ <a href="#">Details: HPCOPULA Procedure</a> ” on page 118:
<b>NORMAL</b>	fits the normal copula.
<b>T</b>	fits the $t$ copula.
<b>CLAYTON</b>	fits the Clayton copula.
<b>FRANK</b>	fits the Frank copula.
<b>GUMBEL</b>	fits the Gumbel copula.

*parameter-value-options*

specify the input parameters that are used to simulate the specified copula. These options must be appropriate for the type of copula specified. You can specify the following *parameter-value-options*:

**CORR=SAS-data-set**

specifies the data set that contains the correlation matrix to use for elliptical copulas. If the correlation matrix is valid but its elements are not submitted in order, then you must provide the variable names in the first column of the matrix, and these names must match the variable names in the VAR statement. See [Output 5.1.1](#) for an example of a correlation matrix input in this form. If the correlation matrix elements are submitted in order, the first column of variable names is not required. You can use this option for normal and  $t$  copulas.

**DF=value**

specifies the degrees of freedom. You can use this option for  $t$  copulas.

**THETA=value**

specifies the parameter value for the Archimedean copulas.

The DEFINE statement is used with the SIMULATE statement.

---

## SIMULATE Statement

**SIMULATE** < *copula-name-list* > /options ;

The SIMULATE statement simulates data from a specified copula model. The copula name specification is the name of a defined copula as specified by *name* in the DEFINE statement.

**NDRAWS=integer**

specifies the number of draws to generate for this simulation. By default, NDRAWS=100.

**OUTUNIFORM=SAS-data-set**

specifies the output data set to contain the result of the simulation in uniform margins. You can use this option when MARGINALS=UNIFORM or MARGINALS=EMPIRICAL. If MARGINALS=EMPIRICAL, then this option enables you to obtain the samples that are simulated from the joint distribution specified by the copula, where all marginal distributions are uniform. The data are not created if you do not specify this option.

**SEED=integer**

specifies the seed for generating random numbers for the simulation. If you do not provide the seed, a random number is used as the seed.

---

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement specifies *performance-options* to control the multithreaded and distributed computing environment and requests detailed performance results of the HPCOPULA procedure. You can also use the PERFORMANCE statement to control whether the HPCOPULA procedure executes in SMP or MPP mode. You can specify the following *performance-options*:

**DETAILS**

requests a table that shows a timing breakdown of the PROC HPCOPULA steps.

**NODES=*n***

specifies the number of nodes in the distributed computing environment, provided that the data are not processed alongside the database.

**NTHREADS=*n***

specifies the number of threads for analytic computations and overrides the SAS system option THREADS | NOTTHREADS. If you do not specify the NTHREADS= option, PROC HPCOPULA creates one thread per CPU for the analytic computations.

For more information about the PERFORMANCE statement, see the section “[PERFORMANCE Statement](#)” on page 39 in Chapter 3, “[Shared Concepts and Topics](#).”

---

## VAR Statement

**VAR** *variables* ;

The VAR statement specifies the variable names in the input data set that is specified by the DATA= option in the PROC HPCOPULA statement. The subset of variables in the data set is used for the copula models in the FIT statement. If there is no input data set, the VAR statement creates the list of variable names for the SIMULATE statement.

---

## Details: HPCOPULA Procedure

---

### Sklar's Theorem

The copula models are tools for studying the dependence structure of multivariate distributions. The usual joint distribution function contains the information both about the marginal behavior of the individual random variables and about the dependence structure between the variables. The copula is introduced to decouple the marginal properties of the random variables and the dependence structures. An  $m$ -dimensional *copula* is a joint distribution function on  $[0, 1]^m$ , where all marginal distributions are standard uniform. The common notation for a copula is  $C(u_1, \dots, u_m)$ .

The Sklar (1959) theorem shows the importance of copulas in modeling multivariate distributions. The first part of the theorem states that a copula can be derived from any joint distribution functions, and the second part asserts the opposite: that any copula can be combined with any set of marginal distributions to result in a multivariate distribution function. The theorem follows:

- Let  $F$  be a joint distribution function, and let  $F_j, j = 1, \dots, m$ , be the marginal distributions. Then there exists a copula  $C : [0, 1]^m \rightarrow [0, 1]$  such that

$$F(x_1, \dots, x_m) = C(F_1(x_1), \dots, F_m(x_m))$$

for all  $x_1, \dots, x_m$  in  $[-\infty, \infty]$ . Moreover, if the margins are continuous, then  $C$  is unique; otherwise  $C$  is uniquely determined on  $\text{Ran} F_1 \times \dots \times \text{Ran} F_m$ , where  $\text{Ran} F_j = F_j([-\infty, \infty])$  is the range of  $F_j$ .

- The converse is also true. That is, if  $C$  is a copula and  $F_1, \dots, F_m$  are univariate distribution functions, then the multivariate function that is defined in the preceding equation is a joint distribution function with marginal distributions  $F_j, j = 1, \dots, m$ .

---

### Dependence Measures

There are three basic types of dependence measures: linear correlation, rank correlation, and tail dependence. Linear correlation is given by

$$\rho \equiv \text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X)}\sqrt{\text{var}(Y)}}$$

The linear correlation coefficient contains very limited information about the joint properties of the variables. A well-known property is that zero correlation does not imply independence, whereas independence implies zero correlation. In addition, there are distinct bivariate distributions that have the same marginal distribution and the same correlation coefficient. These results suggest that caution must be used in interpreting the linear correlation.

Another statistical measure of dependence is rank correlation, which is nonparametric. For example, Kendall's tau is the covariance between the sign statistics  $X_1 - \tilde{X}_1$  and  $X_2 - \tilde{X}_2$ , where  $(\tilde{X}_1, \tilde{X}_2)$  is an independent copy of  $(X_1, X_2)$ :

$$\rho_\tau \equiv E[\text{sign}(X_1 - \tilde{X}_1)(X_2 - \tilde{X}_2)]$$

The sign function (sometimes written as  $\text{sgn}$ ) is defined as

$$\text{sign}(x) = \begin{cases} -1 & \text{if } x \leq 0 \\ 0 & \text{if } x = 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

Spearman's rho is the correlation between the transformed random variables:

$$\rho_S(X_1, X_2) \equiv \rho(F_1(X_1), F_2(X_2))$$

The variables are transformed by their distribution functions so that the transformed variables are uniformly distributed on  $[0, 1]$ . The rank correlations depend only on the copula of the random variables and are indifferent to the marginal distributions. Like linear correlation, rank correlation has its limitations. In particular, different copulas result in the same rank correlation.

A third measure, tail dependence, focuses on only part of the joint properties between the variables. Tail dependence measures the dependence when both variables have extreme values. Formally, they can be defined as the conditional probabilities of quantile exceedances. There are two types of tail dependence:

- Upper tail dependence is defined as

$$\lambda_u(X_1, X_2) \equiv \lim_{q \rightarrow 1^-} P(X_2 > F_2^{-1}(q) | X_1 > F_1^{-1}(q))$$

when the limit exists and  $\lambda_u \in [0, 1]$ . Here  $F_j^{-1}$  is the quantile function (that is, the inverse of the CDF).

- Lower tail dependence is defined symmetrically.

## Normal Copula

Let  $u_j \sim U(0, 1)$  for  $j = 1, \dots, m$ , where  $U(0, 1)$  represents the uniform distribution on the  $[0, 1]$  interval. Let  $\Sigma$  be the correlation matrix, where  $m(m-1)/2$  parameters satisfy the positive semidefiniteness constraint. The normal copula can be written as

$$C_\Sigma(u_1, u_2, \dots, u_m) = \Phi_\Sigma(\Phi^{-1}(u_1), \dots, \Phi^{-1}(u_m))$$

where  $\Phi$  is the distribution function of a standard normal random variable and  $\Phi_\Sigma$  is the  $m$ -variate standard normal distribution with mean vector 0 and covariance matrix  $\Sigma$ . That is, the distribution  $\Phi_\Sigma$  is  $N_m(0, \Sigma)$ .

## Simulation

For the normal copula, the input of the simulation is the correlation matrix  $\Sigma$ . The normal copula can be simulated by the following steps, in which  $\mathbf{U} = (U_1, \dots, U_m)$  denotes one random draw from the copula:

1. Generate a multivariate normal vector  $\mathbf{Z} \sim N(0, \Sigma)$ , where  $\Sigma$  is an  $m$ -dimensional correlation matrix.
2. Transform the vector  $\mathbf{Z}$  into  $\mathbf{U} = (\Phi(Z_1), \dots, \Phi(Z_m))^T$ , where  $\Phi$  is the distribution function of univariate standard normal.

The first step can be achieved by Cholesky decomposition of the correlation matrix  $\Sigma = LL^T$ , where  $L$  is a lower triangular matrix with positive elements on the diagonal. If  $\tilde{\mathbf{Z}} \sim N(0, I)$ , then  $L\tilde{\mathbf{Z}} \sim N(0, \Sigma)$ .

## Student's $t$ copula

Let  $\Theta = \{(\nu, \Sigma) : \nu \in (1, \infty), \Sigma \in \mathbb{R}^{m \times m}\}$ , and let  $t_\nu$  be a univariate  $t$  distribution with  $\nu$  degrees of freedom.

The Student's  $t$  copula can be written as

$$C_\Theta(u_1, u_2, \dots, u_m) = \mathbf{t}_{\nu, \Sigma} \left( t_\nu^{-1}(u_1), t_\nu^{-1}(u_2), \dots, t_\nu^{-1}(u_m) \right)$$

where  $\mathbf{t}_{\nu, \Sigma}$  is the multivariate Student's  $t$  distribution that has a correlation matrix  $\Sigma$  with  $\nu$  degrees of freedom.

## Simulation

The input parameters for the simulation are  $(\nu, \Sigma)$ . The  $t$  copula can be simulated by the following steps:

1. Generate a multivariate vector  $\mathbf{X} \sim t_m(\nu, 0, \Sigma)$  that follows the centered  $t$  distribution with  $\nu$  degrees of freedom and correlation matrix  $\Sigma$ .
2. Transform the vector  $\mathbf{X}$  into  $\mathbf{U} = (t_\nu(X_1), \dots, t_\nu(X_m))^T$ , where  $t_\nu$  is the distribution function of univariate  $t$  distribution with  $\nu$  degrees of freedom.

To simulate centered multivariate  $t$  random variables, you can use the property that  $\mathbf{X} \sim t_m(\nu, 0, \Sigma)$  if  $\mathbf{X} = \sqrt{\nu/s} \mathbf{Z}$ , where  $\mathbf{Z} \sim N(0, \Sigma)$  and the univariate random variable  $s \sim \chi_\nu^2$ .

## Archimedean Copulas

### Overview of Archimedean Copulas

Let function  $\phi : [0, 1] \rightarrow [0, \infty)$  be a strict Archimedean copula generator function, and suppose that its inverse  $\phi^{-1}$  is completely monotonic on  $[0, \infty)$ . A strict generator is a decreasing function  $\phi : [0, 1] \rightarrow [0, \infty)$  that satisfies  $\phi(0) = \infty$  and  $\phi(1) = 0$ . A decreasing function  $f(t) : [a, b] \rightarrow (-\infty, \infty)$  is completely monotonic if it satisfies

$$(-1)^k \frac{d^k}{dt^k} f(t) \geq 0, k \in \mathbb{N}, t \in (a, b)$$

An Archimedean copula is defined as follows:

$$C(u_1, u_2, \dots, u_m) = \phi^{-1} \left( \phi(u_1) + \dots + \phi(u_m) \right)$$

The Archimedean copulas available in the HPCOPULA procedure are the Clayton copula, the Frank copula, and the Gumbel copula.



### Clayton Copula

Let the generator function  $\phi(u) = \theta^{-1} (u^{-\theta} - 1)$ . A Clayton copula is defined as

$$C_\theta(u_1, u_2, \dots, u_m) = \left[ \sum_{i=1}^m u_i^{-\theta} - m + 1 \right]^{-1/\theta}$$

where  $\theta > 0$ .

### Frank Copula

Let the generator function be

$$\phi(u) = -\log \left[ \frac{\exp(-\theta u) - 1}{\exp(-\theta) - 1} \right]$$

A Frank copula is defined as

$$C_\theta(u_1, u_2, \dots, u_m) = \frac{1}{\theta} \log \left\{ 1 + \frac{\prod_{i=1}^m [\exp(-\theta u_i) - 1]}{[\exp(-\theta) - 1]^{m-1}} \right\}$$

where  $\theta \in (-\infty, \infty) \setminus \{0\}$  for  $m = 2$  and  $\theta > 0$  for  $m \geq 3$ .

### Gumbel Copula

Let the generator function  $\phi(u) = (-\log u)^\theta$ . A Gumbel copula is defined as

$$C_\theta(u_1, u_2, \dots, u_m) = \exp \left\{ - \left[ \sum_{i=1}^m (-\log u_i)^\theta \right]^{1/\theta} \right\}$$

where  $\theta > 1$ .

### Simulation

Suppose that the generator of the Archimedean copula is  $\phi$ . Then the simulation method that uses a Laplace-Stieltjes transformation of the distribution function is given by Marshall and Olkin (1988), where  $\tilde{F}(t) = \int_0^\infty e^{-tx} dF(x)$ :

1. Generate a random variable  $V$  that has the distribution function  $F$  such that  $\tilde{F}(t) = \phi^{-1}(t)$ .
2. Draw samples from the independent uniform random variables  $X_1, \dots, X_m$ .
3. Return  $\mathbf{U} = (\tilde{F}(-\log(X_1)/V), \dots, \tilde{F}(-\log(X_m)/V))^T$ .

The Laplace-Stieltjes transformations are as follows:

- For the Clayton copula,  $\tilde{F} = (1 + t)^{-1/\theta}$ , and the distribution function  $F$  is associated with a gamma random variable that has a shape parameter of  $\theta^{-1}$  and a scale parameter of 1.
- For the Gumbel copula,  $\tilde{F} = \exp(-t^{1/\theta})$ , and  $F$  is the distribution function of the stable variable  $\text{St}(\theta^{-1}, 1, \gamma, 0)$ , where  $\gamma = [\cos(\pi/(2\theta))]^\theta$ .

- For the Frank copula where  $\theta > 0$ ,  $\tilde{F} = -\log\{1 - \exp(-t)[1 - \exp(-\theta)]\}/\theta$ , and  $F$  is a discrete probability function  $P(V = k) = (1 - \exp(-\theta))^k / (k\theta)$ . This probability function is related to a logarithmic random variable that has a parameter value of  $1 - e^{-\theta}$ .

For more information about simulating a random variable from a stable distribution, see Theorem 1.19 in Nolan (2010). For more information about simulating a random variable from a logarithmic series, see Chapter 10.5 in Devroye (1986).

For a Frank copula where  $m = 2$  and  $\theta < 0$ , the simulation can be done through conditional distributions as follows:

1 Draw independent  $v_1, v_2$  from a uniform distribution.

2 Let  $u_1 = v_1$ .

3 Let  $u_2 = -\frac{1}{\theta} \log \left( 1 + \frac{v_2(1-e^{-\theta})}{v_2(e^{-\theta v_1}-1)-e^{-\theta v_1}} \right)$ .

---

## OUTUNIFORM= Data Sets

The number of columns and the names of columns in OUTUNIFORM= data sets match the number and names of the *variables* in the VAR statement.

---

## Examples: HPCOPULA Procedure

---

### Example 5.1: Simulating Default Times

Suppose the correlation structure that is required for a normal copula function is already known. For example, the correlation structure can be estimated from the historical data on default times in some industries, but this estimation is not within the scope of this example. The correlation structure is saved in a SAS data set called `lnparm`. The following statements and their output in [Output 5.1.1](#) show that the correlation parameter is set at 0.8:

```
proc print data = lnparm;
run;
```

**Output 5.1.1** Copula Correlation Matrix

	Obs	Y1	Y2
	1	1.0	0.8
	2	0.8	1.0

The following statements use PROC HPCOPULA to simulate the data:

```

option set=GRIDHOST="%GRIDHOST";
option set=GRIDINSTALLLOC="%GRIDINSTALLLOC";

/* simulate the data from bivariate normal copula */
proc hpcopula;
  var Y1-Y2;
  define cop normal (corr=inparm);
  simulate cop /
    ndraws      = 1000000
    seed        = 1234
    outuniform  = normal_unifdata;
  PERFORMANCE nodes=4 nthreads=4 details
    host="%GRIDHOST" install="%GRIDINSTALLLOC";
run;

```

The VAR statement specifies the list of variables that contains the simulated data. The DEFINE statement assigns the name COP and specifies a normal copula that reads the correlation matrix from the Inparm data set. The SIMULATE statement refers to the COP label that is defined in the VAR statement and specifies several options: the NDRAWS= option specifies a sample size, the SEED= option specifies 1234 as the random number generator seed, and the OUTUNIFORM=NORMAL\_UNIFDATA option names the output data set to contain the result of simulation in uniforms. The PERFORMANCE statement requests that the analytic computations be performed on four nodes in the distributed computing environment and four threads on each node. [Output 5.1.2](#) shows the run time of this particular simulation experiment.

**Output 5.1.2** Run-Time Performance

Performance Information		
Host Node	<< your grid host >>	
Install Location	<< your grid install location >>	
Execution Mode	Distributed	
Grid Mode	Symmetric	
Number of Compute Nodes	4	
Number of Threads per Node	4	
Procedure Task Timing		
Task	Seconds	Percent
Simulation of Model	0.07	0.22%
Writing of output data	31.70	99.78%

The following DATA step transforms the variables from zero-one uniformly distributed to nonnegative exponentially distributed with parameter 0.5 and adds three indicator variables to the data set: SURVIVE1 and SURVIVE2 are equal to 1 if company 1 or company 2, respectively, has remained in business for more than three years, and SURVIVE is equal to 1 if both companies survived the same period together.

```

/* default time has exponential marginal distribution with parameter 0.5 */
data default;
  set normal_unifdata;
  array arr{2} Y1-Y2;
  array time{2} time1-time2;
  array surv{2} survive1-survive2;
  lambda = 0.5;
  do i=1 to 2;
    time[i] = -log(1-arr[i])/lambda;
    surv[i] = 0;
    if (time[i] >3) then surv[i]=1;
  end;
  survive = 0;
  if (time1 >3) && (time2 >3) then survive = 1;
run;

```

The first analysis step is to look at correlations between survival times of the two companies. You can perform this step by using the CORR procedure as follows:

```

proc corr data = default pearson kendall;
  var time1 time2;
run;

```

Output 5.1.3 shows the output of this code. The output contains some descriptive statistics and two measures of correlation: Pearson and Kendall. Both measures indicate high and statistically significant dependence between the life spans of the two companies.

**Output 5.1.3** Default Time Descriptive Statistics and Correlations

The CORR Procedure						
2 Variables:      time1      time2						
Simple Statistics						
Variable	N	Mean	Std Dev	Median	Minimum	Maximum
time1	1000000	2.00042	1.99724	1.38664	1.78961E-6	28.39277
time2	1000000	2.00190	2.00064	1.38787	2.24931E-6	30.50949
Pearson Correlation Coefficients, N = 1000000						
Prob >  r  under H0: Rho=0						
		time1		time2		
	time1	1.00000		0.76950		
				<.0001		
	time2	0.76950		1.00000		
		<.0001				

**Output 5.1.3** *continued*

Kendall Tau b Correlation Coefficients, N = 1000000			
Prob >  tau  under H0: Tau=0			
	time1	time2	
time1	1.00000	0.58998	
		<.0001	
time2	0.58998	1.00000	
	<.0001		

The second and final step is to empirically estimate the default probabilities of the two companies. This is done by using the FREQ procedure as follows:

```
proc freq data=default;
  table survive survive1-survive2;
run;
```

The results are shown in [Output 5.1.4](#).

**Output 5.1.4** Probabilities of Default

The FREQ Procedure					
survive	Frequency	Percent	Cumulative Frequency	Cumulative Percent	
0	852314	85.23	852314	85.23	
1	147686	14.77	1000000	100.00	
survive1	Frequency	Percent	Cumulative Frequency	Cumulative Percent	
0	776565	77.66	776565	77.66	
1	223435	22.34	1000000	100.00	
survive2	Frequency	Percent	Cumulative Frequency	Cumulative Percent	
0	776382	77.64	776382	77.64	
1	223618	22.36	1000000	100.00	

[Output 5.1.4](#) shows that the empirical default probabilities are 75% and 78%. Assuming that these companies are independent yields the probability estimate that both companies default during the period of three years as  $0.75 \times 0.78 = 0.59$  (59%). Comparing this naive estimate with the much higher actual 83% joint default probability illustrates that neglecting the correlation between the two companies significantly underestimates the probability of default.

---

## References

Devroye, L. (1986), *Non-uniform Random Variate Generation*, New York: Springer-Verlag.

URL <http://luc.devroye.org/rnbookindex.html>

Marshall, A. W. and Olkin, I. (1988), “Families of Multivariate Distributions,” *Journal of the American Statistical Association*, 83, 834–841.

Nolan, J. P. (2010), *Stable Distributions: Models for Heavy Tailed Data*, Boston: Birkhäuser.

Sklar, A. (1959), “Fonctions de répartition à  $n$  dimensions et leurs marges,” *Publications de l’Institut de Statistique de L’Université de Paris*, 8, 229–231.

# Chapter 6

## The HPCOUNTREG Procedure

### Contents

---

Overview: HPCOUNTREG Procedure . . . . .	<b>128</b>
PROC HPCOUNTREG Features . . . . .	128
Getting Started: HPCOUNTREG Procedure . . . . .	<b>128</b>
Syntax: HPCOUNTREG Procedure . . . . .	<b>131</b>
Functional Summary . . . . .	132
PROC HPCOUNTREG Statement . . . . .	134
BOUNDS Statement . . . . .	137
BY Statement . . . . .	138
FREQ Statement . . . . .	138
INIT Statement . . . . .	138
MODEL Statement . . . . .	138
OUTPUT Statement . . . . .	139
PERFORMANCE Statement . . . . .	140
RESTRICT Statement . . . . .	141
WEIGHT Statement . . . . .	141
ZEROMODEL Statement . . . . .	142
Details: HPCOUNTREG Procedure . . . . .	<b>143</b>
Missing Values . . . . .	143
Poisson Regression . . . . .	143
Negative Binomial Regression . . . . .	144
Zero-Inflated Count Regression Overview . . . . .	146
Zero-Inflated Poisson Regression . . . . .	146
Zero-Inflated Negative Binomial Regression . . . . .	148
Computational Resources . . . . .	150
Covariance Matrix Types . . . . .	150
Displayed Output . . . . .	150
OUTPUT OUT= Data Set . . . . .	152
OUTEST= Data Set . . . . .	152
ODS Table Names . . . . .	152
Examples: The HPCOUNTREG Procedure . . . . .	<b>153</b>
Example 6.1: High-Performance Zero-Inflated Poisson Model . . . . .	153
References . . . . .	<b>157</b>

---

---

## Overview: HPCOUNTREG Procedure

The HPCOUNTREG procedure is a high-performance version of the COUNTREG procedure in SAS/ETS software. Like the COUNTREG procedure, the HPCOUNTREG procedure fits regression models in which the dependent variable takes on nonnegative integer or count values. Unlike the COUNTREG procedure, which can be run only on an individual workstation, the HPCOUNTREG procedure takes advantage of a computing environment that enables it to distribute the optimization task among one or more nodes. In addition, each node can use one or more threads to carry out the optimization on its subset of the data. When several nodes are employed, with each node using several threads to carry out its part of the work, the result is a highly parallel computation that provides a dramatic gain in performance.

The HPCOUNTREG procedure enables you to read and write data in distributed form and perform analyses in distributed mode and single-machine mode. For information about how to affect the execution mode of SAS high-performance analytical procedures, see the section “[Processing Modes](#)” on page 12 in Chapter 3, “[Shared Concepts and Topics](#).”

The HPCOUNTREG procedure is specifically designed to operate in the high-performance distributed environment. By default, PROC HPCOUNTREG performs computations in multiple threads.

---

## PROC HPCOUNTREG Features

The HPCOUNTREG procedure estimates the parameters of a count regression model by maximum likelihood techniques. The following list summarizes some basic features of the HPCOUNTREG procedure:

- can perform analysis on a massively parallel high-performance appliance
- reads input data in parallel and writes output data in parallel when the data source is the appliance database
- is highly multithreaded during all phases of analytic execution
- performs maximum likelihood estimation
- supports multiple link functions
- uses the [WEIGHT](#) statement for weighted analysis
- uses the [FREQ](#) statement for grouped analysis
- uses the [OUTPUT](#) statement to produce a data set that contains predicted probabilities and other observationwise statistics

---

## Getting Started: HPCOUNTREG Procedure

Except for its ability to operate in the high-performance distributed environment, the HPCOUNTREG procedure is similar in use to other regression model procedures in the SAS System. For example, the following statements are used to estimate a Poisson regression model:



```
proc hpcountreg data=one ;
  model y = x / dist=poisson ;
run;
```

The response variable *y* is numeric and has nonnegative integer values.

This section illustrates two simple examples that use PROC HPCOUNTREG. The data are taken from Long (1997). This study examines how factors such as gender (*fem*), marital status (*mar*), number of young children (*kid5*), prestige of the graduate program (*phd*), and number of articles published by a scientist's mentor (*ment*) affect the number of articles (*art*) published by the scientist.

The first 10 observations are shown in [Figure 6.1](#).

**Figure 6.1** Article Count Data

Obs	art	fem	mar	kid5	phd	ment
1	3	0	1	2	1.38000	8.0000
2	0	0	0	0	4.29000	7.0000
3	4	0	0	0	3.85000	47.0000
4	1	0	1	1	3.59000	19.0000
5	1	0	1	0	1.81000	0.0000
6	1	0	1	1	3.59000	6.0000
7	0	0	1	1	2.12000	10.0000
8	0	0	1	0	4.29000	2.0000
9	3	0	1	2	2.58000	2.0000
10	3	0	1	1	1.80000	4.0000

The following SAS statements estimate the Poisson regression model. The model is executed in the distributed computing environment with two threads and four nodes.

```
/*-- Poisson Regression --*/
proc hpcountreg data=long97data;
  model art = fem mar kid5 phd ment / dist=poisson method=quanew;
  performance nthreads=2 nodes=4 details;
run;
```

The “Model Fit Summary” table that is shown in [Figure 6.2](#) lists several details about the model. By default, the HPCOUNTREG procedure uses the Newton-Raphson optimization technique. The maximum log-likelihood value is shown, in addition to two information measures—Akaike’s information criterion (AIC) and Schwarz’s Bayesian information criterion (SBC)—which can be used to compare competing Poisson models. Smaller values of these criteria indicate better models.

**Figure 6.2** Estimation Summary Table for a Poisson Regression

The HPCOUNTREG Procedure	
Model Fit Summary	
Dependent Variable	art
Number of Observations	915
Data Set	WORK.LONG97DATA
Model	Poisson
Log Likelihood	-1651
Maximum Absolute Gradient	0.0002080
Number of Iterations	13
Optimization Method	Quasi-Newton
AIC	3314
SBC	3343

Figure 6.3 shows the parameter estimates of the model and their standard errors. All covariates are significant predictors of the number of articles, except for the prestige of the program (phd), which has a  $p$ -value of 0.6271.

**Figure 6.3** Parameter Estimates of Poisson Regression

Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	0.3046	0.1030	2.96	0.0031
fem	1	-0.2246	0.05461	-4.11	<.0001
mar	1	0.1552	0.06137	2.53	0.0114
kid5	1	-0.1849	0.04013	-4.61	<.0001
phd	1	0.01282	0.02640	0.49	0.6271
ment	1	0.02554	0.002006	12.73	<.0001

To allow for variance greater than the mean, you can fit the negative binomial model instead of the Poisson model by specifying the DIST=NEGBIN option, as shown in the following statements. Whereas the Poisson model requires that the conditional mean and conditional variance be equal, the negative binomial model allows for overdispersion, in which the conditional variance can exceed the conditional mean.

```

/*-- Negative Binomial Regression --*/
proc hpcountreg data=long97data;
  model art = fem mar kid5 phd ment / dist=negbin(p=2) method=quanew;
  performance nthreads=2 nodes=4 details;
run;

```

Figure 6.4 shows the fit summary and Figure 6.5 shows the parameter estimates.

**Figure 6.4** Estimation Summary Table for a Negative Binomial Regression

The HPCOUNTREG Procedure	
Model Fit Summary	
Dependent Variable	art
Number of Observations	915
Data Set	WORK.LONG97DATA
Model	NegBin
Log Likelihood	-1561
Maximum Absolute Gradient	0.0000666
Number of Iterations	16
Optimization Method	Quasi-Newton
AIC	3136
SBC	3170

**Figure 6.5** Parameter Estimates of Negative Binomial Regression

Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	0.2561	0.1386	1.85	0.0645
fem	1	-0.2164	0.07267	-2.98	0.0029
mar	1	0.1505	0.08211	1.83	0.0668
kid5	1	-0.1764	0.05306	-3.32	0.0009
phd	1	0.01527	0.03604	0.42	0.6718
ment	1	0.02908	0.003470	8.38	<.0001
_Alpha	1	0.4416	0.05297	8.34	<.0001

The parameter estimate for `_Alpha` of 0.4416 is an estimate of the dispersion parameter in the negative binomial distribution. A  $t$  test for the hypothesis  $H_0 : \alpha = 0$  is provided. It is highly significant, indicating overdispersion ( $p < 0.0001$ ).

The null hypothesis  $H_0 : \alpha = 0$  can be also tested against the alternative  $\alpha > 0$  by using the likelihood ratio test, as described by Cameron and Trivedi (1998, pp. 45, 77–78). The likelihood ratio test statistic is equal to  $-2(\mathcal{L}_P - \mathcal{L}_{NB}) = -2(-1651 + 1561) = 180$ , which is highly significant, providing strong evidence of overdispersion.

## Syntax: HPCOUNTREG Procedure

The following statements are available in the HPCOUNTREG procedure. Items within angle brackets (< >) or square brackets ([ ]) are optional.

```

PROC HPCOUNTREG <options> ;
  BOUNDS bound1 [ , bound2 ... ] ;
  BY variables ;
  FREQ freq-variable ;
  INIT initialization1 < , initialization2 ... > ;
  MODEL dependent-variable = regressors </ options> ;
  OUTPUT <output-options> ;
  PERFORMANCE performance-options ;
  RESTRICT restriction1 [ , restriction2 ... ] ;
  WEIGHT variable </ option> ;
  ZEROMODEL dependent-variable ~ zero-inflated-regressors </ options> ;

```

There can be only one MODEL statement. The ZEROMODEL statement, if used, must appear after the MODEL statement. If a FREQ or WEIGHT statement is specified more than once, the variable specified in the first instance is used.

## Functional Summary

Table 6.1 summarizes the statements and options used with the HPCOUNTREG procedure.

**Table 6.1** PROC HPCOUNTREG Functional Summary

Description	Statement	Option
<b>Data Set Options</b>		
Specifies the input data set	PROC HPCOUN- TREG	DATA=
Writes parameter estimates to an output data set	PROC HPCOUN- TREG	OUTEST=
Writes estimates to an output data set	OUTPUT	OUT=
Specifies BY-group processing	BY	
Specifies an optional frequency variable	FREQ	
Specifies an optional weight variable	WEIGHT	
<b>Printing Control Options</b>		
Prints the correlation matrix of the estimates	PROC HPCOUN- TREG	CORRB
Prints the covariance matrix of the estimates	PROC HPCOUN- TREG	COVB
Suppresses the normal printed output	PROC HPCOUN- TREG	NOPRINT
Requests all printing options	PROC HPCOUN- TREG	PRINTALL
<b>Options to Control the Optimization Process</b>		
Specifies maximum number of iterations allowed	PROC HPCOUN- TREG	MAXITER=

Description	Statement	Option
Selects the iterative minimization method to use	PROC HPCOUN- TREG	METHOD=
Specifies maximum number of iterations allowed	PROC HPCOUN- TREG	MAXITER=
Specifies maximum number of function calls	PROC HPCOUN- TREG	MAXFUNC=
Specifies the upper limit of CPU time in seconds	PROC HPCOUN- TREG	MAXTIME=
Specifies absolute function convergence criterion	PROC HPCOUN- TREG	ABSCONV=
Specifies absolute function convergence criterion	PROC HPCOUN- TREG	ABSFCNV=
Specifies absolute gradient convergence criterion	PROC HPCOUN- TREG	ABSGCONV=
Specifies relative function convergence criterion	PROC HPCOUN- TREG	FCONV=
Specifies relative gradient convergence criterion	PROC HPCOUN- TREG	GCONV=
Specifies absolute parameter convergence criterion	PROC HPCOUN- TREG	ABSXCONV=
Specifies matrix singularity criterion	PROC HPCOUN- TREG	SINGULAR=
Sets boundary restrictions on parameters	BOUNDS	
Sets initial values for parameters	INIT	
Sets linear restrictions on parameters	RESTRICT	
<b>Model Estimation Options</b>		
Specifies the type of model	PROC HPCOUN- TREG	DIST=
Specifies the type of covariance matrix	PROC HPCOUN- TREG	COVEST=
Suppresses the intercept parameter	MODEL	NOINT
Specifies the offset variable	MODEL	OFFSET=
Specifies the zero-inflated offset variable	ZEROMODEL	OFFSET=
Specifies the zero-inflated link function	ZEROMODEL	LINK=
<b>Output Control Options</b>		
Includes covariances in the OUTEST= data set	PROC HPCOUN- TREG	COVOUT
Includes correlations in the OUTEST= data set	PROC HPCOUN- TREG	CORROUT
Outputs SAS variables to the output data set	OUTPUT	COPYVAR=
Outputs probability of the actual value	OUTPUT	PROB=
Outputs expected value of response variable	OUTPUT	PRED=
Outputs estimates of $\mathbf{XBeta} = \mathbf{x}_i' \boldsymbol{\beta}$	OUTPUT	XBETA=
Outputs estimates of $\mathbf{ZGamma} = \mathbf{z}_i' \boldsymbol{\gamma}$	OUTPUT	ZGAMMA=

Description	Statement	Option
Outputs probability of a zero value as a result of the zero-generating process	OUTPUT	PROBZERO=
<b>Performance Options</b>		
Requests a table that shows a timing breakdown	PERFORMANCE	DETAILS
Specifies the number of threads to use	PERFORMANCE	NTHREADS=
Specifies the number of nodes to use on the SAS appliance	PERFORMANCE	NODES=

## PROC HPCOUNTREG Statement

**PROC HPCOUNTREG** *<options>* ;

The following *options* can be used in the PROC HPCOUNTREG statement.

### Input Data Set Option

**DATA=SAS-data-set**

specifies the input SAS data set. If the DATA= option is not specified, PROC HPCOUNTREG uses the most recently created SAS data set.

### Output Data Set Options

**OUTEST=SAS-data-set**

writes the parameter estimates to the specified output data set.

**CORROUT**

writes the correlation matrix for the parameter estimates to the OUTEST= data set. This option is valid only if the OUTEST= option is specified.

**COVOUT**

writes the covariance matrix for the parameter estimates to the OUTEST= data set. This option is valid only if the OUTEST= option is specified.

### Printing Options

You can specify the following options in either the PROC HPCOUNTREG statement or the MODEL statement:

**CORRB**

prints the correlation matrix of the parameter estimates.

**COVB**

prints the covariance matrix of the parameter estimates.

**NOPRINT**

suppresses all printed output.

**PRINTALL**

requests all printing options.

## Estimation Control Options

You can specify the following options in either the PROC HPCOUNTREG statement or the MODEL statement:

**COVEST=***value*

specifies the type of covariance matrix for the parameter estimates.

The default is COVEST=HESSIAN. You can specify the following *values*:

HESSIAN	specifies the covariance from the Hessian matrix.
OP	specifies the covariance from the outer product matrix.
QML	specifies the covariance from the outer product and Hessian matrices.

## Optimization Control Options

PROC HPCOUNTREG uses the nonlinear optimization (NLO) subsystem to perform nonlinear optimization tasks. You can specify the following *options* in either the PROC HPCOUNTREG statement or the MODEL statement.

**ABSCONV=***r***ABSTOL=***r*

specifies an absolute function value convergence criterion by which minimization stops when  $f(\theta^{(k)}) \leq r$ . The default value of *r* is the negative square root of the largest double-precision value, which serves only as a protection against overflows.

**ABSFCONV=***r***ABSFTOL=***r*

specifies an absolute function difference convergence criterion by which minimization stops when the function value has a small change in successive iterations:

$$|f(\theta^{(k-1)}) - f(\theta^{(k)})| \leq r$$

The default is 0.

**ABSGCONV=***r***ABSGTOL=***r*

specifies an absolute gradient convergence criterion. Optimization stops when the maximum absolute gradient element is small:

$$\max_j |g_j(\theta^{(k)})| \leq r$$

The default is 1E-5.

**ABSXCONV=*r*****ABSXTOL=*r***

specifies an absolute parameter convergence criterion. Optimization stops when the Euclidean distance between successive parameter vectors is small:

$$\| \theta^{(k)} - \theta^{(k-1)} \|_2 \leq r$$

The default is 0.

**FCONV=*r*****FTOL=*r***

specifies a relative function convergence criterion. Optimization stops when a relative change of the function value in successive iterations is small:

$$\frac{|f(\theta^{(k)}) - f(\theta^{(k-1)})|}{|f(\theta^{(k-1)})|} \leq r$$

The default value is  $2\epsilon$ , where  $\epsilon$  denotes the machine precision constant, which is the smallest double-precision floating-point number such that  $1 + \epsilon > 1$ .

**GCONV=*r*****GTOL=*r***

specifies a relative gradient convergence criterion. For all techniques except CONGRA, optimization stops when the normalized predicted function reduction is small:

$$\frac{g(\theta^{(k)})^T [H^{(k)}]^{-1} g(\theta^{(k)})}{|f(\theta^{(k)})|} \leq r$$

For the CONGRA technique (where a reliable Hessian estimate  $H$  is not available), the following criterion is used:

$$\frac{\|g(\theta^{(k)})\|_2^2 \|s(\theta^{(k)})\|_2}{\|g(\theta^{(k)}) - g(\theta^{(k-1)})\|_2 |f(\theta^{(k)})|} \leq r$$

The default is 1E-8.

**MAXFUNC=*i*****MAXFU=*i***

specifies the maximum number of function calls in the optimization process. The default is 1,000.

The optimization can terminate only after completing a full iteration. Therefore, the number of function calls that are actually performed can exceed the number of calls that are specified by this option.

**MAXITER=*i*****MAXIT=*i***

specifies the maximum number of iterations in the optimization process. The default is 200.

**MAXTIME=*r***

specifies an upper limit of  $r$  seconds of CPU time for the optimization process. The default value is the largest floating-point double representation of your computer. The time that is specified by this option is checked only once at the end of each iteration. Therefore, the actual run time can be much longer than  $r$ . The actual run time includes the remaining time needed to finish the iteration and the time needed to generate the output of the results.



### METHOD=*value*

specifies the iterative minimization method to use. The default is METHOD=NEWRAP. You can specify the following *values*:

<b>CONGRA</b>	specifies the conjugate-gradient method.
<b>DBLDOG</b>	specifies the double-dogleg method.
<b>NEWRAP</b>	specifies the Newton-Raphson method (this is the default).
<b>NONE</b>	specifies that no optimization be performed beyond using the ordinary least squares method to compute the parameter estimates.
<b>NRRIDG</b>	specifies the Newton-Raphson Ridge method.
<b>QUANEW</b>	specifies the quasi-Newton method.
<b>TRUREG</b>	specifies the trust region method.

### SINGULAR=*r*

specifies the general singularity criterion that is applied by the HPCOUNTREG procedure in sweeps and inversions. The default is 1E-8.

---

## BOUNDS Statement

**BOUNDS** *bound1* [, *bound2* ...] ;

The BOUNDS statement imposes simple boundary constraints on the parameter estimates. You can specify any number of BOUNDS statements.

Each *bound* is composed of parameter names, constants, and inequality operators as follows:

*item operator item* [ *operator item* [ *operator item* ... ] ]

Each *item* is a constant, a parameter name, or a list of parameter names. Each *operator* is <, >, <=, or >=. Parameter names are as shown in the Effect column of the “Parameter Estimates” table.

You can use both the BOUNDS statement and the RESTRICT statement to impose boundary constraints. However, the BOUNDS statement provides a simpler syntax for specifying these kinds of constraints. For more information, see the section “[RESTRICT Statement](#)” on page 141.

The following BOUNDS statement illustrates the use of parameter lists to specify boundary constraints. It constrains the estimates of the parameter for *z* to be negative, the parameters for *x1* through *x10* to be between 0 and 1, and the parameter for *x1* in the zero-inflation model to be less than 1.

```
bounds z < 0,
       0 < x1-x10 < 1,
       Inf_x1 < 1;
```

---

## BY Statement

**BY** *variables* ;

A BY statement can be used with PROC HPCOUNTREG to obtain separate analyses on observations in groups defined by the BY variables. When a BY statement appears, the input data set should be sorted in order of the BY variables.

BY statement processing is not supported when the HPCOUNTREG procedure runs alongside the database or alongside the Hadoop Distributed File System (HDFS). These modes are used if the input data are stored in a database or HDFS and the grid host is the appliance that houses the data.

---

## FREQ Statement

**FREQ** *freq-variable* ;

The FREQ statement identifies a variable (*freq-variable*) that contains the frequency of occurrence of each observation. PROC HPCOUNTREG treats each observation as if it appears  $n$  times, where  $n$  is the value of *freq-variable* for the observation. If the value for the observation is not an integer, it is truncated to an integer. If the value is less than 1 or missing, the observation is not used in the model fitting. When the FREQ statement is not specified, each observation is assigned a frequency of 1.

---

## INIT Statement

**INIT** *initialization1* < , *initialization2* ... > ;

The INIT statement sets initial values for parameters in the optimization.

Each *initialization* is written as a parameter or parameter list, followed by an optional equal sign (=), followed by a number:

*parameter* <=> *number*

Parameter names are as shown in the Effect column of the “Parameter Estimates” table.

---

## MODEL Statement

**MODEL** *dependent-variable* = *regressors* </ *options* > ;

The MODEL statement specifies the dependent variable and independent regressor variables for the regression model. The dependent count variable should take only nonnegative integer values from the input data set. PROC HPCOUNTREG rounds any positive noninteger count value to the nearest integer. PROC HPCOUNTREG discards any observation that has a negative count.

Only one MODEL statement can be specified. You can specify the following *options* in the MODEL statement after a slash (/).

**DIST=***value*

specifies a type of model to be analyzed. You can specify the following *values*:

POISSON | P specifies the Poisson regression model.

NEGBIN(P=1) specifies the negative binomial regression model that uses a linear variance function.

NEGBIN(P=2) | NEGBIN specifies the negative binomial regression model that uses a quadratic variance function.

ZIPOISSON | ZIP specifies zero-inflated Poisson regression.

ZINEGBIN | ZINB specifies zero-inflated negative binomial regression.

You can also specify the DIST option in the HPCOUNTREG statement.

**NOINT**

suppresses the intercept parameter.

**OFFSET=***offset-variable*

specifies a variable in the input data set to be used as an offset variable. The *offset-variable* is used to allow the observational units to vary across observations. For example, when the number of shipping accidents could be measured across different time periods or the number of students who participate in an activity could be reported across different class sizes, the observational units need to be adjusted to a common denominator by using the offset variable. The offset variable appears as a covariate in the model with its parameter restricted to 1. The offset variable cannot be the response variable, the zero-inflation offset variable (if any), or any of the explanatory variables. The “Model Fit Summary” table gives the name of the data set variable that is used as the offset variable; it is labeled “Offset.”

## Printing Options

You can specify the following options in either the PROC HPCOUNTREG statement or the MODEL statement:

**CORRB**

prints the correlation matrix of the parameter estimates.

**COVB**

prints the covariance matrix of the parameter estimates.

**NOPRINT**

suppresses all printed output.

**PRINTALL**

requests all printing options.

---

## OUTPUT Statement

**OUTPUT** < *output-options* > ;

The OUTPUT statement creates a new SAS data set that includes variables created by the *output-options*. These variables include the estimates of  $\mathbf{x}'_i\boldsymbol{\beta}$ , the expected value of the response variable, and the probability

of the response variable taking on the current value. Furthermore, if a zero-inflated model was fit, you can request that the output data set contain the estimates of  $\mathbf{z}_i' \boldsymbol{\gamma}$  and the probability that the response is zero as a result of the zero-generating process. These statistics can be computed for all observations in which the regressors are not missing, even if the response is missing. By adding observations that have missing response values to the input data set, you can compute these statistics for new observations or for settings of the regressors that are not present in the data without affecting the model fit.

You can specify only one OUTPUT statement. You can specify the following *output-options*:

**OUT=SAS-data-set**

names the output data set

**COPYVAR=SAS-variable-names**

**COPYVARS=SAS-variable-names**

adds SAS variables to the output data set.

**PRED=name**

names the variable to contain the predicted value of the response variable.

**PROB=name**

names the variable to contain the probability that the response variable will take the actual value,  $\Pr(Y = y_i)$ .

**PROBCOUNT(value1 < value2 ... >)**

outputs the probability that the response variable will take particular values. Each value should be a nonnegative integer. Nonintegers are rounded to the nearest integer. For *value*, you can also specify a list of the form X TO Y BY Z. For example, PROBCOUNT(0 1 2 TO 10 BY 2 15) requests predicted probabilities for the counts 0, 1, 2, 4, 5, 6, 8, 10, and 15. This option is not available for the fixed- and random-effects panel models.

**PROBZERO=name**

names the variable to contain the value of  $\varphi_i$ , which is the probability that the response variable will take the value of 0 as a result of the zero-generating process. This variable is written to the output file only if the model is zero-inflated.

**XBETA=name**

names the variable to contain estimates of  $\mathbf{x}_i' \boldsymbol{\beta}$ .

**ZGAMMA=name**

names the variable to contain estimates of  $\mathbf{z}_i' \boldsymbol{\gamma}$ .

---

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement specifies options to control the multithreaded and distributed computing environment and requests detailed results about the performance characteristics of the HPCOUNTREG procedure. You can also use the PERFORMANCE statement to control whether the HPCOUNTREG procedure executes in single-machine or distributed mode. The most commonly used *performance-options* in the PERFORMANCE statement are as follows:

**DETAILS**

requests a table that shows a timing breakdown of the procedure steps.

**NODES=*n***

specifies the number of nodes in the distributed computing environment, provided that the data are not processed alongside the database.

**NTHREADS=*n***

specifies the number of threads for analytic computations and overrides the SAS system option `THREADS` | `NOTHREADS`. If you do not specify the `NTHREADS=` option, PROC HPCOUNTREG creates one thread per CPU for the analytic computations.

For more information about the PERFORMANCE statement for high-performance analytical procedures, see the section “[PERFORMANCE Statement](#)” on page 39 of Chapter 3, “[Shared Concepts and Topics](#).”

---

## RESTRICT Statement

**RESTRICT** *restriction1* [, *restriction2* ...] ;

The RESTRICT statement imposes linear restrictions on the parameter estimates. You can specify any number of RESTRICT statements.

Each *restriction* is written as an expression, followed by an equality operator (=) or an inequality operator (<, >, <=, >=) and then by a second expression, as follows:

*expression operator expression*

The *operator* can be =, <, >, <=, or >=.

Restriction expressions can be composed of parameter names, constants, and the following operators: times (\*), plus (+), and minus (−). Parameter names are as shown in the Effect column of the “Parameter Estimates” table. The restriction expressions must be a linear function of the variables.

Lagrange multipliers are reported in the “Parameter Estimates” table for all the active linear constraints. They are identified by the names Restrict1, Restrict2, and so on. The probabilities of these Lagrange multipliers are computed using a beta distribution (LaMotte 1994). Nonactive (nonbinding) restrictions have no effect on the estimation results and are not noted in the output.

The following RESTRICT statement constrains the negative binomial dispersion parameter  $\alpha$  to 1, which restricts the conditional variance to be  $\mu + \mu^2$ :

```
restrict _Alpha = 1;
```

---

## WEIGHT Statement

**WEIGHT** *variable* </option> ;

The WEIGHT statement specifies a variable to supply weighting values to use for each observation in estimating parameters. The log likelihood for each observation is multiplied by the corresponding weight variable value.

If the weight of an observation is nonpositive, that observation is not used in the estimation.

The following *option* can be added to the WEIGHT statement after a slash (/).

### NONNORMALIZE

does not normalize the weights. (By default, the weights are normalized so that they add up to the actual sample size. The weights  $w_i$  are normalized by multiplying them by  $\frac{n}{\sum_{i=1}^n w_i}$ , where  $n$  is the sample size.) If the weights are required to be used as they are, then specify the NONNORMALIZE option.

---

## ZEROMODEL Statement

**ZEROMODEL** *dependent-variable* ~ *zero-inflated-regressors* < / *options* > ;

The ZEROMODEL statement is required if either ZIP or ZINB is specified in the DIST= option in the MODEL statement. If ZIP or ZINB is specified, then the ZEROMODEL statement must follow the MODEL statement. The dependent variable in the ZEROMODEL statement must be the same as the dependent variable in the MODEL statement.

The zero-inflated (ZI) regressors appear in the equation that determines the probability ( $\varphi_i$ ) of a zero count. Each of these  $q$  variables has a parameter to be estimated in the regression. For example, let  $\mathbf{z}'_i$  be the  $i$ th observation's  $1 \times (q + 1)$  vector of values of the  $q$  ZI explanatory variables ( $w_0$  is set to 1 for the intercept term). Then  $\varphi_i$  is a function of  $\mathbf{z}'_i \boldsymbol{\gamma}$ , where  $\boldsymbol{\gamma}$  is the  $(q + 1) \times 1$  vector of parameters to be estimated. (The zero-inflated intercept is  $\gamma_0$ ; the coefficients for the  $q$  zero-inflated covariates are  $\gamma_1, \dots, \gamma_q$ .) If  $q$  is equal to 0 (no ZI explanatory variables are provided), then only the intercept term  $\gamma_0$  is estimated. The “Parameter Estimates” table in the displayed output shows the estimates for the ZI intercept and ZI explanatory variables; they are labeled with the prefix “Inf\_”. For example, the ZI intercept is labeled “Inf\_intercept”. If you specify Age (a variable in your data set) as a ZI explanatory variable, then the “Parameter Estimates” table labels the corresponding parameter estimate “Inf\_Age”.

You can specify the following *options* in the ZEROMODEL statement after a slash (/):

### LINK=*value*

specifies the distribution function used to compute probability of zeros. The supported distribution functions are as follows:

LOGISTIC      specifies logistic distribution.

NORMAL        specifies standard normal distribution.

If this option is omitted, then the default ZI link function is logistic.

### OFFSET=*zero-inflated-offset-variable*

specifies a variable in the input data set to be used as a zero-inflated (ZI) offset variable. The ZI offset variable *zero-inflated-offset-variable* is included as a term, with coefficient restricted to 1, in the equation that determines the probability ( $\varphi_i$ ) of a zero count and represents an adjustment to a common observational unit. The ZI offset variable cannot be the response variable, the offset variable (if any), or any of the explanatory variables. The name of the data set variable that is used as the ZI offset variable is displayed in the “Model Fit Summary” table, where it is labeled as “Inf\_offset”.

## Details: HPCOUNTREG Procedure

### Missing Values

Any observations in the input data set that have a missing value for one or more of the regressors are ignored by PROC HPCOUNTREG and not used in the model fit. PROC HPCOUNTREG rounds any positive noninteger count values to the nearest integer and ignores any observations that have a negative count.

If the input data set contains any observations that have missing response values but nonmissing regressors, PROC HPCOUNTREG can compute several statistics and store them in an output data set by using the OUTPUT statement. For example, you can request that the output data set contain the estimates of  $\mathbf{x}_i' \boldsymbol{\beta}$ , the expected value of the response variable, and the probability that the response variable will take the current value. Furthermore, if a zero-inflated model was fit, you can request that the output data set contain the estimates of  $\mathbf{z}_i' \boldsymbol{\gamma}$ , and the probability that the response is 0 as a result of the zero-generating process. Note that the presence of such observations (that have missing response values) does not affect the model fit.

### Poisson Regression

The most widely used model for count data analysis is Poisson regression. Poisson regression assumes that  $y_i$ , given the vector of covariates  $\mathbf{x}_i$ , is independently Poisson distributed with

$$P(Y_i = y_i | \mathbf{x}_i) = \frac{e^{-\mu_i} \mu_i^{y_i}}{y_i!}, \quad y_i = 0, 1, 2, \dots$$

and the mean parameter—that is, the mean number of events per period—is given by

$$\mu_i = \exp(\mathbf{x}_i' \boldsymbol{\beta})$$

where  $\boldsymbol{\beta}$  is a  $(k + 1) \times 1$  parameter vector. (The intercept is  $\beta_0$ ; the coefficients for the  $k$  regressors are  $\beta_1, \dots, \beta_k$ .) Taking the exponential of  $\mathbf{x}_i' \boldsymbol{\beta}$  ensures that the mean parameter  $\mu_i$  is nonnegative. It can be shown that the conditional mean is given by

$$E(y_i | \mathbf{x}_i) = \mu_i = \exp(\mathbf{x}_i' \boldsymbol{\beta})$$

Note that the conditional variance of the count random variable is equal to the conditional mean in the Poisson regression model:

$$V(y_i | \mathbf{x}_i) = E(y_i | \mathbf{x}_i) = \mu_i$$

The equality of the conditional mean and variance of  $y_i$  is known as *equidispersion*.

The standard estimator for the Poisson model is the maximum likelihood estimator (MLE). Because the observations are independent, the log-likelihood function is written as

$$\mathcal{L} = \sum_{i=1}^N (-\mu_i + y_i \ln \mu_i - \ln y_i!) = \sum_{i=1}^N (-e^{\mathbf{x}_i' \boldsymbol{\beta}} + y_i \mathbf{x}_i' \boldsymbol{\beta} - \ln y_i!)$$

For more information about the Poisson regression model, see the section “Poisson Regression” (Chapter 11, *SAS/ETS User’s Guide*).

The Poisson model has been criticized for its restrictive property that the conditional variance equals the conditional mean. Real-life data are often characterized by *overdispersion*—that is, the variance exceeds the mean. Allowing for overdispersion can improve model predictions because the Poisson restriction of equal mean and variance results in the underprediction of zeros when overdispersion exists. The most commonly used model that accounts for overdispersion is the negative binomial model.

## Negative Binomial Regression

The Poisson regression model can be generalized by introducing an unobserved heterogeneity term for observation  $i$ . Thus, the individuals are assumed to differ randomly in a manner that is not fully accounted for by the observed covariates. This is formulated as

$$E(y_i | \mathbf{x}_i, \tau_i) = \mu_i \tau_i = e^{\mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i}$$

where the unobserved heterogeneity term  $\tau_i = e^{\epsilon_i}$  is independent of the vector of regressors  $\mathbf{x}_i$ . Then the distribution of  $y_i$  conditional on  $\mathbf{x}_i$  and  $\tau_i$  is Poisson with conditional mean and conditional variance  $\mu_i \tau_i$ :

$$f(y_i | \mathbf{x}_i, \tau_i) = \frac{\exp(-\mu_i \tau_i) (\mu_i \tau_i)^{y_i}}{y_i!}$$

Let  $g(\tau_i)$  be the probability density function of  $\tau_i$ . Then, the distribution  $f(y_i | \mathbf{x}_i)$  (no longer conditional on  $\tau_i$ ) is obtained by integrating  $f(y_i | \mathbf{x}_i, \tau_i)$  with respect to  $\tau_i$ :

$$f(y_i | \mathbf{x}_i) = \int_0^\infty f(y_i | \mathbf{x}_i, \tau_i) g(\tau_i) d\tau_i$$

An analytical solution to this integral exists when  $\tau_i$  is assumed to follow a gamma distribution. This solution is the negative binomial distribution. If the model contains a constant term, then in order to identify the mean of the distribution, it is necessary to assume that  $E(e^{\epsilon_i}) = E(\tau_i) = 1$ . Thus, it is assumed that  $\tau_i$  follows a  $\text{gamma}(\theta, \theta)$  distribution with  $E(\tau_i) = 1$  and  $V(\tau_i) = 1/\theta$ ,

$$g(\tau_i) = \frac{\theta^\theta}{\Gamma(\theta)} \tau_i^{\theta-1} \exp(-\theta \tau_i)$$

where  $\Gamma(x) = \int_0^\infty z^{x-1} \exp(-z) dz$  is the gamma function and  $\theta$  is a positive parameter. Then, the density of  $y_i$  given  $\mathbf{x}_i$  is derived as

$$\begin{aligned} f(y_i | \mathbf{x}_i) &= \int_0^\infty f(y_i | \mathbf{x}_i, \tau_i) g(\tau_i) d\tau_i \\ &= \frac{\theta^\theta \mu_i^{y_i}}{y_i! \Gamma(\theta)} \int_0^\infty e^{-(\mu_i + \theta) \tau_i} \tau_i^{\theta + y_i - 1} d\tau_i \\ &= \frac{\theta^\theta \mu_i^{y_i} \Gamma(y_i + \theta)}{y_i! \Gamma(\theta) (\theta + \mu_i)^{\theta + y_i}} \\ &= \frac{\Gamma(y_i + \theta)}{y_i! \Gamma(\theta)} \left( \frac{\theta}{\theta + \mu_i} \right)^\theta \left( \frac{\mu_i}{\theta + \mu_i} \right)^{y_i} \end{aligned}$$



If you make the substitution  $\alpha = \frac{1}{\theta}$  ( $\alpha > 0$ ), the negative binomial distribution can then be rewritten as

$$f(y_i | \mathbf{x}_i) = \frac{\Gamma(y_i + \alpha^{-1})}{y_i! \Gamma(\alpha^{-1})} \left( \frac{\alpha^{-1}}{\alpha^{-1} + \mu_i} \right)^{\alpha^{-1}} \left( \frac{\mu_i}{\alpha^{-1} + \mu_i} \right)^{y_i}, \quad y_i = 0, 1, 2, \dots$$

Thus, the negative binomial distribution is derived as a gamma mixture of Poisson random variables. It has the conditional mean

$$E(y_i | \mathbf{x}_i) = \mu_i = e^{\mathbf{x}_i' \boldsymbol{\beta}}$$

and the conditional variance

$$V(y_i | \mathbf{x}_i) = \mu_i \left[ 1 + \frac{1}{\theta} \mu_i \right] = \mu_i [1 + \alpha \mu_i] > E(y_i | \mathbf{x}_i)$$

The conditional variance of the negative binomial distribution exceeds the conditional mean. Overdispersion results from neglected unobserved heterogeneity. The negative binomial model with variance function  $V(y_i | \mathbf{x}_i) = \mu_i + \alpha \mu_i^2$ , which is quadratic in the mean, is referred to as the NEGBIN2 model (Cameron and Trivedi 1986). To estimate this model, specify `DIST=NEGBIN(P=2)` in the `MODEL` statement. The Poisson distribution is a special case of the negative binomial distribution where  $\alpha = 0$ . A test of the Poisson distribution can be carried out by testing the hypothesis that  $\alpha = \frac{1}{\theta_i} = 0$ . A Wald test of this hypothesis is provided (it is the reported  $t$  statistic for the estimated  $\alpha$  in the negative binomial model).

The log-likelihood function of the negative binomial regression model (NEGBIN2) is given by

$$\begin{aligned} \mathcal{L} = \sum_{i=1}^N \left\{ \sum_{j=0}^{y_i-1} \ln(j + \alpha^{-1}) - \ln(y_i!) \right. \\ \left. - (y_i + \alpha^{-1}) \ln(1 + \alpha \exp(\mathbf{x}_i' \boldsymbol{\beta})) + y_i \ln(\alpha) + y_i \mathbf{x}_i' \boldsymbol{\beta} \right\} \end{aligned}$$

where use of the following fact is made if  $y$  is an integer:

$$\Gamma(y + a) / \Gamma(a) = \prod_{j=0}^{y-1} (j + a)$$

Cameron and Trivedi (1986) consider a general class of negative binomial models that have mean  $\mu_i$  and variance function  $\mu_i + \alpha \mu_i^p$ . The NEGBIN2 model, with  $p = 2$ , is the standard formulation of the negative binomial model. Models that have other values of  $p$ ,  $-\infty < p < \infty$ , have the same density  $f(y_i | \mathbf{x}_i)$ , except that  $\alpha^{-1}$  is replaced everywhere by  $\alpha^{-1} \mu_i^{2-p}$ . The negative binomial model NEGBIN1, which sets  $p = 1$ , has the variance function  $V(y_i | \mathbf{x}_i) = \mu_i + \alpha \mu_i$ , which is linear in the mean. To estimate this model, specify `DIST=NEGBIN(P=1)` in the `MODEL` statement.

The log-likelihood function of the NEGBIN1 regression model is given by

$$\begin{aligned} \mathcal{L} = \sum_{i=1}^N \left\{ \sum_{j=0}^{y_i-1} \ln(j + \alpha^{-1} \exp(\mathbf{x}_i' \boldsymbol{\beta})) \right. \\ \left. - \ln(y_i!) - (y_i + \alpha^{-1} \exp(\mathbf{x}_i' \boldsymbol{\beta})) \ln(1 + \alpha) + y_i \ln(\alpha) \right\} \end{aligned}$$

For more information about the negative binomial regression model, see the section “Negative Binomial Regression” (Chapter 11, *SAS/ETS User's Guide*).

## Zero-Inflated Count Regression Overview

The main motivation for using zero-inflated count models is that real-life data frequently display overdispersion and excess zeros. Zero-inflated count models provide a way to both model the excess zeros and allow for overdispersion. In particular, there are two possible data generation processes for each observation. The result of a Bernoulli trial is used to determine which of the two processes to use. For observation  $i$ , Process 1 is chosen with probability  $\varphi_i$  and Process 2 with probability  $1 - \varphi_i$ . Process 1 generates only zero counts. Process 2 generates counts from either a Poisson or a negative binomial model. In general,

$$y_i \sim \begin{cases} 0 & \text{with probability } \varphi_i \\ g(y_i) & \text{with probability } 1 - \varphi_i \end{cases}$$

Therefore, the probability of  $\{Y_i = y_i\}$  can be described as

$$\begin{aligned} P(y_i = 0 | \mathbf{x}_i) &= \varphi_i + (1 - \varphi_i)g(0) \\ P(y_i | \mathbf{x}_i) &= (1 - \varphi_i)g(y_i), \quad y_i > 0 \end{aligned}$$

where  $g(y_i)$  follows either the Poisson or the negative binomial distribution.

If the probability  $\varphi_i$  depends on the characteristics of observation  $i$ , then  $\varphi_i$  is written as a function of  $\mathbf{z}_i' \boldsymbol{\gamma}$ , where  $\mathbf{z}_i'$  is the  $1 \times (q + 1)$  vector of zero-inflated covariates and  $\boldsymbol{\gamma}$  is the  $(q + 1) \times 1$  vector of zero-inflated coefficients to be estimated. (The zero-inflated intercept is  $\gamma_0$ ; the coefficients for the  $q$  zero-inflated covariates are  $\gamma_1, \dots, \gamma_q$ .) The function  $F$  that relates the product  $\mathbf{z}_i' \boldsymbol{\gamma}$  (which is a scalar) to the probability  $\varphi_i$  is called the zero-inflated link function,

$$\varphi_i = F_i = F(\mathbf{z}_i' \boldsymbol{\gamma})$$

In the HPCOUNTREG procedure, the zero-inflated covariates are indicated in the ZEROMODEL statement. Furthermore, the zero-inflated link function  $F$  can be specified as either the logistic function,

$$F(\mathbf{z}_i' \boldsymbol{\gamma}) = \Lambda(\mathbf{z}_i' \boldsymbol{\gamma}) = \frac{\exp(\mathbf{z}_i' \boldsymbol{\gamma})}{1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma})}$$

or the standard normal cumulative distribution function (also called the probit function),

$$F(\mathbf{z}_i' \boldsymbol{\gamma}) = \Phi(\mathbf{z}_i' \boldsymbol{\gamma}) = \int_0^{\mathbf{z}_i' \boldsymbol{\gamma}} \frac{1}{\sqrt{2\pi}} \exp(-u^2/2) du$$

The zero-inflated link function is indicated by using the LINK= option in the ZEROMODEL statement. The default ZI link function is the logistic function.

## Zero-Inflated Poisson Regression

In the zero-inflated Poisson (ZIP) regression model, the data generation process that is referred to earlier as Process 2 is

$$g(y_i) = \frac{\exp(-\mu_i) \mu_i^{y_i}}{y_i!}$$

where  $\mu_i = e^{\mathbf{x}_i' \boldsymbol{\beta}}$ . Thus the ZIP model is defined as

$$\begin{aligned} P(y_i = 0 | \mathbf{x}_i, \mathbf{z}_i) &= F_i + (1 - F_i) \exp(-\mu_i) \\ P(y_i | \mathbf{x}_i, \mathbf{z}_i) &= (1 - F_i) \frac{\exp(-\mu_i) \mu_i^{y_i}}{y_i!}, \quad y_i > 0 \end{aligned}$$

The conditional expectation and conditional variance of  $y_i$  are given by

$$\begin{aligned} E(y_i | \mathbf{x}_i, \mathbf{z}_i) &= \mu_i (1 - F_i) \\ V(y_i | \mathbf{x}_i, \mathbf{z}_i) &= E(y_i | \mathbf{x}_i, \mathbf{z}_i) (1 + \mu_i F_i) \end{aligned}$$

Note that the ZIP model (in addition to the ZINB model) exhibits overdispersion because  $V(y_i | \mathbf{x}_i, \mathbf{z}_i) > E(y_i | \mathbf{x}_i, \mathbf{z}_i)$ .

In general, the log-likelihood function of the ZIP model is

$$\mathcal{L} = \sum_{i=1}^N \ln [P(y_i | \mathbf{x}_i, \mathbf{z}_i)]$$

After a specific link function (either logistic or standard normal) for the probability  $\varphi_i$  is chosen, it is possible to write the exact expressions for the log-likelihood function and the gradient.

### ZIP Model with Logistic Link Function

First, consider the ZIP model in which the probability  $\varphi_i$  is expressed by a logistic link function, namely

$$\varphi_i = \frac{\exp(\mathbf{z}_i' \boldsymbol{\gamma})}{1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma})}$$

The log-likelihood function is

$$\begin{aligned} \mathcal{L} &= \sum_{\{i: y_i=0\}} \ln [\exp(\mathbf{z}_i' \boldsymbol{\gamma}) + \exp(-\exp(\mathbf{x}_i' \boldsymbol{\beta}))] \\ &\quad + \sum_{\{i: y_i>0\}} \left[ y_i \mathbf{x}_i' \boldsymbol{\beta} - \exp(\mathbf{x}_i' \boldsymbol{\beta}) - \sum_{k=2}^{y_i} \ln(k) \right] \\ &\quad - \sum_{i=1}^N \ln [1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma})] \end{aligned}$$

### ZIP Model with Standard Normal Link Function

Next, consider the ZIP model in which the probability  $\varphi_i$  is expressed by a standard normal link function:  $\varphi_i = \Phi(\mathbf{z}_i' \boldsymbol{\gamma})$ . The log-likelihood function is

$$\begin{aligned} \mathcal{L} = & \sum_{\{i: y_i=0\}} \ln \{ \Phi(\mathbf{z}_i' \boldsymbol{\gamma}) + [1 - \Phi(\mathbf{z}_i' \boldsymbol{\gamma})] \exp(-\exp(\mathbf{x}_i' \boldsymbol{\beta})) \} \\ & + \sum_{\{i: y_i>0\}} \left\{ \ln [(1 - \Phi(\mathbf{z}_i' \boldsymbol{\gamma}))] - \exp(\mathbf{x}_i' \boldsymbol{\beta}) + y_i \mathbf{x}_i' \boldsymbol{\beta} - \sum_{k=2}^{y_i} \ln(k) \right\} \end{aligned}$$

For more information about the zero-inflated Poisson regression model, see the section “Zero-Inflated Poisson Regression” (Chapter 11, *SAS/ETS User's Guide*).

---

### Zero-Inflated Negative Binomial Regression

The zero-inflated negative binomial (ZINB) model in PROC HPCOUNTREG is based on the negative binomial model that has a quadratic variance function (when DIST=NEGBIN in the MODEL or PROC HPCOUNTREG statement). The ZINB model is obtained by specifying a negative binomial distribution for the data generation process referred to earlier as Process 2:

$$g(y_i) = \frac{\Gamma(y_i + \alpha^{-1})}{y_i! \Gamma(\alpha^{-1})} \left( \frac{\alpha^{-1}}{\alpha^{-1} + \mu_i} \right)^{\alpha^{-1}} \left( \frac{\mu_i}{\alpha^{-1} + \mu_i} \right)^{y_i}$$

Thus the ZINB model is defined to be

$$\begin{aligned} P(y_i = 0 | \mathbf{x}_i, \mathbf{z}_i) &= F_i + (1 - F_i) (1 + \alpha \mu_i)^{-\alpha^{-1}} \\ P(y_i | \mathbf{x}_i, \mathbf{z}_i) &= (1 - F_i) \frac{\Gamma(y_i + \alpha^{-1})}{y_i! \Gamma(\alpha^{-1})} \left( \frac{\alpha^{-1}}{\alpha^{-1} + \mu_i} \right)^{\alpha^{-1}} \\ &\quad \times \left( \frac{\mu_i}{\alpha^{-1} + \mu_i} \right)^{y_i}, \quad y_i > 0 \end{aligned}$$

In this case, the conditional expectation ( $E$ ) and conditional variance ( $V$ ) of  $y_i$  are

$$E(y_i | \mathbf{x}_i, \mathbf{z}_i) = \mu_i (1 - F_i)$$

$$V(y_i | \mathbf{x}_i, \mathbf{z}_i) = E(y_i | \mathbf{x}_i, \mathbf{z}_i) [1 + \mu_i (F_i + \alpha)]$$

Like the ZIP model, the ZINB model exhibits overdispersion because the conditional variance exceeds the conditional mean.

### ZINB Model with Logistic Link Function

In this model, the probability  $\varphi_i$  is given by the logistic function, namely

$$\varphi_i = \frac{\exp(\mathbf{z}_i' \boldsymbol{\gamma})}{1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma})}$$

The log-likelihood function is

$$\begin{aligned} \mathcal{L} = & \sum_{\{i: y_i=0\}} \ln \left[ \exp(\mathbf{z}_i' \boldsymbol{\gamma}) + (1 + \alpha \exp(\mathbf{x}_i' \boldsymbol{\beta}))^{-\alpha^{-1}} \right] \\ & + \sum_{\{i: y_i > 0\}} \sum_{j=0}^{y_i-1} \ln(j + \alpha^{-1}) \\ & + \sum_{\{i: y_i > 0\}} \{ -\ln(y_i!) - (y_i + \alpha^{-1}) \ln(1 + \alpha \exp(\mathbf{x}_i' \boldsymbol{\beta})) + y_i \ln(\alpha) + y_i \mathbf{x}_i' \boldsymbol{\beta} \} \\ & - \sum_{i=1}^N \ln [1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma})] \end{aligned}$$

### ZINB Model with Standard Normal Link Function

For this model, the probability  $\varphi_i$  is expressed by the standard normal distribution function (probit function):  $\varphi_i = \Phi(\mathbf{z}_i' \boldsymbol{\gamma})$ . The log-likelihood function is

$$\begin{aligned} \mathcal{L} = & \sum_{\{i: y_i=0\}} \ln \left\{ \Phi(\mathbf{z}_i' \boldsymbol{\gamma}) + [1 - \Phi(\mathbf{z}_i' \boldsymbol{\gamma})] (1 + \alpha \exp(\mathbf{x}_i' \boldsymbol{\beta}))^{-\alpha^{-1}} \right\} \\ & + \sum_{\{i: y_i > 0\}} \ln [1 - \Phi(\mathbf{z}_i' \boldsymbol{\gamma})] \\ & + \sum_{\{i: y_i > 0\}} \sum_{j=0}^{y_i-1} \{ \ln(j + \alpha^{-1}) \} \\ & - \sum_{\{i: y_i > 0\}} \ln(y_i!) \\ & - \sum_{\{i: y_i > 0\}} (y_i + \alpha^{-1}) \ln(1 + \alpha \exp(\mathbf{x}_i' \boldsymbol{\beta})) \\ & + \sum_{\{i: y_i > 0\}} y_i \ln(\alpha) \\ & + \sum_{\{i: y_i > 0\}} y_i \mathbf{x}_i' \boldsymbol{\beta} \end{aligned}$$

For more information about the zero-inflated negative binomial regression model, see the section “Zero-Inflated Negative Binomial Regression” (Chapter 11, *SAS/ETS User's Guide*).

## Computational Resources

The time and memory that PROC HPCOUNTREG requires are proportional to the number of parameters in the model and the number of observations in the data set being analyzed. Less time and memory are required for smaller models and fewer observations. When PROC HPCOUNTREG is run in the high-performance distributed environment, the amount of time required is also affected by the number of nodes and the number of threads per node as specified in the PERFORMANCE statement.

The method that is chosen to calculate the variance-covariance matrix and the optimization method also affect the time and memory resources. All optimization methods available through the METHOD= option have similar memory use requirements. The processing time might differ for each method, depending on the number of iterations and functional calls needed. The data set is read into memory to save processing time. If not enough memory is available to hold the data, the HPCOUNTREG procedure stores the data in a utility file on disk and rereads the data as needed from this file, substantially increasing the execution time of the procedure. The gradient and the variance-covariance matrix must be held in memory. If the model has  $p$  parameters including the intercept, then at least  $8 * (p + p * (p + 1)/2)$  bytes of memory are needed. The processing time is also a function of the number of iterations needed to converge to a solution for the model parameters. The number of iterations that are needed cannot be known in advance. You can use the MAXITER= option to limit the number of iterations that PROC HPCOUNTREG executes. You can alter the convergence criteria by using the nonlinear optimization options available in the PROC HPCOUNTREG statement. For a list of all the nonlinear optimization options, see “[Optimization Control Options](#)” on page 135.

---

## Covariance Matrix Types

The COVEST= option in the PROC HPCOUNTREG statement enables you to specify the estimation method for the covariance matrix. COVEST=HESSIAN estimates the covariance matrix that is based on the inverse of the Hessian matrix; COVEST=OP uses the outer product of gradients; and COVEST=QML produces the covariance matrix that is based on both the Hessian and outer product matrices. Although all three methods produce asymptotically equivalent results, they differ in computational intensity and produce results that might differ in finite samples. The COVEST=OP option provides the covariance matrix that is typically the easiest to compute. In some cases, the OP approximation is considered more efficient than the Hessian or QML approximation because it contains fewer random elements. The QML approximation is computationally the most complex because it requires both the outer product of gradients and the Hessian matrix. In most cases, the OP or Hessian approximation is preferred to QML. The need for QML approximation arises in cases where the model is misspecified and the information matrix equality does not hold. The default is COVEST=HESSIAN.

---

## Displayed Output

PROC HPCOUNTREG produces the following displayed output.

## Model Fit Summary

The “Model Fit Summary” table contains the following information:

- dependent (count) variable name
- number of observations used
- number of missing values in data set, if any
- data set name
- type of model that was fit
- offset variable name, if any
- zero-inflated link function, if any
- zero-inflated offset variable name, if any
- log-likelihood value at solution
- maximum absolute gradient at solution
- number of iterations
- AIC value at solution (smaller value indicates better fit)
- SBC value at solution (smaller value indicates better fit)

A line in the “Model Fit Summary” table indicates whether the algorithm successfully converged.

