

# SAS/STAT® 14.2 User's Guide The SURVEYPHREG Procedure

This document is an individual chapter from SAS/STAT® 14.2 User's Guide.

The correct bibliographic citation for this manual is as follows: SAS Institute Inc. 2016. SAS/STAT® 14.2 User's Guide. Cary, NC: SAS Institute Inc.

## SAS/STAT<sup>®</sup> 14.2 User's Guide

Copyright © 2016, 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, NC 27513-2414

#### November 2016

 $SAS^{(0)}$  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. (1) indicates USA registration.

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

SAS software may be provided with certain third-party software, including but not limited to open-source software, which is licensed under its applicable third-party software license agreement. For license information about third-party software distributed with SAS software, refer to http://support.sas.com/thirdpartylicenses.

# Chapter 115 The SURVEYPHREG Procedure

# Contents

Overview: SURVEYPHREG Procedure	9338
Getting Started: SURVEYPHREG Procedure	9339
Syntax: SURVEYPHREG Procedure	9343
PROC SURVEYPHREG Statement	9344
BY Statement	9350
CLASS Statement	9350
CLUSTER Statement	9353
DOMAIN Statement	9353
ESTIMATE Statement	9354
FREQ Statement	9355
LSMEANS Statement	9355
LSMESTIMATE Statement	9356
MODEL Statement	9358
NLOPTIONS Statement	9362
OUTPUT Statement	9362
Programming Statements	9363
REPWEIGHTS Statement	9365
SLICE Statement	9367
STORE Statement	9367
STRATA Statement	9367
TEST Statement	9368
WEIGHT Statement	9368
Details: SURVEYPHREG Procedure	9369
Notation and Estimation	9369
Failure Time Distribution	9370
Time and CLASS Variables Usage	9371
Partial Likelihood Function for the Cox Model	9374
Specifying the Sample Design	9375
Missing Values	9377
Variance Estimation	9379
Taylor Series Linearization	9380
Bootstrap Method	9381
Balanced Repeated Replication (BRR) Method	9381
Jackknife Method	9383
Replicate Weights Method	9385
Degrees of Freedom	9386

Variance Adjustment Factors	86
Variance Ratios and Standard Error Ratios	87
Domain Analysis	88
Hypothesis Tests, Confidence Intervals, and Residuals	89
Testing the Global Null Hypothesis	89
Model Fit Statistics	89
Contrasts	90
Confidence Intervals	91
Hazard Ratios	91
Residuals	91
Output Data Sets	93
Displayed Output	95
ODS Table Names	98
ODS Graphics	99
Examples: SURVEYPHREG Procedure	99
Example 115.1: Analysis of Clustered Data	99
Example 115.2: Stratification, Clustering, and Unequal Weights	01
Example 115.3: Domain Analysis	06
Example 115.4: Variance Estimation by Using Replicate Weights	10
Example 115.5: A Test of the Proportional Hazards Assumption by Using the Pro-	
gramming Statements	12
References	13

# **Overview: SURVEYPHREG Procedure**

The SURVEYPHREG procedure performs regression analysis based on the Cox proportional hazards model for sample survey data. Cox's semiparametric model is widely used in the analysis of survival data to estimate hazard rates when adequate explanatory variables are available. The procedure provides design-based variance estimates, confidence intervals, and hypothesis tests concerning the parameters and model effects. See Chapter 3, "Introduction to Statistical Modeling with SAS/STAT Software," and Chapter 14, "Introduction to Survey Procedures," for an introduction to the basic concepts of survey data analysis; see Chapter 13, "Introduction to Survival Analysis Procedures," for an introduction to the basic concepts of survey analysis.

The survival time of each member of a finite population is assumed to follow its own hazard function,  $\lambda_i(t)$ , expressed as

 $\lambda_i(t) = \lambda(t; \mathbf{Z}_i(t)) = \lambda_0(t) \exp(\mathbf{Z}'_i(t)\boldsymbol{\beta})$ 

where  $\lambda_0(t)$  is an arbitrary and unspecified baseline hazard function,  $\mathbf{Z}_i(t)$  is the vector of explanatory variables for the *i*th population unit at time *t*, and  $\boldsymbol{\beta}$  is the vector of unknown regression parameters.

The finite population regression parameter  $\beta_N$  is defined as the maximizer of the partial log likelihood when the entire finite population is observed. The SURVEYPHREG procedure produces a sample-based estimate  $\hat{\beta}$  of the proportional hazards regression parameters  $\beta_N$  for the finite population by maximizing the partial pseudo-log-likelihood  $l_{\pi}(\boldsymbol{\beta}; \mathbf{Z}_{i}(t), t_{i})$  based on observed covariates  $\mathbf{Z}_{i}(t)$  and observed survival time  $t_{i}$ . The procedure also produces an estimate of the sampling variance  $V(\boldsymbol{\beta}|\mathcal{F}_{N})$ , which assumes that the values of the finite population  $\mathcal{F}_{N}$  are fixed. For statistical inference, PROC SURVEYPHREG incorporates complex survey sample designs, including designs with stratification, clustering, and unequal weighting.

The procedure also allows time-dependent explanatory variables. An explanatory variable is time-dependent if its value for any given individual can change over time. Time-dependent variables have many useful applications in survival analysis. You can include time-dependent variables such as blood pressure or blood chemistry measures that vary with time during the course of a study. You can also use time-dependent variables to test the validity of the proportional hazards model.

Several optimization techniques are available in SURVEYPHREG to maximize the log likelihood. Hazard ratio estimates can also be obtained along with parameter estimates. Sampling errors of the regression parameters and hazard ratios are computed by using either the Taylor series (linearization) method or one of the replication (resampling) methods that are based on complex sample designs (Binder 1983; Wolter 2007; Särndal, Swensson, and Wretman 1992; Binder 1992; Lohr 2010; Fuller 2009). These variance estimators essentially assume the finite population as fixed and estimate the variability due to the random sample selection mechanism.

The remaining sections of this chapter contain information about how to use PROC SURVEYPHREG, information about the underlying statistical methodology, and some applications of the procedure. The section "Getting Started: SURVEYPHREG Procedure" on page 9339 introduces PROC SURVEYPHREG with an example. The section "Syntax: SURVEYPHREG Procedure" on page 9343 describes the syntax of the procedure. The section "Details: SURVEYPHREG Procedure" on page 9369 summarizes the statistical techniques employed in PROC SURVEYPHREG. The section "Examples: SURVEYPHREG Procedure" on page 9399 includes some additional examples of useful applications. Experienced SAS/STAT software users might decide to proceed to the "Syntax" section, while other users might choose to read both the "Getting Started" and "Examples" sections before proceeding to "Syntax" and "Details."

# Getting Started: SURVEYPHREG Procedure

This section uses a data set that is obtained by stratified random sampling from a simulated finite population to illustrate some of the basic features of PROC SURVEYPHREG.

Suppose the library system for a small county wants to study the length of time that books are borrowed over a specified study period, adjusting for the age of the borrower and accounting for the fact that some books are never returned. Suppose there are 10 branch libraries in the county. Assume that a list of 11,617 (simulated) transactions is available for the study period October 1, 2008, to December 31, 2008, and assume that this list can be used as the sampling frame. A stratified random sample with replacement is used to select 100 transactions, where branch libraries are the strata. The total number of transactions within branches range from 510 to 2,011 for the study period. The total sample size of 100 transactions is allocated proportionally across branches based on the number of transactions. For each selected transaction, telephone interviews were conducted to find out additional characteristics of the borrower. The data set LibrarySurvey contains the following variables for all units (transactions) in the sample:

- Branch, the library branch from which the book was borrowed
- SampleWeight, the survey sampling weight for the transaction

- CheckOut, the date the book was borrowed
- CheckIn, the date the book was returned, with a missing value if the book was not returned by December 31, 2008
- Age, the age of the borrower

```
data LibrarySurvey;
   input Branch
                        2.
         SamplingWeight 7.2
        CheckOut date10.
        CheckIn
                      date10.
        Age;
  datalines;
 1 103.60 08NOV2008 13NOV2008 18
 1 103.60 010CT2008 070CT2008 30
 1 103.60 05NOV2008 06NOV2008 73
 1 103.60 250CT2008 260CT2008 53
1 103.60 09NOV2008 10NOV2008 55
2 127.50 10DEC2008 15DEC2008 39
                           . 33
 2 127.50 19DEC2008
2 127.50 26NOV2008 27NOV2008 41
 2 127.50 03NOV2008 07NOV2008 33
   ... more lines ...
10 118.35 14NOV2008 17NOV2008 29
10 118.35 11DEC2008 13DEC2008 35
10 118.35 21NOV2008 23NOV2008 46
data LibrarySurvey;
  set LibrarySurvey;
  Returned = (CheckIn ^= .);
  if (Returned) then
      lenBorrow = CheckIn
                                         - CheckOut;
   else
  lenBorrow = input('31Dec2008',date9.) - CheckOut;
run;
```

PROC SURVEYPHREG can be used to estimate the regression parameters of a proportional hazards model and the design-based variance of the estimated coefficients. The design-based variance is useful when the finite population is considered fixed, as in this example. See Lohr (2010) and Särndal, Swensson, and Wretman (1992) for details.

The following statements request a proportional hazards regression of lenBorrow on Age with Returned as the censor indicator. A transaction is considered to be censored if its check-in date is missing. The WEIGHT statement specifies the sampling weight variable (SamplingWeight), and the STRATA statement specifies the stratification variable (Branch).

```
proc surveyphreg data = LibrarySurvey;
  weight SamplingWeight;
   strata Branch;
   model lenBorrow*Returned(0) = Age;
run;
```

Summary information about the model, number of observations, survey design, censored values, and variance estimation method are shown in Output 115.1. The "Model Information" table summarizes the model you fit. The "Number of Observations" table displays the number of observations read and used by the procedure. This table also displays the sum of weights read and used. The sum of weights read (11,616.79) can be used as an estimator of the population size, and the sum of weights used can be used as an estimator of the respondent size in the population. The "Design Summary" table displays survey design information such as stratification and clustering. This example implements a stratified design with 10 strata. The "Censored Summary" and "Weighted Censored Summary" tables display the (weighted) number of censored and event units. Weighted counts can be used as estimators of the corresponding finite population quantities. For example, Output 115.1 shows that 10% of the sampled units are censored and an estimated 10.05% of the population units are censored.

#### Figure 115.1 Summary Statistics

# The SURVEYPHREG Procedure

Model	Informatio	n
Data Set	WORK.LIB	RARYSURVEY
Dependent Variable	lenBorrow	
Censoring Variable	Returned	
Censoring Value(s)	0	
Weight Variable	SamplingW	/eight
Stratum Variable	Branch	
Ties Handling	BRESLOW	/
Number of Observ	vations Rea	ad 100
Number of Observ	vations Use	ed 100
Sum of Weights R	Read	11616.79
Sum of Weights U	Jsed	11616.79
Desig Number o	n Summary of Strata	/ 10
Summary of th and Cen	ne Number Isored Valu	of Event es
		Percent
Total Event C	Censored C	Censored
100 90	10	10.00
Summary of the Event and (	Weighted I Censored V	Number of alues
		Percent
Total Even	t Censore	d Censored
11616.79 10449.2	2 1167.5	7 10.05

#### Figure 115.1 continued

Variance Estimation Method Taylor Series

Parameter estimates and their standard errors are shown in Output 115.2. The estimated regression coefficient is highly significant with a value of 0.062, indicating a positive association between age and the length of time books are borrowed (recall that these are simulated data). In this example, the procedure uses the STRATA and WEIGHT statements to incorporate stratification and unequal weighting, respectively, into variance estimation. The degrees of freedom are calculated as the number of sampling units (100) minus the number of strata (10). Note that the estimated variance reported in Output 115.2 ignores the finite population correction (*fpc*). You can use the TOTAL= or RATE= option in the PROC statement to include an *fpc* in your variance estimator.

Figure 115.2 Weighted Estimates and Their Standard Errors

Analysis of Maximum Likelihood Estimates						
Standard Hazard			Hazard			
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio
Age	90	0.061593	0.008366	7.36	<.0001	1.064

# Syntax: SURVEYPHREG Procedure

The following statements are available in the SURVEYPHREG procedure. Items within <> are optional.

**PROC SURVEYPHREG** < options > : BY variables; CLASS variable < (options) > < ... variable < (options) > > < / options > ; **CLUSTER** variables; **DOMAIN** variables < variable\* variable variable\* variable\* variable ... > ; **ESTIMATE** <'label' > estimate-specification < / options > ; FREQ variable : LSMEANS < model-effects > < / options > ; LSMESTIMATE model-effect Ismestimate-specification < / options>; **MODEL** response < \* censor(list) > = effects < / options > ; NLOPTIONS < options > ; **OUTPUT** < **OUT**=SAS-data-set> < keyword=name ... keyword=name > < / options>; Programming statements ; **REPWEIGHTS** variables < / options > : SLICE model-effect < / options > ; STRATA variables < / option > ; STORE < OUT= >item-store-name < / LABEL='label' > ; **TEST** < model-effects > < / options > ; WEIGHT variable;

The PROC SURVEYPHREG and MODEL statements are required. The CLASS statement, if present, must precede the MODEL statement.

The following sections describe the PROC SURVEYPHREG statement and then describe the other statements in alphabetical order.

The ESTIMATE, LSMEANS, LSMESTIMATE, SLICE, STORE, and TEST statements are also available in other procedures. Summary descriptions of functionality and syntax for these statements are provided in this chapter, and you can find full documentation about them in Chapter 19, "Shared Concepts and Topics."

# **PROC SURVEYPHREG Statement**

#### **PROC SURVEYPHREG** < options> ;

The PROC SURVEYPHREG statement invokes the SURVEYPHREG procedure. It also identifies the data set to be analyzed. Table 115.1 summarizes the *options* available in the PROC SURVEYPHREG statement.

Description
Names the input SAS data set
Treats missing values as a valid category
Suppresses all displayed output
Uses missing observations specified as not missing completely at random
Specifies the sort order of CLASS variables
Specifies the sampling rate
Specifies the total number of primary sampling units
Specifies the variance estimation method

Table 115.1 PROC SURVEYPHREG Statement Optic	ons
--	-----

You can specify the following options in the PROC SURVEYPHREG statement:

# DATA=SAS-data-set

names the SAS data set that contains the data to be analyzed. If you omit the DATA= option, the procedure uses the most recently created SAS data set.

# MISSING

treats missing values as a valid (nonmissing) category for all categorical variables, which include CLASS, STRATA, CLUSTER, and DOMAIN variables. By default, if you do not specify the MISSING option, an observation is excluded from the analysis if it has a missing value for any of these categorical variables. For more information, see the section "Missing Values" on page 9377.

#### NOPRINT

suppresses all displayed output. Note that this option temporarily disables the Output Delivery System (ODS); see Chapter 20, "Using the Output Delivery System," for more information.

#### NOMCAR

includes observations with missing values of the analysis variables that are specified in the MODEL statement as *not missing completely at random* (NOMCAR) for Taylor series variance estimation. When you specify the NOMCAR option, PROC SURVEYPHREG computes variance estimates by analyzing the nonmissing values as a domain (subpopulation), where the entire population includes both nonmissing and missing domains. See the section "Missing Values" on page 9377 for details.

By default, PROC SURVEYPHREG excludes an observation from analyses (and the corresponding variance computations) if that observation has a missing value for any of the variables in the MODEL statement. Note that if you specify the MISSING option for classification variables, then the procedure treats the missing values as a valid nonmissing level.

The NOMCAR option applies only to Taylor series variance estimation. Other replication methods do not use the NOMCAR option.

# ORDER=DATA | FORMATTED | FREQ | INTERNAL

specifies the sort order for the levels of the classification variables (which are specified in the CLASS statement).