## Parameter Estimates

The “Parameter Estimates” table in the displayed output gives the estimates for the ZI intercept and ZI explanatory variables; they are labeled with the prefix “Inf\_”. For example, the ZI intercept is labeled “Inf\_intercept”. If you specify “Age” (a variable in your data set) as a ZI explanatory variable, then the “Parameter Estimates” table labels the corresponding parameter estimate “Inf\_Age”. If you do not list any ZI explanatory variables (for the ZI option VAR=), then only the intercept term is estimated.

“\_Alpha” is the negative binomial dispersion parameter. The  $t$  statistic that is given for “\_Alpha” is a test of overdispersion.

## Covariance of Parameter Estimates

If you specify the COVB option in the PROC HPCOUNTREG or MODEL statement, the HPCOUNTREG procedure displays the estimated covariance matrix, which is defined as the inverse of the information matrix at the final iteration.

## Correlation of Parameter Estimates

If you specify the CORRB option in the PROC HPCOUNTREG or MODEL statement, the HPCOUNTREG procedure displays the estimated correlation matrix, which is based on the Hessian matrix used at the final iteration.

---

## OUTPUT OUT= Data Set

The OUTPUT statement creates a new SAS data set that contains various estimates that you specify. You can request that the output data set contain the estimates of  $\mathbf{x}_i' \boldsymbol{\beta}$ , the expected value of the response variable, and the probability that the response variable will take the current value. Furthermore, if a zero-inflated model is fit, you can request that the output data set contain the estimates of  $\mathbf{z}_i' \boldsymbol{\gamma}$  and the probability that the response is 0 as a result of the zero-generating process. These statistics can be computed for all observations in which the regressors are not missing, even if the response is missing. By adding observations with missing response values to the input data set, you can compute these statistics for new observations or for settings of the regressors that are not present in the data without affecting the model fit. Because of potential space limitations on the client workstation, the data set that is created by the OUTPUT statement does not contain the variables in the input data set.

---

## OUTEST= Data Set

The OUTEST= data set is made up of at least two rows: the first row (with `_TYPE_='PARM'`) contains each of the parameter estimates in the model, and the second row (with `_TYPE_='STD'`) contains the standard errors for the parameter estimates in the model.

If you use the COVOUT option in the PROC HPCOUNTREG statement, the OUTEST= data set also contains the covariance matrix for the parameter estimates. The covariance matrix appears in the observations with `_TYPE_='COV'`, and the `_NAME_` variable labels the rows with the parameter names.

---

## ODS Table Names

PROC HPCOUNTREG assigns a name to each table that it creates. You can use these names to denote the table when you use the Output Delivery System (ODS) to select tables and create output data sets. These table names are listed in [Table 6.2](#).

**Table 6.2** ODS Tables Produced in PROC HPCOUNTREG

ODS Table Name	Description	Option
<b>ODS Tables Created by the MODEL Statement</b>		
FitSummary	Summary of nonlinear estimation	Default
ConvergenceStatus	Convergence status	Default
ParameterEstimates	Parameter estimates	Default
CovB	Covariance of parameter estimates	COVB
CorrB	Correlation of parameter estimates	CORRB



---

## Examples: The HPCOUNTREG Procedure

---

### Example 6.1: High-Performance Zero-Inflated Poisson Model

This example shows the use of the HPCOUNTREG procedure with an emphasis on large data set processing and the performance improvements that are achieved by executing in the high-performance distributed environment.

The following DATA step generates one million replicates from the zero-inflated Poisson (ZIP) model. The model contains seven variables and three variables that correspond to the zero-inflated process.

```
data simulate;
  call streaminit(12345);
  array vars x1-x7;
  array zero_vars z1-z3;

  array parms{7} (.3 .4 .2 .4 -.3 -.5 -.3);
  array zero_parms{3} (-.6 .3 .2);

  intercept=2;
  z_intercept=-1;
  theta=0.5;

  do i=1 to 1000000;
    sum_xb=0;
    sum_gz=0;
    do j=1 to 7;
      vars[j]=rand('NORMAL',0,1);
      sum_xb=sum_xb+parms[j]*vars[j];
    end;
    mu=exp(intercept+sum_xb);
    y_p=rand('POISSON', mu);

    do j=1 to 3;
      zero_vars[j]=rand('NORMAL',0,1);
      sum_gz = sum_gz+zero_parms[j]*zero_vars[j];
    end;
    z_gamma = z_intercept+sum_gz;
    pzero = cdf('LOGISTIC',z_gamma);
    cut=rand('UNIFORM');
    if cut<pzero then y_p=0;
    output;
  end;
  keep y_p x1-x7 z1-z3;
run;
```

The following statements estimate a zero-inflated Poisson model.

```

option set=GRIDHOST("&GRIDHOST");
option set=GRIDINSTALLLOC("&GRIDINSTALLLOC");

proc hpcountreg data=simulate dist=zip;
  performance nthreads=2 nodes=1 details
    host=&GRIDHOST install=&GRIDINSTALLLOC;
  model y_p=x1-x7;
  zeromodel y_p ~ z1-z3;
run;

```

The model is executed in the distributed computing environment on two threads and only one node. These settings are used to obtain a hypothetical environment that might resemble running the HPCOUNTREG procedure on a desktop workstation with a dual-core CPU. To run these statements successfully, you need to set the macro variables GRIDHOST and GRIDINSTALLLOC to resolve to appropriate values, or you can replace the references to the macro variables in the example with the appropriate values. [Output 6.1.1](#) shows the “Performance Information” table for this hypothetical scenario.

**Output 6.1.1** Performance Information with One Node and One Thread

Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	1
Number of Threads per Node	2

[Output 6.1.2](#) shows the results for the zero-inflated Poisson model. The “Model Fit Summary” table shows detailed information about the model and indicates that all one million observations were used to fit the model. All parameter estimates in the “Parameter Estimates” table are highly significant and correspond to their theoretical values set during the data generating process. The optimization of the model that contains one million observations took 41.88 seconds.

**Output 6.1.2** Zero-Inflated Poisson Model Execution on One Node and Two Threads

Model Fit Summary	
Dependent Variable	Y_P
Number of Observations	1000000
Data Set	WORK.SIMULATE
Model	ZIP
ZI Link Function	Logistic
Log Likelihood	-2215238
Maximum Absolute Gradient	2.0586E-8
Number of Iterations	7
Optimization Method	Newton-Raphson
AIC	4430500
SBC	4430642

**Output 6.1.2** *continued*

Convergence criterion (FCONV=2.220446E-16) satisfied.					
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	2.0005	0.000492	4069.80	<.0001
x1	1	0.2995	0.000352	850.17	<.0001
x2	1	0.3998	0.000353	1132.23	<.0001
x3	1	0.2008	0.000352	570.27	<.0001
x4	1	0.3994	0.000353	1132.85	<.0001
x5	1	-0.2995	0.000353	-848.95	<.0001
x6	1	-0.5000	0.000353	-1414.9	<.0001
x7	1	-0.3002	0.000352	-852.14	<.0001
Inf_Intercept	1	-0.9993	0.002521	-396.45	<.0001
Inf_z1	1	-0.6024	0.002585	-233.02	<.0001
Inf_z2	1	0.2976	0.002454	121.25	<.0001
Inf_z3	1	0.1974	0.002430	81.20	<.0001
Procedure Task Timing					
Task			Seconds	Percent	
Reading and Levelizing Data			0.40	0.97%	
Communication to Client			0.07	0.17%	
Optimization			41.29	98.57%	
Post-Optimization			0.12	0.29%	

In the following statements, the PERFORMANCE statement is modified to use a grid with 10 nodes, with each node capable of spawning eight threads:

```
proc hpcountreg data=simulate dist=zip;
  performance nthreads=8 nodes=10 details
               host="%GRIDHOST" install="%GRIDINSTALLLOC";
  model y_p=x1-x7;
  zeromodel y_p ~ z1-z3;
run;
```

Because the two models being estimated are identical, it is reasonable to expect that [Output 6.1.2](#) and [Output 6.1.3](#) would show the same results. However, you can see a significant difference in performance between the two models. The second model, which was run on a grid that used 10 nodes with eight threads each, took only 1.89 seconds instead of 41.88 seconds to optimize.

In certain circumstances, you might observe slight numerical differences in the results, depending on the number of nodes and threads involved. This happens because the order in which partial results are accumulated can make a difference in the final result, owing to the limits of numerical precision and the propagation of error in numerical computations.

**Output 6.1.3** Zero-Inflated Poisson Model Execution on 10 Nodes with Eight Threads Each

## The HPCOUNTREG Procedure

## Model Fit Summary

Dependent Variable	Y_P
Number of Observations	1000000
Data Set	WORK.SIMULATE
Model	ZIP
ZI Link Function	Logistic
Log Likelihood	-2215238
Maximum Absolute Gradient	2.0608E-8
Number of Iterations	7
Optimization Method	Newton-Raphson
AIC	4430500
SBC	4430642

Convergence criterion (FCONV=2.220446E-16) satisfied.

## Parameter Estimates

Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	2.0005	0.000492	4069.80	<.0001
x1	1	0.2995	0.000352	850.17	<.0001
x2	1	0.3998	0.000353	1132.23	<.0001
x3	1	0.2008	0.000352	570.27	<.0001
x4	1	0.3994	0.000353	1132.85	<.0001
x5	1	-0.2995	0.000353	-848.95	<.0001
x6	1	-0.5000	0.000353	-1414.9	<.0001
x7	1	-0.3002	0.000352	-852.14	<.0001
Inf_Intercept	1	-0.9993	0.002521	-396.45	<.0001
Inf_z1	1	-0.6024	0.002585	-233.02	<.0001
Inf_z2	1	0.2976	0.002454	121.25	<.0001
Inf_z3	1	0.1974	0.002430	81.20	<.0001

## Procedure Task Timing

Task	Seconds	Percent
Reading and Levelizing Data	0.02	0.98%
Communication to Client	0.05	2.76%
Optimization	1.70	89.87%
Post-Optimization	0.12	6.39%

As this example suggests, increasing the number of nodes and the number of threads per node improves performance significantly. When you use the parallelism afforded by a high-performance distributed environment, you can see an even more dramatic reduction in the time required for the optimization as the number of observations in the data set increases. When the data set is extremely large, the computations might not even be possible in some cases, given the typical memory resources and computational constraints of a desktop computer. Under such circumstances the high-performance distributed environment becomes a necessity.

---

## References

- Cameron, A. C. and Trivedi, P. K. (1986), “Econometric Models Based on Count Data: Comparisons and Applications of Some Estimators and Some Tests,” *Journal of Applied Econometrics*, 1, 29–53.
- Cameron, A. C. and Trivedi, P. K. (1998), *Regression Analysis of Count Data*, Cambridge: Cambridge University Press.
- LaMotte, L. R. (1994), “A Note on the Role of Independence in  $t$  Statistics Constructed from Linear Statistics in Regression Models,” *The American Statistician*, 48, 238–240.
- Long, J. S. (1997), *Regression Models for Categorical and Limited Dependent Variables*, Thousand Oaks, CA: Sage Publications.



## Chapter 7

# The HPPANEL Procedure (Experimental)

### Contents

---

Overview: HPPANEL Procedure . . . . .	<b>160</b>
Getting Started: HPPANEL Procedure . . . . .	<b>161</b>
Syntax: HPPANEL Procedure . . . . .	<b>163</b>
Functional Summary . . . . .	163
PROC HPPANEL Statement . . . . .	164
ID Statement . . . . .	165
MODEL Statement . . . . .	165
OUTPUT Statement . . . . .	166
PERFORMANCE Statement . . . . .	167
RESTRICT Statement . . . . .	167
TEST Statement . . . . .	168
Details: HPPANEL Procedure . . . . .	<b>169</b>
Specifying the Input Data . . . . .	169
Specifying the Regression Model . . . . .	169
Specifying the Number of Nodes and Number of Threads . . . . .	169
Unbalanced Data . . . . .	170
One-Way Fixed-Effects Model . . . . .	170
Two-Way Fixed-Effects Model . . . . .	171
Balanced Panels . . . . .	172
Unbalanced Panels . . . . .	173
One-Way Random-Effects Model . . . . .	174
Two-Way Random-Effects Model . . . . .	174
Linear Hypothesis Testing . . . . .	176
Specification Tests . . . . .	177
OUTPUT OUT= Data Set . . . . .	177
OUTEST= Data Set . . . . .	177
Printed Output . . . . .	178
ODS Table Names . . . . .	179
Example: HPPANEL Procedure . . . . .	<b>180</b>
Example 7.1: One-Way Random-Effects High-Performance Model . . . . .	180
References . . . . .	<b>184</b>

---

---

## Overview: HPPANEL Procedure

The HPPANEL procedure is a high-performance version of the PANEL procedure in SAS/ETS software. Both procedures analyze a class of linear econometric models that commonly arise when time series and cross-sectional data are combined (pooled). This type of data on time series cross-sectional bases is often referred to as panel data. Typical examples of panel data include observations over time about households, countries, firms, trade, and so on. For example, in the case of survey data about household income, the panel is created by repeatedly surveying the same households in different time periods (years).

Unlike the PANEL procedure (which can be run only on an individual workstation), the HPPANEL procedure takes advantage of a computing environment that enables it to distribute the optimization task among one or more nodes. Running on one node is called single-machine, and running on more than one node is called distributed mode. In addition, each node (whether in single-machine mode or in distributed mode) can use one or more threads to carry out the optimization on its subset of the data. When several nodes are used and each node uses several threads to carry out its part of the work, the result is a highly parallel computation that provides a dramatic gain in performance.

**NOTE:** Distributed mode requires SAS High-Performance Econometrics.

You can use the HPPANEL procedure to read and write data in distributed form and perform analyses in distributed mode or in single-machine mode. For more information about how to affect the execution mode of SAS high-performance analytical procedures, see the section “[Processing Modes](#)” on page 12 in Chapter 3, “[Shared Concepts and Topics](#).”

The HPPANEL procedure is specifically designed to operate in the high-performance distributed mode. By default, PROC HPPANEL performs computations in multiple threads.

The panel data models can be grouped into several categories that depend on the structure of the error term. The HPPANEL procedure uses the following error structures and the corresponding methods to analyze data:

- one-way and two-way models
- fixed-effects and random-effects models

A one-way model depends only on the cross section to which the observation belongs. A two-way model depends on both the cross section and the time period to which the observation belongs.

Apart from the possible one-way or two-way nature of the effect, the other dimension of difference between the possible specifications is the nature of the cross-sectional or time-series effect. The models are referred to as fixed-effects models if the effects are nonrandom and as random-effects models otherwise.



If the effects are fixed, the models are essentially regression models that have dummy variables that correspond to the specified effects. For fixed-effects models, ordinary least squares (OLS) estimation is the best linear unbiased estimator. Random-effects models use a two-stage approach: In the first stage, variance components are calculated by using methods described by Fuller and Battese (1974); Wansbeek and Kapteyn (1989); Wallace and Hussain (1969); Nerlove (1971). In the second stage, variance components are used to standardize the data, and ordinary least squares (OLS) regression is performed.

---

## Getting Started: HPPANEL Procedure

The following statements use the cost function data from Greene (1990) to estimate the variance components model. The variable `Production` is the log of output in millions of kilowatt-hours, and the variable `Cost` is the log of cost in millions of dollars. See Greene (1990) for details.

```
data greene;
  input firm year production cost @@;
datalines;
1 1955    5.36598    1.14867  1 1960    6.03787    1.45185
1 1965    6.37673    1.52257  1 1970    6.93245    1.76627
2 1955    6.54535    1.35041  2 1960    6.69827    1.71109
2 1965    7.40245    2.09519  2 1970    7.82644    2.39480
3 1955    8.07153    2.94628  3 1960    8.47679    3.25967
... more lines ...
```

You decide to fit the following model to the data,

$$C_{it} = \text{Intercept} + \beta P_{it} + v_i + e_t + \epsilon_{it} \text{ for } i = 1, \dots, N \text{ and } t = 1, \dots, T$$

where  $C_{it}$  and  $P_{it}$  represent the cost and production; and  $v_i$ ,  $e_t$ , and  $\epsilon_{it}$  are the cross-sectional, time series, and error variance components, respectively.

If you assume that the time and cross-sectional effects are random, four possible estimators are left for the variance components. The following statements choose the Fuller-Battese method to fit this model:

```
proc hppanel data=greene;
  model cost = production / rantwo vcomp = fb;
  id firm year;
  performance nodes=0 nthreads=2;
run;
```

The output of the HPPANEL procedure is shown in [Output 7.1](#).

**Figure 7.1** Two-Way Random Effects Results

The HPPANEL Procedure					
Model Information					
Data Source		WORK.GREENE			
Response Variable		cost			
Model		RANTWO			
Variance Component		FULLER			
Execution Mode		Single-Machine			
Fit Statistics					
Sum of Squared Error				0.34808	
Degree of Freedom				22.00000	
Mean Squared Error				0.01582	
Root Mean Squared Error				0.12579	
R-Square				0.81362	
Variance Component Estimates					
Variance Component for Cross Sections				0.0469	
Variance Component for Time Series				0.00906	
Variance Component for Error				0.00875	
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	-2.99992	0.64778	-4.63	<.0001
production	1	0.74660	0.07618	9.80	<.0001

Printed first is the model description, which reports the method used for estimation and the method used for estimating error components. Printed next is the fit statistics table, and then the variance components estimates. Finally, the table of regression parameter estimates shows the estimates, standard errors, and  $t$  tests.

## Syntax: HPPANEL Procedure

The following statements are available in the HPPANEL procedure:

```

PROC HPPANEL options ;
  ID cross-section-id time-series-id ;
  MODEL response = regressors </options> ;
  RESTRICT equation1 < , equation2 ... > ;
  TEST equation < , equation2 ... > < /options> ;
  OUTPUT OUT=SAS-data-set < output-options > ;
  PERFORMANCE < performance-options > ;

```

The ID and MODEL statements are required.

The following sections provide a functional summary of statements and options, then describe the PROC HPPANEL statement, and then describe the other statements in alphabetical order.

## Functional Summary

Table 7.1 summarizes the statements and options that you can use in the HPPANEL procedure.

**Table 7.1** Functional Summary

Description	Statement	Option
<b>Data Set Options</b>		
Includes correlations in the OUTEST= data set	HPPANEL	CORROUT
Includes covariances in the OUTEST= data set	HPPANEL	COVOUT
Specifies the input data set	HPPANEL	DATA=
Specifies the name of an output SAS data set	OUTPUT	OUT=
Writes parameter estimates to an output data set	HPPANEL	OUTEST=
<b>Variable Role Options</b>		
Specifies the cross-sectional and time ID variables	ID	
<b>Printing Control Options</b>		
Prints correlations of the estimates	MODEL	CORRB
Prints covariances of the estimates	MODEL	COVB
Suppresses printed output	MODEL	NOPRINT
Prints fixed effects	MODEL	PRINTFIXED
Performs tests of linear hypotheses	TEST	
<b>Model Estimation Options</b>		
Requests the one-way fixed-effects model	MODEL	FIXONE
Requests the one-way fixed-effects model with respect to time	MODEL	FIXONETIME

Description	Statement	Option
Requests the two-way fixed-effects model	MODEL	FIXTWO
Suppresses the intercept term	MODEL	NOINT
Requests the one-way random-effects model	MODEL	RANONE
Requests the two-way random-effects model	MODEL	RANTWO
Specifies the method for the variance components estimator	MODEL	VCOMP=
Specifies linear equality restrictions on the parameters	RESTRICT	
Specifies which tests to perform	TEST	WALD, LM, LR

## PROC HPPANEL Statement

**PROC HPPANEL** *options* ;

The HPPANEL statement invokes the HPPANEL procedure.

You can specify the following *options*:

**DATA=SAS-data-set**

names the input data set. Only one observation is allowed for each cross section and time period. If you omit the DATA= option, PROC HPPANEL uses the most recently created SAS data set.

**OUTEST=SAS-data-set**

names an output data set to contain the parameter estimates. When the OUTEST= option is not specified, the OUTEST= data set is not created. For more information about the structure of the OUTEST= data set, see the section “[OUTEST= Data Set](#)” on page 177.

**OUTCOV**

**COVOUT**

writes the standard errors and covariance matrix of the parameter estimates to the OUTEST= data set. For more information, see the section “[OUTEST= Data Set](#)” on page 177.

**OUTCORR**

**CORROUT**

writes the correlation matrix of the parameter estimates to the OUTEST= data set. For more information, see the section “[OUTEST= Data Set](#)” on page 177.

In addition, you can specify any of the following MODEL statement options in the PROC HPPANEL statement: CORRB, COVB, FIXONE, FIXONETIME, FIXTWO, RANONE, RANTWO, NOINT, NOPRINT, PRINTFIXED, and VCOMP=. Specifying these options in the PROC HPPANEL statement is equivalent to specifying them in the MODEL statement. For a complete description of each of these options, see the section “[MODEL Statement](#)” on page 165.

---

## ID Statement

**ID** *cross-section-id time-series-id* ;

The ID statement specifies variables in the input data set that identify the cross section and the time period for each observation. The ID statement is required. Unlike the PANEL procedure, the HPPANEL procedure does not require the data set to be sorted.

---

## MODEL Statement

**MODEL** *response = regressors* < / *options* > ;

The MODEL statement specifies the regression model and the error structure that are assumed for the regression residuals. The *response* variable is regressed on the independent variables (*regressors*). You can specify only one MODEL statement and only one *response*.

The error structure is specified by the FIXONE, FIXONETIME, FIXTWO, RANONE, and RANTWO options.

You can specify the following *options* after a slash (/).

### **CORRB**

prints the matrix of estimated correlations between the parameter estimates.

### **COVB**

prints the matrix of estimated covariances between the parameter estimates.

### **FIXONE**

requests that a one-way fixed-effects model be estimated, where the one-way model corresponding to cross-sectional effects only.

### **FIXONETIME**

requests that a one-way fixed-effects model be estimated, where the one-way model corresponding to time effects only.

### **FIXTWO**

requests that a two-way fixed-effects model be estimated.

### **NOINT**

suppresses the intercept parameter from the model.

### **NOPRINT**

suppresses the normal printed output.

### **PRINTFIXED**

prints the fixed effects.

### **RANONE**

requests that a one-way random-effects model be estimated.

**RANTWO**

requests that a two-way random-effects model be estimated.

**VCOMP=FB | NL | WH | WK**

specifies the type of variance component estimator to use.

For more information about these estimators, see the sections “[One-Way Random-Effects Model](#)” on page 174 and “[Two-Way Random-Effects Model](#)” on page 174.

You can specify the following values:

<b>FB</b>	requests the Fuller-Battese estimator.
<b>WK</b>	requests the Wansbeek-Kapteyn estimator.
<b>WH</b>	requests the Wallace-Hussain estimator.
<b>NERLOVE</b>	requests the Nerlove estimator.

By default, VCOMP=WK for both balanced and unbalanced data.

---

## OUTPUT Statement

**OUTPUT OUT=SAS-data-set** < output-options > ;

The OUTPUT statement creates a new SAS data set to contain variables that are specified by the COPYVAR option, the cross-sectional ID (\_CSID\_), and the time period (\_TSID\_). This data set also contains the predicted value and the residual if they are specified by *output-options*. When the response values are missing for the observation, all output estimates except the residual are still computed as long as none of the explanatory variables are missing. You can specify only one OUTPUT statement.

You must specify the OUT= option:

**OUT=SAS-data-set**

names the output data set.

You can specify one or more of the following *output-options*:

**COPYVAR=(SAS-variable-names)**

**COPYVARS=(SAS-variable-names)**

adds SAS variables to the output data set.

**PREDICTED**

outputs estimates of predicted dependent variables.

**RESIDUAL**

outputs estimates of residuals.

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement specifies *performance-options* to control the multithreaded and distributed computing environment and requests detailed performance results of the HPPANEL procedure. You can also use the PERFORMANCE statement to control whether the HPPANEL procedure executes in single-machine or distributed mode. You can specify the following *performance-options*:

### DETAILS

requests a table that shows a timing breakdown of the procedure steps.

### NODES=*n*

specifies the number of nodes in the distributed computing environment, provided that the data are not processed alongside the database.

### NTHREADS=*n*

specifies the number of threads for analytic computations and overrides the SAS system option THREADS | NOTTHREADS. If you do not specify the NTHREADS= option, PROC HPPANEL creates one thread per CPU for the analytic computations.

The PERFORMANCE statement is documented further in the section “[PERFORMANCE Statement](#)” on page 39 in Chapter 3, “[Shared Concepts and Topics](#).”

## RESTRICT Statement

**RESTRICT** *equation1* < ,*equation2*... > ;

The RESTRICT statement specifies linear equality restrictions on the parameters in the MODEL statement. There can be as many unique restrictions as the number of parameters in the MODEL statement. Multiple RESTRICT statements are understood as joint restrictions on the model’s parameters.

Currently, PROC HPPANEL only supports linear equality restrictions. Restriction expressions can be composed only of algebraic operations that involve the addition symbol (+), subtraction symbol (–), and multiplication symbol (\*).

The following statements illustrate the use of the RESTRICT statement:

```
proc hppanel;
  id csid tsid;
  model y = x1 x2 x3;
  restrict x1 = 0, x2 * .5 + 2 * x3 = 0;
  restrict x2 = 0, intercept = 0;
run;
```

A RESTRICT statement cannot include a division sign in its formulation. As in the preceding example, you can obtain restrictions on the intercept by using the keyword INTERCEPT.

## TEST Statement

**TEST** *equation1* <*,equation2...>* / *options* ;

The TEST statement performs Wald, LaGrange multiplier, and likelihood ratio tests of linear hypotheses about the regression parameters in the MODEL statement. Each *equation* specifies a linear hypothesis to be tested. Currently, only linear equality restrictions and tests are permitted in PROC HPPANEL. Test expressions can be composed only of algebraic operations that involve the addition symbol (+), subtraction symbol (−), and multiplication symbol (\*). All hypotheses in one TEST statement are tested jointly. Variable names in the equations must correspond to regressors in the preceding MODEL statement, and each name represents the coefficient of the corresponding regressor. In the equality restrictions, you can use the keyword INTERCEPT to refer to the coefficient of the intercept.

You can specify the following *options* after the slash (/):

### ALL

specifies Wald, LaGrange multiplier, and likelihood ratio tests.

### WALD

specifies the Wald test.

### LM

specifies the LaGrange multiplier test.

### LR

specifies the likelihood ratio test.

By default, the Wald test is performed.

The following statements illustrate the use of the TEST statement:

```
proc hppanel;
  id csid tsid;
  model y = x1 x2 x3;
  test x1 = 0, x2 * .5 + 2 * x3 = 0;
  test intercept = 0, x3 = 0;
run;
```

The first test investigates the joint hypothesis that

$$\beta_1 = 0$$

and

$$0.5\beta_2 + 2\beta_3 = 0$$



---

## Details: HPPANEL Procedure

---

### Specifying the Input Data

The HPPANEL procedure is similar to other regression procedures in SAS. Suppose you want to regress the variable *Y* on regressors *X1* and *X2*. Cross sections are identified by the variable *State*, and time periods are identified by the variable *Date*. Unlike the PANEL procedure, the HPPANEL procedure does not require the data set to be sorted. To invoke the HPPANEL procedure, you must specify the cross section and time series variables in an ID statement. The following statements show the correct syntax:

```
proc hppanel data=a;
  id state date;
  model y = x1 x2;
  performance nodes=2 nthreads=4;
run;
```

---

### Specifying the Regression Model

The MODEL statement in PROC HPPANEL is specified like the MODEL statement in other SAS regression procedures: the dependent variable is listed first, followed by an equal sign, followed by the list of regressor variables, as shown in the following statements:

```
proc hppanel data=a;
  id state date;
  model y = x1 x2;
  performance nodes=2 nthreads=4;
run;
```

---

### Specifying the Number of Nodes and Number of Threads

The PERFORMANCE statement in PROC HPPANEL is specified like the PERFORMANCE statement in other SAS high-performance procedures. The following statements execute the model in the distributed computing environment with two threads and four nodes:

```
proc hppanel data=a;
  id state date;
  model y = x1 x2;
  performance nodes=2 nthreads=4;
run;
```

The major advantage of using PROC HPPANEL is that you can incorporate a model for the structure of the random errors. It is important to consider what type of error structure model is appropriate for your data and to specify the corresponding option in the MODEL statement.

The error structure options supported by the HPPANEL procedure are FIXONE, FIXONETIME, FIXTWO, RANONE, and RANTWO. For more information about these methods and the error structures they assume, see the following sections. The following statements fit a Fuller-Battese one-way random-effects model:

```
proc hppanel data=a;
  id state date;
  model y = x1 x2 / ranone vcomp=fb;
  performance nodes=0 nthreads=1;
run;
```

To aid in model specification within this class of models, PROC HPPANEL provides one specification test statistic, the Hausman  $m$  statistic, which provides information about the appropriateness of the random-effects specification. The  $m$  statistic is based on the idea that, under the null hypothesis of no correlation between the effects variables and the regressors, ordinary least squares (OLS) and generalized least squares (GLS) are consistent. However, OLS is inefficient. Hence, a test can be based on the result that the covariance between an efficient estimator and its difference from an inefficient estimator is 0. Rejection of the null hypothesis might suggest that the fixed-effects model is more appropriate.

The HPPANEL procedure also provides the Buse R-square measure. This number is interpreted as a measure of the proportion of the transformed sum of squares of the dependent variable that is attributable to the influence of the independent variables. For OLS estimation, the Buse R-square measure is equivalent to the usual R-square measure.

---

## Unbalanced Data

The HPPANEL procedure can process data that have different numbers of time series observations across different cross sections. The missing time series observations are recognized by the absence of time series ID variable values in some of the cross sections in the input data set. Moreover, if an observation that has a particular time series ID value and cross-sectional ID value is present in the input data set but one or more of the model variables are missing, that time series point is treated as missing for that cross section.

---

## One-Way Fixed-Effects Model

The specification for the one-way fixed-effects model is

$$u_{it} = \gamma_i + \epsilon_{it}$$

where the  $\gamma_i$  are nonrandom parameters to be estimated.

Let  $\mathbf{Q}_0 = \text{diag}(\mathbf{E}_{T_i})$ , with  $\bar{\mathbf{J}}_{T_i} = \mathbf{J}_{T_i} / T_i$  and  $\mathbf{E}_{T_i} = \mathbf{I}_{T_i} - \bar{\mathbf{J}}_{T_i}$ , where  $\mathbf{J}_{T_i}$  is a matrix of  $T_i$  ones.

The matrix  $\mathbf{Q}_0$  represents the within transformation. In the one-way model, the within transformation is the conversion of the raw data to deviations from a cross section's mean. The vector  $\tilde{\mathbf{x}}_{it}$  is a row of the general matrix  $\mathbf{X}_s$ , where the subscripted  $s$  implies that the constant (column of ones) is missing.

Let  $\tilde{\mathbf{X}}_s = \mathbf{Q}_0 \mathbf{X}_s$  and  $\tilde{\mathbf{y}} = \mathbf{Q}_0 \mathbf{y}$ . The estimator of the slope coefficients is given by

$$\tilde{\beta}_s = (\tilde{\mathbf{X}}_s' \tilde{\mathbf{X}}_s)^{-1} \tilde{\mathbf{X}}_s' \tilde{\mathbf{y}}$$

After the slope estimates have been calculated, the estimation of an intercept or the cross-sectional fixed effects is handled as follows. First, you obtain the cross-sectional effects:

$$\gamma_i = \bar{y}_{i\cdot} - \tilde{\beta}_s \bar{x}_{i\cdot} \quad \text{for } i = 1 \dots N$$

If the NOINT option is specified, then the dummy variables' coefficients are set equal to the fixed effects. If you want an intercept, then the  $i$ th dummy variable is obtained from the following expression:

$$D_i = \gamma_i - \gamma_N \quad \text{for } i = 1 \dots N - 1$$

The intercept is the  $N$ th fixed effect  $\gamma_N$ .

The within-model sum of squared errors is

$$\text{SSE} = \sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \gamma_i - \mathbf{X}_s \tilde{\beta}_s)^2$$

The estimated error variance can be written as

$$\hat{\sigma}_\epsilon^2 = \text{SSE} / (M - N - (K - 1))$$

Alternatively, an equivalent way to express the error variance is

$$\hat{\sigma}_\epsilon^2 = \tilde{\mathbf{u}}' \mathbf{Q}_0 \tilde{\mathbf{u}} / (M - N - (K - 1))$$

where the residuals  $\tilde{\mathbf{u}}$  are given by  $\tilde{\mathbf{u}} = (\mathbf{I}_M - \mathbf{j}_M \mathbf{j}_M' / M)(\mathbf{y} - \mathbf{X}_s \tilde{\beta}_s)$  if there is an intercept and by  $\tilde{\mathbf{u}} = (\mathbf{y} - \mathbf{X}_s \tilde{\beta}_s)$  if there is not. The drawback is that the formula changes (but the results do not) with the inclusion of a constant.

The variance covariance matrix of  $\tilde{\beta}_s$  is given by

$$\text{Var}[\tilde{\beta}_s] = \hat{\sigma}_\epsilon^2 (\tilde{\mathbf{X}}_s' \tilde{\mathbf{X}}_s)^{-1}$$

The covariance of the dummy variables and the dummy variables with the  $\tilde{\beta}_s$  depends on whether the intercept is included in the model. For more information, see the section “One-Way Fixed-Effects Model” (Chapter 20, *SAS/ETS User's Guide*).

Alternatively, the FIXONETIME model option estimates a one-way model in which the heterogeneity comes from time effects. This option is analogous to re-sorting the data by time and then by cross section, and then running a FIXONE model. The advantage of using the FIXONETIME option is that sorting is avoided and the model remains labeled correctly.

## Two-Way Fixed-Effects Model

The specification for the two-way fixed-effects model is

$$u_{it} = \gamma_i + \alpha_t + \epsilon_{it}$$

where the  $\gamma_i$  and  $\alpha_t$  are nonrandom parameters to be estimated.

If you do not specify the NOINT option (which suppresses the intercept) in the MODEL statement, the estimates for the fixed effects are reported under the restriction that  $\gamma_N = 0$  and  $\alpha_T = 0$ . If you specify the NOINT option to suppress the intercept, only the restriction  $\alpha_T = 0$  is imposed.

## Balanced Panels

Assume that the data are balanced (for example, all cross sections have  $T$  observations). Then you can write

$$\tilde{y}_{it} = y_{it} - \bar{y}_{i\cdot} - \bar{y}_{\cdot t} + \bar{\bar{y}}$$

$$\tilde{\mathbf{x}}_{it} = \mathbf{x}_{it} - \bar{\mathbf{x}}_{i\cdot} - \bar{\mathbf{x}}_{\cdot t} + \bar{\bar{\mathbf{x}}}$$

where the symbols are as follows:

- $y_{it}$  and  $\mathbf{x}_{it}$  are the dependent variable (a scalar) and the explanatory variables (a vector whose columns are the explanatory variables, not including a constant), respectively
- $\bar{y}_{i\cdot}$  and  $\bar{\mathbf{x}}_{i\cdot}$  are cross section means
- $\bar{y}_{\cdot t}$  and  $\bar{\mathbf{x}}_{\cdot t}$  are time means
- $\bar{\bar{y}}$  and  $\bar{\bar{\mathbf{x}}}$  are the overall means

The two-way fixed-effects model is simply a regression of  $\tilde{y}_{it}$  on  $\tilde{\mathbf{x}}_{it}$ . Therefore, the two-way  $\beta$  is given by

$$\tilde{\beta}_s = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \tilde{\mathbf{y}}$$

The following calculations of cross-sectional dummy variables, time dummy variables, and intercepts are similar to how they are calculated in the one-way model:

First, you obtain the net cross-sectional and time effects. Denote the cross-sectional effects by  $\gamma$  and the time effects by  $\alpha$ . These effects are calculated from the following relations:

$$\hat{\gamma}_i = (\bar{y}_{i\cdot} - \bar{\bar{y}}) - \tilde{\beta}_s (\bar{x}_{i\cdot} - \bar{\bar{x}})$$

$$\hat{\alpha}_t = (\bar{y}_{\cdot t} - \bar{\bar{y}}) - \tilde{\beta}_s (\bar{x}_{\cdot t} - \bar{\bar{x}})$$

Use the superscript C and T to denote the cross-sectional dummy variables and time dummy variables, respectively. Under the NOINT option, the following equations produce the dummy variables:

$$D_i^C = \hat{\gamma}_i + \hat{\alpha}_T$$

$$D_t^T = \hat{\alpha}_t - \hat{\alpha}_T$$

When an intercept is specified, the equations for dummy variables and intercept are

$$D_i^C = \hat{\gamma}_i - \hat{\gamma}_N$$

$$D_t^T = \hat{\alpha}_t - \hat{\alpha}_T$$

$$\text{Intercept} = \hat{\gamma}_N + \hat{\alpha}_T$$

The sum of squared errors is

$$\text{SSE} = \sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \gamma_i - \alpha_t - \mathbf{X}_s \tilde{\beta}_s)^2$$

The estimated error variance is

$$\hat{\sigma}_\epsilon^2 = \text{SSE}/(M - N - T - (K - 1))$$

With or without a constant, the covariance matrix of  $\tilde{\beta}_s$  is given by

$$\text{Var}[\tilde{\beta}_s] = \hat{\sigma}_\epsilon^2 (\tilde{\mathbf{X}}_s' \tilde{\mathbf{X}}_s)^{-1}$$

For information about the covariance matrix that is related to dummy variables, see the section “Two-Way Fixed-Effects Model” (Chapter 20, *SAS/ETS User's Guide*).

## Unbalanced Panels

Let  $\mathbf{X}_*$  and  $\mathbf{y}_*$  be the independent and dependent variables, respectively, that are arranged by time and by cross section within each time period. (Note that the input data set that the PANEL procedure uses must be sorted by cross section and then by time within each cross section.) Let  $M_t$  be the number of cross sections that are observed in year  $t$ , and let  $\sum_t M_t = M$ . Let  $\mathbf{D}_t$  be the  $M_t \times N$  matrix that is obtained from the  $N \times N$  identity matrix from which rows that correspond to cross sections that are not observed at time  $t$  have been omitted. Consider

$$\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2)$$

where  $\mathbf{Z}_1 = (\mathbf{D}_1', \mathbf{D}_2', \dots, \mathbf{D}_T')'$  and  $\mathbf{Z}_2 = \text{diag}(\mathbf{D}_1 \mathbf{j}_N, \mathbf{D}_2 \mathbf{j}_N, \dots, \mathbf{D}_T \mathbf{j}_N)$ . The matrix  $\mathbf{Z}$  contains the dummy variable structure for the two-way model.

Let

$$\begin{aligned} \Delta_N &= \mathbf{Z}_1' \mathbf{Z}_1 \\ \Delta_T &= \mathbf{Z}_2' \mathbf{Z}_2 \\ \mathbf{A} &= \mathbf{Z}_2' \mathbf{Z}_1 \\ \bar{\mathbf{Z}} &= \mathbf{Z}_2 - \mathbf{Z}_1 \Delta_N^{-1} \mathbf{A}' \\ \mathbf{Q} &= \Delta_T - \mathbf{A} \Delta_N^{-1} \mathbf{A}' \\ \mathbf{P} &= (\mathbf{I}_M - \mathbf{Z}_1 \Delta_N^{-1} \mathbf{Z}_1') - \bar{\mathbf{Z}} \mathbf{Q}^{-1} \bar{\mathbf{Z}}' \end{aligned}$$

The estimate of the regression slope coefficients is given by

$$\tilde{\beta}_s = (\mathbf{X}_{*s}' \mathbf{P} \mathbf{X}_{*s})^{-1} \mathbf{X}_{*s}' \mathbf{P} \mathbf{y}_*$$

where  $\mathbf{X}_{*s}$  is the  $\mathbf{X}_*$  matrix without the vector of 1s.

The estimator of the error variance is

$$\hat{\sigma}_\epsilon^2 = \tilde{\mathbf{u}}' \mathbf{P} \tilde{\mathbf{u}} / (M - T - N + 1 - (K - 1))$$

where the residuals are given by  $\tilde{\mathbf{u}} = (\mathbf{I}_M - \mathbf{j}_M \mathbf{j}_M' / M)(\mathbf{y}_* - \mathbf{X}_{*s} \tilde{\beta}_s)$  if there is an intercept in the model and by  $\tilde{\mathbf{u}} = \mathbf{y}_* - \mathbf{X}_{*s} \tilde{\beta}_s$  if there is no intercept.

The actual implementation is quite different from the theory. For more information, see the section “Two-Way Fixed-Effects Model” (Chapter 20, *SAS/ETS User's Guide*).

## One-Way Random-Effects Model

The specification for the one-way random-effects model is

$$u_{it} = v_i + \epsilon_{it}$$

Let  $\mathbf{Z}_0 = \text{diag}(\mathbf{J}_{T_i})$ ,  $\mathbf{P}_0 = \text{diag}(\bar{\mathbf{J}}_{T_i})$ , and  $\mathbf{Q}_0 = \text{diag}(\mathbf{E}_{T_i})$ , with  $\bar{\mathbf{J}}_{T_i} = \mathbf{J}_{T_i}/T_i$  and  $\mathbf{E}_{T_i} = \mathbf{I}_{T_i} - \bar{\mathbf{J}}_{T_i}$ . Define  $\tilde{\mathbf{X}}_s = \mathbf{Q}_0 \mathbf{X}_s$ . Also define  $\tilde{\mathbf{y}} = \mathbf{Q}_0 \mathbf{y}$  and  $\mathbf{J}$  as a vector of 1s whose length is  $T_i$ .

In the one-way model, estimation proceeds in a two-step fashion. First, you obtain estimates of the variance of the  $\sigma_\epsilon^2$  and  $\sigma_v^2$ . There are multiple ways to derive these estimates; PROC HPPANEL provides four options. For more information, see the section “One-Way Random-Effects Model” (Chapter 20, *SAS/ETS User's Guide*).

After the variance components are calculated from any method, the next task is to estimate the regression model of interest. For each individual, you form a weight ( $\theta_i$ ),

$$\theta_i = 1 - \sigma_\epsilon/w_i$$

$$w_i^2 = T_i \sigma_v^2 + \sigma_\epsilon^2$$

where  $T_i$  is the  $i$ th cross section's time observations.

Taking the  $\theta_i$ , you form the partial deviations,

$$\tilde{y}_{it} = y_{it} - \theta_i \bar{y}_i.$$

$$\tilde{x}_{it} = x_{it} - \theta_i \bar{x}_i.$$

where  $\bar{y}_i$  and  $\bar{x}_i$  are cross section means of the dependent variable and independent variables (including the constant if any), respectively.

The random-effects  $\beta$  is then the result of simple OLS on the transformed data.

## Two-Way Random-Effects Model

The specification for the two-way random-effects model is

$$u_{it} = v_i + e_t + \epsilon_{it}$$

As it does for the one-way random-effects model, the HPPANEL procedure provides four options for variance component estimators. However, unbalanced panels present some special concerns that do not occur for one-way random-effects models.

Let  $\mathbf{X}_*$  and  $\mathbf{y}_*$  be the independent and dependent variables that are arranged by time and by cross section within each time period. (Note that the input data set that the PANEL procedure uses must be sorted by cross section and then by time within each cross section.) Let  $M_t$  be the number of cross sections that are observed in time  $t$ , and let  $\sum_t M_t = M$ . Let  $\mathbf{D}_t$  be the  $M_t \times N$  matrix that is obtained from the  $N \times N$  identity matrix from which rows that correspond to cross sections that are not observed at time  $t$  have been omitted. Consider

$$\mathbf{Z} = (\mathbf{Z}_1, \mathbf{Z}_2)$$

where  $\mathbf{Z}_1 = (\mathbf{D}'_1, \mathbf{D}'_2, \dots, \mathbf{D}'_T)'$  and  $\mathbf{Z}_2 = \text{diag}(\mathbf{D}_1\mathbf{j}_N, \mathbf{D}_2\mathbf{j}_N, \dots, \mathbf{D}_T\mathbf{j}_N)$ .

The matrix  $\mathbf{Z}$  contains the dummy variable structure for the two-way model.

For notational ease, let

$$\Delta_N = \mathbf{Z}'_1 \mathbf{Z}_1$$

$$\Delta_T = \mathbf{Z}'_2 \mathbf{Z}_2$$

$$\mathbf{A} = \mathbf{Z}'_2 \mathbf{Z}_1$$

$$\bar{\mathbf{Z}} = \mathbf{Z}_2 - \mathbf{Z}_1 \Delta_N^{-1} \mathbf{A}'$$

$$\bar{\Delta}_1 = \mathbf{I}_M - \mathbf{Z}_1 \Delta_N^{-1} \mathbf{Z}'_1$$

$$\bar{\Delta}_2 = \mathbf{I}_M - \mathbf{Z}_2 \Delta_T^{-1} \mathbf{Z}'_2$$

$$\mathbf{Q} = \Delta_T - \mathbf{A} \Delta_N^{-1} \mathbf{A}'$$

$$\mathbf{P} = (\mathbf{I}_M - \mathbf{Z}_1 \Delta_N^{-1} \mathbf{Z}'_1) - \bar{\mathbf{Z}} \mathbf{Q}^{-1} \bar{\mathbf{Z}}'$$

PROC HPPANEL provides four methods to estimate the variance components. For more information, see the section “Two-Way Random-Effects Model” (Chapter 20, *SAS/ETS User's Guide*).

After the estimates of the variance components are calculated, you can proceed to the final estimation. If the panel is balanced, partial mean deviations are used as follows

$$\tilde{y}_{it} = y_{it} - \theta_1 \bar{y}_{i\cdot} - \theta_2 \bar{y}_{\cdot t} + \theta_3 \bar{y}_{\cdot\cdot}$$

$$\tilde{x}_{it} = x_{it} - \theta_1 \bar{x}_{i\cdot} - \theta_2 \bar{x}_{\cdot t} + \theta_3 \bar{x}_{\cdot\cdot}$$

The  $\theta$  estimates are obtained from

$$\theta_1 = 1 - \frac{\sigma_\epsilon}{\sqrt{T\sigma_v^2 + \sigma_\epsilon^2}}$$

$$\theta_2 = 1 - \frac{\sigma_\epsilon}{\sqrt{N\sigma_e^2 + \sigma_\epsilon^2}}$$

$$\theta_3 = \theta_1 + \theta_2 + \frac{\sigma_\epsilon}{\sqrt{T\sigma_v^2 + N\sigma_e^2 + \sigma_\epsilon^2}} - 1$$

With these partial deviations, PROC HPPANEL uses OLS on the transformed series (including an intercept if you want).

The case of an unbalanced panel is somewhat more complicated. Wansbeek and Kapteyn show that the inverse of  $\Omega$  can be written as

$$\sigma_\epsilon^2 \Omega^{-1} = \mathbf{V} - \mathbf{V} \mathbf{Z}_2 \tilde{\mathbf{P}}^{-1} \mathbf{Z}'_2 \mathbf{V}$$

with the following:

$$\mathbf{V} = \mathbf{I}_M - \mathbf{Z}_1 \tilde{\Delta}_N^{-1} \mathbf{Z}'_1$$

$$\tilde{\mathbf{P}} = \tilde{\Delta}_T - \mathbf{A} \tilde{\Delta}_N^{-1} \mathbf{A}'$$

$$\tilde{\Delta}_N = \Delta_N + \left( \frac{\sigma_\epsilon^2}{\sigma_v^2} \right) \mathbf{I}_N$$

$$\tilde{\Delta}_T = \Delta_T + \left( \frac{\sigma_\epsilon^2}{\sigma_e^2} \right) \mathbf{I}_T$$

By using the inverse of the covariance matrix of the error, it becomes possible to complete GLS on the unbalanced panel.

## Linear Hypothesis Testing

For a linear hypothesis of the form  $\mathbf{R}\beta = \mathbf{r}$ , where  $\mathbf{R}$  is  $J \times K$  and  $\mathbf{r}$  is  $J \times 1$ , the  $F$ -statistic with  $J, M - K$  degrees of freedom is computed as

$$(\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}\hat{\mathbf{V}}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})$$

However, it is also possible to write the  $F$  statistic as

$$F = \frac{(\hat{\mathbf{u}}_*' \hat{\mathbf{u}}_* - \hat{\mathbf{u}}' \hat{\mathbf{u}}) / J}{\hat{\mathbf{u}}' \hat{\mathbf{u}} / (M - K)}$$

where

- $\hat{\mathbf{u}}_*$  is the residual vector from the restricted regression
- $\hat{\mathbf{u}}$  is the residual vector from the unrestricted regression
- $J$  is the number of restrictions
- $M - K$  are the degrees of freedom,  $M$  is the number of observations, and  $K$  is the number of parameters in the model

The Wald, likelihood ratio (LR), and LaGrange multiplier (LM) tests are all related to the  $F$  test. You use this relationship of the  $F$  test to the likelihood ratio and LaGrange multiplier tests. The Wald test is calculated from its definition.

The Wald test statistic is

$$W = (\mathbf{R}\hat{\beta} - \mathbf{r})' [\mathbf{R}\hat{\mathbf{V}}\mathbf{R}']^{-1} (\mathbf{R}\hat{\beta} - \mathbf{r})$$

The likelihood ratio is

$$\text{LR} = M \ln \left[ 1 + \frac{1}{M - K} JF \right]$$

The LaGrange multiplier test statistic is

$$\text{LM} = M \left[ \frac{JF}{M - K + JF} \right]$$

where  $JF$  represents the number of restrictions multiplied by the result of the  $F$  test.

The distribution of these test statistics is the  $\chi^2$  distribution whose degrees of freedom equal the number of restrictions imposed ( $J$ ). The three tests are asymptotically equivalent, but they have differing small-sample properties. Greene (2000, p. 392) and Davidson and MacKinnon (1993, pp. 456–458) discuss the small-sample properties of these statistics.



## Specification Tests

The HPPANEL procedure outputs one specification test for random effects: the Hausman (1978) specification test ( $m$  statistic) can be used to test hypotheses in terms of bias or inconsistency of an estimator. This test was also proposed by Wu (1973) and further extended in Hausman and Taylor (1982). Hausman's  $m$  statistic is as follows.

Consider two estimators,  $\hat{\beta}_a$  and  $\hat{\beta}_b$ , which under the null hypothesis are both consistent, but only  $\hat{\beta}_a$  is asymptotically efficient. Under the alternative hypothesis, only  $\hat{\beta}_b$  is consistent. The  $m$  statistic is

$$m = (\hat{\beta}_b - \hat{\beta}_a)' (\hat{S}_b - \hat{S}_a)^{-1} (\hat{\beta}_b - \hat{\beta}_a)$$

where  $\hat{S}_b$  and  $\hat{S}_a$  are consistent estimates of the asymptotic covariance matrices of  $\hat{\beta}_b$  and  $\hat{\beta}_a$ . Then  $m$  is distributed as  $\chi^2$  with  $k$  degrees of freedom, where  $k$  is the dimension of  $\hat{\beta}_a$  and  $\hat{\beta}_b$ .

In the random-effects specification, the null hypothesis of no correlation between effects and regressors implies that the OLS estimates of the slope parameters are consistent and inefficient but the GLS estimates of the slope parameters are consistent and efficient. This facilitates a Hausman specification test. The reported degrees of freedom for the  $\chi^2$  statistic are equal to the number of slope parameters. If the null hypothesis holds, the random-effects specification should be used.

## OUTPUT OUT= Data Set

PROC HPPANEL writes the initial data of the estimated model, predicted values, and residuals to an output data set when the OUT= option is specified in the OUTPUT statement. The OUT= data set contains the following variables:

<code>_CSID_</code>	is the value of the cross section ID. The variable name is the one specified in the id statement.
<code>_TSID_</code>	is the value of the time period in the dynamic model. The variable name is the one specified in the id statement.
Regressors	are the values of regressor variables that are specified in the COPYVAR option.
Pred	is the predicted value of dependent variable. This column is output only if the PRED option is specified.
Resid	is the residual from the regression. This column is output only if the RESIDUAL option is specified.

## OUTEST= Data Set

PROC HPPANEL writes the parameter estimates to an output data set when the OUTEST= option is specified in the PROC HPPANEL statement. The OUTEST= data set contains the following variables in the PROC statement:

<code>_METHOD_</code>	is a character variable that identifies the estimation method.
<code>_TYPE_</code>	is a character variable that identifies the type of observation. Values of the <code>_TYPE_</code> variable are CORRB, COVB, CSPARMS, STD, and the type of model estimated. The CORRB observation contains correlations of the parameter estimates; the COVB observation contains covariances of the parameter estimates; the STD observation indicates the row of standard deviations of the corresponding coefficients; and the type of model estimated observation contains the parameter estimates.
<code>_NAME_</code>	is a character variable that contains the name of a regressor variable for COVB and CORRB observations and is left blank for other observations. The <code>_NAME_</code> variable is used in conjunction with the <code>_TYPE_</code> values COVB and CORRB to identify rows of the correlation or covariance matrix.
<code>_DEPVAR_</code>	is a character variable that contains the name of the response variable.
<code>_MSE_</code>	is the mean square error of the transformed model.
<code>_VARCS_</code>	is the variance component estimate due to cross sections. The <code>_VARCS_</code> variable is included in the OUTEST= data set when the RANONE option is specified in the MODEL or PROC HPPANEL statement.
<code>_VARTS_</code>	is the variance component estimate due to time series. The <code>_VARTS_</code> variable is included in the OUTEST= data set when the RANTWO option is specified in the MODEL or PROC HPPANEL statement.
<code>_VARERR_</code>	is the variance component estimate due to error. The <code>_VARERR_</code> variable is included in the OUTEST= data set when the RANONE or RANTWO option is specified in the MODEL or PROC HPPANEL statement.
Intercept	is the intercept parameter estimate. (The intercept is missing for models when the NOINT option is specified in the MODEL statement.)
Regressors	are the regressor variables that are specified in the MODEL statement. The regressor variables in the OUTEST= data set contain the corresponding parameter estimates, and the corresponding covariance or correlation matrix elements for <code>_TYPE_=COVB</code> and <code>_TYPE_=CORRB</code> observations.

---

## Printed Output

The printed output from PROC HPPANEL includes the following:

- the model information, which includes the data source, the dependent variable name, the estimation method used, the execution mode, and for random-effects model analysis, the variance component estimation method.
- the number of observations
- the fit statistics, which include the sum of squared error (SSE), the degree of freedom for error (DFE), the mean square error (MSE), the root mean square error (RMSE), and the R-square
- the error components estimates for random-effects model

- the Hausman test statistics, which include the degree of freedom (DF), the test statistics, and the  $p$ -value.
- the regression parameter estimates and analysis, which include for each regressor the name of the regressor, the degrees of freedom, the parameter estimate, the standard error of the estimate, a  $t$  statistic for testing whether the estimate is significantly different from 0, and the significance probability of the  $t$  statistic

Optionally, PROC HPPANEL prints the following:

- the covariance and correlation of the resulting regression parameter estimates
- the WALD, LR, and LM test statistics for linear equality restrictions that are specified in the TEST statements
- the timing breakdown of the procedure steps

## ODS Table Names

PROC HPPANEL assigns a name to each table it creates. You can use these names to refer to the table when you use the Output Delivery System (ODS) to select tables and create output data sets. These names are listed in [Table 7.2](#).

**Table 7.2** ODS Tables Produced in PROC HPPANEL

ODS Table Name	Description	Option
<b>ODS Tables Created by the MODEL Statement</b>		
ModelInfo	Model information	Default
PerformanceInfo	Performance information	Default
Nobs	Number of observations	Default
FitStatistics	Fit statistics	Default
ParameterEstimates	Parameter estimates	Default
CovB	Covariance of parameter estimates	COVB
CorrB	Correlations of parameter estimates	CORRB
RandomEffectsTest	Hausman test for random effects	RANONE, RANTWO
<b>ODS Tables Created by the TEST Statement</b>		
TestResults	Test results	
<b>ODS Tables Created by the PERFORMANCE Statement</b>		
Timing	Timing Table	

## Example: HPPANEL Procedure

### Example 7.1: One-Way Random-Effects High-Performance Model

This example shows the use of the one-way random effects model that is available in the HPPANEL procedure with an emphasis on processing a large data set and on the performance improvements that are achieved by executing in a high-performance distributed environment.

The following DATA step generates 5 million replications from a one-way panel data that includes 50,000 cross sections and 100 time periods:

```
data hppan_ex01 (keep = cs ts y x1-x10);
  retain seed1 55371 seed2 97335 seed3 19412;
  array x[10];
  label y = 'dependent var.';
  label x1='first independent var.';
  label x2='second independent var.';
  label x3='third independent var.';
  int = 1;
  do cs = 1 to 50000;
    dummy = 10000*rannor( seed3 );
    do ts = 1 to 100;
      /*- generate regressors and compute the structural */
      /*- part of the dependent variable */
      y = 5; /* intercept */
      do k = 1 to 10;
        x[k] = (cs + ts ) * (0.001*ranuni( k ) + 1) ;
        y = y + x[k] * k;
      end;

      /*- add an error term, such that e ~ N(0,100) -----*/
      y = y + 10000*rannor( seed2 );
      /*- add a random effect, such that e ~ N(0,100) -----*/
      y = y + dummy;
      output;
    end;
  end;
run;
```

The model is executed in the distributed computing environment with one thread and only one node. These settings are used to obtain a hypothetical environment that might resemble running the HPPANEL procedure on a desktop workstation with a single-core CPU. To run the following statements successfully, you need to set the macro variables GRIDHOST and GRIDINSTALLLOC to resolve to appropriate values, or you can replace the references to the macro variables in the example with the appropriate values.

```
option set=GRIDHOST="%GRIDHOST";
option set=GRIDINSTALLLOC="%GRIDINSTALLLOC";

proc hppanel data=hppan_ex01 ranone;
  id cs ts;
  model y = x1-x10;
  performance nodes = 1 threads = 1 details
    host="%GRIDHOST" install="%GRIDINSTALLLOC";
run;
```

In [Output 7.1.1](#), the “Performance Information” table shows that the model was estimated on the grid that is defined in a macro variable named GRIDHOST in a distributed environment on only one node with one thread. The grid install location is defined in a macro variable named GRIDINSTALLLOC.

**Output 7.1.1** Grid Information with One Node and One Thread

Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	1
Number of Threads per Node	1

[Output 7.1.2](#) shows the results for the one-way random effects model. The “Model Information” table shows detailed information about the model. The “Number of Observations” table indicates that all 5 million observations were used to fit the model. All parameter estimates in the “Parameter Estimates” table are highly significant and correspond to the theoretical values that were set for them during the data generating process. In the “Timing” table, you can see that for 5 million observations, computing moments took 5706.25 seconds, and the cross-product accumulation took 272.95 seconds.

**Output 7.1.2** One-Way Random Effects Model

Model Information					
Data Source	WORK.HPPAN_EX01				
Response Variable	Y				
Model	RANONE				
Variance Component	WANSBEEK				
Execution Mode	Distributed				
Fit Statistics					
Sum of Squared Error	5.00008E14				
Degree of Freedom	4999989				
Mean Squared Error	100001811				
Root Mean Squared Error	10000				
R-Square	0.98318				
Variance Component Estimates					
Variance Component for Cross Sections	1.0704E8				
Variance Component for Error	1.0007E8				
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t
Intercept	1	27.06229	93.06534	0.29	0.7712
x1	1	0.44857	0.51089	0.88	0.3799
x2	1	2.18393	0.51098	4.27	<.0001
x3	1	2.70052	0.51099	5.28	<.0001
x4	1	4.49262	0.51100	8.79	<.0001
x5	1	5.54728	0.51076	10.86	<.0001
x6	1	6.50872	0.51088	12.74	<.0001
x7	1	6.54937	0.51098	12.82	<.0001
x8	1	7.09160	0.51090	13.88	<.0001
x9	1	8.64988	0.51092	16.93	<.0001
x10	1	10.82664	0.51051	21.21	<.0001
Procedure Task Timing					
Task		Seconds		Percent	
Data Read and Variable Levelization		2.00		0.03%	
Communication to Client		0.00		0.00%	
Computing Moments		5706.25		95.40%	
Cross-Product Accumulation		272.95		4.56%	

In the following statements, the PERFORMANCE statement is modified to request a grid that has 10 nodes, where each node spawns one thread:

```
proc hppanel data=hppan_ex01 ranone;
  id cs ts;
  model y = x1-x10;
  performance nodes = 10 threads = 1 details
    host="&GRIDHOST" install="&GRIDINSTALLLOC";
run;
```

In [Output 7.1.3](#), the “Performance Information” table shows that the model was estimated on the grid that is defined in a macro variable named GRIDHOST in a distributed environment on 10 nodes with one thread each. The grid install location is defined in a macro variable named GRIDINSTALLLOC.

**Output 7.1.3** Grid Information for 10 Nodes with One Thread Each

Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	10
Number of Threads per Node	1

Although the two models are identical, estimating the model took only 14 minutes for the second implementation, which was run on a grid that used 10 nodes with one thread each, instead of 1 hour and 40 minutes for the first implementation.

**Output 7.1.4** Timing Information for 10 Nodes with One Thread Each

Procedure Task Timing		
Task	Seconds	Percent
Data Read and Variable Levelization	0.48	0.06%
Communication to Client	0.00	0.00%
Computing Moments	784.07	96.34%
Cross-Product Accumulation	29.34	3.60%

---

## References

- Davidson, R. and MacKinnon, J. G. (1993), *Estimation and Inference in Econometrics*, New York: Oxford University Press.
- Fuller, W. A. and Battese, G. E. (1974), “Estimation of Linear Models with Crossed-Error Structure,” *Journal of Econometrics*, 2, 67–78.
- Greene, W. H. (1990), *Econometric Analysis*, New York: Macmillan.
- Greene, W. H. (2000), *Econometric Analysis*, 4th Edition, Upper Saddle River, NJ: Prentice-Hall.
- Hausman, J. A. (1978), “Specification Tests in Econometrics,” *Econometrica*, 46, 1251–1271.
- Hausman, J. A. and Taylor, W. E. (1982), “A Generalized Specification Test,” *Economics Letters*, 8, 239–245.
- Nerlove, M. (1971), “Further Evidence on the Estimation of Dynamic Relations from a Time Series of Cross Sections,” *Econometrica*, 39, 359–382.
- Wallace, T. and Hussain, A. (1969), “The Use of Error Components Model in Combining Cross Section with Time Series Data,” *Econometrica*, 37, 55–72.
- Wansbeek, T. and Kapteyn, A. (1989), “Estimation of the Error-Components Model with Incomplete Panels,” *Journal of Econometrics*, 41, 341–361.
- Wu, D. M. (1973), “Alternative Tests of Independence between Stochastic Regressors and Disturbances,” *Econometrica*, 41, 733–750.



# Chapter 8

## The HPQLIM Procedure

### Contents

---

Overview: HPQLIM Procedure . . . . .	<b>186</b>
PROC HPQLIM Features . . . . .	187
Getting Started: HPQLIM Procedure . . . . .	<b>187</b>
Syntax: HPQLIM Procedure . . . . .	<b>189</b>
Functional Summary . . . . .	189
PROC HPQLIM Statement . . . . .	192
BAYES Statement . . . . .	196
BOUNDS Statement . . . . .	200
BY Statement . . . . .	200
ENDOGENOUS Statement . . . . .	201
FREQ Statement . . . . .	203
HETERO Statement . . . . .	203
INIT Statement . . . . .	204
MODEL Statement . . . . .	204
OUTPUT Statement . . . . .	206
PERFORMANCE Statement . . . . .	208
PRIOR Statement . . . . .	208
RESTRICT Statement . . . . .	209
TEST Statement . . . . .	209
WEIGHT Statement . . . . .	210
Details: HPQLIM Procedure . . . . .	<b>211</b>
Ordinal Discrete Choice Modeling . . . . .	211
Limited Dependent Variable Models . . . . .	212
Stochastic Frontier Production and Cost Models . . . . .	213
Heteroscedasticity . . . . .	214
Tests on Parameters . . . . .	215
Bayesian Analysis . . . . .	216
Prior Distributions . . . . .	217
Output to SAS Data Set . . . . .	220
OUTEST= Data Set . . . . .	223
Naming . . . . .	224
ODS Table Names . . . . .	225
ODS Graphics . . . . .	226
Examples: The HPQLIM Procedure . . . . .	<b>226</b>
Example 8.1: High-Performance Model with Censoring . . . . .	226
Example 8.2: Bayesian High-Performance Model with Censoring . . . . .	231
References . . . . .	<b>235</b>

---

## Overview: HPQLIM Procedure

The HPQLIM (high-performance qualitative and limited dependent variable model) procedure is a high-performance version of the QLIM procedure in SAS/ETS software, which analyzes univariate limited dependent variable models in which dependent variables are observed only in a limited range of values. Unlike the QLIM procedure, which can be run only on an individual workstation, the HPQLIM procedure takes advantage of a computing environment that enables it to distribute the optimization task to one or more nodes. In addition, each node can use one or more threads to perform the optimization on its subset of the data. When several nodes are used and each node uses several threads to carry out its part of the work, the result is a highly parallel computation that provides a dramatic gain in performance.

With the HPQLIM procedure you can read and write data in distributed form and perform analyses in distributed mode and single-machine mode. For more information about how to affect the execution mode of SAS high-performance analytical procedures, see the section “[Processing Modes](#)” on page 12 in Chapter 3, “[Shared Concepts and Topics](#).”

The HPQLIM procedure is specifically designed to operate in the high-performance distributed environment. It can use maximum likelihood or Bayesian methods. In both cases, the likelihood evaluation is performed in a distributed environment. By default, PROC HPQLIM uses multiple threads to perform computations.

The HPQLIM procedure is similar in use to the other SAS procedures that support regression or simultaneous equations models. For example, the standard model with censoring or truncation is estimated by specifying the endogenous variable to be truncated or censored. When the data are limited by specific values or variables, the limits of the dependent variable can be specified with the CENSORED or TRUNCATED option in the ENDOGENOUS or MODEL statement. For example, the two-limit censored model requires two variables: one that contains the lower (bottom) bound and one that contains the upper (top) bound. The following statements execute the model in the distributed computing environment with two threads and four nodes:

```
proc hpqlim data=a;
  model y = x1 x2 x3;
  endogenous y ~ censored(lb=bottom ub=top);
  performance nthreads=2 nodes=4 details;
run;
```

The bounds can be numbers if they are fixed for all observations in the data set. For example, the standard Tobit model can be specified as follows:

```
proc hpqlim data=a;
  model y = x1 x2 x3;
  endogenous y ~ censored(lb=0);
  performance nthreads=2 nodes=4 details;
run;
```

---

## PROC HPQLIM Features

The HPQLIM procedure supports the following models:

- linear regression models with heteroscedasticity
- Tobit models (censored and truncated) with heteroscedasticity
- stochastic frontier production and cost models

In linear regression models with heteroscedasticity, the assumption that error variance is constant across observations is relaxed. The HPQLIM procedure allows for a number of different linear and nonlinear variance specifications.

The HPQLIM procedure also offers a class of models in which the dependent variable is censored or truncated from below or above or both. When a continuous dependent variable is observed only within a certain range, and values outside this range are not available, the HPQLIM procedure offers a class of models that adjust for truncation. In some cases, the dependent variable is continuous only in a certain range, and all values outside this range are reported as being on its boundary. For example, if it is not possible to observe negative values, the value of the dependent variable is reported as equal to 0. Because the data are censored, ordinary least squares (OLS) results are inconsistent, and it cannot be guaranteed that the predicted values from the model will fall in the appropriate region.

Stochastic frontier production and cost models allow for random shocks of the production or cost. They include a systematic positive component in the error term that adjusts for technical or cost inefficiency.

The HPQLIM procedure can use maximum likelihood or Bayesian methods. Initial starting values for the nonlinear optimizations are typically calculated by OLS. Initial values for the Bayesian sampling are typically calculated by maximum likelihood.

---

## Getting Started: HPQLIM Procedure

This example illustrates the use of the HPQLIM procedure. The data were originally published by Mroz (1987), and the following statements show a subset of the Mroz (1987) data set:

```

title1 'Estimating a Tobit model';

data subset;
  input Hours Yrs_Ed Yrs_Exp @@;
  if Hours eq 0 then Lower=.;
  else Lower=Hours;
datalines;
0 8 9 0 8 12 0 9 10 0 10 15 0 11 4 0 11 6
1000 12 1 1960 12 29 0 13 3 2100 13 36
3686 14 11 1920 14 38 0 15 14 1728 16 3
1568 16 19 1316 17 7 0 17 15
;

```

In these data, Hours is the number of hours that the wife worked outside the household in a given year, Yrs\_Ed is the years of education, and Yrs\_Exp is the years of work experience.

By the nature of the data it is clear that there are a number of women who committed some positive number of hours to outside work ( $y_i > 0$  is observed). There are also a number of women who did not work outside the home at all ( $y_i = 0$  is observed). This yields the following model:

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

$$y_i = \begin{cases} y_i^* & \text{if } y_i^* > 0 \\ 0 & \text{if } y_i^* \leq 0 \end{cases}$$

where  $\epsilon_i \sim iidN(0, \sigma^2)$  and the set of explanatory variables is denoted by  $\mathbf{x}_i$ . The following statements fit a Tobit model to the hours worked with years of education and years of work experience as covariates:

```
/*-- Tobit Model --*/
proc hpqlim data=subset;
  model hours = yrs_ed yrs_exp;
  endogenous hours ~ censored(lb=0);
  performance nthreads=2 nodes=4 details;
run;
```

The output of the HPQLIM procedure is shown in [Output 8.1](#).

**Figure 8.1** Tobit Analysis Results

Estimating a Tobit model					
The HPQLIM Procedure					
Model Fit Summary					
Number of Endogenous Variables					1
Endogenous Variable				Hours	
Number of Observations				17	
Log Likelihood				-74.93700	
Maximum Absolute Gradient				1.18953E-6	
Number of Iterations				23	
Optimization Method				Quasi-Newton	
AIC				157.87400	
Schwarz Criterion				161.20685	
Parameter Estimates					
Parameter	DF	Estimate	Standard Error	t Value	Approx Pr >  t
Intercept	1	-5598.295129	27.692220	-202.16	<.0001
Yrs_Ed	1	373.123254	53.988877	6.91	<.0001
Yrs_Exp	1	63.336247	36.551299	1.73	0.0831
_Sigma	1	1582.859635	390.076480	4.06	<.0001

The “Parameter Estimates” table contains four rows. The first three rows correspond to the vector estimate of the regression coefficients  $\boldsymbol{\beta}$ . The last row is called `_Sigma`, which corresponds to the estimate of the error variance  $\sigma$ .

## Syntax: HPQLIM Procedure

The following statements are available in the HPQLIM procedure:

```

PROC HPQLIM options ;
  BAYES <options> ;
  BOUNDS bound1 < , bound2 ... > ;
  BY variables ;
  FREQ variable ;
  ENDOGENOUS variables ~ options ;
  HETERO dependent variables ~ exogenous variables / options ;
  INIT initvalue1 < , initvalue2 ... > ;
  MODEL dependent variables = regressors / options ;
  OUTPUT options ;
  PRIOR variables ~ distributions ;
  RESTRICT restriction1 < , restriction2 ... > ;
  TEST options ;
  WEIGHT variable ;

```

One MODEL statement is required. If a FREQ or WEIGHT statement is specified more than once, the variable that is specified in the first instance is used.

## Functional Summary

Table 8.1 summarizes the statements and options used with the HPQLIM procedure.

**Table 8.1** PROC HPQLIM Functional Summary

Description	Statement	Option
<b>Data Set Options</b>		
Specifies the input data set	PROC HPQLIM	DATA=
Writes parameter estimates to an output data set	PROC HPQLIM	OUTEST=
Writes predictions to an output data set	OUTPUT	OUT=
<b>Declaring the Role of Variables</b>		
Specifies BY-group processing	BY	
Specifies a frequency variable	FREQ	
Specifies a weight variable	WEIGHT	NONNORMALIZE
<b>Printing Control Options</b>		
Requests all printing options	PROC HPQLIM	PRINTALL
Prints the correlation matrix of the estimates	PROC HPQLIM	CORRB
Prints the covariance matrix of the estimates	PROC HPQLIM	COVB
Suppresses the normal printed output	PROC HPQLIM	NOPRINT

Table 8.1 *continued*

Description	Statement	Option
<b>Plotting Options</b>		
Displays plots	PROC HPQLIM	PLOTS=
<b>Optimization Process Control Options</b>		
Selects the iterative minimization method to use	PROC HPQLIM	METHOD=
Specifies the maximum number of iterations allowed	PROC HPQLIM	MAXITER=
Specifies the maximum number of function calls	PROC HPQLIM	MAXFUNC=
Specifies the upper limit of CPU time in seconds	PROC HPQLIM	MAXTIME=
Specifies an absolute convergence criterion	PROC HPQLIM	ABSCONV=
Specifies an absolute function convergence criterion	PROC HPQLIM	ABSFCNV=
Specifies an absolute gradient convergence criterion	PROC HPQLIM	ABSGCONV=
Specifies a relative function convergence criterion	PROC HPQLIM	FCONV=
Specifies a relative gradient convergence criterion	PROC HPQLIM	GCONV=
Specifies an absolute parameter convergence criterion	PROC HPQLIM	ABSXCONV=
Specifies a matrix singularity criterion	PROC HPQLIM	SINGULAR=
Sets boundary restrictions on parameters	BOUNDS	
Sets initial values for parameters	INIT	
Sets linear restrictions on parameters	RESTRICT	
<b>Model Estimation Options</b>		
Suppresses the intercept parameter	MODEL	NOINT
Specifies the method to calculate parameter covariance	PROC HPQLIM	COVEST=
<b>Bayesian MCMC Options</b>		
Specifies the initial values of the MCMC	INIT	
Specifies the maximum number of tuning phases	BAYES	MAXTUNE=
Specifies the minimum number of tuning phases	BAYES	MINTUNE=
Specifies the number of burn-in iterations	BAYES	NBI=
Specifies the number of iterations during the sampling phase	BAYES	NMC=
Specifies the number of iterations during the tuning phase	BAYES	NTU=
Controls options for constructing the initial proposal covariance matrix	BAYES	PROPCOV
Specifies the sampling scheme	BAYES	SAMPLING=
Specifies the random number generator seed	BAYES	SEED=
Controls the thinning of the Markov chain	BAYES	THIN=
<b>Bayesian Summary Statistics and Convergence Diagnostic Options</b>		
Displays convergence diagnostics	BAYES	DIAGNOSTICS=
Displays summary statistics of the posterior samples	BAYES	STATISTICS=

**Table 8.1** *continued*

Description	Statement	Option
<b>Bayesian Prior and Posterior Sample Options</b>		
Specifies a SAS data set for the posterior samples	BAYES	OUTPOST=
<b>Bayesian Analysis Options</b>		
Specifies the normal prior distribution	PRIOR	NORMAL(MEAN=, VAR=)
Specifies the gamma prior distribution	PRIOR	GAMMA(SHAPE=, SCALE=)
Specifies the inverse gamma prior distribution	PRIOR	IGAMMA(SHAPE=, SCALE=)
Specifies the uniform prior distribution	PRIOR	UNIFORM(MIN=, MAX=)
Specifies the beta prior distribution	PRIOR	BETA(SHAPE1=, SHAPE2=, MIN=, MAX=)
Specifies the <i>t</i> prior distribution	PRIOR	T(LOCATION=, DF=)
<b>Endogenous Variable Options</b>		
Specifies a discrete variable	ENDOGENOUS	DISCRETE()
Specifies a censored variable	ENDOGENOUS	CENSORED()
Specifies a truncated variable	ENDOGENOUS	TRUNCATED()
Specifies a stochastic frontier variable	ENDOGENOUS	FRONTIER()
<b>Heteroscedasticity Model Options</b>		
Specifies the function for heteroscedasticity models	HETERO	LINK=
Squares the function for heteroscedasticity models	HETERO	SQUARE
Specifies no constant for heteroscedasticity models	HETERO	NOCONST
<b>Output Control Options</b>		
Outputs predicted values	OUTPUT	PREDICTED
Outputs the structured part	OUTPUT	XBETA
Outputs residuals	OUTPUT	RESIDUAL
Outputs the error standard deviation	OUTPUT	ERRSTD
Outputs marginal effects	OUTPUT	MARGINAL
Outputs probability for the current response	OUTPUT	PROB
Outputs probability for all responses	OUTPUT	PROBALL
Outputs the expected value	OUTPUT	EXPECTED
Outputs the conditional expected value	OUTPUT	CONDITIONAL
Outputs inverse Mills ratio	OUTPUT	MILLS
Outputs technical efficiency measures	OUTPUT	TE1 TE2
Includes covariances in the OUTEST= data set	PROC HPQLIM	COVOUT
Includes correlations in the OUTEST= data set	PROC HPQLIM	CORROUT

Table 8.1 *continued*

Description	Statement	Option
<b>Test Request Options</b>		
Requests Wald, Lagrange multiplier, and likelihood ratio tests	TEST	ALL
Requests the Wald test	TEST	WALD
Requests the Lagrange multiplier test	TEST	LM
Requests the likelihood ratio test	TEST	LR

## PROC HPQLIM Statement

**PROC HPQLIM** *options* ;

The PROC HPQLIM statement invokes the HPQLIM procedure. You can specify the following *options*.

### Data Set Options

**DATA**=*SAS-data-set*

specifies the input SAS data set. If this option is not specified, PROC HPQLIM uses the most recently created SAS data set.

### Output Data Set Options

**OUTEST**=*SAS-data-set*

writes the parameter estimates to an output data set.

**COVOUT**

writes the covariance matrix for the parameter estimates to the OUTEST= data set. This option is valid only if the OUTEST= option is specified.

**CORROUT**

writes the correlation matrix for the parameter estimates to the OUTEST= data set. This option is valid only if the OUTEST= option is specified.

### Printing Options

**NOPRINT**

suppresses the normal printed output but does not suppress error listings. If this option is specified, then any other print option is turned off.

**PRINTALL**

turns on all the printing options. The options that are set by PRINTALL are COVB and CORRB.

**CORRB**

prints the correlation matrix of the parameter estimates.



**COVB**

prints the covariance matrix of the parameter estimates.

**Model Estimation Options****COVEST=covariance-option**

specifies the method for calculating the covariance matrix of parameter estimates. You can specify the following *covariance-options*.

<b>OP</b>	specifies the covariance from the outer product matrix.
<b>HESSIAN</b>	specifies the covariance from the inverse Hessian matrix.
<b>QML</b>	specifies the covariance from the outer product and Hessian matrices (the quasi-maximum likelihood estimates).

The default is COVEST=HESSIAN.

**Optimization Control Options**

PROC HPQLIM uses the nonlinear optimization (NLO) subsystem to perform nonlinear optimization tasks. You can specify the following *options*:

**ABSCONV=r****ABSTOL=r**

specifies an absolute function value convergence criterion by which minimization stops when  $f(\theta^{(k)}) \leq r$ . The default value of  $r$  is the negative square root of the largest double-precision value, which serves only as a protection against overflows.

**ABSFCNV=r****ABSFTOL=r**

specifies an absolute function difference convergence criterion by which minimization stops when the function value has a small change in successive iterations:

$$|f(\theta^{(k-1)}) - f(\theta^{(k)})| \leq r$$

The default value is  $r = 0$ .

**ABSGCONV=r****ABSGTOL=r**

specifies an absolute gradient convergence criterion. Optimization stops when the maximum absolute gradient element is small:

$$\max_j |g_j(\theta^{(k)})| \leq r$$

The default value is  $r=1\text{E}-5$ .

**ABSXCONV=r****ABSXTOL=r**

specifies an absolute parameter convergence criterion. Optimization stops when the Euclidean distance between successive parameter vectors is small:

$$\|\theta^{(k)} - \theta^{(k-1)}\|_2 \leq r$$

The default is 0.

**FCONV=*r*****FTOL=*r***

specifies a relative function convergence criterion. Optimization stops when a relative change of the function value in successive iterations is small:

$$\frac{|f(\theta^{(k)}) - f(\theta^{(k-1)})|}{|f(\theta^{(k-1)})|} \leq r$$

The default value is  $r = 2\epsilon$ , where  $\epsilon$  denotes the machine precision constant, which is the smallest double-precision floating-point number such that  $1 + \epsilon > 1$ .

**GCONV=*r*****GTOL=*r***

specifies a relative gradient convergence criterion. For all techniques except CONGRA, optimization stops when the normalized predicted function reduction is small:

$$\frac{g(\theta^{(k)})^T [H^{(k)}]^{-1} g(\theta^{(k)})}{|f(\theta^{(k)})|} \leq r$$

For the CONGRA technique (where a reliable Hessian estimate  $H$  is not available), the following criterion is used:

$$\frac{\|g(\theta^{(k)})\|_2^2 \|s(\theta^{(k)})\|_2}{\|g(\theta^{(k)}) - g(\theta^{(k-1)})\|_2 |f(\theta^{(k)})|} \leq r$$

The default value is  $r = 1\text{E-}8$ .

**MAXFUNC=*i*****MAXFU=*i***

specifies the maximum number of function calls in the optimization process. The default is 1,000.

The optimization can terminate only after completing a full iteration. Therefore, the number of function calls that are actually performed can exceed the number of calls that are specified by this option.

**MAXITER=*i*****MAXIT=*i***

specifies the maximum number of iterations in the optimization process. The default is 200.

**MAXTIME=*r***

specifies an upper limit of  $r$  seconds of CPU time for the optimization process. The default value is the largest floating-point double representation of your computer. The time that is specified by this option is checked only once at the end of each iteration. Therefore, the actual running time can be much longer than  $r$ . The actual running time includes the remaining time needed to finish the iteration and the time needed to generate the output of the results.

**METHOD=*value***

specifies the iterative minimization method to use. The default is METHOD=NEWRAP. You can specify the following *values*:

**CONGRA** specifies the conjugate-gradient method.

**DBLDOG** specifies the double dogleg method.

<b>NONE</b>	specifies that no optimization be performed beyond using the ordinary least squares method to compute the parameter estimates.
<b>NEWRAP</b>	specifies the Newton-Raphson method (the default).
<b>NRRIDG</b>	specifies the Newton-Raphson Ridge method.
<b>QUANEW</b>	specifies the quasi-Newton method.
<b>TRUREG</b>	specifies the trust region method.

**SINGULAR=*r***

specifies the general singularity criterion that is applied by the HPQLIM procedure in sweeps and inversions. The default for the optimization is 1E-8.

## Plotting Options

**PLOTS**< (*global-plot-options*) > = *plot-request* | (*plot-requests*)

controls the display of plots. By default, the plots are displayed in panels unless the UNPACK *global-plot-option* is specified. When you specify only one *plot-request*, you can omit the parentheses around it.

### Global Plot Options

You can specify the following *global-plot-options*:

**ONLY**

displays only the requested plot.

**UNPACKPANEL****UNPACK**

specifies that all paneled plots be unpacked, meaning that each plot in a panel is displayed separately.

### Plot Requests

You can specify the following *plot-requests*:

**ALL**

specifies all types of available plots.

**AUTOCORR**< (**LAGS**=*n*) >

displays the autocorrelation function plots for the parameters. The optional LAGS= suboption specifies the number (up to lag *n*) of autocorrelations to be plotted in the autocorrelation function plot. If this suboption is not specified, autocorrelations are plotted up to lag 50. This *plot-request* is available only for Bayesian analysis.

**BAYESDIAG**

is equivalent to specifying the TRACE, AUTOCORR, and DENSITY *plot-requests*.

**DENSITY**< (**FRINGE**) >

displays the kernel density plots for the parameters. If you specify the FRINGE suboption, a fringe plot is created on the X axis of the kernel density plot. This *plot-request* is available only for Bayesian analysis.

**NONE**

suppresses all diagnostic plots.

**TRACE<(SMOOTH)>**

displays the trace plots for the parameters. The SMOOTH suboption displays a fitted penalized B-spline curve for each plot. This *plot-request* is available only for Bayesian analysis.

**BAYES Statement**

**BAYES** <options> ;

The BAYES statement controls the Metropolis sampling scheme that is used to obtain samples from the posterior distribution of the underlying model and data.

**DIAGNOSTICS=ALL | NONE** | (*keyword-list*)

**DIAG=ALL | NONE** | (*keyword-list*)

controls which diagnostics are produced. All the following diagnostics are produced when you specify **DIAGNOSTICS=ALL**. If you do not want any of these diagnostics, specify **DIAGNOSTICS=NONE**. If you want some but not all of the diagnostics, or if you want to change certain settings of these diagnostics, specify one or more of the following keywords. The default is **DIAGNOSTICS=NONE**.

**AUTOCORR** <(**LAGS=numeric-list**)>

computes the autocorrelations at lags that are specified in the *numeric-list*. Elements in the *numeric-list* are truncated to integers, and repeated values are removed. If the **LAGS=** option is not specified, autocorrelations of lags 1, 5, and 10 are computed.

**ESS**

computes Carlin's estimate of the effective sample size, the correlation time, and the efficiency of the chain for each parameter.

**GEWEKE** <(*geweke-options*)>

computes the Geweke spectral density diagnostics, which are essentially a two-sample *t* test between the first  $f_1$  portion and the last  $f_2$  portion of the chain. The defaults are  $f_1 = 0.1$  and  $f_2 = 0.5$ , but you can choose other fractions by using the following *geweke-options*:

**FRAC1=value**

specifies the fraction  $f_1$  for the first window.

**FRAC2=value**

specifies the fraction  $f_2$  for the second window.

**HEIDELBERGER** <(*heidel-options*)>

computes for each variable the Heidelberg and Welch diagnostic, which consists of a stationarity test of the null hypothesis that the sample values form a stationary process. If the stationarity test is not rejected, a halfwidth test is then carried out. Optionally, you can specify one or more of the following *heidel-options*:

**EPS=value**

specifies a positive number  $\epsilon$  such that if the halfwidth is less than  $\epsilon$  times the sample mean of the retained iterates, the halfwidth test is passed.

**HALPHA=value**

specifies the  $\alpha$  level ( $0 < \alpha < 1$ ) for the halfwidth test.

**SALPHA=value**

specifies the  $\alpha$  level ( $0 < \alpha < 1$ ) for the stationarity test.

**MCSE****MCERROR**

computes the Monte Carlo standard error for each parameter. The Monte Carlo standard error, which measures the simulation accuracy, is the standard error of the posterior mean estimate and is calculated as the posterior standard deviation divided by the square root of the effective sample size.

**RAFTERY<(raftery-options)>**

computes the Raftery and Lewis diagnostics, which evaluate the accuracy of the estimated quantile ( $\hat{\theta}_Q$  for a given  $Q \in (0, 1)$ ) of a chain.  $\hat{\theta}_Q$  can achieve any degree of accuracy when the chain is allowed to run for a long time. The computation stops when the estimated probability  $\hat{P}_Q = \Pr(\theta \leq \hat{\theta}_Q)$  reaches within  $\pm R$  of the value  $Q$  with probability  $S$ ; that is,  $\Pr(Q - R \leq \hat{P}_Q \leq Q + R) = S$ . The following *raftery-options* enable you to specify  $Q$ ,  $R$ ,  $S$ , and a precision level  $\epsilon$  for the test:

**QUANTILE | Q=value**

specifies the order (a value between 0 and 1) of the quantile of interest. The default is 0.025.

**ACCURACY | R=value**

specifies a small positive number as the margin of error for measuring the accuracy of the estimation of the quantile. The default is 0.005.

**PROBABILITY | S=value**

specifies the probability of attaining the accuracy of the estimation of the quantile. The default is 0.95.

**EPSILON | EPS=value**

specifies the tolerance level (a small positive number) for the stationary test. The default is 0.001.

**MINTUNE=number**

specifies the minimum number of tuning phases. The default is 2.

**MAXTUNE=number**

specifies the maximum number of tuning phases. The default is 24.

**NBI=number**

specifies the number of burn-in iterations before the chains are saved. The default is 1,000.

**NMC=number**

specifies the number of iterations after the burn-in. The default is 1,000.

**NTU=number**

specifies the number of samples for each tuning phase. The default is 500.

**OUTPOST=SAS-data-set**

names the SAS data set to contain the posterior samples. Alternatively, you can create the output data set by specifying an ODS OUTPUT statement as follows:

**ODS OUTPUT POSTERIORSAMPLE = < SAS-data-set > ;**

**PROPCOV=value**

specifies the method that is used in constructing the initial covariance matrix for the Metropolis-Hastings algorithm. The QUANEW and NMSIMP methods find numerically approximated covariance matrices at the optimum of the posterior density function with respect to all continuous parameters. The tuning phase starts at the optimized values; in some problems, this can greatly increase convergence performance. If the approximated covariance matrix is not positive definite, then an identity matrix is used instead. You can specify the following *values*:

**CONGRA**

performs a conjugate-gradient optimization.

**DBLDOG**

performs a version of double-dogleg optimization.

**NEWRAP**

performs a Newton-Raphson optimization that combines a line-search algorithm with ridging.

**NMSIMP**

performs a Nelder-Mead simplex optimization.

**NRRIDG**

performs a Newton-Raphson optimization with ridging.

**QUANEW**

performs a quasi-Newton optimization.

**TRUREG**

performs a trust-region optimization.

**SAMPLING=MULTIMETROPOLIS | UNIMETROPOLIS**

specifies how to sample from the posterior distribution. **SAMPLING=MULTIMETROPOLIS** implements a Metropolis sampling scheme on a single block that contains all the parameters of the model. **SAMPLING=UNIMETROPOLIS** implements a Metropolis sampling scheme on multiple blocks, one for each parameter of the model. The default is **SAMPLING=MULTIMETROPOLIS**.

**SEED=number**

specifies an integer seed in the range 1 to  $2^{31} - 1$  for the random number generator in the simulation. Specifying a seed enables you to reproduce identical Markov chains for the same specification. If you do not specify the **SEED=** option, or if you specify a nonpositive seed, a random seed is derived from the time of day.

**STATISTICS** <(global-options)> = **ALL** | **NONE** | *keyword* | (*keyword-list*)

**STATS** <(global-options)> = **ALL** | **NONE** | *keyword* | (*keyword-list*)

controls the number of posterior statistics that are produced. Specifying **STATISTICS=ALL** is equivalent to specifying **STATISTICS=(CORR COV INTERVAL PRIOR SUMMARY)**. If you do not want any posterior statistics, specify **STATISTICS=NONE**. The default is **STATISTICS=(SUMMARY INTERVAL)**. You can specify the following *global-options*:

**ALPHA**=*value* < ,*value*>...< ,*value*>

controls the probabilities of the credible intervals. The *value*, which must be between 0 and 1, produces a pair of 100(1-*value*)% equal-tail and highest posterior density (HPD) intervals for each parameter. The default is **ALPHA=0.05**, which yields the 95% credible intervals for each parameter.

**PERCENT**=*value* < ,*value*>...< ,*value*>

requests the percentile points of the posterior samples. The *value* must be between 0 and 100. The default is **PERCENT=25, 50, 75**, which yields the 25th, 50th, and 75th percentile points, respectively, for each parameter.

You can specify the following *keywords*:

#### **CORR**

produces the posterior correlation matrix.

#### **COV**

produces the posterior covariance matrix.

#### **INTERVAL**

produces equal-tail credible intervals and HPD intervals. The default is to produce the 95% equal-tail credible intervals and 95% HPD intervals, but you can use the **ALPHA= *global-option*** to request intervals of any probabilities.

#### **NONE**

suppresses printing of all summary statistics.

#### **PRIOR**

produces a summary table of the prior distributions that are used in the Bayesian analysis.

#### **SUMMARY**

produces the means, standard deviations, and percentile points (25th, 50th, and 75th) for the posterior samples. You can use the **PERCENT= *global-option*** to request specific percentile points.

**THIN**=*number*

**THINNING**=*number*

controls the thinning of the Markov chain. Only one in every  $k$  samples is used when **THIN**= $k$ . If **NBI**= $n_0$  and **NMC**= $n$ , the number of samples that are retained is

$$\left[ \frac{n_0 + n}{k} \right] - \left[ \frac{n_0}{k} \right]$$

where  $[a]$  represents the integer part of the number  $a$ . The default is **THIN**=1.

## BOUNDS Statement

**BOUNDS** *bound1* < , *bound2* ... > ;

The BOUNDS statement imposes simple boundary constraints on the parameter estimates. BOUNDS statement constraints refer to the parameters that are estimated by the HPQLIM procedure. You can specify any number of BOUNDS statements.

Each *bound* is composed of parameters, constants, and inequality operators. Parameters that are associated with regressor variables are referred to by the names of the corresponding regressor variables. Specify each bound as follows:

*item operator item* < *operator item* < *operator item* ... > >

Each *item* is a constant, the name of a parameter, or a list of parameter names. For more information about how parameters are named in the HPQLIM procedure, see the section “[Naming of Parameters](#)” on page 224. Each *operator* is <, >, <=, or >=.

You can use both the BOUNDS statement and the RESTRICT statement to impose boundary constraints; however, the BOUNDS statement provides a simpler syntax for specifying these types of constraints. For more information, see the section “[RESTRICT Statement](#)” on page 209.

The following BOUNDS statement constrains the estimates of the parameters that are associated with the variable *ttime* and the variables *x1* through *x10* to be between 0 and 1. The following example illustrates the use of parameter lists to specify boundary constraints.

```
bounds 0 < ttime x1-x10 < 1;
```

The following BOUNDS statement constrains the estimates of the correlation (*\_RHO*) and sigma (*\_SIGMA*) in the bivariate model:

```
bounds _rho >= 0, _sigma.y1 > 1, _sigma.y2 < 5;
```

## BY Statement

**BY** *variables* ;

A BY statement can be used with PROC HPQLIM to obtain separate analyses on observations in groups defined by the BY variables.

BY statement processing is not supported when the HPQLIM procedure runs alongside the database or alongside the Hadoop Distributed File System (HDFS). These modes are used if the input data are stored in a database or HDFS and the grid host is the appliance that houses the data.



## ENDOGENOUS Statement

**ENDOGENOUS** *variables ~ options ;*

The ENDOGENOUS statement specifies the type of dependent variables that appear on the left-hand side of the equation. The listed endogenous variables refer to the dependent variables that appear on the left-hand side of the equation. Currently, no right-hand-side endogeneity is handled in PROC HPQLIM. All variables that appear on the right-hand side of the equation are treated as exogenous.

### Discrete Variable Options

**DISCRETE** *<(discrete-options)>*

specifies that the endogenous variables in this statement be discrete. You can specify the following *discrete-options*:

**DISTRIBUTION**=*distribution-type*

**DIST**=*distribution-type*

**D**=*distribution-type*

specifies the cumulative distribution function that is used to model the response probabilities. You can specify the following *distribution-types*:

**LOGISTIC**      specifies the logistic distribution for the logit model.

**NORMAL**        specifies the normal distribution for the probit model.

By default, **DISTRIBUTION**=**NORMAL**.

**ORDER**=**DATA** | **FORMATTED** | **FREQ** | **INTERNAL**

specifies the sort order for the levels of the discrete variables that are specified in the ENDOGENOUS statement. This ordering determines which parameters in the model correspond to each level in the data. You can specify the following sort orders:

**DATA**            sorts levels by order of appearance in the input data set.

**FORMATTED**    sorts levels by formatted value. The sort order is machine-dependent.

**FREQ**            sorts levels by descending frequency count; levels that have the most observations come first in the order.

**INTERNAL**       sorts levels by unformatted value. The sort order is machine-dependent.

By default, **ORDER**=**FORMATTED**. For more information about sort order, see the chapter on the SORT procedure in the *Base SAS Procedures Guide*.

### Censored Variable Options

**CENSORED** *(censored-options)*

specifies that the endogenous variables in this statement be censored. You can specify the following *censored-options*:

**LB=***value* | *variable*

**LOWERBOUND=***value* | *variable*

specifies the lower bound of the censored variables. If *value* is missing or the value in *variable* is missing, no lower bound is set. By default, no lower bound is set.

**UB=***value* | *variable*

**UPPERBOUND=***value* | *variable*

specifies the upper bound of the censored variables. If *value* is missing or the value in *variable* is missing, no upper bound is set. By default, no upper bound is set.

### Truncated Variable Options

**TRUNCATED** (*truncated-options*)

You can specify the following *truncated-options*:

**LB=***value* | *variable*

**LOWERBOUND=***value* | *variable*

specifies the lower bound of the truncated variables. If *value* is missing or the value in *variable* is missing, no lower bound is set. By default, no lower bound is set.

**UB=***value* | *variable*

**UPPERBOUND=***value* | *variable*

specifies the upper bound of the truncated variables. If *value* is missing or the value in *variable* is missing, no upper bound is set. By default, no upper bound is set.

### Stochastic Frontier Variable Options

**FRONTIER** <(*frontier-options*)>

You can specify the following *frontier-options*:

**TYPE=**HALF | EXPONENTIAL | TRUNCATED

specifies the model type.

**HALF**

specifies half-normal model.

**EXPONENTIAL**

specifies exponential model.

**TRUNCATED**

specifies truncated normal model.

**PRODUCTION**

specifies that the estimated model be a production function.

**COST**

specifies that the estimated model be a cost function.

If neither PRODUCTION nor COST is specified, a production function is estimated by default.

---

## FREQ Statement

**FREQ** *variable* ;

The FREQ statement identifies a variable that contains the frequency of occurrence of each observation. PROC HPQLIM treats each observation as if it appeared  $n$  times, where  $n$  is the value of the FREQ variable for the observation. If the frequency value is not an integer, it is truncated to an integer. If the frequency value is less than 1 or missing, the observation is not used in the model fitting. When the FREQ statement is not specified, each observation is assigned a frequency of 1. If you specify more than one FREQ statement, then the first FREQ statement is used.

---

## HETERO Statement

**HETERO** *dependent variables ~ exogenous variables </ options >* ;

The HETERO statement specifies variables that are related to the heteroscedasticity of the residuals and the way that these variables are used to model the error variance. PROC HPQLIM supports the following heteroscedastic regression model:

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

$$\epsilon_i \sim N(0, \sigma_i^2)$$

For more information about the specification of functional forms, see the section “[Heteroscedasticity](#)” on page 214. The following *options* specify the functional forms of heteroscedasticity:

### LINK=EXP | LINEAR

specifies the functional form.

#### EXP

specifies the exponential link function:

$$\sigma_i^2 = \sigma^2(1 + \exp(\mathbf{z}_i' \boldsymbol{\gamma}))$$

#### LINEAR

specifies the linear link function:

$$\sigma_i^2 = \sigma^2(1 + \mathbf{z}_i' \boldsymbol{\gamma})$$

The default is LINK=EXP.

### NOCONST

specifies that there be no constant in the linear or exponential heteroscedasticity model:

$$\sigma_i^2 = \sigma^2(\mathbf{z}_i' \boldsymbol{\gamma})$$

$$\sigma_i^2 = \sigma^2 \exp(\mathbf{z}_i' \boldsymbol{\gamma})$$

This option is ignored if you do not specify the LINK= option.

**SQUARE**

estimates the model by using the square of the linear heteroscedasticity function. For example, you can specify the following heteroscedasticity function:

$$\sigma_i^2 = \sigma^2(1 + (\mathbf{z}_i' \boldsymbol{\gamma})^2)$$

```
model y = x1 x2 / censored(lb=0);
hetero y ~ z1 / link=linear square;
```

The SQUARE option does not apply to the exponential heteroscedasticity function because the square of an exponential function of  $\mathbf{z}_i' \boldsymbol{\gamma}$  is the same as the exponential of  $2\mathbf{z}_i' \boldsymbol{\gamma}$ . Hence, the only difference is that all  $\boldsymbol{\gamma}$  estimates are divided by two.

This option is ignored if you do not specify the LINK= option. You cannot use the HETERO statement within a Bayesian framework.

**INIT Statement**

**INIT** *initvalue1* < , *initvalue2* ... > ;

The INIT statement sets initial values for parameters in the optimization. You can specify any number of INIT statements.

Each *initvalue* is written as a parameter or parameter list, followed by an optional equality operator (=), followed by a number:

*parameter* <=> *number*

**MODEL Statement**

**MODEL** *dependent* = *regressors* < / *options* > ;

The MODEL statement specifies the dependent variable and independent regressor variables for the regression model.

You can specify the following *option* after a slash (/).

**NOINT**

suppresses the intercept parameter.

You can also specify the following endogenous variable options, which are the same as the options that are specified in the ENDOGENOUS statement. If an endogenous variable option is specified in both the MODEL statement and the ENDOGENOUS statement, the option in the ENDOGENOUS statement is used.

## Discrete Variable Options

**DISCRETE** <(discrete-options) >

specifies that the endogenous variables in this statement be discrete. You can specify the following *discrete-options*:

**DISTRIBUTION**=*distribution-type*

**DIST**=*distribution-type*

**D**=*distribution-type*

specifies the cumulative distribution function that is used to model the response probabilities. You can specify the following *distribution-types*:

**LOGISTIC**      specifies the logistic distribution for the logit model.

**NORMAL**        specifies the normal distribution for the probit model.

By default, DISTRIBUTION=NORMAL.

**ORDER=DATA | FORMATTED | FREQ | INTERNAL**

specifies the sort order for the levels of the discrete variables that are specified in the ENDOGENOUS statement. This ordering determines which parameters in the model correspond to each level in the data. You can specify the following sort orders:

**DATA**            sorts levels by order of appearance in the input data set.

**FORMATTED**     sorts levels by formatted value. The sort order is machine-dependent.

**FREQ**            sorts levels by descending frequency count; levels that have the most observations come first in the order.

**INTERNAL**       sorts levels by unformatted value. The sort order is machine-dependent.

By default, ORDER=FORMATTED. For more information about sort order, see the chapter on the SORT procedure in the *Base SAS Procedures Guide*.

## Censored Variable Options

**CENSORED** <(censored-options) >

specifies that the endogenous variables in this statement be censored. You can specify the following *censored-options*:

**LB**=*value* | *variable*

**LOWERBOUND**=*value* | *variable*

specifies the lower bound of the censored variables. If *value* is missing or the value in *variable* is missing, no lower bound is set. By default, no lower bound is set.

**UB**=*value* | *variable*

**UPPERBOUND**=*value* | *variable*

specifies the upper bound of the censored variables. If *value* is missing or the value in *variable* is missing, no upper bound is set. By default, no upper bound is set.

## Truncated Variable Options

### TRUNCATED <(truncated-options)>

You can specify the following *truncated-options*:

**LB=***value* | *variable*

**LOWERBOUND=***value* | *variable*

specifies the lower bound of the truncated variables. If *value* is missing or the value in *variable* is missing, no lower bound is set. By default, no lower bound is set.

**UB=***value* | *variable*

**UPPERBOUND=***value* | *variable*

specifies the upper bound of the truncated variables. If *value* is missing or the value in *variable* is missing, no upper bound is set. By default, no upper bound is set.

## Stochastic Frontier Variable Options

### FRONTIER <(frontier-options)>

You can specify the following *frontier-options*:

**TYPE=**HALF | EXPONENTIAL | TRUNCATED

specifies the model type.

**HALF**

specifies a half-normal model.

**EXPONENTIAL**

specifies an exponential model.

**TRUNCATED**

specifies a truncated normal model.

**PRODUCTION**

specifies that the estimated model be a production function.

**COST**

specifies that the estimated model be a cost function.

If neither PRODUCTION nor COST is specified, a production function is estimated by default.

---

## OUTPUT Statement

**OUTPUT** **OUT=***SAS-data-set* <*output-options*> ;

The OUTPUT statement creates a new SAS data set to contain variables that are specified with the COPYVAR option and the following data if they are specified by *output-options*: estimates of  $\mathbf{x}'\boldsymbol{\beta}$ , predicted value, residual, marginal effects, probability, standard deviation of the error, expected value, conditional expected value, technical efficiency measures, and inverse Mills ratio. When the response values are missing for the observation, all output estimates except the residual are still computed as long as none of the explanatory

variables are missing. This enables you to compute these statistics for prediction. You can specify only one OUTPUT statement.

You must specify the OUT= option:

**OUT=SAS-data-set**

names the output data set.

**COPYVAR=SAS-variable-names**

**COPYVARS=(SAS-variable-names)**

adds SAS variables to the output data set

You can specify one or more of the following *output-options*:

**CONDITIONAL**

outputs estimates of conditional expected values of continuous endogenous variables.

**ERRSTD**

outputs estimates of  $\sigma_j$ , the standard deviation of the error term.

**EXPECTED**

outputs estimates of expected values of continuous endogenous variables.

**MARGINAL**

outputs marginal effects.

**MILLS**

outputs estimates of inverse Mills ratios of censored or truncated continuous, binary discrete, and selection endogenous variables.

**PREDICTED**

outputs estimates of predicted endogenous variables.

**PROB**

outputs estimates of probability of discrete endogenous variables taking the current observed responses.

**PROBALL**

outputs estimates of probability of discrete endogenous variables for all possible responses.

**RESIDUAL**

outputs estimates of residuals of continuous endogenous variables.

**XBETA**

outputs estimates of  $\mathbf{x}'\boldsymbol{\beta}$ .

**TE1**

outputs estimates of technical efficiency for each producer in the stochastic frontier model that is suggested by Battese and Coelli (1988).

**TE2**

outputs estimates of technical efficiency for each producer in the stochastic frontier model that is suggested by Jondrow et al. (1982).

---

## PERFORMANCE Statement

**PERFORMANCE** < *performance-options* > ;

The PERFORMANCE statement specifies *performance-options* to control the multithreaded and distributed computing environment and requests detailed performance results of the HPQLIM procedure. You can also use the PERFORMANCE statement to control whether the HPQLIM procedure executes in single-machine or distributed mode. You can specify the following *performance-options*:

### DETAILS

requests a table that shows a timing breakdown of the procedure steps.

### NODES=*n*

specifies the number of nodes in the distributed computing environment, provided that the data are not processed alongside the database.

### NTHREADS=*n*

specifies the number of threads for analytic computations and overrides the SAS system option THREADS | NOTHEADS. If you do not specify the NTHREADS= option, PROC HPQLIM creates one thread per CPU for the analytic computations.

The PERFORMANCE statement is documented further in the section “[PERFORMANCE Statement](#)” on page 39 in Chapter 3, “[Shared Concepts and Topics](#).”

---

## PRIOR Statement

**PRIOR** \_REGRESSORS | *parameter-list* ~ *distribution* ;

The PRIOR statement specifies the prior distribution of the model parameters. You must specify one parameter or a list of parameters, a tilde ~, and then a distribution with its parameters. Multiple PRIOR statements are allowed.

You can specify the following *distributions*:

### NORMAL(MEAN= $\mu$ , VAR= $\sigma^2$ )

specifies a normal distribution with the parameters MEAN and VAR.

### GAMMA(SHAPE=*a*, SCALE=*b*)

specifies a gamma distribution with the parameters SHAPE and SCALE.

### IGAMMA(SHAPE=*a*, SCALE=*b*)

specifies an inverse gamma distribution with the parameters SHAPE and SCALE.

### UNIFORM(MIN=*m*, MAX=*M*)

specifies a uniform distribution that is defined between MIN and MAX.

### BETA(SHAPE1=*a*, SHAPE2=*b*, MIN=*m*, MAX=*M*)

specifies a beta distribution with the parameters SHAPE1 and SHAPE2 and defined between MIN and MAX.



**T(LOCATION= $\mu$ , DF= $\nu$ )**

specifies a noncentral  $t$  distribution with DF degrees of freedom and a location parameter equal to LOCATION.

For more information about how to specify *distributions*, see the section “[Standard Distributions](#)” on page 218.

You can specify the special keyword REGRESSORS to select all the parameters that are used in the linear regression component of the model.

---

## RESTRICT Statement

**RESTRICT** *restriction1* <, *restriction2* ... > ;

The RESTRICT statement imposes linear restrictions on the parameter estimates. You can specify any number of RESTRICT statements, but the number of restrictions that are imposed is limited by the number of regressors.

Each *restriction* is written as an expression, followed by an equality operator (=) or an inequality operator (<, >, <=, >=), followed by a second expression:

*expression operator expression*

The *operator* can be =, <, >, <=, or >=. The *operator* and second *expression* are optional.

Restriction expressions can be composed of parameter names; multiplication (\*), addition (+), and subtraction (−) operators; and constants. Parameters that are named in restriction expressions must be among the parameters that are estimated by the model. Parameters that are associated with a regressor variable are referred to by the name of the corresponding regressor variable. The restriction expressions must be a linear function of the parameters.

The following statements illustrate the use of the RESTRICT statement:

```
proc hpqlim data=one;
  model y = x1-x10 / censored(lb=0);
  restrict x1*x2 <= x2 + x3;
run;
```

---

## TEST Statement

<'label':> **TEST** <'string':> *equation* <, *equation* ... > / *options* ;

The TEST statement performs Wald, Lagrange multiplier, and likelihood ratio tests of linear hypotheses about the regression parameters in the preceding MODEL statement. Each equation specifies a linear hypothesis to be tested. All hypotheses in one TEST statement are tested jointly. Variable names in the equations must correspond to regressors in the preceding MODEL statement, and each name represents the coefficient of the corresponding regressor. Use the keyword INTERCEPT for a test that includes a constant.

You can specify the following *options* after the slash (/):

**ALL**

requests Wald, Lagrange multiplier, and likelihood ratio tests.

**LM**

requests the Lagrange multiplier test.

**LR**

requests the likelihood ratio test.

**WALD**

requests the Wald test.

The following statements illustrate the use of the TEST statement (note the use of the INTERCEPT keyword in the second TEST statement):

```
proc hpqlim;
  model y = x1 x2 x3;
  test x1 = 0, x2 * .5 + 2 * x3 = 0;
  test _int: test intercept = 0, x3 = 0;
run;
```

The first TEST statement investigates the joint hypothesis that

$$\beta_1 = 0$$

and

$$0.5\beta_2 + 2\beta_3 = 0$$

Only linear equality restrictions and tests are permitted in PROC HPQLIM. Test expressions can be composed only of algebraic operations that involve the addition symbol (+), subtraction symbol (–), and multiplication symbol (\*).

The TEST statement accepts labels that are reproduced in the printed output. You can label a TEST statement in two ways: you can specify a label followed by a colon before the TEST keyword, or you can specify a quoted string after the TEST keyword. If you specify both a label before the TEST keyword and a quoted string after the keyword, PROC HPQLIM uses the label that precedes the colon. If no label or quoted string is specified, PROC HPQLIM labels the test automatically.

---

## WEIGHT Statement

**WEIGHT** *variable* *</ option>* ;

The WEIGHT statement specifies a variable that supplies weighting values to use for each observation in estimating parameters. The log likelihood for each observation is multiplied by the corresponding weight variable value.

If the weight of an observation is nonpositive, that observation is not used in the estimation.

You can add the following *option* after a slash (/):

**NONNORMALIZE**

specifies that the weights must be used as is. When this option is not specified, the weights are normalized so that they add up to the actual sample size. Weights  $w_i$  are normalized by multiplying them by  $\frac{n}{\sum_{i=1}^n w_i}$ , where  $n$  is the sample size.

---

## Details: HPQLIM Procedure

---

### Ordinal Discrete Choice Modeling

#### Binary Probit and Logit Model

The binary choice model is

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

where the value of the latent dependent variable,  $y_i^*$ , is observed only as follows:

$$\begin{aligned} y_i &= 1 && \text{if } y_i^* > 0 \\ &= 0 && \text{otherwise} \end{aligned}$$

The disturbance,  $\epsilon_i$ , of the probit model has a standard normal distribution with the distribution function (CDF)

$$\Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} \exp(-t^2/2) dt$$

The disturbance of the logit model has a standard logistic distribution with the distribution function (CDF)

$$\Lambda(x) = \frac{\exp(x)}{1 + \exp(x)} = \frac{1}{1 + \exp(-x)}$$

The binary discrete choice model has the following probability that the event  $\{y_i = 1\}$  occurs:

$$P(y_i = 1) = F(\mathbf{x}_i' \boldsymbol{\beta}) = \begin{cases} \Phi(\mathbf{x}_i' \boldsymbol{\beta}) & \text{(probit)} \\ \Lambda(\mathbf{x}_i' \boldsymbol{\beta}) & \text{(logit)} \end{cases}$$

For more information, see the section “Ordinal Discrete Choice Modeling” (Chapter 22, *SAS/ETS User's Guide*).

#### Ordinal Probit/Logit

When the dependent variable is observed in sequence with  $M$  categories, binary discrete choice modeling is not appropriate for data analysis. McKelvey and Zavoina (1975) propose the ordinal (or ordered) probit model.

Consider the regression equation

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

where error disturbances,  $\epsilon_i$ , have the distribution function  $F$ . The unobserved continuous random variable,  $y_i^*$ , is identified as  $M$  categories. Suppose there are  $M + 1$  real numbers,  $\mu_0, \dots, \mu_M$ , where  $\mu_0 = -\infty$ ,  $\mu_1 = 0$ ,  $\mu_M = \infty$ , and  $\mu_0 \leq \mu_1 \leq \dots \leq \mu_M$ . Define

$$R_{i,j} = \mu_j - \mathbf{x}_i' \boldsymbol{\beta}$$

The probability that the unobserved dependent variable is contained in the  $j$ th category can be written as

$$P[\mu_{j-1} < y_i^* \leq \mu_j] = F(R_{i,j}) - F(R_{i,j-1})$$

For more information, see the section “Ordinal Discrete Choice Modeling” (Chapter 22, *SAS/ETS User's Guide*).

## Limited Dependent Variable Models

### Censored Regression Models

When the dependent variable is censored, values in a certain range are all transformed to a single value. For example, the standard Tobit model can be defined as

$$y_i^* = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$$

$$y_i = \begin{cases} y_i^* & \text{if } y_i^* > 0 \\ 0 & \text{if } y_i^* \leq 0 \end{cases}$$

where  $\epsilon_i \sim iidN(0, \sigma^2)$ .

The Tobit model can be generalized to handle observation-by-observation censoring. The censored model on both the lower and upper limits can be defined as

$$y_i = \begin{cases} R_i & \text{if } y_i^* \geq R_i \\ y_i^* & \text{if } L_i < y_i^* < R_i \\ L_i & \text{if } y_i^* \leq L_i \end{cases}$$

You can see Chapter 22.7, “Censored Regression Models” (*SAS/ETS User's Guide*), for more details.

### Truncated Regression Models

In a truncated model, the observed sample is a subset of the population where the dependent variable falls within a certain range. For example, when neither a dependent variable nor exogenous variables are observed for  $y_i^* \leq 0$ , the truncated regression model can be specified as

$$\ell = \sum_{i \in \{y_i > 0\}} \left\{ -\ln \Phi(\mathbf{x}_i' \boldsymbol{\beta} / \sigma) + \ln \left[ \frac{\phi((y_i - \mathbf{x}_i' \boldsymbol{\beta}) / \sigma)}{\sigma} \right] \right\}$$

For more information, see the section “Truncated Regression Models” (Chapter 22, *SAS/ETS User's Guide*).

## Stochastic Frontier Production and Cost Models

Stochastic frontier production models were first developed by Aigner, Lovell, and Schmidt (1977); Meeusen and van den Broeck (1977). Specification of these models allow for random shocks of the production or cost but also include a term for technical or cost inefficiency. Assuming that the production function takes a log-linear Cobb-Douglas form, the stochastic frontier production model can be written as

$$\ln(y_i) = \beta_0 + \sum_n \beta_n \ln(x_{ni}) + \epsilon_i$$

where  $\epsilon_i = v_i - u_i$ . The  $v_i$  term represents the stochastic error component, and the  $u_i$  term represents the nonnegative, technical inefficiency error component. The  $v_i$  error component is assumed to be distributed iid normal and independent from  $u_i$ . If  $u_i > 0$ , the error term  $\epsilon_i$  is negatively skewed and represents technical inefficiency. If  $u_i < 0$ , the error term  $\epsilon_i$  is positively skewed and represents cost inefficiency. PROC HPQLIM models the  $u_i$  error component as a half-normal, exponential, or truncated normal distribution.

### The Normal-Half-Normal Model

When  $v_i$  is iid  $N(0, \sigma_v^2)$  in a normal-half-normal model,  $u_i$  is iid  $N^+(0, \sigma_u^2)$ , with  $v_i$  and  $u_i$  independent of each other. Given the independence of error terms, the joint density of  $v$  and  $u$  can be written as

$$f(u, v) = \frac{2}{2\pi\sigma_u\sigma_v} \exp \left\{ -\frac{u^2}{2\sigma_u^2} - \frac{v^2}{2\sigma_v^2} \right\}$$

Substituting  $v = \epsilon + u$  into the preceding equation and integrating  $u$  out gives

$$f(\epsilon) = \frac{2}{\sigma} \phi \left( \frac{\epsilon}{\sigma} \right) \Phi \left( -\frac{\epsilon\lambda}{\sigma} \right)$$

where  $\lambda = \sigma_u/\sigma_v$  and  $\sigma = \sqrt{\sigma_u^2 + \sigma_v^2}$ .

In the case of a stochastic frontier cost model,  $v = \epsilon - u$  and

$$f(\epsilon) = \frac{2}{\sigma} \phi \left( \frac{\epsilon}{\sigma} \right) \Phi \left( \frac{\epsilon\lambda}{\sigma} \right)$$

For more information, see the section “Stochastic Frontier Production and Cost Models” (Chapter 22, *SAS/ETS User's Guide*).

### The Normal-Exponential Model

Under the normal-exponential model,  $v_i$  is iid  $N(0, \sigma_v^2)$  and  $u_i$  is iid exponential. Given the independence of error term components  $u_i$  and  $v_i$ , the joint density of  $v$  and  $u$  can be written as

$$f(u, v) = \frac{1}{\sqrt{2\pi}\sigma_u\sigma_v} \exp \left\{ -\frac{u}{\sigma_u} - \frac{v^2}{2\sigma_v^2} \right\}$$

The marginal density function of  $\epsilon$  for the production function is

$$f(\epsilon) = \left( \frac{1}{\sigma_u} \right) \Phi \left( -\frac{\epsilon}{\sigma_v} - \frac{\sigma_v}{\sigma_u} \right) \exp \left\{ \frac{\epsilon}{\sigma_u} + \frac{\sigma_v^2}{2\sigma_u^2} \right\}$$

The marginal density function for the cost function is equal to

$$f(\epsilon) = \left( \frac{1}{\sigma_u} \right) \Phi \left( \frac{\epsilon}{\sigma_v} - \frac{\sigma_v}{\sigma_u} \right) \exp \left\{ -\frac{\epsilon}{\sigma_u} + \frac{\sigma_v^2}{2\sigma_u^2} \right\}$$

For more information, see the section “Stochastic Frontier Production and Cost Models” (Chapter 22, *SAS/ETS User's Guide*).

### The Normal–Truncated Normal Model

The normal–truncated normal model is a generalization of the normal-half-normal model that allows the mean of  $u_i$  to differ from zero. Under the normal–truncated normal model, the error term component  $v_i$  is iid  $N^+(0, \sigma_v^2)$  and  $u_i$  is iid  $N(\mu, \sigma_u^2)$ . The joint density of  $v_i$  and  $u_i$  can be written as

$$f(u, v) = \frac{1}{\sqrt{2\pi}\sigma_u\sigma_v\Phi(\mu/\sigma_u)} \exp \left\{ -\frac{(u - \mu)^2}{2\sigma_u^2} - \frac{v^2}{2\sigma_v^2} \right\}$$

The marginal density function of  $\epsilon$  for the production function is

$$f(\epsilon) = \frac{1}{\sigma} \phi \left( \frac{\epsilon + \mu}{\sigma} \right) \Phi \left( \frac{\mu}{\sigma\lambda} - \frac{\epsilon\lambda}{\sigma} \right) \left[ \Phi \left( \frac{\mu}{\sigma_u} \right) \right]^{-1}$$

The marginal density function for the cost function is

$$f(\epsilon) = \frac{1}{\sigma} \phi \left( \frac{\epsilon - \mu}{\sigma} \right) \Phi \left( \frac{\mu}{\sigma\lambda} + \frac{\epsilon\lambda}{\sigma} \right) \left[ \Phi \left( \frac{\mu}{\sigma_u} \right) \right]^{-1}$$

For more information, see the section “Stochastic Frontier Production and Cost Models” (Chapter 22, *SAS/ETS User's Guide*).

For more information about normal-half-normal, normal-exponential, and normal–truncated normal models, see Kumbhakar and Lovell (2000); Coelli, Prasada Rao, and Battese (1998).

---

## Heteroscedasticity

If the variance of regression disturbance,  $(\epsilon_i)$ , is heteroscedastic, the variance can be specified as a function of variables

$$E(\epsilon_i^2) = \sigma_i^2 = f(\mathbf{z}_i' \boldsymbol{\gamma})$$

Table 8.2 shows various functional forms of heteroscedasticity and the corresponding options to request each model.

**Table 8.2** Specification Summary for Modeling Heteroscedasticity

Number	Model	Options
1	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2(1 + \exp(\mathbf{z}'_i \boldsymbol{\gamma}))$	LINK=EXP (default)
2	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2 \exp(\mathbf{z}'_i \boldsymbol{\gamma})$	LINK=EXP NOCONST
3	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2(1 + \sum_{l=1}^L \gamma_l z_{li})$	LINK=LINEAR
4	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2(1 + (\sum_{l=1}^L \gamma_l z_{li})^2)$	LINK=LINEAR SQUARE
5	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2(\sum_{l=1}^L \gamma_l z_{li})$	LINK=LINEAR NOCONST
6	$f(\mathbf{z}'_i \boldsymbol{\gamma}) = \sigma^2((\sum_{l=1}^L \gamma_l z_{li})^2)$	LINK=LINEAR SQUARE NOCONST

In models 3 and 5, variances of some observations might be negative. Although the HPQLIM procedure assigns a large penalty to move the optimization away from such a region, the optimization might not be able to improve the objective function value and might become locked in the region. Signs of such an outcome include extremely small likelihood values or missing standard errors in the estimates. In models 2 and 6, variances are guaranteed to be greater than or equal to zero, but variances of some observations might be very close to 0. In these scenarios, standard errors might be missing. Models 1 and 4 do not have such problems. Variances in these models are always positive and never close to 0.

For more information, see the section “Heteroscedasticity and Box-Cox Transformation” (Chapter 22, *SAS/ETS User's Guide*).

## Tests on Parameters

In general, the tested hypothesis can be written as

$$H_0 : \mathbf{h}(\theta) = 0$$

where  $\mathbf{h}(\theta)$  is an  $r \times 1$  vector-valued function of the parameters  $\theta$  given by the  $r$  expressions that are specified in the TEST statement.

Let  $\hat{V}$  be the estimate of the covariance matrix of  $\hat{\theta}$ . Let  $\hat{\theta}$  be the unconstrained estimate of  $\theta$  and  $\tilde{\theta}$  be the constrained estimate of  $\theta$  such that  $\mathbf{h}(\tilde{\theta}) = 0$ . Let

$$A(\theta) = \partial \mathbf{h}(\theta) / \partial \theta \big|_{\hat{\theta}}$$

Using this notation, the test statistics for the three types of tests are computed as follows.

- The Wald test statistic is defined as

$$W = \mathbf{h}'(\hat{\theta}) \left( A(\hat{\theta}) \hat{V} A'(\hat{\theta}) \right)^{-1} \mathbf{h}(\hat{\theta})$$

- The Lagrange multiplier test statistic is

$$LM = \lambda' A(\tilde{\theta}) \tilde{V} A'(\tilde{\theta}) \lambda$$

where  $\lambda$  is the vector of Lagrange multipliers from the computation of the restricted estimate  $\tilde{\theta}$ .

- The likelihood ratio test statistic is

$$LR = 2 \left( L(\hat{\theta}) - L(\tilde{\theta}) \right)$$

where  $\tilde{\theta}$  represents the constrained estimate of  $\theta$  and  $L$  is the concentrated log-likelihood value.

The following statements use the TEST statement to perform a likelihood ratio test:

```
proc hpqlim;
  model y = x1 x2 x3;
  test x1 = 0, x2 * .5 + 2 * x3 = 0 /lr;
run;
```

For more information, see the section “Tests on Parameters” (Chapter 22, *SAS/ETS User’s Guide*).

---

## Bayesian Analysis

To perform Bayesian analysis, you must specify a BAYES statement. Unless otherwise stated, all options that are described in this section are options in the BAYES statement.

By default, PROC HPQLIM uses the random walk Metropolis algorithm to obtain posterior samples. For the implementation details of the Metropolis algorithm in PROC HPQLIM, such as the blocking of the parameters and tuning of the covariance matrices, see the sections “[Blocking of Parameters](#)” on page 216 and “[Tuning the Proposal Distribution](#)” on page 216.

The Bayes theorem states that

$$p(\theta|y) \propto \pi(\theta)L(y|\theta)$$

where  $\theta$  is a parameter or a vector of parameters and  $\pi(\theta)$  is the product of the prior densities that are specified in the **PRIOR** statement. The term  $L(y|\theta)$  is the likelihood that is associated with the **MODEL** statement.

## Blocking of Parameters

In a multivariate parameter model, all the parameters are updated in one single block (by default or when you specify the **SAMPLING=MULTIMETROPOLIS** option). This can be inefficient, especially when parameters have vastly different scales. As an alternative, you can update the parameters one at a time (by specifying **SAMPLING=UNIMETROPOLIS**).

## Tuning the Proposal Distribution

One key factor in achieving high efficiency of a Metropolis-based Markov chain is finding a good proposal distribution for each block of parameters. This process is called tuning. The tuning phase consists of a number of loops that are controlled by the options **MINTUNE** and **MAXTUNE**. The **MINTUNE=** option controls the minimum number of tuning loops and has a default value of 2. The **MAXTUNE=** option controls the maximum number of tuning loops and has a default value of 24. Each loop repeats the number of times specified by the **NTU=** option, which has a default of 500. At the end of every loop, PROC HPQLIM



examines the acceptance probability for each block. The acceptance probability is the percentage of NTU proposed values that have been accepted. If this probability does not fall within the acceptance tolerance range (see the following section), the proposal distribution is modified before the next tuning loop.

A good proposal distribution should resemble the actual posterior distribution of the parameters. Large sample theory states that the posterior distribution of the parameters approaches a multivariate normal distribution (Gelman et al. 2004, Appendix B; Schervish 1995, Section 7.4). That is why a normal proposal distribution often works well in practice. The default proposal distribution in PROC HPQLIM is the normal distribution.

You can see Chapter 22.7, “Bayesian Analysis” (*SAS/ETS User’s Guide*), for more details.

## Initial Values of the Markov Chains

You can assign initial values to any parameters. For more information, see the **INIT** statement. If you use the optimization **PROPCOV=** option, PROC HPQLIM starts the tuning at the optimized values. This option overwrites the provided initial values.

## Prior Distributions

The **PRIOR** statement specifies the prior distribution of the model parameters. You must specify one parameter or a list of parameters, a tilde  $\sim$ , and then a distribution with its parameters. You can specify multiple **PRIOR** statements to define independent priors. Parameters that are associated with a regressor variable are referred to by the name of the corresponding regressor variable.

You can specify the special keyword **\_REGRESSORS** to consider all the regressors of a model. If multiple **PRIOR** statements affect the same parameter, the last **PRIOR** statement prevails. For example, in a regression with two regressors (X1, X2), the following statements imply that the prior on X1 is **NORMAL**(MEAN=0, VAR=1), the prior on X2 is **GAMMA**(SHAPE=3, SCALE=4).

```
...
prior _Regressors ~ uniform(min=0, max=1);
prior X1 X2 ~ gamma(shape=3, scale=4);
prior X1 ~ normal(mean=0, var=1);
...
```

If a parameter is not associated with a **PRIOR** statement or if some of the prior hyperparameters are missing, then the default choices in Table 8.3 are considered.

**Table 8.3** Default Values for Prior Distributions

PRIOR Distribution	Hyperparameter <sub>1</sub>	Hyperparameter <sub>2</sub>	Min	Max	Parameters Default Choice
NORMAL	MEAN=0	VAR=1E6	$-\infty$	$\infty$	Regression-Location-Threshold
IGAMMA	SHAPE=2.000001	SCALE=1	$> 0$	$\infty$	Scale
GAMMA	SHAPE=1	SCALE=1	0	$\infty$	
UNIFORM			$-\infty$	$\infty$	
BETA	SHAPE1=1	SHAPE2=1	$-\infty$	$\infty$	
T	LOCATION=0	DF=3	$-\infty$	$\infty$	

For density specification, see the section “**Standard Distributions**” on page 218.

## Standard Distributions

Table 8.4 through Table 8.9 show all the distribution density functions that PROC HPQLIM recognizes. You specify these distribution densities in the **PRIOR** statement.

**Table 8.4 Beta Distribution**

PRIOR statement	BETA(SHAPE1= $a$ , SHAPE2= $b$ , MIN= $m$ , MAX= $M$ )
	Note: Commonly $m = 0$ and $M = 1$ .
Density	$\frac{(\theta-m)^{a-1}(M-\theta)^{b-1}}{B(a,b)(M-m)^{a+b-1}}$
Parameter restriction	$a > 0, b > 0, -\infty < m < M < \infty$
Range	$\begin{cases} [m, M] & \text{when } a = 1, b = 1 \\ [m, M) & \text{when } a = 1, b \neq 1 \\ (m, M] & \text{when } a \neq 1, b = 1 \\ (m, M) & \text{otherwise} \end{cases}$
Mean	$\frac{a}{a+b} \times (M - m) + m$
Variance	$\frac{ab}{(a+b)^2(a+b+1)} \times (M - m)^2$
Mode	$\begin{cases} \frac{a-1}{a+b-2} \times M + \frac{b-1}{a+b-2} \times m & a > 1, b > 1 \\ m \text{ and } M & a < 1, b < 1 \\ m & \begin{cases} a < 1, b \geq 1 \\ a = 1, b > 1 \end{cases} \\ M & \begin{cases} a \geq 1, b < 1 \\ a > 1, b = 1 \end{cases} \\ \text{not unique} & a = b = 1 \end{cases}$
Defaults	SHAPE1=SHAPE2=1, MIN $\rightarrow -\infty$ , MAX $\rightarrow \infty$

**Table 8.5 Gamma Distribution**

PRIOR statement	GAMMA(SHAPE= $a$ , SCALE= $b$ )
Density	$\frac{1}{b^a \Gamma(a)} \theta^{a-1} e^{-\theta/b}$
Parameter restriction	$a > 0, b > 0$
Range	$[0, \infty)$
Mean	$ab$
Variance	$ab^2$
Mode	$(a - 1)b$
Defaults	SHAPE=SCALE=1

**Table 8.6 Inverse Gamma Distribution**

PRIOR statement	IGAMMA(SHAPE= $a$ , SCALE= $b$ )
Density	$\frac{b^a}{\Gamma(a)} \theta^{-(a+1)} e^{-b/\theta}$
Parameter restriction	$a > 0, b > 0$
Range	$0 < \theta < \infty$
Mean	$\frac{b}{a-1}, \quad a > 1$
Variance	$\frac{b^2}{(a-1)^2(a-2)}, \quad a > 2$
Mode	$\frac{b}{a+1}$
Defaults	SHAPE=2.000001, SCALE=1

**Table 8.7 Normal Distribution**

PRIOR statement	NORMAL(MEAN= $\mu$ , VAR= $\sigma^2$ )
Density	$\frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(\theta-\mu)^2}{2\sigma^2}\right)$
Parameter restriction	$\sigma^2 > 0$
Range	$-\infty < \theta < \infty$
Mean	$\mu$
Variance	$\sigma^2$
Mode	$\mu$
Defaults	MEAN=0, VAR=1000000

**Table 8.8  $t$  Distribution**

PRIOR statement	T(LOCATION= $\mu$ , DF= $\nu$ )
Density	$\frac{\Gamma(\frac{\nu+1}{2})}{\Gamma(\frac{\nu}{2})\sqrt{\pi\nu}} \left[1 + \frac{(\theta-\mu)^2}{\nu}\right]^{-\frac{\nu+1}{2}}$
Parameter restriction	$\nu > 0$
Range	$-\infty < \theta < \infty$
Mean	$\mu$ , for $\nu > 1$
Variance	$\frac{\nu}{\nu-2}$ , for $\nu > 2$
Mode	$\mu$
Defaults	LOCATION=0, DF=3

**Table 8.9 Uniform Distribution**

PRIOR statement	UNIFORM(MIN= $m$ , MAX= $M$ )
Density	$\frac{1}{M-m}$

Parameter restriction	$-\infty < m < M < \infty$
Range	$\theta \in [m, M]$
Mean	$\frac{m+M}{2}$
Variance	$\frac{(M-m)^2}{12}$
Mode	Not unique
Defaults	MIN $\rightarrow -\infty$ , MAX $\rightarrow \infty$

---

## Output to SAS Data Set

### XBeta, Predicted, and Residual

Xbeta is the structural part on the right-hand side of the model. The predicted value is the predicted dependent variable value. For censored variables, if the predicted value is outside the boundaries, it is reported as the closest boundary. The residual is defined only for continuous variables and is defined as

$$\text{Residual} = \text{Observed} - \text{Predicted}$$

### Error Standard Deviation

The error standard deviation is  $\sigma_i$  in the model. It varies only when the HETERO statement is used.

### Marginal Effects

A marginal effect is defined as a contribution of one control variable to the response variable. For a binary choice model with two response categories,  $\mu_0 = -\infty$  and  $\mu_1 = 0$ ,  $\mu_2 = \infty$ . For an ordinal response model with  $M$  response categories ( $\mu_0, \dots, \mu_M$ ), define

$$R_{i,j} = \mu_j - \mathbf{x}_i' \boldsymbol{\beta}$$

The probability that the unobserved dependent variable is contained in the  $j$ th category can be written as

$$P[\mu_{j-1} < y_i^* \leq \mu_j] = F(R_{i,j}) - F(R_{i,j-1})$$

The marginal effect of changes in the regressors on the probability of  $y_i = j$  is then

$$\frac{\partial \text{Prob}[y_i = j]}{\partial \mathbf{x}} = [f(\mu_{j-1} - \mathbf{x}_i' \boldsymbol{\beta}) - f(\mu_j - \mathbf{x}_i' \boldsymbol{\beta})] \boldsymbol{\beta}$$

where  $f(x) = \frac{dF(x)}{dx}$ . In particular,

$$f(x) = \frac{dF(x)}{dx} = \begin{cases} \frac{1}{\sqrt{2\pi}} e^{-x^2/2} & (\text{probit}) \\ \frac{e^{-x}}{[1+e^{(-x)}]^2} & (\text{logit}) \end{cases}$$

The marginal effects in the truncated regression model are

$$\frac{\partial E[y_i | L_i < y_i^* < R_i]}{\partial \mathbf{x}} = \boldsymbol{\beta} \left[ 1 - \frac{(\phi(a_i) - \phi(b_i))^2}{(\Phi(b_i) - \Phi(a_i))^2} + \frac{a_i \phi(a_i) - b_i \phi(b_i)}{\Phi(b_i) - \Phi(a_i)} \right]$$

where  $a_i = \frac{L_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma_i}$  and  $b_i = \frac{R_i - \mathbf{x}_i' \boldsymbol{\beta}}{\sigma_i}$ .

The marginal effects in the censored regression model are

$$\frac{\partial E[y|\mathbf{x}_i]}{\partial \mathbf{x}} = \boldsymbol{\beta} \times \text{Prob}[L_i < y_i^* < R_i]$$

### Expected and Conditionally Expected Values

The expected value is the unconditional expectation of the dependent variable. For a censored variable, it is

$$E[y_i] = \Phi(a_i)L_i + (\mathbf{x}_i' \boldsymbol{\beta} + \lambda \sigma_i)(\Phi(b_i) - \Phi(a_i)) + (1 - \Phi(b_i))R_i$$

For a left-censored variable ( $R_i = \infty$ ), this formula is

$$E[y_i] = \Phi(a_i)L_i + (\mathbf{x}_i' \boldsymbol{\beta} + \lambda \sigma_i)(1 - \Phi(a_i))$$

where  $\lambda = \frac{\phi(a_i)}{1 - \Phi(a_i)}$ .

For a right-censored variable ( $L_i = -\infty$ ), this formula is

$$E[y_i] = (\mathbf{x}_i' \boldsymbol{\beta} + \lambda \sigma_i)\Phi(b_i) + (1 - \Phi(b_i))R_i$$

where  $\lambda = -\frac{\phi(b_i)}{\Phi(b_i)}$ .

For a noncensored variable, this formula is

$$E[y_i] = \mathbf{x}_i' \boldsymbol{\beta}$$

The conditional expected value is the expectation when the variable is inside the boundaries:

$$E[y_i | L_i < y_i < R_i] = \mathbf{x}_i' \boldsymbol{\beta} + \lambda \sigma_i$$

### Technical Efficiency

Technical efficiency for each producer is computed only for stochastic frontier models.

In general, the stochastic production frontier can be written as

$$y_i = f(x_i; \boldsymbol{\beta}) \exp\{v_i\} TE_i$$

where  $y_i$  denotes producer  $i$ 's actual output,  $f(\cdot)$  is the deterministic part of the production frontier,  $\exp\{v_i\}$  is a producer-specific error term, and  $TE_i$  is the technical efficiency coefficient, which can be written as

$$TE_i = \frac{y_i}{f(x_i; \boldsymbol{\beta}) \exp\{v_i\}}$$

For a Cobb-Douglas production function,  $TE_i = \exp\{-u_i\}$ . For more information, see the section “[Stochastic Frontier Production and Cost Models](#)” on page 213.

The cost frontier can be written in general as

$$E_i = c(y_i, w_i; \boldsymbol{\beta}) \exp\{v_i\} / CE_i$$

where  $w_i$  denotes producer  $i$ 's input prices,  $c(\cdot)$  is the deterministic part of the cost frontier,  $\exp\{v_i\}$  is a producer-specific error term, and  $CE_i$  is the cost efficiency coefficient, which can be written as

$$CE_i = \frac{c(x_i, w_i; \beta) \exp\{v_i\}}{E_i}$$

For a Cobb-Douglas cost function,  $CE_i = \exp\{-u_i\}$ . For more information, see the section “[Stochastic Frontier Production and Cost Models](#)” on page 213. Hence, both technical and cost efficiency coefficients are the same. The estimates of technical efficiency are provided in the following subsections.

### Normal-Half-Normal Model

Define  $\mu_* = -\epsilon\sigma_u^2/\sigma^2$  and  $\sigma_*^2 = \sigma_u^2\sigma_v^2/\sigma^2$ . Then, as shown by Jondrow et al. (1982), conditional density is as follows:

$$f(u|\epsilon) = \frac{f(u, \epsilon)}{f(\epsilon)} = \frac{1}{\sqrt{2\pi}\sigma_*} \exp\left\{-\frac{(u - \mu_*)^2}{2\sigma_*^2}\right\} \Bigg/ \left[1 - \Phi\left(-\frac{\mu_*}{\sigma_*}\right)\right]$$

Hence,  $f(u|\epsilon)$  is the density for  $N^+(\mu_*, \sigma_*^2)$ .

From this result, it follows that the estimate of technical efficiency (Battese and Coelli 1988) is

$$TE1_i = E(\exp\{-u_i\}|\epsilon_i) = \left[\frac{1 - \Phi(\sigma_* - \mu_{*i}/\sigma_*)}{1 - \Phi(-\mu_{*i}/\sigma_*)}\right] \exp\left\{-\mu_{*i} + \frac{1}{2}\sigma_*^2\right\}$$

The second version of the estimate (Jondrow et al. 1982) is

$$TE2_i = \exp\{-E(u_i|\epsilon_i)\}$$

where

$$E(u_i|\epsilon_i) = \mu_{*i} + \sigma_* \left[\frac{\phi(-\mu_{*i}/\sigma_*)}{1 - \Phi(-\mu_{*i}/\sigma_*)}\right] = \sigma_* \left[\frac{\phi(\epsilon_i \lambda / \sigma)}{1 - \Phi(\epsilon_i \lambda / \sigma)} - \left(\frac{\epsilon_i \lambda}{\sigma}\right)\right]$$

### Normal-Exponential Model

Define  $A = -\tilde{\mu}/\sigma_v$  and  $\tilde{\mu} = -\epsilon - \sigma_v^2/\sigma_u$ . Then, as shown by Kumbhakar and Lovell (2000), conditional density is as follows:

$$f(u|\epsilon) = \frac{1}{\sqrt{2\pi}\sigma_v\Phi(-\tilde{\mu}/\sigma_v)} \exp\left\{-\frac{(u - \tilde{\mu})^2}{2\sigma^2}\right\}$$

Hence,  $f(u|\epsilon)$  is the density for  $N^+(\tilde{\mu}, \sigma_v^2)$ .

From this result, it follows that the estimate of technical efficiency is

$$TE1_i = E(\exp\{-u_i\}|\epsilon_i) = \left[\frac{1 - \Phi(\sigma_v - \tilde{\mu}_i/\sigma_v)}{1 - \Phi(-\tilde{\mu}_i/\sigma_v)}\right] \exp\left\{-\tilde{\mu}_i + \frac{1}{2}\sigma_v^2\right\}$$

The second version of the estimate is

$$TE2_i = \exp\{-E(u_i|\epsilon_i)\}$$

where

$$E(u_i|\epsilon_i) = \tilde{\mu}_i + \sigma_v \left[\frac{\phi(-\tilde{\mu}_i/\sigma_v)}{1 - \Phi(-\tilde{\mu}_i/\sigma_v)}\right] = \sigma_v \left[\frac{\phi(A)}{\Phi(-A)} - A\right]$$

### Normal-Truncated Normal Model

Define  $\tilde{\mu} = (-\sigma_u^2 \epsilon_i + \mu \sigma_v^2) / \sigma^2$  and  $\sigma_*^2 = \sigma_u^2 \sigma_v^2 / \sigma^2$ . Then, as shown by Kumbhakar and Lovell (2000), conditional density is as follows:

$$f(u|\epsilon) = \frac{1}{\sqrt{2\pi}\sigma_*[1 - \Phi(-\tilde{\mu}/\sigma_*)]} \exp \left\{ -\frac{(u - \tilde{\mu})^2}{2\sigma_*^2} \right\}$$

Hence,  $f(u|\epsilon)$  is the density for  $N^+(\tilde{\mu}, \sigma_*^2)$ .

From this result, it follows that the estimate of technical efficiency is

$$TE1_i = E(\exp\{-u_i\}|\epsilon_i) = \frac{1 - \Phi(\sigma_* - \tilde{\mu}_i/\sigma_*)}{1 - \Phi(-\tilde{\mu}_i/\sigma_*)} \exp \left\{ -\tilde{\mu}_i + \frac{1}{2}\sigma_*^2 \right\}$$

The second version of the estimate is

$$TE2_i = \exp\{-E(u_i|\epsilon_i)\}$$

where

$$E(u_i|\epsilon_i) = \tilde{\mu}_i + \sigma_* \left[ \frac{\phi(\tilde{\mu}_i/\sigma_*)}{1 - \Phi(-\tilde{\mu}_i/\sigma_*)} \right]$$

---

## OUTEST= Data Set

The OUTEST= data set contains all the parameters that are estimated by a MODEL statement. Each parameter contains the estimate for the corresponding parameter in the corresponding model. In addition, the OUTEST= data set contains the following variables:

<code>_NAME_</code>	indicates the name of the independent variable.
<code>_TYPE_</code>	indicates the type of observation. PARM indicates the row of coefficients; STD indicates the row of standard deviations of the corresponding coefficients.
<code>_STATUS_</code>	indicates the convergence status for optimization.

The rest of the columns correspond to the explanatory variables.

The OUTEST= data set contains one observation for the MODEL statement, which shows the parameter estimates for that model. If you specify the COVOUT option in the PROC HPQLIM statement, the OUTEST= data set includes additional observations for the MODEL statement, which show the rows of the covariance matrix of parameter estimates. For covariance observations, the value of the `_TYPE_` variable is COV, and the `_NAME_` variable identifies the parameter that is associated with that row of the covariance matrix. If you specify the CORROUT option in the PROC HPQLIM statement, the OUTEST= data set includes additional observations for the MODEL statement, which show the rows of the correlation matrix of parameter estimates. For correlation observations, the value of the `_TYPE_` variable is CORR, and the `_NAME_` variable identifies the parameter that is associated with that row of the correlation matrix.

## Naming

### Naming of Parameters

The parameters are named in the same way as in other SAS procedures such as the REG and PROBIT procedures. The constant in the regression equation is called Intercept. The coefficients of independent variables are named by the independent variables. The standard deviation of the errors is called `_Sigma`. If the HETERO statement is included, the coefficients of the independent variables in the HETERO statement are called `_H.x`, where *x* is the name of the independent variable.

### Naming of Output Variables

Table 8.10 shows the *options* in the OUTPUT statement, with the corresponding variable names and their explanations.

**Table 8.10** OUTPUT Statement Options

<i>output-option</i>	<b>Variable Name</b>	<b>Explanation</b>
CONDITIONAL	CEXPCT_y	Conditional expected value of <i>y</i> , conditioned on the truncation
ERRSTD	ERRSTD_y	Standard deviation of error term
EXPECTED	EXPCT_y	Unconditional expected value of <i>y</i>
MARGINAL	MEFF_x	Marginal effect of <i>x</i> on <i>y</i> ( $\frac{\partial y}{\partial x}$ ) with single equation
PREDICTED	P_y	Predicted value of <i>y</i>
RESIDUAL	RESID_y	Residual of <i>y</i> , ( <i>y</i> – PredictedY)
PROB	PROB_y	Probability that <i>y</i> is taking the observed value in this observation (discrete <i>y</i> only)
PROBALL	PROB <sub><i>i</i></sub> _y	Probability that <i>y</i> is taking the <i>i</i> th value (discrete <i>y</i> only)
MILLS	MILLS_y	Inverse Mills ratio for <i>y</i>
TE1	TE1	Technical efficiency estimate for each producer proposed by Battese and Coelli (1988)
TE2	TE2	Technical efficiency estimate for each producer proposed by Jondrow et al. (1982)
XBETA	XBETA_y	Structure part ( $\mathbf{x}'\boldsymbol{\beta}$ ) of <i>y</i> equation

If you prefer to name the output variables differently, you can use the RENAME option in the data set. For example, the following statements rename the residual of *y* as *Resid*:

```
proc hpqlim data=one;
  model y = x1-x10 / censored;
  output out=outds(rename=(resid_y=resid)) residual;
run;
```



## ODS Table Names

PROC HPQLIM assigns a name to each table that it creates. You can use these names to refer to the table when you use the Output Delivery System (ODS) to select tables and create output data sets. These names are listed in Table 8.11.

**Table 8.11** ODS Tables Produced in PROC HPQLIM

ODS Table Name	Description	Option
<b>ODS Tables Created by the MODEL Statement and TEST Statement</b>		
ResponseProfile	Response profile	Default
FitSummary	Summary of nonlinear estimation	Default
ParameterEstimates	Parameter estimates	Default
SummaryContResponse	Summary of continuous response	Default
CovB	Covariance of parameter estimates	COVB
CorrB	Correlation of parameter estimates	CORRB
<b>ODS Tables Created by the BAYES Statement</b>		
AutoCorr	Autocorrelation statistics for each parameter	Default
Corr	Correlation matrix of the posterior samples	STATS=COR
Cov	Covariance matrix of the posterior samples	STATS=COV
ESS	Effective sample size for each parameter	Default
MCSE	Monte Carlo standard error for each parameter	Default
Geweke	Geweke diagnostics for each parameter	Default
Heidelberger	Heidelberger-Welch diagnostics for each parameter	DIAGNOSTICS=HEIDEL
PostIntervals	Equal-tail and HPD intervals for each parameter	Default
PosteriorSample	Posterior samples	(ODS output data set only)
PostSummaries	Posterior summaries	Default
PriorSummaries	Prior summaries	STATS=PRIOR
Raftery	Raftery-Lewis diagnostics for each parameter	DIAGNOSTICS=RAFTERY
<b>ODS Tables Created by the TEST Statement</b>		
TestResults	Test results	Default

## ODS Graphics

You can use a name to reference every graph that is produced through ODS Graphics. The names of the graphs that PROC HPQLIM generates are listed in [Table 8.12](#).

**Table 8.12** Graphs Produced by PROC HPQLIM When a BAYES Statement Is Included

ODS Graph Name	Plot Description	Statement and Option
<b>Bayesian Diagnostic Plots</b>		
ADPanel	Autocorrelation function and density panel	PLOTS=(AUTOCORR DENSITY)
AutocorrPanel	Autocorrelation function panel	PLOTS=AUTOCORR
AutocorrPlot	Autocorrelation function plot	PLOTS(UNPACK)=AUTOCORR
DensityPanel	Density panel	PLOTS=DENSITY
DensityPlot	Density plot	PLOTS(UNPACK)=DENSITY
TAPanel	Trace and autocorrelation function panel	PLOTS=(TRACE AUTOCORR)
TADPanel	Trace, density, and autocorrelation function panel	PLOTS=(TRACE AUTOCORR DENSITY)
TDPanel	Trace and density panel	PLOTS=(TRACE DENSITY)
TracePanel	Trace panel	PLOTS=TRACE
TracePlot	Trace plot	PLOTS(UNPACK)=TRACE

## Examples: The HPQLIM Procedure

### Example 8.1: High-Performance Model with Censoring

This example shows the use of the HPQLIM procedure with an emphasis on processing a large data set and on the performance improvements that are achieved by executing in the high-performance distributed environment.

The following DATA step generates 5 million replicates from a censored model. The model contains seven variables.

```
data simulate;
  call streaminit(12345);
  array vars x1-x7;
  array parms{7} (3 4 2 4 -3 -5 -3);

  intercept=2;

  do i=1 to 5000000;
    sum_xb=0;
    do j=1 to 7;
```

```

vars[j]=rand('NORMAL',0,1);
sum_xb=sum_xb+parms[j]*vars[j];
end;
y=intercept+sum_xb+400*rand('NORMAL',0,1);
if y>400 then y=400;
if y<0 then y=0;
output;
end;
keep y x1-x7;
run;

```

The following statements estimate a censored model. The model is executed in the distributed computing environment with two threads and only one node. These settings are used to obtain a hypothetical environment that might resemble running the HPQLIM procedure on a desktop workstation with a dual-core CPU. To run these statements successfully, you need to set the macro variables GRIDHOST and GRIDINSTALLLOC to resolve to appropriate values, or you can replace the references to the macro variables in the example with the appropriate values.

```

option set=GRIDHOST="%GRIDHOST";
option set=GRIDINSTALLLOC="%GRIDINSTALLLOC";

proc hpqlim data=simulate ;
  performance nthreads=2 nodes=1 details
              host="%GRIDHOST" install="%GRIDINSTALLLOC";
  model y=x1-x7 /censored(lb=0 ub=400);
run;

```

Output 8.1.1 shows that the censored model was estimated on the grid, defined in a macro variable named GRIDHOST, in a distributed environment on only one node with two threads.