This option applies to the levels for all classification variables, except when you use the (default) ORDER=FORMATTED option with numeric classification variables that have no explicit format. In that case, the levels of such variables are ordered by their internal value.

The ORDER= option can take the following values:

Value of ORDER=	Levels Sorted By
DATA	Order of appearance in the input data set
FORMATTED	External formatted value, except for numeric variables with no explicit format, which are sorted by their unformatted (internal) value
FREQ	Descending frequency count; levels with the most observations come first in the order
INTERNAL	Unformatted value

By default, ORDER=FORMATTED. For ORDER=FORMATTED and ORDER=INTERNAL, the sort order is machine-dependent.

For more information about sort order, see the chapter on the SORT procedure in the Base SAS Procedures Guide and the discussion of BY-group processing in SAS Language Reference: Concepts.

#### RATE=value | SAS-data-set

#### R=value | SAS-data-set

specifies the sampling rate, which PROC SURVEYPHREG uses to compute a finite population correction for Taylor series variance estimation. This option is ignored for bootstrap, BRR, and jackknife variance estimation.

If your sample design has multiple stages, you should specify the *first-stage sampling rate*, which is the ratio of the number of primary sampling units (PSUs) that are selected to the total number of PSUs in the population.

You can specify the sampling rate in either of the following ways:

- *value* specifies a nonnegative number to use for a nonstratified design or for a stratified design that has the same sampling rate in each stratum.
- SAS-data-set specifies a SAS-data-set that contains the stratification variables and the sampling rates for a stratified design that has different sampling rates in the strata. You must provide the sampling rates in the data set variable named \_RATE\_.

The sampling rates must be nonnegative numbers. You can specify *value* as a number between 0 and 1. Or you can specify *value* in percentage form as a number between 1 and 100, and PROC SURVEYPHREG converts that number to a proportion. The procedure treats the value 1 as 100% instead of 1%.

For more information, see the section "Population Totals and Sampling Rates" on page 9376.

If you do not specify the RATE= or TOTAL= option, then the Taylor series variance estimation does not include a finite population correction. You cannot specify both the TOTAL= option and the RATE= option in the same PROC SURVEYPHREG statement.

# TOTAL=value | SAS-data-set

N=value | SAS-data-set

specifies the total number of primary sampling units (PSUs) in the population, which PROC SUR-VEYPHREG uses to compute a finite population correction for Taylor series variance estimation. This option is ignored for bootstrap, BRR, and jackknife variance estimation.

You can specify the total number of PSUs in either of the following ways:

*value* specifies a positive number to use for a nonstratified design or for a stratified design that has the same population total in each stratum.

SAS-data-set specifies a SAS-data-set that contains the stratification variables and the population totals for a stratified design that has different population totals in the strata. You must provide the stratum totals in the data set variable named \_TOTAL\_.

The stratum totals must be positive numbers.

For more information, see the section "Population Totals and Sampling Rates" on page 9376.

If you do not specify the TOTAL= or RATE= option, then the Taylor series variance estimation does not include a finite population correction. You cannot specify both the TOTAL= option and the RATE= option in the same PROC SURVEYPHREG statement.

# VARMETHOD=method < (method-options) >

specifies the variance estimation *method*. PROC SURVEYPHREG provides the Taylor series method and balanced repeated replication (BRR), jackknife, and bootstrap replication (resampling) methods.

Table 115.2 summarizes the available *methods* and *method-options*.

method	Variance Estimation Method	method-options
BOOTSTRAP BRR	Bootstrap Balanced repeated replication	None CENTER=FULLSAMPLE   REPLICATES DETAILS FAY <=value > HADAMARD=SAS-data-set OUTWEIGHTS=SAS-data-set PRINTH REPS=number
JACKKNIFE	Jackknife	CENTER=FULLSAMPLE   REPLICATES DETAILS OUTJKCOEFS= <i>SAS-data-set</i> OUTWEIGHTS= <i>SAS-data-set</i>
TAYLOR	Taylor series linearization	None

Table 115.2 Variance Estimation Options

For VARMETHOD=BRR and VARMETHOD=JACKKNIFE, you can specify *method-options* in parentheses following the *method*.

By default, VARMETHOD=JACKKNIFE if you also specify a REPWEIGHTS statement; otherwise, VARMETHOD=TAYLOR by default.

You can specify the following *methods*:

#### BOOTSTRAP

requests variance estimation by the bootstrap method. When you specify this option, you must also provide bootstrap replicate weights by using a REPWEIGHTS statement; PROC SURVEYPHREG does not create bootstrap weights. For more information, see the section "Bootstrap Method" on page 9381.

#### BRR < (method-options) >

requests variance estimation by balanced repeated replication (BRR). The BRR method requires a stratified sample design with two primary sampling units (PSUs) in each stratum. If you specify the VARMETHOD=BRR option, you must also specify a STRATA statement unless you provide replicate weights with a REPWEIGHTS statement. See the section "Balanced Repeated Replication (BRR) Method" on page 9381 for details.

You can specify the following *method-options* in parentheses after the VARMETHOD=BRR option:

# **CENTER=FULLSAMPLE | REPLICATES**

defines how to compute the deviations for the BRR method. CENTER=FULLSAMPLE is the default, which computes the deviations of the replicate estimates from the full sample estimate. Alternatively, you can specify CENTER=REPLICATES to compute the deviations of the replicate estimates from the average of the replicate estimates. See the section "Balanced Repeated Replication (BRR) Method" on page 9381 for details.

#### DETAILS

displays the maximum likelihood estimates of model parameters for replicate samples when the replicate parameter estimates are available. A replicate sample might not provide useful parameter estimates (replicate estimates), for reasons such as nonconvergence of the optimization or inestimability of some parameters in that replicate sample.

#### FAY <=value>

requests Fay's method, which is a modification of the BRR method. See the section "Fay's BRR Method" on page 9382 for details.

You can specify the *value* of the Fay coefficient, which is used in converting the original sampling weights to replicate weights. The Fay coefficient must be a nonnegative number less than 1. By default, the value of the Fay coefficient equals 0.5.

# HADAMARD=SAS-data-set

# H=SAS-data-set

names a SAS data set that contains the Hadamard matrix for BRR replicate construction. If you do not provide a Hadamard matrix with the HADAMARD= *method-option*, PROC SURVEYPHREG generates an appropriate Hadamard matrix for replicate construction. See the sections "Balanced Repeated Replication (BRR) Method" on page 9381 and "Hadamard Matrix" on page 9383 for details.

If a Hadamard matrix of a given dimension exists, it is not necessarily unique. Therefore, if you want to use a specific Hadamard matrix, you must provide the matrix as a SAS data set in the HADAMARD= *method-option*.

In the HADAMARD= input data set, each variable corresponds to a column of the Hadamard matrix, and each observation corresponds to a row of the matrix. You can use any variable names in the HADAMARD= data set. All values in the data set must equal either 1 or -1. You must ensure that the matrix you provide is indeed a Hadamard matrix—that is, A'A = RI, where A is the Hadamard matrix of dimension R and I is an identity matrix. PROC SURVEYPHREG does not check the validity of the Hadamard matrix that you provide.

The HADAMARD= input data set must contain at least H variables, where H denotes the number of first-stage strata in your design. If the data set contains more than H variables, PROC SURVEYPHREG uses only the first H variables. Similarly, the HADAMARD= input data set must contain at least H observations.

If you do not specify the REPS= *method-option*, then the number of replicates is equal to the number of observations in the HADAMARD= input data set. If you specify the number of replicates—for example, REPS=*nreps*—then the first *nreps* observations in the HADAMARD= data set are used to construct the replicates.

You can specify the PRINTH *method-option* to display the Hadamard matrix that PROC SURVEYPHREG uses to construct replicates for BRR variance estimation.

# OUTWEIGHTS=SAS-data-set

names an output SAS data set to store the replicate weights that PROC SURVEYPHREG creates for BRR variance estimation. For more information about replicate weights, see the section "Balanced Repeated Replication (BRR) Method" on page 9381. For more information about the contents of the OUTWEIGHTS= data set, see the section "Replicate Weights Output Data Set" on page 9394.

The OUTWEIGHTS= *method-option* is not available when you provide replicate weights by using a REPWEIGHTS statement.

# PRINTH

displays the Hadamard matrix that is used to construct replicates for BRR variance estimation. When you provide the Hadamard matrix in the HADAMARD= *method-option*, PROC SURVEYPHREG displays only the rows and columns that are actually used to construct replicates. For more information, see the sections "Balanced Repeated Replication (BRR) Method" on page 9381 and "Hadamard Matrix" on page 9383.

The PRINTH *method-option* is not available when you provide replicate weights by using a REPWEIGHTS statement, because PROC SURVEYPHREG does not use a Hadamard matrix in this case.

#### **REPS**=number

specifies the number of replicates for BRR variance estimation. The value of *number* must be an integer greater than 1.

If you do not provide a Hadamard matrix by using the HADAMARD= *method-option*, the number of replicates should be greater than the number of strata and should be a multiple of

4. For more information, see the section "Balanced Repeated Replication (BRR) Method" on page 9381. If a Hadamard matrix cannot be constructed for the REPS= value that you specify, the value is increased until a Hadamard matrix of that dimension can be constructed. Therefore, it is possible for the actual number of replicates to be larger than the REPS= value that you specify.

If you provide a Hadamard matrix by using the HADAMARD= *method-option*, the value of REPS= must not be greater than the number of rows in the Hadamard matrix. If you provide a Hadamard matrix and do not specify the REPS= *method-option*, the number of replicates equals the number of rows in the Hadamard matrix.

If you do not specify the REPS= or HADAMARD= *method-option* and do not include a REPWEIGHTS statement, the number of replicates equals the smallest multiple of 4 that is greater than the number of strata.

If you provide replicate weights with a REPWEIGHTS statement, the procedure does not use the REPS= *method-option*. With a REPWEIGHTS statement, the number of replicates equals the number of REPWEIGHTS variables.

# JACKKNIFE | JK < (method-options) >

requests variance estimation by the delete-1 jackknife method. See the section "Jackknife Method" on page 9383 for details. If you provide replicate weights with a REPWEIGHTS statement, VARMETHOD=JACKKNIFE is the default variance estimation method. The JACKKNIFE method requires at least two primary sampling units (PSUs) in each stratum for stratified designs unless you provide replicate weights with a REPWEIGHTS statement.

You can specify the following method-options in parentheses following VARMETHOD=JACKKNIFE:

#### **CENTER=FULLSAMPLE | REPLICATES**

defines how to compute the deviations for the jackknife method. CENTER=FULLSAMPLE is the default, which computes the deviations of the replicate estimates from the full sample estimate. Alternatively, you can specify CENTER=REPLICATES to compute the deviations of the replicate estimates from the average of the replicate estimates. See the section "Jackknife Method" on page 9383 for details.

# DETAILS

displays the maximum likelihood estimates of model parameters for replicate samples when the replicate parameter estimates are available. A replicate sample might not provide useful parameter estimates (replicate estimates), for reasons such as nonconvergence of the optimization or inestimability of some parameters in that replicate sample.

#### OUTWEIGHTS=SAS-data-set

names an output SAS data set that contains replicate weights. See the section "Jackknife Method" on page 9383 for more information about replicate weights. See the section "Replicate Weights Output Data Set" on page 9394 for more details about the contents of the OUTWEIGHTS= data set.

The OUTWEIGHTS= *method-option* is not available when you provide replicate weights with a REPWEIGHTS statement.

#### OUTJKCOEFS=SAS-data-set

names an output SAS data set that contains jackknife coefficients. See the section "Jackknife Coefficients Output Data Set" on page 9394 for more details about the contents of the OUTJKCOEFS= data set.

# TAYLOR

requests Taylor series variance estimation. This is the default method if you do not specify the VARMETHOD= option or a REPWEIGHTS statement. See the section "Taylor Series Linearization" on page 9380 for more information.

# **BY Statement**

#### BY variables;

You can specify a BY statement with PROC SURVEYPHREG to obtain separate analyses of observations in groups that are defined by the BY variables. When a BY statement appears, the procedure expects the input data set to be sorted in order of the BY variables. If you specify more than one BY statement, only the last one specified is used.

If your input data set is not sorted in ascending order, use one of the following alternatives:

- Sort the data by using the SORT procedure with a similar BY statement.
- Specify the NOTSORTED or DESCENDING option in the BY statement for the SURVEYPHREG procedure. The NOTSORTED option does not mean that the data are unsorted but rather that the data are arranged in groups (according to values of the BY variables) and that these groups are not necessarily in alphabetical or increasing numeric order.
- Create an index on the BY variables by using the DATASETS procedure (in Base SAS software).

Note that using a BY statement provides completely separate analyses of the BY groups. It does not provide a domain (subpopulation) analysis, where the number of sampling units in the subpopulation is not known at the time the survey is designed. For such an analysis use the DOMAIN statement.

For more information about BY-group processing, see the discussion in SAS Language Reference: Concepts. For more information about the DATASETS procedure, see the discussion in the Base SAS Procedures Guide.

# **CLASS Statement**

# CLASS variable < (options) > ... < variable < (options) >> < / options >;

The CLASS statement names the classification variables to be used as explanatory variables in the analysis.

The CLASS statement must precede the MODEL statement. Most *options* can be specified either as individual variable *options* or as global *options*. You can specify *options* for each variable by enclosing the options in parentheses after the variable name. You can also specify global *options* for the CLASS statement by placing the *options* after a slash (/). Global *options* are applied to all the variables specified in the CLASS statement. If you specify more than one CLASS statement, the global *options* specified in any one CLASS statement apply to all CLASS statements. However, individual CLASS variable *options* override the global *options*. The following *options* are available:

# DESCENDING

# DESC

reverses the sort order of the classification variable. If both the DESCENDING and ORDER= options are specified, PROC SURVEYPHREG orders the categories according to the ORDER= option and then reverses that order.

# MISSING

treats missing values (".", .\_, .A, ..., .Z for numeric variables and blanks for character variables) as valid values for the CLASS variable.

# ORDER=DATA | FORMATTED | FREQ | INTERNAL

specifies the sort order for the levels of classification variables. This ordering determines which parameters in the model correspond to each level in the data, so the ORDER= option can be useful when you use the CONTRAST statement. By default, ORDER=FORMATTED. For ORDER=FORMATTED and ORDER=INTERNAL, the sort order is machine-dependent. When ORDER=FORMATTED is in effect for numeric variables for which you have supplied no explicit format, the levels are ordered by their internal values.

The following table shows how PROC SURVEYPHREG interprets values of the ORDER= option.

Value of ORDER=	Levels Sorted By
DATA	Order of appearance in the input data set
FORMATTED	External formatted values, except for numeric
	variables with no explicit format, which are sorted
	by their unformatted (internal) values
FREQ	Descending frequency count; levels with more
	observations come earlier in the order
INTERNAL	Unformatted value

For more information about sort order, see the chapter on the SORT procedure in the Base SAS Procedures Guide and the discussion of BY-group processing in SAS Language Reference: Concepts.

# **PARAM**=keyword

specifies the parameterization method for the classification variable or variables. If the PARAM= option is not specified together with any individual CLASS variable, then by default, PARAM=GLM. Otherwise, the default is PARAM=EFFECT. You can specify any of the *keywords* shown in the following table.

Design matrix columns are created from CLASS variables according to the corresponding coding schemes:

Value of PARAM=	Coding
EFFECT	Effect coding
GLM	Less-than-full-rank reference cell coding (this <i>keyword</i> can be used only in a global option)
ORDINAL THERMOMETER	Cumulative parameterization for an ordinal CLASS variable
POLYNOMIAL POLY	Polynomial coding
REFERENCE REF	Reference cell coding
ORTHEFFECT	Orthogonalizes PARAM=EFFECT coding
ORTHORDINAL ORTHOTHERM	Orthogonalizes PARAM=ORDINAL coding
ORTHPOLY	Orthogonalizes PARAM=POLYNOMIAL coding
ORTHREF	Orthogonalizes PARAM=REFERENCE coding

All parameterizations are full rank, except for the GLM parameterization. The REF= option in the CLASS statement determines the reference level for EFFECT and REFERENCE coding and for their orthogonal parameterizations. It also indirectly determines the reference level for a singular GLM parameterization through the order of levels.

If PARAM=ORTHPOLY or PARAM=POLY and the classification variable is numeric, then the ORDER= option in the CLASS statement is ignored, and the internal unformatted values are used. See the section "Other Parameterizations" on page 389 in Chapter 19, "Shared Concepts and Topics," for further details.

# REF='level' | keyword

specifies the reference level for PARAM=EFFECT, PARAM=REFERENCE, and their orthogonalizations. For PARAM=GLM, the REF= option specifies a level of the classification variable to be put at the end of the list of levels. This level thus corresponds to the reference level in the usual interpretation of the linear estimates with a singular parameterization.

For an individual variable REF= option (but not for a global REF= option), you can specify the *level* of the variable to use as the reference level. Specify the formatted value of the variable if a format is assigned. For a global or individual variable REF= option, you can use one of the following *keywords*. The default is REF=LAST.

- **FIRST** designates the first ordered level as reference.
- **LAST** designates the last ordered level as reference.

# TRUNCATE<=n>

specifies the length n of CLASS variable values to use in determining CLASS variable levels. The default is to use the full formatted length of the CLASS variable. If you specify TRUNCATE without the length n, the first 16 characters of the formatted values are used. When formatted values are longer

than 16 characters, you can use this option to revert to the levels as determined in releases before SAS 9. The TRUNCATE option is available only as a global option.

# **CLUSTER Statement**

# CLUSTER variables ;

The CLUSTER statement names variables that identify the first-stage clusters in a clustered sample design. First-stage clusters are also known as primary sampling units (PSUs). The combinations of categories of CLUSTER variables define the clusters in the sample. If a STRATA statement is specified, clusters are nested within strata.

If your sample design has clustering at multiple stages, you should specify only the first-stage clusters (PSUs) in the CLUSTER statement. For more information, see the section "Specifying the Sample Design" on page 9375.

If you provide replicate weights for replication variance estimation by specifying a REPWEIGHTS statement, you do not need to specify a CLUSTER statement.

The CLUSTER *variables* are one or more variables in the DATA= input data set. These variables can be either character or numeric, but the procedure treats them as categorical variables. The formatted values of the CLUSTER variables determine the CLUSTER variable levels. Thus, you can use formats to group values into levels. For more information, see the discussion of the FORMAT procedure in the *Base SAS Procedures Guide* and the discussions of the FORMAT statement and SAS formats in *SAS Formats and Informats: Reference*.

You can use multiple CLUSTER statements to specify CLUSTER variables. The procedure uses variables from all CLUSTER statements to create clusters. Cluster variables must not occur in the CLASS statement.

# **DOMAIN Statement**

**DOMAIN** variables < variable\* variable variable\* variable\* variable ... > ;

The DOMAIN statement requests analysis for domains (subpopulations), in addition to analysis for the entire study population. The DOMAIN statement names the variables that identify domains, which are called domain variables.

It is common practice to compute statistics for domains. The formation of these domains might not be known at the design stage. Therefore, the sample sizes for the domains are often random. Use a DOMAIN statement to incorporate this variability into the variance estimation.

Note that a DOMAIN statement is different from a BY statement. In a BY statement, you treat the sample sizes as fixed in each subpopulation, and you perform analysis within each BY group independently.

Use the DOMAIN statement on the entire data set to perform a domain analysis. Creating a new data set from a single domain and analyzing that with PROC SURVEYPHREG yields inappropriate estimates of variance.

A domain variable can be either character or numeric. The procedure treats domain variables as categorical variables. If a variable appears by itself in a DOMAIN statement, each level of this variable determines a domain in the study population. If two or more variables are joined by asterisks (\*), then every possible

combination of levels of these variables determines a domain. The procedure performs a descriptive analysis within each domain that is defined by the domain variables. Domain variables must not occur in the CLASS statement.

The formatted values of the domain variables determine the categorical variable levels. Thus, you can use formats to group values into levels. For more information, see the FORMAT procedure in the *Base SAS Procedures Guide* and the FORMAT statement and SAS formats in the *SAS Formats and Informats: Reference.* 

# **ESTIMATE Statement**

The ESTIMATE statement provides a mechanism for obtaining custom hypothesis tests. Estimates are formed as linear estimable functions of the form  $L\beta$ . You can perform hypothesis tests for the estimable functions, construct confidence limits, and obtain specific nonlinear transformations.

Table 115.3 summarizes the options available in the ESTIMATE statement.

Option	Description			
Construction and Com	Construction and Computation of Estimable Functions			
DIVISOR=	Specifies a list of values to divide the coefficients			
NOFILL	Suppresses the automatic fill-in of coefficients for higher-order effects			
SINGULAR=	Tunes the estimability checking difference			
Degrees of Freedom an	d <i>p</i> -values			
ADJUST=	Determines the method for multiple comparison adjustment of estimates			
ALPHA= $\alpha$	Determines the confidence level $(1 - \alpha)$			
LOWER	Performs one-sided, lower-tailed inference			
STEPDOWN	Adjusts multiplicity-corrected <i>p</i> -values further in a step-down fashion			
TESTVALUE=	Specifies values under the null hypothesis for tests			
UPPER	Performs one-sided, upper-tailed inference			
Statistical Output				
CL	Constructs confidence limits			
CORR	Displays the correlation matrix of estimates			
COV	Displays the covariance matrix of estimates			
E	Prints the L matrix			
JOINT	Produces a joint $F$ or chi-square test for the estimable functions			
PLOTS=	Requests ODS statistical graphics if the analysis is sampling-based			

 Table 115.3
 ESTIMATE Statement Options

Table 115.3	continued
Option	Description
SEED=	Specifies the seed for computations that depend on random numbers
Generalized Linear Mod	eling
CATEGORY=	Specifies how to construct estimable functions with multinomial
	data
EXP	Exponentiates and displays estimates
ILINK	Computes and displays estimates and standard errors on the inverse
	linked scale

For details about the syntax of the ESTIMATE statement, see the section "ESTIMATE Statement" on page 442 in Chapter 19, "Shared Concepts and Topics."

# **FREQ Statement**

# FREQ variable;

The FREQ statement names a numeric *variable* that provides a frequency for each observation in the input data set. PROC SURVEYPHREG treats each observation as if it appears *n* times, where *n* is the value of the FREQ variable for the observation. If not an integer, the frequency value is truncated to an integer. If the frequency value is missing, the observation is not used in the analysis. The FREQ statement allows one frequency variable.

If you use the FREQ statement and request the jackknife or BRR variance estimator by specifying the VARMETHOD=JACKKNIFE or VARMETHOD=BRR option in the PROC SURVEYPHREG statement, then you must identify the primary sampling units with a CLUSTER statement unless you also provide replicate weights with a REPWEIGHTS statement.

# LSMEANS Statement

LSMEANS < model-effects > < / options > ;

The LSMEANS statement computes and compares least squares means (LS-means) of fixed effects. LS-means are *predicted population margins*—that is, they estimate the marginal means over a balanced population. In a sense, LS-means are to unbalanced designs as class and subclass arithmetic means are to balanced designs.

Table 115.4 summarizes the options available in the LSMEANS statement.

Option	Description
Construction and Comp	utation of LS-Means
AT	Modifies the covariate value in computing LS-means
BYLEVEL	Computes separate margins
DIFF	Requests differences of LS-means
OM=	Specifies the weighting scheme for LS-means computation as
	determined by the input data set
SINGULAR=	Tunes estimability checking
Degrees of Freedom and	<i>p</i> -values
ADJUST=	Determines the method for multiple-comparison adjustment of
	LS-means differences
ALPHA= $\alpha$	Determines the confidence level $(1 - \alpha)$
STEPDOWN	Adjusts multiple-comparison p-values further in a step-down
	fashion
Statistical Output	
CL	Constructs confidence limits for means and mean differences
CORR	Displays the correlation matrix of LS-means
COV	Displays the covariance matrix of LS-means
E	Prints the L matrix
LINES	Produces a "Lines" display for pairwise LS-means differences
MEANS	Prints the LS-means
PLOTS=	Requests graphs of means and mean comparisons
SEED=	Specifies the seed for computations that depend on random
	numbers
Generalized Linear Mod	leling
EXP	Exponentiates and displays estimates of LS-means or LS-means
	differences
ILINK	Computes and displays estimates and standard errors of LS-means
	(but not differences) on the inverse linked scale
ODDSRATIO	Reports (simple) differences of least squares means in terms of
	odds ratios if permitted by the link function

 Table 115.4
 LSMEANS Statement Options

For details about the syntax of the LSMEANS statement, see the section "LSMEANS Statement" on page 458 in Chapter 19, "Shared Concepts and Topics."

# LSMESTIMATE Statement

LSMESTIMATE model-effect <'label' > values < divisor=n> <, ... <'label' > values < divisor=n>> </ options>; The LSMESTIMATE statement provides a mechanism for obtaining custom hypothesis tests among least squares means.

Table 115.5 summarizes the *options* available in the LSMESTIMATE statement.

Option	Description
Construction and Comp	outation of LS-Means
AT	Modifies covariate values in computing LS-means
BYLEVEL	Computes separate margins
DIVISOR=	Specifies a list of values to divide the coefficients
OM=	Specifies the weighting scheme for LS-means computation as
	determined by a data set
SINGULAR=	Tunes estimability checking
Degrees of Freedom and	<i>p</i> -values
ADJUST=	Determines the method for multiple-comparison adjustment of
	LS-means differences
ALPHA=α	Determines the confidence level $(1 - \alpha)$
LOWER	Performs one-sided, lower-tailed inference
STEPDOWN	Adjusts multiple-comparison p-values further in a step-down
	fashion
TESTVALUE=	Specifies values under the null hypothesis for tests
UPPER	Performs one-sided, upper-tailed inference
Statistical Output	
CL	Constructs confidence limits for means and mean differences
CORR	Displays the correlation matrix of LS-means
COV	Displays the covariance matrix of LS-means
E	Prints the L matrix
ELSM	Prints the K matrix
JOINT	Produces a joint F or chi-square test for the LS-means and
	LS-means differences
PLOTS=	Requests graphs of means and mean comparisons
SEED=	Specifies the seed for computations that depend on random
	numbers
Generalized Linear Mod	leling
CATEGORY=	Specifies how to construct estimable functions with multinomial
	data
EXP	Exponentiates and displays LS-means estimates
ILINK	Computes and displays estimates and standard errors of LS-means
	(but not differences) on the inverse linked scale

 Table 115.5
 LSMESTIMATE Statement Options

For details about the syntax of the LSMESTIMATE statement, see the section "LSMESTIMATE Statement" on page 477 in Chapter 19, "Shared Concepts and Topics."

# **MODEL Statement**

#### **MODEL** response < \* censor (list) > = effects < / options > ;

The MODEL statement identifies the variable to be used as the failure time variable, the optional censoring variable, and the explanatory effects, including covariates, main effects, and interactions; see the section "Specification of Effects" on page 3670 in Chapter 47, "The GLM Procedure," for more information. A note of caution: specifying the effect T\*A in the MODEL statement, where T is the time variable and A is a CLASS variable, does not make the effect time-dependent. You must specify exactly one MODEL statement.

The MODEL statement allows one response variable. In the MODEL statement, the failure time variable precedes the equal sign. This can optionally be followed by an asterisk, the name of the censoring variable, and a list of censoring values (separated by blanks or commas if there is more than one) enclosed in parentheses. If the censoring variable takes on one of these values, the corresponding failure time is considered to be censored. The variables following the equal sign are the explanatory variables (sometimes called independent variables or covariates) for the model.

The censoring variable must be numeric. The failure time variable must contain nonnegative values. Any observation with a negative failure time is excluded from the analysis, as is any observation with a missing value for any of the variables listed in the MODEL statement. See "Missing Values" on page 9377 for details.

Table 115.6 summarizes the *options* available in the MODEL statement, which can be specified after a slash (/).

Option	Description
ALPHA=	Specifies $\alpha$ for the $100(1 - \alpha)\%$ confidence limits
CLPARM	Computes confidence limits for regression parameters
COVB	Displays covariance matrix
DF=	Specifies the denominator degrees of freedom
HESS	Displays the Hessian matrix
INVHESS	Displays the inverse of the Hessian matrix
RISKLIMITS	Computes confidence limits for the exponentials of the
	regression parameters
SERATIO=	Computes the ratio of two standard errors for the
	regression coefficients
SINGULAR=	Specifies tolerance for testing singularity
TIES=	Specifies the method of handling ties in failure times
VADJUST=	Specifies a variance adjustment factor
VARRATIO=	Computes the ratio of two variances for the regression
	coefficients

## Table 115.6 MODEL Statement Options

#### ALPHA= $\alpha$

sets the level of the confidence limits for the estimated regression parameters and the hazard ratios. The value of *alpha* must be between 0 and 1, and the default is 0.05. A confidence level of  $\alpha$  produces  $100(1 - \alpha)\%$  confidence limits. The default of ALPHA=0.05 produces 95% confidence limits.

The ALPHA= option has no effect unless you also specify the CLPARM or RISKLIMITS option.

#### CLPARM

produces confidence limits for regression parameters of Cox proportional hazards models. You can specify the confidence coefficient by using the ALPHA= option. Classification main effects that use parameterizations other than REF, EFFECT, or GLM are ignored. For more information, see the section "Confidence Intervals" on page 9391.

# COVB

displays the estimated covariance matrix of the parameter estimates.

#### DF=value | keyword < (value) >

specifies the denominator degrees of freedom for hypothesis tests, specifies the degrees of freedom for confidence limits, and requests adjustments to the Wald test statistics. If you specify a *value*, it must be a nonnegative number.

In the description that follows, *d* denotes the usual degrees of freedom computed from the survey data by using the number of strata, clusters, or replicate weights. For more information, see the section "Degrees of Freedom" on page 9386.

By default, DF=PARMADJ when you use the Taylor series linearized variance estimator, and DF=DESIGN when you use the replication variance estimator. Alternatively, you can specify a nonnegative *value* for the degrees of freedom, or you can specify one of the following *keywords*:

#### **ALLREPS**

computes the denominator degrees of freedom for replication methods by using the total number of replicate samples. By default, PROC SURVEYPHREG computes the denominator degrees of freedom based on the number of replicate samples that are used. Some replicate samples might not be usable, in the sense that they cannot be used for variance estimation because of factors such as inestimability or nonconvergence. These replicate samples are not accounted for in the denominator degrees of freedom unless you specify DF=ALLREPS. For more information, see the section "Degrees of Freedom" on page 9386.