**Output 8.1.1** Censored Model with One Node and Two Threads: Performance Table

Estimating a Tobit model	
Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	1
Number of Threads per Node	2

Output 8.1.2 shows the estimation results for the censored model. The “Model Fit Summary” table shows detailed information about the model and indicates that all 5 million observations were used to fit the model. All parameter estimates in the “Parameter Estimates” table are highly significant and correspond to their theoretical values that were set during the data generating process. The optimization of the model with 5 million observations took 43 seconds.

**Output 8.1.2** Censored Model with One Node and Two Threads: Summary

Model Information							
Data Source		WORK.SIMULATE					
Response Variable		y					
Optimization Technique		Quasi-Newton					
Number of Observations							
Number of Observations Read		5000000					
Number of Observations Used		5000000					
Summary Statistics of Continuous Responses							
Variable	Mean	Standard Error	Type	Lower Bound	Upper Bound	N Obs Lower Bound	N Obs Upper Bound
y	127.0	159.491090	Censored	0	400.0	249E4	8E5
Convergence criterion (FCONV=2.220446E-16) satisfied.							
Model Fit Summary							
Number of Endogenous Variables		1					
Endogenous Variable		y					
Number of Observations		5000000					
Log Likelihood		-15268972					
Maximum Absolute Gradient		0.0003291					
Number of Iterations		11					
Optimization Method		Quasi-Newton					
AIC		30537962					
Schwarz Criterion		30538083					
Parameter Estimates							
Parameter	DF	Estimate	Standard Error	t Value	Approx Pr >  t		
Intercept	1	2.220379	0.222201	9.99	<.0001		
x1	1	3.055533	0.201620	15.15	<.0001		
x2	1	4.000176	0.201570	19.85	<.0001		
x3	1	1.852740	0.201555	9.19	<.0001		
x4	1	4.170266	0.201533	20.69	<.0001		
x5	1	-3.010679	0.201458	-14.94	<.0001		
x6	1	-5.176016	0.201541	-25.68	<.0001		
x7	1	-2.695948	0.201671	-13.37	<.0001		
_Sigma	1	399.997845	0.261930	1527.12	<.0001		

**Output 8.1.2** *continued*

Procedure Task Timing		
Task	Seconds	Percent
Reading and Levelizing Data	0.94	2.12%
Communication to Client	0.03	0.06%
Optimization	43.33	97.82%
Post-optimization	0.00	0.00%

In the following statements, the PERFORMANCE statement is modified to use a grid with 10 nodes, with each node capable of spawning eight threads:

```
proc hpqlim data=simulate ;
  performance nthreads=8 nodes=10 details
             host="&GRIDHOST" install="&GRIDINSTALLLOC";
  model y=x1-x7 /censored(lb=0 ub=400);
run;
```

The second model which was run on a grid with 10 nodes and eight threads each ([Output 8.1.3](#)) took only 1.9 seconds instead of 43 seconds to optimize.

**Output 8.1.3** Censored Model on Ten Nodes with Eight Threads Each: Performance Table

Estimating a Tobit model	
Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	10
Number of Threads per Node	8

Because the two models being estimated are identical, it is reasonable to expect that [Output 8.1.2](#) and [Output 8.1.4](#) would show the same results except for the performance. However, in certain circumstances, you might observe slight numerical differences in the results (depending on the number of nodes and threads) because the order in which partial results are accumulated, the limits of numerical precision, and the propagation of error in numerical computations can make a difference in the final result.

**Output 8.1.4** Censored Model on Ten Nodes with Eight Threads Each: Summary

Model Information	
Data Source	WORK.SIMULATE
Response Variable	y
Optimization Technique	Quasi-Newton

## Output 8.1.4 continued

Number of Observations							
Number of Observations Read				5000000			
Number of Observations Used				5000000			
Summary Statistics of Continuous Responses							
Variable	Mean	Standard Error	Type	Lower Bound	Upper Bound	N Obs Lower Bound	N Obs Upper Bound
y	127.0	159.491090	Censored	0	400.0	249E4	8E5
Convergence criterion (FCONV=2.220446E-16) satisfied.							
Model Fit Summary							
Number of Endogenous Variables				1			
Endogenous Variable				y			
Number of Observations				5000000			
Log Likelihood				-15268972			
Maximum Absolute Gradient				0.0008332			
Number of Iterations				10			
Optimization Method				Quasi-Newton			
AIC				30537962			
Schwarz Criterion				30538083			
Parameter Estimates							
Parameter	DF	Estimate	Standard Error	t Value	Approx Pr >  t		
Intercept	1	2.220358	0.222201	9.99	<.0001		
x1	1	3.055491	0.201620	15.15	<.0001		
x2	1	4.000196	0.201570	19.85	<.0001		
x3	1	1.852735	0.201555	9.19	<.0001		
x4	1	4.170323	0.201533	20.69	<.0001		
x5	1	-3.010670	0.201458	-14.94	<.0001		
x6	1	-5.176019	0.201541	-25.68	<.0001		
x7	1	-2.695886	0.201671	-13.37	<.0001		
_Sigma	1	399.997846	0.261930	1527.12	<.0001		
Procedure Task Timing							
Task				Seconds	Percent		
Reading and Levelizing Data				0.07	3.20%		
Communication to Client				0.17	7.84%		
Optimization				1.91	88.96%		
Post-optimization				0.00	0.00%		

As this example suggests, increasing the number of nodes and the number of threads per node improves performance significantly. When you use the parallelism that a high-performance distributed environment affords, you can see an even more dramatic reduction in the time required for the optimization as the number of observations in the data set increases. When the data set is extremely large, the computations might not even be possible with the typical memory resources and computational constraints of a desktop computer. Under such circumstances the high-performance distributed environment becomes a necessity.

---

## Example 8.2: Bayesian High-Performance Model with Censoring

This example shows the use of the Bayesian analysis available in the HPQLIM procedure with an emphasis on processing a large data set and on the performance improvements that are achieved by executing in a high-performance distributed environment.

The model and the data set are the same as in [Example 8.1](#), and the priors are set to the defaults.

The model is executed in the distributed computing environment with two threads and only one node. These settings are used to obtain a hypothetical environment that might resemble running the HPQLIM procedure on a desktop workstation with a dual-core CPU. To run the following statements successfully, you need to set the macro variables GRIDHOST and GRIDINSTALLLOC to resolve to appropriate values, or you can replace the references to the macro variables in the example with the appropriate values.

```
option set=GRIDHOST="%GRIDHOST";
option set=GRIDINSTALLLOC="%GRIDINSTALLLOC";

proc hpqlim data=simulate ;
  bayes nbi=10000 nmc=30000;
  performance nthreads=2 nodes=1 details
             host="%GRIDHOST" install="%GRIDINSTALLLOC";
  model y=x1-x7 /censored(lb=0 ub=400);
  %*;      ods output PerformanceInfo=perfInfo;
  %*;      ods output Timing=time;
run;
```

[Output 8.2.1](#) shows a summary of the posterior distribution that is associated with the censored model when you use diffuse prior distributions.

**Output 8.2.1** Posterior Summary for Bayesian Censored Model

Estimating a Tobit model						
The HPQLIM Procedure						
Posterior Summaries						
Parameter	N	Mean	Standard Deviation	Percentiles 25%	50%	75%
Intercept	30000	2.2240	0.2206	2.0730	2.2209	2.3725
x1	30000	3.0489	0.2006	2.9127	3.0419	3.1803
x2	30000	3.9984	0.1978	3.8667	4.0014	4.1284
x3	30000	1.8443	0.2018	1.7069	1.8456	1.9822
x4	30000	4.1748	0.2000	4.0419	4.1753	4.3104
x5	30000	-3.0096	0.1998	-3.1447	-3.0114	-2.8736
x6	30000	-5.1686	0.2003	-5.3041	-5.1680	-5.0348
x7	30000	-2.6953	0.2099	-2.8375	-2.6955	-2.5545
_Sigma	30000	400.0	0.2615	399.8	400.0	400.2

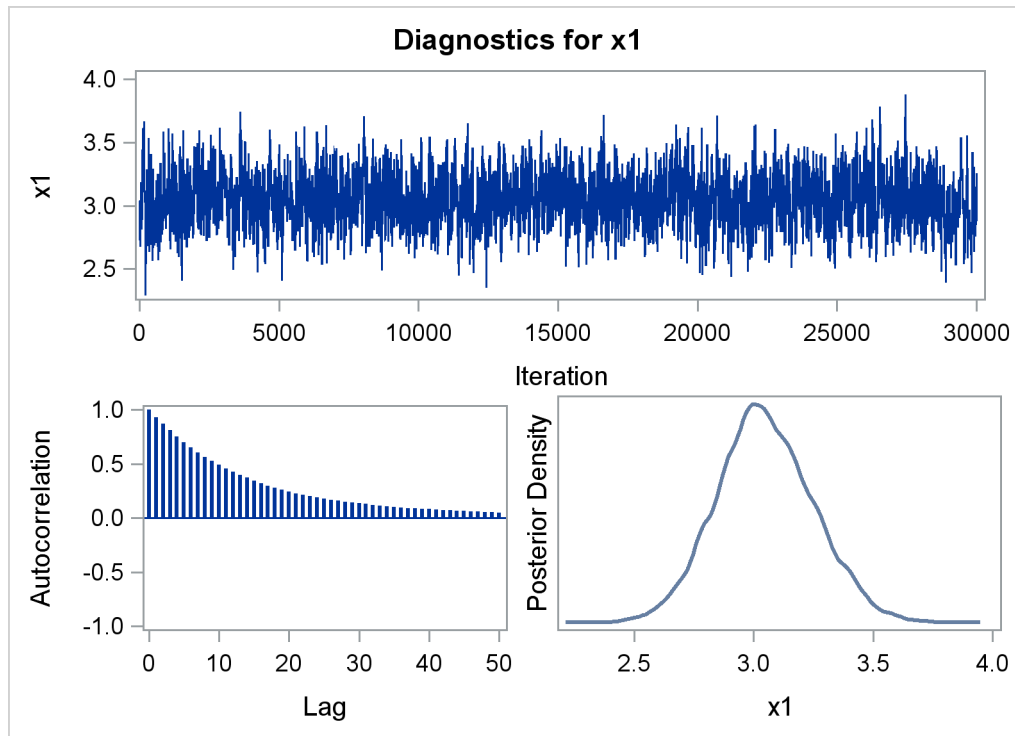
Output 8.2.2 show a summary of the performance when you use a distributed computing environment with one node and two threads.

**Output 8.2.2** Performance Analysis for Bayesian Censored Model on One Node with Two Threads

Estimating a Tobit model		
Performance Information		
Host Node	<< your grid host >>	
Install Location	<< your grid install location >>	
Execution Mode	Distributed	
Grid Mode	Symmetric	
Number of Compute Nodes	1	
Number of Threads per Node	2	
Estimating a Tobit model		
Procedure Task Timing		
Task	Seconds	Percent
Reading and Levelizing Data	0.95	0.00%
Communication to Client	0.15	0.00%
Bayesian Analysis: Likelihood for MCMC	30819.21	99.86%
Bayesian Analysis: MCMC	0.22	0.00%
Optimization	43.41	0.14%
Post-optimization	0.00	0.00%

Finally, Output 8.2.3 shows the diagnostic and summary plots that are associated with *X1*.



**Output 8.2.3** Bayesian Diagnostic and Summary Plots for x1

In the following statements, the PERFORMANCE statement is modified to use a grid with 10 nodes, where each node spawns eight threads:

```
option set=GRIDHOST="%GRIDHOST";
option set=GRIDINSTALLLOC="%GRIDINSTALLLOC";

proc hpqlim data=simulate ;
  bayes nbi=10000 nmc=30000;
    performance nthreads=8 nodes=10 details
      host="%GRIDHOST" install="%GRIDINSTALLLOC";
  model y=x1-x7 /censored(lb=0 ub=400);
  %*;    ods output PerformanceInfo=perfInfo;
  %*;    ods output Timing=time;
run;
```

The two models are identical, but the second implementation, which was run on a grid that used 10 nodes with eight threads each, took only 21 minutes instead of 8.5 hours to sample from the same posterior distribution.

**Output 8.2.4** Performance Analysis for Bayesian Censored Model on Ten Nodes with Eight Threads Each

Estimating a Tobit model		
Performance Information		
Host Node	<< your grid host >>	
Install Location	<< your grid install location >>	
Execution Mode	Distributed	
Grid Mode	Symmetric	
Number of Compute Nodes	10	
Number of Threads per Node	8	
Estimating a Tobit model		
Procedure Task Timing		
Task	Seconds	Percent
Reading and Levelizing Data	0.07	0.01%
Communication to Client	0.23	0.02%
Bayesian Analysis: Likelihood for MCMC	1242.54	99.84%
Bayesian Analysis: MCMC	0.17	0.01%
Optimization	1.49	0.12%
Post-optimization	0.00	0.00%

---

## References

- Aigner, C., Lovell, C. A. K., and Schmidt, P. (1977), “Formulation and Estimation of Stochastic Frontier Production Function Models,” *Journal of Econometrics*, 6, 21–37.
- Battese, G. E. and Coelli, T. J. (1988), “Prediction of Firm-Level Technical Efficiencies with a Generalized Frontier Production Function and Panel Data,” *Journal of Econometrics*, 38, 387–399.
- Christensen, L. R. and Greene, W. H. (1976), “Economies of Scale in U.S. Electric Power Generation,” *Journal of Political Economy*, 84, 655–676.
- Coelli, T. J., Prasada Rao, D. S., and Battese, G. E. (1998), *An Introduction to Efficiency and Productivity Analysis*, London: Kluwer Academic.
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004), *Bayesian Data Analysis*, 2nd Edition, London: Chapman & Hall.
- Jondrow, J., Lovell, C. A. K., Materov, I. S., and Schmidt, P. (1982), “On the Estimation of Technical Efficiency in the Stochastic Frontier Production Function Model,” *Journal of Econometrics*, 19, 233–238.
- Kumbhakar, S. C. and Lovell, C. A. K. (2000), *Stochastic Frontier Analysis*, New York: Cambridge University Press.
- McKelvey, R. D. and Zavoina, W. (1975), “A Statistical Model for the Analysis of Ordinal Level Dependent Variables,” *Journal of Mathematical Sociology*, 4, 103–120.
- Meeusen, W. and van den Broeck, J. (1977), “Efficiency Estimation from Cobb-Douglas Production Functions with Composed Error,” *International Economic Review*, 18, 435–444.
- Mroz, T. A. (1987), “The Sensitivity of an Empirical Model of Married Women’s Work to Economic and Statistical Assumptions,” *Econometrica*, 55, 765–799.
- Schervish, M. J. (1995), *Theory of Statistics*, New York: Springer-Verlag.
- Wooldridge, J. M. (2002), *Econometric Analysis of Cross Section of Panel Data*, Cambridge, MA: MIT Press.



# Chapter 9

## The HPSEVERITY Procedure

### Contents

---

Overview: HPSEVERITY Procedure . . . . .	<b>238</b>
Getting Started: HPSEVERITY Procedure . . . . .	<b>240</b>
A Simple Example of Fitting Predefined Distributions . . . . .	240
An Example with Left-Truncation and Right-Censoring . . . . .	243
An Example of Modeling Regression Effects . . . . .	247
Syntax: HPSEVERITY Procedure . . . . .	<b>251</b>
Functional Summary . . . . .	252
PROC HPSEVERITY Statement . . . . .	254
BY Statement . . . . .	261
DIST Statement . . . . .	261
LOSS Statement . . . . .	263
NLOPTIONS Statement . . . . .	265
OUTSCORELIB Statement (Experimental) . . . . .	269
PERFORMANCE Statement . . . . .	271
SCALEMODEL Statement . . . . .	272
WEIGHT Statement . . . . .	273
Programming Statements . . . . .	274
Details: HPSEVERITY Procedure . . . . .	<b>274</b>
Predefined Distributions . . . . .	274
Censoring and Truncation . . . . .	284
Parameter Estimation Method . . . . .	285
Details of Optimization Algorithms . . . . .	287
Parameter Initialization . . . . .	289
Estimating Regression Effects . . . . .	290
Empirical Distribution Function Estimation Methods . . . . .	295
Statistics of Fit . . . . .	301
Distributed and Multithreaded Computation . . . . .	305
Defining a Severity Distribution Model with the FCMP Procedure . . . . .	308
Predefined Utility Functions . . . . .	320
Scoring Functions (Experimental) . . . . .	325
Custom Objective Functions . . . . .	332
Input Data Sets . . . . .	335
Output Data Sets . . . . .	336
Displayed Output . . . . .	339
Examples: HPSEVERITY Procedure . . . . .	<b>342</b>
Example 9.1: Defining a Model for Gaussian Distribution . . . . .	342

Example 9.2: Defining a Model for the Gaussian Distribution with a Scale Parameter	346
Example 9.3: Defining a Model for Mixed-Tail Distributions . . . . .	352
Example 9.4: Fitting a Scaled Tweedie Model with Regressors . . . . .	359
Example 9.5: Fitting Distributions to Interval-Censored Data . . . . .	363
Example 9.6: Benefits of Distributed and Multithreaded Computing . . . . .	365
Example 9.7: Estimating Parameters Using Cramér-von Mises Estimator . . . . .	371
Example 9.8: Defining a Finite Mixture Model That Has a Scale Parameter . . . . .	374
Example 9.9: Predicting Mean and Value-at-Risk by Using Scoring Functions . . . . .	380
References . . . . .	385

---

## Overview: HPSEVERITY Procedure

PROC HPSEVERITY estimates parameters of any arbitrary continuous probability distribution that is used to model the magnitude (severity) of a continuous-valued event of interest. Some examples of such events are loss amounts paid by an insurance company and demand of a product as depicted by its sales. PROC HPSEVERITY is especially useful when the severity of an event does not follow typical distributions (such as the normal distribution) that are often assumed by standard statistical methods.

PROC HPSEVERITY runs in either single-machine mode or distributed mode. **NOTE:** Distributed mode requires SAS High-Performance Econometrics.

PROC HPSEVERITY provides a default set of probability distribution models that includes the Burr, exponential, gamma, generalized Pareto, inverse Gaussian (Wald), lognormal, Pareto, Tweedie, and Weibull distributions. In the simplest form, you can estimate the parameters of any of these distributions by using a list of severity values that are recorded in a SAS data set. You can optionally group the values by a set of BY variables. PROC HPSEVERITY computes the estimates of the model parameters, their standard errors, and their covariance structure by using the maximum likelihood method for each of the BY groups.

PROC HPSEVERITY can fit multiple distributions at the same time and choose the best distribution according to a selection criterion that you specify. You can use seven different statistics of fit as selection criteria. They are log likelihood, Akaike's information criterion (AIC), corrected Akaike's information criterion (AICC), Schwarz Bayesian information criterion (BIC), Kolmogorov-Smirnov statistic (KS), Anderson-Darling statistic (AD), and Cramér-von Mises statistic (CvM).

You can request the procedure to output the status of the estimation process, the parameter estimates and their standard errors, the estimated covariance structure of the parameters, and the statistics of fit.

The following key features make PROC HPSEVERITY unique among SAS procedures that can estimate continuous probability distributions:

- It enables you to fit a distribution model when the severity values are truncated or censored or both. You can specify any combination of the following types of censoring and truncation effects: left-censoring, right-censoring, left-truncation, or right-truncation. This is especially useful in applications with an insurance-type model where a severity (loss) is reported and recorded only if it is greater than the deductible amount (left-truncation) and where a severity value greater than or equal to the policy limit is recorded at the limit (right-censoring). Another useful application is that of interval-censored data,

where you know both the lower limit (right-censoring) and upper limit (left-censoring) on the severity, but you do not know the exact value.

It uses an appropriate estimator of the empirical distribution function (EDF). EDF is required to compute the KS, AD, and CvM statistics-of-fit. The procedure also provides the EDF estimates to your custom parameter initialization method. When you specify truncation or censoring, the EDF is estimated by using either Kaplan-Meier's product-limit estimator or Turnbull's estimator. The former is used by default when you specify only one form of censoring effect (right-censoring or left-censoring), whereas the latter is used by default when you specify both left-censoring and right-censoring effects.

- It enables you to define any arbitrary continuous parametric distribution model and to estimate its parameters. You just need to define the key components of the distribution, such as its probability density function (PDF) and cumulative distribution function (CDF), as a set of functions and subroutines written with the FCMP procedure, which is part of Base SAS software. As long as the functions and subroutines follow certain rules, the HPSEVERITY procedure can fit the distribution model defined by them.
- It can model the effect of exogenous or regressor variables on a probability distribution, as long as the distribution has a scale parameter. A linear combination of the regressor variables is assumed to affect the scale parameter via an exponential link function.

If a distribution does not have a scale parameter, then either it needs to have another parameter that can be derived from a scale parameter by using a supported transformation or it needs to be reparameterized to have a scale parameter. If neither of these is possible, then regression effects cannot be modeled.

- It enables you to specify your own objective function to be optimized for estimating the parameters of a model. You can write SAS programming statements to specify the contribution of each observation to the objective function. You can use keyword functions such as `_PDF_` and `_CDF_` to generalize the objective function to any distribution. If you do not specify your own objective function, then the parameters of a model are estimated by maximizing the likelihood function of the data.
- It enables you to create scoring functions that offer a convenient way to evaluate any distribution function, such as PDF, CDF, QUANTILE, or your custom distribution function, for a fitted model on new observations.

Because the HPSEVERITY procedure is a high-performance analytical procedure, it also does the following:

- enables you to run in distributed mode on a cluster of machines that distribute the data and the computations
- enables you to run in single-machine mode on the server where SAS is installed
- exploits all the available cores and concurrent threads, regardless of execution mode

For more information, see the section “[Processing Modes](#)” on page 12 in Chapter 3, “[Shared Concepts and Topics](#).”

## Getting Started: HPSEVERITY Procedure

This section outlines the use of the HPSEVERITY procedure to fit continuous probability distribution models. Three examples illustrate different features of the procedure.

### A Simple Example of Fitting Predefined Distributions

The simplest way to use PROC HPSEVERITY is to fit all the predefined distributions to a set of values and let the procedure identify the best fitting distribution.

Consider a lognormal distribution, whose probability density function (PDF)  $f$  and cumulative distribution function (CDF)  $F$  are as follows, respectively, where  $\Phi$  denotes the CDF of the standard normal distribution:

$$f(x; \mu, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2} \quad \text{and} \quad F(x; \mu, \sigma) = \Phi\left(\frac{\log(x)-\mu}{\sigma}\right)$$

The following DATA step statements simulate a sample from a lognormal distribution with population parameters  $\mu = 1.5$  and  $\sigma = 0.25$ , and store the sample in the variable Y of a data set Work.Test\_sev1:

```
/*----- Simple Lognormal Example -----*/
data test_sev1(keep=y label='Simple Lognormal Sample');
  call streaminit(45678);
  label y='Response Variable';
  Mu = 1.5;
  Sigma = 0.25;
  do n = 1 to 100;
    y = exp(Mu) * rand('LOGNORMAL')**Sigma;
    output;
  end;
run;
```

The following statements fit all the predefined distribution models to the values of Y and identify the best distribution according to the corrected Akaike's information criterion (AICC):

```
proc hpseverity data=test_sev1 crit=aicc;
  loss y;
  dist _predefined_;
run;
```

The PROC HPSEVERITY statement specifies the input data set along with the model selection criterion, the LOSS statement specifies the variable to be modeled, and the DIST statement with the `_PREDEFINED_` keyword specifies that all the predefined distribution models be fitted.

Some of the default output displayed by this step is shown in [Figure 9.1](#) through [Figure 9.3](#). First, information about the input data set is displayed followed by the “Model Selection” table, as shown in [Figure 9.1](#). The model selection table displays the convergence status, the value of the selection criterion, and the selection status for each of the candidate models. The Converged column indicates whether the estimation process for a given distribution model has converged, might have converged, or failed. The Selected column indicates whether a given distribution has the best fit for the data according to the selection criterion. For this example, the lognormal distribution model is selected, because it has the lowest value for the selection criterion.



**Figure 9.1** Data Set Information and Model Selection Table

The HPSEVERITY Procedure			
Input Data Set			
Name	WORK.TEST_SEV1		
Label	Simple Lognormal Sample		
Model Selection			
Distribution	Converged	AICC	Selected
Burr	Yes	322.50845	No
Exp	Yes	508.12287	No
Gamma	Yes	320.50264	No
Igauss	Yes	319.61652	No
Logn	Yes	319.56579	Yes
Pareto	Yes	510.28172	No
Gpd	Yes	510.20576	No
Weibull	Yes	334.82373	No

Next, the estimation information for each of the candidate models is displayed. The information for the lognormal model, which is the best fitting model, is shown in Figure 9.2. The first table displays a summary of the distribution. The second table displays the convergence status. This is followed by a summary of the optimization process which indicates the technique used, the number of iterations, the number of times the objective function was evaluated, and the log likelihood attained at the end of the optimization. Since the model with lognormal distribution has converged, PROC HPSEVERITY displays its statistics of fit and parameter estimates. The estimates of  $\mu=1.49605$  and  $\sigma=0.26243$  are quite close to the population parameters of  $\mu=1.5$  and  $\sigma=0.25$  from which the sample was generated. The  $p$ -value for each estimate indicates the rejection of the null hypothesis that the estimate is 0, implying that both the estimates are significantly different from 0.

**Figure 9.2** Estimation Details for the Lognormal Model

The HPSEVERITY Procedure		
Logn Distribution		
Distribution Information		
Name	Logn	
Description	Lognormal Distribution	
Distribution Parameters	2	
Convergence Status		
Convergence criterion (GCONV=1E-8) satisfied.		

Figure 9.2 continued

Optimization Summary	
Optimization Technique	Trust Region
Iterations	2
Function Calls	8
Log Likelihood	-157.72104

Fit Statistics	
-2 Log Likelihood	315.44208
AIC	319.44208
AICC	319.56579
BIC	324.65242
Kolmogorov-Smirnov	0.50641
Anderson-Darling	0.31240
Cramer-von Mises	0.04353

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu	1.49605	0.02651	56.43	<.0001
Sigma	0.26243	0.01874	14.00	<.0001

The parameter estimates of the Burr distribution are shown in Figure 9.3. These estimates are used in the next example.

Figure 9.3 Parameter Estimates for the Burr Model

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Theta	4.62348	0.46181	10.01	<.0001
Alpha	1.15706	0.47493	2.44	0.0167
Gamma	6.41227	0.99039	6.47	<.0001

## An Example with Left-Truncation and Right-Censoring

PROC HPSEVERITY enables you to specify that the response variable values are left-truncated or right-censored. The following DATA step expands the data set of the previous example to simulate a scenario that is typically encountered by an automobile insurance company. The values of the variable Y represent the loss values on claims that are reported to an auto insurance company. The variable THRESHOLD records the deductible on the insurance policy. If the actual value of Y is less than or equal to the deductible, then it is unobservable and does not get recorded. In other words, THRESHOLD specifies the left-truncation of Y. LIMIT records the policy limit. If the value of Y is equal to or greater than the recorded value, then the observation is right-censored.

```

/*----- Lognormal Model with left-truncation and censoring -----*/
data test_sev2(keep=y threshold limit
               label='A Lognormal Sample With Censoring and Truncation');
  set test_sev1;
  label y='Censored & Truncated Response';
  if _n_ = 1 then call streaminit(45679);

  /* make about 20% of the observations left-truncated */
  if (rand('UNIFORM') < 0.2) then
    threshold = y * (1 - rand('UNIFORM'));
  else
    threshold = .;
  /* make about 15% of the observations right-censored */
  iscens = (rand('UNIFORM') < 0.15);
  if (iscens) then
    limit = y;
  else
    limit = .;
run;

```

The following statements use the AICC criterion to analyze which of the four predefined distributions (lognormal, Burr, gamma, and Weibull) has the best fit for the data:

```

proc hpseverity data=test_sev2 crit=aicc print=all ;
  loss y / lt=threshold rc=limit;

  dist logn burr gamma weibull;
  performance nthreads=2;
run;

```

The LOSS statement specifies the left-truncation and right-censoring variables. The DIST statement specifies the candidate distributions. The PRINT= option in the PROC HPSEVERITY statement requests that all the displayed output be prepared. The NTHREADS option in the PERFORMANCE statement specifies that two threads of computation be used. The option is shown here just for illustration. You should use it only when you want to restrict the procedure to use a different number of threads than the value of the CPUCOUNT= system option, which usually defaults to the number of physical CPU cores available on your machine, thereby allowing the procedure to fully utilize the computational power of your machine.

Some of the key results prepared by PROC HPSEVERITY are shown in Figure 9.4 through Figure 9.7. In addition to the estimates of the range, mean, and standard deviation of Y, the “Descriptive Statistics for Variable y” table shown in Figure 9.4 also indicates the number of observations that are left-truncated or right-censored. The “Model Selection” table in Figure 9.4 shows that models with all the candidate distributions have converged and that the Logn (lognormal) model has the best fit for the data according to the AICC criterion.

**Figure 9.4** Summary Results for the Truncated and Censored Data

The HPSEVERITY Procedure			
Input Data Set			
Name	WORK.TEST_SEV2		
Label	A Lognormal Sample With Censoring and Truncation		
Descriptive Statistics for y			
Observations			100
Observations Used for Estimation			100
Minimum			2.30264
Maximum			8.34116
Mean			4.62007
Standard Deviation			1.23627
Left Truncated Observations			23
Right Censored Observations			14
Model Selection			
Distribution	Converged	AICC	Selected
Logn	Yes	298.92672	Yes
Burr	Yes	302.66229	No
Gamma	Yes	299.45293	No
Weibull	Yes	309.26779	No

PROC HPSEVERITY also prepares a table that shows all the fit statistics for all the candidate models. It is useful to see which model would be the best fit according to each of the criteria. The “All Fit Statistics” table prepared for this example is shown in Figure 9.5. It indicates that the lognormal model is chosen by all the criteria.

**Figure 9.5** Comparing All Statistics of Fit for the Truncated and Censored Data

All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Logn	294.80301*	298.80301*	298.92672*	304.01335*	0.51824*
Burr	296.41229	302.41229	302.66229	310.22780	0.66984
Gamma	295.32921	299.32921	299.45293	304.53955	0.62511
Weibull	305.14408	309.14408	309.26779	314.35442	0.93307
Note: The asterisk (*) marks the best model according to each column's criterion.					

All Fit Statistics		
Distribution	AD	CvM
Logn	0.34736*	0.05159*
Burr	0.36712	0.05726
Gamma	0.42921	0.05526
Weibull	1.40699	0.17465
Note: The asterisk (*) marks the best model according to each column's criterion.		

### Specifying Initial Values for Parameters

All the predefined distributions have parameter initialization functions built into them. For the current example, Figure 9.6 shows the initial values that are obtained by the predefined method for the Burr distribution. It also shows the summary of the optimization process and the final parameter estimates.

**Figure 9.6** Burr Model Summary for the Truncated and Censored Data

Initial Parameter Values and Bounds			
Parameter	Initial Value	Lower Bound	Upper Bound
Theta	4.78102	1.05367E-8	Infty
Alpha	2.00000	1.05367E-8	Infty
Gamma	2.00000	1.05367E-8	Infty
Optimization Summary			
Optimization Technique	Trust	Region	
Iterations		8	
Function Calls		23	
Log Likelihood		-148.20614	

Figure 9.6 continued

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Theta	4.76980	0.62492	7.63	<.0001
Alpha	1.16363	0.58859	1.98	0.0509
Gamma	5.94081	1.05004	5.66	<.0001

You can specify a different set of initial values if estimates are available from fitting the distribution to similar data. For this example, the parameters of the Burr distribution can be initialized with the final parameter estimates of the Burr distribution that were obtained in the first example (shown in Figure 9.3). One of the ways in which you can specify the initial values is as follows:

```
/*----- Specifying initial values using INIT= option -----*/
proc hpseverity data=test_sev2 crit=aicc print=all;
  loss y / lt=threshold rc=limit;

  dist burr(init=(theta=4.62348 alpha=1.15706 gamma=6.41227));
  performance nthreads=2;
run;
```

The names of the parameters that are specified in the INIT option must match the parameter names in the definition of the distribution. The results obtained with these initial values are shown in Figure 9.7. These results indicate that new set of initial values causes the optimizer to reach the same solution with fewer iterations and function evaluations as compared to the default initialization.

Figure 9.7 Burr Model Optimization Summary for the Truncated and Censored Data

The HPSEVERITY Procedure				
Burr Distribution				
Optimization Summary				
Optimization Technique	Trust Region			
Iterations	5			
Function Calls	16			
Log Likelihood	-148.20614			
Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Theta	4.76980	0.62492	7.63	<.0001
Alpha	1.16363	0.58859	1.98	0.0509
Gamma	5.94081	1.05004	5.66	<.0001

## An Example of Modeling Regression Effects

Consider a scenario in which the magnitude of the response variable might be affected by some regressor (exogenous or independent) variables. The HPSEVERITY procedure enables you to model the effect of such variables on the distribution of the response variable via an exponential link function. In particular, if you have  $k$  random regressor variables denoted by  $x_j$  ( $j = 1, \dots, k$ ), then the distribution of the response variable  $Y$  is assumed to have the form

$$Y \sim \exp\left(\sum_{j=1}^k \beta_j x_j\right) \cdot \mathcal{F}(\Theta)$$

where  $\mathcal{F}$  denotes the distribution of  $Y$  with parameters  $\Theta$  and  $\beta_j$  ( $j = 1, \dots, k$ ) denote the regression parameters (coefficients).

For the effective distribution of  $Y$  to be a valid distribution from the same parametric family as  $\mathcal{F}$ , it is necessary for  $\mathcal{F}$  to have a scale parameter. The effective distribution of  $Y$  can be written as

$$Y \sim \mathcal{F}(\theta, \Omega)$$

where  $\theta$  denotes the scale parameter and  $\Omega$  denotes the set of nonscale parameters. The scale  $\theta$  is affected by the regressors as

$$\theta = \theta_0 \cdot \exp\left(\sum_{j=1}^k \beta_j x_j\right)$$

where  $\theta_0$  denotes a *base* value of the scale parameter.

Given this form of the model, PROC HPSEVERITY allows a distribution to be a candidate for modeling regression effects only if it has an untransformed or a log-transformed scale parameter.

All the predefined distributions, except the lognormal distribution, have a direct scale parameter (that is, a parameter that is a scale parameter without any transformation). For the lognormal distribution, the parameter  $\mu$  is a log-transformed scale parameter. This can be verified by replacing  $\mu$  with a parameter  $\theta = e^\mu$ , which results in the following expressions for the PDF  $f$  and the CDF  $F$  in terms of  $\theta$  and  $\sigma$ , respectively, where  $\Phi$  denotes the CDF of the standard normal distribution:

$$f(x; \theta, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\log(x) - \log(\theta)}{\sigma}\right)^2} \quad \text{and} \quad F(x; \theta, \sigma) = \Phi\left(\frac{\log(x) - \log(\theta)}{\sigma}\right)$$

With this parameterization, the PDF satisfies the  $f(x; \theta, \sigma) = \frac{1}{\theta} f\left(\frac{x}{\theta}; 1, \sigma\right)$  condition and the CDF satisfies the  $F(x; \theta, \sigma) = F\left(\frac{x}{\theta}; 1, \sigma\right)$  condition. This makes  $\theta$  a scale parameter. Hence,  $\mu = \log(\theta)$  is a log-transformed scale parameter and the lognormal distribution is eligible for modeling regression effects.

The following DATA step simulates a lognormal sample whose scale is decided by the values of the three regressors X1, X2, and X3 as follows:

$$\mu = \log(\theta) = 1 + 0.75 X1 - X2 + 0.25 X3$$

```

/*----- Lognormal Model with Regressors -----*/
data test_sev3(keep=y x1-x3
                label='A Lognormal Sample Affected by Regressors');
    array x{*} x1-x3;
    array b{4} _TEMPORARY_ (1 0.75 -1 0.25);
    call streaminit(45678);
    label y='Response Influenced by Regressors';
    Sigma = 0.25;
    do n = 1 to 100;
        Mu = b(1); /* log of base value of scale */
        do i = 1 to dim(x);
            x(i) = rand('UNIFORM');
            Mu = Mu + b(i+1) * x(i);
        end;
        y = exp(Mu) * rand('LOGNORMAL')**Sigma;
        output;
    end;
run;

```

The following PROC HPSEVERITY step fits the lognormal, Burr, and gamma distribution models to this data. The regressors are specified in the SCALEMODEL statement.

```

proc hpseverity data=test_sev3 crit=aicc print=all;
    loss y;
    scalemodel x1-x3;

    dist logn burr gamma;
run;

```

Some of the key results prepared by PROC HPSEVERITY are shown in [Figure 9.8](#) through [Figure 9.12](#). The descriptive statistics of all the variables are shown in [Figure 9.8](#).

**Figure 9.8** Summary Results for the Regression Example

The HPSEVERITY Procedure	
Input Data Set	
Name	WORK.TEST_SEV3
Label	A Lognormal Sample Affected by Regressors
Descriptive Statistics for y	
Observations	100
Observations Used for Estimation	100
Minimum	1.17863
Maximum	6.65269
Mean	2.99859
Standard Deviation	1.12845



**Figure 9.8** *continued*

Descriptive Statistics for Regressors					
Variable	N	Minimum	Maximum	Mean	Standard Deviation
x1	100	0.0005115	0.97971	0.51689	0.28206
x2	100	0.01883	0.99937	0.47345	0.28885
x3	100	0.00255	0.97558	0.48301	0.29709

The comparison of the fit statistics of all the models is shown in [Figure 9.9](#). It indicates that the lognormal model is the best model according to each of the likelihood-based statistics.

**Figure 9.9** Comparison of Statistics of Fit for the Regression Example

All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Logn	187.49609*	197.49609*	198.13439*	210.52194*	1.97544
Burr	190.69154	202.69154	203.59476	218.32256	2.09334
Gamma	188.91483	198.91483	199.55313	211.94069	1.94472*
Note: The asterisk (*) marks the best model according to each column's criterion.					

All Fit Statistics		
Distribution	AD	CvM
Logn	17.24618	1.21665
Burr	13.93436*	1.28529
Gamma	15.84787	1.17617*
Note: The asterisk (*) marks the best model according to each column's criterion.		

The distribution information and the convergence results of the lognormal model are shown in [Figure 9.10](#). The iteration history gives you a summary of how the optimizer is traversing the surface of the log-likelihood function in its attempt to reach the optimum. Both the change in the log likelihood and the maximum gradient of the objective function with respect to any of the parameters typically approach 0 if the optimizer converges.

**Figure 9.10** Convergence Results for the Lognormal Model with Regressors

The HPSEVERITY Procedure				
Logn Distribution				
Distribution Information				
Name				Logn
Description	Lognormal Distribution			
Distribution Parameters				2
Regression Parameters				3
Convergence Status				
Convergence criterion (GCONV=1E-8) satisfied.				
Optimization Iteration History				
Iter	Function Calls	-Log Likelihood	Change	Maximum Gradient
0	2	93.75285		6.16002
1	4	93.74805	-0.0048055	0.11031
2	6	93.74805	-1.5017E-6	0.00003376
3	10	93.74805	-1.279E-13	3.1051E-12
Optimization Summary				
Optimization Technique		Trust Region		
Iterations		3		
Function Calls		10		
Log Likelihood		-93.74805		

The final parameter estimates of the lognormal model are shown in [Figure 9.11](#). All the estimates are significantly different from 0. The estimate that is reported for the parameter  $Mu$  is the base value for the log-transformed scale parameter  $\mu$ . Let  $x_i$  ( $1 \leq i \leq 3$ ) denote the observed value for regressor  $X_i$ . If the lognormal distribution is chosen to model  $Y$ , then the effective value of the parameter  $\mu$  varies with the observed values of regressors as

$$\mu = 1.04047 + 0.65221 x_1 - 0.91116 x_2 + 0.16243 x_3$$

These estimated coefficients are reasonably close to the population parameters (that is, within one or two standard errors).

**Figure 9.11** Parameter Estimates for the Lognormal Model with Regressors

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu	1.04047	0.07614	13.66	<.0001
Sigma	0.22177	0.01609	13.78	<.0001
x1	0.65221	0.08167	7.99	<.0001
x2	-0.91116	0.07946	-11.47	<.0001
x3	0.16243	0.07782	2.09	0.0395

The estimates of the gamma distribution model, which is the next best model according to the fit statistics, are shown in Figure 9.12. The estimate that is reported for the parameter *Theta* is the base value for the scale parameter  $\theta$ . If the gamma distribution is chosen to model  $Y$ , then the effective value of the scale parameter is  $\theta = 0.14293 \exp(0.64562 x_1 - 0.89831 x_2 + 0.14901 x_3)$ .

**Figure 9.12** Parameter Estimates for the Gamma Model with Regressors

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Theta	0.14293	0.02329	6.14	<.0001
Alpha	20.37726	2.93277	6.95	<.0001
x1	0.64562	0.08224	7.85	<.0001
x2	-0.89831	0.07962	-11.28	<.0001
x3	0.14901	0.07870	1.89	0.0613

## Syntax: HPSEVERITY Procedure

The following statements are available in the HPSEVERITY procedure:

```

PROC HPSEVERITY options ;
  BY variable-list ;
  LOSS < response-variable > < / censoring-truncation-options > ;
  WEIGHT weight-variable ;
  SCALEMODEL regressor-variable-list < / scalemodel-options > ;
  DIST distribution-name-or-keyword < ( distribution-option ) < distribution-name-or-keyword
    < ( distribution-option ) > > ... > < / preprocess-options > ;
  OUTSCORELIB < OUTLIB = > fcmp-library-name options ;
  NLOPTIONS options ;
  PERFORMANCE options ;
  Programming statements ;

```

## Functional Summary

Table 9.1 summarizes the statements and options that control the HPSEVERITY procedure.

**Table 9.1** HPSEVERITY Functional Summary

Description	Statement	Option
<b>Statements</b>		
Specifies BY-group processing	BY	
Specifies the response variable to model along with censoring and truncation effects	LOSS	
Specifies the weight variable	WEIGHT	
Specifies the regression variables to model	SCALEMODEL	
Specifies distributions to fit	DIST	
Specifies the library to write scoring functions to	OUTSCORELIB	
Specifies optimization options	NLOPTIONS	
Specifies performance options	PERFORMANCE	
Specifies programming statements that define an objective function	Programming statements	
<b>Data Set Options</b>		
Specifies that the OUTEST= data set contain covariance estimates	PROC HPSEVERITY	COVOUT
Specifies the input data set	PROC HPSEVERITY	DATA=
Specifies the input data set for parameter estimates	PROC HPSEVERITY	INEST=
Specifies the output data set for parameter estimates	PROC HPSEVERITY	OUTEST=
Specifies the output data set for model information	PROC HPSEVERITY	OUTMODELINFO=
Specifies the output data set for statistics of fit	PROC HPSEVERITY	OUTSTAT=
<b>Data Interpretation Options</b>		
Specifies left-censoring	LOSS	LEFTCENSORED=
Specifies left-truncation	LOSS	LEFTTRUNCATED=
Specifies right-censoring	LOSS	RIGHTCENSORED=
Specifies right-truncation	LOSS	RIGHTTRUNCATED=
<b>Model Estimation Options</b>		
Specifies the model selection criterion	PROC HPSEVERITY	CRITERION=
Specifies the method for computing mixture distribution	SCALEMODEL	DFMIXTURE=
Specifies initial values for model parameters	DIST	INIT=
Specifies the objective function symbol	PROC HPSEVERITY	OBJECTIVE=
Specifies the offset variable in the scale regression model	SCALEMODEL	OFFSET=
Specifies the optimization technique	NLOPTIONS	TECHNIQUE=
Specifies the denominator for computing covariance estimates	PROC HPSEVERITY	VARDEF=

Table 9.1 *continued*

Description	Statement	Option
<b>Optimization Termination Criteria</b>		
Absolute function value criterion	NLOPTIONS	ABSCONV=
Absolute function difference convergence criterion	NLOPTIONS	ABSFCNV=
Absolute gradient convergence criterion	NLOPTIONS	ABSGCONV=
Absolute parameter convergence criterion	NLOPTIONS	ABSXCONV=
Relative function convergence criterion	NLOPTIONS	FCONV=
Another relative function convergence criterion	NLOPTIONS	FCONV2=
Used in FCONV, GCONV criterion	NLOPTIONS	FSIZE=
Relative gradient convergence criterion	NLOPTIONS	GCONV=
Maximum number of function calls	NLOPTIONS	MAXFUNC=
Maximum number of iterations	NLOPTIONS	MAXITER=
Upper limit on CPU time in seconds	NLOPTIONS	MAXTIME=
Minimum number of iterations	NLOPTIONS	MINITER=
Relative parameter convergence criterion	NLOPTIONS	XCONV=
Used in XCONV criterion	NLOPTIONS	XSIZE=
<b>Empirical Distribution Function (EDF)</b>		
<b>Estimation Options</b>		
Specifies the nonparametric method of CDF estimation	PROC HPSEVERITY	EMPIRICALCDF=
Specifies the sample to be used for computing the EDF estimates	PROC HPSEVERITY	INITSAMPLE
<b>EMPIRICALCDF=MODIFIEDKM Options</b>		
Specifies the $\alpha$ value for the lower bound on risk set size	PROC HPSEVERITY	ALPHA=
Specifies the $c$ value for the lower bound on risk set size	PROC HPSEVERITY	C=
Specifies the absolute lower bound on risk set size	PROC HPSEVERITY	RSLB=
<b>EMPIRICALCDF=TURNBULL Options</b>		
Specifies that the final EDF estimates be maximum likelihood estimates	PROC HPSEVERITY	ENSUREMLE
Specifies the relative convergence criterion	PROC HPSEVERITY	EPS=
Specifies the maximum number of iterations	PROC HPSEVERITY	MAXITER=
Specifies the threshold below which an EDF estimate is deemed to be 0	PROC HPSEVERITY	ZEROPROB=
<b>Scoring Function Generation Options</b>		
Specifies that scoring functions of all models be written to one package	OUTSCORELIB	COMMONPACKAGE
Specifies the output data set for BY-group identifiers	OUTSCORELIB	OUTBYID=
Specifies the output library for scoring functions	OUTSCORELIB	OUTLIB=

Table 9.1 *continued*

Description	Statement	Option
<b>Displayed Output Options</b>		
Specifies that distributions be listed to the log without estimating any models that use them	DIST	LISTONLY
Suppresses all displayed output	PROC HPSEVERITY	NOPRINT
Specifies which output to display	PROC HPSEVERITY	PRINT=
Specifies the verbosity of messages printed to the log	PROC HPSEVERITY	VERBOSE=
Specifies that distributions be validated without estimating any models that use them	DIST	VALIDATEONLY

## PROC HPSEVERITY Statement

**PROC HPSEVERITY** *options* ;

The PROC HPSEVERITY statement invokes the procedure. You can specify two types of *options* in the PROC HPSEVERITY statement. One set of *options* controls input and output. The other set of *options* controls the model estimation and selection process.

The following *options* control the input data sets used by PROC HPSEVERITY and various forms of output generated by PROC HPSEVERITY. The *options* are listed in alphabetical order:

### COVOUT

specifies that the OUTEST= data set contain the estimate of the covariance structure of the parameters. This option has no effect if you do not specify the OUTEST= option. For more information about how the covariance is reported in the OUTEST= data set, see the section “OUTEST= Data Set” on page 336.

### DATA=SAS-data-set

names the input data set. If you do not specify the DATA= option, then the most recently created SAS data set is used.

### INEST=SAS-data-set

names the input data set that contains the initial values of the parameter estimates to start the optimization process. The initial values that you specify in the INIT= option in the DIST statement take precedence over any initial values that you specify in the INEST= data set. For more information about the variables in this data set, see the section “INEST= Data Set” on page 335.

**INITSAMPLE** (*initsample-option*)**INITSAMPLE** (*initsample-option* ... *initsample-option*)

specifies that a sample of the input data be used for initializing the distribution parameters. If you specify more than one *initsample-option*, then separate them with spaces.

When you do not specify initial values for the distribution parameters, PROC HPSEVERITY needs to compute the empirical distribution function (EDF) estimates as part of the default method for parameter initialization. The EDF estimation process can be expensive, especially when you specify censoring or truncation effects for the loss variable. Furthermore, it is not amenable to parallelism due to the sequential nature of the algorithm for truncation effects. You can use the INITSAMPLE option to specify that only a fraction of the input data be used in order to reduce the time taken to compute the EDF estimates. PROC HPSEVERITY uses the uniform random sampling method to select the sample, the size and randomness of which is controlled by the following *initsample-options*:

**FRACTION=***number*

specifies the fraction, between 0 and 1, of the input data to be used for sampling.

**SEED=***number*

specifies the seed to be used for the uniform random number generator. This option enables you to select the same sample from the same input data across different runs of PROC HPSEVERITY, which can be useful for replicating the results across different runs. If you do not specify the seed value, PROC HPSEVERITY generates a seed that is based on the system clock.

**SIZE=***number*

specifies the size of the sample. If the data are distributed across different nodes, then this size applies to the sample that is prepared at each node. For example, let the input data set of size 100,000 observations be distributed across 10 nodes such that each node has 10,000 observations. If you specify SIZE=1000, then each node computes a local EDF estimate by using a sample of size 1,000 selected randomly from its 10,000 observations. If you specify both of the SIZE= and FRACTION= options, then the value that you specify in the SIZE= option is used and the FRACTION= option is ignored.

If you do not specify the INITSAMPLE option, then a uniform random sample of at most 10,000 observations is used for EDF estimation.

**NOPRINT**

turns off all displayed output. If you specify this option, then any value that you specify for the PRINT= option is ignored.

**OUTEST=***SAS-data-set*

names the output data set to contain estimates of the parameter values and their standard errors for each model whose parameter estimation process converges. For more information about the variables in this data set, see the section “[OUTEST= Data Set](#)” on page 336.

**OUTMODELINFO=***SAS-data-set*

names the output data set to contain the information about each candidate distribution. For more information about the variables in this data set, see the section “[OUTMODELINFO= Data Set](#)” on page 337.

**OUTSTAT=SAS-data-set**

names the output data set to contain the values of statistics of fit for each model whose parameter estimation process converges. For more information about the variables in this data set, see the section “[OUTSTAT= Data Set](#)” on page 338.

**PRINT** <(*global-display-option*)> <=*display-option*>**PRINT** <(*global-display-option*)> <= (*display-option* . . . *display-option*) >

specifies the desired displayed output. If you specify more than one *display-option*, then separate them with spaces and enclose them in parentheses.

The following *global-display-option* is available:

**ONLY**

turns off the default displayed output and displays only the requested output.

The following *display-options* are available:

**ALL**

displays all the output.

**ALLFITSTATS**

displays the comparison of all the statistics of fit for all the models in one table. The table does not include the models whose parameter estimation process does not converge.

**CONVSTATUS**

displays the convergence status of the parameter estimation process.

**DESCSTATS**

displays the descriptive statistics for the response variable. If you specify the SCALEMODEL statement, then this option also displays the descriptive statistics for the regressor variables.

**DISTINFO**

displays the information about each specified distribution. For each distribution, the information includes the name, description, validity status, and number of distribution parameters.

**ESTIMATES****PARMEST**

displays the final estimates of parameters. The estimates are not displayed for models whose parameter estimation process does not converge.

**ESTIMATIONDETAILS**

displays the details of the estimation process for all the models in one table.

**INITIALVALUES**

displays the initial values and bounds used for estimating each model.

**NLOHISTORY**

displays the iteration history of the nonlinear optimization process used for estimating the parameters.



**NLOSUMMARY**

displays the summary of the nonlinear optimization process used for estimating the parameters.

**NONE**

displays none of the output. If you specify this option, then it overrides all other display options. The default displayed output is also suppressed.

**SELECTION****SELECT**

displays the model selection table.

**STATISTICS****FITSTATS**

displays the statistics of fit for each model. The statistics of fit are not displayed for models whose parameter estimation process does not converge.

If you do not specify the PRINT= option or if you do not specify the ONLY *global-display-option*, then the default displayed output is equivalent to specifying PRINT=(SELECTION CONVSTATUS NLOSUMMARY STATISTICS ESTIMATES).

**VARDEF=option**

specifies the denominator to use for computing the covariance estimates. You can specify one of the following values for *option*:

**DF**

specifies that the number of nonmissing observations minus the model degrees of freedom (number of parameters) be used.

**N**

specifies that the number of nonmissing observations be used.

For more information about the covariance estimation, see the section “[Estimating Covariance and Standard Errors](#)” on page 287.

**VERBOSE=verbosity-level**

specifies the amount of messages printed to the SAS log by PROC HPSEVERITY. A higher number prints messages with the same or more detail.

The following *options* control the model estimation and selection process:

**CRITERION=criterion-option****CRITERIA=criterion-option****CRIT=criterion-option**

specifies the model selection criterion.

If you specify two or more candidate models for estimation, then the one with the best value for the selection criterion is chosen as the best model. If you specify the OUTSTAT= data set, then the best model’s observation has a value of 1 for the `_SELECTED_` variable.

You can specify one of the following *criterion-options*:

**AD**

specifies the Anderson-Darling (AD) statistic value, which is computed by using the empirical distribution function (EDF) estimate, as the selection criterion. A lower value is deemed better.

**AIC**

specifies Akaike's information criterion (AIC) as the selection criterion. A lower value is deemed better.

**AICC**

specifies the finite-sample corrected Akaike's information criterion (AICC) as the selection criterion. A lower value is deemed better.

**BIC**

specifies the Schwarz Bayesian information criterion (BIC) as the selection criterion. A lower value is deemed better.

**CUSTOM**

specifies the custom objective function as the selection criterion. You can specify this only if you also specify the **OBJECTIVE=** option. A lower value is deemed better.

**CVM**

specifies the Cramér-von Mises (CvM) statistic value, which is computed by using the empirical distribution function (EDF) estimate, as the selection criterion. A lower value is deemed better.

**KS**

specifies the Kolmogorov-Smirnov (KS) statistic value, which is computed by using the empirical distribution function (EDF) estimate, as the selection criterion. A lower value is deemed better.

**LOGLIKELIHOOD**

**LL**

specifies  $-2 * \log(L)$  as the selection criterion, where  $L$  is the likelihood of the data. A lower value is deemed better. This is the default.

For more information about these *criterion-options*, see the section “[Statistics of Fit](#)” on page 301.

**EMPIRICALCDF | EDF=method**

specifies the method to use for computing the nonparametric or empirical estimate of the cumulative distribution function of the data. You can specify one of the following values for *method*:

**AUTOMATIC | AUTO**

specifies that the method be chosen automatically based on the data specification.

If you do not specify any censoring or truncation, then the standard empirical estimation method (STANDARD) is chosen. If you specify both right-censoring and left-censoring, then Turnbull's estimation method (TURNBULL) is chosen. For all other combinations of censoring and truncation, the Kaplan-Meier method (KAPLANMEIER) is chosen.

**KAPLANMEIER | KM**

specifies that the product limit estimator proposed by Kaplan and Meier (1958) be used. Specification of this method has no effect when you specify both right-censoring and left-censoring.

**MODIFIEDKM | MKM <(options)>**

specifies that the modified product limit estimator be used. Specification of this method has no effect when you specify both right-censoring and left-censoring.

This method allows Kaplan-Meier's product limit estimates to be more robust by ignoring the contributions to the estimate due to small risk-set sizes. The risk set is the set of observations at the risk of failing, where an observation is said to fail if it has not been processed yet and might experience censoring or truncation. You can specify the minimum risk-set size that makes it eligible to be included in the estimation either as an absolute lower bound on the size (RSLB= option) or a relative lower bound determined by the formula  $cn^\alpha$  proposed by Lai and Ying (1991). You can specify the values of  $c$  and  $\alpha$  by using the C= and ALPHA= options, respectively. By default, the relative lower bound is used with values of  $c = 1$  and  $\alpha = 0.5$ . However, you can modify the default by using the following *options*:

**ALPHA | A=number**

specifies the value to use for  $\alpha$  when the lower bound on the risk set size is defined as  $cn^\alpha$ . This value must satisfy  $0 < \alpha < 1$ .

**C=number**

specifies the value to use for  $c$  when the lower bound on the risk set size is defined as  $cn^\alpha$ . This value must satisfy  $c > 0$ .

**RSLB=number**

specifies the absolute lower bound on the risk set size to be included in the estimate.

**NOTURNBULL**

specifies that the method be chosen automatically based on the data specification and that Turnbull's method not be used. This option is the default.

This method first replaces each left-censored or interval-censored observation with an uncensored observation. If the resulting set of observations has any truncated or right-censored observations, then the Kaplan-Meier method (KAPLANMEIER) is chosen. Otherwise, the standard empirical estimation method (STANDARD) is chosen. The observations are modified only for the purpose of computing the EDF estimates; the modification does not affect the parameter estimation process.

**STANDARD | STD**

specifies that the standard empirical estimation method be used. If you specify both right-censoring and left-censoring, then the specification of this method has no effect. If you specify any other combination of censoring or truncation effects, then this method ignores such effects, and can thus result in estimates that are more biased than those obtained with other methods that are more suitable for censored or truncated data.

**TURNBULL | EM <(options)>**

specifies that the Turnbull's method be used. This method is used when you specify both right-censoring and left-censoring. An iterative expectation-maximization (EM) algorithm proposed

by Turnbull (1976) is used to compute the empirical estimates. If you also specify truncation, then the modification suggested by Frydman (1994) is used.

This method is used if you specify both right-censoring and left-censoring and if you explicitly specify the `EMPIRICALCDF=TURNBULL` option.

You can modify the default behavior of the EM algorithm by using the following *options*:

**ENSUREMLE**

specifies that the final EDF estimates be maximum likelihood estimates. The Kuhn-Tucker conditions are computed for the likelihood maximization problem and checked to ensure that EM algorithm converges to maximum likelihood estimates. The method generalizes the method proposed by Gentleman and Geyer (1994) by taking into account any truncation information that you might specify.

**EPS=number**

specifies the maximum relative error to be allowed between estimates of two consecutive iterations. This criterion is used to check the convergence of the algorithm. If you do not specify this option, then PROC HPSEVERITY uses a default value of 1.0E–8.

**MAXITER=number**

specifies the maximum number of iterations to attempt to find the empirical estimates. If you do not specify this option, then PROC HPSEVERITY uses a default value of 500.

**ZEROPROB=number**

specifies the threshold below which an empirical estimate of the probability is considered zero. This option is used to decide if the final estimate is a maximum likelihood estimate. This option does not have an effect if you do not specify the ENSUREMLE option. If you specify the ENSUREMLE option, but do not specify this option, then PROC HPSEVERITY uses a default value of 1.0E–8.

For more information about each of the methods, see the section “[Empirical Distribution Function Estimation Methods](#)” on page 295.

**OBJECTIVE=symbol-name**

names the symbol that represents the objective function in the SAS programming statements that you specify. For each model to be estimated, PROC HPSEVERITY executes the programming statements to compute the value of this symbol for each observation. The values are added across all observations to obtain the value of the objective function. The optimization algorithm estimates the model parameters such that the objective function value is *minimized*. A separate optimization problem is solved for each candidate distribution. If you specify a BY statement, then a separate optimization problem is solved for each candidate distribution within each BY group.

For more information about writing SAS programming statements to define your own objective function, see the section “[Custom Objective Functions](#)” on page 332.

## BY Statement

**BY** *variable-list* ;

A BY statement can be used in the HPSEVERITY procedure to process the input data set in groups of observations defined by the BY variables.

If you specify the BY statement, then PROC HPSEVERITY expects the input data set to be sorted in the order of the BY variables unless you specify the NOTSORTED option.

The BY statement is always supported in the single-machine mode of execution. For the distributed mode, it is supported only when the DATA= data set resides on the client machine. In other words, the BY statement is supported only in the client-data (or local-data) mode of the distributed computing model and not for any of the alongside modes, such as the alongside-the-database or alongside-HDFS mode.

## DIST Statement

**DIST** *distribution-name-or-keyword* < (*distribution-option*) > < *distribution-name-or-keyword* < (*distribution-option*) > > ... ;

The DIST statement specifies candidate distributions to be estimated by the HPSEVERITY procedure. You can specify multiple DIST statements, and each statement can contain one or more distribution specifications.

For your convenience, PROC HPSEVERITY provides the following 10 different predefined distributions (the name in the parentheses is the name to use in the DIST statement): Burr (BURR), exponential (EXP), gamma (GAMMA), generalized Pareto (GPD), inverse Gaussian or Wald (IGAUSS), lognormal (LOGN), Pareto (PARETO), Tweedie (TWEEDIE), scaled Tweedie (STWEEDIE), and Weibull (WEIBULL). These are described in detail in the section “[Predefined Distributions](#)” on page 274.

You can specify any of the predefined distributions or any distribution that you have defined. If a distribution that you specify is not a predefined distribution, then you must submit the CMPLIB= system option with appropriate libraries before you submit the PROC HPSEVERITY step to enable the procedure to find the functions associated with your distribution. The predefined distributions are defined in the Sashelp.Svrtldist library. However, you are not required to specify this library in the CMPLIB= system option. For more information about defining your own distributions, see the section “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308.

As a convenience, you can also use a shortcut keyword to indicate a list of distributions. You can specify one or more of the following keywords:

### **ALL**

specifies all the predefined distributions and the distributions that you have defined in the libraries that you specify in the CMPLIB= system option. In addition to the eight predefined distributions included by the **PREDEFINED** keyword, this list also includes the Tweedie and scaled Tweedie distributions that are defined in the Sashelp.Svrtldist library.

### **PREDEFINED**

specifies the list of eight predefined distributions: BURR, EXP, GAMMA, GPD, IGAUSS, LOGN, PARETO, and WEIBULL. Although the TWEEDIE and STWEEDIE distributions are available in the Sashelp.Svrtldist library along with these eight distributions, they are not included by this keyword. If

you want to fit the TWEEDIE and STWEEDIE distributions, then you must specify them explicitly or use the `_ALL_` keyword.

### **`_USER_`**

specifies the list of all the distributions that you have defined in the libraries that you specify in the `CMPLIB=` system option. This list does not include the distributions defined in the `Sashelp.Svrtldist` library, even if you specify the `Sashelp.Svrtldist` library in the `CMPLIB=` option.

The use of these keywords, especially `_ALL_`, can result in a large list of distributions, which might take a longer time to estimate. A warning is printed to the SAS log if the number of total distribution models to estimate exceeds 10.

The following *distribution-option* values can be used in the `DIST` statement for a distribution name that is not a shortcut keyword:

### **`INIT=(name=value ... name=value)`**

specifies the initial values to be used for the distribution parameters to start the parameter estimation process. You must specify the values by parameter names and the parameter names must match the names used in the model definition. For example, let a model `M`'s definition contain a `M_PDF` function with following signature:

```
function M_PDF(x, alpha, beta);
```

For this model, the names `alpha` and `beta` must be used for the `INIT` option. The names are case-insensitive. If you do not specify initial values for some parameters in the `INIT` statement, then a default value of 0.001 is assumed for those parameters. If you specify an incorrect parameter, PROC HPSEVERITY prints a warning to the SAS log and does not fit the model. All specified values must be nonmissing.

If you are modeling regression effects, then the initial value of the first distribution parameter (`alpha` in the preceding example) should be the initial *base* value of the scale parameter or log-transformed scale parameter. For more information, see the section “[Estimating Regression Effects](#)” on page 290.

The use of `INIT=` option is one of the three methods available for initializing the parameters. For more information, see the section “[Parameter Initialization](#)” on page 289. If none of the initialization methods is used, then PROC HPSEVERITY initializes all parameters to 0.001.

You can specify the following *preprocess-options* in the `DIST` statement:

### **`LISTONLY`**

specifies that the list of all candidate distributions be printed to the SAS log without doing any further processing on them. This option is especially useful when you use a shortcut keyword to include a list of distributions. It enables you to find out which distributions are included by the keyword.

### **`VALIDATEONLY`**

specifies that all candidate distributions be checked for validity without doing any further processing on them. If a distribution is invalid, the reason for invalidity is written to the SAS log. If all distributions are valid, then the distribution information is written to the SAS log. The information includes name, description, validity status (valid or invalid), and number of distribution parameters. The information is not written to the SAS log if you specify an `OUTMODELINFO=` data set or the `PRINT=DISTINFO` or `PRINT=ALL` option in the PROC HPSEVERITY statement. This option is especially useful

when you specify your own distributions or when you specify the `_USER_` or `_ALL_` keywords in the `DIST` statement. It enables you to check whether your custom distribution definitions satisfy PROC HPSEVERITY's requirements for the specified modeling task. It is recommended that you specify the `SCALEMODEL` statement if you intend to fit a model with regression effects, because the `SCALEMODEL` statement instructs PROC HPSEVERITY to perform additional checks to validate whether regression effects can be modeled on each candidate distribution.

## LOSS Statement

**LOSS** < *response-variable-name* > < / *censoring-truncation-options* > ;

The `LOSS` statement specifies the name of the response or loss variable whose distribution needs to be modeled. You can also specify additional options to indicate any truncation or censoring of the response. The specification of response variable is optional if you specify at least one type of censoring. You must specify a response variable if you do not specify any censoring. If you specify more than one `LOSS` statement, then the first statement is used.

All the analysis variables that you specify in this statement must be present in the input data set that you specify by using the `DATA=` option in the PROC HPSEVERITY statement. The response variable is expected to have nonmissing values. If the variable has a missing value in an observation, then a warning is written to the SAS log and that observation is ignored.

The following *censoring-truncation-options* can be used in the `LOSS` statement:

**LEFTCENSORED**=*variable-name*

**LC**=*variable-name*

**LEFTCENSORED**=*number*

**LC**=*number*

specifies the left-censoring variable or a global left-censoring limit.

You can use the *variable-name* argument to specify a data set variable that contains the left-censoring limit. If the value of this variable is missing, then PROC HPSEVERITY assumes that such observations are not left-censored.

Alternatively, you can use the *number* argument to specify a left-censoring limit value that applies to all the observations in the data set. This limit must be a nonzero positive number.

By the definition of left-censoring, an exact value of the response is not known when it is less than or equal to the left-censoring limit. If you specify the response variable and the value of that variable is less than or equal to the value of the left-censoring limit for some observations, then PROC HPSEVERITY treats such observations as left-censored and the value of the response variable is ignored. If you specify the response variable and the value of that variable is greater than the value of the left-censoring limit for some observations, then PROC HPSEVERITY assumes that such observations are not left-censored and the value of the left-censoring limit is ignored.

If you specify both right-censoring and left-censoring limits, then the left-censoring limit must be greater than or equal to the right-censoring limit. If both limits are identical, then the observation is assumed to be uncensored.

For more information about left-censoring, see the section “[Censoring and Truncation](#)” on page 284.

**LEFTTRUNCATED** | **LT=***variable-name*

**LT=***variable-name*

**LEFTTRUNCATED=***number*

**LT=***number*

specifies the left-truncation variable or a global left-truncation threshold.

You can use the *variable-name* argument to specify a data set variable that contains the left-truncation threshold. If the value of this variable is missing or 0 for some observations, then PROC HPSEVERITY assumes that such observations are not left-truncated.

Alternatively, you can use the *number* argument to specify a left-truncation threshold that applies to all the observations in the data set. This threshold must be a nonzero positive number.

It is assumed that the response variable contains the observed values. By the definition of left-truncation, you can observe only a value that is greater than the left-truncation threshold. If a response variable value is less than or equal to the left-truncation threshold, a warning is printed to the SAS log, and the observation is ignored. For more information about left-truncation, see the section “[Censoring and Truncation](#)” on page 284.

**RIGHTCENSORED** | **RC=***variable-name*

**RC=***variable-name*

**RIGHTCENSORED=***number*

**RC=***number*

specifies the right-censoring variable or a global right-censoring limit.

You can use the *variable-name* argument to specify a data set variable that contains the right-censoring limit. If the value of this variable is missing, then PROC HPSEVERITY assumes that such observations are not right-censored.

Alternatively, you can use the *number* argument to specify a right-censoring limit value that applies to all the observations in the data set. This limit must be a nonzero positive number.

By the definition of right-censoring, an exact value of the response is not known when it is greater than or equal to the right-censoring limit. If you specify the response variable and the value of that variable is greater than or equal to the value of the right-censoring limit for some observations, then PROC HPSEVERITY treats such observations as right-censored and the value of the response variable is ignored. If you specify the response variable and the value of that variable is less than the value of the right-censoring limit for some observations, then PROC HPSEVERITY assumes that such observations are not right-censored and the value of the right-censoring limit is ignored.

If you specify both right-censoring and left-censoring limits, then the left-censoring limit must be greater than or equal to the right-censoring limit. If both limits are identical, then the observation is assumed to be uncensored.

For more information about right-censoring, see the section “[Censoring and Truncation](#)” on page 284.



**RIGHTTRUNCATED** | **RT=***variable-name*

**RT=***variable-name*

**RIGHTTRUNCATED=***number*

**RT=***number*

specifies the right-truncation variable or a global right-truncation threshold.

You can use the *variable-name* argument to specify a data set variable that contains the right-truncation threshold. If the value of this variable is missing for some observations, then PROC HPSEVERITY assumes that such observations are not right-truncated.

Alternatively, you can use the *number* argument to specify a right-truncation threshold that applies to all the observations in the data set. This threshold must be a nonzero positive number.

It is assumed that the response variable contains the observed values. By the definition of right-truncation, you can observe only a value that is less than or equal to the right-truncation threshold. If a response variable value is greater than the right-truncation threshold, a warning is printed to the SAS log, and the observation is ignored. For more information about right-truncation, see the section [“Censoring and Truncation”](#) on page 284.

---

## NLOPTIONS Statement

**NLOPTIONS** *options* ;

The HPSEVERITY procedure uses the nonlinear optimization (NLO) subsystem to perform nonlinear optimization of the likelihood function to obtain the estimates of distribution and regression parameters. You can use the NLOPTIONS statement to control different aspects of this optimization process. If you specify more than one NLOPTIONS statement, then the first statement is used.

For most problems, the default settings of the optimization process are adequate. However, in some cases it might be useful to change the optimization technique or to change the maximum number of iterations. The following statement uses the MAXITER= option to set the maximum number of iterations to 200 and uses the TECH= option to change the optimization technique to the double-dogleg optimization (DBLDOG) rather than the default technique, the trust region optimization (TRUREG), used in the HPSEVERITY procedure:

```
nloptions tech=dbldog maxiter=200;
```

The following *options* values can be used in the NLOPTIONS statement:

**ABSCONV=***r*

**ABSTOL=***r*

specifies an absolute function value convergence criterion. For minimization, termination requires  $f(\theta^{(k)}) \leq r$ . The default value of *r* is the negative square root of the largest double-precision value, which serves only as a protection against overflows.

**ABSFCNV= $r$** **ABSFTOL= $r$** 

specifies an absolute function difference convergence criterion. For all techniques except NMSIMP, termination requires a small change of the function value in successive iterations:

$$|f(\theta^{(k-1)}) - f(\theta^{(k)})| \leq r$$

The same formula is used for the NMSIMP technique, but  $\theta^{(k)}$  is defined as the vertex with the lowest function value, and  $\theta^{(k-1)}$  is defined as the vertex with the highest function value in the simplex. The default value is  $r = 0$ .

**ABSGCNV= $r$** **ABSGTOL= $r$** 

specifies an absolute gradient convergence criterion. Termination requires the maximum absolute gradient element to be small:

$$\max_j |g_j(\theta^{(k)})| \leq r$$

This criterion is not used by the NMSIMP technique. The default value is  $r=1\text{E-}5$ .

**ABSXCNV= $r$** **ABSXTOL= $r$** 

specifies an absolute parameter convergence criterion. For all techniques except NMSIMP, termination requires a small Euclidean distance between successive parameter vectors,

$$\|\theta^{(k)} - \theta^{(k-1)}\|_2 \leq r$$

For the NMSIMP technique, termination requires either a small length  $\alpha^{(k)}$  of the vertices of a restart simplex,

$$\alpha^{(k)} \leq r$$

or a small simplex size,

$$\delta^{(k)} \leq r$$

where the simplex size  $\delta^{(k)}$  is defined as the L1 distance from the simplex vertex  $\xi^{(k)}$  with the smallest function value to the other simplex points  $\theta_l^{(k)} \neq \xi^{(k)}$ :

$$\delta^{(k)} = \sum \|\theta_l^{(k)} - \xi^{(k)}\|_1$$

The default is  $r=1\text{E-}8$  for the NMSIMP technique and  $r = 0$  otherwise.

**FCONV= $r$** **FTOL= $r$** 

specifies a relative function convergence criterion. For all techniques except NMSIMP, termination requires a small relative change of the function value in successive iterations,

$$\frac{|f(\theta^{(k)}) - f(\theta^{(k-1)})|}{\max(|f(\theta^{(k-1)})|, \text{FSIZE})} \leq r$$

where FSIZE is defined by the FSIZE= option. The same formula is used for the NMSIMP technique, but  $\theta^{(k)}$  is defined as the vertex with the lowest function value, and  $\theta^{(k-1)}$  is defined as the vertex with the highest function value in the simplex.

The default value is  $r = 2\epsilon$ , where  $\epsilon$  denotes the machine precision constant, which is the smallest double-precision floating-point number such that  $1 + \epsilon > 1$ .

**FCONV2=r****FTOL2=r**

specifies another function convergence criterion.

For all techniques except NMSIMP, termination requires a small predicted reduction of the objective function:

$$df^{(k)} \approx f(\theta^{(k)}) - f(\theta^{(k)} + s^{(k)})$$

The predicted reduction

$$\begin{aligned} df^{(k)} &= -g^{(k)T} s^{(k)} - \frac{1}{2} s^{(k)T} H^{(k)} s^{(k)} \\ &= -\frac{1}{2} s^{(k)T} g^{(k)} \\ &\leq r \end{aligned}$$

is computed by approximating the objective function  $f$  by the first two terms of the Taylor series and substituting the Newton step

$$s^{(k)} = -[H^{(k)}]^{-1} g^{(k)}$$

For the NMSIMP technique, termination requires a small standard deviation of the function values of the  $p + 1$  simplex vertices  $\theta_l^{(k)}$ ,  $l = 0, \dots, p$ ,  $\sqrt{\frac{1}{p+1} \sum_l [f(\theta_l^{(k)}) - \bar{f}(\theta^{(k)})]^2} \leq r$  where  $\bar{f}(\theta^{(k)}) = \frac{1}{p+1} \sum_l f(\theta_l^{(k)})$ . If there are  $p_{act}$  boundary constraints active at  $\theta^{(k)}$ , the mean and standard deviation are computed only for the  $p + 1 - p_{act}$  unconstrained vertices.

The default value is  $r=1E-6$  for the NMSIMP technique and  $r = 0$  otherwise.

**FSIZE=r**

specifies the FSIZE parameter of the relative function and relative gradient termination criteria. The default value is  $r = 0$ . For more information, see the FCONV= and GCONV= options.