# DESIGN

computes the denominator degrees of freedom as *d*. When you specify DF=DESIGN, the corresponding Wald *F* statistics do not account for the number of parameters in the model. This option is useful if you do not want to apply the adjustment described in Korn and Graubard (1999, p. 93). For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# **DESIGN** (value)

computes the denominator degrees of freedom as *value*. When you specify DF=DESIGN (*value*), the corresponding Wald F statistics do not account for the number of parameters in the model. This option is useful if you do not want to apply the adjustment described in Korn and Graubard (1999, p. 93) and you want to specify the denominator degrees of freedom. You might want to specify a denominator degrees of freedom other than d for reasons such as missing values or domain estimation for relatively small domains. For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# DESIGNADJ

computes the denominator degrees of freedom as d. When you specify DF=DESIGNADJ, the corresponding Wald F statistics account for the number of parameters in the model. This option is useful if you are fitting a model that has many parameters relative to d but you want to use d as the denominator degrees of freedom. For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# NONE

specifies the denominator degrees of freedom to be infinite. This option is useful if you want to compute chi-square tests and normal confidence intervals. For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# PARMADJ

computes the denominator degrees of freedom as d minus the number of nonsingular parameters plus 1. When you specify DF=PARMADJ, the corresponding Wald F statistics account for the number of parameters in the model. This option is useful if you are fitting a model that has many parameters relative to d. For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# PARMADJ (value)

computes the denominator degrees of freedom as *value*. When you specify DF=PARMADJ (*value*), the corresponding Wald *F* statistics account for the number of parameters in the model. This option is useful if you are fitting a model with that has parameters relative to *d* and you want to specify the denominator degrees of freedom. You might want to specify the denominator degrees of freedom for reasons such as missing values or domain estimation for relatively small domains. For more information, see the section "Testing the Global Null Hypothesis" on page 9389.

# HESS

displays the last evaluation of the Hessian matrix.

# INVHESS

displays the inverse of the Hessian matrix that is evaluated at the estimated regression parameters.

# RISKLIMITS

# RL

produces confidence limits for hazard ratios and related quantities. For more information, see the section "Hazard Ratios" on page 9391. You can specify the confidence coefficient by using the ALPHA= option. You must take great care with any interpretation of the estimates and their confidence limits if interaction effects are involved in the model or if parameterizations other than REF, EFFECT, or GLM are used.

# SERATIO=ALL | MODEL | IND

computes the ratio of two standard errors for the regression parameters. The standard error in the numerator uses the complete design information that you specify. You can specify the following options to compute different standard errors for the denominator:

# ALL

requests both MODEL and IND standard error ratios.

# MODEL

computes the standard errors in the denominator as the square root of the diagonals of the inverse Hessian matrix evaluated at the estimated regression parameters. For more information, see the section "Variance Ratios and Standard Error Ratios" on page 9387.

## IND

computes the standard errors in the denominator by ignoring stratification and clustering. For more information, see the section "Variance Ratios and Standard Error Ratios" on page 9387.

#### SINGULAR=value

specifies the singularity criterion for determining linear dependencies in the set of explanatory variables. The default value is  $10^{-12}$ .

# **TIES**=method

specifies how to handle ties in the failure time. You can specify the following methods:

#### BRESLOW

uses the approximate partial likelihood of Breslow (1974).

# EFRON

uses the approximate partial likelihood of Efron (1977).

If there are no ties, both methods result in the same likelihood and yield identical estimates. By default, TIES=BRESLOW, which is the most efficient method when there are no ties.

### VADJUST=DF | PARMADJ | NONE | AVGREPSS

specifies variance adjustment factors. You can specify the following keywords:

# DF

# PARMADJ

requests the degrees-of-freedom adjustment (n-1)/(n-p) in the computation of the matrix **G** for the Taylor series linearization variance estimation.

# NONE

excludes the degrees-of-freedom adjustment (n-1)/(n-p) from the computation of the matrix G for the Taylor series linearization variance estimation. By default, VADJUST=NONE.

#### **AVGREPSS**

use the average sum of squares from all the usable replicate samples for the unusable replicates. This option is applicable only for the jackknife replication method. VADJUST=AVGREPSS multiplies the default jackknife variance estimator by the factor  $R/R_a$ , where  $R_a$  is the number of usable replicates and R is the total number of replicates. For more information, see the section "Variance Adjustment Factors" on page 9386.

#### VARRATIO=ALL | MODEL | IND

computes the ratio of two variances for the regression parameters. The variance in the numerator uses the complete design information. You can specify the following options to compute different variances for the denominator:

# ALL

requests both MODEL and IND variance ratios.

# MODEL

computes the variances in the denominator as the diagonals of the inverse Hessian matrix evaluated at the estimated regression parameters. For more information, see the section "Variance Ratios and Standard Error Ratios" on page 9387.

#### IND

computes the variances in the denominator by ignoring stratification and clustering. For more information, see the section "Variance Ratios and Standard Error Ratios" on page 9387.

# **NLOPTIONS Statement**

# **NLOPTIONS** < options > ;

The NLOPTIONS statement specifies details of the nonlinear optimization used by PROC SURVEYPHREG to maximize the log likelihood function. By default, the procedure uses the Newton-Raphson optimization technique. For more information about the NLOPTIONS statement, see the section "NLOPTIONS Statement" on page 489 in Chapter 19, "Shared Concepts and Topics."

# **OUTPUT Statement**

**OUTPUT** < **OUT**=SAS-data-set> < keyword=name ... keyword=name > < / options> ;

The OUTPUT statement creates a new SAS data set that contains statistics that are calculated for each observation unit. These statistics can include the estimated linear predictor  $(\mathbf{z}'_{j}\hat{\boldsymbol{\beta}})$  and its standard error, residuals, and influence statistics. In addition, this data set includes all the variables from the DATA= input data set.

Only score residuals are available in the OUTPUT data set if the model contains a time-dependent variable that is defined by means of programming statements.

The following list explains specifications in the OUTPUT statement:

#### OUT=SAS-data-set

names the output data set. If you omit the OUT= option, the OUTPUT data set is named by using the DATA*n* convention. See the section "OUT= Data Set for the OUTPUT statement" on page 9394 for more information.

# keyword=name

specifies the statistics to include in the OUTPUT data set and names the new variables that contain the statistics. Specify a *keyword* for each desired statistic (see the following list of *keywords*), and optionally an equal sign with either a variable or a list of variables in parentheses to contain the statistics. If you specify a *keyword* without a variable name, then the procedure uses default names. The *keywords* that accept a list of variables are RESSCH, RESSCO, and WTRESSCH. For these *keywords*, you can specify as many names in *name* as the number of explanatory variables in the MODEL statement. If you specify *k* names and *k* is less than the total number of explanatory variables, only the first *k* names are taken from the list; the procedure assigns default names for the rest of the statistics. The *keywords* and the corresponding statistics are as follows:

## ATRISK

specifies the number of subjects at risk at the observation time  $\tau_i$ .

#### RESDEV

specifies the deviance residual  $\hat{D}_j$ . This is a transform of the martingale residual to achieve a more symmetric distribution.

## RESMART

specifies the martingale residual  $\hat{M}_j$ . The residual at the observation time  $\tau_j$  can be interpreted as the difference over  $[0, \tau_j]$  in the observed number of events minus the expected number of events given by the model.

#### RESSCH

specifies the Schoenfeld residuals. These residuals are useful in assessing the proportional hazards assumption.

# RESSCO

specifies the score residuals. These residuals are a decomposition of the first partial derivative of the log likelihood. They can be used to assess the leverage that is exerted by each subject in the parameter estimation. They are also useful in constructing design-based variance estimators.

# STDXBETA

specifies the standard error of the estimated linear predictor,  $\sqrt{\mathbf{z}'_j \hat{\mathbf{V}}(\hat{\boldsymbol{\beta}}|F)}\mathbf{z}_j$ .

#### WTATRISK

specifies the weighted number of subjects at risk at the observation time  $\tau_i$ .

# XBETA

specifies the estimate of the linear predictor,  $\mathbf{z}'_{i} \hat{\boldsymbol{\beta}}$ .

# **Programming Statements**

Programming statements are used to create or modify the values of the explanatory variables in the MODEL statement. They are especially useful in fitting models with time-dependent explanatory variables. Programming statements can also be used to create explanatory variables that are not time-dependent. PROC SURVEYPHREG programming statements cannot be used to create or modify the values of the response variable, the censoring variable, the frequency variable, the WEIGHT variable, the CLASS variables, the STRATA variables, the CLUSTER variables, or the DOMAIN variables.

The following DATA step statements are available in PROC SURVEYPHREG:

```
ABORT;
ARRAY arrayname < [ dimensions ] > < $ > < variables-and-constants >;
CALL name < (expression < , expression ... >) >;
DELETE:
DO < variable = expression < TO expression > < BY expression > >
   <, expression < TO expression > < BY expression >> ...
   < WHILE expression > < UNTIL expression >;
END;
GOTO statement-label;
IF expression;
IF expression THEN program-statement;
             ELSE program-statement;
variable = expression;
variable + expression;
LINK statement-label;
PUT < variable > < = > ...;
RETURN;
SELECT < (expression) >;
STOP:
SUBSTR(variable, index, length)= expression;
WHEN (expression)program-statement;
       OTHERWISE program-statement;
```

By default, the PUT statement in PROC SURVEYPHREG writes results to the Output window instead of the Log window. If you want the results of the PUT statements to go to the Log window, add the following statement before the PUT statement:

# FILE LOG;

DATA step functions are also available. Use these programming statements the same way you use them in the DATA step. For detailed information, see the SAS Functions and CALL Routines: Reference.

Consider the following example of using programming statements in PROC SURVEYPHREG. Suppose blood pressure is measured at multiple times during the course of a study that investigates the effect of blood pressure on some survival time. By treating the blood pressure as a time-dependent explanatory variable, you can use the value of the most recent blood pressure at each specific point of time in the modeling process rather than using the initial blood pressure or the final blood pressure. The values of the following variables are recorded for each patient, if they are available. Otherwise, the variables contain missing values.

Time	survival time
Censor	censoring indicator (with 0 as the censoring value)
BP0	blood pressure on entry to the study
T1	time 1
BP1	blood pressure at T1
T2	time 2
BP2	blood pressure at T2
WT	design weight
PSU	identification of primary sampling units

The following programming statements create a variable BP. At each time T, the value of BP is the blood pressure reading for that time, if available. Otherwise, it is the last blood pressure reading.

```
proc surveyphreg;
  weight WT;
  model Time*Censor(0)=BP;
  cluster PSU;
  BP = BP0;
  if Time>=T1 and T1^=. then BP=BP1;
  if Time>=T2 and T2^=. then BP=BP2;
run;
```

# **REPWEIGHTS Statement**

# **REPWEIGHTS** variables < / options > ;

The REPWEIGHTS statement names *variables* that provide replicate weights for replication variance estimation, which you request with the VARMETHOD=BOOTSTRAP, VARMETHOD=BRR, or VARMETHOD=JACKKNIFE option in the PROC SURVEYPHREG statement. If you do not provide a REPWEIGHTS statement for VARMETHOD=BRR or VARMETHOD=JACKKNIFE, then PROC SUR-VEYPHREG constructs replicate weights for the analysis. For more information, see the sections "Balanced Repeated Replication (BRR) Method" on page 9381 and "Jackknife Method" on page 9383. For VARMETHOD=BOOTSTRAP, you must specify the REPWEIGHTS statement to provide replicate weights. For more information, see the section "Bootstrap Method" on page 9381.

Each REPWEIGHTS variable should contain the weights for a single replicate, and the number of replicates equals the number of REPWEIGHTS variables. The REPWEIGHTS variables must be numeric, and the variable values must be nonnegative numbers.

If you provide replicate weights with a REPWEIGHTS statement, you do not need to specify a CLUSTER or STRATA statement. If you use a REPWEIGHTS statement and do not specify the VARMETHOD= option in the PROC SURVEYPHREG statement, the procedure uses VARMETHOD=JACKKNIFE by default.

If you specify a REPWEIGHTS statement but do not include a WEIGHT statement, PROC SURVEYPHREG uses the average of each observation's replicate weights as the observation's weight.

You can specify the following options in the REPWEIGHTS statement after a slash (/):

#### JKCOEFS=jackknife-coefficient-specification

specifies jackknife coefficients for VARMETHOD=JACKKNIFE. The default value for the jackknife coefficient is (R - 1)/R, where *R* is the total number of replicates. You can specify an alternative value with one of the following three forms:

# JKCOEFS=value

specifies a single jackknife coefficient for all replicates. The coefficient *value* must be a nonnegative number.

# JKCOEFS=(values)

specifies jackknife coefficients for VARMETHOD=JACKKNIFE, where each coefficient corresponds to an individual replicate identified by a REPWEIGHTS variable. You can separate *values* with blanks or commas. The coefficient *values* must be nonnegative numbers. The number of *values* must equal the number of replicate weight variables named in the REPWEIGHTS statement. List these values in the same order in which you list the corresponding replicate weight variables in the REPWEIGHTS statement.

# JKCOEFS=SAS-data-set

names a SAS data set that contains the jackknife coefficients for VARMETHOD=JACKKNIFE. You provide the jackknife coefficients in the JKCOEFS= data set variable JKCoefficient. Each coefficient value must be a nonnegative number. The observations in the JKCOEFS= data set should correspond to the replicates that are identified by the REPWEIGHTS variables. Arrange the coefficients or observations in the JKCOEFS= data set in the same order in which you list the corresponding replicate weight variables in the REPWEIGHTS statement. The number of observations in the JKCOEFS= data set must not be less than the number of REPWEIGHTS variables.

See the section "Jackknife Method" on page 9383 for details about jackknife coefficients.

# REPCOEFS=replication-coefficient-specification

specifies replicate coefficients for replication methods. When you specify VARMETHOD=JACKKNIFE, the default value for the replicate coefficient is (R - 1)/R, where *R* is the total number of replicates. When you specify VARMETHOD=BOOTSTRAP or VARMETHOD=BRR, the default value for the replicate coefficient is 1/R.

For VARMETHOD=BOOTSTRAP or VARMETHOD=JACKKNIFE, you can specify one of the following three *replication-coefficient-specifications*:

# **REPCOEFS=**value

specifies a single replicate coefficient for all replicates, where *value* must be a nonnegative number.

# **REPCOEFS=(**values)

specifies a list of replicate coefficients, where each coefficient corresponds to an individual replicate that is identified by a REPWEIGHTS variable. You can separate *values* with blanks or commas. The coefficient *values* must be nonnegative numbers. The number of *values* must equal the number of replicate weight variables specified in the REPWEIGHTS statement. List these values in the same order in which you list the corresponding replicate weight variables in the REPWEIGHTS statement.

# REPCOEFS=SAS-data-set

names a *SAS-data-set* that contains the replicate coefficients. You must provide the replicate coefficients in a variable named Coefficient or RepCoefficient in the *SAS-data-set*. Each coefficient value must be a nonnegative number. The observations in the *SAS-data-set* should correspond to the replicates that are identified by the variables that are specified in the REPWEIGHTS statement. Arrange the coefficients or observations in the *SAS-data-set* in the same order in which you list the corresponding replicate weight variables in the REPWEIGHTS statement. The number of observations in the *SAS-data-set* must not be less than the number of variables specified in the REPWEIGHTS statement.

For more information about replication coefficients, see the section "Replicate Weights Method" on page 9385.

# **SLICE Statement**

SLICE model-effect < / options > ;

The SLICE statement provides a general mechanism for performing a partitioned analysis of the LS-means for an interaction. This analysis is also known as an analysis of simple effects.

The SLICE statement uses the same *options* as the LSMEANS statement, which are summarized in Table 19.21. For details about the syntax of the SLICE statement, see the section "SLICE Statement" on page 506 in Chapter 19, "Shared Concepts and Topics."

# **STORE Statement**

STORE < OUT= >item-store-name < / LABEL='label' > ;

The STORE statement requests that the procedure save the context and results of the statistical analysis. The resulting item store has a binary file format that cannot be modified. The contents of the item store can be processed with the PLM procedure. For details about the syntax of the STORE statement, see the section "STORE Statement" on page 509 in Chapter 19, "Shared Concepts and Topics."

# **STRATA Statement**

# **STRATA** variables < / option>;

The STRATA statement names variables that form the strata in a stratified sample design. The combinations of levels of STRATA variables define the strata in the sample, where strata are nonoverlapping subgroups that were sampled independently.