**GCONV=r****GTOL=r**

specifies a relative gradient convergence criterion. For all techniques except CONGRA and NMSIMP, termination requires that the normalized predicted function reduction is small,

$$\frac{g(\theta^{(k)})^T [H^{(k)}]^{-1} g(\theta^{(k)})}{\max(|f(\theta^{(k)})|, \text{FSIZE})} \leq r$$

where FSIZE is defined by the FSIZE= option. For the CONGRA technique (where a reliable Hessian estimate  $H$  is not available), the following criterion is used:

$$\frac{\|g(\theta^{(k)})\|_2^2 \|s(\theta^{(k)})\|_2}{\|g(\theta^{(k)}) - g(\theta^{(k-1)})\|_2 \max(|f(\theta^{(k)})|, \text{FSIZE})} \leq r$$

This criterion is not used by the NMSIMP technique. The default value is  $r = 1E - 8$ .

**MAXFUNC=*i*****MAXFU=*i***

specifies the maximum number *i* of function calls in the optimization process. The default values are

- TRUREG, NRRIDG, NEWRAP: 125
- QUANEW, DBLDOG: 500
- CONGRA: 1000
- NMSIMP: 3000

Note that the optimization can terminate only after completing a full iteration. Therefore, the number of function calls that is actually performed can exceed the number that you specify by the MAXFUNC= option.

**MAXITER=*i*****MAXIT=*i***

specifies the maximum number *i* of iterations in the optimization process. The default values are

- TRUREG, NRRIDG, NEWRAP: 50
- QUANEW, DBLDOG: 200
- CONGRA: 400
- NMSIMP: 1000

These default values are also valid when you specify a missing value for *i*.

**MAXTIME=*r***

specifies an upper limit of *r* seconds of CPU time for the optimization process. The default value is the largest floating-point double representation of your computer. The time that you specify in the MAXTIME= option is checked only once at the end of each iteration. Therefore, the actual running time can be much longer than that you specify by the MAXTIME= option. The actual running time includes the rest of the time needed to finish the iteration and the time needed to generate the output of the results.

**MINITER=*i*****MINIT=*i***

specifies the minimum number of iterations. The default value is 0. If you request more iterations than are actually needed for convergence to a stationary point, the optimization algorithms can behave strangely. For example, the effect of rounding errors can prevent the algorithm from continuing for the required number of iterations.

**TECHNIQUE=*name*****TECH=*name***

specifies the optimization technique. Valid values for *name* are as follows:

CONGRA	performs a conjugate-gradient optimization.
DBLDOG	performs a version of double-dogleg optimization.
NMSIMP	performs a Nelder-Mead simplex optimization.

NONE	does not perform any optimization. This option can be used as follows: <ul style="list-style-type: none"> <li>to perform a grid search without optimization</li> <li>to compute estimates and predictions that cannot be obtained efficiently with any of the optimization techniques</li> </ul>
NEWRAP	performs a Newton-Raphson optimization that combines a line-search algorithm with ridging.
NRRIDG	performs a Newton-Raphson optimization with ridging.
QUANEW	performs a quasi-Newton optimization.
TRUREG	performs a trust region optimization. This is the default estimation method.

For more information about optimization algorithms, see the section “[Details of Optimization Algorithms](#)” on page 287.

**XCONV=*r***

**XTOL=*r***

specifies the relative parameter convergence criterion. For all techniques except NMSIMP, termination requires a small relative parameter change in subsequent iterations:

$$\frac{\max_j |\theta_j^{(k)} - \theta_j^{(k-1)}|}{\max(|\theta_j^{(k)}|, |\theta_j^{(k-1)}|, \text{XSIZE})} \leq r$$

For the NMSIMP technique, the same formula is used, but  $\theta_j^{(k)}$  is defined as the vertex that has the lowest function value and  $\theta_j^{(k-1)}$  is defined as the vertex that has the highest function value in the simplex. The default value is  $r=1\text{E}-8$  for the NMSIMP technique and  $r = 0$  otherwise.

**XSIZE=*r***

specifies the XSIZE parameter of the relative parameter termination criterion. The value of  $r$  must be greater than or equal to 0; the default is  $r = 0$ . For more information, see the XCONV= option.

---

## OUTSCORELIB Statement (Experimental)

**OUTSCORELIB <OUTLIB=> *fcmp-library-name options* ;**

The OUTSCORELIB statement specifies the library to write scoring functions to. Scoring functions enable you to easily execute a distribution function on the fitted parameters of the distribution without going through a potentially complex process of extracting the fitted parameter estimates from other output such as the OUTEST= data set that is created by PROC HPSEVERITY.

You must specify the following option as the first option in the statement:

**OUTLIB=*fcmp-library-name***

names the FCMP library to contain the scoring functions. PROC HPSEVERITY writes the scoring functions to the FCMP library named *fcmp-library-name*. If a library or data set named *fcmp-library-name* already exists, PROC HPSEVERITY deletes it before proceeding.

This option is similar to the OUTLIB= option that you would specify in a PROC FCMP statement, except that *fcmp-library-name* must be a two-level name whereas the OUTLIB= option in the PROC FCMP statement requires a three-level name. The third level of a three-level name specifies the package to which the functions belong. You do not need to specify the package name in the *fcmp-library-name*, because PROC HPSEVERITY automatically creates the package for you. By default, a separate package is created for each distribution that has not failed to converge. Each package is named for a distribution. For example, if you define and fit a distribution named *mydist*, and if *mydist* does not fail to converge, then PROC HPSEVERITY creates a package named *mydist* in the OUTLIB= library that you specify. Further, let the definition of the *mydist* distribution contain three distribution functions, *mydist\_PDF*(*x*,*Parm1*,*Parm2*), *mydist\_LOGCDF*(*x*,*Parm1*,*Parm2*), and *mydist\_XYZ*(*x*,*Parm1*,*Parm2*). If you specify the OUTSCORELIB statement

```
outscorelib outlib=sasuser.scorefunc;
```

then the Sasuser.Scorefunc library contains the following three functions in a package named *mydist*: *SEV\_PDF*(*x*), *SEV\_LOGCDF*(*x*), and *SEV\_XYZ*(*x*).

The key feature of scoring functions is that they do not require the parameter arguments (*Parm1* and *Parm2* in this example). The fitted parameter estimates are encoded inside the scoring function so that you can compute or score the value of each function for a given value of the loss variable without having to know or extract the parameter estimates through some other means.

For convenience, you can omit the OUTLIB= portion of the specification and just specify the name, as in the following example:

```
outscorelib sasuser.scorefunc;
```

When the HPSEVERITY procedure executes successfully, the *fcmp-library-name* is appended to the CMPLIB system option, so you can immediately start using the scoring functions in a DATA step or PROC FCMP step.

You can specify the following *options* in the OUTSCORELIB statement:

### COMMONPACKAGE

#### ONEPACKAGE

requests that only one common package be created to contain all the scoring functions.

If you specify this option, then all the scoring functions are created in a package called *sevfit*. For each distribution function that has the name *distribution\_suffix*, the name of the corresponding scoring function is formed as *SEV\_suffix\_distribution*. For example, the scoring function of the distribution function 'MYDIST\_BAR' is named 'SEV\_BAR\_MYDIST'.

If you do not specify this option, then all scoring functions for a distribution are created in a package that has the same name as the distribution, and for each distribution function that has the name *distribution\_suffix*, the name of the corresponding scoring function is formed as *SEV\_suffix*. For example, the scoring function of the distribution function 'MYDIST\_BAR' is named 'SEV\_BAR'.

#### OUTBYID=SAS-data-set

names the output data set to contain the unique identifier for each BY group. This unique identifier is used as part of the name of the package or scoring function for each distribution. This is a required option when you specify a BY statement in PROC HPSEVERITY.

The OUTBYID= data set contains one observation per BY group and a variable named `_ID_` in addition to the BY variables that you specify in the BY statement. The `_ID_` variable contains the unique identifier for each BY group. The identifier of the BY group is the decimal representation of the sequence number of the BY group. The first BY group has an identifier of 1, the second BY group has an identifier of 2, the tenth BY group has an identifier of 10, and so on.

If you do not specify the COMMONPACKAGE option in the OUTSCORELIB statement, then for each distribution, PROC HPSEVERITY creates as many packages as the number of BY groups. The unique BY-group identifier is used as a suffix for the package name. For example, if your DATA= data set has three BY groups and if you specify the OUTSCORELIB statement

```
outscorelib outlib=sasuser.byscorefunc outbyid=sasuser.byid;
```

then for the distribution ‘MYDIST’, the Sasuser.Byscorefunc library contains the three packages ‘MYDIST1’, ‘MYDIST2’, and ‘MYDIST3’, and each package contains one scoring function named ‘SEV\_BAR’ for each distribution function named ‘MYDIST\_BAR’.

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, PROC HPSEVERITY creates as many versions of the distribution function as the number of BY groups. The unique BY-group identifier is used as a suffix for the function name. Extending the previous example, if you specify the OUTSCORELIB statement with the COMMONPACKAGE option,

```
outscorelib outlib=sasuser.byscorefunc outbyid=sasuser.byid commonpackage;
```

then for the distribution function ‘MYDIST\_BAR’ of the distribution ‘MYDIST’, the Sasuser.Byscorefunc library contains the following three scoring functions: ‘SEV\_BAR\_MYDIST1’, ‘SEV\_BAR\_MYDIST2’, and ‘SEV\_BAR\_MYDIST3’. All the scoring functions are created in one common package named *sevfit*.

For both the preceding examples, the Sasuser.Byid data set contains three observations, one for each BY group. The value of the `_ID_` variable is 1 for the first BY group, 2 for the second BY group, and 3 for the third BY group.

For more information about scoring functions, see the section “[Scoring Functions \(Experimental\)](#)” on page 325.

---

## PERFORMANCE Statement

### PERFORMANCE options ;

The PERFORMANCE statement defines performance parameters for distributed and multithreaded computing, passes variables that describe the distributed computing environment, and requests detailed results about the performance characteristics of PROC HPSEVERITY.

You can also use the PERFORMANCE statement to control whether a high-performance analytical procedure executes in single-machine or distributed mode.

The PERFORMANCE statement is documented further in the section “[PERFORMANCE Statement](#)” on page 39 of Chapter 3, “[Shared Concepts and Topics](#).”

## SCALEMODEL Statement

**SCALEMODEL** *regressor-variable-list* < / *scalemodel-options* > ;

The SCALEMODEL statement specifies regression variables. All the variables that you specify in this statement must be present in the input data set that you specify by using the DATA= option in the PROC HPSEVERITY statement. The scale parameter of each candidate distribution is linked to a linear combination of these regression variables along with an intercept. If a distribution does not have a scale parameter, then a model based on that distribution is not estimated. If you specify more than one SCALEMODEL statement, then the first statement is used.

The regressor variables are expected to have nonmissing values. If any of the variables has a missing value in an observation, then a warning is written to the SAS log and that observation is ignored.

For more information about modeling regression effects, see the section “[Estimating Regression Effects](#)” on page 290.

You can specify the following *scalemodel-options* in the SCALEMODEL statement:

**DFMIXTURE=***method-name* < (*method-options*) >

specifies the method for computing representative estimates of the cumulative distribution function (CDF).

When you specify regression variables, the scale of the distribution depends on the values of the regressors. For a given distribution family, each observation in the input data set implies a different scaled version of the distribution. To compute estimates of CDF that are comparable across different distribution families, PROC HPSEVERITY needs to construct a single representative distribution from all such distributions. You can specify one of the following *method-name* values to specify the method that is used to construct the representative distribution. For more information about each of the methods, see the section “[CDF Estimates with Regression Effects](#)” on page 293.

### FULL

specifies that the representative distribution be the mixture of  $N$  distributions such that each distribution has a scale value that is implied by each of the  $N$  observations that are used for estimation. This method is the slowest.

### MEAN

specifies that the representative distribution be the one-point mixture of the distribution whose scale value is the mean of the  $N$  scale values that are implied by the  $N$  observations that are used for estimation. If you do not specify the DFMIXTURE= option, then this method is used by default. This is also the fastest method.

### QUANTILE < (K=q) >

specifies that the representative distribution be the mixture of a fixed number of distributions whose scale values are the quantiles from the sample of  $N$  scale values that are implied by the  $N$  observations in the current BY group (or in the entire DATA= data set if you do not specify the BY statement).



You can use the **K=** option to specify the number of distributions in the mixture. If you specify **K=q**, then the mixture contains  $(q - 1)$  distributions such that each distribution has as its scale one of the  $(q - 1)$ -quantiles.

If you do not specify the **K=** option, then PROC HPSEVERITY uses the default of 2, which implies the use of a one-point mixture with a distribution whose scale value is the median of all scale values.

**RANDOM** <(random-method-options)>

specifies that the representative distribution be the mixture of a fixed number of distributions whose scale values are the scale values that are implied by a randomly chosen subset of the set of all observations in the current BY group (or in the entire DATA= data set if you do not specify the BY statement). The same subset of observations is used for each distribution family.

You can specify the following *random-method-options* to specify how the subset is chosen:

**K=r**

specifies the number of distributions to include in the mixture. If you do not specify this option, then PROC HPSEVERITY uses the default of 15.

**SEED=number**

specifies the seed that is used to generate the uniform random sample of observation indices. If you do not specify this option, then PROC HPSEVERITY generates a seed internally that is based on the current value of the system clock.

**OFFSET=offset-variable-name**

specifies the name of the offset variable in the scale regression model. An offset variable is a regressor variable whose regression coefficient is known to be 1. For more information, see the section “[Offset Variable](#)” on page 291.

---

## WEIGHT Statement

**WEIGHT** *variable-name* ;

The WEIGHT statement specifies the name of a variable whose values represent the weight of each observation. PROC HPSEVERITY associates a weight of  $w$  to each observation, where  $w$  is the value of the WEIGHT variable for the observation. If the weight value is missing or less than or equal to 0, then the observation is ignored and a warning is written to the SAS log. When you do not specify the WEIGHT statement, each observation is assigned a weight of 1. If you specify more than one WEIGHT statement, then the last statement is used.

The weights are normalized so that they add up to the actual sample size. In particular, the weight of each observation is multiplied by  $\frac{N}{\sum_{i=1}^N w_i}$ , where  $N$  is the sample size.

---

## Programming Statements

You can use a series of programming statements that use variables in the input data set that you specify in DATA= option in the PROC HPSEVERITY statement to assign a value to an objective function symbol. You must specify the objective function symbol by using the **OBJECTIVE=** option in the PROC HPSEVERITY statement. If you do not specify the **OBJECTIVE=** option in the PROC HPSEVERITY statement, then the programming statements are ignored and models are estimated using the maximum likelihood method.

You can use most DATA step statements and functions in your program. Any additional functions, restrictions, and differences are listed in the section “[Custom Objective Functions](#)” on page 332.

---

## Details: HPSEVERITY Procedure

---

### Predefined Distributions

PROC HPSEVERITY assumes the following model for the response variable  $Y$

$$Y \sim \mathcal{F}(\Theta)$$

where  $\mathcal{F}$  is a continuous probability distribution with parameters  $\Theta$ . The model hypothesizes that the observed response is generated from a stochastic process that is governed by the distribution  $\mathcal{F}$ . This model is usually referred to as the error model. Given a representative input sample of response variable values, PROC HPSEVERITY estimates the model parameters for any distribution  $\mathcal{F}$  and computes the statistics of fit for each model. This enables you to find the distribution that is most likely to generate the observed sample.

A set of predefined distributions is provided with the HPSEVERITY procedure. A summary of the distributions is provided in [Table 9.2](#). For each distribution, the table lists the name of the distribution that should be used in the DIST statement, the parameters of the distribution along with their bounds, and the mathematical expressions for the probability density function (PDF) and cumulative distribution function (CDF) of the distribution.

All the predefined distributions, except LOGN and TWEEDIE, are parameterized such that their first parameter is the scale parameter. For LOGN, the first parameter  $\mu$  is a log-transformed scale parameter. TWEEDIE does not have a scale parameter. The presence of scale parameter or a log-transformed scale parameter enables you to use all of the predefined distributions, except TWEEDIE, as a candidate for estimating regression effects.

A distribution model is associated with each predefined distribution. You can also define your own distribution model, which is a set of functions and subroutines that you define by using the FCMP procedure. For more information, see the section “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308.

**Table 9.2** Predefined HPSEVERITY Distributions

Name	Distribution	Parameters	PDF ( $f$ ) and CDF ( $F$ )
BURR	Burr	$\theta > 0, \alpha > 0,$ $\gamma > 0$	$f(x) = \frac{\alpha \gamma z^\gamma}{x(1+z^\gamma)^{(\alpha+1)}}$ $F(x) = 1 - \left(\frac{1}{1+z^\gamma}\right)^\alpha$
EXP	Exponential	$\theta > 0$	$f(x) = \frac{1}{\theta} e^{-z}$ $F(x) = 1 - e^{-z}$
GAMMA	Gamma	$\theta > 0, \alpha > 0$	$f(x) = \frac{z^\alpha e^{-z}}{x \Gamma(\alpha)}$ $F(x) = \frac{\gamma(\alpha, z)}{\Gamma(\alpha)}$
GPD	Generalized Pareto	$\theta > 0, \xi > 0$	$f(x) = \frac{1}{\theta} (1 + \xi z)^{-1-1/\xi}$ $F(x) = 1 - (1 + \xi z)^{-1/\xi}$
IGAUSS	Inverse Gaussian (Wald)	$\theta > 0, \alpha > 0$	$f(x) = \frac{1}{\theta} \sqrt{\frac{\alpha}{2\pi z^3}} e^{-\frac{\alpha(z-1)^2}{2z}}$ $F(x) = \Phi\left((z-1)\sqrt{\frac{\alpha}{z}}\right) + \Phi\left(-(z+1)\sqrt{\frac{\alpha}{z}}\right) e^{2\alpha}$
LOGN	Lognormal	$\mu$ (no bounds), $\sigma > 0$	$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2}$ $F(x) = \Phi\left(\frac{\log(x)-\mu}{\sigma}\right)$
PARETO	Pareto	$\theta > 0, \alpha > 0$	$f(x) = \frac{\alpha \theta^\alpha}{(x+\theta)^{\alpha+1}}$ $F(x) = 1 - \left(\frac{\theta}{x+\theta}\right)^\alpha$
TWEEDIE	Tweedie <sup>6</sup>	$p > 1, \mu > 0,$ $\phi > 0$	$f(x) = a(x, \phi) \exp\left[\frac{1}{\phi} \left(\frac{x\mu^{1-p}}{1-p} - \kappa(\mu, p)\right)\right]$ $F(x) = \int_0^x f(t) dt$
STWEEDIE	Scaled Tweedie <sup>6</sup>	$\theta > 0, \lambda > 0,$ $1 < p < 2$	$f(x) = a(x, \theta, \lambda, p) \exp\left(-\frac{x}{\theta} - \lambda\right)$ $F(x) = \int_0^x f(t) dt$
WEIBULL	Weibull	$\theta > 0, \tau > 0$	$f(x) = \frac{1}{x} \tau z^\tau e^{-z^\tau}$ $F(x) = 1 - e^{-z^\tau}$

Notes:

1.  $z = x/\theta$ , wherever  $z$  is used.
2.  $\theta$  denotes the scale parameter for all the distributions. For LOGN,  $\log(\theta) = \mu$ .
3. Parameters are listed in the order in which they are defined in the distribution model.
4.  $\gamma(a, b) = \int_0^b t^{a-1} e^{-t} dt$  is the lower incomplete gamma function.
5.  $\Phi(y) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{y}{\sqrt{2}}\right)\right)$  is the standard normal CDF.
6. For more information, see the section “[Tweedie Distributions](#)” on page 276.

## Tweedie Distributions

Tweedie distributions are a special case of the exponential dispersion family (Jørgensen 1987) with a property that the variance of the distribution is equal to  $\phi\mu^p$ , where  $\mu$  is the mean of the distribution,  $\phi$  is a dispersion parameter, and  $p$  is an index parameter as discovered by Tweedie (1984). The distribution is defined for all values of  $p$  except for values of  $p$  in the open interval  $(0, 1)$ . Many important known distributions are a special case of Tweedie distributions including normal ( $p=0$ ), Poisson ( $p=1$ ), gamma ( $p=2$ ), and the inverse Gaussian ( $p=3$ ). Apart from these special cases, the probability density function (PDF) of the Tweedie distribution does not have an analytic expression. For  $p > 1$ , it has the form (Dunn and Smyth 2005),

$$f(x; \mu, \phi, p) = a(x, \phi) \exp \left[ \frac{1}{\phi} \left( \frac{x\mu^{1-p}}{1-p} - \kappa(\mu, p) \right) \right]$$

where  $\kappa(\mu, p) = \mu^{2-p}/(2-p)$  for  $p \neq 2$  and  $\kappa(\mu, p) = \log(\mu)$  for  $p = 2$ . The function  $a(x, \phi)$  does not have an analytical expression. It is typically evaluated using series expansion methods described in Dunn and Smyth (2005).

For  $1 < p < 2$ , the Tweedie distribution is a compound Poisson-gamma mixture distribution, which is the distribution of  $S$  defined as

$$S = \sum_{i=1}^N X_i$$

where  $N \sim \text{Poisson}(\lambda)$  and  $X_i \sim \text{gamma}(\alpha, \theta)$  are independent and identically distributed gamma random variables with shape parameter  $\alpha$  and scale parameter  $\theta$ . At  $X = 0$ , the density is a probability mass that is governed by the Poisson distribution, and for values of  $X > 0$ , it is a mixture of gamma variates with Poisson mixing probability. The parameters  $\lambda$ ,  $\alpha$ , and  $\theta$  are related to the natural parameters  $\mu$ ,  $\phi$ , and  $p$  of the Tweedie distribution as

$$\begin{aligned} \lambda &= \frac{\mu^{2-p}}{\phi(2-p)} \\ \alpha &= \frac{2-p}{p-1} \\ \theta &= \phi(p-1)\mu^{p-1} \end{aligned}$$

The mean of a Tweedie distribution is positive for  $p > 1$ .

Two predefined versions of the Tweedie distribution are provided with the HPSEVERITY procedure. The first version, named TWEEDIE and defined for  $p > 1$ , has the natural parameterization with parameters  $\mu$ ,  $\phi$ , and  $p$ . The second version, named STWEEDIE and defined for  $1 < p < 2$ , is the version with a scale parameter. It corresponds to the compound Poisson-gamma distribution with gamma scale parameter  $\theta$ , Poisson mean parameter  $\lambda$ , and the index parameter  $p$ . The index parameter decides the shape parameter  $\alpha$  of the gamma distribution as

$$\alpha = \frac{2-p}{p-1}$$

The parameters  $\theta$  and  $\lambda$  of the STWEEDIE distribution are related to the parameters  $\mu$  and  $\phi$  of the TWEEDIE distribution as

$$\begin{aligned} \mu &= \lambda\theta\alpha \\ \phi &= \frac{(\lambda\theta\alpha)^{2-p}}{\lambda(2-p)} = \frac{\theta}{(p-1)(\lambda\theta\alpha)^{p-1}} \end{aligned}$$

You can fit either version when there are no regression variables. Each version has its own merits. If you fit the TWEEDIE version, you have the direct estimate of the overall mean of the distribution. If you are interested in the most practical range of the index parameter  $1 < p < 2$ , then you can fit the STWEEDIE version, which provides you direct estimates of the Poisson and gamma components that comprise the distribution (an estimate of the gamma shape parameter  $\alpha$  is easily obtained from the estimate of  $p$ ).

If you want to estimate the effect of exogenous (regression) variables on the distribution, then you must use the STWEEDIE version, because PROC HPSEVERITY requires a distribution to have a scale parameter in order to estimate regression effects. For more information, see the section “[Estimating Regression Effects](#)” on page 290. The gamma scale parameter  $\theta$  is the scale parameter of the STWEEDIE distribution. If you are interested in determining the effect of regression variables on the mean of the distribution, you can do so by first fitting the STWEEDIE distribution to determine the effect of the regression variables on the scale parameter  $\theta$ . Then, you can easily estimate how the mean of the distribution  $\mu$  is affected by the regression variables using the relationship  $\mu = c\theta$ , where  $c = \lambda\alpha = \lambda(2 - p)/(p - 1)$ . The estimates of the regression parameters remain the same, whereas the estimate of the intercept parameter is adjusted by the estimates of the  $\lambda$  and  $p$  parameters.

### Parameter Initialization for Predefined Distributions

The parameters are initialized by using the method of moments for all the distributions, except for the gamma and the Weibull distributions. For the gamma distribution, approximate maximum likelihood estimates are used. For the Weibull distribution, the method of percentile matching is used.

Given  $n$  observations of the severity value  $y_i$  ( $1 \leq i \leq n$ ), the estimate of  $k$ th raw moment is denoted by  $m_k$  and computed as

$$m_k = \frac{1}{n} \sum_{i=1}^n y_i^k$$

The 100 $p$ th percentile is denoted by  $\pi_p$  ( $0 \leq p \leq 1$ ). By definition,  $\pi_p$  satisfies

$$F(\pi_p-) \leq p \leq F(\pi_p)$$

where  $F(\pi_p-) = \lim_{h \downarrow 0} F(\pi_p - h)$ . PROC HPSEVERITY uses the following practical method of computing  $\pi_p$ . Let  $\hat{F}_n(y)$  denote the empirical distribution function (EDF) estimate at a severity value  $y$ . Let  $y_p^-$  and  $y_p^+$  denote two consecutive values in the ascending sequence of  $y$  values such that  $\hat{F}_n(y_p^-) < p$  and  $\hat{F}_n(y_p^+) \geq p$ . Then, the estimate  $\hat{\pi}_p$  is computed as

$$\hat{\pi}_p = y_p^- + \frac{p - \hat{F}_n(y_p^-)}{\hat{F}_n(y_p^+) - \hat{F}_n(y_p^-)}(y_p^+ - y_p^-)$$

Let  $\epsilon$  denote the smallest double-precision floating-point number such that  $1 + \epsilon > 1$ . This machine precision constant can be obtained by using the CONSTANT function in Base SAS software.

The details of how parameters are initialized for each predefined distribution are as follows:

**BURR** The parameters are initialized by using the method of moments. The  $k$ th raw moment of the Burr distribution is:

$$E[X^k] = \frac{\theta^k \Gamma(1 + k/\gamma) \Gamma(\alpha - k/\gamma)}{\Gamma(\alpha)}, \quad -\gamma < k < \alpha\gamma$$

Three moment equations  $E[X^k] = m_k$  ( $k = 1, 2, 3$ ) need to be solved for initializing the three parameters of the distribution. In order to get an approximate closed form solution, the second shape parameter  $\hat{\gamma}$  is initialized to a value of 2. If  $2m_3 - 3m_1m_2 > 0$ , then simplifying and solving the moment equations yields the following feasible set of initial values:

$$\hat{\theta} = \sqrt{\frac{m_2m_3}{2m_3 - 3m_1m_2}}, \quad \hat{\alpha} = 1 + \frac{m_3}{2m_3 - 3m_1m_2}, \quad \hat{\gamma} = 2$$

If  $2m_3 - 3m_1m_2 < \epsilon$ , then the parameters are initialized as follows:

$$\hat{\theta} = \sqrt{m_2}, \quad \hat{\alpha} = 2, \quad \hat{\gamma} = 2$$

EXP The parameters are initialized by using the method of moments. The  $k$ th raw moment of the exponential distribution is:

$$E[X^k] = \theta^k \Gamma(k + 1), \quad k > -1$$

Solving  $E[X] = m_1$  yields the initial value of  $\hat{\theta} = m_1$ .

GAMMA The parameter  $\alpha$  is initialized by using its *approximate* maximum likelihood (ML) estimate. For a set of  $n$  independent and identically distributed observations  $y_i$  ( $1 \leq i \leq n$ ) drawn from a gamma distribution, the log likelihood  $l$  is defined as follows:

$$\begin{aligned} l &= \sum_{i=1}^n \log \left( y_i^{\alpha-1} \frac{e^{-y_i/\theta}}{\theta^\alpha \Gamma(\alpha)} \right) \\ &= (\alpha - 1) \sum_{i=1}^n \log(y_i) - \frac{1}{\theta} \sum_{i=1}^n y_i - n\alpha \log(\theta) - n \log(\Gamma(\alpha)) \end{aligned}$$

Using a shorter notation of  $\sum$  to denote  $\sum_{i=1}^n$  and solving the equation  $\partial l / \partial \theta = 0$  yields the following ML estimate of  $\theta$ :

$$\hat{\theta} = \frac{\sum y_i}{n\alpha} = \frac{m_1}{\alpha}$$

Substituting this estimate in the expression of  $l$  and simplifying gives

$$l = (\alpha - 1) \sum \log(y_i) - n\alpha - n\alpha \log(m_1) + n\alpha \log(\alpha) - n \log(\Gamma(\alpha))$$

Let  $d$  be defined as follows:

$$d = \log(m_1) - \frac{1}{n} \sum \log(y_i)$$

Solving the equation  $\partial l / \partial \alpha = 0$  yields the following expression in terms of the digamma function,  $\psi(\alpha)$ :

$$\log(\alpha) - \psi(\alpha) = d$$

The digamma function can be approximated as follows:

$$\hat{\psi}(\alpha) \approx \log(\alpha) - \frac{1}{\alpha} \left( 0.5 + \frac{1}{12\alpha + 2} \right)$$

This approximation is within 1.4% of the true value for all the values of  $\alpha > 0$  except when  $\alpha$  is arbitrarily close to the positive root of the digamma function (which is approximately 1.461632). Even for the values of  $\alpha$  that are close to the positive root, the absolute error between true and approximate values is still acceptable ( $|\hat{\psi}(\alpha) - \psi(\alpha)| < 0.005$  for  $\alpha > 1.07$ ). Solving the equation that arises from this approximation yields the following estimate of  $\alpha$ :

$$\hat{\alpha} = \frac{3 - d + \sqrt{(d - 3)^2 + 24d}}{12d}$$

If this approximate ML estimate is infeasible, then the method of moments is used. The  $k$ th raw moment of the gamma distribution is:

$$E[X^k] = \theta^k \frac{\Gamma(\alpha + k)}{\Gamma(\alpha)}, \quad k > -\alpha$$

Solving  $E[X] = m_1$  and  $E[X^2] = m_2$  yields the following initial value for  $\alpha$ :

$$\hat{\alpha} = \frac{m_1^2}{m_2 - m_1^2}$$

If  $m_2 - m_1^2 < \epsilon$  (almost zero sample variance), then  $\alpha$  is initialized as follows:

$$\hat{\alpha} = 1$$

After computing the estimate of  $\alpha$ , the estimate of  $\theta$  is computed as follows:

$$\hat{\theta} = \frac{m_1}{\hat{\alpha}}$$

Both the maximum likelihood method and the method of moments arrive at the same relationship between  $\hat{\alpha}$  and  $\hat{\theta}$ .

GPD

The parameters are initialized by using the method of moments. Notice that for  $\xi > 0$ , the CDF of the generalized Pareto distribution (GPD) is:

$$\begin{aligned} F(x) &= 1 - \left(1 + \frac{\xi x}{\theta}\right)^{-1/\xi} \\ &= 1 - \left(\frac{\theta/\xi}{x + \theta/\xi}\right)^{1/\xi} \end{aligned}$$

This is equivalent to a Pareto distribution with scale parameter  $\theta_1 = \theta/\xi$  and shape parameter  $\alpha = 1/\xi$ . Using this relationship, the parameter initialization method used for the PARETO distribution is used to get the following initial values for the parameters of the GPD distribution:

$$\hat{\theta} = \frac{m_1 m_2}{2(m_2 - m_1^2)}, \quad \hat{\xi} = \frac{m_2 - 2m_1^2}{2(m_2 - m_1^2)}$$

If  $m_2 - m_1^2 < \epsilon$  (almost zero sample variance) or  $m_2 - 2m_1^2 < \epsilon$ , then the parameters are initialized as follows:

$$\hat{\theta} = \frac{m_1}{2}, \quad \hat{\xi} = \frac{1}{2}$$

IGAUSS The parameters are initialized by using the method of moments. The standard parameterization of the inverse Gaussian distribution (also known as the Wald distribution), in terms of the location parameter  $\mu$  and shape parameter  $\lambda$ , is as follows (Klugman, Panjer, and Willmot 1998, p. 583):

$$f(x) = \sqrt{\frac{\lambda}{2\pi x^3}} \exp\left(\frac{-\lambda(x - \mu)^2}{2\mu^2 x}\right)$$

$$F(x) = \Phi\left(\left(\frac{x}{\mu} - 1\right) \sqrt{\frac{\lambda}{x}}\right) + \Phi\left(-\left(\frac{x}{\mu} + 1\right) \sqrt{\frac{\lambda}{x}}\right) \exp\left(\frac{2\lambda}{\mu}\right)$$

For this parameterization, it is known that the mean is  $E[X] = \mu$  and the variance is  $Var[X] = \mu^3/\lambda$ , which yields the second raw moment as  $E[X^2] = \mu^2(1 + \mu/\lambda)$  (computed by using  $E[X^2] = Var[X] + (E[X])^2$ ).

The predefined IGAUSS distribution in PROC HPSEVERITY uses the following alternate parameterization to allow the distribution to have a scale parameter,  $\theta$ :

$$f(x) = \sqrt{\frac{\alpha\theta}{2\pi x^3}} \exp\left(\frac{-\alpha(x - \theta)^2}{2x\theta}\right)$$

$$F(x) = \Phi\left(\left(\frac{x}{\theta} - 1\right) \sqrt{\frac{\alpha\theta}{x}}\right) + \Phi\left(-\left(\frac{x}{\theta} + 1\right) \sqrt{\frac{\alpha\theta}{x}}\right) \exp(2\alpha)$$

The parameters  $\theta$  (scale) and  $\alpha$  (shape) of this alternate form are related to the parameters  $\mu$  and  $\lambda$  of the preceding form such that  $\theta = \mu$  and  $\alpha = \lambda/\mu$ . Using this relationship, the first and second raw moments of the IGAUSS distribution are:

$$E[X] = \theta$$

$$E[X^2] = \theta^2 \left(1 + \frac{1}{\alpha}\right)$$

Solving  $E[X] = m_1$  and  $E[X^2] = m_2$  yields the following initial values:

$$\hat{\theta} = m_1, \quad \hat{\alpha} = \frac{m_1^2}{m_2 - m_1^2}$$

If  $m_2 - m_1^2 < \epsilon$  (almost zero sample variance), then the parameters are initialized as follows:

$$\hat{\theta} = m_1, \quad \hat{\alpha} = 1$$



**LOGN** The parameters are initialized by using the method of moments. The  $k$ th raw moment of the lognormal distribution is:

$$E[X^k] = \exp\left(k\mu + \frac{k^2\sigma^2}{2}\right)$$

Solving  $E[X] = m_1$  and  $E[X^2] = m_2$  yields the following initial values:

$$\hat{\mu} = 2\log(m_1) - \frac{\log(m_2)}{2}, \quad \hat{\sigma} = \sqrt{\log(m_2) - 2\log(m_1)}$$

**PARETO** The parameters are initialized by using the method of moments. The  $k$ th raw moment of the Pareto distribution is:

$$E[X^k] = \frac{\theta^k \Gamma(k+1) \Gamma(\alpha-k)}{\Gamma(\alpha)}, \quad -1 < k < \alpha$$

Solving  $E[X] = m_1$  and  $E[X^2] = m_2$  yields the following initial values:

$$\hat{\theta} = \frac{m_1 m_2}{m_2 - 2m_1^2}, \quad \hat{\alpha} = \frac{2(m_2 - m_1^2)}{m_2 - 2m_1^2}$$

If  $m_2 - m_1^2 < \epsilon$  (almost zero sample variance) or  $m_2 - 2m_1^2 < \epsilon$ , then the parameters are initialized as follows:

$$\hat{\theta} = m_1, \quad \hat{\alpha} = 2$$

**TWEEDIE** The parameter  $p$  is initialized by assuming that the sample is generated from a gamma distribution with shape parameter  $\alpha$  and by computing  $\hat{p} = \frac{\hat{\alpha}+2}{\hat{\alpha}+1}$ . The initial value  $\hat{\alpha}$  is obtained from using the method previously described for the GAMMA distribution. The parameter  $\mu$  is the mean of the distribution. Hence, it is initialized to the sample mean as

$$\hat{\mu} = m_1$$

Variance of a Tweedie distribution is equal to  $\phi\mu^p$ . Thus, the sample variance is used to initialize the value of  $\phi$  as

$$\hat{\phi} = \frac{m_2 - m_1^2}{\hat{\mu}^{\hat{p}}}$$

**STWEEDIE** STWEEDIE is a compound Poisson-gamma mixture distribution with mean  $\mu = \lambda\theta\alpha$ , where  $\alpha$  is the shape parameter of the gamma random variables in the mixture and the parameter  $p$  is determined solely by  $\alpha$ . First, the parameter  $p$  is initialized by assuming that the sample is generated from a gamma distribution with shape parameter  $\alpha$  and by computing  $\hat{p} = \frac{\hat{\alpha}+2}{\hat{\alpha}+1}$ . The initial value  $\hat{\alpha}$  is obtained from using the method previously described for the GAMMA distribution. As done for initializing the parameters of the TWEEDIE distribution, the sample mean and variance are used to compute the values  $\hat{\mu}$  and  $\hat{\phi}$  as

$$\hat{\mu} = m_1$$

$$\hat{\phi} = \frac{m_2 - m_1^2}{\hat{\mu}^{\hat{p}}}$$

Based on the relationship between the parameters of TWEEDIE and STWEEDIE distributions described in the section “[Tweedie Distributions](#)” on page 276, values of  $\theta$  and  $\lambda$  are initialized as

$$\hat{\theta} = \hat{\phi}(\hat{p} - 1)\hat{\mu}^{p-1}$$

$$\hat{\lambda} = \frac{\hat{\mu}}{\hat{\theta}\hat{\alpha}}$$

**WEIBULL** The parameters are initialized by using the percentile matching method. Let  $q1$  and  $q3$  denote the estimates of the 25th and 75th percentiles, respectively. Using the formula for the CDF of Weibull distribution, they can be written as

$$1 - \exp(-(q1/\theta)^\tau) = 0.25$$

$$1 - \exp(-(q3/\theta)^\tau) = 0.75$$

Simplifying and solving these two equations yields the following initial values:

$$\hat{\theta} = \exp\left(\frac{r \log(q1) - \log(q3)}{r - 1}\right), \quad \hat{\tau} = \frac{\log(\log(4))}{\log(q3) - \log(\hat{\theta})}$$

where  $r = \log(\log(4))/\log(\log(4/3))$ . These initial values agree with those suggested in Klugman, Panjer, and Willmot (1998).

A summary of the initial values of all the parameters for all the predefined distributions is given in [Table 9.3](#). The table also provides the names of the parameters to use in the **INIT=** option in the DIST statement if you want to provide a different initial value.

**Table 9.3** Parameter Initialization for Predefined Distributions

Distribution	Parameter	Name for INIT option	Default Initial Value
BURR	$\theta$	theta	$\sqrt{\frac{m_2 m_3}{2m_3 - 3m_1 m_2}}$
	$\alpha$	alpha	$1 + \frac{m_3}{2m_3 - 3m_1 m_2}$
	$\gamma$	gamma	2
EXP	$\theta$	theta	$m_1$
GAMMA	$\theta$	theta	$m_1/\alpha$
	$\alpha$	alpha	$\frac{3-d+\sqrt{(d-3)^2+24d}}{12d}$
GPD	$\theta$	theta	$m_1 m_2 / (2(m_2 - m_1^2))$
	$\xi$	xi	$(m_2 - 2m_1^2) / (2(m_2 - m_1^2))$
IGAUSS	$\theta$	theta	$m_1$
	$\alpha$	alpha	$m_1^2 / (m_2 - m_1^2)$
LOGN	$\mu$	mu	$2 \log(m_1) - \log(m_2)/2$
	$\sigma$	sigma	$\sqrt{\log(m_2) - 2 \log(m_1)}$
PARETO	$\theta$	theta	$m_1 m_2 / (m_2 - 2m_1^2)$
	$\alpha$	alpha	$2(m_2 - m_1^2) / (m_2 - 2m_1^2)$
TWEEDIE	$\mu$	mu	$m_1$
	$\phi$	phi	$(m_2 - m_1^2) / m_1^p$
	$p$	p	$(\alpha + 2) / (\alpha + 1)$
			where $\alpha = \frac{3-d+\sqrt{(d-3)^2+24d}}{12d}$
STWEEDIE	$\theta$	theta	$(m_2 - m_1^2)(p - 1) / m_1$
	$\lambda$	lambda	$m_1^2 / (\alpha(m_2 - m_1^2)(p - 1))$
	$p$	p	$(\alpha + 2) / (\alpha + 1)$
			where $\alpha = \frac{3-d+\sqrt{(d-3)^2+24d}}{12d}$
WEIBULL	$\theta$	theta	$\exp\left(\frac{r \log(q1) - \log(q3)}{r-1}\right)$
	$\tau$	tau	$\log(\log(4)) / (\log(q3) - \log(\hat{\theta}))$

Notes:

- $m_k$  denotes the  $k$ th raw moment
- $d = \log(m_1) - (\sum \log(y_i)) / n$
- $q1$  and  $q3$  denote the 25th and 75th percentiles, respectively
- $r = \log(\log(4)) / \log(\log(4/3))$

## Censoring and Truncation

One of the key features of PROC HPSEVERITY is that it enables you to specify whether the severity event's magnitude is observable and if it is observable, then whether the exact value of the magnitude is known. If an event is unobservable when the magnitude is in certain intervals, then it is referred to as a truncation effect. If the exact magnitude of the event is not known, but it is known to have a value in a certain interval, then it is referred to as a censoring effect.

PROC HPSEVERITY allows a severity event to be subject to any combination of the following four censoring and truncation effects:

- **Left-truncation:** An event is said to be left-truncated if it is observed only when  $Y > T^l$ , where  $Y$  denotes the random variable for the magnitude and  $T^l$  denotes a random variable for the truncation threshold. You can specify left-truncation using the `LEFTTRUNCATED=` option in the LOSS statement.
- **Right-truncation:** An event is said to be right-truncated if it is observed only when  $Y \leq T^r$ , where  $Y$  denotes the random variable for the magnitude and  $T^r$  denotes a random variable for the truncation threshold. You can specify right-truncation using the `RIGHTTRUNCATED=` option in the LOSS statement.
- **Left-censoring:** An event is said to be left-censored if it is known that the magnitude is  $Y \leq C^l$ , but the exact value of  $Y$  is not known.  $C^l$  is a random variable for the censoring limit. You can specify left-censoring using the `LEFTCENSORED=` option in the LOSS statement.
- **Right-censoring:** An event is said to be right-censored if it is known that the magnitude is  $Y > C^r$ , but the exact value of  $Y$  is not known.  $C^r$  is a random variable for the censoring limit. You can specify right-censoring using the `RIGHTCENSORED=` option in the LOSS statement.

For each effect, you can specify a different threshold or limit for each observation or specify a single threshold or limit that applies to all the observations.

If all the four types of effects are present on an event, then the following relationship holds:  $T^l < C^r \leq C^l \leq T^r$ . PROC HPSEVERITY checks these relationships and write a warning to the SAS log if any is violated.

If you specify the response variable in the LOSS statement, then PROC HPSEVERITY also checks whether each observation satisfies the definitions of the specified censoring and truncation effects. If you specify left-truncation, then PROC HPSEVERITY ignores observations where  $Y \leq T^l$ , because such observations are not observable by definition. Similarly, if you specify right-truncation, then PROC HPSEVERITY ignores observations where  $Y > T^r$ . If you specify left-censoring, then PROC HPSEVERITY treats an observation with  $Y > C^l$  as uncensored and ignores the value of  $C^l$ . The observations with  $Y \leq C^l$  are considered as left-censored, and the value of  $Y$  is ignored. If you specify right-censoring, then PROC HPSEVERITY treats an observation with  $Y \leq C^r$  as uncensored and ignores the value of  $C^r$ . The observations with  $Y > C^r$  are considered as right-censored, and the value of  $Y$  is ignored. If you specify both left-censoring and right-censoring, it is referred to as interval-censoring. If  $C^r < C^l$  is satisfied for an observation, then it is considered as interval-censored and the value of the response variable is ignored. If  $C^r = C^l$  for an observation, then PROC HPSEVERITY assumes that observation to be uncensored. If all the observations in a data set are censored in some form, then the specification of the response variable in the LOSS statement is optional, because the actual value of the response variable is not required for the purposes of estimating a model.

Specification of censoring and truncation affects the likelihood of the data (see the section “[Likelihood Function](#)” on page 286) and how the empirical distribution function (EDF) is estimated (see the section “[Empirical Distribution Function Estimation Methods](#)” on page 295).

### Truncation and Conditional CDF Estimates

If you specify left-truncation or right-truncation, then the EDF estimates that are computed by all methods except the STANDARD method are conditional on the truncation information. See the section “[EDF Estimates and Truncation](#)” on page 300 for more information. In such cases, PROC HPSEVERITY uses conditional estimates of the CDF when it computes the EDF-based statistics of fit.

Let  $t_{\min}^l = \min_i \{t_i^l\}$  be the smallest value of the left-truncation threshold ( $t_i^l$  is the left-truncation threshold for observation  $i$ ) and  $t_{\max}^r = \max_i \{t_i^r\}$  be the largest value of the right-truncation threshold ( $t_i^r$  is the right-truncation threshold for observation  $i$ ). If  $\hat{F}(y)$  denotes the unconditional estimate of the CDF at  $y$ , then the conditional estimate  $\hat{F}^c(y)$  is computed as follows:

- If an observation is both left-truncated and right-truncated, then

$$\hat{F}^c(y) = \frac{\hat{F}(y) - \hat{F}(t_{\min}^l)}{\hat{F}(t_{\max}^r) - \hat{F}(t_{\min}^l)}$$

- If an observation is left-truncated but not right-truncated, then

$$\hat{F}^c(y) = \frac{\hat{F}(y) - \hat{F}(t_{\min}^l)}{1 - \hat{F}(t_{\min}^l)}$$

- If an observation is right-truncated but not left-truncated, then

$$\hat{F}^c(y) = \frac{\hat{F}(y)}{\hat{F}(t_{\max}^r)}$$

If you specify regressors, then  $\hat{F}(y)$ ,  $\hat{F}(t_{\min}^l)$ , and  $\hat{F}(t_{\max}^r)$  are all computed from a mixture distribution, as described in the section “[CDF Estimates with Regression Effects](#)” on page 293.

---

## Parameter Estimation Method

If you do not specify a custom objective function by specifying programming statements and the **OBJECTIVE=** option in the PROC HPSEVERITY statement, then PROC HPSEVERITY uses the maximum likelihood (ML) method to estimate the parameters of each model. A nonlinear optimization process is used to maximize the log of the likelihood function. If you specify a custom objective function, then PROC HPSEVERITY uses a nonlinear optimization algorithm to estimate the parameters of each model that minimize the value of your specified objective function. For more information, see the section “[Custom Objective Functions](#)” on page 332.

## Likelihood Function

Let  $f_{\Theta}(x)$  and  $F_{\Theta}(x)$  denote the PDF and CDF, respectively, evaluated at  $x$  for a set of parameter values  $\Theta$ . Let  $Y$  denote the random response variable, and let  $y$  denote its value recorded in an observation in the input data set. Let  $T^l$  and  $T^r$  denote the random variables for the left-truncation and right-truncation threshold, respectively, and let  $t^l$  and  $t^r$  denote their values for an observation, respectively. If there is no left-truncation, then  $t^l = \tau^l$ , where  $\tau^l$  is the smallest value in the support of the distribution; so  $F(t^l) = 0$ . If there is no right-truncation, then  $t^r = \tau_h$ , where  $\tau_h$  is the largest value in the support of the distribution; so  $F(t^r) = 1$ . Let  $C^l$  and  $C^r$  denote the random variables for the left-censoring and right-censoring limit, respectively, and let  $c^l$  and  $c^r$  denote their values for an observation, respectively. If there is no left-censoring, then  $c^l = \tau_h$ ; so  $F(c^l) = 1$ . If there is no right-censoring, then  $c^r = \tau^l$ ; so  $F(c^r) = 0$ .

The set of input observations can be categorized into the following four subsets within each BY group:

- $E$  is the set of uncensored and untruncated observations. The likelihood of an observation in  $E$  is

$$l_E = \Pr(Y = y) = f_{\Theta}(y)$$

- $E_t$  is the set of uncensored observations that are truncated. The likelihood of an observation in  $E_t$  is

$$l_{E_t} = \Pr(Y = y | t^l < Y \leq t^r) = \frac{f_{\Theta}(y)}{F_{\Theta}(t^r) - F_{\Theta}(t^l)}$$

- $C$  is the set of censored observations that are not truncated. The likelihood of an observation  $C$  is

$$l_C = \Pr(c^r < Y \leq c^l) = F_{\Theta}(c^l) - F_{\Theta}(c^r)$$

- $C_t$  is the set of censored observations that are truncated. The likelihood of an observation  $C_t$  is

$$l_{C_t} = \Pr(c^r < Y \leq c^l | t^l < Y \leq t^r) = \frac{F_{\Theta}(c^l) - F_{\Theta}(c^r)}{F_{\Theta}(t^r) - F_{\Theta}(t^l)}$$

Note that  $(E \cup E_t) \cap (C \cup C_t) = \emptyset$ . Also, the sets  $E_t$  and  $C_t$  are empty when you do not specify truncation, and the sets  $C$  and  $C_t$  are empty when you do not specify censoring.

Given this, the likelihood of the data  $L$  is as follows:

$$L = \prod_E f_{\Theta}(y) \prod_{E_t} \frac{f_{\Theta}(y)}{F_{\Theta}(t^r) - F_{\Theta}(t^l)} \prod_C F_{\Theta}(c^l) - F_{\Theta}(c^r) \prod_{C_t} \frac{F_{\Theta}(c^l) - F_{\Theta}(c^r)}{F_{\Theta}(t^r) - F_{\Theta}(t^l)}$$

The maximum likelihood procedure used by PROC HPSEVERITY finds an optimal set of parameter values  $\hat{\Theta}$  that maximizes  $\log(L)$  subject to the boundary constraints on parameter values. For a distribution *dist*, you can specify such boundary constraints by using the *dist\_LOWERBOUNDS* and *dist\_UPPERBOUNDS* subroutines. For more information, see the section “Defining a Severity Distribution Model with the FCMP Procedure” on page 308. Some aspects of the optimization process can be controlled by using the NLOPTIONS statement.

## Estimating Covariance and Standard Errors

PROC HPSEVERITY computes an estimate of the covariance matrix of the parameters by using the asymptotic theory of the maximum likelihood estimators (MLE). If  $N$  denotes the number of observations used for estimating a parameter vector  $\theta$ , then the theory states that as  $N \rightarrow \infty$ , the distribution of  $\hat{\theta}$ , the estimate of  $\theta$ , converges to a normal distribution with mean  $\theta$  and covariance  $\hat{C}$  such that  $\mathbf{I}(\theta) \cdot \hat{C} \rightarrow 1$ , where  $\mathbf{I}(\theta) = -E [\nabla^2 \log(L(\theta))]$  is the information matrix for the likelihood of the data,  $L(\theta)$ . The covariance estimate is obtained by using the inverse of the information matrix.

In particular, if  $\mathbf{G} = \nabla^2 \log(-L(\theta))$  denotes the Hessian matrix of the negative of log likelihood, then the covariance estimate is computed as

$$\hat{C} = \frac{N}{d} \mathbf{G}^{-1}$$

where  $d$  is a denominator that is determined by the **VARDEF=** option. If **VARDEF=N**, then  $d = N$ , which yields the asymptotic covariance estimate. If **VARDEF=DF**, then  $d = N - k$ , where  $k$  is number of parameters (the model's degrees of freedom). The **VARDEF=DF** option is the default, because it attempts to correct the potential bias introduced by the finite sample.

The standard error  $s_i$  of the parameter  $\theta_i$  is computed as the square root of the  $i$ th diagonal element of the estimated covariance matrix; that is,  $s_i = \sqrt{\hat{C}_{ii}}$ .

If you specify a custom objective function, then the covariance matrix of the parameters is still computed by inverting the information matrix, except that the Hessian matrix  $\mathbf{G}$  is computed as  $\mathbf{G} = \nabla^2 \log(U(\theta))$ , where  $U$  denotes your custom objective function that is minimized by the optimizer.

Covariance and standard error estimates might not be available if the Hessian matrix is found to be singular at the end of the optimization process. This can especially happen if the optimization process stops without converging.

---

## Details of Optimization Algorithms

### Overview

There are several optimization techniques available. You can choose a particular optimizer with the **TECH=name** option in the PROC statement or **NLOPTIONS** statement.

**Table 9.4** Optimization Techniques

Algorithm	TECH=
Trust region Method	TRUREG
Newton-Raphson method with line search	NEWRAP
Newton-Raphson method with ridging	NRRIDG
Quasi-Newton methods (DBFGS, DDFP, BFGS, DFP)	QUANEW
Double-dogleg method (DBFGS, DDFP)	DBLDOG
Conjugate gradient methods (PB, FR, PR, CD)	CONGRA
Nelder-Mead simplex method	NMSIMP

No algorithm for optimizing general nonlinear functions exists that always finds the global optimum for a general nonlinear minimization problem in a reasonable amount of time. Since no single optimization

technique is invariably superior to others, NLO provides a variety of optimization techniques that work well in various circumstances. However, you can devise problems for which none of the techniques in NLO can find the correct solution. Moreover, nonlinear optimization can be computationally expensive in terms of time and memory, so you must be careful when matching an algorithm to a problem.

All optimization techniques in NLO use  $O(n^2)$  memory except the conjugate gradient methods, which use only  $O(n)$  of memory and are designed to optimize problems with many parameters. These iterative techniques require repeated computation of the following:

- the function value (optimization criterion)
- the gradient vector (first-order partial derivatives)
- for some techniques, the (approximate) Hessian matrix (second-order partial derivatives)

However, since each of the optimizers requires different derivatives, some computational efficiencies can be gained. Table 9.5 shows which derivatives are required for each optimization technique. (FOD means that first-order derivatives or the gradient is computed; SOD means that second-order derivatives or the Hessian is computed.)

**Table 9.5** Optimization Computations

Algorithm	FOD	SOD
TRUREG	x	x
NEWRAP	x	x
NRRIDG	x	x
QUANEW	x	-
DBLDOG	x	-
CONGRA	x	-
NMSIMP	-	-

Each optimization method uses one or more convergence criteria that determine when it has converged. The various termination criteria are listed and described in the previous section. An algorithm is considered to have converged when any one of the convergence criteria is satisfied. For example, under the default settings, the QUANEW algorithm converges if  $\text{ABSGCONV} < 1\text{E-}5$ ,  $\text{FCONV} < 10^{-\text{FDIGITS}}$ , or  $\text{GCONV} < 1\text{E-}8$ .

## Choosing an Optimization Algorithm

The factors for choosing a particular optimization technique for a particular problem are complex and might involve trial and error.

For many optimization problems, computing the gradient takes more computer time than computing the function value, and computing the Hessian sometimes takes *much* more computer time and memory than computing the gradient, especially when there are many decision variables. Unfortunately, optimization techniques that do not use some kind of Hessian approximation usually require many more iterations than techniques that do use a Hessian matrix. As a result the total run time of these techniques is often longer. Techniques that do not use the Hessian also tend to be less reliable. For example, they can more easily terminate at stationary points rather than at global optima.



A few general remarks about the various optimization techniques follow:

- The second-derivative methods TRUREG, NEWRAP, and NRRIDG are best for small problems where the Hessian matrix is not expensive to compute. Sometimes the NRRIDG algorithm can be faster than the TRUREG algorithm, but TRUREG can be more stable. The NRRIDG algorithm requires only one matrix with  $n(n + 1)/2$  double words; TRUREG and NEWRAP require two such matrices.
- The first-derivative methods QUANEW and DBLDOG are best for medium-sized problems where the objective function and the gradient are much faster to evaluate than the Hessian. The QUANEW and DBLDOG algorithms generally require more iterations than TRUREG, NRRIDG, and NEWRAP, but each iteration can be much faster. The QUANEW and DBLDOG algorithms require only the gradient to update an approximate Hessian, and they require slightly less memory than TRUREG or NEWRAP (essentially one matrix with  $n(n + 1)/2$  double words). QUANEW is the default optimization method.
- The first-derivative method CONGRA is best for large problems where the objective function and the gradient can be computed much faster than the Hessian and where too much memory is required to store the (approximate) Hessian. The CONGRA algorithm generally requires more iterations than QUANEW or DBLDOG, but each iteration can be much faster. Since CONGRA requires only a factor of  $n$  double-word memory, many large applications can be solved only by CONGRA.
- The no-derivative method NMSIMP is best for small problems where derivatives are not continuous or are very difficult to compute.

---

## Parameter Initialization

PROC HPSEVERITY enables you to initialize parameters of a model in different ways. A model can have two kinds of parameters: distribution parameters and regression parameters.

The distribution parameters can be initialized by using one of the following three methods:

INIT= option	You can use the <a href="#">INIT=</a> option in the DIST statement.
INEST= data set	You can use the <a href="#">INEST=</a> data set.
PARMINIT subroutine	You can define a <a href="#">dist_PARMINIT</a> subroutine in the distribution model. For more information, see the section “ <a href="#">Defining a Severity Distribution Model with the FCMP Procedure</a> ” on page 308.

Note that only one of the initialization methods is used. You cannot combine them. They are used in the following order:

- The method that uses the [INIT=](#) option takes the highest precedence. If you use the [INIT=](#) option to provide an initial value for at least one parameter, then other initialization methods ([INEST=](#) and [PARMINIT](#)) are not used. If you specify initial values for some but not all the parameters by using the [INIT=](#) option, then the uninitialized parameters are initialized to the default value of 0.001.

If you use this option and if you specify the regression effects, then the value of the first distribution parameter must be related to the initial value for the *base* value of the scale or log-transformed scale parameter. For more information, see the section “[Estimating Regression Effects](#)” on page 290.

- The method that uses the INEST= data set takes the second precedence. If you specify a nonmissing value for even one distribution parameter, then the PARMINIT method is not used and any uninitialized parameters are initialized to the default value of 0.001.
- If none of the distribution parameters are initialized by using the INIT= option or the INEST= data set, but the distribution model defines a PARMINIT subroutine, then PROC HPSEVERITY invokes that subroutine with appropriate inputs to initialize the parameters. If the PARMINIT subroutine returns missing values for some parameters, then those parameters are initialized to the default value of 0.001.
- If none of the initialization methods are used, each distribution parameter is initialized to the default value of 0.001.

For more information about regression models and initialization of regression parameters, see the section “[Estimating Regression Effects](#)” on page 290.

### PARMINIT-Based Parameter Initialization Method and Distributed Data

If you specify a distributed mode of execution for the procedure, then the input data are distributed across the computational nodes. For more information about the distributed computing model, see the section “[Distributed and Multithreaded Computation](#)” on page 305. If the PARMINIT subroutine is used for initializing the distribution parameters, then PROC HPSEVERITY invokes that subroutine on each computational node with the data that are local to that node. The EDF estimates that are supplied to the PARMINIT subroutine are also computed using the local data. The initial values of the parameters that are supplied to the optimizer are the average of the local estimates that are computed on each node. This approach works well if the data are distributed randomly across nodes. If you distribute the data on the appliance before you run the procedure (alongside-the-database model), then you should try to make the distribution as random as possible in order to increase the chances of computing good initial values. If you specify a data set that is not distributed before you run the procedure, then PROC HPSEVERITY distributes the data for you by sending the first observation to the first node, the second observation to the second node, and so on. If the order of observations is random, then this method ensures random distribution of data across the computational nodes.

---

## Estimating Regression Effects

The HPSEVERITY procedure enables you to estimate the effects of regressor (exogenous) variables while fitting a distribution if the distribution has a scale parameter or a log-transformed scale parameter.

Let  $x_j$ ,  $j = 1, \dots, k$ , denote the  $k$  regressor variables. Let  $\beta_j$  denote the regression parameter that corresponds to the regressor  $x_j$ . If you do not specify regression effects, then the model for the response variable  $Y$  is of the form

$$Y \sim \mathcal{F}(\Theta)$$

where  $\mathcal{F}$  is the distribution of  $Y$  with parameters  $\Theta$ . This model is usually referred to as the error model. The regression effects are modeled by extending the error model to the following form:

$$Y \sim \exp\left(\sum_{j=1}^k \beta_j x_j\right) \cdot \mathcal{F}(\Theta)$$

Under this model, the distribution of  $Y$  is valid and belongs to the same parametric family as  $\mathcal{F}$  if and only if  $\mathcal{F}$  has a scale parameter. Let  $\theta$  denote the scale parameter and  $\Omega$  denote the set of nonscale distribution parameters of  $\mathcal{F}$ . Then the model can be rewritten as

$$Y \sim \mathcal{F}(\theta, \Omega)$$

such that  $\theta$  is modeled by the regressors as

$$\theta = \theta_0 \cdot \exp\left(\sum_{j=1}^k \beta_j x_j\right)$$

where  $\theta_0$  is the *base* value of the scale parameter. Thus, the scale regression model consists of the following parameters:  $\theta_0$ ,  $\Omega$ , and  $\beta_j$  ( $j = 1, \dots, k$ ).

Given this form of the model, distributions without a scale parameter cannot be considered when regression effects are to be modeled. If a distribution does not have a direct scale parameter, then PROC HPSEVERITY accepts it only if it has a log-transformed scale parameter — that is, if it has a parameter  $p = \log(\theta)$ .

### Offset Variable

You can specify that an offset variable be included in the scale regression model by specifying it in the **OFFSET=** option of the SCALEMODEL statement. The offset variable is a regressor whose regression coefficient is known to be 1. If  $x_o$  denotes the offset variable, then the scale regression model becomes

$$\theta = \theta_0 \cdot \exp(x_o + \sum_{j=1}^k \beta_j x_j)$$

The regression coefficient of the offset variable is fixed at 1 and not estimated, so it is not reported in the ParameterEstimates ODS table. However, if you specify the OUTEST= data set, then the regression coefficient is added as a variable to that data set. The value of the offset variable in OUTEST= data set is equal to 1 for the estimates row (\_TYPE\_='EST') and is equal to a special missing value (.F) for the standard error (\_TYPE\_='STDERR') and covariance (\_TYPE\_='COV') rows.

An offset variable is useful to model the scale parameter per unit of some measure of exposure. For example, in the automobile insurance context, measure of exposure can be the number of car-years insured or the total number of miles driven by a fleet of cars at a rental car company. For worker's compensation insurance, if you want to model the expected loss per enterprise, then you can use the number of employees or total employee salary as the measure of exposure. For epidemiological data, measure of exposure can be the number of people who are exposed to a certain pathogen when you are modeling the loss associated with an epidemic. In general, if  $e$  denotes the value of the exposure measure and if you specify  $x_o = \log(e)$  as the offset variable, then you are modeling the effect of other regressors ( $x_j$ ) on the size of the scale of the distribution *per unit of exposure*.

Another use for an offset variable is when you have a priori knowledge of the effect of some exogenous variables that cannot be included in the SCALEMODEL statement. You can model the combined effect of such variables as an offset variable in order to correct for the omitted variable bias.

## Parameter Initialization for Regression Models

The regression parameters are initialized either by using the values that you specify or by the default method.

- If you provide initial values for the regression parameters, then you must provide valid, nonmissing initial values for  $\theta_0$  and  $\beta_j$  parameters for all  $j$ .

You can specify the initial value for  $\theta_0$  by using either the INEST= data set or the INIT= option in the DIST statement. If the distribution has a direct scale parameter (no transformation), then the initial value for the first parameter of the distribution is used as an initial value for  $\theta_0$ . If the distribution has a log-transformed scale parameter, then the initial value for the first parameter of the distribution is used as an initial value for  $\log(\theta_0)$ .

You can use only the INEST= data set to specify the initial values for  $\beta_j$ . The INEST= data set must contain nonmissing initial values for all the regressors that you specify in the SCALEMODEL statement. The only missing value allowed is the special missing value .R, which indicates that the regressor is linearly dependent on other regressors. If you specify .R for a regressor for one distribution in a BY group, you must specify it so for all the distributions in that BY group.

- If you do not specify valid initial values for  $\theta_0$  or  $\beta_j$  parameters for all  $j$ , then PROC HPSEVERITY initializes those parameters by using the following method:

Let a random variable  $Y$  be distributed as  $\mathcal{F}(\theta, \Omega)$ , where  $\theta$  is the scale parameter. By the definition of the scale parameter, a random variable  $W = Y/\theta$  is distributed as  $\mathcal{G}(\Omega)$  such that  $\mathcal{G}(\Omega) = \mathcal{F}(1, \Omega)$ . Given a random error term  $e$  that is generated from a distribution  $\mathcal{G}(\Omega)$ , a value  $y$  from the distribution of  $Y$  can be generated as

$$y = \theta \cdot e$$

Taking the logarithm of both sides and using the relationship of  $\theta$  with the regressors yields:

$$\log(y) = \log(\theta_0) + \sum_{j=1}^k \beta_j x_j + \log(e)$$

PROC HPSEVERITY makes use of the preceding relationship to initialize parameters of a regression model with distribution *dist* as follows:

1. The following linear regression problem is solved to obtain initial estimates of  $\beta_0$  and  $\beta_j$ :

$$\log(y) = \beta_0 + \sum_{j=1}^k \beta_j x_j$$

The estimates of  $\beta_j$  ( $j = 1, \dots, k$ ) in the solution of this regression problem are used to initialize the respective regression parameters of the model. The estimate of  $\beta_0$  is later used to initialize the value of  $\theta_0$ .

The results of this regression are also used to detect whether any regressors are linearly dependent on the other regressors. If any such regressors are found, then a warning is written to the SAS log and the corresponding regressor is eliminated from further analysis. The estimates for linearly dependent regressors are denoted by a special missing value of .R in the OUTEST= data set and in any displayed output.

2. Let  $s_0$  denote the initial value of the scale parameter.

If the distribution model of *dist* does not contain the *dist\_PARMINIT* subroutine, then  $s_0$  and all the nonscale distribution parameters are initialized to the default value of 0.001.

However, it is strongly recommended that each distribution's model contain the *dist\_PARMINIT* subroutine. For more information, see the section “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308. If that subroutine is defined, then  $s_0$  is initialized as follows: Each input value  $y_i$  of the response variable is transformed to its scale-normalized version  $w_i$  as

$$w_i = \frac{y_i}{\exp(\beta_0 + \sum_{j=1}^k \beta_j x_{ij})}$$

where  $x_{ij}$  denotes the value of  $j$ th regressor in the  $i$ th input observation. These  $w_i$  values are used to compute the input arguments for the *dist\_PARMINIT* subroutine. The values that are computed by the subroutine for nonscale parameters are used as their respective initial values. If the distribution has an untransformed scale parameter, then  $s_0$  is set to the value of the scale parameter that is computed by the subroutine. If the distribution has a log-transformed scale parameter  $P$ , then  $s_0$  is computed as  $s_0 = \exp(l_0)$ , where  $l_0$  is the value of  $P$  computed by the subroutine.

3. The value of  $\theta_0$  is initialized as

$$\theta_0 = s_0 \cdot \exp(\beta_0)$$

## Reporting Estimates of Regression Parameters

When you request estimates to be written to the output (either ODS displayed output or in the OUTEST= data set), the estimate of the base value of the first distribution parameter is reported. If the first parameter is the log-transformed scale parameter, then the estimate of  $\log(\theta_0)$  is reported; otherwise, the estimate of  $\theta_0$  is reported. The transform of the first parameter of a distribution *dist* is controlled by the *dist\_SCALETRANSFORM* function that is defined for it.

## CDF Estimates with Regression Effects

When regression effects are estimated, the estimate of the scale parameter depends on the values of the regressors and the estimates of the regression parameters. This dependency results in a potentially different distribution for each observation. PROC HPSEVERITY needs to compute the estimates of the cumulative distribution function (CDF) to compute the EDF-based statistics of fit that compare the nonparametric and parametric estimates of the distribution function. To make the CDF estimates comparable across distributions and comparable to the empirical distribution function (EDF), PROC HPSEVERITY computes the CDF estimates from a representative distribution. The *representative distribution* is a mixture of a certain number of distributions, where each distribution differs only in the value of the scale parameter. You can specify the number of distributions in the mixture and how their scale values are chosen by using the *DFMIXTURE=* option in the SCALEMODEL statement.

Let  $N$  denote the number of observations used for EDF estimation,  $K$  denote the number of components in the mixture distribution,  $s_k$  denote the scale parameter of the  $k$ th mixture component, and  $d_k$  denote the weight associated with  $k$ th mixture component.

Let  $F(y; s_k, \hat{\Omega})$  denote the CDF of the  $k$ th component distribution, where  $\hat{\Omega}$  denotes the set of estimates of all parameters of the distribution other than the scale parameter. Then, the CDF estimates,  $F^*(y)$ , of the mixture distribution at  $y$  are computed as

$$F^*(y) = \frac{1}{D} \sum_{k=1}^K d_k F(y; s_k, \hat{\Omega})$$

where  $D$  is the normalization factor ( $D = \sum_{k=1}^K d_k$ ). The  $F^*(y)$  values are used for computing the EDF-based statistics of fit.

The scale values  $s_k$  for the  $K$  mixture components are derived from the set  $\{\hat{\theta}_i\}$  ( $i = 1 \dots N$ ) of  $N$  scale values, where  $\hat{\theta}_i$  denotes the estimate of the scale parameter due to observation  $i$ . It is computed as

$$\hat{\theta}_i = \hat{\theta}_0 \cdot \exp\left(\sum_{j=1}^k \hat{\beta}_j x_{ij}\right)$$

where  $\hat{\theta}_0$  is an estimate of the base value of the scale parameter,  $\hat{\beta}_j$  are the estimates of regression coefficients, and  $x_{ij}$  is the value of regressor  $j$  in observation  $i$ .

Let  $w_i$  denote the weight of observation  $i$ . If you specify the WEIGHT statement, then the weight is equal to the value of the specified weight variable for the corresponding observation in the DATA= data set; otherwise, the weight is set to 1.

You can specify one of the following *method-names* in the **DFMIXTURE=** option in the SCALEMODEL statement to specify the method of choosing  $K$  and the corresponding  $s_k$  and  $d_k$  values:

**FULL** In this method, there are as many mixture components as the number of observations that are used for estimation. In other words,  $K = N$ ,  $s_k = \hat{\theta}_k$ , and  $d_k = w_k$  ( $k = 1 \dots N$ ). This is the slowest method, because it requires  $O(N)$  computations to compute the mixture CDF  $F^*(y_i)$  of one observation. For  $N$  observations, the computational complexity in terms of number of CDF evaluations is  $O(N^2)$ . Even for moderately large values of  $N$ , the time taken to compute the mixture CDF can significantly exceed the time taken to estimate the model parameters. So, it is recommended that you use the FULL method only for small data sets.

**MEAN** In this method, the mixture contains only one distribution, whose scale value is the mean of the scale values that are implied by all the observations. In other words,  $s_1$  is computed as

$$s_1 = \frac{1}{W} \sum_{i=1}^N w_i \hat{\theta}_i$$

where  $W$  is the total weight ( $W = \sum_{i=1}^N w_i$ ).

This method is the fastest because it requires only one CDF evaluation per observation. The computational complexity is  $O(N)$  for  $N$  observations.

If you do not specify the **DFMIXTURE=** option in the SCALEMODEL statement, then this is the default method.

- QUANTILE** In this method, a certain number of quantiles are chosen from the set of all scale values. If you specify a value of  $q$  for the **K=** option when specifying this method, then  $K = q - 1$  and  $s_k$  are set to be the  $(q - 1)$   $q$ -quantiles from the set  $\{\hat{\theta}_i\}$  ( $i = 1 \dots N$ ). The weight of each of the components ( $d_k$ ) is assumed to be 1 for this method.
- The default value of  $q$  is 2, which implies a one-point mixture that has a distribution whose scale value is equal to the median scale value.
- For this method, PROC HPSEVERITY needs to sort the  $N$  scale values in the set  $\{\hat{\theta}_i\}$ ; the sorting requires  $O(N \log(N))$  computations. Then, computing the mixture estimate of one observation requires  $(q - 1)$  CDF evaluations. Hence, the computational complexity of this method is  $O(qN) + O(N \log(N))$  for computing a mixture CDF of  $N$  observations. For  $q \ll N$ , the QUANTILE method is significantly faster than the FULL method.
- RANDOM** In this method, a uniform random sample of observations is chosen, and the mixture contains the distributions that are implied by those observations. If you specify a value of  $r$  for the **K=** option when specifying this method, then the size of the sample is  $r$ . Hence,  $K = r$ . If  $l_j$  denotes the index of  $j$ th observation in the sample ( $j = 1 \dots r$ ), such that  $1 \leq l_j \leq N$ , then the scale of  $k$ th component distribution in the mixture is  $s_k = \hat{\theta}_{l_k}$  and the weight associated with it is  $d_k = w_{l_k}$ .
- You can also specify the seed to be used for generating the random sample by using the **SEED=** option for this method. The same sample of observations is used for all models.
- Computing a mixture estimate of one observation requires  $r$  CDF evaluations. Hence, the computational complexity of this method is  $O(rN)$  for computing a mixture CDF of  $N$  observations. For  $r \ll N$ , the RANDOM method is significantly faster than the FULL method.

---

## Empirical Distribution Function Estimation Methods

The empirical distribution function (EDF) is a nonparametric estimate of the cumulative distribution function (CDF) of the distribution. PROC HPSEVERITY computes EDF estimates for two purposes: to send the estimates to a distribution's **PARMINIT** subroutine in order to initialize the distribution parameters, and to compute the EDF-based statistics of fit.

To reduce the time that it takes to compute the EDF estimates, you can use the **INITSAMPLE** option to specify that only a fraction of the input data be used. If you do not specify the **INITSAMPLE** option and the data set has more than 10,000 valid observations, then a uniform random sample of at most 10,000 observations is used for EDF estimation.

In the distributed mode of execution, in which data are distributed across the grid nodes, the EDF estimates are computed on each node by using the portion of the input data that is located on that node. These local EDF estimates are an approximation of the global EDF estimates, which would be computed by using the entire input data set. PROC HPSEVERITY does not compute global EDF estimates. Let  $X$  denote a quantity that depends on the EDF estimates.  $X$  can be either an EDF-based initial value of a distribution parameter or an EDF-based statistic of fit. PROC HPSEVERITY estimates  $X$  as follows: First, each grid node  $k$  computes an estimate  $X_k$  by using the local EDF estimates that are computed on that node. Then, the estimate  $\hat{X}$  of  $X$  is computed as an average of all the  $X_k$  values; that is,  $\hat{X} = \sum_{i=1}^K X_k$ , where  $K$  denotes the total number of nodes where the data reside.

This section describes the methods that are used for computing EDF estimates.



## Notation

Let there be a set of  $N$  observations, each containing a quintuplet of values  $(y_i, t_i^l, t_i^r, c_i^r, c_i^l)$ ,  $i = 1, \dots, N$ , where  $y_i$  is the value of the response variable,  $t_i^l$  is the value of the left-truncation threshold,  $t_i^r$  is the value of the right-truncation threshold, and  $c_i^r$  is the value of the right-censoring limit,  $c_i^l$  is the value of the left-censoring limit.

If an observation is not left-truncated, then  $t_i^l = \tau^l$ , where  $\tau^l$  is the smallest value in the support of the distribution; so  $F(t_i^l) = 0$ . If an observation is not right-truncated, then  $t_i^r = \tau_h$ , where  $\tau_h$  is the largest value in the support of the distribution; so  $F(t_i^r) = 1$ . If an observation is not right-censored, then  $c_i^r = \tau^l$ ; so  $F(c_i^r) = 0$ . If an observation is not left-censored, then  $c_i^l = \tau_h$ ; so  $F(c_i^l) = 1$ .