If your sample design has stratification at multiple stages, you should identify only the first-stage strata in the STRATA statement. See the section "Specifying the Sample Design" on page 9375 for more information.

If you provide replicate weights for replication variance estimation in a REPWEIGHTS statement, you do not need to specify a STRATA statement.

The STRATA *variables* are one or more variables in the DATA= input data set. These variables can be either character or numeric, but the procedure treats them as categorical variables. The formatted values of the STRATA variables determine the STRATA variable levels. Thus, you can use formats to group values into levels. See the discussion of the FORMAT procedure in the *Base SAS Procedures Guide* and the discussions of the FORMAT statement and SAS formats in the *SAS Formats and Informats: Reference*. Strata variables must not occur in the CLASS statement.

The STRATA statement in PROC SURVEYPHREG is different from the STRATA statement in PROC PHREG (Chapter 86, "The PHREG Procedure"). PROC PHREG fits different baseline hazard functions in different strata, which is useful if the proportional hazards assumption is not satisfied.

You can specify the following option in the STRATA statement after a slash (/):

#### LIST

displays a "Stratum Information" table, which lists all strata together with the corresponding values of the STRATA variables. This table provides the number of observations and the number of clusters in each stratum, as well as the sampling fraction if you specify the RATE= or TOTAL= option.

# **TEST Statement**

**TEST** < model-effects > < / options > ;

The TEST statement enables you to perform F tests for model effects that test Type I, Type II, or Type III hypotheses. See Chapter 15, "The Four Types of Estimable Functions," for details about the construction of Type I, II, and III estimable functions.

Table 115.7 summarizes the options available in the TEST statement.

Option	Description
CHISQ	Requests chi-square tests
DDF=	Specifies denominator degrees of freedom for fixed effects
E	Requests Type I, Type II, and Type III coefficients
E1	Requests Type I coefficients
E2	Requests Type II coefficients
E3	Requests Type III coefficients
HTYPE=	Indicates the type of hypothesis test to perform
INTERCEPT	Adds a row that corresponds to the overall intercept

Table 115.7 TEST Statement Options

For details about the syntax of the TEST statement, see the section "TEST Statement" on page 510 in Chapter 19, "Shared Concepts and Topics."

# **WEIGHT Statement**

#### WEIGHT variable ;

The WEIGHT statement names the variable that contains the sampling weights. This variable must be numeric, and the sampling weights must be positive numbers. If an observation has a weight that is nonpositive or missing, then the procedure omits that observation from the analysis. See the section "Missing Values" on page 9377 for more information. The WEIGHT statement allows one weight variable.

If you do not specify a WEIGHT statement but provide replicate weights with a REPWEIGHTS statement, PROC SURVEYPHREG uses the average of each observation's replicate weights as the observation's weight.

If you specify neither a WEIGHT statement nor a REPWEIGHTS statement, PROC SURVEYPHREG assigns all observations a weight of one.

# **Details: SURVEYPHREG Procedure**

# Notation and Estimation

Let  $U = \{1, 2, ..., N\}$  be the set of indices and let  $\mathcal{F}_N$  be the set of values for a finite population of size N. The survival time of each member of the finite population is assumed to follow its own hazard function,  $\lambda_i(t)$ , expressed as

 $\lambda_i(t) = \lambda(t; \mathbf{Z}_i(t)) = \lambda_0(t) \exp(\mathbf{Z}'_i(t)\boldsymbol{\beta})$ 

where  $\lambda_0(t)$  is an arbitrary and unspecified baseline hazard function,  $\mathbf{Z}_i(t)$  is the vector of explanatory variables for the *i*th unit at time *t*, and  $\boldsymbol{\beta}$  is the vector of unknown regression parameters that are associated with the explanatory variables. The vector  $\boldsymbol{\beta}$  is assumed to be the same for all individuals.

The partial likelihood function introduced by Cox (1972, 1975) eliminates the unknown baseline hazard  $\lambda_0(t)$  and accounts for censored survival times. If the entire population is observed, then this partial likelihood can be used to estimate  $\beta$ . Let  $\beta_N$  be the desired estimator. Assuming a working model with uncorrelated responses,  $\beta_N$  is obtained by maximizing the partial log likelihood,

$$l(\boldsymbol{\beta}) = \sum_{i \in U} \log \left\{ L(\boldsymbol{\beta}; \mathbf{Z}_i(t), t_i) \right\}$$

with respect to  $\boldsymbol{\beta}$ , where  $L(\boldsymbol{\beta}; \mathbf{Z}_i(t), t_i)$  is Cox's partial likelihood function.

Assume that probability sample A is selected from the finite population U and  $\pi_i$  is the selection probability for unit *i*. Further assume that covariates  $\mathbf{Z}_i(t)$  and survival time  $t_i$  are available for every unit in the sample A. An estimator of the finite population log likelihood is

$$l_{\pi}(\boldsymbol{\beta}) = \sum_{i \in A} \pi_i^{-1} \log \left\{ L(\boldsymbol{\beta}; \mathbf{Z}_i(t), t_i) \right\}$$

See "Partial Likelihood Function for the Cox Model" on page 9374 for more details.

A sample-based estimator  $\hat{\beta}$  for the finite population quantity  $\beta_N$  can be obtained by maximizing the partial pseudo-log-likelihood  $l_{\pi}(\beta; \mathbf{Z}_i(t), t_i)$  with respect to  $\beta$ . The design-based variance for  $\hat{\beta}$  is obtained by assuming the set of finite population values  $\mathcal{F}_N$  as fixed. For more information about maximum pseudo-likelihood estimators and other inferential approaches for survey data, see Kish and Frankel (1974); Godambe and Thompson (1986); Pfeffermann (1993), Korn and Graubard (1999, chapter 3), Chambers and Skinner (2003, chapter 2), and Fuller (2009, section 6.5). Maximum pseudo-likelihood estimators and their properties for Cox's proportional hazards model for survey data are discussed in Binder (1990, 1992); Lin and Wei (1989); Lin (2000); Boudreau and Lawless (2006).

Without loss of generality, the rest of this section uses indices for stratified clustered designs. For a stratified clustered sample design, observations are represented by a matrix

 $(\mathbf{w}, \mathbf{t}, \mathbf{\Delta}, \mathbf{Z}) = (w_{hij}, t_{hij}, \Delta_{hij}, \mathbf{z}_{hij})$ 

where

• w denotes the vector of sampling weights

- t denotes the event time variable
- $\Delta$  denotes the event indicator
- **Z** denotes the  $n \times p$  matrix of auxiliary information
- $h = 1, 2, \ldots, H$  is the stratum index
- $i = 1, 2, ..., n_h$  is the cluster index within stratum h
- $j = 1, 2, ..., m_{hi}$  is the unit index within cluster *i* of stratum *h*
- *p* is the total number of parameters
- $n = \sum_{h=1}^{H} \sum_{i=1}^{n_h} m_{hi}$  is the total number of observations in the sample
- $y_{hii}(t) = I(t_{hii} \ge t)$ , where  $I(\cdot)$  is an indicator function
- $n_{hii}(t) = I(t_{hii} \le t)$ , where  $I(\cdot)$  is an indicator function

Let  $\sum_{B} = \sum_{(h,i,j) \in B}$  denote the summation over the set of indices such that the observation unit *j* in PSU *i* and stratum *h* belongs to the index set *B*. Typically, *B* is the set of all population indices that are in the sample, the risk set, or the set of all units with a failure.

The first-stage sampling rate (fraction of PSUs selected for the sample) is denoted by  $f_h$ . The first-stage sampling rate is used in Taylor series variance estimation. You can specify the stratum sampling rates with the RATE= option. Or if you specify population totals with the TOTAL= option, PROC SURVEYPHREG computes  $f_h$  as the ratio of stratum sample size to the stratum total, in terms of PSUs. See the section "Population Totals and Sampling Rates" on page 9376 for details. If you do not specify the RATE= option or the TOTAL= option, then the procedure assumes that the stratum sampling rates  $f_h$  are negligible and does not use a finite population correction when computing variances.

# **Failure Time Distribution**

Let T be a nonnegative random variable that represents the failure time of an individual from a homogeneous superpopulation. The survival distribution function (also known as the survivor function) of T is written as

$$S(t) = \Pr(T \ge t)$$

A mathematically equivalent way of specifying the distribution of T is through its hazard function. The hazard function  $\lambda(t)$  specifies the instantaneous failure rate at t. If T is a continuous random variable,  $\lambda(t)$  is expressed as

$$\lambda(t) = \lim_{\Delta t \to 0^+} \frac{\Pr(t \le T < t + \Delta t \mid T \ge t)}{\Delta t} = \frac{f(t)}{S(t)}$$

where f(t) is the probability density function of T.
# **Time and CLASS Variables Usage**

The following DATA step creates an artificial data set, **Test**, to be used in this section. There are six variables in **Test**: the variable T contains the failure times; the variable Status is the censoring indicator variable with the value 1 for an uncensored failure time and the value 0 for a censored time; the variable A is a categorical variable with values 1, 2, and 3 representing three different categories; the variable MirrorT is an exact copy of T; the variable W is the observation weight; and the variable S is the strata indicator.

#### data Test;

in	put T	Stat	us A	W S	@@;				
Mi	rrorT	' = Т;							
da	talin	es;							
23	1	1	10	1	7	0	1	20	2
23	1	1	10	1	10	1	1	20	2
20	0	1	10	1	13	0	1	20	2
24	1	1	10	1	10	1	1	20	2
18	1	2	10	1	6	1	2	20	2
18	0	2	10	1	6	1	2	20	2
13	0	2	10	1	13	1	2	20	2
9	0	2	10	1	15	1	2	20	2
8	1	3	10	1	6	1	3	20	2
12	0	3	10	1	4	1	3	20	2
11	1	3	10	1	8	1	1	20	2
6	1	3	10	1	7	1	3	20	2
7	1	3	10	1	12	1	3	20	2
9	1	2	10	1	15	1	2	20	2
3	1	2	10	1	14	0	3	20	2
6	1	1	10	1	13	1	2	20	2

#### ;

# Time Variable on the Right Side of the MODEL Statement

The time variable cannot be used explicitly as an explanatory effect in the MODEL statement. The following statements produce an error message:

```
proc surveyphreg data=Test;
  weight W;
   strata S;
   class A;
   model T*Status(0)=T*A;
run;
```

To use the time variable as an explanatory effect, replace T by MirrorT as an effect, which is an exact copy of T, as in the following statements:

```
proc surveyphreg data=Test;
  weight W;
  strata S;
  class A;
  model T*Status(0)=A*MirrorT;
run;
```

Note that neither T\*A nor MirrorT\*A in the MODEL statement is time-dependent. The results of fitting this model are shown in Figure 115.3.

### Figure 115.3 T\*A Effect

### The SURVEYPHREG Procedure

Analysis of Maximum Likelihood Estimates							
		Standard				Hazard	
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio	
MirrorT*A 1	30	-17.560699	0.342527	-51.27	<.0001	0.000	
MirrorT*A 2	30	-17.424235	0.285106	-61.11	<.0001	0.000	
MirrorT*A 3	30	-17.448672	0.321376	-54.29	<.0001	0.000	

# **CLASS Variables and Programming Statements**

In PROC SURVEYPHREG, the levels of CLASS variables are determined by the CLASS statement and the input data and are not affected by user-supplied programming statements. Consider the following statements, which produce the results in Figure 115.4. Variable A is declared as a CLASS variable in the CLASS statement.

```
proc surveyphreg data=Test;
  weight W;
   strata S;
   class A;
   model T*Status(0)=A;
run;
```

Figure 115.4 shows the parameters that correspond to A and their respective regression coefficients estimates.

Figure 115.4 Design Variable and Regression Coefficient Estimates

		C	lass Level formation			
		Class	Levels Va	alues		
		Α	3 1	2 3		
Ar	alys	sis of Maxin	num Likeli	hood Es	stimates	
			Standard			Hazard
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio
A 1	30	-1.162184	0.644483	-1.80	0.0814	0.313
A 2	30	-0.616962	0.513355	-1.20	0.2388	0.540
A 3	30	0				1.000

```
The \ {\it SURVEYPHREG} \ Procedure
```

Now consider the programming statement that attempts to change the value of the CLASS variable A as in the following specification:

```
proc surveyphreg data=Test;
weight W;
strata S;
```

```
class A;
model T*Status(0)=A;
if A=3 then A=2;
run;
```

Results of this analysis are shown in Figure 115.5 and are identical to those in Figure 115.4. The **if A=3 then A=2** programming statement has no effect on the explanatory variable for A, which have already been determined.

Figure 115.5	Design Variable and Regression Coefficient Estimates
	The SURVEYPHREG Procedure

		Class	lass Level			
		Class	Levels va	alues		
		Α	3 1	23		
Ar	nalys	sis of Maxir	num Likeli	hood E	stimates	
			Standard			Hazard
Parameter	DF	Estimate	Error	t Value	• Pr >  t	Ratio
A 1	30	-1.162184	0.644483	-1.80	0.0814	0.313
A 2	30	-0.616962	0.513355	-1.20	0.2388	0.540
A 3	30	0				1.000

Additionally any variable used in a programming statement that has already been declared in the CLASS statement is *not* treated as a collection of the corresponding design variables. Consider the following statements:

```
proc surveyphreg data=Test;
    class A;
    model T*Status(0)=A X;
    X=T*A;
run;
```

The CLASS variable A generates two design variables as explanatory variables. The variable X created by the **X=T\*A** programming statement is a single time-dependent covariate whose values are evaluated using the exact values of A given in the data, not the dummy coded values that represent A. In the data set Test, A has the values of 1, 2, and 3, and these values are multiplied by the values of T to produce X. If A were a character variable with values 'Bird', 'Cat', and 'Dog', the programming statement X=T\*A would have produced an error in the attempt to multiply a number with a character value.

Figure 115.6 Single Time-Dependent Variable X\*A

Analysis of Maximum Likelihood Estimates						
			Hazard			
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio
A 1	31	0.158010	1.182556	0.13	0.8946	1.171
A 2	31	0.008993	0.652504	0.01	0.9891	1.009
A 3	31	0				1.000
х	31	0.092679	0.071328	1.30	0.2034	1.097

#### The SURVEYPHREG Procedure

The following statements are not the same as in the preceding program. If you want to create time-dependent covariates from the values of a CLASS variable, you could use syntax like the following:

```
proc surveyphreg data=Test;
   class A;
   model T*Status(0)=A X1 X2;
   X1= T*(A=1);
   X2= T*(A=2);
run;
```

The Boolean parenthetical expressions (A=1) and (A=2) resolve to a value of 1 or 0, depending on whether the expression is true or false, respectively.

Results of this test are shown in Figure 115.7.

Analysis of Maximum Likelihood Estimates							
		Standard					
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio	
A 1	31	-0.007655	1.221122	-0.01	0.9950	0.992	
A 2	31	-0.881383	1.743507	-0.51	0.6168	0.414	
A 3	31	0				1.000	
X1	31	-0.155220	0.164334	-0.94	0.3522	0.856	
X2	31	0.011554	0.188932	0.06	0.9516	1.012	

The SURVEYPHREG Procedure

Figure 115.7 Simple Test of Proportional Hazards Assumption

In general, when your model contains a categorical explanatory variable that is time-dependent, it might be necessary to use hardcoded dummy variables to represent the categories of the categorical variable.

# Partial Likelihood Function for the Cox Model

Let  $t_{(1)} < t_{(2)} < ... < t_{(K)}$  denote the K distinct, ordered event times. Let  $d_k$  denote the multiplicity of failures at  $t_{(k)}$ ; that is,  $d_k$  is the size of the set  $\mathcal{D}_k$  of individuals that fail at  $t_{(k)}$ . Let  $w_{hij}$  be the weight associated with the *j*th observation unit in the *i*th cluster in stratum h. Using this notation, the pseudo-likelihood functions used in PROC SURVEYPHREG to estimate  $\beta_N$  are described in the following sections.