Let  $w_i$  denote the weight associated with  $i$ th observation. If you specify the **WEIGHT** statement, then  $w_i$  is the normalized value of the weight variable; otherwise, it is set to 1. The weights are normalized such that they sum up to  $N$ .

An indicator function  $I[e]$  takes a value of 1 or 0 if the expression  $e$  is true or false, respectively.

## Estimation Methods

If the response variable is subject to both left-censoring and right-censoring effects and if you do not specify the **EMPIRICALCDF=NOTURNBULL** method, then PROC HPSEVERITY uses the Turnbull's method. This section describes methods other than Turnbull's method. For Turnbull's method, see the next section "Turnbull's EDF Estimation Method" on page 298.

The method descriptions assume that all observations are either uncensored or right-censored; that is, each observation is of the form  $(y_i, t_i^l, t_i^r, \tau^l, \tau_h)$  or  $(y_i, t_i^l, t_i^r, c_i^r, \tau_h)$ .

If all observations are either uncensored or left-censored, then each observation is of the form  $(y_i, t_i^l, t_i^r, \tau_l, c_i^l)$ . It is converted to an observation  $(-y_i, -t_i^r, -t_i^l, -c_i^l, \tau_h)$ ; that is, the signs of all the response variable values are reversed, the new left-truncation threshold is equal to the negative of the original right-truncation threshold, the new right-truncation threshold is equal to the negative of the original left-truncation threshold, and the negative of the original left-censoring limit becomes the new right-censoring limit. With this transformation, each observation is either uncensored or right-censored. The methods described for handling uncensored or right-censored data are now applicable. After the EDF estimates are computed, the observations are transformed back to the original form and EDF estimates are adjusted such  $F_n(y_i) = 1 - F_n(-y_i-)$ , where  $F_n(-y_i-)$  denotes the EDF estimate of the value slightly less than the transformed value  $-y_i$ .

Further, a set of uncensored or right-censored observations can be converted to a set of observations of the form  $(y_i, t_i^l, t_i^r, \delta_i)$ , where  $\delta_i$  is the indicator of right-censoring.  $\delta_i = 0$  indicates a right-censored observation, in which case  $y_i$  is assumed to record the right-censoring limit  $c_i^r$ .  $\delta_i = 1$  indicates an uncensored observation, and  $y_i$  records the exact observed value. In other words,  $\delta_i = I[Y \leq C^r]$  and  $y_i = \min(y_i, c_i^r)$ .

Given this notation, the EDF is estimated as

$$F_n(y) = \begin{cases} 0 & \text{if } y < y^{(1)} \\ \hat{F}_n(y^{(k)}) & \text{if } y^{(k)} \leq y < y^{(k+1)}, k = 1, \dots, N-1 \\ \hat{F}_n(y^{(N)}) & \text{if } y^{(N)} \leq y \end{cases}$$

where  $y^{(k)}$  denotes the  $k$ th order statistic of the set  $\{y_i\}$  and  $\hat{F}_n(y^{(k)})$  is the estimate computed at that value. The definition of  $\hat{F}_n$  depends on the estimation method. You can specify a particular method or let



PROC HPSEVERITY choose an appropriate method by using the **EMPIRICALCDF=** option in the PROC HPSEVERITY statement. Each method computes  $\hat{F}_n$  as follows:

**NOTURNBULL** This is the default method. First, censored observations, if any, are processed as follows:

- An observation that is left-censored but not right-censored is converted to an uncensored observation  $(y_i^u, t_i^l, t_i^r, \tau^l, \tau_h)$ , where  $y_i^u = c_i^l/2$ .
- An observation that is both left-censored and right-censored is converted to an uncensored observation  $(y_i^u, t_i^l, t_i^r, \tau^l, \tau_h)$ , where  $y_i^u = (c_i^r + c_i^l)/2$ .
- An observation that is right-censored but not left-censored is left unchanged.

If the processed set of observations contains any truncated or right-censored observations, the KAPLANMEIER method is used. Otherwise, the STANDARD method is used.

The observations are modified only for the purpose of computing the EDF estimates. The original censoring information is used by the parameter estimation process.

**STANDARD** If you do not specify any censoring or truncation information, then this method is chosen. It is the standard way of computing EDF. The EDF estimate at observation  $i$  is computed as follows:

$$\hat{F}_n(y_i) = \frac{1}{N} \sum_{j=1}^N I[y_j \leq y_i]$$

**KAPLANMEIER** This method is chosen when you specify at least one form of censoring or truncation. The Kaplan-Meier (KM) estimator, also known as the product-limit estimator, was first introduced by Kaplan and Meier (1958) for censored data. Lynden-Bell (1971) derived a similar estimator for left-truncated data. PROC HPSEVERITY uses the definition that combines both censoring and truncation information (Klein and Moeschberger 1997; Lai and Ying 1991).

The EDF estimate at observation  $i$  is computed as

$$\hat{F}_n(y_i) = 1 - \prod_{\tau \leq y_i} \left( 1 - \frac{n_\tau}{R_n(\tau)} \right)$$

where  $n_\tau$  and  $R_n(\tau)$  are defined as follows:

- $n_\tau = \sum_{k=1}^N w_k \cdot I[y_k = \tau \text{ and } \tau \leq t_k^r \text{ and } \delta_k = 1]$ , which is the number of uncensored observations ( $\delta_k = 1$ ) for which the response variable value is equal to  $\tau$  and  $\tau$  is observable according to the right-truncation threshold of that observation ( $\tau \leq t_k^r$ ).
- $R_n(\tau) = \sum_{k=1}^N w_k \cdot I[y_k \geq \tau > t_k^l]$ , which is the size (cardinality) of the risk set at  $\tau$ . The term *risk set* has its origins in survival analysis; it contains the events that are at risk of failure at a given time,  $\tau$ . In other words, it contains the events that have survived up to time  $\tau$  and might fail at or after  $\tau$ . For PROC HPSEVERITY, *time* is equivalent to the magnitude of the event and *failure* is equivalent to an uncensored and observable event, where observable means it satisfies the truncation thresholds.

**MODIFIEDKM**

The product-limit estimator used by the KAPLANMEIER method does not work well if the risk set size becomes very small. For right-censored data, the size can become small towards the right tail. For left-truncated data, the size can become small at the left tail and can remain so for the entire range of data. This was demonstrated by Lai and Ying (1991). They proposed a modification to the estimator that ignores the effects due to small risk set sizes.

The EDF estimate at observation  $i$  is computed as

$$\hat{F}_n(y_i) = 1 - \prod_{\tau \leq y_i} \left( 1 - \frac{n(\tau)}{R_n(\tau)} \cdot I[R_n(\tau) \geq cN^\alpha] \right)$$

where the definitions of  $n(\tau)$  and  $R_n(\tau)$  are identical to those used for the KAPLANMEIER method described previously.

You can specify the values of  $c$  and  $\alpha$  by using the **C=** and **ALPHA=** options. If you do not specify a value for  $c$ , the default value used is  $c = 1$ . If you do not specify a value for  $\alpha$ , the default value used is  $\alpha = 0.5$ .

As an alternative, you can also specify an absolute lower bound, say  $L$ , on the risk set size by using the **RSLB=** option, in which case  $I[R_n(\tau) \geq cN^\alpha]$  is replaced by  $I[R_n(\tau) \geq L]$  in the definition.

### Turnbull's EDF Estimation Method

If the response variable is subject to both left-censoring and right-censoring effects and if you do not specify the NOTURNBULL method, then the HPSEVERITY procedure uses a method proposed by Turnbull (1976) to compute the nonparametric estimates of the cumulative distribution function. The original Turnbull's method is modified using the suggestions made by Frydman (1994) when truncation effects are present.

Let the input data consist of  $N$  observations in the form of quintuplets of values  $(y_i, t_i^l, t_i^r, c_i^r, c_i^l)$ ,  $i = 1, \dots, N$  with notation described in the section “**Notation**” on page 296. For each observation, let  $A_i = (c_i^r, c_i^l]$  be the censoring interval; that is, the response variable value is known to lie in the interval  $A_i$ , but the exact value is not known. If an observation is uncensored, then  $A_i = (y_i - \epsilon, y_i]$  for any arbitrarily small value of  $\epsilon > 0$ . If an observation is censored, then the value  $y_i$  is ignored. Similarly, for each observation, let  $B_i = (t_i^l, t_i^r]$  be the truncation interval; that is, the observation is drawn from a truncated (conditional) distribution  $F(y, B_i) = P(Y \leq y | Y \in B_i)$ .

Two sets,  $L$  and  $R$ , are formed using  $A_i$  and  $B_i$  as follows:

$$\begin{aligned} L &= \{c_i^r, 1 \leq i \leq N\} \cup \{t_i^r, 1 \leq i \leq N\} \\ R &= \{c_i^l, 1 \leq i \leq N\} \cup \{t_i^l, 1 \leq i \leq N\} \end{aligned}$$

The sets  $L$  and  $R$  represent the left endpoints and right endpoints, respectively. A set of disjoint intervals  $C_j = [q_j, p_j]$ ,  $1 \leq j \leq M$  is formed such that  $q_j \in L$  and  $p_j \in R$  and  $q_j \leq p_j$  and  $p_j < q_{j+1}$ . The value of  $M$  is dependent on the nature of censoring and truncation intervals in the input data. Turnbull (1976) showed that the maximum likelihood estimate (MLE) of the EDF can increase only inside intervals  $C_j$ . In other words, the MLE estimate is constant in the interval  $(p_j, q_{j+1})$ . The likelihood is independent of the behavior of  $F_n$  inside any of the intervals  $C_j$ . Let  $s_j$  denote the increase in  $F_n$  inside an interval  $C_j$ . Then,

the EDF estimate is as follows:

$$F_n(y) = \begin{cases} 0 & \text{if } y < q_1 \\ \sum_{k=1}^j s_k & \text{if } p_j < y < q_{j+1}, 1 \leq j \leq M-1 \\ 1 & \text{if } y > p_M \end{cases}$$

PROC HPSEVERITY computes the estimates  $F_n(p_j+) = F_n(q_{j+1}-) = \sum_{k=1}^j s_k$  at points  $p_j$  and  $q_{j+1}$  and computes  $F_n(q_1-) = 0$  at point  $q_1$ , where  $F_n(x+)$  denotes the limiting estimate at a point that is infinitesimally larger than  $x$  when approaching  $x$  from values larger than  $x$  and where  $F_n(x-)$  denotes the limiting estimate at a point that is infinitesimally smaller than  $x$  when approaching  $x$  from values smaller than  $x$ .

PROC HPSEVERITY uses the expectation-maximization (EM) algorithm proposed by Turnbull (1976), who referred to the algorithm as the self-consistency algorithm. By default, the algorithm runs until one of the following criteria is met:

- Relative-error criterion: The maximum relative error between the two consecutive estimates of  $s_j$  falls below a threshold  $\epsilon$ . If  $l$  indicates an index of the current iteration, then this can be formally stated as

$$\arg \max_{1 \leq j \leq M} \left\{ \frac{|s_j^l - s_j^{l-1}|}{s_j^{l-1}} \right\} \leq \epsilon$$

You can control the value of  $\epsilon$  by specifying the **EPS=** suboption of the **EDF=TURNBULL** option in the PROC HPSEVERITY statement. The default value is 1.0E-8.

- Maximum-iteration criterion: The number of iterations exceeds an upper limit that you specify for the **MAXITER=** suboption of the **EDF=TURNBULL** option in the PROC HPSEVERITY statement. The default number of maximum iterations is 500.

The self-consistent estimates obtained in this manner might not be maximum likelihood estimates. Gentleman and Geyer (1994) suggested the use of the Kuhn-Tucker conditions for the maximum likelihood problem to ensure that the estimates are MLE. If you specify the **ENSUREMLE** suboption of the **EDF=TURNBULL** option in the PROC HPSEVERITY statement, then PROC HPSEVERITY computes the Kuhn-Tucker conditions at the end of each iteration to determine whether the estimates  $\{s_j\}$  are MLE. If you do not specify any truncation effects, then the Kuhn-Tucker conditions derived by Gentleman and Geyer (1994) are used. If you specify any truncation effects, then PROC HPSEVERITY uses modified Kuhn-Tucker conditions that account for the truncation effects. An integral part of checking the conditions is to determine whether an estimate  $s_j$  is zero or whether an estimate of the Lagrange multiplier or the reduced gradient associated with the estimate  $s_j$  is zero. PROC HPSEVERITY declares these values to be zero if they are less than or equal to a threshold  $\delta$ . You can control the value of  $\delta$  by specifying the **ZEROPROB=** suboption of the **EDF=TURNBULL** option in the PROC HPSEVERITY statement. The default value is 1.0E-8. The algorithm continues until the Kuhn-Tucker conditions are satisfied or the number of iterations exceeds the upper limit. The relative-error criterion stated previously is not used when you specify the **ENSUREMLE** option.

## EDF Estimates and Truncation

If you specify truncation, then the estimate  $\hat{F}_n(y)$  that is computed by any method other than the STANDARD method is a *conditional* estimate. In other words,  $\hat{F}_n(y) = \Pr(Y \leq y | \tau_G < Y \leq \tau_H)$ , where  $G$  and  $H$  denote the (unknown) distribution functions of the left-truncation threshold variable  $T^l$  and the right-truncation threshold variable  $T^r$ , respectively,  $\tau_G$  denotes the smallest left-truncation threshold with a nonzero cumulative probability, and  $\tau_H$  denotes the largest right-truncation threshold with a nonzero cumulative probability. Formally,  $\tau_G = \inf\{s : G(s) > 0\}$  and  $\tau_H = \sup\{s : H(s) > 0\}$ . For computational purposes, PROC HPSEVERITY estimates  $\tau_G$  and  $\tau_H$  by  $t_{\min}^l$  and  $t_{\max}^r$ , respectively, defined as

$$t_{\min}^l = \min\{t_k^l : 1 \leq k \leq N\}$$

$$t_{\max}^r = \max\{t_k^r : 1 \leq k \leq N\}$$

These estimates of  $t_{\min}^l$  and  $t_{\max}^r$  are used to compute the conditional estimates of the CDF as described in the section “[Truncation and Conditional CDF Estimates](#)” on page 285.

## Supplying EDF Estimates to Functions and Subroutines

The parameter initialization subroutines in distribution models and some predefined utility functions require EDF estimates. For more information, see the sections “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308 and “[Predefined Utility Functions](#)” on page 320.

PROC HPSEVERITY supplies the EDF estimates to these subroutines and functions by using two arrays,  $x$  and  $F$ , the dimension of each array, and a type of the EDF estimates. The type identifies how the EDF estimates are computed and stored. They are as follows:

- Type 1 specifies that EDF estimates are computed using the STANDARD method; that is, the data used for estimation are neither censored nor truncated.
- Type 2 specifies that EDF estimates are computed using the KAPLANMEIER method; that is, the data used for estimation are subject to at least one form of truncation or censoring.
- Type 3 specifies that EDF estimates are computed using the TURNBULL method; that is, the data used for estimation are subject to both left- and right-censoring. The data might or might not be truncated.

For Types 1 and 2, the EDF estimates are stored in arrays  $x$  and  $F$  of dimension  $N$  such that the following holds:

$$F_n(y) = \begin{cases} 0 & \text{if } y < x[1] \\ F[k] & \text{if } x[k] \leq y < x[k+1], k = 1, \dots, N-1 \\ F[N] & \text{if } x[N] \leq y \end{cases}$$

where  $[k]$  denotes  $k$ th element of the array ( $[1]$  denotes the first element of the array).

For Type 3, the EDF estimates are stored in arrays  $x$  and  $F$  of dimension  $N$  such that the following holds:

$$F_n(y) = \begin{cases} 0 & \text{if } y < x[1] \\ \text{undefined} & \text{if } x[2k-1] \leq y < x[2k], k = 1, \dots, (N-1)/2 \\ F[2k] = F[2k+1] & \text{if } x[2k] \leq y < x[2k+1], k = 1, \dots, (N-1)/2 \\ F[N] & \text{if } x[N] \leq y \end{cases}$$

Although the behavior of EDF is theoretically undefined for the interval  $[x[2k-1], x[2k])$ , for computational purposes, all predefined functions and subroutines assume that the EDF increases linearly from  $F[2k-1]$

to  $F[2k]$  in that interval if  $x[2k - 1] < x[2k]$ . If  $x[2k - 1] = x[2k]$ , which can happen when the EDF is estimated from a combination of uncensored and interval-censored data, the predefined functions and subroutines assume that  $F_n(x[2k - 1]) = F_n(x[2k]) = F[2k]$ .

## Statistics of Fit

PROC HPSEVERITY computes and reports various statistics of fit to indicate how well the estimated model fits the data. The statistics belong to two categories: likelihood-based statistics and EDF-based statistics. Neg2LogLike, AIC, AICC, and BIC are likelihood-based statistics, and KS, AD, and CvM are EDF-based statistics.

In the distributed mode of execution, in which data are distributed across the grid nodes, the EDF estimates are computed by using the local data. The EDF-based statistics are computed by using these local EDF estimates. The reported value of each EDF-based statistic is an average of the values of the statistic that are computed by all the grid nodes where the data reside. Also, for large data sets, in both single-machine and distributed modes of execution, the EDF estimates are computed by using a fraction of the input data that is governed by either the **INITSAMPLE** option or the default sample size. Because of this nature of computing the EDF estimates, the EDF-based statistics of fit are an approximation of the values that would have been computed if the entire input data set were used for computing the EDF estimates. So the values that are reported for EDF-based statistics should be used only for comparing different models. The reported values should not be interpreted as true estimates of the corresponding statistics.

The likelihood-based statistics are reported for the entire input data in both single-machine and distributed modes of execution.

The following subsections provide definitions of each category of statistics.

### Likelihood-Based Statistics of Fit

Let  $y_i, i = 1, \dots, N$ , denote the response variable values. Let  $L$  be the likelihood as defined in the section “**Likelihood Function**” on page 286. Let  $p$  denote the number of model parameters that are estimated. Note that  $p = p_d + (k - k_r)$ , where  $p_d$  is the number of distribution parameters,  $k$  is the number of regressors that you specify in the SCALEMODEL statement, and  $k_r$  is the number of regressors found to be linearly dependent (redundant) on other regressors. Given this notation, the likelihood-based statistics are defined as follows:

Neg2LogLike      The log likelihood is reported as

$$\text{Neg2LogLike} = -2 \log(L)$$

The multiplying factor  $-2$  makes it easy to compare it to the other likelihood-based statistics. A model that has a smaller value of Neg2LogLike is deemed better.

AIC      Akaike’s information criterion (AIC) is defined as

$$\text{AIC} = -2 \log(L) + 2p$$

A model that has a smaller AIC value is deemed better.

AICC The corrected Akaike's information criterion (AICC) is defined as

$$\text{AICC} = -2 \log(L) + \frac{2Np}{N - p - 1}$$

A model that has a smaller AICC value is deemed better. It corrects the finite-sample bias that AIC has when  $N$  is small compared to  $p$ . AICC is related to AIC as

$$\text{AICC} = \text{AIC} + \frac{2p(p + 1)}{N - p - 1}$$

As  $N$  becomes large compared to  $p$ , AICC converges to AIC. AICC is usually recommended over AIC as a model selection criterion.

BIC The Schwarz Bayesian information criterion (BIC) is defined as

$$\text{BIC} = -2 \log(L) + p \log(N)$$

A model that has a smaller BIC value is deemed better.

## EDF-Based Statistics

This class of statistics is based on the difference between the estimate of the cumulative distribution function (CDF) and the estimate of the empirical distribution function (EDF). A model that has a smaller value of the chosen EDF-based statistic is deemed better.

Let  $y_i, i = 1, \dots, N$  denote the sample of  $N$  values of the response variable. Let  $r_i = \sum_{j=1}^N I[y_j \leq y_i]$  denote the number of observations with a value less than or equal to  $y_i$ , where  $I$  is an indicator function. Let  $F_n(y_i)$  denote the EDF estimate that is computed by using the method that you specify in the **EMPIRICALCDF=** option. Let  $Z_i = \hat{F}(y_i)$  denote the estimate of the CDF. Let  $F_n(Z_i)$  denote the EDF estimate of  $Z_i$  values that are computed using the same method that is used to compute the EDF of  $y_i$  values. Using the probability integral transformation, if  $F(y)$  is the true distribution of the random variable  $Y$ , then the random variable  $Z = F(y)$  is uniformly distributed between 0 and 1 (D'Agostino and Stephens 1986, Ch. 4). Thus, comparing  $F_n(y_i)$  with  $\hat{F}(y_i)$  is equivalent to comparing  $F_n(Z_i)$  with  $\hat{F}(Z_i) = Z_i$  (uniform distribution).

Note the following two points regarding which CDF estimates are used for computing the test statistics:

- If you specify regressor variables, then the CDF estimates  $Z_i$  that are used for computing the EDF test statistics are from a mixture distribution. See the section “[CDF Estimates with Regression Effects](#)” on page 293 for more information.
- If the EDF estimates are conditional because of the truncation information, then each unconditional estimate  $Z_i$  is converted to a conditional estimate using the method described in the section “[Truncation and Conditional CDF Estimates](#)” on page 285.

In the following, it is assumed that  $Z_i$  denotes an appropriate estimate of the CDF if you specify any truncation or regression effects. Given this, the EDF-based statistics of fit are defined as follows:

KS The Kolmogorov-Smirnov (KS) statistic computes the largest vertical distance between the CDF and the EDF. It is formally defined as follows:

$$KS = \sup_y |F_n(y) - F(y)|$$

If the STANDARD method is used to compute the EDF, then the following formula is used:

$$\begin{aligned} D^+ &= \max_i \left( \frac{r_i}{N} - Z_i \right) \\ D^- &= \max_i \left( Z_i - \frac{r_{i-1}}{N} \right) \\ KS &= \sqrt{N} \max(D^+, D^-) + \frac{0.19}{\sqrt{N}} \end{aligned}$$

Note that  $r_0$  is assumed to be 0.

If the method used to compute the EDF is any method other than the STANDARD method, then the following formula is used:

$$\begin{aligned} D^+ &= \max_i (F_n(Z_i) - Z_i), \text{ if } F_n(Z_i) \geq Z_i \\ D^- &= \max_i (Z_i - F_n(Z_i)), \text{ if } F_n(Z_i) < Z_i \\ KS &= \sqrt{N} \max(D^+, D^-) + \frac{0.19}{\sqrt{N}} \end{aligned}$$

AD The Anderson-Darling (AD) statistic is a quadratic EDF statistic that is proportional to the expected value of the weighted squared difference between the EDF and CDF. It is formally defined as follows:

$$AD = N \int_{-\infty}^{\infty} \frac{(F_n(y) - F(y))^2}{F(y)(1 - F(y))} dF(y)$$

If the STANDARD method is used to compute the EDF, then the following formula is used:

$$AD = -N - \frac{1}{N} \sum_{i=1}^N [(2r_i - 1) \log(Z_i) + (2N + 1 - 2r_i) \log(1 - Z_i)]$$

If the method used to compute the EDF is any method other than the STANDARD method, then the statistic can be computed by using the following two pieces of information:

- If the EDF estimates are computed using the KAPLANMEIER or MODIFIEDKM methods, then EDF is a step function such that the estimate  $F_n(z)$  is a constant equal to  $F_n(Z_{i-1})$  in interval  $[Z_{i-1}, Z_i]$ . If the EDF estimates are computed using the TURNBULL method, then there are two types of intervals: one in which the EDF curve is constant and the other in which the EDF curve is theoretically undefined. For computational purposes, it is assumed that the EDF curve is linear for the latter type of the interval. For each method, the EDF estimate  $F_n(y)$  at  $y$  can be written as

$$F_n(z) = F_n(Z_{i-1}) + S_i(z - Z_{i-1}), \text{ for } z \in [Z_{i-1}, Z_i]$$

where  $S_i$  is the slope of the line defined as

$$S_i = \frac{F_n(Z_i) - F_n(Z_{i-1})}{Z_i - Z_{i-1}}$$

For the KAPLANMEIER or MODIFIEDKM method,  $S_i = 0$  in each interval.

- Using the probability integral transform  $z = F(y)$ , the formula simplifies to

$$AD = N \int_{-\infty}^{\infty} \frac{(F_n(z) - z)^2}{z(1-z)} dz$$

The computation formula can then be derived from the following approximation:

$$\begin{aligned} AD &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} \frac{(F_n(z) - z)^2}{z(1-z)} dz \\ &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} \frac{(F_n(Z_{i-1}) + S_i(z - Z_{i-1}) - z)^2}{z(1-z)} dz \\ &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} \frac{(P_i - Q_i z)^2}{z(1-z)} dz \end{aligned}$$

where  $P_i = F_n(Z_{i-1}) - S_i Z_{i-1}$ ,  $Q_i = 1 - S_i$ , and  $K$  is the number of points at which the EDF estimate are computed. For the TURNBULL method,  $K = 2k$  for some  $k$ .

Assuming  $Z_0 = 0$ ,  $Z_{K+1} = 1$ ,  $F_n(0) = 0$ , and  $F_n(Z_K) = 1$  yields the following computation formula:

$$\begin{aligned} AD &= -N(Z_1 + \log(1 - Z_1) + \log(Z_K) + (1 - Z_K)) \\ &\quad + N \sum_{i=2}^K [P_i^2 A_i - (Q_i - P_i)^2 B_i - Q_i^2 C_i] \end{aligned}$$

where  $A_i = \log(Z_i) - \log(Z_{i-1})$ ,  $B_i = \log(1 - Z_i) - \log(1 - Z_{i-1})$ , and  $C_i = Z_i - Z_{i-1}$ .

If EDF estimates are computed using the KAPLANMEIER or MODIFIEDKM method, then  $P_i = F_n(Z_{i-1})$  and  $Q_i = 1$ , which simplifies the formula as

$$\begin{aligned} AD &= -N(1 + \log(1 - Z_1) + \log(Z_K)) \\ &\quad + N \sum_{i=2}^K [F_n(Z_{i-1})^2 A_i - (1 - F_n(Z_{i-1}))^2 B_i] \end{aligned}$$

**CvM** The Cramér-von Mises (CvM) statistic is a quadratic EDF statistic that is proportional to the expected value of the squared difference between the EDF and CDF. It is formally defined as follows:

$$CvM = N \int_{-\infty}^{\infty} (F_n(y) - F(y))^2 dF(y)$$

If the STANDARD method is used to compute the EDF, then the following formula is used:

$$CvM = \frac{1}{12N} + \sum_{i=1}^N \left( Z_i - \frac{(2r_i - 1)}{2N} \right)^2$$

If the method used to compute the EDF is any method other than the STANDARD method, then the statistic can be computed by using the following two pieces of information:



- As described previously for the AD statistic, the EDF estimates are assumed to be piecewise linear such that the estimate  $F_n(y)$  at  $y$  is

$$F_n(z) = F_n(Z_{i-1}) + S_i(z - Z_{i-1}), \text{ for } z \in [Z_{i-1}, Z_i]$$

where  $S_i$  is the slope of the line defined as

$$S_i = \frac{F_n(Z_i) - F_n(Z_{i-1})}{Z_i - Z_{i-1}}$$

For the KAPLANMEIER or MODIFIEDKM method,  $S_i = 0$  in each interval.

- Using the probability integral transform  $z = F(y)$ , the formula simplifies to:

$$\text{CvM} = N \int_{-\infty}^{\infty} (F_n(z) - z)^2 dz$$

The computation formula can then be derived from the following approximation:

$$\begin{aligned} \text{CvM} &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} (F_n(z) - z)^2 dz \\ &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} (F_n(Z_{i-1}) + S_i(z - Z_{i-1}) - z)^2 dz \\ &= N \sum_{i=1}^{K+1} \int_{Z_{i-1}}^{Z_i} (P_i - Q_i z)^2 dz \end{aligned}$$

where  $P_i = F_n(Z_{i-1}) - S_i Z_{i-1}$ ,  $Q_i = 1 - S_i$ , and  $K$  is the number of points at which the EDF estimate are computed. For the TURNBULL method,  $K = 2k$  for some  $k$ .

Assuming  $Z_0 = 0$ ,  $Z_{K+1} = 1$ , and  $F_n(0) = 0$  yields the following computation formula:

$$\text{CvM} = N \frac{Z_1^3}{3} + N \sum_{i=2}^{K+1} \left[ P_i^2 A_i - P_i Q_i B_i - \frac{Q_i^2}{3} C_i \right]$$

where  $A_i = Z_i - Z_{i-1}$ ,  $B_i = Z_i^2 - Z_{i-1}^2$ , and  $C_i = Z_i^3 - Z_{i-1}^3$ .

If EDF estimates are computed using the KAPLANMEIER or MODIFIEDKM method, then  $P_i = F_n(Z_{i-1})$  and  $Q_i = 1$ , which simplifies the formula as

$$\text{CvM} = \frac{N}{3} + N \sum_{i=2}^{K+1} [F_n(Z_{i-1})^2 (Z_i - Z_{i-1}) - F_n(Z_{i-1})(Z_i^2 - Z_{i-1}^2)]$$

which is similar to the formula proposed by Koziol and Green (1976).

---

## Distributed and Multithreaded Computation

PROC HPSEVERITY makes an attempt to use all the computational resources that you specify in the PERFORMANCE statement in order to complete the assigned tasks as fast as possible. This section describes the distributed and multithreading computing methods that PROC HPSEVERITY uses.

## Distributed Computing

Distributed computing refers to the organization of computation work into multiple tasks that are processed on different nodes; a node is one of the machines that constitute the grid. The number of nodes that PROC HPSEVERITY uses is determined by the distributed processing execution mode. If you specify the client-data (or local-data) mode of execution, then the number of nodes is determined by the **NODES=** option in the **PERFORMANCE** statement. If you are using the alongside-the-database mode of execution, then PROC HPSEVERITY determines the number of nodes internally by using the information that is associated with the **DATA=** data set and the grid information that you specify either in the **PERFORMANCE** statement or in the grid environment variables. For more information about distributed processing modes, see the section “Processing Modes” on page 12.

In the client-data model, PROC HPSEVERITY distributes the input data across the number of nodes that you specify by sending the first observation to the first node, the second observation to the second node, and so on.

In the alongside-the-database model, PROC HPSEVERITY uses the existing distributed organization of the data. You do not need to specify the **NODES=** option.

The number of nodes that are used for distributed computing is displayed in the “Performance Information” table, which is part of the default output.

## Multithreading

Threading refers to the organization of computational work into multiple tasks (processing units that can be scheduled by the operating system). A task is associated with a thread. Multithreading refers to the concurrent execution of threads. When multithreading is possible, you can achieve more substantial performance gains than you can with sequential (single-threaded) execution.

The number of threads the HPSEVERITY procedure spawns is determined by the number of CPUs on a machine. You can control the number of CPUs in the following ways:

- You can use the **CPUCOUNT=** SAS system option to specify the CPU count. For example, if you specify the following statement, then PROC HPSEVERITY schedules threads as if it were executing on a system that had four CPUs, regardless of the actual CPU count:

```
options cpucount=4;
```

You can use this specification only in single-machine mode, and it does not take effect if the **THREADS** system option is turned off.

The default value of the **CPUCOUNT=** system option might not equal the number of all the logical CPU cores available on your machine, such as those available because of hyperthreading. To allow PROC HPSEVERITY to use all the logical cores in single-machine mode, specify the following **OPTIONS** statement:

```
options cpucount=actual;
```

- You can specify the **NTHREADS=** option in the **PERFORMANCE** statement. This specification overrides the **THREADS** and **CPUCOUNT=** system options. Specify **NTHREADS=1** to force single-threaded execution.

If you do not specify the `NTHREADS=` option and the `THREADS` system option is turned on, then the number of threads that are used in distributed mode is equal to the total number of logical CPU cores available on each node of the grid, and the number of threads used in single-machine mode is determined by the `CPUCOUNT=` system option.

If you do not specify the `NTHREADS=` option and the `THREADS` system option is turned off, then only one thread of execution is used in both single-machine and distributed modes.

The number of threads per machine is displayed in the “Performance Information” table, which is part of the default output.

Performance improvement is not always guaranteed when you use more threads, for several reasons: the increased cost of communication and synchronization among threads might offset the reduced cost of computation, the hyperthreading feature of the processor might not be very efficient for floating-point computations, and other applications might be running on the machine.

## Combining the Power of Distributed and Multithreading Computing

The `HPSEVERITY` procedure combines the powers of distributed and multithreading paradigms by using a data-parallel model. In particular, the distributed tasks are defined by dividing the data among multiple nodes, and within one node, the multithreading tasks are defined by further dividing the local data among the threads. For example, if the input data set has 10,000 observations and you are running on a grid that has five nodes, then each node processes 2,000 observations (this assumes that if you specify an alongside-the-database model, then you have equally and randomly divided the input data among the nodes). Further, if each node has eight CPUs, then 250 observations are associated with each thread within the node. All computations that require access to the data are then distributed and multithreaded.

Note that in single-machine mode (see the section “[Processing Modes](#)” on page 12), only multithreading is available.

When you specify more than one candidate distribution model, for some tasks `PROC HPSEVERITY` exploits the independence among models by processing multiple models in parallel on a single node such that each model is assigned to one of the threads executing in parallel. When a thread finishes processing the assigned model, it starts processing the next unprocessed model, if one exists.

The computations that take advantage of the distributed and multithreaded model include the following:

- **Validation and preparation of data:** In this stage, the observations in the input data set are validated and transformed, if necessary. The summary statistics of the data are prepared. Because each observation is independent, the computations can be distributed among nodes and among threads within nodes without significant communication overhead.
- **Initialization of distribution parameters:** In this stage, the parallelism is achieved by initializing multiple models in parallel. The only computational step that is not fully parallelized in this release is the step of computing empirical distribution function (EDF) estimates, which are required when `PROC HPSEVERITY` needs to invoke a distribution’s `PARMINIT` subroutine to initialize distribution parameters. The EDF estimation step is not amenable to full-fledged parallelism because it requires sequential access to sorted data, especially when the loss variable is modified by truncation effects. When the data are distributed across nodes, the EDF computations take place on local data and the `PARMINIT` function is invoked on the local data by using the local EDF estimates. The initial values

that are supplied to the nonlinear optimizer are computed by averaging the local estimates of the distribution parameters that are returned by the PARMINIT functions on each node.

- Initialization of regression parameters (if you specify the SCALEMODEL statement): In this stage, if you do not specify initial values for the regression parameters by using the INEST= data set, then PROC HPSEVERITY initializes those parameters by solving a linear regression problem  $\log(y) = \beta_0 + \sum_{j=1}^k \beta_j x_j$ . For more information, see the section “[Parameter Initialization for Regression Models](#)” on page 292. The most computationally intensive step is the formation of the crossproducts matrix. PROC HPSEVERITY exploits the parallelism by observing the fact that the contribution to the crossproducts matrix due to one observation is independent from the contribution due to another observation. Each node computes the contribution of its local data to each entry of the crossproducts matrix. Within each node, each thread computes the contribution of its chunk of data to each entry of the crossproducts matrix. On each node, the contributions from all the threads are added up to form the contribution due to all of the local data. The partial crossproducts matrices are then gathered from all nodes on a central node, which sums them up to form the final crossproducts matrix.
- Optimization: In this stage, the nonlinear optimizer iterates over the parameter space in search of the optimal set of parameters. In each iteration, it evaluates the objective function along with the gradient and Hessian of the objective function, if needed by the optimization method. Within one iteration, for the current estimates of the parameters, each observation’s contribution to the objective function, gradient, and Hessian is independent of another observation. This enables PROC HPSEVERITY to fully exploit the distributed and multithreaded paradigms to efficiently parallelize each iteration of the algorithm.

---

## Defining a Severity Distribution Model with the FCMP Procedure

A *severity distribution model* consists of a set of functions and subroutines that are defined using the FCMP procedure. The FCMP procedure is part of Base SAS software. Each function or subroutine must be named as `<distribution-name>_<keyword>`, where *distribution-name* is the identifying short name of the distribution and *keyword* identifies one of the functions or subroutines. The total length of the name should not exceed 32. Each function or subroutine must have a specific signature, which consists of the number of arguments, sequence and types of arguments, and return value type. The summary of all the recognized function and subroutine names and their expected behavior is given in [Table 9.6](#).

Consider the following points when you define a distribution model:

- When you define a function or subroutine requiring parameter arguments, the names and order of those arguments must be the same. Arguments other than the parameter arguments can have any name, but they must satisfy the requirements on their type and order.
- When the HPSEVERITY procedure invokes any function or subroutine, it provides the necessary input values according to the specified signature, and expects the function or subroutine to prepare the output and return it according to the specification of the return values in the signature.
- You can use most of the SAS programming statements and SAS functions that you can use in a DATA step for defining the FCMP functions and subroutines. However, there are a few differences in the capabilities of the DATA step and the FCMP procedure. To learn more, see the documentation of the FCMP procedure in the *Base SAS Procedures Guide*.

- You must specify either the PDF or the LOGPDF function. Similarly, you must specify either the CDF or the LOGCDF function. All other functions are optional, except when necessary for correct definition of the distribution. It is strongly recommended that you define the PARMINIT subroutine to provide a good set of initial values for the parameters. The information provided by PROC HPSEVERITY to the PARMINIT subroutine enables you to use popular initialization approaches based on the method of moments and the method of percentile matching, but you can implement any algorithm to initialize the parameters by using the values of the response variable and the estimate of its empirical distribution function.
- The LOWERBOUNDS subroutines should be defined if the lower bound on at least one distribution parameter is different from the default lower bound of 0. If you define a LOWERBOUNDS subroutine but do not set a lower bound for some parameter inside the subroutine, then that parameter is assumed to have no lower bound (or a lower bound of  $-\infty$ ). Hence, it is recommended that you explicitly return the lower bound for each parameter when you define the LOWERBOUNDS subroutine.
- The UPPERBOUNDS subroutines should be defined if the upper bound on at least one distribution parameter is different from the default upper bound of  $\infty$ . If you define an UPPERBOUNDS subroutine but do not set an upper bound for some parameter inside the subroutine, then that parameter is assumed to have no upper bound (or a upper bound of  $\infty$ ). Hence, it is recommended that you explicitly return the upper bound for each parameter when you define the UPPERBOUNDS subroutine.
- If you want to use the distribution in a model with regression effects, then make sure that the first parameter of the distribution is the scale parameter itself or a log-transformed scale parameter. If the first parameter is a log-transformed scale parameter, then you must define the SCALETRANSFORM function.
- In general, it is not necessary to define the gradient and Hessian functions, because the HPSEVERITY procedure uses an internal system to evaluate the required derivatives. The internal system typically computes the derivatives analytically. But it might not be able to do so if your function definitions use other functions that it cannot differentiate analytically. In such cases, derivatives are approximated using a finite difference method and a note is written to the SAS log to indicate the components that are differentiated using such approximations. PROC HPSEVERITY does reasonably well with these finite difference approximations. But, if you know of a way to compute the derivatives of such components analytically, then you should define the gradient and Hessian functions.

In order to use your distribution with PROC HPSEVERITY, you need to record the FCMP library that contains the functions and subroutines for your distribution and other FCMP libraries that contain FCMP functions or subroutines used within your distribution's functions and subroutines. Specify all those libraries in the CMPLIB= system option by using the OPTIONS global statement. For more information about the OPTIONS statement, see *SAS Statements: Reference*. For more information about the CMPLIB= system option, see *SAS System Options: Reference*.

Each predefined distribution mentioned in the section “[Predefined Distributions](#)” on page 274 has a distribution model associated with it. The functions and subroutines of all those models are available in the Sashelp.Svrtldist library. The order of the parameters in the signatures of the functions and subroutines is the same as listed in [Table 9.2](#). You do not need to use the CMPLIB= option in order to use the predefined distributions with PROC HPSEVERITY. However, if you need to use the functions or subroutines of the predefined distributions in SAS statements other than the PROC HPSEVERITY step (such as in a DATA step), then specify the Sashelp.Svrtldist library in the CMPLIB= system option by using the OPTIONS global statement prior to using them.

Table 9.6 shows functions and subroutines that define a distribution model, and subsections after the table provide more detail. The functions are listed in alphabetical order of the keyword suffix.

**Table 9.6** List of Functions and Subroutines That Define a Distribution Model

Name	Type	Required	Expected to Return
<i>dist_CDF</i>	Function	YES <sup>1</sup>	Cumulative distribution function value
<i>dist_CDFGRADIENT</i>	Subroutine	NO	Gradient of the CDF
<i>dist_CDFHESSIAN</i>	Subroutine	NO	Hessian of the CDF
<i>dist_CONSTANTPARM</i>	Subroutine	NO	Constant parameters
<i>dist_DESCRIPTION</i>	Function	NO	Description of the distribution
<i>dist_LOGCDF</i>	Function	YES <sup>1</sup>	Log of cumulative distribution function value
<i>dist_LOGCDFGRADIENT</i>	Subroutine	NO	Gradient of the LOGCDF
<i>dist_LOGCDFHESSIAN</i>	Subroutine	NO	Hessian of the LOGCDF
<i>dist_LOGPDF</i>	Function	YES <sup>2</sup>	Log of probability density function value
<i>dist_LOGPDFGRADIENT</i>	Subroutine	NO	Gradient of the LOGPDF
<i>dist_LOGPDFHESSIAN</i>	Subroutine	NO	Hessian of the LOGPDF
<i>dist_LOGSDF</i>	Function	NO	Log of survival function value
<i>dist_LOGSDFGRADIENT</i>	Subroutine	NO	Gradient of the LOGSDF
<i>dist_LOGSDFHESSIAN</i>	Subroutine	NO	Hessian of the LOGSDF
<i>dist_LOWERBOUNDS</i>	Subroutine	NO	Lower bounds on parameters
<i>dist_PARMINIT</i>	Subroutine	NO	Initial values for parameters
<i>dist_PDF</i>	Function	YES <sup>2</sup>	Probability density function value
<i>dist_PDFGRADIENT</i>	Subroutine	NO	Gradient of the PDF
<i>dist_PDFHESSIAN</i>	Subroutine	NO	Hessian of the PDF
<i>dist_SCALETRANSFORM</i>	Function	NO	Type of relationship between the first distribution parameter and the scale parameter
<i>dist_SDF</i>	Function	NO	Survival function value
<i>dist_SDFGRADIENT</i>	Subroutine	NO	Gradient of the SDF
<i>dist_SDFHESSIAN</i>	Subroutine	NO	Hessian of the SDF
<i>dist_UPPERBOUNDS</i>	Subroutine	NO	Upper bounds on parameters

Notes:

1. Either the *dist\_CDF* or the *dist\_LOGCDF* function must be defined.
2. Either the *dist\_PDF* or the *dist\_LOGPDF* function must be defined.

The signature syntax and semantics of each function or subroutine are as follows:

**dist\_CDF**

defines a function that returns the value of the cumulative distribution function (CDF) of the distribution at the specified values of the random variable and distribution parameters.

- *Type*: Function
- *Required*: YES
- *Number of arguments*:  $m + 1$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments*:

$x$  Numeric value of the random variable at which the CDF value should be evaluated

$p_1$  Numeric value of the first parameter

$p_2$  Numeric value of the second parameter

.....

$p_m$  Numeric value of the  $m$ th parameter

- *Return value*: Numeric value that contains the CDF value  $F(x; p_1, p_2, \dots, p_m)$

If you want to consider this distribution as a candidate distribution when you estimate a response variable model with regression effects, then the first parameter of this distribution must be a scale parameter or log-transformed scale parameter. In other words, if the distribution has a scale parameter, then the following equation must be satisfied:

$$F(x; p_1, p_2, \dots, p_m) = F\left(\frac{x}{p_1}; 1, p_2, \dots, p_m\right)$$

If the distribution has a log-transformed scale parameter, then the following equation must be satisfied:

$$F(x; p_1, p_2, \dots, p_m) = F\left(\frac{x}{\exp(p_1)}; 0, p_2, \dots, p_m\right)$$

Here is a sample structure of the function for a distribution named 'FOO':

```
function FOO_CDF(x, P1, P2);
  /* Code to compute CDF by using x, P1, and P2 */

  F = <computed CDF>;
  return (F);
endsub;
```

**dist\_CONSTANTPARM**

defines a subroutine that specifies constant parameters. A parameter is *constant* if it is required for defining a distribution but is not subject to optimization in PROC HPSEVERITY. Constant parameters are required to be part of the model in order to compute the PDF or the CDF of the distribution. Typically, values of these parameters are known a priori or estimated using some means other than the maximum likelihood method used by PROC HPSEVERITY. You can estimate them inside the *dist\_PARMINIT* subroutine. Once initialized, the parameters remain constant in the context of PROC HPSEVERITY; that is, they retain their initial value. PROC HPSEVERITY estimates only the nonconstant parameters.

- *Type*: Subroutine
- *Required*: NO
- *Number of arguments*:  $k$ , where  $k$  is the number of constant parameters
- *Sequence and type of arguments*:

constant parameter 1    Name of the first constant parameter

.....

constant parameter  $k$     Name of the  $k$ th constant parameter

- *Return value*: None

Here is a sample structure of the subroutine for a distribution named 'FOO' that has P3 and P5 as its constant parameters, assuming that distribution has at least three parameters:

```
subroutine FOO_CONSTANTPARM(p5, p3);
endsub;
```

Note the following points when you specify the constant parameters:

- At least one distribution parameter must be free to be optimized; that is, if a distribution has total  $m$  parameters, then  $k$  must be strictly less than  $m$ .
- If you want to use this distribution for modeling regression effects, then the first parameter must not be a constant parameter.
- The order of arguments in the signature of this subroutine does not matter as long as each argument's name matches the name of one of the parameters that are defined in the signature of the *dist\_PDF* function.
- The constant parameters must be specified in signatures of all the functions and subroutines that accept distribution parameters as their arguments.
- You must provide a nonmissing initial value for each constant parameter by using one of the supported parameter initialization methods.

#### *dist\_***DESCRIPTION**

defines a function that returns a description of the distribution.

- *Type*: Function
- *Required*: NO
- *Number of arguments*: None
- *Sequence and type of arguments*: Not applicable
- *Return value*: Character value containing a description of the distribution

Here is a sample structure of the function for a distribution named 'FOO':

```
function FOO_DESCRIPTION() $48;
  length desc $48;
  desc = "A model for a continuous distribution named foo";
  return (desc);
endsub;
```



There is no restriction on the length of the description (the length of 48 used in the previous example is for illustration purposes only). However, if the length is greater than 256, then only the first 256 characters appear in the displayed output and in the `_DESCRIPTION_` variable of the `OUTMODELINFO=` data set. Hence, the recommended length of the description is less than or equal to 256.

#### `dist_`**LOGcore**

defines a function that returns the natural logarithm of the specified *core* function of the distribution at the specified values of the random variable and distribution parameters. The *core* keyword can be PDF, CDF, or SDF.

- *Type*: Function
- *Required*: YES only if *core* is PDF or CDF and you have not defined that *core* function; otherwise, NO
- *Number of arguments*:  $m + 1$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments*:
  - x    Numeric value of the random variable at which the natural logarithm of the *core* function should be evaluated
  - p1   Numeric value of the first parameter
  - p2   Numeric value of the second parameter
  - .....
  - pm   Numeric value of the  $m$ th parameter
- *Return value*: Numeric value that contains the natural logarithm of the *core* function

Here is a sample structure of the function for the core function PDF of a distribution named 'FOO':

```
function FOO_LOGPDF(x, P1, P2);
  /* Code to compute LOGPDF by using x, P1, and P2 */

  l = <computed LOGPDF>;
  return (l);
endsub;
```

#### `dist_`**LOWERBOUNDS**

defines a subroutine that returns lower bounds for the parameters of the distribution. If this subroutine is not defined for a given distribution, then the HPSEVERITY procedure assumes a lower bound of 0 for each parameter. If a lower bound of  $l_i$  is returned for a parameter  $p_i$ , then the HPSEVERITY procedure assumes that  $l_i < p_i$  (strict inequality). If a missing value is returned for some parameter, then the HPSEVERITY procedure assumes that there is no lower bound for that parameter (equivalent to a lower bound of  $-\infty$ ).

- *Type*: Subroutine
- *Required*: NO
- *Number of arguments*:  $m$ , where  $m$  is the number of distribution parameters

- *Sequence and type of arguments:*

- |       |   |
|-------|---|
| p1    | Output argument that returns the lower bound on the first parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.       |
| p2    | Output argument that returns the lower bound on the second parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.      |
| ..... |   |
| pm    | Output argument that returns the lower bound on the <i>m</i> th parameter. You must specify this in the OUTARGS statement inside the subroutine's definition. |
- *Return value:* The results, lower bounds on parameter values, should be returned in the parameter arguments of the subroutine.

Here is a sample structure of the subroutine for a distribution named 'FOO':

```
subroutine FOO_LOWERBOUNDS(p1, p2);
  outargs p1, p2;

  p1 = <lower bound for P1>;
  p2 = <lower bound for P2>;
endsub;
```

#### *dist\_***PARMINIT**

defines a subroutine that returns the initial values for the distribution's parameters given an empirical distribution function (EDF) estimate.

- *Type:* Subroutine
- *Required:* NO
- *Number of arguments:*  $m + 4$ , where *m* is the number of distribution parameters
- *Sequence and type of arguments:*

- |       |   |
|-------|---|
| dim   | Input numeric value that contains the dimension of the x, nx, and F array arguments.  |
| x{*}  | Input numeric array of dimension <i>dim</i> that contains values of the random variables at which the EDF estimate is available. It can be assumed that x contains values in an increasing order. In other words, if $i < j$ , then $x[i] < x[j]$ . |
| nx{*} | Input numeric array of dimension <i>dim</i> . Each $nx[i]$ contains the number of observations in the original data that have the value $x[i]$ .  |
| F{*}  | Input numeric array of dimension <i>dim</i> . Each $F[i]$ contains the EDF estimate for $x[i]$ . This estimate is computed by the HPSEVERITY procedure based on the options that you specify in the LOSS statement.                                 |
| Ftype | Input numeric value that contains the type of the EDF estimate that is stored in x and F. See the section “ <a href="#">Supplying EDF Estimates to Functions and Subroutines</a> ” on page 300 for definition of types.                             |
| p1    | Output argument that returns the initial value of the first parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.   |

- p2*      Output argument that returns the initial value of the second parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.
- .....
- pm*      Output argument that returns the initial value of the *m*th parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.
- *Return value:* The results, initial values of the parameters, should be returned in the parameter arguments of the subroutine.

Here is a sample structure of the subroutine for a distribution named 'FOO':

```
subroutine FOO_PARMINIT(dim, x{*}, nx{*}, F{*}, Ftype, p1, p2);
  outargs p1, p2;

  /* Code to initialize values of P1 and P2 by using
     dim, x, nx, and F */

  p1 = <initial value for p1>;
  p2 = <initial value for p2>;
endsub;
```

#### *dist\_PDF*

defines a function that returns the value of the probability density function (PDF) of the distribution at the specified values of the random variable and distribution parameters.

- *Type:* Function
- *Required:* YES
- *Number of arguments:*  $m + 1$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments:*

*x*      Numeric value of the random variable at which the PDF value should be evaluated

*p1*      Numeric value of the first parameter

*p2*      Numeric value of the second parameter

.....

*pm*      Numeric value of the *m*th parameter

- *Return value:* Numeric value that contains the PDF value  $f(x; p_1, p_2, \dots, p_m)$

If you want to consider this distribution as a candidate distribution when you estimate a response variable model with regression effects, then the first parameter of this distribution must be a scale parameter or log-transformed scale parameter. In other words, if the distribution has a scale parameter, then the following equation must be satisfied:

$$f(x; p_1, p_2, \dots, p_m) = \frac{1}{p_1} f\left(\frac{x}{p_1}; 1, p_2, \dots, p_m\right)$$

If the distribution has a log-transformed scale parameter, then the following equation must be satisfied:

$$f(x; p_1, p_2, \dots, p_m) = \frac{1}{\exp(p_1)} f\left(\frac{x}{\exp(p_1)}; 0, p_2, \dots, p_m\right)$$

Here is a sample structure of the function for a distribution named 'FOO':

```
function FOO_PDF(x, P1, P2);
  /* Code to compute PDF by using x, P1, and P2 */

  f = <computed PDF>;
  return (f);
endsub;
```

### *dist\_***SCALETRANSFORM**

defines a function that returns a keyword to identify the transform that needs to be applied to the scale parameter to convert it to the first parameter of the distribution.

If you want to use this distribution for modeling regression effects, then the first parameter of this distribution must be a scale parameter. However, for some distributions, a typical or convenient parameterization might not have a scale parameter, but one of the parameters can be a simple transform of the scale parameter. As an example, consider a typical parameterization of the lognormal distribution with two parameters, location  $\mu$  and shape  $\sigma$ , for which the PDF is defined as follows:

$$f(x; \mu, \sigma) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\log(x)-\mu}{\sigma}\right)^2}$$

You can reparameterize this distribution to contain a parameter  $\theta$  instead of the parameter  $\mu$  such that  $\mu = \log(\theta)$ . The parameter  $\theta$  would then be a scale parameter. However, if you want to specify the distribution in terms of  $\mu$  and  $\sigma$  (which is a more recognized form of the lognormal distribution) and still allow it as a candidate distribution for estimating regression effects, then instead of writing another distribution with parameters  $\theta$  and  $\sigma$ , you can simply define the distribution with  $\mu$  as the first parameter and specify that it is the logarithm of the scale parameter.

- *Type*: Function
- *Required*: NO
- *Number of arguments*: None
- *Sequence and type of arguments*: Not applicable
- *Return value*: Character value that contains one of the following keywords:

LOG	specifies that the first parameter is the logarithm of the scale parameter.
IDENTITY	specifies that the first parameter is a scale parameter without any transformation.

If you do not specify this function, then the IDENTITY transform is assumed.

Here is a sample structure of the function for a distribution named 'FOO':

```
function FOO_SCALETRANSFORM() $8;
  length xform $8;
  xform = "IDENTITY";
  return (xform);
endsub;
```

**dist\_SDF**

defines a function that returns the value of the survival distribution function (SDF) of the distribution at the specified values of the random variable and distribution parameters.

- *Type*: Function
- *Required*: NO
- *Number of arguments*:  $m + 1$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments*:

$x$     Numeric value of the random variable at which the SDF value should be evaluated

$p_1$    Numeric value of the first parameter

$p_2$    Numeric value of the second parameter

.....

$p_m$    Numeric value of the  $m$ th parameter

- *Return value*: Numeric value that contains the SDF value  $S(x; p_1, p_2, \dots, p_m)$

If you want to consider this distribution as a candidate distribution when estimating a response variable model with regression effects, then the first parameter of this distribution must be a scale parameter or log-transformed scale parameter. In other words, if the distribution has a scale parameter, then the following equation must be satisfied:

$$S(x; p_1, p_2, \dots, p_m) = S\left(\frac{x}{p_1}; 1, p_2, \dots, p_m\right)$$

If the distribution has a log-transformed scale parameter, then the following equation must be satisfied:

$$S(x; p_1, p_2, \dots, p_m) = S\left(\frac{x}{\exp(p_1)}; 0, p_2, \dots, p_m\right)$$

Here is a sample structure of the function for a distribution named 'FOO':

```
function FOO_SDF(x, P1, P2);
  /* Code to compute SDF by using x, P1, and P2 */

  S = <computed SDF>;
  return (S);
endsub;
```

**dist\_UPPERBOUNDS**

defines a subroutine that returns upper bounds for the parameters of the distribution. If this subroutine is not defined for a given distribution, then the HPSEVERITY procedure assumes that there is no upper bound for any of the parameters. If an upper bound of  $u_i$  is returned for a parameter  $p_i$ , then the HPSEVERITY procedure assumes that  $p_i < u_i$  (strict inequality). If a missing value is returned for some parameter, then the HPSEVERITY procedure assumes that there is no upper bound for that parameter (equivalent to an upper bound of  $\infty$ ).

- *Type*: Subroutine
- *Required*: NO
- *Number of arguments*:  $m$ , where  $m$  is the number of distribution parameters

- *Sequence and type of arguments:*

p1	Output argument that returns the upper bound on the first parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.
p2	Output argument that returns the upper bound on the second parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.
....	
pm	Output argument that returns the upper bound on the <i>m</i> th parameter. You must specify this in the OUTARGS statement inside the subroutine's definition.

- *Return value:* The results, upper bounds on parameter values, should be returned in the parameter arguments of the subroutine.

Here is a sample structure of the subroutine for a distribution named 'FOO':

```
subroutine FOO_UPPERBOUNDS(p1, p2);
  outargs p1, p2;

  p1 = <upper bound for P1>;
  p2 = <upper bound for P2>;
endsub;
```

#### *dist\_core***GRADIENT**

defines a subroutine that returns the gradient vector of the specified *core* function of the distribution at the specified values of the random variable and distribution parameters. The *core* keyword can be PDF, CDF, SDF, LOGPDF, LOGCDF, or LOGSDF.

- *Type:* Subroutine
- *Required:* NO
- *Number of arguments:*  $m + 2$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments:*

x	Numeric value of the random variable at which the gradient should be evaluated
p1	Numeric value of the first parameter
p2	Numeric value of the second parameter
....	
pm	Numeric value of the <i>m</i> th parameter
grad{*}	Output numeric array of size <i>m</i> that contains the gradient vector evaluated at the specified values. If <i>h</i> denotes the value of the <i>core</i> function, then the expected order of the values in the array is as follows: $\frac{\partial h}{\partial p_1} \frac{\partial h}{\partial p_2} \dots \frac{\partial h}{\partial p_m}$

- *Return value:* Numeric array that contains the gradient evaluated at *x* for the parameter values  $(p_1, p_2, \dots, p_m)$

Here is a sample structure of the function for the core function CDF of a distribution named 'FOO':

```
subroutine FOO_CDFGRADIENT(x, P1, P2, grad{*});
  outargs grad;

  /* Code to compute gradient by using x, P1, and P2 */
  grad[1] = <partial derivative of CDF w.r.t. P1
             evaluated at x, P1, P2>;
  grad[2] = <partial derivative of CDF w.r.t. P2
             evaluated at x, P1, P2>;
endsub;
```

#### *dist\_core*HESIAN

defines a subroutine that returns the Hessian matrix of the specified *core* function of the distribution at the specified values of the random variable and distribution parameters. The *core* keyword can be PDF, CDF, SDF, LOGPDF, LOGCDF, or LOGSDF.

- *Type*: Subroutine
- *Required*: NO
- *Number of arguments*:  $m + 2$ , where  $m$  is the number of distribution parameters
- *Sequence and type of arguments*:

x	Numeric value of the random variable at which the Hessian matrix should be evaluated
p1	Numeric value of the first parameter
p2	Numeric value of the second parameter
....	
pm	Numeric value of the $m$ th parameter
hess{*}	Output numeric array of size $m(m + 1)/2$ that contains the lower triangular portion of the Hessian matrix in a packed vector form, evaluated at the specified values. If $h$ denotes the value of the <i>core</i> function, then the expected order of the values in the array is as follows: $\frac{\partial^2 h}{\partial p_1^2} \mid \frac{\partial^2 h}{\partial p_1 \partial p_2} \frac{\partial^2 h}{\partial p_2^2} \mid \cdots \mid \frac{\partial^2 h}{\partial p_1 \partial p_m} \frac{\partial^2 h}{\partial p_2 \partial p_m} \cdots \frac{\partial^2 h}{\partial p_m^2}$

- *Return value*: Numeric array that contains the lower triangular portion of the Hessian matrix evaluated at  $x$  for the parameter values  $(p_1, p_2, \dots, p_m)$

Here is a sample structure of the subroutine for the core function LOGSDF of a distribution named 'FOO':

```
subroutine FOO_LOGSDFHESIAN(x, P1, P2, hess{*});
  outargs hess;

  /* Code to compute Hessian by using x, P1, and P2 */
  hess[1] = <second order partial derivative of LOGSDF
             w.r.t. P1 evaluated at x, P1, P2>;
  hess[2] = <second order partial derivative of LOGSDF
             w.r.t. P1 and P2 evaluated at x, P1, P2>;
  hess[3] = <second order partial derivative of LOGSDF
             w.r.t. P2 evaluated at x, P1, P2>;
endsub;
```

## Predefined Utility Functions

The following predefined utility functions are provided with the HPSEVERITY procedure and are available in the Sashelp.Svrtldist library:

### SVRTUTIL\_EDF:

This function computes the empirical distribution function (EDF) estimate at the specified value of the random variable given the EDF estimate for a sample.

- *Type:* Function
- *Signature:* SVRTUTIL\_EDF(y, n, x{\*}, F{\*}, Ftype)
- *Argument Description:*

y	Value of the random variable at which the EDF estimate is desired.
n	Dimension of the $x$ and $F$ input arrays.
x{*}	Input numeric array of dimension $n$ that contains values of the random variable observed in the sample. These values are sorted in nondecreasing order.
F{*}	Input numeric array of dimension $n$ in which each $F[i]$ contains the EDF estimate for $x[i]$ . These values must be sorted in nondecreasing order.
Ftype	Type of the empirical estimate that is stored in the $x$ and $F$ arrays. See the section <a href="#">“Supplying EDF Estimates to Functions and Subroutines”</a> on page 300 for definition of types.
- *Return value:* The EDF estimate at  $y$ .

The type of the sample EDF estimate determines how the EDF estimate at  $y$  is computed. For more information, see the section [“Supplying EDF Estimates to Functions and Subroutines”](#) on page 300.

### SVRTUTIL\_EMPLIMMOMENT:

This function computes the empirical estimate of the limited moment of specified order for the specified upper limit, given the EDF estimate for a sample.

- *Type:* Function
- *Signature:* SVRTUTIL\_EMPLIMMOMENT(k, u, n, x{\*}, F{\*}, Ftype)
- *Argument Description:*

k	Order of the desired empirical limited moment.
u	Upper limit on the value of the random variable to be used in the computation of the desired empirical limited moment.
n	Dimension of the $x$ and $F$ input arrays.
x{*}	Input numeric array of dimension $n$ that contains values of the random variable observed in the sample. These values are sorted in nondecreasing order.
F{*}	Input numeric array of dimension $n$ in which each $F[i]$ contains the EDF estimate for $x[i]$ . These values must be sorted in nondecreasing order.



**Ftype** Type of the empirical estimate that is stored in the  $x$  and  $F$  arrays. See the section “[Supplying EDF Estimates to Functions and Subroutines](#)” on page 300 for definition of types.

- *Return value:* The desired empirical limited moment.

The empirical limited moment is computed by using the empirical estimate of the CDF. If  $F_n(x)$  denotes the EDF at  $x$ , then the empirical limited moment of order  $k$  with upper limit  $u$  is defined as

$$E_n[(X \wedge u)^k] = k \int_0^u (1 - F_n(x))x^{k-1} dx$$

The SVRTUTIL\_EMPLIMMOMENT function uses the piecewise linear nature of  $F_n(x)$  as described in the section “[Supplying EDF Estimates to Functions and Subroutines](#)” on page 300 to compute the integration.

#### SVRTUTIL\_HILLCUTOFF:

This function computes an estimate of the value where the right tail of a distribution is expected to begin. The function implements the algorithm described in Danielsson et al. 2001. The description of the algorithm uses the following notation:

- $n$  Number of observations in the original sample.
- $B$  Number of bootstrap samples to draw.
- $m_1$  Size of the bootstrap sample in the first step of the algorithm ( $m_1 < n$ ).
- $x_{(i)}^{j,m}$   $i$ th order statistic of  $j$ th bootstrap sample of size  $m$  ( $1 \leq i \leq m, 1 \leq j \leq B$ ).
- $x_{(i)}$   $i$ th order statistic of the original sample ( $1 \leq i \leq n$ ).

Given the input sample  $x$  and values of  $B$  and  $m_1$ , the steps of the algorithm are as follows:

1. Take  $B$  bootstrap samples of size  $m_1$  from the original sample.
2. Find the integer  $k_1$  that minimizes the bootstrap estimate of the mean squared error:

$$k_1 = \arg \min_{1 \leq k < m_1} Q(m_1, k)$$

3. Take  $B$  bootstrap samples of size  $m_2 = m_1^2/n$  from the original sample.
4. Find the integer  $k_2$  that minimizes the bootstrap estimate of the mean squared error:

$$k_2 = \arg \min_{1 \leq k < m_2} Q(m_2, k)$$

5. Compute the integer  $k_{\text{opt}}$ , which is used for computing the cutoff point:

$$k_{\text{opt}} = \frac{k_1^2}{k_2} \left( \frac{\log(k_1)}{2 \log(m_1) - \log(k_1)} \right)^{2 - 2 \log(k_1) / \log(m_1)}$$

6. Set the cutoff point equal to  $x_{(k_{\text{opt}}+1)}$ .

The bootstrap estimate of the mean squared error is computed as

$$Q(m, k) = \frac{1}{B} \sum_{j=1}^B \text{MSE}_j(m, k)$$

The mean squared error of  $j$ th bootstrap sample is computed as

$$\text{MSE}_j(m, k) = (M_j(m, k) - 2(\gamma_j(m, k))^2)^2$$

where  $M_j(m, k)$  is a control variate proposed by Danielsson et al. 2001,

$$M_j(m, k) = \frac{1}{k} \sum_{i=1}^k \left( \log(x_{(m-i+1)}^{j,m}) - \log(x_{(m-k)}^{j,m}) \right)^2$$

and  $\gamma_j(m, k)$  is the Hill's estimator of the tail index (Hill 1975),

$$\gamma_j(m, k) = \frac{1}{k} \sum_{i=1}^k \log(x_{(m-i+1)}^{j,m}) - \log(x_{(m-k)}^{j,m})$$

This algorithm has two tuning parameters,  $B$  and  $m_1$ . The number of bootstrap samples  $B$  is chosen based on the availability of computational resources. The optimal value of  $m_1$  is chosen such that the following ratio,  $R(m_1)$ , is minimized:

$$R(m_1) = \frac{(Q(m_1, k_1))^2}{Q(m_2, k_2)}$$

The SVRTUTIL\_HILLCUTOFF utility function implements the preceding algorithm. It uses the grid search method to compute the optimal value of  $m_1$ .

- *Type:* Function
- *Signature:* SVRTUTIL\_HILLCUTOFF( $n$ ,  $x\{*\}$ ,  $b$ ,  $s$ ,  $status$ )
- *Argument Description:*

$n$	Dimension of the array $x$ .
$x\{*\}$	Input numeric array of dimension $n$ that contains the sample.
$b$	Number of bootstrap samples used to estimate the mean squared error. If $b$ is less than 10, then a default value of 50 is used.
$s$	Approximate number of steps used to search the optimal value of $m_1$ in the range $[n^{0.75}, n - 1]$ . If $s$ is less than or equal to 1, then a default value of 10 is used.
$status$	Output argument that contains the status of the algorithm. If the algorithm succeeds in computing a valid cutoff point, then <i>status</i> is set to 0. If the algorithm fails, then <i>status</i> is set to 1.
- *Return value:* The cutoff value where the right tail is estimated to start. If the size of the input sample is inadequate ( $n \leq 5$ ), then a missing value is returned and *status* is set to a missing value. If the algorithm fails to estimate a valid cutoff value (*status* = 1), then the fifth largest value in the input sample is returned.

**SVRTUTIL\_PERCENTILE:**

This function computes the specified empirical percentile given the EDF estimates.

- *Type:* Function
- *Signature:* SVRTUTIL\_PERCENTILE(p, n, x{ \* }, F{ \* }, Ftype)
- *Argument Description:*

p	Desired percentile. The value must be in the interval (0,1). The function returns the 100p <sup>th</sup> percentile.
n	Dimension of the x and F input arrays.
x{ * }	Input numeric array of dimension n that contains values of the random variable observed in the sample. These values are sorted in nondecreasing order.
F{ * }	Input numeric array of dimension n in which each F[i] contains the EDF estimate for x[i]. These values must be sorted in nondecreasing order.
Ftype	Type of the empirical estimate that is stored in the x and F arrays. See the section “ <a href="#">Supplying EDF Estimates to Functions and Subroutines</a> ” on page 300 for definition of types.
- *Return value:* The 100p<sup>th</sup> percentile of the input sample.

The method used to compute the percentile depends on the type of the EDF estimate (Ftype argument).

- Ftype = 1      Smoothed empirical estimates are computed using the method described in Klugman, Panjer, and Willmot (1998). Let  $\lfloor x \rfloor$  denote the greatest integer less than or equal to x. Define  $g = \lfloor p(n+1) \rfloor$  and  $h = p(n+1) - g$ . Then the empirical percentile  $\hat{\pi}_p$  is defined as

$$\hat{\pi}_p = (1-h)x[g] + hx[g+1]$$

This method does not work if  $p < 1/(n+1)$  or  $p > n/(n+1)$ . If  $p < 1/(n+1)$ , then the function returns  $\hat{\pi}_p = x[1]/2$ , which assumes that the EDF is 0 in the interval  $[0, x[1])$ . If  $p > n/(n+1)$ , then  $\hat{\pi}_p = x[n]$ .

- Ftype = 2      If  $p < F[1]$ , then  $\hat{\pi}_p = x[1]/2$ , which assumes that the EDF is 0 in the interval  $[0, x[1])$ . If  $|p - F[i]| < \epsilon$  for some value of i and  $i < n$ , then  $\hat{\pi}_p$  is computed as

$$\hat{\pi}_p = \frac{x[i] + x[i+1]}{2}$$

where  $\epsilon$  is a machine-precision constant as returned by the SAS function CONSTANT(‘MACEPS’). If  $F[i-1] < p < F[i]$ , then  $\hat{\pi}_p$  is computed as

$$\hat{\pi}_p = x[i]$$

If  $p \geq F[n]$ , then  $\hat{\pi}_p = x[n]$ .

**SVRTUTIL\_RAWMOMENTS:**

This subroutine computes the raw moments of a sample.

- *Type:* Subroutine
- *Signature:* SVRTUTIL\_RAWMOMENTS(n, x{ \* }, nx{ \* }, nRaw, raw{ \* })

- *Argument Description:*

$n$	Dimension of the $x$ and $nx$ input arrays.
$x\{*\}$	Input numeric array of dimension $n$ that contains distinct values of the random variable that are observed in the sample.
$nx\{*\}$	Input numeric array of dimension $n$ in which each $nx[i]$ contains the number of observations in the sample that have the value $x[i]$ .
$nRaw$	Desired number of raw moments. The output array <i>raw</i> contains the first $nRaw$ raw moments.
$raw\{*\}$	Output array of raw moments. The $k$ th element in the array ( $raw\{k\}$ ) contains the $k$ th raw moment, where $1 \leq k \leq nRaw$ .

- *Return value:* Numeric array *raw* that contains the first  $nRaw$  raw moments. The array contains missing values if the sample has no observations (that is, if all the values in the  $nx$  array add up to zero).

## SVRTUTIL\_SORT:

This function sorts the given array of numeric values in an ascending or descending order.

- *Type:* Subroutine
- *Signature:* SVRTUTIL\_SORT( $n$ ,  $x\{*\}$ , *flag*)
- *Argument Description:*

$n$	Dimension of the input array $x$ .
$x\{*\}$	Numeric array that contains the values to be sorted at input. The subroutine uses the same array to return the sorted values.
<i>flag</i>	A numeric value that controls the sort order. If <i>flag</i> is 0, then the values are sorted in an ascending order. If <i>flag</i> has any value other than 0, then the values are sorted in descending order.

- *Return value:* Numeric array  $x$ , which is sorted in place (that is, the sorted array is stored in the same storage area occupied by the input array  $x$ ).

You can use the following predefined functions when you use the FCMP procedure to define functions and subroutines. They are summarized here for your information. For more information, see the FCMP procedure documentation in *Base SAS Procedures Guide*.

## INVCDF:

This function computes the quantile from any continuous probability distribution by numerically inverting the CDF of that distribution. You need to specify the CDF function of the distribution, the values of its parameters, and the cumulative probability to compute the quantile.

## LIMMOMENT:

This function computes the limited moment of order  $k$  with upper limit  $u$  for any continuous probability distribution. The limited moment is defined as

$$\begin{aligned}
 E[(X \wedge u)^k] &= \int_0^u x^k f(x) dx + \int_u^\infty u^k f(x) dx \\
 &= \int_0^u x^k f(x) dx + u^k (1 - F(u))
 \end{aligned}$$

where  $f(x)$  and  $F(x)$  denote the PDF and the CDF of the distribution, respectively. The LIMMOMENT function uses the following alternate definition, which can be derived using integration-by-parts:

$$E[(X \wedge u)^k] = k \int_0^u (1 - F(x))x^{k-1} dx$$

You need to specify the CDF function of the distribution, the values of its parameters, and the values of  $k$  and  $u$  to compute the limited moment.

---

## Scoring Functions (Experimental)

Scoring refers to the act of evaluating a distribution function, such as LOGPDF, SDF, or QUANTILE, on an observation by using the fitted parameter estimates of that distribution. You can do scoring in a DATA step by using the `OUTEST=` data set that you create with PROC HPSEVERITY. However, that approach requires some cumbersome programming. In order to simplify the scoring process, you can specify that PROC HPSEVERITY create scoring functions for each fitted distribution.

As an example, assume that you have fitted the Pareto distribution by using PROC HPSEVERITY and that it converges. Further assume that you want to use the fitted distribution to evaluate the probability of observing a loss value greater than some set of regulatory limits  $\{L\}$  that are encoded in a data set. You can simplify this scoring process as follows. First, in the PROC HPSEVERITY step that fits your distributions, you create the scoring functions library by specifying the `OUTSCORELIB` statement as illustrated in the following steps:

```
proc hpseverity data=input;
    loss lossclaim;
    dist pareto;
    outscorelib outlib=sasuser.fitdist;
run;
```

Upon successful completion, if the Pareto distribution model has converged, then the `Sasuser.Fitdist` library contains the `SEV_SDF` scoring function in addition to other scoring functions, such as `SEV_PDF`, `SEV_LOGPDF`, and so on. Further, PROC HPSEVERITY also sets the `CMPLIB` system option to include the `Sasuser.Fitdist` library. If the set of limits  $\{L\}$  is recorded in the variable `Limit` in the scoring data set `Work.Limits`, then you can submit the following DATA step to compute the probability of seeing a loss greater than each limit:

```
data prob;
    set work.limits;
    exceedance_probability = sev_sdf(limit);
run;
```

Without the use of scoring functions, you can still perform this scoring task, but the DATA step that you need to write to accomplish it becomes more complicated and less flexible. For example, you would need to read the parameter estimates from some output created by PROC HPSEVERITY. To do that, you would need to know the parameter names, which are different for different distributions; this in turn would require you to write a specific DATA step for each distribution or to write a SAS macro. With the use of scoring functions, you can accomplish that task much more easily.

If you fit multiple distributions, then you can specify the `COMMONPACKAGE` option in the `OUTSCORELIB` statement as follows:

```
proc hpseverity data=input;
  loss lossclaim;
  dist exp pareto weibull;
  outscorelib outlib=sasuser.fitdist commonpackage;
run;
```

The preceding step creates scoring functions such as *SEV\_SDF\_Exp*, *SEV\_SDF\_Pareto*, and *SEV\_SDF\_Weibull*. You can use them to compare the probabilities of exceeding the limit for different distributions by using the following DATA step:

```
data prob;
  set work.limits;
  exceedance_exp = sev_sdf_exp(limit);
  exceedance_pareto = sev_sdf_pareto(limit);
  exceedance_weibull = sev_sdf_weibull(limit);
run;
```

## Formal Description

PROC HPSEVERITY creates a scoring function for each distribution function. A distribution function is defined as any function named *dist\_suffix*, where *dist* is the name of a distribution that you specify in the DIST statement and the function's signature is identical to the signature of the required distribution function such as *dist\_CDF* or *dist\_LOGCDF*. For example, for the function 'FOO\_BAR' to be a distribution function, you must specify the distribution 'FOO' in the DIST statement and you must define 'FOO\_BAR' in the following manner if the distribution 'FOO' has parameters named 'P1' and 'P2':

```
function FOO_BAR(y, P1, P2);
  /* Code to compute BAR by using y, P1, and P2 */
  R = <computed BAR>;
  return (R);
endsub;
```

For more information about the signature that defines a distribution function, see the description of the *dist\_CDF* function in the section “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308.

The name and package of the scoring function of a distribution function depend on whether you specify the `COMMONPACKAGE` option in the `OUTSCORELIB` statement.

When you do not specify the `COMMONPACKAGE` option, the scoring function that corresponds to the distribution function *dist\_suffix* is named *SEV\_suffix*, where *SEV\_* is the standard prefix of all scoring functions. The scoring function is created in a package named *dist*. Each scoring function accepts only one argument, the value of the loss variable, and returns the same value as the value returned by the corresponding distribution function for the final estimates of the distribution's parameters. For example, for the preceding 'FOO\_BAR' distribution function, the scoring function named 'SEV\_BAR' is created in the package named 'FOO' and 'SEV\_BAR' has the following signature:

```
function SEV_BAR(y);
  /* returns value of FOO_BAR for the supplied value
    of y and fitted values of P1, P2 */
endsub;
```

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, then the scoring function that corresponds to the distribution function *dist\_suffix* is named *SEV\_suffix\_dist*, where *SEV\_* is the standard prefix of all scoring functions. The scoring function is created in a package named *sevfit*. For example, for the preceding ‘FOO\_BAR’ distribution function, if you specify the COMMONPACKAGE option, the scoring function named ‘SEV\_BAR\_FOO’ is created in the *sevfit* package and ‘SEV\_BAR\_FOO’ has the following signature:

```
function SEV_BAR_FOO(y);
  /* returns value of FOO_BAR for the supplied value
    of y and fitted values of P1, P2 */
endsub;
```

## Scoring Functions for the Scale Regression Model

If you use the SCALEMODEL statement to specify a scale regression model, then the estimate of the scale parameter or the log-transformed scale parameter of the distribution depends on the values of the regressors. So PROC HPSEVERITY creates a scoring function that has the following signature, where  $x\{*\}$  represents the array of regressors:

```
function SEV_BAR(y, x{*});
  /* returns value of FOO_BAR for the supplied value of x and fitted values of P1, P2 */
endsub;
```

As an illustration of using this form, assume that you submit the following PROC HPSEVERITY step to create the scoring library Sasuser.Scalescore:

```
proc hpseverity data=input;
  loss lossclaim;
  scalemodel x1-x3;
  dist pareto;
  outscorelib outlib=sasuser.scalescore;
run;
```

Your scoring data set must contain all the regressors that you specify in the SCALEMODEL statement. You can submit the following DATA step to score observations by using the scale regression model:

```
data prob;
  array regvals{*} x1-x3;
  set work.limits;
  exceedance_probability = sev_sdf(limit, regvals);
run;
```

PROC HPSEVERITY creates two utility functions, SEV\_NUMREG and SEV\_REGNAME, in the OUTLIB= library that return the number of regressors and name of a given regressor, respectively. They are described in detail in the next section. These utility functions are useful when you do not have easy access to the regressor names in the SCALEMODEL statement.

You can use the utility functions as follows:

```
data prob;
  array regvals{10} _temporary_;
  set work.limits;
  do i = 1 to sev_numreg();
    regvals(i) = input(vvaluex(sev_regname(i)), best12.);
  end;
  exceedance_probability = sev_sdf(limit, regvals);
run;
```

The dimension of the regressor values array that you supply to the scoring function must be equal to  $K + L$ , where  $K$  is the number of regressors that you specify in the SCALEMODEL statement irrespective of whether PROC HPSEVERITY deems any of those regressors to be redundant.  $L$  is 1 if you specify an **OFFSET=** variable in the SCALEMODEL statement, and 0 otherwise.

## Utility Functions and Subroutines in the OUTLIB= Library

In addition to creating the scoring functions for all distribution functions, PROC HPSEVERITY creates the following utility functions and subroutines in the OUTLIB= library.

### SEV\_NUMPARM | SEV\_NUMPARM\_dist

is a function that returns the number of distribution parameters and has the following signature:

- *Type:* Function
- *Number of arguments:* 0
- *Sequence and type of arguments:* Not applicable
- *Return value:* Numeric value that contains the number of distribution parameters

If you do not specify the COMMONPACKAGE option in the OUTSCORELIB statement, then a function named SEV\_NUMPARM is created in the package of each distribution. Here is a sample structure of the code that PROC HPSEVERITY uses to define the function:

```
function SEV_NUMPARM();
  n = <number of distribution parameters>;
  return (n);
endsub;
```

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, then for each distribution *dist*, the function named SEV\_NUMPARM\_dist is created in the *sevfit* package. SEV\_NUMPARM\_dist has the same structure as the SEV\_NUMPARM function that is described previously.

### SEV\_PARMEST | SEV\_PARMEST\_dist

is a subroutine that returns the estimate and standard error of a specified distribution parameter and has the following signature:

- *Type:* Subroutine
- *Number of arguments:* 3



- *Sequence and type of arguments:*

*index* specifies the numeric value of the index of the distribution parameter for which you want the information. The value of *index* must be in the interval  $[1, m]$ , where  $m$  is the number of parameters in the distribution to which this subroutine belongs.

*est* specifies the output argument that returns the estimate of the requested parameter.

*stderr* specifies the output argument that returns the standard error of the requested parameter.

- *Return value:* Estimate and standard error of the requested distribution parameter that are returned in the output arguments *est* and *stderr*, respectively

If you do not specify the COMMONPACKAGE option in the OUTSCORELIB statement, then a subroutine named SEV\_PARMEST is created in the package of each distribution. Here is a sample structure of the code that PROC HPSEVERITY uses to define the subroutine:

```
subroutine SEV_PARMEST(index, est, stderr);
  outargs est, stderr;
  est = <value of the estimate for the distribution parameter
        at position 'index'>;
  stderr = <value of the standard error for distribution parameter
           at position 'index'>;
endsub;
```

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, then for each distribution *dist*, the subroutine named SEV\_PARMEST\_*dist* is created in the *sevfit* package. SEV\_PARMEST\_*dist* has the same structure as the SEV\_PARMEST subroutine that is described previously.

If you use the SCALEMODEL statement to specify a scale regression model, then for *index*=1, the returned estimates are of  $\theta_0$ , the base value of the scale parameter, or  $\log(\theta_0)$  if the distribution has a log-scale parameter. For more information about  $\theta_0$ , see the section “[Estimating Regression Effects](#)” on page 290.

### SEV\_PARMNAME | SEV\_PARMNAME\_*dist*

is a function that returns the name of a specified distribution parameter and has the following signature:

- *Type:* Function
- *Number of arguments:* 1
- *Sequence and type of arguments:*

*index* specifies the numeric value of the index of the distribution parameter for which you want the information. The value of *index* must be in the interval  $[1, m]$ , where  $m$  is the number of parameters in the distribution to which this function belongs.

- *Return value:* Character value that contains the name of the distribution parameter that appears at the position *index* in the distribution’s definition

If you do not specify the COMMONPACKAGE option in the OUTSCORELIB statement, then a function named SEV\_PARMNAME is created in the package of each distribution.

Here is a sample structure of the code that PROC HPSEVERITY uses to define the function:

```
function SEV_PARMNAME(index) $32;
    name = <name of the distribution parameter at position 'index'>;
    return (name);
endsub;
```

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, then for each distribution *dist*, a function named SEV\_PARMNAME\_*dist* is created in the *sevfit* package. SEV\_PARMNAME\_*dist* has the same structure as the SEV\_PARMNAME function that is described previously.

If you use the SCALEMODEL statement to specify a scale regression model, then the following helper functions and subroutines are also created in the OUTLIB= library.

### SEV\_NUMREG

is a function that returns the number of regressors and has the following signature:

- *Type*: Function
- *Number of arguments*: 0
- *Sequence and type of arguments*: Not applicable
- *Return value*: Numeric value that contains the number of regressors that you specify in the SCALEMODEL statement. If you specify an OFFSET= variable in the SCALEMODEL statement, then the returned value is equal to 1 plus the number of regressors that you specify in the SCALEMODEL statement.

Here is a sample structure of the code that PROC HPSEVERITY uses to define the function:

```
function SEV_NUMREG();
    m = <number of regressors>;
    if (<offset variable is specified>) then m = m + 1;
    return (m);
endsub;
```

This function does not depend on any distribution, so it is always created in the *sevfit* package.

### SEV\_REGEST | SEV\_REGEST\_*dist*

is a subroutine that returns the estimate and standard error of a specified regression parameter and has the following signature:

- *Type*: Subroutine
- *Number of arguments*: 3
- *Sequence and type of arguments*:

*index* specifies the numeric value of the index of the regression parameter for which you want the information. The value of *index* must be in the interval  $[1, K]$ , where  $K$  is the number of regressors as returned by the SEV\_NUMREG function. If you specify an OFFSET= variable in the SCALEMODEL statement, then an *index* value of  $K$  corresponds to the offset variable.

*est* specifies the output argument that returns the estimate of the requested regression parameter.

*stderr* specifies the output argument that returns the standard error of the requested regression parameter.

- *Return value:* Estimate and standard error of the requested regression parameter that are returned in the output arguments *est* and *stderr*, respectively

If you do not specify the COMMONPACKAGE option in the OUTSCORELIB statement, then a subroutine named SEV\_REGEST is created in the package of each distribution. Here is a sample structure of the code that PROC HPSEVERITY uses to define the subroutine:

```
subroutine SEV_REGEST(index, est, stderr);
  outargs est, stderr;
  est = <value of the estimate for the regression parameter
        at position 'index'>;
  stderr = <value of the standard error for regression parameter
           at position 'index'>;
endsub;
```

If you specify the COMMONPACKAGE option in the OUTSCORELIB statement, then for each distribution *dist*, the subroutine named SEV\_REGEST\_*dist* is created in the *sevfit* package. SEV\_REGEST\_*dist* has the same structure as the SEV\_REGEST subroutine that is described previously.

If the regressor that corresponds to the specified *index* value is a redundant regressor, the returned values of both *est* and *stderr* are equal to the special missing value of .R. If you specify an OFFSET= variable in the SCALEMODEL statement and if the *index* value corresponds to the offset variable — that is, it is equal to the value that the SEV\_NUMREG function returns — then the returned value of *est* is equal to 1 and the returned value of *stderr* is equal to the special missing value of .F.

## SEV\_REGNAME

is a function that returns the name of a specified regressor and has the following signature:

- *Type:* Function
- *Number of arguments:* 1
- *Sequence and type of arguments:*

*index* specifies the numeric value of the index of the regressor for which you want the name. The value of *index* must be in the interval  $[1, K]$ , where  $K$  is the number of regressors as returned by the SEV\_NUMREG function. If you specify an OFFSET= variable in the SCALEMODEL statement, then an *index* value of  $K$  corresponds to the offset variable.

- *Return value:* Character value that contains the name of the regressor that appears at the position *index* in the SCALEMODEL statement. If you specify an OFFSET= variable in the SCALEMODEL statement, then for an *index* value of  $K$ , the returned value contains the name of the offset variable.

Here is a sample structure of the code that PROC HPSEVERITY uses to define the function:

```
function SEV_REGNAME(index) $32;
    name = <name of regressor at position 'index'>;
    return (name);
endsub;
```

This function does not depend on any distribution, so it is always created in the *sevfit* package.

---

## Custom Objective Functions

You can use a series of programming statements that use variables in the DATA= data set to assign a value to an objective function symbol. You must specify the objective function symbol by using the **OBJECTIVE=** option in the PROC HPSEVERITY statement.

The objective function can be programmed such that it is applicable to any distribution that is used in the model. For that purpose, PROC HPSEVERITY recognizes the following *keyword* functions in the programming statements:

<code>_PDF_(x)</code>	returns the probability density function (PDF) of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_CDF_(x)</code>	returns the cumulative distribution function (CDF) of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_SDF_(x)</code>	returns the survival distribution function (SDF) of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_LOGPDF_(x)</code>	returns the natural logarithm of the PDF of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_LOGCDF_(x)</code>	returns the natural logarithm of the CDF of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_LOGSDF_(x)</code>	returns the natural logarithm of the SDF of a distribution evaluated at the current value of a data set variable <i>x</i> .
<code>_EDF_(x)</code>	returns the empirical distribution function (EDF) estimate evaluated at the current value of a data set variable <i>x</i> . Internally, PROC HPSEVERITY computes the estimate using the SVRTUTIL_EDF function as described in the section “ <a href="#">Predefined Utility Functions</a> ” on page 320. The EDF estimate that is required by the SVRTUTIL_EDF function is computed by using the response variable values in the current BY group or in the entire input data set if you do not specify the BY statement.
<code>_EMPLIMMOMENT_(k, u)</code>	returns the empirical limited moment of order <i>k</i> evaluated at the current value of a data set variable <i>u</i> that represents the upper limit of the limited moment. The order <i>k</i> can also be a data set variable. Internally, PROC HPSEVERITY computes the moment using the SVRTUTIL_EMPLIMMOMENT function as described in the section “ <a href="#">Predefined Utility Functions</a> ” on page 320. The EDF estimate that is required by the SVRTUTIL_EMPLIMMOMENT function is computed by using the response variable values in the current BY group or in the entire input data set if you do not specify the BY statement.

`_LIMMOMENT_(k, u)`

returns the limited moment of order  $k$  evaluated at the current value of a data set variable  $u$  that represents the upper limit of the limited moment. The order  $k$  can be a data set variable or a constant. Internally, for each candidate distribution, PROC HPSEVERITY computes the moment using the LIMMOMENT function as described in the section “[Predefined Utility Functions](#)” on page 320.

All the preceding functions are right-hand side functions. They act as placeholders for distribution-specific functions, with the exception of `_EDF_` and `_EMPLIMMOMENT_` functions. As an example, let the data set `Work.Test` contain a response variable  $Y$  and a left-truncation threshold variable  $T$ . The following statements use the values in this data set to fit a model with distribution  $D$  such that the parameters of the model minimize the value of the objective function symbol `MYOBJ`:

```
options cmplib=(work.mydist);
proc hpseverity data=work.test objective=myobj;
  loss y / lt=t;

  myobj = -_LOGPDF_(y);
  if (not(missing(t))) then
    myobj = myobj + log(1-_CDF_(t));

  dist d;
run;
```

The symbol `MYOBJ` is designated as an objective function symbol by using the `OBJECTIVE=` option in the PROC HPSEVERITY statement. The response variable  $Y$  and left-truncation variable  $T$  are specified in the LOSS statement. The distribution  $D$  is specified in the DIST statement. The remaining statements constitute a program that computes the value of the `MYOBJ` symbol.

Let the distribution  $D$  have parameters  $P1$  and  $P2$ . In order to estimate the model for this distribution, PROC HPSEVERITY internally converts the generic program to the following program specific to distribution  $D$ :

```
myobj = -D_LOGPDF(y, p1, p2);
if (not(missing(t))) then
  myobj = myobj + log(1-D_CDF(t, p1, p2));
```

Note that the generic keyword functions `_LOGPDF_` and `_CDF_` have been replaced with distribution-specific functions `D_LOGPDF` and `D_CDF`, respectively, with appropriate distribution parameters. The `D_LOGPDF` and `D_CDF` functions must have been defined previously and are assumed to be available in the `Work.Mydist` library that you specify in the `CMPLIB=` option.

The program is executed for each observation in `Work.Test` to compute the value of `MYOBJ` by using the values of variables  $Y$  and  $T$  in that observation and internally computed values of the model parameters  $P1$  and  $P2$ . The values of `MYOBJ` are then added over all the observations of the data set or over all the observations of the current BY group if you specify the BY statement. The resulting aggregate value is the value of the objective function, and it is supplied to the optimizer. If the optimizer requires derivatives of the objective function, then PROC HPSEVERITY automatically differentiates `MYOBJ` with respect to the parameters  $P1$  and  $P2$ . The optimizer iterates over various combinations of the values of parameters  $P1$  and  $P2$ , each time computing a new value of the objective function and the needed derivatives of it, until it finds a combination that minimizes the objective function.

Note the following points when you define your own program to compute the custom objective function:

- The value of the objective function is always minimized by PROC HPSEVERITY. If you want to maximize the value of a certain objective, then add a statement that assigns the negated value of the maximization objective to the objective function symbol that you specify in the **OBJECTIVE=** option. Minimization of the negated objective is equivalent to the maximization of the original objective.
- The contributions of individual observations are always added to compute the overall objective function in a given iteration of the optimizer. If you specify the **WEIGHT** statement, then the contribution of each observation is weighted by multiplying it with the normalized value of the weight variable for that observation.
- If you are fitting multiple distributions in one PROC HPSEVERITY step and use any of the keyword functions in your program, then it is recommended that you do not explicitly use the parameters of any of the specified distributions in your programming statements.
- If you use a specific keyword function in your programming statements, then the corresponding distribution functions must be defined in a library that you specify in the **CMPLIB=** system option or in `Sashelp.Svrtlist`, the predefined functions library. In the preceding example, it is assumed that the functions `D_LOGPDF` and `D_CDF` are defined in the `Work.Mydist` library that is specified in the **CMPLIB=** option.
- You can use most DATA step statements and functions in your program. The DATA step file and the data set I/O statements (for example, **INPUT**, **FILE**, **SET**, and **MERGE**) are not available. However, some functionality of the **PUT** statement is supported. See the section “PROC FCMP and DATA Step Differences” in *Base SAS Procedures Guide* for more information. In addition to the differences listed in that section, the following differences exist:
  - Only numeric-valued variables can be used in PROC HPSEVERITY programming statements. This restriction also implies that you cannot use SAS functions or call routines that require character-valued arguments, unless you pass those arguments as constant (literal) strings or characters.
  - You cannot use functions that create lagged versions of a variable in PROC HPSEVERITY programming statements. If you need lagged versions, then you can use a DATA step prior to the PROC HPSEVERITY step to add those versions to the input data set.
- When coding your programming statements, avoid defining variables that begin with an underscore (`_`), because they might conflict with internal variables created by PROC HPSEVERITY.

## Custom Objective Functions and Regression Effects

If you specify regressors by using the **SCALEMODEL** statement, then PROC HPSEVERITY automatically adds a statement prior to your programming statements to compute the value of the scale parameter or the log-transformed scale parameter of the distribution using the values of the regression variables and internally created regression parameters. For example, if you specify three regressors `x1`, `x2`, and `x3` in the **SCALEMODEL** statement, then for a model that contains the distribution `D` with scale parameter `S`, PROC HPSEVERITY prepends your program with a statement that is equivalent to the following statement:

```
S = _SEVTHETA0 * exp(_SEVBETA1 * x1 + _SEVBETA2 * x2 + _SEVBETA3 * x3);
```

If a model contains a distribution D1 with a log-transformed scale parameter M, PROC HPSEVERITY prepends your program with a statement that is equivalent to the following statement:

$$M = \_SEVTHETA0 + \_SEVBETA1 * x1 + \_SEVBETA2 * x2 + \_SEVBETA3 * x3;$$

The `\_SEVTHETA0`, `\_SEVBETA1`, `\_SEVBETA2`, and `\_SEVBETA3` are the internal regression parameters associated with the intercept and the regressors `x1`, `x2`, and `x3`, respectively.

Since the names of the internal regression parameters start with a prefix `\_SEV`, if you use a variable in your program with a name that begins with `\_SEV`, then PROC HPSEVERITY writes an error message to the SAS log and stops processing.

---

## Input Data Sets

PROC HPSEVERITY accepts `DATA=` and `INEST=` data sets as input data sets. This section details the information they are expected to contain.

### DATA= Data Set

The `DATA=` data set is expected to contain the values of the analysis variables that you specify in the [LOSS statement](#) and the [SCALEMODEL statement](#).

If you specify the `BY` statement, then the `DATA=` data set must contain all the `BY` variables that you specify in the `BY` statement and the data set must be sorted by the `BY` variables unless you specify the `NOTSORTED` option in the `BY` statement.

### INEST= Data Set

The `INEST=` data set is expected to contain the initial values of the parameters for the parameter estimation process.

If you specify the `BY` statement, then the `INEST=` data set must contain all the `BY` variables that you specify in the `BY` statement. If you do not specify the `NOTSORTED` option in the `BY` statement, then the `INEST=` data set must be sorted by the `BY` variables. However, it is not required to contain all the `BY` groups present in the `DATA=` data set. For the `BY` groups that are not present in the `INEST=` data set, the default parameter initialization method is used. If you specify the `NOTSORTED` option in the `BY` statement, then the `INEST=` data set must contain all the `BY` groups that are present in the `DATA=` data set and they must appear in the same order as they appear in the `DATA=` data set.

In addition to any variables that you specify in the `BY` statement, the data set must contain the following variables:

<code>\_MODEL\_</code>	identifying name of the distribution for which the estimates are provided.
<code>\_TYPE\_</code>	type of the estimate. The value of this variable must be <code>EST</code> for an observation to be valid.
<code>&lt;Parameter 1&gt; ... &lt;Parameter M&gt;</code>	<i>M</i> variables, named after the parameters of all candidate distributions, that contain initial values of the respective parameters. <i>M</i> is the cardinality of the union of parameter name sets from all candidate distributions. In an observation, estimates are read only

from variables for parameters that correspond to the distribution that is indicated by the `_MODEL_` variable.

If you specify a missing value for some parameters, then default initial values are used unless the parameter is initialized by using the `INIT=` option in the `DIST` statement. If you want to use the `dist_PARMINIT` subroutine for initializing the parameters of a model, then you should either not specify the model in the `INEST=` data set or specify missing values for all the distribution parameters in the `INEST=` data set and not use the `INIT=` option in the `DIST` statement.

If you specify regressors, then the initial value that you provide for the first parameter of each distribution must be the base value of the scale or log-transformed scale parameter. For more information, see the section “[Estimating Regression Effects](#)” on page 290.

<Regressor 1> ... <Regressor K>

If you specify  $K$  regressors in the `SCALEMODEL` statement, then the `INEST=` data set must contain  $K$  variables that are named for each regressor. The variables contain initial values of the respective regression coefficients. If a regressor is linearly dependent on other regressors for a given BY group, then you can indicate this by providing a special missing value of `.R` for the respective variable. In a given BY group, if you mark a variable as linearly dependent for one model, then you must mark that variable as linearly dependent for all the models. Similarly, in a given BY group, if you do not mark a variable as linearly dependent for one model, then you must not mark that variable as linearly dependent for all the models.

---

## Output Data Sets

PROC HPSEVERITY writes the `OUTEST=`, `OUTMODELINFO=`, and `OUTSTAT=` data sets when requested by their respective options. The data sets and their contents are described in the following sections.

### OUTEST= Data Set

The `OUTEST=` data set records the estimates of the model parameters. It also contains estimates of their standard errors and optionally their covariance structure. If you specify BY variables, then the data are organized in BY groups and the data set contains variables that you specify in the BY statement.

If you do not specify the `COVOUT` option, then the data set contains the following variables:

<code>_MODEL_</code>	identifying name of the distribution model. The observation contains information about this distribution.				
<code>_TYPE_</code>	type of the estimates reported in this observation. It can take one of the following two values: <table data-bbox="422 1690 1128 1780"> <tr> <td><code>EST</code></td><td>point estimates of model parameters</td></tr> <tr> <td><code>STDERR</code></td><td>standard error estimates of model parameters</td></tr> </table>	<code>EST</code>	point estimates of model parameters	<code>STDERR</code>	standard error estimates of model parameters
<code>EST</code>	point estimates of model parameters				
<code>STDERR</code>	standard error estimates of model parameters				
<code>_STATUS_</code>	status of the reported estimates. The possible values are listed in the section “ <a href="#">_STATUS_ Variable Values</a> ” on page 339.				



## &lt;Parameter 1&gt; ... &lt;Parameter M&gt;

$M$  variables, named after the parameters of all candidate distributions, that contain estimates of the respective parameters.  $M$  is the cardinality of the union of parameter name sets from all candidate distributions. In an observation, estimates are populated only for parameters that correspond to the distribution that is indicated by the `_MODEL_` variable. If `_TYPE_` is EST, then the estimates are missing if the model does not converge. If `_TYPE_` is STDERR, then the estimates are missing if covariance estimates cannot be obtained.

If you specify regressors, then the estimate that is reported for the first parameter of each distribution is the estimate of the base value of the scale or log-transformed scale parameter. For more information, see the section “[Estimating Regression Effects](#)” on page 290.

## &lt;Regressor 1&gt; ... &lt;Regressor K&gt;

If you specify  $K$  regressors in the [SCALEMODEL](#) statement, then the OUTEST= data set contains  $K$  variables that are named for each regressor. The variables contain estimates for their respective regression coefficients. If a regressor is deemed to be linearly dependent on other regressors for a given BY group, then a warning message is written to the SAS log and a special missing value of .R is written in the respective variable. If `_TYPE_` is EST, then the estimates are missing if the model does not converge. If `_TYPE_` is STDERR, then the estimates are missing if covariance estimates cannot be obtained.

## &lt;Offset Variable&gt;

If you specify an OFFSET= variable in the SCALEMODEL statement, then the OUTEST= data set contains a variable that is named after the offset variable. If `_TYPE_` is EST, then the value of this variable is 1. If `_TYPE_` is STDERR, then the value of this variable is a special missing value of .F.

If you specify the COVOUT option in the PROC HPSEVERITY statement, then the OUTEST= data set contains additional observations that contain the estimates of the covariance structure. Given the symmetric nature of the covariance structure, only the lower triangular portion is reported. In addition to the variables listed and described previously, the data set contains the following variables that are either new or have a modified description:

<code>_TYPE_</code>	type of the estimates reported in this observation. For observations that contain rows of the covariance structure, the value is COV.
<code>_STATUS_</code>	status of the reported estimates. For observations that contain rows of the covariance structure, the status is 0 if covariance estimation was successful. If estimation fails, the status is 1 and a single observation is reported with <code>_TYPE_=COV</code> and missing values for all the parameter variables.
<code>_NAME_</code>	name of the parameter for the row of covariance matrix reported in the current observation.

**OUTMODELINFO= Data Set**

The OUTMODELINFO= data set records the information about each candidate distribution that you specify in the DIST statement. It contains the following variables:

<code>_MODEL_</code>	identifying name of the distribution model. The observation contains information about this distribution.
<code>_DEPVAR_</code>	name of the loss variable.
<code>_DESCRIPTION_</code>	descriptive name of the model. This has a nonmissing value only if the DESCRIPTION function has been defined for this model.
<code>_VALID_</code>	validity of the distribution definition. This has a value of 1 for valid definitions and a value of 0 for invalid definitions. If the definition is invalid, then PROC HPSEVERITY writes the reason for invalidity to the SAS log.
<code>_PARMNAME1 ... _PARMNAMEM</code>	$M$ variables that contain names of parameters of the distribution model, where $M$ is the maximum number of parameters across all the specified distribution models. For a given distribution with $m$ parameters, values of variables <code>_PARMNAME<math>j</math></code> ( $j > m$ ) are missing.

## OUTSTAT= Data Set

The OUTSTAT= data set records statistics of fit and model selection information. If you specify BY variables, then the data are organized in BY groups and the data set contains variables that you specify in the BY statement. The data set contains the following variables:

<code>_MODEL_</code>	identifying name of the distribution model. The observation contains information about this distribution.
<code>_NMODELPRM_</code>	number of parameters in the distribution.
<code>_NESTPRM_</code>	number of estimated parameters. This includes the regression parameters, if you specify any regressors.
<code>_NOBS_</code>	number of nonmissing observations used for parameter estimation.
<code>_STATUS_</code>	status of the parameter estimation process for this model. The possible values are listed in the section “ <a href="#">_STATUS_ Variable Values</a> ” on page 339.
<code>_SELECTED_</code>	indicator of the best distribution model. If the value is 1, then this model is the best model for the current BY group according to the specified model selection criterion. This value is missing if the parameter estimation process does not converge for this model.
Neg2LogLike	value of the log likelihood, multiplied by $-2$ , that is attained at the end of the parameter estimation process. This value is missing if the parameter estimation process does not converge for this model.
AIC	value of the Akaike’s information criterion (AIC) that is attained at the end of the parameter estimation process. This value is missing if the parameter estimation process does not converge for this model.
AICC	value of the corrected Akaike’s information criterion (AICC) that is attained at the end of the parameter estimation process. This value is missing if the parameter estimation process does not converge for this model.
BIC	value of the Schwarz Bayesian information criterion (BIC) that is attained at the end of the parameter estimation process. This value is missing if the parameter estimation process does not converge for this model.

KS	value of the Kolmogorov-Smirnov (KS) statistic that is attained at the end of the parameter estimation process. This value is missing if parameter estimation process does not converge for this model.
AD	value of the Anderson-Darling (AD) statistic that is attained at the end of the parameter estimation process. This value is missing if parameter estimation process does not converge for this model.
CVM	value of the Cramér-von Mises (CvM) statistic that is attained at the end of the parameter estimation process. This value is missing if parameter estimation process does not converge for this model.

## **\_STATUS\_ Variable Values**

The `_STATUS_` variable in the `OUTEST=` and `OUTSTAT=` data sets contains a value that indicates the status of the parameter estimation process for the respective distribution model. The variable can take the following values in the `OUTEST=` data set for `_TYPE_=EST` observations and in the `OUTSTAT=` data set:

- 0      The parameter estimation process converged for this model.
- 301    The parameter estimation process might not have converged for this model because there is no improvement in the objective function value. This might indicate that the initial values of the parameters are optimal, or you can try different convergence criteria in the `NLOPTIONS` statement.
- 302    The parameter estimation process might not have converged for this model because the number of iterations exceeded the maximum allowed value. You can try setting a larger value for the `MAXITER=` options in the `NLOPTIONS` statement.
- 303    The parameter estimation process might not have converged for this model because the number of objective function evaluations exceeded the maximum allowed value. You can try setting a larger value for the `MAXFUNC=` options in the `NLOPTIONS` statement.
- 304    The parameter estimation process might not have converged for this model because the time taken by the process exceeded the maximum allowed value. You can try setting a larger value for the `MAXTIME=` option in the `NLOPTIONS` statement.
- 400    The parameter estimation process did not converge for this model.

The `_STATUS_` variable can take the following values in the `OUTEST=` data set for `_TYPE_=STDERR` and `_TYPE_=COV` observations:

- 0      The covariance and standard error estimates are available and valid.
- 1      The covariance and standard error estimates are not available, because the process of computing covariance estimates failed.

---

## **Displayed Output**

The `HPSEVERITY` procedure optionally produces displayed output by using the Output Delivery System (ODS). All output is controlled by the `PRINT=` option in the `PROC HPSEVERITY` statement. [Table 9.7](#) relates the ODS tables to `PRINT=` options.

**Table 9.7** ODS Tables Produced in PROC HPSEVERITY

ODS Table Name	Description	Option
AllFitStatistics	Statistics of fit for all the distribution models	PRINT=ALLFITSTATS
ConvergenceStatus	Convergence status of parameter estimation process	PRINT=CONVSTATUS
DescStats	Descriptive statistics for the response variable	PRINT=DESCSTATS
DistributionInfo	Distribution information	PRINT=DISTINFO
EstimationDetails	Details of the estimation process for all the distribution models	PRINT=ESTIMATIONDETAILS
InitialValues	Initial parameter values and bounds	PRINT=INITIALVALUES
IterationHistory	Optimization iteration history	PRINT=NLOHISTORY
ModelSelection	Model selection summary	PRINT=SELECTION
OptimizationSummary	Optimization summary	PRINT=NLOSUMMARY
ParameterEstimates	Final parameter estimates	PRINT=ESTIMATES
PerformanceInfo	Execution environment information that pertains to the computational performance	Default
RegDescStats	Descriptive statistics for the regressor variables	PRINT=DESCSTATS
StatisticsOfFit	Statistics of fit	PRINT=STATISTICS
Timing	Timing information for various computational stages of the procedure	DETAILS (PERFORMANCE statement)

### PRINT= Option

If you do not specify the PRINT= option, then by default PROC HPSEVERITY produces ModelSelection, PerformanceInfo, ConvergenceStatus, OptimizationSummary, StatisticsOfFit, and ParameterEstimates ODS tables.

You can specify the following values for the PRINT= option:

#### PRINT=ALLFITSTATS

displays the comparison of all the statistics of fit for all the models in one table. The table does not include the models whose parameter estimation process does not converge. If all the models fail to converge, then this table is not produced. If the table contains more than one model, then the best model according to each statistic is indicated with an asterisk (\*) in that statistic's column.

#### PRINT=CONVSTATUS

displays the convergence status of the parameter estimation process.

**PRINT=DESCSTATS**

displays the descriptive statistics for the response variable. If you specify regressor variables in the SCALEMODEL statement, then a separate table with descriptive statistics of regressors is also displayed.

**PRINT=DISTINFO**

displays the information about all the candidate distribution. It includes the name, the description, the number of distribution parameters, and whether the distribution is valid for the specified modeling task.

**PRINT=ESTIMATES**

displays the final estimates of parameters. The estimates are not displayed for models whose parameter estimation process does not converge.

**PRINT=ESTIMATIONDETAILS**

displays the comparative details of the estimation process that is used to fit each candidate distribution. If you specify the DETAILS option in the PERFORMANCE statement, then this table contains a column that indicates the time taken to estimate each candidate model.

**PRINT=INITIALVALUES**

displays the initial values and bounds used for estimating each model.

**PRINT=NLOHISTORY**

displays the iteration history of the nonlinear optimization process used for estimating the parameters.

**PRINT=NLOSUMMARY**

displays the summary of the nonlinear optimization process used for estimating the parameters.

**PRINT=SELECTION**

displays the model selection table. The table shows the convergence status of each candidate model, and the value of the selection criterion along with an indication of the selected model.

**PRINT=STATISTICS**

displays the statistics of fit for each model. The statistics of fit are not displayed for models whose parameter estimation process does not converge.

## Performance Information

The “Performance Information” table is produced by default. It displays information about the execution mode. For single-machine mode, the table displays the number of threads used. For distributed mode, the table displays the grid mode (symmetric or asymmetric), the number of compute nodes, and the number of threads per node.

If you specify the DETAILS option in the PERFORMANCE statement, the procedure also produces a “Timing” table in which elapsed times (absolute and relative) for the main tasks of the procedure are displayed.

## Examples: HPSEVERITY Procedure

### Example 9.1: Defining a Model for Gaussian Distribution

Suppose you want to fit a distribution model other than one of the predefined ones available to you. Suppose you want to define a model for the Gaussian distribution with the following typical parameterization of the PDF ( $f$ ) and CDF ( $F$ ):

$$f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

$$F(x; \mu, \sigma) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{x - \mu}{\sigma \sqrt{2}}\right)\right)$$

For PROC HPSEVERITY, a *distribution model* consists of a set of functions and subroutines that are defined with the FCMP procedure. Each function and subroutine should be written following certain rules. For more information, see the section “[Defining a Severity Distribution Model with the FCMP Procedure](#)” on page 308.

**NOTE:** The Gaussian distribution is not a commonly used severity distribution. It is used in this example primarily to illustrate the process of defining your own distribution models. Although the distribution has a support over the entire real line, you can fit the distribution with PROC HPSEVERITY only if the input sample contains nonnegative values.

The following SAS statements define a distribution model named NORMAL for the Gaussian distribution. The OUTLIB= option in the PROC FCMP statement stores the compiled versions of the functions and subroutines in the ‘models’ package of the Work.Sevexmpl library. The LIBRARY= option in the PROC FCMP statement enables this PROC FCMP step to use the SVRTUTIL\_RAWMOMENTS utility subroutine that is available in the Sashelp.Svrtdist library. The subroutine is described in the section “[Predefined Utility Functions](#)” on page 320.

```

/*----- Define Normal Distribution with PROC FCMP -----*/
proc fcmp library=sashelp.svrtdist outlib=work.sevexmpl.models;
  function normal_pdf(x,Mu,Sigma);
    /* Mu      : Location */
    /* Sigma   : Standard Deviation */
    return ( exp(-(x-Mu)**2/(2 * Sigma**2)) /
              (Sigma * sqrt(2*constant('PI')) ) );
  endsub;

  function normal_cdf(x,Mu,Sigma);
    /* Mu      : Location */
    /* Sigma   : Standard Deviation */
    z = (x-Mu)/Sigma;
    return (0.5 + 0.5*erf(z/sqrt(2)));
  endsub;

  subroutine normal_parminit(dim, x[*], nx[*], F[*], Ftype, Mu, Sigma);
    outargs Mu, Sigma;
    array m[2] / nosymbols;

```

```

/* Compute estimates by using method of moments */
call svrtutil_rawmoments(dim, x, nx, 2, m);
Mu    = m[1];
Sigma = sqrt(m[2] - m[1]**2);
endsub;

subroutine normal_lowerbounds(Mu, Sigma);
outargs Mu, Sigma;
Mu = .; /* Mu has no lower bound */
Sigma = 0; /* Sigma > 0 */
endsub;
quit;

```

The statements define the two functions required of any distribution model (NORMAL\_PDF and NORMAL\_CDF) and two optional subroutines (NORMAL\_PARMINIT and NORMAL\_LOWERBOUNDS). The name of each function or subroutine must follow a specific structure. It should start with the model's short or identifying name, which is 'NORMAL' in this case, followed by an underscore '\_', followed by a keyword suffix such as 'PDF'. Each function or subroutine has a specific purpose. For more information about all the functions and subroutines that you can define for a distribution model, see the section [“Defining a Severity Distribution Model with the FCMP Procedure”](#) on page 308. Following is the description of each function and subroutine defined in this example:

- The PDF and CDF suffixes define functions that return the probability density function and cumulative distribution function values, respectively, given the values of the random variable and the distribution parameters.
- The PARMINIT suffix defines a subroutine that returns the initial values for the parameters by using the sample data or the empirical distribution function (EDF) estimate computed from it. In this example, the parameters are initialized by using the method of moments. Hence, you do not need to use the EDF estimates, which are available in the F array. The first two raw moments of the Gaussian distribution are as follows:

$$E[x] = \mu, \quad E[x^2] = \mu^2 + \sigma^2$$

Given the sample estimates,  $m_1$  and  $m_2$ , of these two raw moments, you can solve the equations  $E[x] = m_1$  and  $E[x^2] = m_2$  to get the following estimates for the parameters:  $\hat{\mu} = m_1$  and  $\hat{\sigma} = \sqrt{m_2 - m_1^2}$ . The NORMAL\_PARMINIT subroutine implements this solution. It uses the SVRTUTIL\_RAWMOMENTS utility subroutine to compute the first two raw moments.

- The LOWERBOUNDS suffix defines a subroutine that returns the lower bounds on the parameters. PROC HPSEVERITY assumes a default lower bound of 0 for all the parameters when a LOWERBOUNDS subroutine is not defined. For the parameter  $\mu$  ( $Mu$ ), there is no lower bound, so you need to define the NORMAL\_LOWERBOUNDS subroutine. It is recommended that you assign bounds for all the parameters when you define the LOWERBOUNDS subroutine or its counterpart, the UPPERBOUNDS subroutine. Any unassigned value is returned as a missing value, which PROC HPSEVERITY interprets to mean that the parameter is unbounded, and that might not be what you want.

You can now use this distribution model with PROC HPSEVERITY. Let the following DATA step statements simulate a normal sample with  $\mu = 10$  and  $\sigma = 2.5$ :

```

/*----- Simulate a Normal sample -----*/
data testnorm(keep=y);
  call streaminit(12345);
  do i=1 to 100;
    y = rand('NORMAL', 10, 2.5);
    output;
  end;
run;

```

Prior to using your distribution with PROC HPSEVERITY, you must communicate the location of the library that contains the definition of the distribution and the locations of libraries that contain any functions and subroutines used by your distribution model. The following OPTIONS statement sets the CMPLIB= system option to include the FCMP library Work.Sevexmpl in the search path used by PROC HPSEVERITY to find FCMP functions and subroutines.

```

/*--- Set the search path for functions defined with PROC FCMP ---*/
options cmplib=(work.sevexmpl);

```

Now, you are ready to fit the NORMAL distribution model with PROC HPSEVERITY. The following statements fit the model to the values of Y in the Work.Testnorm data set:

```

/*--- Fit models with PROC HPSEVERITY ---*/
proc hpseverity data=testnorm print=all;
  loss y;
  dist Normal;
run;

```

The DIST statement specifies the identifying name of the distribution model, which is 'NORMAL'. Neither is the INEST= option specified in the PROC HPSEVERITY statement nor is the INIT= option specified in the DIST statement. So, PROC HPSEVERITY initializes the parameters by invoking the NORMAL\_PARMINIT subroutine.

Some of the results prepared by the preceding PROC HPSEVERITY step are shown in [Output 9.1.1](#) and [Output 9.1.2](#). The descriptive statistics of variable Y and the “Model Selection” table, which includes just the normal distribution, are shown in [Output 9.1.1](#).

**Output 9.1.1** Summary of Results for Fitting the Normal Distribution

The HPSEVERITY Procedure	
Input Data Set	
Name	WORK.TESTNORM
Descriptive Statistics for y	
Observations	100
Observations Used for Estimation	100
Minimum	3.88249
Maximum	16.00864
Mean	10.02059
Standard Deviation	2.37730



**Output 9.1.1** *continued*

Model Selection			
Distribution	Converged	-2 Log Likelihood	Selected
Normal	Yes	455.97541	Yes

The initial values for the parameters, the optimization summary, and the final parameter estimates are shown in [Output 9.1.2](#). No iterations are required to arrive at the final parameter estimates, which are identical to the initial values. This confirms the fact that the maximum likelihood estimates for the Gaussian distribution are identical to the estimates obtained by the method of moments that was used to initialize the parameters in the NORMAL\_PARMINIT subroutine.

**Output 9.1.2** Details of the Fitted Normal Distribution Model

The HPSEVERITY Procedure				
Normal Distribution				
Distribution Information				
Name	Normal			
Distribution Parameters	2			
Initial Parameter Values and Bounds				
Parameter	Initial Value	Lower Bound	Upper Bound	
Mu	10.02059	-Infty	Infty	
Sigma	2.36538	1.05367E-8	Infty	
Optimization Summary				
Optimization Technique	Trust Region			
Iterations	0			
Function Calls	4			
Log Likelihood	-227.98770			
Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu	10.02059	0.23894	41.94	<.0001
Sigma	2.36538	0.16896	14.00	<.0001

The NORMAL distribution defined and illustrated here has no scale parameter, because all the following inequalities are true:

$$f(x; \mu, \sigma) \neq \frac{1}{\mu} f\left(\frac{x}{\mu}; 1, \sigma\right)$$

$$f(x; \mu, \sigma) \neq \frac{1}{\sigma} f\left(\frac{x}{\sigma}; \mu, 1\right)$$

$$F(x; \mu, \sigma) \neq F\left(\frac{x}{\mu}; 1, \sigma\right)$$

$$F(x; \mu, \sigma) \neq F\left(\frac{x}{\sigma}; \mu, 1\right)$$

This implies that you cannot estimate the effect of regressors on a model for the response variable based on this distribution.

---

### Example 9.2: Defining a Model for the Gaussian Distribution with a Scale Parameter

If you want to estimate the effects of regressors, then the model needs to be parameterized to have a scale parameter. Although this might not be always possible, it is possible for the Gaussian distribution by replacing the location parameter  $\mu$  with another parameter,  $\alpha = \mu/\sigma$ , and defining the PDF ( $f$ ) and the CDF ( $F$ ) as follows:

$$f(x; \sigma, \alpha) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{1}{2}\left(\frac{x}{\sigma} - \alpha\right)^2\right)$$

$$F(x; \sigma, \alpha) = \frac{1}{2} \left(1 + \operatorname{erf}\left(\frac{1}{\sqrt{2}}\left(\frac{x}{\sigma} - \alpha\right)\right)\right)$$

You can verify that  $\sigma$  is the scale parameter, because both of the following equalities are true:

$$f(x; \sigma, \alpha) = \frac{1}{\sigma} f\left(\frac{x}{\sigma}; 1, \alpha\right)$$

$$F(x; \sigma, \alpha) = F\left(\frac{x}{\sigma}; 1, \alpha\right)$$

**NOTE:** The Gaussian distribution is not a commonly used severity distribution. It is used in this example primarily to illustrate the concept of parameterizing a distribution such that it has a scale parameter. Although the distribution has a support over the entire real line, you can fit the distribution with PROC HPSEVERITY only if the input sample contains nonnegative values.

The following statements use the alternate parameterization to define a new model named NORMAL\_S. The definition is stored in the Work.Sevexmpl library.

```

/*----- Define Normal Distribution With Scale Parameter -----*/
proc fcmp library=sashelp.svrtldist outlib=work.sevexmpl.models;
  function normal_s_pdf(x, Sigma, Alpha);
    /* Sigma : Scale & Standard Deviation */
    /* Alpha : Scaled mean */
    return ( exp(-(x/Sigma - Alpha)**2/2) /
              (Sigma * sqrt(2*constant('PI')) ) );
  endsub;

  function normal_s_cdf(x, Sigma, Alpha);
    /* Sigma : Scale & Standard Deviation */
    /* Alpha : Scaled mean */
    z = x/Sigma - Alpha;
    return (0.5 + 0.5*erf(z/sqrt(2)));
  endsub;

  subroutine normal_s_parinit(dim, x[*], nx[*], F[*], Ftype, Sigma, Alpha);
    outargs Sigma, Alpha;
    array m[2] / nosymbols;

    /* Compute estimates by using method of moments */
    call svrtutil_rawmoments(dim, x, nx, 2, m);
    Sigma = sqrt(m[2] - m[1]**2);
    Alpha = m[1]/Sigma;
  endsub;

  subroutine normal_s_lowerbounds(Sigma, Alpha);
    outargs Sigma, Alpha;
    Alpha = .; /* Alpha has no lower bound */
    Sigma = 0; /* Sigma > 0 */
  endsub;
quit;

```

An important point to note is that the scale parameter *Sigma* is the first distribution parameter (after the 'x' argument) listed in the signatures of NORMAL\_S\_PDF and NORMAL\_S\_CDF functions. *Sigma* is also the first distribution parameter listed in the signatures of other subroutines. This is required by PROC HPSEVERITY, so that it can identify which is the scale parameter. When you specify regressor variables, PROC HPSEVERITY checks whether the first parameter of each candidate distribution is a scale parameter (or a log-transformed scale parameter if [dist\\_SCALETRANSFORM](#) subroutine is defined for the distribution with LOG as the transform). If it is not, then an appropriate message is written the SAS log and that distribution is not fitted.

Let the following DATA step statements simulate a sample from the normal distribution where the parameter  $\sigma$  is affected by the regressors as follows:

$$\sigma = \exp(1 + 0.5 X_1 + 0.75 X_3 - 2 X_4 + X_5)$$

The sample is simulated such that the regressor  $X_2$  is linearly dependent on regressors  $X_1$  and  $X_3$ .

```

/*--- Simulate a Normal sample affected by Regressors ---*/
data testnorm_reg(keep=y x1-x5 Sigma);
  array x{*} x1-x5;
  array b{6} _TEMPORARY_ (1 0.5 . 0.75 -2 1);
  call streaminit(34567);
  label y='Normal Response Influenced by Regressors';

  do n = 1 to 100;
    /* simulate regressors */
    do i = 1 to dim(x);
      x(i) = rand('UNIFORM');
    end;
    /* make x2 linearly dependent on x1 */
    x(2) = 5 * x(1);

    /* compute log of the scale parameter */
    logSigma = b(1);
    do i = 1 to dim(x);
      if (i ne 2) then
        logSigma = logSigma + b(i+1) * x(i);
      end;

    Sigma = exp(logSigma);
    y = rand('NORMAL', 25, Sigma);

    output;
  end;
run;

```

The following statements use PROC HPSEVERITY to fit the NORMAL\_S distribution model along with some of the predefined distributions to the simulated sample:

```

/*--- Set the search path for functions defined with PROC FCMP ---*/
options cmplib=(work.sevexmpl);

/*----- Fit models with PROC HPSEVERITY -----*/
proc hpseverity data=testnorm_reg print=all;
  loss y;
  scalemodel x1-x5;
  dist Normal_s burr logn pareto weibull;
run;

```

The “Model Selection” table in [Output 9.2.1](#) indicates that all the models, except the Burr distribution model, have converged. Also, only three models, Normal\_s, Burr, and Weibull, seem to have a good fit for the data. The table that compares all the fit statistics indicates that Normal\_s model is the best according to the likelihood-based statistics.

**Output 9.2.1** Summary of Results for Fitting the Normal Distribution with Regressors

The HPSEVERITY Procedure					
Input Data Set					
Name	WORK.TESTNORM_REG				
Model Selection					
Distribution	Converged	-2 Log Likelihood	Selected		
Normal_s	Yes	603.95786	Yes		
Burr	Maybe	612.81685	No		
Logn	Yes	749.20125	No		
Pareto	Yes	841.07022	No		
Weibull	Yes	612.77496	No		
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Normal_s	603.95786*	615.95786*	616.86108*	631.58888*	1.52388
Burr	612.81685	626.81685	628.03424	645.05304	1.50448*
Logn	749.20125	761.20125	762.10448	776.83227	2.88110
Pareto	841.07022	853.07022	853.97345	868.70124	4.83810
Weibull	612.77496	624.77496	625.67819	640.40598	1.50490
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
Normal_s	4.00152		0.70769		
Burr	3.90072*		0.63399*		
Logn	16.20558		3.04825		
Pareto	31.60568		6.84046		
Weibull	3.90559		0.63458		
Note: The asterisk (*) marks the best model according to each column's criterion.					

This prompts you to further evaluate why the model with Burr distribution has not converged. The initial values, convergence status, and the optimization summary for the Burr distribution are shown in [Output 9.2.2](#). The initial values table indicates that the regressor X2 is redundant, which is expected. More importantly, the convergence status indicates that it requires more than 50 iterations. PROC HPSEVERITY enables you to change several settings of the optimizer by using the [NLOPTIONS](#) statement. In this case, you can increase the limit of 50 on the iterations, change the convergence criterion, or change the technique to something other than the default trust-region technique.

**Output 9.2.2** Details of the Fitted Burr Distribution Model

The HPSEVERITY Procedure			
Burr Distribution			
Distribution Information			
Name	Burr		
Description	Burr Distribution		
Distribution Parameters	3		
Regression Parameters	4		
Initial Parameter Values and Bounds			
Parameter	Initial Value	Lower Bound	Upper Bound
Theta	25.75198	1.05367E-8	Infty
Alpha	2.00000	1.05367E-8	Infty
Gamma	2.00000	1.05367E-8	Infty
x1	0.07345	-709.78271	709.78271
x2	Redundant		
x3	-0.14056	-709.78271	709.78271
x4	0.27064	-709.78271	709.78271
x5	-0.23230	-709.78271	709.78271
Convergence Status			
Needs more than 50 iterations.			
Optimization Summary			
Optimization Technique	Trust Region		
Iterations	50		
Function Calls	137		
Log Likelihood	-306.40842		

The following PROC HPSEVERITY step uses the NLOPTIONS statement to change the convergence criterion and the limits on the iterations and function evaluations, exclude the lognormal and Pareto distributions that have been confirmed previously to fit the data poorly, and exclude the redundant regressor X2 from the model:

```

/*--- Refit and compare models with higher limit on iterations ---*/
proc hpseverity data=testnorm_reg print=all;
  loss y;
  scalemodel x1 x3-x5;
  dist Normal_s burr weibull;
  nloptions absfconv=2.0e-5 maxiter=100 maxfunc=500;
run;

```

The results shown in [Output 9.2.3](#) indicate that the Burr distribution has now converged and that the Burr and Weibull distributions have an almost identical fit for the data. The NORMAL\_S distribution is still the best distribution according to the likelihood-based criteria.

**Output 9.2.3** Summary of Results after Changing Maximum Number of Iterations

The HPSEVERITY Procedure					
Input Data Set					
Name	WORK.TESTNORM_REG				
Model Selection					
Distribution	Converged	-2 Log Likelihood	Selected		
Normal_s	Yes	603.95786	Yes		
Burr	Yes	612.79276	No		
Weibull	Yes	612.77496	No		
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Normal_s	603.95786*	615.95786*	616.86108*	631.58888*	1.52388
Burr	612.79276	626.79276	628.01015	645.02895	1.50472*
Weibull	612.77496	624.77496	625.67819	640.40598	1.50490
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
Normal_s	4.00152		0.70769		
Burr	3.90351*		0.63433*		
Weibull	3.90559		0.63458		
Note: The asterisk (*) marks the best model according to each column's criterion.					

### Example 9.3: Defining a Model for Mixed-Tail Distributions

In some applications, a few severity values tend to be extreme as compared to the typical values. The extreme values represent the worst case scenarios and cannot be discarded as outliers. Instead, their distribution must be modeled to prepare for their occurrences. In such cases, it is often useful to fit one distribution to the non-extreme values and another distribution to the extreme values. The *mixed-tail* distribution mixes two distributions: one for the *body* region, which contains the non-extreme values, and another for the *tail* region, which contains the extreme values. The tail distribution is usually a generalized Pareto distribution (GPD), because it is usually good for modeling the conditional excess severity above a threshold. The body distribution can be any distribution. The following definitions are used in describing a generic formulation of a mixed-tail distribution:

$g(x)$	PDF of the body distribution
$G(x)$	CDF of the body distribution
$h(x)$	PDF of the tail distribution
$H(x)$	CDF of the tail distribution
$\theta$	scale parameter for the body distribution
$\Omega$	set of nonscale parameters for the body distribution
$\xi$	shape parameter for the GPD tail distribution
$x_r$	normalized value of the response variable at which the tail starts
$p_n$	mixing probability

Given these notations, the PDF  $f(x)$  and the CDF  $F(x)$  of the mixed-tail distribution are defined as

$$f(x) = \begin{cases} \frac{p_n}{G(x_b)} g(x) & \text{if } x \leq x_b \\ (1 - p_n) h(x - x_b) & \text{if } x > x_b \end{cases}$$

$$F(x) = \begin{cases} \frac{p_n}{G(x_b)} G(x) & \text{if } x \leq x_b \\ p_n + (1 - p_n) H(x - x_b) & \text{if } x > x_b \end{cases}$$

where  $x_b = \theta x_r$  is the value of the response variable at which the tail starts.



These definitions indicate the following:

- The body distribution is conditional on  $X \leq x_b$ , where  $X$  denotes the random response variable.
- The tail distribution is the generalized Pareto distribution of the  $(X - x_b)$  values.
- The probability that a response variable value belongs to the body is  $p_n$ . Consequently the probability that the value belongs to the tail is  $(1 - p_n)$ .

The parameters of this distribution are  $\theta$ ,  $\Omega$ ,  $\xi$ ,  $x_r$ , and  $p_n$ . The scale of the GPD tail distribution  $\theta_t$  is computed as

$$\theta_t = \frac{G(x_b; \theta, \Omega)}{g(x_b; \theta, \Omega)} \frac{(1 - p_n)}{p_n} = \theta \frac{G(x_r; \theta = 1, \Omega)}{g(x_r; \theta = 1, \Omega)} \frac{(1 - p_n)}{p_n}$$

The parameter  $x_r$  is usually estimated using a tail index estimation algorithm. One such algorithm is the Hill's algorithm (Danielsson et al. 2001), which is implemented by the predefined utility function SVRTUTIL\_HILLCUTOFF available to you in the Sashelp.Svrtdist library. The algorithm and the utility function are described in detail in the section “[Predefined Utility Functions](#)” on page 320. The function computes an estimate of  $x_b$ , which can be used to compute an estimate of  $x_r$  because  $x_r = x_b / \hat{\theta}$ , where  $\hat{\theta}$  is the estimate of the scale parameter of the body distribution.

The parameter  $p_n$  is usually determined by the domain expert based on the fraction of losses that are expected to belong to the tail.

The following SAS statements define the LOGNGPD distribution model for a mixed-tail distribution with the lognormal distribution as the body distribution and GPD as the tail distribution:

```

/*----- Define Lognormal Body-GPD Tail Mixed Distribution -----*/
proc fcmp library=sashelp.svtrdist outlib=work.sevexmpl.models;
  function LOGNGPD_DESCRIPTION() $256;
    length desc $256;
    desc1 = "Lognormal Body-GPD Tail Distribution.";
    desc2 = " Mu, Sigma, and Xi are free parameters.";
    desc3 = " Xr and Pn are constant parameters.";
    desc = desc1 || desc2 || desc3;
    return(desc);
  endsub;

  function LOGNGPD_SCALETRANSFORM() $3;
    length xform $3;
    xform = "LOG";
    return (xform);
  endsub;

  subroutine LOGNGPD_CONSTANTPARM(Xr,Pn);
  endsub;

  function LOGNGPD_PDF(x, Mu,Sigma,Xi,Xr,Pn);
    cutoff = exp(Mu) * Xr;
    p = CDF('LOGN',cutoff, Mu, Sigma);
    if (x < cutoff + constant('MACEPS')) then do;
      return ((Pn/p)*PDF('LOGN', x, Mu, Sigma));
    end;
    else do;
      gpd_scale = p*((1-Pn)/Pn)/PDF('LOGN', cutoff, Mu, Sigma);
      h = (1+Xi*(x-cutoff)/gpd_scale)**(-1-(1/Xi))/gpd_scale;
      return ((1-Pn)*h);
    end;
  endsub;

  function LOGNGPD_CDF(x, Mu,Sigma,Xi,Xr,Pn);
    cutoff = exp(Mu) * Xr;
    p = CDF('LOGN',cutoff, Mu, Sigma);
    if (x < cutoff + constant('MACEPS')) then do;
      return ((Pn/p)*CDF('LOGN', x, Mu, Sigma));
    end;
    else do;
      gpd_scale = p*((1-Pn)/Pn)/PDF('LOGN', cutoff, Mu, Sigma);
      H = 1 - (1 + Xi*((x-cutoff)/gpd_scale))**(-1/Xi);
      return (Pn + (1-Pn)*H);
    end;
  endsub;

  subroutine LOGNGPD_PARMINIT(dim,x[*],nx[*],F[*],Ftype,
    Mu,Sigma,Xi,Xr,Pn);
    outargs Mu,Sigma,Xi,Xr,Pn;
    array xe[1] / nosymbols;
    array nxe[1] / nosymbols;

    eps = constant('MACEPS');

```

```

Pn = 0.8; /* Set mixing probability */
_status_ = .;
call streaminit(56789);
Xb = svrtutil_hillcutoff(dim, x, 100, 25, _status_);
if (missing(_status_) or _status_ = 1) then
    Xb = svrtutil_percentile(Pn, dim, x, F, Ftype);

/* prepare arrays for excess values */
i = 1;
do while (i <= dim and x[i] < Xb+eps);
    i = i + 1;
end;
dime = dim-i+1;
call dynamic_array(xe, dime);
call dynamic_array(nxe, dime);
j = 1;
do while(i <= dim);
    xe[j] = x[i] - Xb;
    nxe[j] = nx[i];
    i = i + 1;
    j = j + 1;
end;

/* Initialize lognormal parameters */
call logn_parminit(dim, x, nx, F, Ftype, Mu, Sigma);
if (not(missing(Mu))) then
    Xr = Xb/exp(Mu);
else
    Xr = .;

/* Initialize GPD's shape parameter using excess values */
call gpd_parminit(dime, xe, nxe, F, Ftype, theta_gpd, Xi);
endsub;

subroutine LOGNGPD_LOWERBOUNDS(Mu, Sigma, Xi, Xr, Pn);
    outargs Mu, Sigma, Xi, Xr, Pn;

    Mu = .; /* Mu has no lower bound */
    Sigma = 0; /* Sigma > 0 */
    Xi = 0; /* Xi > 0 */
endsub;
quit;

```

Note the following points about the LOGNGPD definition:

- The parameters  $x_r$  and  $p_n$  are not estimated with the maximum likelihood method used by PROC HPSEVERITY, so you need to specify them as *constant* parameters by defining the [dist\\_CONSTANTPARM](#) subroutine. The signature of LOGNGPD\_CONSTANTPARM subroutine lists only the constant parameters  $X_r$  and  $P_n$ .
- The parameter  $x_r$  is estimated by first using the SVRTUTIL\_HILLCUTOFF utility function to compute an estimate of the cutoff point  $\hat{x}_b$  and then computing  $x_r = \hat{x}_b / e^{\hat{\mu}}$ . If SVRTUTIL\_HILLCUTOFF

fails to compute a valid estimate, then the SVRTUTIL\_PERCENTILE utility function is used to set  $\hat{x}_b$  to the  $p_n$ th percentile of the data. The parameter  $p_n$  is fixed to 0.8.

- The Sashelp.Svrtdist library is specified with the LIBRARY= option in the PROC FCMP statement to enable the LOGNGPD\_PARMINIT subroutine to use the predefined utility functions (SVRTUTIL\_HILLCUTOFF and SVRTUTIL\_PERCENTILE) and parameter initialization subroutines (LOGN\_PARMINIT and GPD\_PARMINIT).
- The LOGNGPD\_LOWERBOUNDS subroutine defines the lower bounds for all parameters. This subroutine is required because the parameter  $\mu$  has a non-default lower bound. The bounds for  $\sigma$  and  $\xi$  must be specified. If they are not specified, they are returned as missing values, which PROC HPSEVERITY interprets as having no lower bound. You need not specify any bounds for the constant parameters  $X_r$  and  $P_n$ , because they are not subject to optimization.

The following DATA step statements simulate a sample from a mixed-tail distribution with a lognormal body and GPD tail. The parameter  $p_n$  is fixed to 0.8, the same value used in the LOGNGPD\_PARMINIT subroutine defined previously.

```
/*----- Simulate a sample for the mixed-tail distribution -----*/
data testmixdist(keep=y label='Lognormal Body-GPD Tail Sample');
  call streaminit(45678);
  label y='Response Variable';
  N = 100;
  Mu = 1.5;
  Sigma = 0.25;
  Xi = 1.5;
  Pn = 0.8;

  /* Generate data comprising the lognormal body */
  Nbody = N*Pn;
  do i=1 to Nbody;
    y = exp(Mu) * rand('LOGNORMAL')**Sigma;
    output;
  end;

  /* Generate data comprising the GPD tail */
  cutoff = quantile('LOGNORMAL', Pn, Mu, Sigma);
  gpd_scale = (1-Pn) / pdf('LOGNORMAL', cutoff, Mu, Sigma);
  do i=Nbody+1 to N;
    y = cutoff + ((1-rand('UNIFORM'))**(-Xi) - 1)*gpd_scale/Xi;
    output;
  end;
run;
```

The following statements use PROC HPSEVERITY to fit the LOGNGPD distribution model to the simulated sample. They also fit three other predefined distributions (BURR, LOGN, and GPD). The final parameter estimates are written to the Work.Parmest data set.

```
/*--- Set the search path for functions defined with PROC FCMP ---*/
options cmplib=(work.sevexmpl);
/*----- Fit LOGNGPD model with PROC HPSEVERITY -----*/
proc hpseverity data=testmixdist print=all outest=parmest;
  loss y;
  dist logngpd burr logn gpd;
run;
```

Some of the results prepared by PROC HPSEVERITY are shown in [Output 9.3.1](#) and [Output 9.3.2](#). The “Model Selection” table in [Output 9.3.1](#) indicates that all models converged. The last table in [Output 9.3.1](#) shows that the model with LOGNGPD distribution has the best fit according to almost all the statistics of fit. The Burr distribution model is the closest contender to the LOGNGPD model, but the GPD distribution model fits the data very poorly.

**Output 9.3.1** Summary of Fitting Mixed-Tail Distribution

The HPSEVERITY Procedure					
Input Data Set					
Name	WORK.TESTMIXDIST				
Label	Lognormal Body-GPD Tail Sample				
Model Selection					
Distribution	Converged	-2 Log Likelihood	Selected		
logngpd	Yes	418.78232	Yes		
Burr	Yes	424.93728	No		
Logn	Yes	459.43471	No		
Gpd	Yes	558.13444	No		
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
logngpd	418.78232*	428.78232*	429.42062*	441.80817	0.62140*
Burr	424.93728	430.93728	431.18728	438.75280*	0.71373
Logn	459.43471	463.43471	463.55842	468.64505	1.55267
Gpd	558.13444	562.13444	562.25815	567.34478	3.43470
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
logngpd	0.31670*		0.04972*		
Burr	0.57649		0.07860		
Logn	3.27122		0.48448		
Gpd	16.74156		3.31860		
Note: The asterisk (*) marks the best model according to each column's criterion.					

The detailed results for the LOGNGPD distribution are shown in [Output 9.3.2](#). The initial values table indicates the values computed by LOGNGPD\_PARMINIT subroutine for the  $Xr$  and  $Pn$  parameters. It also uses the bounds columns to indicate the constant parameters. The last table in the figure shows the final parameter estimates. The estimates of all free parameters are significantly different from 0. As expected, the final estimates of the constant parameters  $Xr$  and  $Pn$  have not changed from their initial values.

**Output 9.3.2** Detailed Results for the LOGNGPD Distribution

The HPSEVERITY Procedure				
logngpd Distribution				
Distribution Information				
Name	logngpd			
Description	Lognormal Body-GPD Tail Distribution. Mu, Sigma, and Xi are free parameters. Xr and Pn are constant parameters.			
Distribution Parameters	5			
Initial Parameter Values and Bounds				
Parameter	Initial Value	Lower Bound	Upper Bound	
Mu	1.49954	-Infty	Infty	
Sigma	0.76306	1.05367E-8	Infty	
Xi	0.36661	1.05367E-8	Infty	
Xr	1.27395	Constant	Constant	
Pn	0.80000	Constant	Constant	
Convergence Status				
Convergence criterion (GCONV=1E-8) satisfied.				
Optimization Summary				
Optimization Technique	Trust Region			
Iterations	11			
Function Calls	33			
Failed Function Calls	1			
Log Likelihood	-209.39116			
Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu	1.57921	0.06426	24.57	<.0001
Sigma	0.31868	0.04459	7.15	<.0001
Xi	1.03771	0.38205	2.72	0.0078
Xr	1.27395	Constant	.	.
Pn	0.80000	Constant	.	.

The following SAS statements use the parameter estimates to compute the value where the tail region is estimated to start ( $x_b = e^{\hat{\mu}} \hat{x}_r$ ) and the scale of the GPD tail distribution ( $\theta_t = \frac{G(x_b)}{g(x_b)} \frac{(1-p_n)}{p_n}$ ):

```

/*----- Compute tail cutoff and tail distribution's scale -----*/
data xb_thetat(keep=x_b theta_t);
  set parmest(where=(_MODEL_='logngpd' and _TYPE_='EST'));
  x_b = exp(Mu) * Xr;
  theta_t = (CDF('LOGN',x_b,Mu,Sigma)/PDF('LOGN',x_b,Mu,Sigma)) *
            ((1-Pn)/Pn);
run;

proc print data=xb_thetat noobs;
run;

```

**Output 9.3.3** Start of the Tail and Scale of the GPD Tail Distribution

x_b	theta_t
6.18005	1.27865

The computed values of  $x_b$  and  $\theta_t$  are shown as `x_b` and `theta_t` in [Output 9.3.3](#). Equipped with this additional derived information, you can now interpret the results of fitting the mixed-tail distribution as follows:

- The tail starts at  $y \approx 6.18$ . The primary benefit of using the scale-normalized cutoff ( $x_r$ ) as the constant parameter instead of using the actual cutoff ( $x_b$ ) is that the absolute cutoff is optimized by virtue of optimizing the scale of the body region ( $\theta = e^{\mu}$ ).
- The values  $y \leq 6.18$  follow the lognormal distribution with parameters  $\mu \approx 1.58$  and  $\sigma \approx 0.32$ . These parameter estimates are reasonably close to the parameters used for simulating the sample.
- The values  $y_t = y - 6.18$  ( $y_t > 0$ ) follow the GPD distribution with scale  $\theta_t \approx 1.28$  and shape  $\xi \approx 1.04$ .

---

## Example 9.4: Fitting a Scaled Tweedie Model with Regressors

The Tweedie distribution is often used in the insurance industry to explain the effect of independent variables (regressors) on the distribution of losses. PROC HPSEVERITY provides a predefined scaled Tweedie distribution (STWEEDIE) that enables you to model the regression effects on the scale parameter. The scale regression model has its own advantages such as the ability to easily account for inflation effects. This example illustrates how that model can be used to evaluate the effect of regressors on the *mean* of the Tweedie distribution, which is useful in problems such rate-making and pure premium modeling.

Assume a Tweedie process, whose mean  $\mu$  is affected by  $k$  regressors  $x_j$ ,  $j = 1, \dots, k$  as follows:

$$\mu = \mu_0 \exp \left( \sum_{j=1}^k \beta_j x_j \right)$$

where  $\mu_0$  represents the base value of the mean (you can think of  $\mu_0$  as  $\exp(\beta_0)$ , where  $\beta_0$  is the intercept). This model for the mean is identical to the popular generalized linear model for the mean with a logarithmic link function.

More interestingly, it parallels the model used by PROC HPSEVERITY for the scale parameter  $\theta$ ,

$$\theta = \theta_0 \exp \left( \sum_{j=1}^k \beta_j x_j \right)$$

where  $\theta_0$  represents the base value of the scale parameter. As described in the section “[Tweedie Distributions](#)” on page 276, for the parameter range  $p \in (1, 2)$ , the mean of the Tweedie distribution is given by

$$\mu = \theta \lambda \frac{2-p}{p-1}$$

where  $\lambda$  is the Poisson mean parameter of the scaled Tweedie distribution. This relationship enables you to use the scale regression model to infer the effect of regressors on the mean of the distribution.

Let the data set `Work.Test_Sevtw` contain a sample generated from a Tweedie distribution with dispersion parameter  $\phi = 0.5$ , index parameter  $p = 1.75$ , and the mean parameter that is affected by three regression variables `x1`, `x2`, and `x3` as follows:

$$\mu = 5 \exp(0.25 x_1 - x_2 + 3 x_3)$$

Thus, the population values of regression parameters are  $\mu_0 = 5$ ,  $\beta_1 = 0.25$ ,  $\beta_2 = -1$ , and  $\beta_3 = 3$ . You can find the code used to generate the sample in the PROC HPSEVERITY sample program `hsever04.sas`.

The following PROC HPSEVERITY step uses the sample in `Work.Test_Sevtw` data set to estimate the parameters of the scale regression model for the predefined scaled Tweedie distribution (STWEEDIE) with the dual quasi-Newton (QUANEW) optimization technique:

```
/*--- Fit the scale parameter version of the Tweedie distribution ---*/
proc hpseverity data=test_sevtw outest=estw covout print=all;
    loss y;
    scalemodel x1-x3;

    dist stweedie;
    nloptions tech=quanew;
run;
```

The dual quasi-Newton technique is used because it requires only the first-order derivatives of the objective function, and it is harder to compute reasonably accurate estimates of the second-order derivatives of Tweedie distribution’s PDF with respect to the parameters.

Some of the key results prepared by PROC HPSEVERITY are shown in [Output 9.4.1](#) and [Output 9.4.2](#). The distribution information and the convergence results are shown in [Output 9.4.1](#).



**Output 9.4.1** Convergence Results for the STWEEDIE Model with Regressors

The HPSEVERITY Procedure	
stweedie Distribution	
Distribution Information	
Name	stweedie
Description	Tweedie Distribution with Scale Parameter
Distribution Parameters	3
Regression Parameters	3
Convergence Status	
Convergence criterion (FCONV=2.220446E-16) satisfied.	
Optimization Summary	
Optimization Technique	Dual Quasi-Newton
Iterations	42
Function Calls	218
Log Likelihood	-1044.3

The final parameter estimates of the STWEEDIE regression model are shown in [Output 9.4.2](#). The estimate that is reported for the parameter Theta is the estimate of the base value  $\theta_0$ . The estimates of regression coefficients  $\beta_1$ ,  $\beta_2$ , and  $\beta_3$  are indicated by the rows of x1, x2, and x3, respectively.

**Output 9.4.2** Parameter Estimates for the STWEEDIE Model with Regressors

Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Theta	0.82888	0.26657	3.11	0.0021
Lambda	16.57174	13.12083	1.26	0.2076
P	1.75440	0.20187	8.69	<.0001
x1	0.27970	0.09876	2.83	0.0049
x2	-0.76715	0.10313	-7.44	<.0001
x3	3.03225	0.10142	29.90	<.0001

If your goal is to explain the effect of regressors on the scale parameter, then the output displayed in [Output 9.4.2](#) is sufficient. But, if you want to compute the effect of regressors on the mean of the distribution, then you need to do some postprocessing. Using the relationship between  $\mu$  and  $\theta$ ,  $\mu$  can be written in terms of the parameters of the STWEEDIE model as

$$\mu = \theta_0 \exp \left( \sum_{j=1}^k \beta_j x_j \right) \lambda^{\frac{2-p}{p-1}}$$

This shows that the parameters  $\beta_j$  are identical for the mean and the scale model, and the base value  $\mu_0$  of the mean model is

$$\mu_0 = \theta_0 \lambda \frac{2-p}{p-1}$$

The estimate of  $\mu_0$  and the standard error associated with it can be computed by using the property of the functions of maximum likelihood estimators (MLE). If  $g(\Omega)$  represents a totally differentiable function of parameters  $\Omega$ , then the MLE of  $g$  has an asymptotic normal distribution with mean  $g(\hat{\Omega})$  and covariance  $C = (\partial g)' \Sigma (\partial g)$ , where  $\hat{\Omega}$  is the MLE of  $\Omega$ ,  $\Sigma$  is the estimate of covariance matrix of  $\Omega$ , and  $\partial g$  is the gradient vector of  $g$  with respect to  $\Omega$  evaluated at  $\hat{\Omega}$ . For  $\mu_0$ , the function is  $g(\Omega) = \theta_0 \lambda (2-p)/(p-1)$ . The gradient vector is

$$\begin{aligned} \partial g &= \left( \frac{\partial g}{\partial \theta_0} \quad \frac{\partial g}{\partial \lambda} \quad \frac{\partial g}{\partial p} \quad \frac{\partial g}{\partial \beta_1} \cdots \frac{\partial g}{\partial \beta_k} \right) \\ &= \left( \frac{\mu_0}{\theta_0} \quad \frac{\mu_0}{\lambda} \quad \frac{-\mu_0}{(p-1)(2-p)} \quad 0 \cdots 0 \right) \end{aligned}$$

You can write a DATA step that implements these computations by using the parameter and covariance estimates prepared by PROC HPSEVERITY step. The DATA step program is available in the sample program *hsevev04.sas*. The estimates of  $\mu_0$  prepared by that program are shown in [Output 9.4.3](#). These estimates and the estimates of  $\beta_j$  as shown in [Output 9.4.2](#) are reasonably close (that is, within one or two standard errors) to the parameters of the population from which the sample in Work.Test\_Sevtw data set was drawn.