#### **Continuous Time Scale**

Let  $\mathcal{R}_k$  denote the risk set just before the *k*th ordered event time  $t_{(k)}$ .

The Breslow likelihood is expressed as

$$L_{\text{Breslow}}(\boldsymbol{\beta}) = \prod_{k=1}^{K} \frac{\exp\left(\boldsymbol{\beta}' \sum_{\mathcal{D}_{k}} w_{hij} \mathbf{Z}_{hij}(t)\right)}{\left\{\sum_{\mathcal{R}_{k}} w_{hij} \exp(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t))\right\}^{\sum_{\mathcal{D}_{k}} w_{hij}}}$$

The Efron likelihood is expressed as

$$L_{\text{Efron}}(\boldsymbol{\beta}) = \prod_{k=1}^{K} \frac{\exp\left(\boldsymbol{\beta}' \sum_{\mathcal{D}_{k}} w_{hij} \mathbf{Z}_{hij}(t)\right)}{\left\{\phi(\boldsymbol{\beta}, \mathbf{Z}, \mathbf{w}, k)\right\}^{\frac{1}{d_{k}} \sum_{\mathcal{D}_{k}} w_{hij}}}$$

where  $\phi(\boldsymbol{\beta}, \mathbf{Z}, \mathbf{w}, k)$  is

$$\phi(\boldsymbol{\beta}, \mathbf{Z}, \mathbf{w}, k) = \prod_{l=1}^{d_k} \left\{ \sum_{\mathcal{R}_k} w_{hij} \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) - \frac{l-1}{d_k} \sum_{\mathcal{D}_k} w_{hij} \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) \right\}$$

# Specifying the Sample Design

PROC SURVEYPHREG produces statistics that are based on the sample design used to obtain the survey data. PROC SURVEYPHREG can be used for single-stage or multistage designs, with or without stratification, and with or without unequal weighting. To analyze your survey data with PROC SURVEYPHREG, you need to provide sample design information for the procedure. This information can include design (or variance) strata, clusters, and sampling weights. You provide sample design information with the STRATA, CLUSTER, and WEIGHT statements, and with the RATE= or TOTAL= option in the PROC SURVEYPHREG statement.

If you provide replicate weights for replication variance estimation, you do not need to specify a STRATA or CLUSTER statement. Otherwise, you should specify STRATA and CLUSTER statements whenever your design includes stratification and clustering.

When there are clusters (PSUs) in the sample design, the procedure estimates variance by using the PSUs, as described in the section "Variance Estimation" on page 9379. For a multistage sample design, PROC SURVEYPHREG uses only the first stage of the sample design for variance estimation. So, the required input includes only first-stage cluster (PSU) and first-stage stratum identification. You do not need to input design information about any additional stages of sampling.

#### Stratification

If your sample design is stratified at the first stage of sampling, use the STRATA statement to name the variables that form the strata. The combinations of categories of STRATA variables define the strata in the sample, where strata are nonoverlapping subgroups that were sampled independently. If your sample design has stratification at multiple stages, you should identify only the first-stage strata in the STRATA statement.

If you use a REPWEIGHTS statement to provide replicate weights for replication variance estimation, you do not need to specify a STRATA statement. Otherwise, you should specify a STRATA statement whenever your design includes stratification. If you do not specify a STRATA statement or a REPWEIGHTS statement, then PROC SURVEYPHREG assumes there is no stratification at the first stage. In other words, in this case, the procedure assumes that all observation units are in the same stratum.

#### Clustering

If your sample design selects clusters at the first stage of sampling, use the CLUSTER statement to name the variables that identify the first-stage clusters, which are also called primary sampling units (PSUs). The combinations of categories of CLUSTER variables define the clusters in the sample. If there is a STRATA statement, clusters are nested within strata. If your sample design has clustering at multiple stages, you should

specify only the first-stage clusters (PSUs) in the CLUSTER statement. PROC SURVEYPHREG assumes that each cluster that is defined by the CLUSTER statement variables represents a PSU in the sample.

If you use a REPWEIGHTS statement to provide replicate weights for replication variance estimation, you do not need to specify a CLUSTER statement. Otherwise, you should specify a CLUSTER statement whenever your design includes clustering at the first stage of sampling. If you do not specify a CLUSTER statement, then PROC SURVEYPHREG treats each observation as a PSU.

# Weighting

If your sample design includes unequal weighting, use the WEIGHT statement to name the variable that contains the sampling weights. Sampling weights must be positive numbers. If an observation has a weight that is nonpositive or missing, then the procedure omits that observation from the analysis. See the section "Missing Values" on page 9377 for more information.

If you do not specify a WEIGHT statement but include a REPWEIGHTS statement, PROC SURVEYPHREG uses the average of each observation's replicate weights as the observation's weight. If you specify neither a WEIGHT statement nor a REPWEIGHTS statement, PROC SURVEYPHREG assumes all observations have a weight of one.

# **Population Totals and Sampling Rates**

To include a finite population correction (fpc) in Taylor series variance estimation, you can specify either the sampling rate or the population total by using the RATE= or TOTAL= option, respectively, in the PROC SURVEYPHREG statement. You cannot specify both of these options in the same PROC SURVEYPHREG statement. The RATE= and TOTAL= options apply only to Taylor series variance estimation. The procedure does not use a finite population correction for bootstrap, BRR, or jackknife variance estimation.

If you do not specify the RATE= or TOTAL= option, the Taylor series variance estimation does not include a finite population correction. For fairly small sampling fractions, this correction is often ignored. See Cochran (1977) and Kish (1965) for more information.

If your design has multiple stages of selection and you are specifying the RATE= option, you should input the first-stage sampling rate, which is the ratio of the number of PSUs in the sample to the total number of PSUs in the study population. If you are specifying the TOTAL= option for a multistage design, you should input the total number of PSUs in the study population.

For a nonstratified sample design, or for a stratified sample design with the same sampling rate or the same population total in all strata, you can use the RATE=*value* or TOTAL=*value* option. If your sample design is stratified with different sampling rates or population totals in different strata, use the RATE=*SAS-data-set* or TOTAL=*SAS-data-set* option to name a SAS data set that contains the stratum sampling rates or totals. This data set is called a *secondary data set*, as opposed to the *primary data set* that you specify with the DATA= option.

The secondary data set must contain all the stratification variables listed in the STRATA statement and all the variables in the BY statement. Furthermore, the BY groups must appear in the same order as in the primary data set. If there are formats associated with the STRATA variables and the BY variables, then the formats must be consistent in the primary and the secondary data sets. If you specify the TOTAL=SAS-data-set option, the secondary data set must have a variable named \_TOTAL\_ that contains the stratum population totals. If you specify the RATE=SAS-data-set option, the secondary data set must have a variable named \_TOTAL\_ that contains the stratum population totals. If you specify the RATE=SAS-data-set option, the secondary data set must have a variable named \_RATE\_ that contains the stratum sampling rates. If the secondary data set contains more than one observation

for any one stratum, the procedure uses the first value of \_TOTAL\_ or \_RATE\_ for that stratum and ignores the rest.

The *value* in the RATE= option or the values of \_RATE\_ in the secondary data set must be nonnegative numbers. You can specify *value* as a number between 0 and 1. Or you can specify *value* in percentage form as a number between 1 and 100, and PROC SURVEYPHREG converts that number to a proportion. The procedure treats the value 1 as 100% instead of 1%.

If you specify the TOTAL=*value* option, *value* must not be less than the sample size. If you provide stratum population totals in a secondary data set, these values must not be less than the corresponding stratum sample sizes.

# **Missing Values**

Missing values in your survey data can compromise the quality of your survey results. Some missing values for survey data are because of nonresponses. An observation whose response to every survey item is available is called a *complete respondent*, and an observation whose response to one or more survey items are missing is called an *incomplete respondent*. If the complete respondents are different from the incomplete respondents with regard to a survey effect or outcome, then survey estimates will be biased and will not accurately represent the survey population. There are a variety of techniques in sample design and survey operations that can reduce nonresponse. After data collection is complete, you can use imputation to replace missing values with acceptable values, and you can use sampling weight adjustments to compensate for nonresponse. You should complete this data preparation and adjustment before you analyze your data with PROC SURVEYPHREG. For more details, see Cochran (1977); Kalton and Kasprzyk (1986); Brick and Kalton (1996).

### **WEIGHT Variable**

If an observation has a missing value or a nonpositive value for the WEIGHT variable, then PROC SUR-VEYPHREG excludes that observation from the analysis.

#### **REPWEIGHTS Variables**

If you provide replicate weights in a REPWEIGHTS statement for replication variance estimation, all REPWEIGHTS variable values must be nonmissing. Similarly, if you provide jackknife coefficients in the JKCOEFS= option in the REPWEIGHTS statement, all values of the JKCoefficient variable must be nonmissing. The procedure does not perform the analysis when any replicate weight or jackknife coefficient value is missing.

### **CLASS, STRATA, CLUSTER, and DOMAIN Variables**

An observation is excluded from the analysis if it has a missing value for any CLASS, STRATA, CLUSTER, or DOMAIN variable, unless you specify the MISSING option in the PROC SURVEYPHREG statement. If you specify the MISSING option, the procedure treats missing values as a valid (nonmissing) category for all categorical variables, which include STRATA variables, CLUSTER variables, CLASS variables, and DOMAIN variables.

# **Analysis Variables**

By default, PROC SURVEYPHREG excludes an observation from the likelihood estimation and all associated analyses if the observation has a missing value for any of the variables in the MODEL statement, unless you specify the MISSING or NOMCAR option in the PROC SURVEYPHREG statement. When the procedure excludes observations with missing values from analyses, it displays the total frequency of observations used in the NObs table.

If you specify time-dependent covariates by using programming statements, the procedure computes the values of the covariates for all observations in the risk set at every event time. If an observation contains missing values for any of the time-dependent covariates at a given event time, then the observation is not used at that event time. However, that same observation can be used at some other event times where it contains no missing values. Therefore, an observation with missing time-dependent covariates can be used at some event times but ignored at other event times, depending on whether any of the corresponding time-dependent covariates are missing.

If you specify the MISSING option, the procedure treats missing levels as a valid (nonmissing) level for each categorical analysis variable.

If you specify the NOMCAR option for Taylor series variance estimation, the procedure includes observations with missing values of analysis variables in the variance computations.

# **The NOMCAR Option**

When you specify the NOMCAR option, PROC SURVEYPHREG computes variance estimates by analyzing the nonmissing values for variables in the regression model as a domain or subpopulation, where the entire population includes both nonmissing and missing domains. By default, if an observation contains missing values for the dependent variable or for any variable used in the independent effects, the observation is excluded from the analysis. See the section "Missing Values" on page 9377 for more information.

Note that the NOMCAR option has no effect on categorical predictors when you specify the MISSING option, which treats missing values as a valid nonmissing level. The NOMCAR option does not affect the inclusion of observations that have missing values in the WEIGHT, FREQ, CLUSTER, STRATA, or DOMAIN variables. Observations that have missing values of the WEIGHT and FREQ variables are always excluded from the analysis. Observations that have missing values of the CLUSTER, DOMAIN, or STRATA variables are excluded unless you specify the MISSING option.

The NOMCAR option applies only to Taylor series variance estimation. The replication methods, which you request by specifying VARMETHOD=BOOTSTRAP, VARMETHOD=BRR, or VARMETHOD=JACKKNIFE, do not use the NOMCAR option.

# **Degrees of Freedom**

PROC SURVEYPHREG uses the degrees of freedom of the variance estimator to obtain t confidence limits and Wald-type F tests. The procedure computes the degrees of freedom based on the variance estimation method, the sample design, and the number of estimable parameters. For more information, see the section "Degrees of Freedom" on page 9386. This section describes how missing values can affect the computation of the degrees of freedom.

#### **Taylor Series Variance Estimation**

The degrees of freedom can depend on the number of clusters, the number of strata, and the number of observations. For Taylor series variance estimation, these numbers are based on the observations that are included in the analysis. These numbers do not count observations that are excluded from the analysis because they have missing values. If all values in a stratum are excluded from the analysis as missing values, then that stratum is called an *empty stratum*. Empty strata are not counted in the total number of strata for the analysis. Similarly, empty clusters and missing observations are not included in the totals counts of clusters and observations that are used to compute the degrees of freedom for the analysis.

If you specify the MISSING option, missing values are treated as valid nonmissing levels and are included in computing the degrees of freedom. If you specify the NOMCAR option for Taylor series variance estimation, observations that have missing values for variables in the regression model are included in computing the degrees of freedom.

#### **Replicate-Based Variance Estimation**

For BRR or jackknife variance estimation, by default PROC SURVEYPHREG computes the degrees of freedom by using all valid observations in the input data set. A valid observation is an observation that has a positive value of the WEIGHT variable and nonmissing values of the STRATA and CLUSTER variables unless you specify the MISSING option.

# Variance Estimation

PROC SURVEYPHREG uses the Taylor series method or replication (resampling) methods to estimate sampling errors of estimators that are based on complex sample designs (Fuller 1975; Särndal, Swensson, and Wretman 1992; Wolter 2007; Rust 1985; Dippo, Fay, and Morganstein 1984; Rao and Shao 1999, 1996; and Binder 1992). You can use the VARMETHOD= option in the PROC statement to specify the variance estimation method. By default, PROC SURVEYPHREG uses the Taylor series method.

However, replication methods have recently gained popularity for estimating variances in complex survey data analysis. One reason for this popularity is the relative simplicity of replication-based estimates, especially for nonlinear estimators; another is that modern computational capacity has made replication methods feasible for practical survey analysis.

Replication methods draw multiple replicates (also called subsamples) from a full sample according to a specific resampling scheme. The most commonly used resampling schemes are the *balanced repeated replication* (BRR) method and the *jackknife* method. For each replicate, the original weights are modified for the PSUs in the replicates to create replicate weights. The parameters of interest are estimated by using the replicate weights for each replicate. Then the variances of parameters of interest are estimated by the variability among the estimates derived from these replicates. The procedure automatically creates replicate weights based on the replication method you specify; alternatively you can use the REPWEIGHTS statement to provide your own replicate weights for variance estimation.

The following sections provide details about how the variance-covariance matrix of the estimated regression coefficients is estimated for each variance estimation method.

#### **Taylor Series Linearization**

The Taylor series linearization method is the default variance estimation method used by PROC SUR-VEYPHREG. See the section "Notation and Estimation" on page 9369 for definitions of the notation used in this section. Let

$$S^{(r)}(\boldsymbol{\beta},t) = \sum_{A} w_{hij} y_{hij}(t) \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) \mathbf{Z}_{hij}^{\otimes r}(t)$$

where r = 0, 1. Let A be the set of indices in the selected sample. Let

$$\mathbf{a}^{\bigotimes r} = \begin{cases} \mathbf{a}\mathbf{a}^{\prime} & , \ r = 1\\ I_{\dim(\mathbf{a})} & , \ r = 0 \end{cases}$$

and let  $I_{\dim(\mathbf{a})}$  be the identity matrix of appropriate dimension.