**Output 9.4.3** Estimate of the Base Value Mu0 of the Mean Parameter

Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu0	4.47179	0.42225	10.5904	0

Another effect of using the scaled Tweedie distribution to model the regression effects is that the regressors also affect the variance  $V$  of the Tweedie distribution. The variance is related to the mean as  $V = \phi \mu^p$ , where  $\phi$  is the dispersion parameter. Using the relationship between the parameters TWEEDIE and STWEEDIE distributions as described in the section “[Tweedie Distributions](#)” on page 276, the regression model for the dispersion parameter is

$$\begin{aligned} \log(\phi) &= (2-p) \log(\mu) - \log(\lambda(2-p)) \\ &= ((2-p) \log(\mu_0) - \log(\lambda(2-p))) + (2-p) \sum_{j=1}^k \beta_j x_j \end{aligned}$$

Subsequently, the regression model for the variance is

$$\begin{aligned} \log(V) &= 2 \log(\mu) - \log(\lambda(2-p)) \\ &= (2 \log(\mu_0) - \log(\lambda(2-p))) + 2 \sum_{j=1}^k \beta_j x_j \end{aligned}$$

In summary, PROC HPSEVERITY enables you to estimate regression effects on various parameters and statistics of the Tweedie model.

## Example 9.5: Fitting Distributions to Interval-Censored Data

In some applications, the data available for modeling might not be exact. A commonly encountered scenario is the use of grouped data from an external agency, which for several reasons, including privacy, does not provide information about individual loss events. The losses are grouped into disjoint bins, and you know only the range and number of values in each bin. Each group is essentially interval-censored, because you know that a loss magnitude is in certain interval, but you do not know the exact magnitude. This example illustrates how you can use PROC HPSEVERITY to model such data.

The following DATA step generates sample grouped data for dental insurance claims, which is taken from Klugman, Panjer, and Willmot (1998):

```
/* Grouped dental insurance claims data
   (Klugman, Panjer, and Willmot, 1998) */
data gdental;
    input lowerbd upperbd count @@;
    datalines;
0 25 30 25 50 31 50 100 57 100 150 42 150 250 65 250 500 84
500 1000 45 1000 1500 10 1500 2500 11 2500 4000 3
;
run;
```

The following PROC HPSEVERITY step fits all the predefined distributions to the data in Work.Gdental data set:

```
/* Fit all predefined distributions */
proc hpseverity data=gdental edf=turnbull print=all criterion=aicc;
    loss / rc=lowerbd lc=upperbd;
    weight count;
    dist _predef_;
    performance nthreads=1;
run;
```

The EDF= option in the PROC HPSEVERITY statement specifies that the Turnbull's method be used for EDF estimation. The LOSS statement specifies the left and right boundaries of each group as the right-censoring and left-censoring limits, respectively. The variable count records the number of losses in each group and is specified in the WEIGHT statement. Note that no response variable is specified in the LOSS statement, which is allowed as long as each observation in the input data set is censored. The PERFORMANCE statement specifies that just one thread of execution be used, to minimize the overhead associated with multithreading, because the input data set is very small.

Some of the key results prepared by PROC HPSEVERITY are shown in [Output 9.5.1](#). According to the “Model Selection” table in [Output 9.5.1](#), all distribution models have converged. The “All Fit Statistics” table in [Output 9.5.1](#) indicates that the exponential distribution (EXP) has the best fit for data according to a majority of the likelihood-based statistics.

**Output 9.5.1** Statistics of Fit for Interval-Censored Data

The HPSEVERITY Procedure					
Input Data Set					
Name	WORK.GDENTAL				
Model Selection					
Distribution	Converged	AICC	Selected		
Burr	Yes	51.41112	No		
Exp	Yes	44.64768	Yes		
Gamma	Yes	47.63969	No		
Igauss	Yes	48.05874	No		
Logn	Yes	47.34027	No		
Pareto	Yes	47.16908	No		
Gpd	Yes	47.16908	No		
Weibull	Yes	47.47700	No		
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
Burr	41.41112*	47.41112	51.41112	48.31888	0.08974*
Exp	42.14768	44.14768*	44.64768*	44.45026*	0.26412
Gamma	41.92541	45.92541	47.63969	46.53058	0.19569
Igauss	42.34445	46.34445	48.05874	46.94962	0.34514
Logn	41.62598	45.62598	47.34027	46.23115	0.16853
Pareto	41.45480	45.45480	47.16908	46.05997	0.11423
Gpd	41.45480	45.45480	47.16908	46.05997	0.11423
Weibull	41.76272	45.76272	47.47700	46.36789	0.17238
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
Burr	0.00103*		0.0000816*		
Exp	0.09936		0.01866		
Gamma	0.04608		0.00759		
Igauss	0.12301		0.02562		
Logn	0.01884		0.00333		
Pareto	0.00739		0.0009084		
Gpd	0.00739		0.0009084		
Weibull	0.03293		0.00472		
Note: The asterisk (*) marks the best model according to each column's criterion.					

## Example 9.6: Benefits of Distributed and Multithreaded Computing

One of the key features of the HPSEVERITY procedure is that it takes advantage of the distributed and multithreaded computing machinery in order to solve a given problem faster. This example illustrates the benefits of using multithreading and distributed computing.

The example uses a simulated data set `Work.Largedata`, which contains 10,000,000 observations, some of which are right-censored or left-truncated. The losses are affected by three external effects. The DATA step program that generates this data set is available in the accompanying sample program `hsevex06.sas`.

The following PROC HPSEVERITY step fits all the predefined distributions to the data in `Work.Largedata` data set on the client machine with just one thread of computation:

```
/* Fit all predefined distributions without any multithreading or
   distributed computing */
proc hpseverity data=largedata criterion=aicc initsample(size=20000);
  loss y / lt=threshold rc=limit;
  scalemodel x1-x3;
  dist _predef_;
  performance nthreads=1 bufsize=1000000 details;
run;
```

The `NTHREADS=1` option in the `PERFORMANCE` statement specifies that just one thread of computation be used. The absence of the `NODES=` option in the `PERFORMANCE` statement specifies that single-machine mode of execution be used. That is, this step does not use any multithreading or distributed computing. The `BUFSIZE=` option in the `PERFORMANCE` statement specifies the number of observations to read at one time. Specifying a larger value tends to decrease the time it takes to load the data. The `DETAILS` option in the performance statement enables reporting of the timing information. The `INITSAMPLE` option in the PROC HPSEVERITY statement specifies that a uniform random sample of maximum 20,000 observations be used for parameter initialization.

The “Performance Information” and “Procedure Task Timing” tables that PROC HPSEVERITY prepares are shown in [Output 9.6.1](#). The “Performance Information” table contains the information about the execution environment. The “Procedure Task Timing” table indicates the total time and relative time taken by each of the four main steps of PROC HPSEVERITY. As that table shows, it takes around 35 minutes for the task of estimating parameters, which is usually the most time-consuming of all the tasks.

**Output 9.6.1** Performance for Single-Machine Mode with No Multithreading

The HPSEVERITY Procedure	
Performance Information	
Execution Mode	Single-Machine
Number of Threads	1

**Output 9.6.1** *continued*

Procedure Task Timing		
Task	Seconds	Percent
Load and Prepare Models	3.92	0.18%
Load and Prepare Data	1.87	0.09%
Initialize Parameters	1.25	0.06%
Estimate Parameters	2110.82	99.56%
Compute Fit Statistics	2.22	0.10%

If the grid appliance is not available, you can improve the performance by using multiple threads of computation; this is in fact the default. The following PROC HPSEVERITY step fits all the predefined distributions by using all the logical CPU cores of the machine:

```
/* Specify that all the logical CPU cores on the machine be used */
options cpucount=actual;

/* Fit all predefined distributions with multithreading, but no
distributed computing */
proc hpseverity data=largedata criterion=aicc initsample(size=20000);
  loss y / lt=threshold rc=limit;
  scalemodel x1-x3;
  dist _predef_;
  performance bufsize=1000000 details;
run;
```

When you do not specify the NTHREADS= option in the PERFORMANCE statement, the HPSEVERITY procedure uses the value of the CPUCOUNT= system option to decide the number of threads to use in single-machine mode. Setting the CPUCOUNT= option to ACTUAL before the PROC HPSEVERITY step enables the procedure to use all the logical cores of the machine. The machine that is used to obtain these results (and the earlier results in [Output 9.6.1](#)) has four physical CPU cores, each with a clock speed of 3.4 GHz. Hyperthreading is enabled on the CPUs to yield eight logical CPU cores; this number is confirmed by the “Performance Information” table in [Output 9.6.2](#). The results in the “Procedure Task Timing” table in [Output 9.6.2](#) indicate that the use of multithreading has improved the performance significantly by reducing the time to estimate parameters to around 8 minutes.

**Output 9.6.2** Performance for Single-Machine Mode with Eight Threads

The HPSEVERITY Procedure	
Performance Information	
Execution Mode	Single-Machine
Number of Threads	8

**Output 9.6.2** *continued*

Procedure Task Timing		
Task	Seconds	Percent
Load and Prepare Models	0.69	0.14%
Load and Prepare Data	1.50	0.32%
Initialize Parameters	0.92	0.19%
Estimate Parameters	470.45	99.11%
Compute Fit Statistics	1.11	0.23%

When a grid appliance is available, performance can be further improved by using more than one node in the distributed mode of execution. Large data sets are usually predistributed on the grid appliance that hosts a distributed database. In other words, large problems are best suited for the alongside-the-database model of execution. However, for the purpose of illustration, this example assumes that the data set is available on the client machine and is then distributed to the grid nodes by the HPSEVERITY procedure according to the options that are specified in the PERFORMANCE statement.

The next few PROC HPSEVERITY steps are run on a grid appliance by varying the number of nodes and the number of threads that are used within each node.

You can specify your distributed computing environment by using SAS environment variables or by specifying options in the PERFORMANCE statement, or by a combination of these methods. For example, you can submit the following statements to specify the appliance host (GRIDHOST= SAS environment variable) and the installation location of shared libraries on the appliance (GRIDINSTALLLOC= SAS environment variable):

```
/* Set the appliance host and installation location that are
   appropriate for your distributed mode setup */
option set=GRIDHOST      ="&GRIDHOST";
option set=GRIDINSTALLLOC="&GRIDINSTALLLOC";
```

To run the preceding statements successfully, you need to set the macro variables GRIDHOST and GRIDINSTALLLOC to resolve to appropriate values, or you can replace the references to macro variables with the appropriate values. Alternatively, you can specify the HOST= and INSTALL= options in the PERFORMANCE statement; this method is used in the PROC HPSEVERITY steps of this example. You can use other SAS environment variables and PERFORMANCE statement options to describe your distributed computing environment. For more information, see the section “[PERFORMANCE Statement](#)” on page 39.

To establish a reference point for the performance of one CPU of a grid node, the results of using only one node of the grid appliance without any multithreading are presented first. The particular grid appliance that is used to obtain these results has more than sixteen nodes. Each node has 8 dual-core CPUs with a clock speed of 2.7 GHz. The following PROC HPSEVERITY step fits all the predefined distributions to the data in the Work.Largedata data set:

```
/* Fit all predefined distributions on 1 grid node without
   any multithreading */
proc hpseverity data=largedata criterion=aicc initsample(size=20000);
  loss y / lt=threshold rc=limit;
  scalemodel x1-x3;
  dist _predef_;
  performance nodes=1 nthreads=1 details
    host="&GRIDHOST" install="&GRIDINSTALLLOC";
run;
```

The PERFORMANCE statement specifies that only one node be used to fit the models, with only one thread of computation on that node. The “Performance Information” and “Procedure Task Timing” tables that PROC HPSEVERITY prepares are shown in [Output 9.6.3](#). It takes around 37 minutes to complete the task of estimating parameters.

**Output 9.6.3** Performance on One Grid Appliance Node with No Multithreading

Performance Information		
Host Node	<< your grid host >>	
Install Location	<< your grid install location >>	
Execution Mode	Distributed	
Grid Mode	Symmetric	
Number of Compute Nodes	1	
Number of Threads per Node	1	
Procedure Task Timing		
Task	Seconds	Percent
Load and Prepare Models	0.47	0.02%
Load and Prepare Data	0.87	0.04%
Initialize Parameters	1.09	0.05%
Estimate Parameters	2207.86	99.81%
Compute Fit Statistics	1.67	0.08%

The computations and time taken to fit each model are shown in the “Estimation Details” table of [Output 9.6.4](#), which is generated whenever you specify the DETAILS option in the PERFORMANCE statement. This table can be useful for comparing the relative effort required to fit each model and drawing some broader conclusions. For example, even if the Pareto distribution takes a larger number of iterations, function calls, and gradient and Hessian updates than the gamma distribution, it takes less time to complete; this indicates that the individual PDF and CDF computations of the gamma distribution are more expensive than those of the Pareto distribution.



**Output 9.6.4** Estimation Details

The HPSEVERITY Procedure						
Estimation Details						
Distribution	Converged	Iterations	Function Calls	Gradient Updates	Hessian Updates	Time (Seconds)
Burr	Yes	11	28	104	90	321.10
Exp	Yes	4	12	27	20	28.84
Gamma	Yes	6	16	44	35	899.13
Igauss	Yes	4	12	27	20	267.96
Logn	Yes	4	12	27	20	108.75
Pareto	Yes	39	113	902	860	385.40
Gpd	Yes	6	17	44	35	125.68
Weibull	Yes	4	12	27	20	70.99

To obtain the next reference point for performance, the following PROC HPSEVERITY step specifies that 16 computation threads be used on one node of the grid appliance:

```
/* Fit all predefined distributions on 1 grid node with multithreading */
proc hpseverity data=largedata criterion=aicc initsample(size=20000);
  loss y / lt=threshold rc=limit;
  scalemodel x1-x3;
  dist _predef_;
  performance nodes=1 nthreads=16 details
    host="&GRIDHOST" install="&GRIDINSTALLLOC";
run;
```

The performance tables that are prepared by the preceding statements are shown in [Output 9.6.5](#). As the “Procedure Task Timing” table shows, use of multithreading has improved the performance significantly over that of the single-threaded case. Now, it takes around 3 minutes to complete the task of estimating parameters.

**Output 9.6.5** Performance Information with Multithreading but No Distributed Computing

Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	1
Number of Threads per Node	16

**Output 9.6.5** *continued*

The HPSEVERITY Procedure		
Procedure Task Timing		
Task	Seconds	Percent
Load and Prepare Models	0.51	0.27%
Load and Prepare Data	0.47	0.24%
Initialize Parameters	0.97	0.50%
Estimate Parameters	190.35	98.48%
Compute Fit Statistics	0.99	0.51%

You can combine the power of multithreading and distributed computing by specifying that multiple nodes of the grid be used to accomplish the task. The following PROC HPSEVERITY step specifies that 16 nodes of the grid appliance be used:

```
/* Fit all predefined distributions with distributed computing and
   multithreading within each node */
proc hpseverity data=largedata criterion=aicc initsample(size=20000);
  loss y / lt=threshold rc=limit;
  scalemodel x1-x3;
  dist _predef_;
  performance nodes=16 nthreads=16 details
    host=&GRIDHOST install=&GRIDINSTALLLOC";
run;
```

When the DATA= data set is local to the client machine, as it is in this example, you must specify a nonzero value for the NODES= option in the PERFORMANCE statement in order to enable the distributed mode of execution. In other words, for the distributed mode that is not executing alongside the database, omitting the NODES= option is equivalent to specifying NODES=0, which is single-machine mode.

The performance tables that are prepared by the preceding statements are shown in [Output 9.6.6](#). If you compare these tables to the tables in [Output 9.6.3](#) and [Output 9.6.5](#), you see that the task that would have taken a long time with a single thread of execution on a single machine (over half an hour) can be performed in a much shorter time (around 17 seconds) by using the computational resources of the grid appliance to combine the power of multithreaded and distributed computing.

**Output 9.6.6** Performance Information with Distributed Computing and Multithreading

Performance Information	
Host Node	<< your grid host >>
Install Location	<< your grid install location >>
Execution Mode	Distributed
Grid Mode	Symmetric
Number of Compute Nodes	16
Number of Threads per Node	16

**Output 9.6.6** *continued*

The HPSEVERITY Procedure		
Procedure Task Timing		
Task	Seconds	Percent
Load and Prepare Models	0.53	2.78%
Load and Prepare Data	0.03	0.16%
Initialize Parameters	0.78	4.08%
Estimate Parameters	16.99	88.90%
Compute Fit Statistics	0.78	4.07%

The machines that were used to obtain these performance results are relatively modest machines, and PROC HPSEVERITY was executed in a multiuser environment; that is, background processes were running in single-machine mode or other users were using the grid in distributed mode. For time-critical applications, you can use a larger, dedicated grid that consists of more powerful machines to achieve more dramatic performance improvement.

---

**Example 9.7: Estimating Parameters Using Cramér-von Mises Estimator**

PROC HPSEVERITY enables you to estimate model parameters by minimizing your own objective function. This example illustrates how you can use PROC HPSEVERITY to implement the Cramér-von Mises estimator. Let  $F(y_i; \Theta)$  denote the estimate of CDF at  $y_i$  for a distribution with parameters  $\Theta$ , and let  $F_n(y_i)$  denote the empirical estimate of CDF (EDF) at  $y_i$  that is computed from a sample  $y_i$ ,  $1 \leq i \leq N$ . Then, the Cramér-von Mises estimator of the parameters is defined as

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{i=1}^N (F(y_i; \Theta) - F_n(y_i))^2$$

This estimator belongs to the class of minimum distance estimators. It attempts to estimate the parameters such that the squared distance between the CDF and EDF estimates is minimized.

The following PROC HPSEVERITY step uses the Cramér-von Mises estimator to fit four candidate distribution models, including the LOGNGPD mixed-tail distribution model that was defined in “[Example 9.3: Defining a Model for Mixed-Tail Distributions](#)” on page 352. The input sample is the same as is used in that example.

```

/*--- Set the search path for functions defined with PROC FCMP ---*/
options cmplib=(work.sevexmpl);

/*----- Fit LOGNGPD model with PROC HPSEVERITY by using -----
----- the Cramer-von Mises minimum distance estimator -----*/
proc hpseverity data=testmixdist obj=cvmobj print=all;
  loss y;
  dist logngpd burr logn gpd;

  * Cramer-von Mises estimator (minimizes the distance *
  * between parametric and nonparametric estimates)    *;
  cvmobj = _cdf_(y);
  cvmobj = (cvmobj - _edf_(y))**2;
run;

```

The OBJ= option in the PROC HPSEVERITY statement specifies that the objective function cvmobj should be minimized. The programming statements compute the contribution of each observation in the input data set to the objective function cvmobj. The use of keyword functions \_CDF\_ and \_EDF\_ makes the program applicable to all the distributions.

Some of the key results prepared by PROC HPSEVERITY are shown in [Output 9.7.1](#). The “Model Selection” table indicates that all models converged. When you specify a custom objective function, the default selection criterion is the value of the custom objective function. The “All Fit Statistics” table indicates that LOGNGPD is the best distribution according to all the statistics of fit. Comparing the fit statistics of [Output 9.7.1](#) with those of [Output 9.3.1](#) indicates that the use of the Cramér-von Mises estimator has resulted in smaller values for all the EDF-based statistics of fit for all the models, which is expected from a minimum distance estimator.

**Output 9.7.1** Summary of Cramér-von Mises Estimation

The HPSEVERITY Procedure					
Input Data Set					
Name	WORK.TESTMIXDIST				
Label	Lognormal Body-GPD Tail Sample				
Model Selection					
Distribution	Converged	cvmobj	Selected		
logngpd	Yes	0.02694	Yes		
Burr	Yes	0.03325	No		
Logn	Yes	0.03633	No		
Gpd	Yes	2.96090	No		
All Fit Statistics					
Distribution	cvmobj	-2 Log Likelihood	AIC	AICC	BIC
logngpd	0.02694*	419.49635*	429.49635*	430.13464*	442.52220*
Burr	0.03325	436.58823	442.58823	442.83823	450.40374
Logn	0.03633	491.88659	495.88659	496.01030	501.09693
Gpd	2.96090	560.35409	564.35409	564.47780	569.56443
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	KS	AD	CvM		
logngpd	0.51332*	0.21563*	0.03030*		
Burr	0.53084	0.82875	0.03807		
Logn	0.52469	2.08312	0.04173		
Gpd	2.99095	15.51378	2.97806		
Note: The asterisk (*) marks the best model according to each column's criterion.					

### Example 9.8: Defining a Finite Mixture Model That Has a Scale Parameter

A finite mixture model is a stochastic model that postulates that the probability distribution of the data generation process is a mixture of a finite number of probability distributions. For example, when an insurance company analyzes loss data from multiple policies that are underwritten in different geographic regions, some regions might behave similarly, but the distribution that governs some regions might be different from the distribution that governs other regions. Further, it might not be known which regions behave similarly. Also, the larger amounts of losses might follow a different stochastic process from the stochastic process that governs the smaller amounts of losses. It helps to model all policies together in order to pool the data together and exploit any commonalities among the regions, and the use of a finite mixture model can help capture the differences in distributions across regions and ranges of loss amounts.

Formally, if  $f_i$  and  $F_i$  denote the PDF and CDF, respectively, of component distribution  $i$  and  $p_i$  represents the mixing probability that is associated with component  $i$ , then the PDF and CDF of the finite mixture of  $K$  distribution components are

$$f(x; \Theta, \mathbf{p}) = \sum_{i=1}^K p_i f_i(x; \Theta_i)$$

$$F(x; \Theta, \mathbf{p}) = \sum_{i=1}^K p_i F_i(x; \Theta_i)$$

where  $\Theta_i$  denotes the parameters of component distribution  $i$  and  $\Theta$  denotes the parameters of the mixture distribution, which is a union of all the  $\Theta_i$  parameters.  $\mathbf{p}$  denotes the set of mixing probabilities. All mixing probabilities must add up to 1 ( $\sum_{i=1}^K p_i = 1$ ).

You can define the finite mixture of a specific number of components and specific distributions for each of the components by defining the FCMP functions for the PDF and CDF. However, in general, it is not possible to fit a scale regression model by using any finite mixture distribution unless you take special care to ensure that the mixture distribution has a scale parameter. This example provides a formulation of a two-component finite mixture model that has a scale parameter.

To start with, each component distribution must have either a scale parameter or a log-transformed scale parameter. Let  $\theta_1$  and  $\theta_2$  denote the scale parameters of the first and second components, respectively. Let  $p_1 = p$  be the mixing probability, which makes  $p_2 = 1 - p$  by using the constraint on  $\mathbf{p}$ . The PDF of the mixture of these two distributions can be written as

$$f(x; \theta_1, \theta_2, \Phi, p) = \frac{p}{\theta_1} f_1\left(\frac{x}{\theta_1}; \Phi_1\right) + \frac{1-p}{\theta_2} f_2\left(\frac{x}{\theta_2}; \Phi_2\right)$$

where  $\Phi_1$  and  $\Phi_2$  denote the sets of nonscale parameters of the first and second components, respectively, and  $\Phi$  denotes a union of  $\Phi_1$  and  $\Phi_2$ . For the mixture to have the scale parameter  $\theta$ , the PDF must be of the form

$$f(x; \theta, \Phi', p) = \frac{1}{\theta} \left( p f_1\left(\frac{x}{\theta}; \Phi'_1\right) + (1 - p) f_2\left(\frac{x}{\theta}; \Phi'_2\right) \right)$$

where  $\Phi'$ ,  $\Phi'_1$ , and  $\Phi'_2$  denote the modified sets of nonscale parameters. One simple way to achieve this is to make  $\theta_1 = \theta_2 = \theta$  and  $\Phi' = \Phi$ ; that is, you simply equate the scale parameters of both components and keep the set of nonscale parameters unchanged. However, forcing the scale parameters to be equal in both components is restrictive, because the mixture cannot model potential differences in the scales of the two components. A better approach is to tie the scale parameters of the two components by a ratio such that  $\theta_1 = \theta$  and  $\theta_2 = \rho\theta$ . If the ratio parameter  $\rho$  is estimated along with the other parameters, then the mixture distribution becomes flexible enough to model the variations across the scale parameters of individual components.

To summarize, the PDF and CDF are of the following form for the two-component mixture that has a scale parameter:

$$\begin{aligned} f(x; \theta, \rho, \Phi, p) &= \frac{1}{\theta} \left( p f_1\left(\frac{x}{\theta}; \Phi_1\right) + (1 - p) f_2\left(\frac{x}{\theta}; \rho, \Phi_2\right) \right) \\ F(x; \theta, \rho, \Phi, p) &= p F_1\left(\frac{x}{\theta}; \Phi_1\right) + (1 - p) F_2\left(\frac{x}{\theta}; \rho, \Phi_2\right) \end{aligned}$$

This can be generalized to a mixture of  $K$  components by introducing the  $K - 1$  ratio parameters  $\rho_i$  that relate the scale parameters of each of the  $K$  components to the scale parameter  $\theta$  of the mixture distribution as follows:

$$\begin{aligned} \theta_1 &= \theta \\ \theta_i &= \rho_i \theta; i \in [2, K] \end{aligned}$$

In order to illustrate this approach, define a mixture of two lognormal distributions by using the following PDF function:

$$\begin{aligned} f(x; \mu, \sigma_1, p_2, \rho_2, \sigma_2) &= \frac{(1 - p_2)}{\sigma_1 x \sqrt{2\pi}} \exp\left(\frac{-(\log(x) - \mu)^2}{2\sigma_1^2}\right) + \\ &\quad \frac{p_2}{\sigma_2 x \sqrt{2\pi}} \exp\left(\frac{-(\log(x) - \mu - \log(\rho_2))^2}{2\sigma_2^2}\right) \end{aligned}$$

You can verify that  $\mu$  serves as the log of the scale parameter  $\theta$  ( $\mu = \log(\theta)$ ).

The following PROC FCMP steps encode this formulation in a distribution named SLOGNMIX2 for use with PROC HPSEVERITY:

```

/*- Define Mixture of 2 Lognormal Distributions with a Log-Scale Parameter -*/
proc fcmp library=sashelp.svrtldist outlib=work.sevexmpl.models;
  function slognmix2_description() $128;
    return ("Mixture of two lognormals with a log-scale parameter Mu");
  endsub;

  function slognmix2_scaletransform() $8;
    return ("LOG");
  endsub;

  function slognmix2_pdf(x, Mu, Sigma1, p2, Rho2, Sigma2);
    Mu1 = Mu;
    Mu2 = Mu + log(Rho2);
    pdf1 = logn_pdf(x, Mu1, Sigma1);
    pdf2 = logn_pdf(x, Mu2, Sigma2);
    return ((1-p2)*pdf1 + p2*pdf2);
  endsub;

  function slognmix2_cdf(x, Mu, Sigma1, p2, Rho2, Sigma2);
    Mu1 = Mu;
    Mu2 = Mu + log(Rho2);
    cdf1 = logn_cdf(x, Mu1, Sigma1);
    cdf2 = logn_cdf(x, Mu2, Sigma2);
    return ((1-p2)*cdf1 + p2*cdf2);
  endsub;

  subroutine slognmix2_parminit(dim, x[*], nx[*], F[*], Ftype,
                                Mu, Sigma1, p2, Rho2, Sigma2);
    outargs Mu, Sigma1, p2, Rho2, Sigma2;
    array m[1] / nosymbols;
    p2 = 0.5;
    Rho2 = 0.5;
    median = svrtutil_percentile(0.5, dim, x, F, Ftype);
    Mu = log(2*median/1.5);
    call svrtutil_rawmoments(dim, x, nx, 1, m);
    lm1 = log(m[1]);

    /* Search Rho2 that makes log(sample mean) > Mu */
    do while (lm1 <= Mu and Rho2 < 1);
      Rho2 = Rho2 + 0.01;
      Mu = log(2*median/(1+Rho2));
    end;
    if (Rho2 >= 1) then
      /* If Mu cannot be decreased enough to make it less
         than log(sample mean), then revert to Rho2=0.5.
         That will set Sigma1 and possibly Sigma2 to missing.
         PROC HPSEVERITY replaces missing initial values with 0.001. */
      Mu = log(2*median/1.5);

    Sigma1 = sqrt(2.0*(log(m[1])-Mu));
    Sigma2 = sqrt(2.0*(log(m[1])-Mu-log(Rho2)));

```



```

endsub;

subroutine slognmix2_lowerbounds(Mu, Sigma1, p2, Rho2, Sigma2);
  outargs Mu, Sigma1, p2, Rho2, Sigma2;
  Mu = .; /* Mu has no lower bound */
  Sigma1 = 0; /* Sigma1 > 0 */
  p2 = 0; /* p2 > 0 */
  Rho2 = 0; /* Rho2 > 0 */
  Sigma2 = 0; /* Sigma2 > 0 */
endsub;

subroutine slognmix2_upperbounds(Mu, Sigma1, p2, Rho2, Sigma2);
  outargs Mu, Sigma1, p2, Rho2, Sigma2;
  Mu = .; /* Mu has no upper bound */
  Sigma1 = .; /* Sigma1 has no upper bound */
  p2 = 1; /* p2 < 1 */
  Rho2 = 1; /* Rho2 < 1 */
  Sigma2 = .; /* Sigma2 has no upper bound */
endsub;
quit;

```

As shown in previous examples, an important aspect of defining a distribution for use with PROC HPSEVERITY is the definition of the PARMINIT subroutine that initializes the parameters. For mixture distributions, in general, the parameter initialization is a nontrivial task. For a two-component mixture, some simplifying assumptions make the problem easier to handle. For the initialization of SLOGNMIX2, the initial values of  $p_2$  and  $\rho_2$  are fixed at 0.5, and the following two simplifying assumptions are made:

- The median of the mixture is the average of the medians of the two components:

$$F^{-1}(0.5) = (\exp(\mu_1) + \exp(\mu_2))/2 = \exp(\mu)(1 + \rho_2)/2$$

Solution of this equation yields the value of  $\mu$  in terms of  $\rho_2$  and the sample median.

- Each component has the same mean, which implies the following:

$$\exp(\mu + \sigma_1^2/2) = \exp(\mu + \log(\rho_2) + \sigma_2^2/2)$$

If  $X_i$  represents the random variable of component distribution  $i$  and  $X$  represents the random variable of the mixture distribution, then the following equation holds for the raw moment of any order  $k$ :

$$E[X^k] = \sum_{i=1}^K p_i E[X_i^k]$$

This, in conjunction with the assumption on component means, leads to the equations

$$\log(m_1) = \mu + \frac{\sigma_1^2}{2} \qquad \log(m_1) = \mu + \log(\rho_2) + \frac{\sigma_2^2}{2}$$

where  $m_1$  denotes the first raw moment of the sample. Solving these equations leads to the following values of  $\sigma_1$  and  $\sigma_2$ :

$$\sigma_1^2 = 2(\log(m_1) - \mu) \qquad \sigma_2^2 = 2(\log(m_1) - \mu - \log(\rho_2))$$

Note that  $\sigma_1$  has a valid value only if  $\log(m_1) > \mu$ . Among the many possible methods of ensuring this condition, the SLOGNMIX2\_PARMINIT subroutine uses the method of doing a linear search over  $\rho_2$ .

Even when the preceding assumptions are not true for a given problem, they produce reasonable initial values to help guide the nonlinear optimizer to an acceptable optimum if the mixture of two lognormal distributions is indeed a good fit for your input data. This is illustrated by the results of the following steps that fit the SLOGNMIX2 distribution to simulated data, which have different means for the two components (12.18 and 22.76, respectively), and the median of the sample (15.94) is not equal to the average of the medians of the two components (7.39 and 20.09, respectively):

```

/*----- Simulate a lognormal mixture sample -----*/
data testlognmix(keep=y);
  call streaminit(12345);
  Mu1 = 2;
  Sigma1 = 1;
  i = 0;
  do j=1 to 2000;
    y = exp(Mu1) * rand('LOGNORMAL')**Sigma1;
    output;
  end;
  Mu2 = 3;
  Sigma2 = 0.5;
  do j=1 to 3000;
    y = exp(Mu2) * rand('LOGNORMAL')**Sigma2;
    output;
  end;
run;

/*-- Fit and compare scale regression models with 2-component --*/
/*-- lognormal mixture and the standard lognormal distribution --*/
options cmplib=(work.sevexmpl);

proc hpseverity data=testlognmix print=all;
  loss y;
  dist slognmix2 logn;
run;

```

The comparison of the fit statistics of SLOGNMIX2 and LOGN, as shown in [Output 9.8.1](#), confirms that the two-component mixture is certainly a better fit to these data than the single lognormal distribution.

**Output 9.8.1** Comparison of Fitting One versus Two Lognormal Components to Mixture Data

The HPSEVERITY Procedure					
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
slognmix2	38343*	38353*	38353*	38386*	0.52221*
Logn	39073	39077	39077	39090	5.86522
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD	CvM			
slognmix2	0.19843*	0.02728*			
Logn	66.93414	11.72703			
Note: The asterisk (*) marks the best model according to each column's criterion.					

The detailed results for the SLOGNMIX2 distribution are shown in [Output 9.8.2](#). According to the “Initial Parameter Values and Bounds” table, the initial value of  $\rho_2$  is not 0.5, indicating that a linear search was conducted to ensure  $\log(m_1) > \mu$ .

**Output 9.8.2** Detailed Estimation Results for the SLOGNMIX2 Distribution

The HPSEVERITY Procedure			
slognmix2 Distribution			
Distribution Information			
Name	slognmix2		
Description	Mixture of two lognormals with a log-scale parameter Mu		
Distribution Parameters	5		
Initial Parameter Values and Bounds			
Parameter	Initial Value	Lower Bound	Upper Bound
Mu	2.92006	-Infty	Infty
Sigma1	0.10455	1.05367E-8	Infty
P2	0.50000	1.05367E-8	1.00000
Rho2	0.72000	1.05367E-8	1.00000
Sigma2	0.81728	1.05367E-8	Infty

**Output 9.8.2** *continued*

Convergence Status				
Convergence criterion (GCONV=1E-8) satisfied.				
Optimization Summary				
Optimization Technique	Trust Region			
Iterations	7			
Function Calls	18			
Log Likelihood	-19171.5			
Parameter Estimates				
Parameter	Estimate	Standard Error	t Value	Approx Pr >  t
Mu	3.00922	0.01554	193.68	<.0001
Sigma1	0.49516	0.01451	34.13	<.0001
P2	0.40619	0.02600	15.62	<.0001
Rho2	0.37212	0.02038	18.26	<.0001
Sigma2	1.00019	0.02124	47.09	<.0001

By using the relationship that  $\mu_2 = \mu + \log(\rho_2)$ , you can see that the final parameter estimates are indeed close to the true parameter values that were used to simulate the input sample.

---

### Example 9.9: Predicting Mean and Value-at-Risk by Using Scoring Functions

If you work in the risk management department of an insurance company or a bank, then one of your primary applications of severity loss distribution models is to predict the value-at-risk (VaR), such that the probability of experiencing a loss value greater than the VaR is very low. The probability level at which VaR is measured is prescribed by industry regulations such as Basel III and Solvency II. The VaR level is usually specified in terms of  $(1 - \alpha)$ , where  $\alpha \in (0, 1)$  is the probability that a loss value exceeds the VaR. Typical VaR levels are 0.95, 0.975, and 0.995.

In addition to predicting VaR, which is regarded as an estimate of the worst-case loss, businesses are often interested in predicting the average loss by estimating either the mean or median of the distribution.

The estimation of the mean and VaR combined with the scale regression model is very potent tool for analyzing worst-case and average losses for various scenarios. For example, if the regressors that are used in a scale regression model represent some key macroeconomic and operational indicators, which are widely referred to as key risk indicators (KRIs), then you can analyze the VaR and mean loss estimates over various values for the KRIs to get a more comprehensive picture of the risk profile of your organization across various market and internal conditions.

This example illustrates the use of scoring functions to simplify the process of predicting the mean and VaR of scale regression models.

First, the following PROC FCMP steps define the functions to compute the mean for each of the 10 predefined distributions that are available in the SasHELP.SVRTDIST library:

```

/*----- Define distribution functions that compute the mean -----*/
proc fcmp library=sashelp.svrtldist outlib=work.means.scalemod;
  function BURR_MEAN(x, Theta, Alpha, Gamma);
    if not(Alpha * Gamma > 1) then
      return (.); /* first moment does not exist */
    return (Theta*gamma(1 + 1/Gamma)*gamma(Alpha - 1/Gamma)/gamma(Alpha));
  endsub;
  function EXP_MEAN(x, Theta);
    return (Theta);
  endsub;
  function GAMMA_MEAN(x, Theta, Alpha);
    return (Theta*Alpha);
  endsub;
  function GPD_MEAN(x, Theta, Xi);
    if not(Xi < 1) then
      return (.); /* first moment does not exist */
    return (Theta/(1 - Xi));
  endsub;
  function IGAUSS_MEAN(x, Theta, Alpha);
    return (Theta);
  endsub;
  function LOGN_MEAN(x, Mu, Sigma);
    return (exp(Mu + Sigma*Sigma/2.0));
  endsub;
  function PARETO_MEAN(x, Theta, Alpha);
    if not(Alpha > 1) then
      return (.); /* first moment does not exist */
    return (Theta/(Alpha - 1));
  endsub;
  function STWEEDIE_MEAN(x, Theta, Lambda, P);
    return (Theta* Lambda * (2 - P) / (P - 1));
  endsub;
  function TWEEDIE_MEAN(x, P, Mu, Phi);
    return (Mu);
  endsub;
  function WEIBULL_MEAN(x, Theta, Tau);
    return (Theta*gamma(1 + 1/Tau));
  endsub;
quit;

```

The following statements include the Work.Means library in the CMPLIB= system option and submit a PROC HPSEVERITY step to estimate the scale regression models for various distributions by using a lognormal sample in the Work.Test\_sev8 data set:

```
/*----- Fit all distributions and generate scoring functions -----*/
options cmplib=work.means;
proc hpseverity data=test_sev9 outest=est print=all;
    loss y;
    scalemodel x1-x5;
    dist _predefined_ stweedie;
    outscorelib outlib=scorefuncs commonpackage;
run;
```

The SAS statements that simulate the sample in the Work.Test\_sev8 data set are available in the PROC HPSEVERITY sample program *hsevex09.sas*. The OUTLIB= option in the OUTSCORELIB statement requests that the scoring functions be written to the Work.Scorefuncs library, and the COMMONPACKAGE option in the OUTSCORELIB statement requests that all the functions be written to the same package. Upon completion, PROC HPSEVERITY sets the CMPLIB system option to the following value:

```
(work.means sashelp.svrtdist work.scorefuncs)
```

The “All Fit Statistics” table in [Output 9.9.1](#) shows that the lognormal distribution’s scale model is the best and the inverse Gaussian’s scale model is a close second according to the likelihood-based statistics.

You can examine the scoring functions that are written to the Work.Scorefuncs library by using the FCMP Function Editor that is available in the Display Manager session of Base SAS when you select **Solutions**→**Analysis** from the main menu. For example, PROC HPSEVERITY automatically generates and submits the following PROC FCMP statements to define the scoring functions ‘SEV\_MEAN\_LOGN’ and ‘SEV\_QUANTILE\_IGAUSS’:

```
proc fcmp library=(work.means sashelp.svrtdist) outlib=work.scorefuncs.sevfit;
    function SEV_MEAN_LOGN(y, x{*});
        _logscale_=0;
        _logscale_ = _logscale_ + ( 7.64722278930350E-01 * x{1});
        _logscale_ = _logscale_ + ( 2.99209540369860E+00 * x{2});
        _logscale_ = _logscale_ + (-1.00788916253430E+00 * x{3});
        _logscale_ = _logscale_ + ( 2.58883602184890E-01 * x{4});
        _logscale_ = _logscale_ + ( 5.00927479793970E+00 * x{5});
        _logscale_ = _logscale_ + ( 9.95078833050690E-01);
        return (LOGN_MEAN(y, _logscale_, 2.31592981635590E-01));
    endsub;

    function SEV_QUANTILE_IGAUSS(y, x{*});
        _logscale_=0;
        _logscale_ = _logscale_ + ( 7.64581738373520E-01 * x{1});
        _logscale_ = _logscale_ + ( 2.99159055015310E+00 * x{2});
        _logscale_ = _logscale_ + (-1.00793496641510E+00 * x{3});
        _logscale_ = _logscale_ + ( 2.58870460543840E-01 * x{4});
        _logscale_ = _logscale_ + ( 5.00996884646730E+00 * x{5});
        _scale_ = 2.77854870591020E+00 * exp(_logscale_);
        return (IGAUSS_QUANTILE(y, _scale_, 1.81511227238720E+01));
    endsub;
quit;
```

**Output 9.9.1** Comparison of Fitted Scale Models for Mean and VaR Illustration

The HPSEVERITY Procedure					
All Fit Statistics					
Distribution	-2 Log Likelihood	AIC	AICC	BIC	KS
stweddie	460.65754	476.65754	476.95082	510.37440	10.44548
Burr	451.42238	467.42238	467.71565	501.13924	10.32782
Exp	1515	1527	1527	1552	8.85827
Gamma	448.28222	462.28222	462.50986	491.78448	10.42272
Igauss	444.44512	458.44512	458.67276	487.94738	10.33028
Logn	444.43670*	458.43670*	458.66434*	487.93895*	10.37035
Pareto	1515	1529	1529	1559	8.85775*
Gpd	1515	1529	1529	1559	8.85827
Weibull	527.28676	541.28676	541.51440	570.78902	10.48084
Note: The asterisk (*) marks the best model according to each column's criterion.					
All Fit Statistics					
Distribution	AD		CvM		
stweddie	64571		37.07702		
Burr	42254		37.19808		
Exp	29917		23.98267		
Gamma	63712		37.19450		
Igauss	83195		37.30880		
Logn	68631		37.18553		
Pareto	29916*		23.98149*		
Gpd	29917		23.98267		
Weibull	72814		36.36039		
Note: The asterisk (*) marks the best model according to each column's criterion.					

An important point to note is that the *dist\_MEAN* distribution functions are not available in the original definition of any of the distributions in the Sashelp.Svrtldist library. PROC HPSEVERITY detects the availability of those functions in the Work.Means library that is included in the value of the CMPLIB= system option just before submitting the PROC HPSEVERITY step. Each *dist\_MEAN* distribution function has a signature that matches the signature of a distribution function of the respective distribution, so PROC HPSEVERITY creates the corresponding scoring functions. Specifying the COMMONPACKAGE option in the OUTSCORELIB statement causes the name of the scoring function to take the form SEV\_MEAN\_*dist*. You can define any distribution function that has the desired signature to compute an estimate of your choice, include its library in the CMPLIB= system option, and then specify the OUTSCORELIB statement to generate the corresponding scoring functions.

To illustrate the use of scoring functions, let Work.Reginput contain the scoring data, where the values of regressors in each observation define one scenario. Scoring functions make it very easy to compute the mean

and VaR of each distribution's scale model for each of the scenarios, as the following steps illustrate for the lognormal and inverse Gaussian distributions:

```

/*--- Set VaR level ---*/
%let varLevel=0.975;

/*--- Compute scores (mean and var) for the      ---
--- scoring data by using the scoring functions ---*/
data scores;
    array x{*} x1-x5;
    set reginput;

    igauss_mean = sev_mean_igauss(., x);
    igauss_var   = sev_quantile_igauss(&varLevel, x);
    logn_mean    = sev_mean_logn(., x);
    logn_var     = sev_quantile_logn(&varLevel, x);
run;

```

The preceding steps use a VaR level of 97.5%.

The following DATA step accomplishes the same task by reading the parameter estimates that were written to the Work.Est data set by the previous PROC HPSEVERITY step:

```

/*--- Compute scores (mean and var) for the      ---
--- scoring data by using the OUTEST= data set ---*/
data scoresWithoutest(keep=x1-x5 igauss_mean igauss_var logn_mean logn_var);
    array _x{*} x1-x5;
    array _xparmIgauss{5} _temporary_;
    array _xparmLogn{5} _temporary_;

    retain _Theta0_ Alpha0;
    retain _Mu0_ Sigma0;
    *--- read parameter estimates for igauss and logn models ---*;
    if (_n_ = 1) then do;
        set est(where=(upcase(_MODEL_)= 'IGAUSS' and _TYPE_='EST'));
        _Theta0_ = Theta; Alpha0 = Alpha;
        do _i_=1 to dim(_x_);
            if (_x_(_i_) = .R) then _xparmIgauss_(_i_) = 0;
            else _xparmIgauss_(_i_) = _x_(_i_);
        end;

        set est(where=(upcase(_MODEL_)= 'LOGN' and _TYPE_='EST'));
        _Mu0_ = Mu; Sigma0 = Sigma;
        do _i_=1 to dim(_x_);
            if (_x_(_i_) = .R) then _xparmLogn_(_i_) = 0;
            else _xparmLogn_(_i_) = _x_(_i_);
        end;
    end;

    set reginput;

    *--- predict mean and VaR for inverse Gaussian ---*;
    * first compute X'*beta for inverse Gaussian *;
    _xbeta_ = 0.0;

```



```

do _i_ = 1 to dim(_x_);
    _xbeta_ = _xbeta_ + _xparmIgauss_(_i_) * _x_(_i_);
end;
* now compute scale for inverse Gaussian *;
_SCALE_ = _Theta0_ * exp(_xbeta_);
igauss_mean = igauss_mean(., _SCALE_, Alpha0);
igauss_var = igauss_quantile(&varLevel, _SCALE_, Alpha0);

*--- predict mean and VaR for lognormal      ---*;
* first compute X'*beta for lognormal*;
_xbeta_ = 0.0;
do _i_ = 1 to dim(_x_);
    _xbeta_ = _xbeta_ + _xparmLogn_(_i_) * _x_(_i_);
end;
* now compute Mu=log(scale) for lognormal *;
_MU_ = _Mu0_ + _xbeta_;
logn_mean = logn_mean(., _MU_, Sigma0);
logn_var = logn_quantile(&varLevel, _MU_, Sigma0);
run;

```

The “Values Comparison Summary” table in [Output 9.9.2](#) shows that the difference between the estimates that are produced by both methods is within the acceptable machine precision. However, the comparison of the DATA step complexity of each method clearly shows that the method that uses the scoring functions is much easier because it saves a lot of programming effort. Further, new distribution functions, such as the *dist\_MEAN* functions that are illustrated here, are automatically discovered and converted to scoring functions by PROC HPSEVERITY. That enables you to focus your efforts on writing the distribution function that computes your desired score, which needs to be done only once. Then, you can create and use the corresponding scoring functions multiple times with much less effort.

**Output 9.9.2** Comparison of Mean and VaR Estimates of Two Scoring Methods

<p style="text-align: center;">The COMPARE Procedure</p> <p style="text-align: center;">Comparison of WORK.SCORESWITHOUTEST with WORK.SCORES</p> <p style="text-align: center;">(Method=RELATIVE(0.0222), Criterion=1.0E-12)</p> <p>NOTE: All values compared are within the equality criterion used. However, 40 of the values compared are not exactly equal.</p>
---

## References

- D’Agostino, R. B. and Stephens, M., eds. (1986), *Goodness-of-Fit Techniques*, New York: Marcel Dekker.
- Danielsson, J., de Haan, L., Peng, L., and de Vries, C. G. (2001), “Using a Bootstrap Method to Choose the Sample Fraction in Tail Index Estimation,” *Journal of Multivariate Analysis*, 76, 226–248.
- Dunn, P. K. and Smyth, G. K. (2005), “Series Evaluation of Tweedie Exponential Dispersion Model Densities,” *Statistics and Computing*, 15, 267–280.

- Frydman, H. (1994), "A Note on Nonparametric Estimation of the Distribution Function from Interval-Censored and Truncated Observations," *Journal of the Royal Statistical Society, Series B*, 56, 71–74.
- Gentleman, R. and Geyer, C. J. (1994), "Maximum Likelihood for Interval Censored Data: Consistency and Computation," *Biometrika*, 81, 618–623.
- Hill, B. M. (1975), "A Simple General Approach to Inference about the Tail of a Distribution," *Annals of Statistics*, 3, 1163–1173.
- Jørgensen, B. (1987), "Exponential Dispersion Models (with discussion)," *Journal of the Royal Statistical Society, Series B*, 49, 127–162.
- Kaplan, E. L. and Meier, P. (1958), "Nonparametric Estimation from Incomplete Observations," *Journal of the American Statistical Association*, 53, 457–481.
- Klein, J. P. and Moeschberger, M. L. (1997), *Survival Analysis: Techniques for Censored and Truncated Data*, New York: Springer-Verlag.
- Klugman, S. A., Panjer, H. H., and Willmot, G. E. (1998), *Loss Models: From Data to Decisions*, New York: John Wiley & Sons.
- Koziol, J. A. and Green, S. B. (1976), "A Cramér–von Mises Statistic for Randomly Censored Data," *Biometrika*, 63, 466–474.
- Lai, T. L. and Ying, Z. (1991), "Estimating a Distribution Function with Truncated and Censored Data," *Annals of Statistics*, 19, 417–442.
- Lynden-Bell, D. (1971), "A Method of Allowing for Known Observational Selection in Small Samples Applied to 3CR Quasars," *Monthly Notices of the Royal Astronomical Society*, 155, 95–118.
- Turnbull, B. W. (1976), "The Empirical Distribution Function with Arbitrarily Grouped, Censored, and Truncated Data," *Journal of the Royal Statistical Society, Series B*, 38, 290–295.
- Tweedie, M. C. K. (1984), "An Index Which Distinguishes between Some Important Exponential Families," in J. K. Ghosh and J. Roy, eds., *Statistics: Applications and New Directions—Proceedings of the Indian Statistical Institute Golden Jubilee International Conference*, 579–604.

# Subject Index

- bounds on parameter estimates, 200
- BOUNDS statement, 200
- BY groups
  - HPCDM procedure, 68
  - HPCOUNTREG procedure, 138
  - HPSEVERITY procedure, 261
- censored regression models
  - HPQLIM procedure, 212
- censoring and truncation
  - HPSEVERITY procedure, 284
- Clayton copula, 121
- compound distribution modeling
  - HPCDM procedure, 44
- conjugate-gradient
  - optimization methods, 137
- copula
  - Clayton, 121
  - Frank, 121
  - Gumbel, 121
  - normal, 119
  - Student's  $t$ , 120
- covariates
  - heteroscedasticity models, 203
- defining a custom objective function
  - HPSEVERITY procedure, 332
- defining a custom probability distribution
  - HPSEVERITY procedure, 308
- dependence measures, 118
- descriptive statistics
  - HPCDM procedure, 90
- double-dogleg
  - optimization methods, 137
- empirical distribution function
  - HPSEVERITY procedure, 295
- fitting custom probability distributions
  - HPSEVERITY procedure, 261
- Frank copula, 121
- gamma distribution
  - definition of (HPQLIM), 218
  - HPQLIM procedure, 218
- Gaussian distribution
  - definition of (HPQLIM), 219
  - HPQLIM procedure, 219
- Gumbel copula, 121
- heteroscedasticity models
  - covariates, 203
- HPCDM procedure
  - BY groups, 68
  - descriptive statistics, 90
  - ODS graph names, 97
  - ODS table names, 95
  - parameter perturbation analysis, 89
  - scenario analysis, 74
  - simulating aggregate adjusted loss distribution, 82
  - simulating aggregate loss distribution, 75
- HPCOPULA procedure
  - multithreading, 117
  - overview, 113
  - syntax, 115
- HPCOUNTREG procedure
  - bounds on parameter estimates, 137
  - BY groups, 138
  - multithreading, 140
  - output table names, 152
  - restrictions on parameter estimates, 141
  - syntax, 131
- HPPANEL procedure
  - ID variables, 165
  - linear hypothesis testing, 176
  - multithreading, 167
  - one-way fixed-effects model, 170
  - one-way random-effects model, 173
  - output data sets, 177
  - output table names, 179
  - printed output, 178
  - specification tests, 176
  - two-way fixed-effects model, 171
  - two-way random-effects model, 174
- HPQLIM procedure, 186
  - BY groups, 200
  - censored regression models, 212
  - frontier, 213
  - gamma distribution, 218
  - Gaussian distribution, 219
  - heteroscedasticity, 214
  - inverse gamma distribution, 219
  - limited dependent variable models, 212
  - logit, 186
  - multithreading, 208
  - normal distribution, 219
  - ordinal discrete choice modeling, 211
  - output, 220

- output ODS Graphics table names, 226
- output table names, 225
- probit, 186
- selection, 186
- standard distributions, 218
- syntax, 189
- $t$  distribution, 219
- tests on parameters, 215
- Tobit, 186
- truncated regression models, 212
- uniform distribution, 219
- HPSEVERITY procedure
  - BY groups, 261
  - censoring and truncation, 284
  - defining a custom objective function, 332
  - defining a custom probability distribution, 308
  - empirical distribution function, 295
  - ODS table names, 339
  - predefined distributions, 274
  - predefined utility functions, 320
  - scale regression model, 290
  - scoring functions, 325
  - statistics of fit, 301
  - Tweedie distribution, 276
- ID variables
  - HPPANEL procedure, 165
- inverse gamma distribution
  - HPQLIM procedure, 219
- inverse gamma distribution
  - definition of (HPQLIM), 219
- Kaplan-Meier's EDF estimator
  - HPSEVERITY procedure, 297
- Lagrange multiplier test
  - nonlinear hypotheses, 215
- limited dependent variable models
  - HPQLIM procedure, 212
- linear hypothesis testing, 176
  - HPPANEL procedure, 176
- logit
  - HPQLIM procedure, 186
- loss distribution modeling
  - HPSEVERITY procedure, 238
- measure
  - dependence, 118
- modified Kaplan-Meier's EDF estimator
  - HPSEVERITY procedure, 298
- multithreading
  - HPCOPULA procedure, 117
  - HPCOUNTREG procedure, 140
  - HPPANEL procedure, 167
  - HPQLIM procedure, 208
- Newton-Raphson
  - optimization methods, 137, 194
- Newton-Raphson method, 194
- Newton-Raphson Ridge
  - optimization methods, 137
- none
  - optimization methods, 137
- nonlinear hypotheses
  - Lagrange multiplier test, 215
- normal copula, 119
- normal distribution
  - definition of (HPQLIM), 219
  - HPQLIM procedure, 219
- ODS graph names
  - HPCDM procedure, 97
- ODS table names
  - HPCDM procedure, 95
  - HPSEVERITY procedure, 339
- one-way fixed-effects model, 170
  - HPPANEL procedure, 170
- one-way random-effects model, 173
  - HPPANEL procedure, 173
- optimization methods
  - conjugate-gradient, 137
  - double-dogleg, 137
  - Newton-Raphson, 137, 194
  - Newton-Raphson Ridge, 137
  - none, 137
  - quasi-Newton, 137
  - trust region, 137, 194
- ordinal discrete choice modeling
  - HPQLIM procedure, 211
- output data sets
  - HPPANEL procedure, 177
- output ODS Graphics table names
  - HPQLIM procedure, 226
- output table names
  - HPCOUNTREG procedure, 152
  - HPPANEL procedure, 179
  - HPQLIM procedure, 225
- parameter perturbation analysis
  - HPCDM procedure, 89
- printed output
  - HPPANEL procedure, 178
- prior distribution
  - distribution specification (HPQLIM), 208
- probit
  - HPQLIM procedure, 186
- quasi-Newton
  - optimization methods, 137
- quasi-Newton method, 194

- scale regression model
  - HPSEVERITY procedure, 290
- scenario analysis
  - HPCDM procedure, 74
- scoring functions
  - HPSEVERITY procedure, 325
- selection
  - HPQLIM procedure, 186
- simulating aggregate adjusted loss distribution
  - HPCDM procedure, 82
- simulating aggregate loss distribution
  - HPCDM procedure, 75
- Sklar's theorem, 118
- specification tests
  - HPPANEL procedure, 176
- standard distributions
  - HPQLIM procedure, 218
- statistics of fit
  - HPSEVERITY procedure, 301
- Student's  $t$  copula, 120
- $t$  distribution
  - definition of (HPQLIM), 219
  - HPQLIM procedure, 219
- Tobit
  - HPQLIM procedure, 186
- truncated regression models
  - HPQLIM procedure, 212
- trust region
  - optimization methods, 137, 194
- trust region method, 194
- Turnbull's EDF estimator
  - HPSEVERITY procedure, 298
- Tweedie distribution
  - HPSEVERITY procedure, 276
- two-way fixed-effects model, 171
  - HPPANEL procedure, 171
- two-way random-effects model, 174
  - HPPANEL procedure, 174
- uniform distribution
  - definition of (HPQLIM), 219
  - HPQLIM procedure, 219



# Syntax Index

- ADJSAMPLEVAR= option
  - OUTPUT statement (HPCDM), 70
- ADJUSTEDSEVERITY= option
  - PROC HPCDM statement, 63
- ALL option
  - TEST statement (HPPANEL), 168
  - TEST statement (HPQLIM), 210
- BAYES statement
  - HPQLIM procedure, 196
- BETA
  - PRIOR statement (HPQLIM), 208
- BOUNDS statement
  - HPCOUNTREG procedure, 137
  - HPQLIM procedure, 200
- BY statement
  - HPCDM procedure, 68
  - HPCOUNTREG procedure, 138
  - HPQLIM procedure, 200
  - HPSEVERITY procedure, 261
- CENSORED option
  - ENDOGENOUS statement (HPQLIM), 201, 205
- COMMIT= option
  - PERFORMANCE statement (high-performance analytical procedures), 39
- COMMONPACKAGE option
  - OUTSCORELIB statement (HPSEVERITY), 270
- CONDITIONAL
  - OUTPUT statement (HPQLIM), 207
- COPYVAR= option
  - OUTPUT statement (HPCOUNTREG), 140
  - OUTPUT statement (HPPANEL), 166
  - OUTPUT statement (HPQLIM), 207
- CORRB option
  - HPQLIM procedure, 192
  - MODEL statement, 139
  - MODEL statement (HPPANEL), 165
  - PROC HPCOUNTREG statement, 134
- CORROUT option
  - PROC HPCOUNTREG statement, 134
  - PROC HPPANEL statement, 164
  - PROC HPQLIM statement, 192
- COST option
  - ENDOGENOUS statement (HPQLIM), 202, 206
- COUNT= option
  - EXTERNALCOUNTS statement (HPCDM), 69
- COUNTSTORE= option
  - PROC HPCDM statement, 63
- COVB option
  - HPQLIM procedure, 193
  - MODEL statement, 139
  - MODEL statement (HPPANEL), 165
  - PROC HPCOUNTREG statement, 134
- COVEST= option
  - HPQLIM procedure, 193
  - PROC HPCOUNTREG statement, 135
- COVOUT option
  - PROC HPCOUNTREG statement, 134
  - PROC HPPANEL statement, 164
  - PROC HPQLIM statement, 192
  - PROC HPSEVERITY statement, 254
- CRITERION= option
  - PROC HPSEVERITY statement, 257
- DATA= option
  - PROC HPCDM statement, 64
  - PROC HPCOUNTREG statement, 134
  - PROC HPPANEL statement, 164
  - PROC HPQLIM statement, 192
  - PROC HPSEVERITY statement, 254
- DATASERVER= option
  - PERFORMANCE statement (high-performance analytical procedures), 39
- DEFINE statement
  - HPCOPULA procedure, 115
- DETAILS option
  - PERFORMANCE statement (high-performance analytical procedures), 40
  - PERFORMANCE statement (HPCOPULA), 117
  - PERFORMANCE statement (HPCOUNTREG), 141
  - PERFORMANCE statement (HPPANEL), 167
  - PERFORMANCE statement (HPQLIM), 208
- DFMIXTURE= option
  - SCALEMODEL statement (HPSEVERITY), 272
- DIAGNOSTICS= option
  - BAYES statement (HPQLIM), 196
- DISCRETE option
  - ENDOGENOUS statement (HPQLIM), 201, 205
- DIST statement
  - HPSEVERITY procedure, 261
- DIST= option
  - HPCOUNTREG statement (HPCOUNTREG), 139
  - MODEL statement (HPCOUNTREG), 139
- DISTBY statement

- HPCDM procedure, 68
- DISTRIBUTION= option
  - ENDOGENOUS statement (HPQLIM), 201, 205
- EMPIRICALCDF= option
  - PROC HPSEVERITY statement, 258
- ERRSTD
  - OUTPUT statement (HPQLIM), 207
- EXPECTED
  - OUTPUT statement (HPQLIM), 207
- EXTERNALCOUNTS statement
  - HPCDM procedure, 69
- FIXONE option
  - MODEL statement (HPPANEL), 165
- FIXONETIME option
  - MODEL statement (HPPANEL), 165
- FIXTWO option
  - MODEL statement (HPPANEL), 165
- FREQ statement
  - HPCOUNTREG procedure, 138
- FRONTIER option
  - ENDOGENOUS statement (HPQLIM), 202, 206
- GAMMA
  - PRIOR statement (HPQLIM), 208
- GRIDHOST= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- GRIDMODE= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- GRIDTIMEOUT= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- high-performance analytical procedures,
  - PERFORMANCE statement, 39
  - COMMIT= option, 39
  - DATASERVER= option, 39
  - DETAILS option, 40
  - GRIDHOST= option, 40
  - GRIDMODE= option, 40
  - GRIDTIMEOUT= option, 40
  - HOST= option, 40
  - INSTALL= option, 40
  - INSTALLLOC= option, 40
  - LASR= option, 40
  - LASRSERVER= option, 40
  - MODE= option, 40
  - NNODES= option, 41
  - NODES= option, 41
  - NTHREADS= option, 42
  - THREADS= option, 42
  - TIMEOUT= option, 40
- HOST= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- HPCDM procedure, 61
  - DISTBY statement, 68
  - EXTERNALCOUNTS statement, 69
  - OUTPUT statement, 69
  - OUTSUM statement, 70
  - PERFORMANCE statement, 73
  - SEVERITYMODEL statement, 73
  - syntax, 61
- HPCDM procedure, EXTERNALCOUNTS statement
  - COUNT= option, 69
  - ID= option, 69
- HPCDM procedure, OUTPUT statement
  - ADJSAMPLEVAR= option, 70
  - OUT= option, 69
  - PERTURBOUT option, 70
  - SAMPLEVAR= option, 70
- HPCDM procedure, OUTSUM statement
  - OUT= option, 70
  - PCTLNAME= option, 72
  - PCTLNDEC= option, 73
  - PCTLPTS= option, 72
- HPCDM procedure, PROC HPCDM statement, 63
  - ADJUSTEDSEVERITY= option, 63
  - COUNTSTORE= option, 63
  - DATA= option, 64
  - NOPRINT option, 64
  - NPERTURBEDSAMPLES= option, 64
  - NREPLICATES= option, 64
  - PCTLDEF= option, 65
  - PLOTS= option, 65
  - PRINT= option, 66
  - SEED= option, 67
  - SEVERITYEST= option, 67
  - VARDEF= option, 67
- HPCOPULA procedure, 115
  - DEFINE statement, 115
  - PERFORMANCE statement, 117
  - PROC HPCOPULA statement, 115
  - SIMULATE statement, 116
  - syntax, 115
  - VAR Statement, 117
- HPCOPULA procedure, PERFORMANCE statement, 117
- HPCOUNTREG procedure, 131
  - PERFORMANCE statement, 140
  - syntax, 131
- HPCOUNTREG procedure, PERFORMANCE statement, 140
- HPCOUNTREG procedure, WEIGHT statement, 141
- HPPANEL procedure, 163
  - PERFORMANCE statement, 167



- syntax, 163
- HPPANEL procedure, PERFORMANCE statement, 167
- HPQLIM procedure, 189
  - PERFORMANCE statement, 208
  - PRIOR statement, 208
  - syntax, 189
- HPQLIM procedure, FREQ statement, 203
- HPQLIM procedure, PERFORMANCE statement, 208
- HPQLIM procedure, TEST statement, 209
- HPQLIM procedure, WEIGHT statement, 210
- HPSEVERITY procedure, 251
  - DIST statement, 261
  - LOSS statement, 263
  - NLOPTIONS statement, 265
  - OUTSCORELIB statement, 269
  - PERFORMANCE statement, 271
  - SCALEMODEL statement, 272
  - syntax, 251
  - WEIGHT statement, 273
- HPSEVERITY procedure, DIST statement
  - INIT= option, 262
  - LISTONLY option, 262
  - VALIDATEONLY option, 262
- HPSEVERITY procedure, LOSS statement
  - LEFTCENSORED= option, 263
  - LEFTTRUNCATED= option, 264
  - RIGHTCENSORED= option, 264
  - RIGHTTRUNCATED= option, 265
- HPSEVERITY procedure, OUTSCORELIB statement
  - COMMONPACKAGE option, 270
  - OUTBYID= option, 270
  - OUTLIB= option, 269
- HPSEVERITY procedure, PROC HPSEVERITY statement, 254
  - COVOUT option, 254
  - CRITERION= option, 257
  - DATA= option, 254
  - EDF=AUTO option, 258
  - EDF=KAPLANMEIER option, 259
  - EDF=MODIFIEDKM option, 259
  - EDF=NOTURNBULL option, 259
  - EDF=STANDARD option, 259
  - EDF=TURNBULL option, 259
  - EMPIRICALCDF= option, 258
  - INEST= option, 254
  - INITSAMPLE option, 255
  - NOPRINT option, 255
  - OBJECTIVE= option, 260
  - OUTEST= option, 255
  - OUTMODELINFO= option, 255
  - OUTSTAT= option, 256
  - PRINT= option, 256
  - VARDEF= option, 257
  - VERBOSE= option, 257
- HPSEVERITY procedure, SCALEMODEL statement
  - DFMIXTURE= option, 272
  - OFFSET= option, 273
- ID statement
  - HPPANEL procedure, 165
- ID= option
  - EXTERNALCOUNTS statement (HPCDM), 69
- IGAMMA
  - PRIOR statement (HPQLIM), 208
- INEST= option
  - PROC HPSEVERITY statement, 254
- INIT statement
  - HPCOUNTREG procedure, 138
  - HPQLIM procedure, 204
- INIT= option
  - DIST statement (HPSEVERITY), 262
- INITSAMPLE option
  - PROC HPSEVERITY statement, 255
- INSTALL= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- INSTALLLOC= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- LASR= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- LASRSERVER= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- LEFTCENSORED= option
  - LOSS statement (HPSEVERITY), 263
- LEFTTRUNCATED= option
  - LOSS statement (HPSEVERITY), 264
- LISTONLY option
  - DIST statement (HPSEVERITY), 262
- LM option
  - TEST statement (HPPANEL), 168
  - TEST statement (HPQLIM), 210
- LOSS statement
  - HPSEVERITY procedure, 263
- LOWERBOUND= option
  - ENDOGENOUS statement (HPQLIM), 202, 205, 206
- LR option
  - TEST statement (HPPANEL), 168
  - TEST statement (HPQLIM), 210
- MARGINAL
  - OUTPUT statement (HPQLIM), 207
- MAXTUNE= option
  - BAYES statement (HPQLIM), 197

- METHOD= option
  - PROC HPCOUNTREG statement, 137
  - PROC HPQLIM statement, 194
- MILLS
  - OUTPUT statement (HPQLIM), 207
- MINTUNE= option
  - BAYES statement (HPQLIM), 197
- MODE= option
  - PERFORMANCE statement (high-performance analytical procedures), 40
- MODEL statement
  - HPCOUNTREG procedure, 138
  - HPPANEL procedure, 165
  - HPQLIM procedure, 204
- NBI= option
  - BAYES statement (HPQLIM), 197
- NLOPTIONS statement
  - HPSEVERITY procedure, 265
- NMC= option
  - BAYES statement (HPQLIM), 198
- NNODES= option
  - PERFORMANCE statement (high-performance analytical procedures), 41
- NODES option
  - PERFORMRANCE statement (HPCOPULA), 117
  - PERFORMRANCE statement (HPPANEL), 167
  - PERFORMRANCE statement (HPQLIM), 208
- NODES= option
  - PERFORMANCE statement (high-performance analytical procedures), 41
  - PERFORMANCE statement (HPCOUNTREG), 141
- NOINT option
  - MODEL statement (HPCOUNTREG), 139
  - MODEL statement (HPPANEL), 165
  - MODEL statement (HPQLIM), 204
- NONORMALIZE option
  - WEIGHT statement (HPCOUNTREG), 142
  - WEIGHT statement (HPQLIM), 211
- NOPRINT option
  - MODEL statement (HPPANEL), 165
  - PROC HPCDM statement, 64
  - PROC HPCOUNTREG statement, 135, 139
  - PROC HPQLIM statement, 192
  - PROC HPSEVERITY statement, 255
- NORMAL
  - PRIOR statement (HPQLIM), 208
- NPERTURBEDSAMPLES= option
  - PROC HPCDM statement, 64
- NREPLICATES= option
  - PROC HPCDM statement, 64
- NTHREADS option
  - PERFORMRANCE statement (HPCOPULA), 117
- PERFORMRANCE statement (HPPANEL), 167
- PERFORMRANCE statement (HPQLIM), 208
- NTHREADS= option
  - PERFORMANCE statement (high-performance analytical procedures), 42
  - PERFORMANCE statement (HPCOUNTREG), 141
- NTU= option
  - BAYES statement (HPQLIM), 198
- OBJECTIVE= option
  - PROC HPSEVERITY statement, 260
- OFFSET= option
  - MODEL statement (HPCOUNTREG), 139
  - SCALEMODEL statement (HPSEVERITY), 273
- ORDER= option
  - ENDOGENOUS statement (HPQLIM), 201, 205
- OUT= option
  - OUTPUT statement (HPCDM), 69
  - OUTPUT statement (HPCOUNTREG), 140
  - OUTPUT statement (HPPANEL), 166
  - OUTPUT statement (HPQLIM), 207
  - OUTSUM statement (HPCDM), 70
- OUTBYID= option
  - OUTSCORELIB statement (HPSEVERITY), 270
- OUTCORR option
  - PROC HPPANEL statement, 164
- OUTCOV option
  - PROC HPPANEL statement, 164
- OUTEST= option
  - PROC HPCOUNTREG statement, 134
  - PROC HPPANEL statement, 164, 177
  - PROC HPQLIM statement, 192
  - PROC HPSEVERITY statement, 255
- OUTLIB= option
  - OUTSCORELIB statement (HPSEVERITY), 269
- OUTMODELINFO= option
  - PROC HPSEVERITY statement, 255
- OUTPOST= option
  - BAYES statement (HPQLIM), 198
- OUTPUT statement
  - HPCDM procedure, 69
  - HPCOUNTREG procedure, 139
  - HPPANEL procedure, 166
  - HPQLIM procedure, 206
  - PROC HPPANEL statement, 177
- OUTSCORELIB statement
  - HPSEVERITY procedure, 269
- OUTSTAT= option
  - PROC HPSEVERITY statement, 256
- OUTSUM statement
  - HPCDM procedure, 70
- PCTLDEF= option

PROC HPCDM statement, 65  
 PCTLNAME= option  
   OUTSUM statement (HPCDM), 72  
 PCTLNDEC= option  
   OUTSUM statement (HPCDM), 73  
 PCTLPTS= option  
   OUTSUM statement (HPCDM), 72  
 PERFORMANCE statement  
   high-performance analytical procedures, 39  
   HPCDM procedure, 73  
   HPCOPULA procedure, 117  
   HPCOUNTREG procedure, 140  
   HPPANEL procedure, 167  
   HPQLIM procedure, 208  
   HPSEVERITY procedure, 271  
 PERTURBOUT option  
   OUTPUT statement (HPCDM), 70  
 PLOTS option  
   HPQLIM statement (HPQLIM), 195  
 PLOTS= option  
   PROC HPCDM statement, 65  
 PRED= option  
   OUTPUT statement (HPCOUNTREG), 140  
 PREDICTED  
   OUTPUT statement (HPPANEL), 166  
   OUTPUT statement (HPQLIM), 207  
 PRINT= option  
   PROC HPCDM statement, 66  
   PROC HPSEVERITY statement, 256  
 PRINTALL option  
   MODEL statement, 139  
   PROC HPCOUNTREG statement, 135  
   PROC HPQLIM statement, 192  
 PRINTFIXED option  
   MODEL statement (HPPANEL), 165  
 PRIOR statement  
   HPQLIM procedure, 208  
 PROB  
   OUTPUT statement (HPQLIM), 207  
 PROB= option  
   OUTPUT statement (HPCOUNTREG), 140  
 PROBALL  
   OUTPUT statement (HPQLIM), 207  
 PROBCOUNT option  
   OUTPUT statement (HPCOUNTREG), 140  
 PROBZERO= option  
   OUTPUT statement (HPCOUNTREG), 140  
 PROC HPCDM statement, 63, *see* HPCDM procedure  
 PROC HPCOPULA statement  
   HPCOPULA procedure, 115  
 PROC HPPANEL statement, 164  
 PROC HPSEVERITY statement, 254  
 PRODUCTION option  
   ENDOGENOUS statement (HPQLIM), 202, 206

PROPCOV= option  
   BAYES statement (HPQLIM), 198  
 RANONE option  
   MODEL statement (HPPANEL), 165  
 RANTWO option  
   MODEL statement (HPPANEL), 166  
 RESIDUAL  
   OUTPUT statement (HPPANEL), 166  
   OUTPUT statement (HPQLIM), 207  
 RESTRICT statement  
   HPCOUNTREG procedure, 141  
   HPQLIM procedure, 209  
 RIGHTCENSORED= option  
   LOSS statement (HPSEVERITY), 264  
 RIGHTTRUNCATED= option  
   LOSS statement (HPSEVERITY), 265  
 SAMPLEVAR= option  
   OUTPUT statement (HPCDM), 70  
 SAMPLING= option  
   BAYES statement (HPQLIM), 198  
 SCALEMODEL statement  
   HPSEVERITY procedure, 272  
 SEED= option  
   BAYES statement (HPQLIM), 198  
   PROC HPCDM statement, 67  
 SEVERITYEST= option  
   PROC HPCDM statement, 67  
 SEVERITYMODEL statement  
   HPCDM procedure, 73  
 SIMULATE statement  
   HPCOPULA procedure, 116  
 STATISTICS option  
   BAYES statement (HPQLIM), 199  
 T  
   PRIOR statement (HPQLIM), 209  
 TE1  
   OUTPUT statement (HPQLIM), 207  
 TE2  
   OUTPUT statement (HPQLIM), 207  
 THIN= option  
   BAYES statement (HPQLIM), 199  
 THREADS= option  
   PERFORMANCE statement (high-performance analytical procedures), 42  
 TIMEOUT= option  
   PERFORMANCE statement (high-performance analytical procedures), 40  
 TRUNCATED option  
   ENDOGENOUS statement (HPQLIM), 202, 206  
 UNIFORM  
   PRIOR statement (HPQLIM), 208

- UPPERBOUND= option
  - ENDOGENOUS statement (HPQLIM), 202, 205, 206
- VALIDATEONLY option
  - DIST statement (HPSEVERITY), 262
- VAR Statement
  - HPCOPULA procedure, 117
- VARDEF= option
  - PROC HPCDM statement, 67
  - PROC HPSEVERITY statement, 257
- VCOMP= option
  - MODEL statement (HPPANEL), 166
- VERBOSE= option
  - PROC HPSEVERITY statement, 257
- WALD option
  - TEST statement (HPPANEL), 168
  - TEST statement (HPQLIM), 210
- WEIGHT statement
  - HPSEVERITY procedure, 273
- XBETA
  - OUTPUT statement (HPQLIM), 207
- XBETA= option
  - OUTPUT statement (HPCOUNTREG), 140
- ZEROMODEL statement
  - HPCOUNTREG procedure, 142
- ZGAMMA= option
  - OUTPUT statement (HPCOUNTREG), 140