Let 
$$\overline{\mathbf{Z}}(\boldsymbol{\beta}, t) = \frac{S^{(1)}(\boldsymbol{\beta}, t)}{S^{(0)}(\boldsymbol{\beta}, t)}$$
. The score residual for the  $(h, i, j)$  subject is

$$\begin{split} \mathbf{L}_{hij}(\boldsymbol{\beta}) &= \Delta_{hij} \left\{ \mathbf{Z}_{hij}(t_{hij}) - \bar{\mathbf{Z}}(\boldsymbol{\beta}, t_{hij}) \right\} \\ &- \sum_{(h', i', j') \in A} \Delta_{h'i'j'} \frac{w_{h'i'j'} Y_{hij}(t_{h'i'j'}) \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t_{h'i'j'})\right)}{S^{(0)}(\boldsymbol{\beta}, t_{h'i'j'})} \left\{ \mathbf{Z}_{hij}(t_{h'i'j'}) - \bar{\mathbf{Z}}(\boldsymbol{\beta}, t_{h'i'j'}) \right\} \end{split}$$

For TIES=EFRON, the computation of the score residuals is modified to comply with the Efron partial likelihood. See the section "Residuals" on page 9391 for more information.

The Taylor series estimate of the covariance matrix of  $\hat{\beta}$  is

 $\hat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \mathcal{I}^{-1}(\hat{\boldsymbol{\beta}})\mathbf{G}\mathcal{I}^{-1}(\hat{\boldsymbol{\beta}})$ 

where  $\mathcal{I}(\hat{\beta})$  is the observed information matrix and the  $p \times p$  matrix G is defined as

$$\mathbf{G} = \frac{n-1}{n-p} \sum_{h=1}^{H} \frac{n_h (1-f_h)}{n_h - 1} \sum_{i=1}^{n_h} (\mathbf{e}_{hi+} - \bar{\mathbf{e}}_{h..})' (\mathbf{e}_{hi+} - \bar{\mathbf{e}}_{h..})$$

The observed residuals, their sums and means are defined as follows:

$$\mathbf{e}_{hij} = w_{hij} \mathbf{L}_{hij}(\hat{\boldsymbol{\beta}})$$
$$\mathbf{e}_{hi+} = \sum_{j=1}^{m_{hi}} \mathbf{e}_{hij}$$
$$\bar{\mathbf{e}}_{h\cdots} = \frac{1}{n_h} \sum_{i=1}^{n_h} \mathbf{e}_{hi+}$$

The factor (n-1)/(n-p) in the computation of the matrix **G** reduces the small sample bias that is associated with using the estimated function to calculate deviations (Fuller et al. (1989), pp. 77–81). For simple random sampling, this factor contributes to the degrees of freedom correction applied to the residual mean square for ordinary least squares in which *p* parameters are estimated. By default, the procedure uses this adjustment in the variance estimation. If you do not want to use this multiplier in the variance estimator, then specify the VADJUST=NONE option in the MODEL statement.

# **Bootstrap Method**

Bootstrap replicate samples are created by sampling the primary sampling units (PSUs) within each stratum with replacement. The original sampling weights are then adjusted in each replicate to reflect the full sample. These adjusted weights are also called bootstrap replicate weights. For more information, see Rao and Wu (1988) and Wolter (2007, Chapter 5).

Currently, PROC SURVEYPHREG does not generate bootstrap replicate weights. However, if the bootstrap replicate weights are available to you for a survey, then you can use the REPWEIGHTS statement to name the variables that contain the bootstrap replicate weights and specify the VARMETHOD=BOOTSTRAP option in the PROC SURVEYIMPUTE statement. The SURVEYPHREG procedure uses 1/R as the default bootstrap replicate coefficient when you specify the VARMETHOD=BOOTSTRAP option, where *R* is the total number of replicates. Alternatively, you can specify different replicate coefficients by using the REPCOEFS= option in the REPWEIGHTS statement.

For more information, see the section "Replicate Weights Method" on page 9385.

# **Balanced Repeated Replication (BRR) Method**

The balanced repeated replication (BRR) method requires that the full sample be drawn by using a stratified sample design with two primary sampling units (PSUs) per stratum. The BRR method constructs half-sample replicates by deleting one PSU per stratum according to a Hadamard matrix and doubling the original weight of the other PSU in that stratum. Let *H* be the total number of strata. The total number of replicates *R* is the smallest multiple of 4 that is greater than *H*. However, if you prefer a larger number of replicates, you can specify the REPS=*n* method-option. If a  $n \times n$  Hadamard matrix cannot be constructed, the number of replicates is increased until a Hadamard matrix becomes available.

Each replicate is obtained by deleting one PSU per stratum according to a corresponding Hadamard matrix and adjusting the original weights for the remaining PSUs. The new weights are called replicate weights.

Replicates are constructed by using the first *H* columns of the  $R \times R$  Hadamard matrix. The *r*th (r = 1, 2, ..., R) replicate is drawn from the full sample according to the *r*th row of the Hadamard matrix as follows:

- If the (r, h) element of the Hadamard matrix is 1, then the first PSU of stratum h is included in the rth replicate and the second PSU of stratum h is excluded.
- If the (r, h) element of the Hadamard matrix is -1, then the second PSU of stratum *h* is included in the *r*th replicate and the first PSU of stratum *h* is excluded.

The replicate weights of the remaining PSUs in each half sample are then doubled to their original weights. For more detail about the BRR method, see Wolter (2007) and Lohr (2010).

By default, an appropriate Hadamard matrix is generated automatically to create the replicates. You can display the Hadamard matrix by specifying the VARMETHOD=BRR(PRINTH) *method-option*. If you provide a Hadamard matrix by specifying the VARMETHOD=BRR(HADAMARD=) *method-option*, then the replicates are generated according to the provided Hadamard matrix. You can use the VARMETHOD=BRR(OUTWEIGHTS=) *method-option* to store the replicate weights in a SAS data set.

Let  $\hat{\beta}$  be the estimated proportional hazards regression coefficients from the full sample, and let  $\hat{\beta}_r$  be the estimated proportional hazards regression coefficients from the *r*th replicate by using replicate weights.

PROC SURVEYPHREG estimates the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\widehat{\boldsymbol{\beta}}) = \frac{1}{R} \sum_{r=1}^{R} \left( \widehat{\boldsymbol{\beta}}_{r} - \widehat{\boldsymbol{\beta}} \right) \left( \widehat{\boldsymbol{\beta}}_{r} - \widehat{\boldsymbol{\beta}} \right)'$$

with H degrees of freedom, where H is the number of strata.

If you specify the CENTER=REPLICATES *method-option*, then PROC SURVEYPHREG computes the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\widehat{\boldsymbol{\beta}}) = \frac{1}{R} \sum_{r=1}^{R} \left( \widehat{\boldsymbol{\beta}}_{r} - \overline{\widehat{\boldsymbol{\beta}}_{r}} \right) \left( \widehat{\boldsymbol{\beta}}_{r} - \overline{\widehat{\boldsymbol{\beta}}_{r}} \right)'$$

where  $\overline{\hat{\beta}_r}$  is the average of the replicate estimates as follows:

$$\overline{\hat{\beta}_r} = \frac{1}{R} \sum_{r=1}^R \hat{\beta_r}$$

If one or more components of  $\hat{\beta}_r$  cannot be calculated for some replicates, then the variance estimate is computed by using only the replicates for which the proportional hazards regression coefficients can be estimated. Estimability and nonconvergence are the two most common reasons why  $\hat{\beta}_r$  might not be available for a replicate sample even if  $\hat{\beta}$  is defined for the full sample. Let  $R_a$  be the number of replicates where  $\hat{\beta}_r$  is available, and let  $R - R_a$  be the number of replicates where  $\hat{\beta}_r$  is not available. Without loss of generality, assume that  $\hat{\beta}_r$  is available only for the first  $R_a$  replicates; then the BRR variance estimator is

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \frac{1}{R_a} \sum_{r=1}^{R_a} \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

with degrees of freedom equal to the minimum of H and  $R_a$ , where H is the number of strata. Alternatively, you can use the FAY= *method-option* to request Fay's BRR method, as discussed in the following section.

#### Fay's BRR Method

The traditional BRR method constructs half-sample replicates by deleting one PSU per stratum according to a Hadamard matrix and doubling the original weight of the other PSU. Fay's BRR method uses the Fay coefficient,  $\epsilon$  ( $0 \le \epsilon < 1$ ), and instead of deleting one PSU per stratum, it multiplies the original weight by the coefficient  $\epsilon$ . The original weight of the remaining PSU in that stratum is multiplied by  $2 - \epsilon$ . PROC SURVEYPHREG uses  $\epsilon = 0.5$  as the default value; alternatively, you can specify a value for  $\epsilon$  with the FAY= method-option. When  $\epsilon = 0$ , Fay's method becomes the traditional BRR method. For more details, see Dippo, Fay, and Morganstein (1984); Fay (1984, 1989); Judkins (1990). Because the traditional BRR method uses only half of the total sample in every replicate, several replicate estimators ( $\hat{\beta}_r$ ) might be undefined even when the full sample estimator ( $\hat{\beta}$ ) is defined. Fay's BRR method is especially useful for this situation because it uses all the sampled units in every replicate.

Let  $\hat{\beta}$  be the estimated proportional hazards regression coefficients from the full sample, and let  $\hat{\beta}_r$  be the estimated regression coefficients that are obtained from the *r*th replicate by using replicate weights. PROC

SURVEYPHREG estimates the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\widehat{\boldsymbol{\beta}}) = \frac{1}{R(1-\epsilon)^2} \sum_{r=1}^{R} \left(\widehat{\boldsymbol{\beta}}_r - \widehat{\boldsymbol{\beta}}\right) \left(\widehat{\boldsymbol{\beta}}_r - \widehat{\boldsymbol{\beta}}\right)'$$

with H degrees of freedom, where H is the number of strata.

### Hadamard Matrix

PROC SURVEYPHREG uses a Hadamard matrix to construct replicates for BRR variance estimation. You can provide a Hadamard matrix for replicate construction by using the HADAMARD= *method-option* for VARMETHOD=BRR. Otherwise, PROC SURVEYPHREG generates an appropriate Hadamard matrix. You can display the Hadamard matrix by specifying the PRINTH *method-option*.

A Hadamard matrix **A** of dimension *R* is a square matrix that has all elements equal to 1 or -1 such that  $\mathbf{A'A} = R\mathbf{I}$ , where  $\mathbf{I}$  is an identity matrix of appropriate order. The dimension of a Hadamard matrix must equal 1, 2, or a multiple of 4.

For example, the following matrix is a Hadamard matrix of dimension k = 8:

1	1	1	1	1	1	1	1
1	-1	1	-1	1	-1	1	-1
1	1	-1	-1	1	1	-1	-1
1	-1	-1	1	1	-1	-1	1
1	1	1	1	-1	-1	-1	-1
1	-1	1	-1	-1	1	-1	1
1	1	-1	-1	-1	-1	1	1
1	-1	-1	1	-1	1	1	-1

For BRR replicate construction, the dimension of the Hadamard matrix must be at least *H*, where *H* denotes the number of first-stage strata in your design. If a Hadamard matrix of a given dimension exists, it is not necessarily unique. Therefore, if you want to use a specific Hadamard matrix, you must provide the matrix as a SAS data set in the HADAMARD= *method-option*. You must ensure that the matrix that you provide is actually a Hadamard matrix; PROC SURVEYPHREG does not check the validity of your Hadamard matrix.

See the section "Balanced Repeated Replication (BRR) Method" on page 9381 for details about how the Hadamard matrix is used to construct replicates for BRR variance estimation.

# **Jackknife Method**

The jackknife method of variance estimation deletes one PSU at a time from the full sample to create replicates. This method is also known as the delete-1 jackknife method because it deletes exactly one PSU in every replicate. The total number of replicates *R* is the same as the total number of PSUs. In each replicate, the sampling weights of the remaining PSUs are modified by the *jackknife coefficient*  $\alpha_r$ . The modified weights are called replicate weights.

Let PSU *i* in stratum  $h_r$  be omitted for the *r*th replicate; then the jackknife coefficient and replicate weights are computed as

$$\alpha_r = \begin{cases} \frac{n_{h_r} - 1}{n_{h_r}} & \text{for a stratified design} \\ \frac{R - 1}{R} & \text{for designs without stratification} \end{cases}$$

$$w_{hij}^{(r)} = \begin{cases} w_{hij} & \text{if observation unit } j \text{ is not in donor stratum } h_r \\ 0 & \text{if observation unit } j \text{ is in PSU } i \text{ of donor stratum } h_r \\ w_{hij}/\alpha_r & \text{if observation unit } j \text{ is not in PSU } i \text{ but in donor stratum } h_i \end{cases}$$

You can use the VARMETHOD=JACKKNIFE(OUTJKCOEFS=) *method-option* to store the jackknife coefficients in a SAS data set and use the VARMETHOD=JACKKNIFE(OUTWEIGHTS=) *method-option* to store the replicate weights in a SAS data set.

If you provide your own replicate weights with a REPWEIGHTS statement, then you can also provide corresponding jackknife coefficients with the JKCOEFS= option. If you provide replicate weights with a REPWEIGHTS statement but do not provide jackknife coefficients, then the procedure uses (R - 1)/R as the default jackknife coefficient for every replicate, where *R* is the total number of replicates.

Let  $\hat{\beta}$  be the estimated proportional hazards regression coefficients from the full sample, and let  $\hat{\beta}_r$  be the estimated regression coefficients for the *r*th replicate. PROC SURVEYPHREG estimates the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

with R - H degrees of freedom, where R is the number of replicates and H is the number of strata, or R - 1 when there is no stratification.

If you specify the CENTER=REPLICATES *method-option*, then PROC SURVEYPHREG computes the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\widehat{\boldsymbol{\beta}}) = \sum_{r=1}^{R} \alpha_r \left(\widehat{\boldsymbol{\beta}}_r - \overline{\widehat{\boldsymbol{\beta}}_r}\right) \left(\widehat{\boldsymbol{\beta}}_r - \overline{\widehat{\boldsymbol{\beta}}_r}\right)'$$

where  $\overline{\hat{\beta}_r}$  is the average of the replicate estimates as follows:

$$\overline{\hat{\beta}_r} = \frac{1}{R} \sum_{r=1}^R \hat{\beta_r}$$

If one or more components of  $\hat{\beta}_r$  cannot be calculated for some replicates, then the variance estimator uses only the replicates for which the proportional hazards regression coefficients can be estimated. Estimability and nonconvergence are two common reasons why  $\hat{\beta}_r$  might not be available for a replicate sample even if  $\hat{\beta}$ is defined for the full sample. Let  $R_a$  be the number of replicates where  $\hat{\beta}_r$  are available, and let  $R - R_a$  be the number of replicates where  $\hat{\beta}_r$  are not available. Without loss of generality, assume that  $\hat{\beta}_r$  is available only for the first  $R_a$  replicates; then the jackknife variance estimator is

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R_a} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

with  $R_a - H$  degrees of freedom, where *H* is the number of strata. Alternatively, you can use the VAD-JUST=AVGREPSS option in the MODEL statement to use the average sum of squares for the invalid replicate samples. See "Variance Adjustment Factors" on page 9386 for details.

and

#### **Replicate Weights Method**

The replicate weights variance estimation method is a general-purpose variance estimation method that uses the replicate weights and replicate coefficients that you provide by using the REPWEIGHTS statement and the REPCOEFS= option, respectively.

If you provide your own replicate weights in a REPWEIGHTS statement but do not specify replicate coefficients in a REPCOEFS= option, then the default replicate coefficient depends on the VARMETHOD= option in the PROC SURVEYPHREG statement as shown in the following table:

Value of VARMETHOD=	Default Replicate Coefficient
None specified	(R-1)/R
BOOTSTRAP	1/R
BRR	1/R
JACKKNIFE	(R-1)/R

Let  $\hat{\beta}$  be the estimated proportional hazards regression coefficients from the full sample, and let  $\hat{\beta}_r$  and  $\alpha_r$  be the estimated regression coefficients and the replicate coefficient for the *r*th replicate, respectively. PROC SURVEYPHREG estimates the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

with R degrees of freedom, where R is the number of replicates.

If you specify the CENTER=REPLICATES *method-option*, then PROC SURVEYPHREG computes the covariance matrix of  $\hat{\beta}$  by

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \overline{\hat{\boldsymbol{\beta}}_r} \right) \left( \hat{\boldsymbol{\beta}}_r - \overline{\hat{\boldsymbol{\beta}}_r} \right)'$$

where  $\overline{\hat{\beta}_r}$  is the average of the replicate estimates as follows:

$$\overline{\hat{\beta}_r} = \frac{1}{R} \sum_{r=1}^R \hat{\beta_r}$$

If one or more components of  $\hat{\beta}_r$  cannot be calculated for some replicates, then the variance estimator uses only the replicates for which the proportional hazards regression coefficients can be estimated. Estimability and nonconvergence are two common reasons why  $\hat{\beta}_r$  might not be available for a replicate sample even if  $\hat{\beta}$ is defined for the full sample. Let  $R_a$  be the number of replicates where  $\hat{\beta}_r$  are available, and let  $R - R_a$  be the number of replicates where  $\hat{\beta}_r$  are not available. Without loss of generality, assume that  $\hat{\beta}_r$  are available only for the first  $R_a$  replicates; then the jackknife variance estimator is

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R_a} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

with  $R_a$  degrees of freedom. Alternatively, you can use the VADJUST=AVGREPSS option in the MODEL statement to use the average sum of squares for the invalid replicate samples. For more information, see "Variance Adjustment Factors" on page 9386.

# **Degrees of Freedom**

PROC SURVEYPHREG uses the degrees of freedom of the variance estimator to obtain t confidence limits and Wald-type F tests. The procedure computes the degrees of freedom based on the variance estimation method, the sample design, and the number of estimable parameters. Alternatively, you can specify the degrees of freedom by using the DF= option in the MODEL statement. This section describes how PROC SURVEYPHREG computes different values of the degrees of freedom based on the variance estimation method and the sample design. For more information about how degrees of freedom depend on the number of estimable parameters and the DF= option in the MODEL statement, see the section "Hypothesis Tests, Confidence Intervals, and Residuals" on page 9389.

For Taylor series variance estimation, PROC SURVEYPHREG calculates the degrees of freedom (df) as the number of clusters minus the number of strata. If the CLUSTER statement is not specified, then the procedure treats each observation as a cluster. If the STRATA statement is not specified, then the procedure assumes that all observations are in the same stratum. These numbers are based on the observations that are included in the analysis. These numbers do not count observations that are excluded from the analysis because they have missing values. For more information, see the section "Missing Values" on page 9377. If you specify the MISSING option in the CLASS statement, missing values are treated as valid nonmissing levels and are included in computing the degrees of freedom. If you specify the NOMCAR option for Taylor series variance estimation, observations that have missing values of the analysis variables are included in computing the degrees of freedom.

If you provide replicate weights by using the REPWEIGHTS statement, the degrees of freedom are equal to the number of replicates used, which is the number of REPWEIGHTS variables that provide replicate estimates. Alternatively, you can specify DF=ALLREPS in the MODEL statement to specify that *df* equals the number of replicates.

For BRR variance estimation (when you do not use the REPWEIGHTS statement), PROC SURVEYPHREG calculates the degrees of freedom as the number of strata. The procedure bases the number of strata on all valid observations in the data set. If some replicate samples are not usable, in the sense that they cannot be used for parameter estimation because of factors such as nonconvergence or inestimability, then *df* equals the minimum of the number of strata and the number of replicates used. Alternatively, you can specify DF=ALLREPS in the MODEL statement to specify that *df* equals the number of strata.

For jackknife variance estimation (when you do not use the REPWEIGHTS statement), PROC SURVEYPHREG calculates the degrees of freedom as the number of clusters minus the number of strata. If you do not specify the CLUSTER statement, then the procedure treats each observation as a cluster. If you do not specify the STRATA statement, then the procedure assumes that all observations are in the same stratum. For jackknife variance estimation, PROC SURVEYPHREG bases the number of strata and clusters on all valid observations in the data set. If some replicate samples are not usable, in the sense that they cannot be used for parameter estimation because of factors such as nonconvergence or inestimability, then *df* equals the number of clusters (or observations if no CLUSTER statement is specified) minus the number of strata (or 1 if no STRATA statement is specified) minus the number of replicate samples that are not used. Alternatively, you can specify DF=ALLREPS in the MODEL statement to specify that *df* equals the number of clusters minus the number of strata.

# **Variance Adjustment Factors**

PROC SURVEYPHREG provides options for adjusting the default variance estimators. VADJUST=NONE and VADJUST=DF are available for the Taylor series linearization variance estimator. VADJUST=AVGREPSS is available for the jackknife replication variance estimators.

For models with large number of parameters, it is reasonable to adjust the Taylor series linearized variance estimator by the number of estimable parameters in the analysis model. Fuller et al. (1989, pp. 77–81) use an adjustment factor (n - 1)/(n - p) to estimate the linearized variance for regression coefficients, where *n* is the total number of observation units and *p* is the number of estimable parameters in the analysis model. By default, PROC SURVEYPHREG uses this adjustment in the computation of the matrix **G** for the Taylor series linearization variance estimation. If you do not want to use this adjustment, then specify VADJUST=NONE.

Variance adjustment factors can be useful for replication variance estimations, especially if some replicate samples are not usable. A replicate sample might not provide useful parameter estimates (replicate estimates) for reasons such as nonconvergence of the optimization or inestimability of some parameters in that subsample. For example, consider the jackknife variance estimator with R replicates. Suppose that only  $R_a(< R)$  replicates are used to obtain replicate estimates and  $R - R_a$  replicates cannot be used due to, say, nonconvergence of the optimization. Without loss of generality, assume that the first  $R_a$  replicates are used. By default SURVEYPHREG uses

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = \sum_{r=1}^{R_a} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)'$$

as the jackknife variance estimator. An alternative estimator is

$$\begin{aligned} \widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) &= \sum_{r=1}^{R_a} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)' + (R - R_a) \left\{ \frac{1}{R_a} \sum_{r=1}^{R_a} \alpha_r \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right) \left( \hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}} \right)' \right\} \\ &= \frac{R}{R_a} \sum_{r=1}^{R_a} \alpha_r (\hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}}) (\hat{\boldsymbol{\beta}}_r - \hat{\boldsymbol{\beta}})' \end{aligned}$$

which uses the average replicate sum of squares for the  $R - R_a$  unusable replicate samples. If you specify the VADJUST=AVGREPSS option, PROC SURVEYPHREG uses the second variance estimator for the jackknife replication method. Note that you can specify the FAY *method-option* for the BRR method to avoid nonconvergence of the optimization or inestimability of some parameters in subsamples.

#### Variance Ratios and Standard Error Ratios

PROC SURVEYPHREG provides options to compute different variance ratios and standard error ratios.

If you specify the VARRATIO=MODEL option, then the procedure computes the variance ratio of the estimated regression parameter  $\hat{\beta}_j$  as  $\frac{\hat{V}(\hat{\beta}_j)}{\hat{V}_M(\hat{\beta}_j)}$ , where  $\hat{V}(\hat{\beta}_j)$ , the estimated variance of  $\hat{\beta}_j$ , uses the complete design information and  $\hat{V}_M(\hat{\beta}_j)$  is the *j*th diagonal element of the observed information matrix  $\mathcal{I}^{-1}(\hat{\beta})$ . If you specify the VARRATIO=IND option, then the procedure computes the variance ratio of the estimated regression parameter  $\hat{\beta}_j$  as  $\frac{\hat{V}(\hat{\beta}_j)}{\hat{V}_{\text{ND}}(\hat{\beta}_j)}$ , where  $\hat{V}(\hat{\beta}_j)$ , the estimated variance of  $\hat{\beta}_j$ , uses the complete design

information and  $\hat{V}_{\text{IND}}(\hat{\beta}_j)$  is the *j*th diagonal element of  $\hat{V}_{\text{IND}}(\hat{\beta})$ .  $\hat{V}_{\text{IND}}(\hat{\beta})$  is the sandwich variance estimator, which ignores the strata and the clusters and is computed as

$$\hat{V}_{\text{IND}}(\hat{\boldsymbol{\beta}}) = \mathcal{I}^{-1}(\hat{\boldsymbol{\beta}}) \left\{ \frac{n}{n-1} (1-f) \sum_{h} \sum_{i} \sum_{j} (\mathbf{e}_{hij} - \bar{\mathbf{e}}_{\cdots})' (\mathbf{e}_{hij} - \bar{\mathbf{e}}_{\cdots}) \right\} \mathcal{I}^{-1}(\hat{\boldsymbol{\beta}})$$

where  $e_{hij}$  are the weighted score residuals, f is the overall sampling fraction, and n is the number of observation units. The three sums are over the observation units (j) across the PSUs (i) and the strata (h).

For Taylor series variance estimation, PROC SURVEYPHREG determines the value of f, the overall sampling fraction, based on the RATE= or TOTAL= option. If you do not specify either of these options, PROC SURVEYPHREG assumes that the value of f is negligible and does not use a finite population correction in the analysis. If you specify RATE=value, PROC SURVEYPHREG uses value as the overall sampling fraction f. If you specify TOTAL=value, PROC SURVEYPHREG computes f as the ratio of the number of PSUs in the sample to the specified total.

If you specify stratum sampling rates by using the RATE=SAS-data-set option, then PROC SURVEYPHREG computes stratum totals based on these stratum sampling rates and the number of sample PSUs in each stratum. The procedure sums the stratum totals to form the overall total and then computes f as the ratio of the number of sample PSUs to the overall total. Alternatively, if you specify stratum totals with the TOTAL=SAS-data-set option, then PROC SURVEYPHREG sums these totals to compute the overall total. The overall sampling fraction f is then computed as the ratio of the number of sample PSUs to the overall total.

The replication methods do not use the finite population correction factor (1 - f) in the denominator.

Standard error ratios are computed as the square root of the variance ratios.

# **Domain Analysis**

*Domain analysis* refers to the computation of statistics for domains (subpopulations). Formation of subpopulations can be unrelated to the sample design, and so the domain sample sizes can actually be random variables. Domain analysis takes this variability into account to compute variance estimates for estimated model parameters. Domain analysis is also known as subgroup analysis, subpopulation analysis, and subdomain analysis. For more information about domain analysis, see Lohr (2010); Särndal, Swensson, and Wretman (1992); Cochran (1977).

To request domain analysis with PROC SURVEYPHREG, use the DOMAIN statement. If your domains are formed by more than one variable, you can specify DomainVariable\_1 \* DomainVariable\_2 in the DOMAIN statement. If you use the DOMAIN statement, the procedure performs separate analyses for all domains, in addition to the overall analysis.

Including the domain variables in a DOMAIN statement request provides a different analysis from that obtained by using a BY statement, which provides completely separate analyses of the BY groups. The BY statement can also be used to analyze the data set by subgroups, but it is critical to note that this does *not* account for random sample sizes that often occur for domain analyses. The BY statement is appropriate only when the number of units in each subgroup is known with certainty. For example, the BY statement can be used to obtain stratum level estimates when you have fixed sample sizes for the strata. When the subgroup sample size is random, include the domain variables in DOMAIN statement.

# Hypothesis Tests, Confidence Intervals, and Residuals

### **Testing the Global Null Hypothesis**

The following statistics are available to test the global null hypothesis  $H_0$ :  $\beta = 0$ . Let *d* be the usual degrees of freedom computed from the survey data by using the number of strata, clusters, or replicate weights; and let *p* be the number of estimable parameters in the null hypothesis  $H_0$ . For more information about computing *d*, see the section "Degrees of Freedom" on page 9386.

The likelihood ratio test is expressed as

$$\chi^2_{\rm LR} = 2 \left[ \log \left\{ L(\hat{\pmb{\beta}}) \right\} - \log \left\{ L(\mathbf{0}) \right\} \right]$$

where  $L(\cdot)$  denotes the partial pseudo-likelihood described in the section "Partial Likelihood Function for the Cox Model" on page 9374 and  $\hat{\beta}$  denotes the estimated regression parameters. The *p*-value is computed by using a chi-square distribution with *p* degrees of freedom. The likelihood ratio statistic is sensitive to the scaling of the weights. The usual assumptions that are required for a likelihood ratio test do not hold for the pseudo-likelihood that is used by PROC SURVEYPHREG (Rao, Scott, and Skinner 1998), leading to other methods of testing the global null hypothesis, such as the Wald test discussed in the following paragraph.

The Wald test uses the variance estimator that accounts for complex sampling such as stratification, clustering, and unequal weighting. Let  $Q = \hat{\beta}' \left[ \hat{V}(\hat{\beta}) \right]^{-1} \hat{\beta}$ , where  $\hat{\beta}$  is the estimated regression parameters and  $\hat{V}(\hat{\beta})$  is the estimated covariance matrix for  $\hat{\beta}$ . You can request the Wald tests that are described in the following table by using the DF= option in the MODEL statement.

Value of DF=	Test Request	Test Statistic	Numerator Degrees of Freedom	Denominator Degrees of Freedom
NONE	Chi-square	Q	р	$\infty$
V	Customized F	vQ/pd	р	V
DESIGN	Unadjusted F	Q/p	р	d
DESIGN (V)	Unadjusted F	Q/p	р	V
PARMADJ	Adjusted F	(d-p+1)Q/pd	р	<i>d</i> – <i>p</i> +1
PARMADJ (V)	Adjusted F	(v-p+1)Q/pv	р	<i>v-p</i> +1
DESIGNADJ	Adjusted F	(d-p+1)Q/pd	р	d

# **Model Fit Statistics**

Suppose the model contains *p* estimable parameters. Then the following two criteria are displayed for model fit statistics:

• -2 log likelihood:

 $-2 \operatorname{Log} L = -2 \log(L(\hat{\beta}))$ 

where L(.) is a partial pseudo-likelihood function for the corresponding TIES= option as described in the section "Partial Likelihood Function for the Cox Model" on page 9374, and  $\hat{\beta}$  is the maximum pseudo-log-likelihood estimate of the proportional hazards regression coefficients.

• Akaike's information criterion (AIC):

$$AIC = -2 \log L + 2p$$

The AIC statistic provides a different way of adjusting the log-likelihood statistic for the number of estimable parameters in the model.

Neither of these criteria is adjusted for the complex sample design, and both criteria are sensitive to the scale of the weights.

# Contrasts

For a testable hypothesis  $H_0$ :  $L\beta = 0$ , you can request different Wald tests by using the DF= option in the MODEL statement.

Let

$$Q = (\mathbf{L}^* \hat{\boldsymbol{\beta}})' (\mathbf{L}^{*'} \widehat{\mathbf{V}} \mathbf{L}^*)^{-1} (\mathbf{L}^* \hat{\boldsymbol{\beta}})$$

where L is a contrast vector or matrix that you specify,  $\hat{\beta}$  is the vector of regression parameters,  $\hat{\beta}$  is the estimated regression coefficients,  $\hat{V}$  is the estimated covariance matrix of  $\hat{\beta}$ , and L<sup>\*</sup> is a matrix such that the following are true:

- L\* has the same number of columns as L.
- L\* has full row rank.
- The rank of L\* equals the rank of the L matrix.
- All rows of L\* are estimable functions.
- The Wald *F* statistic that is computed by using the L\* matrix is equivalent to the Wald *F* statistic computed by using the L matrix.

If L is a full-rank matrix and all rows of L are estimable functions, then  $L^*$  is the same as L. It is possible that such an  $L^*$  matrix cannot be constructed for a given set of linear contrasts, in which case the contrasts are not testable. Let *r* be the rank of L. The following table describes the Wald tests available in PROC SURVEYPHREG.

Value of DF=	Test Request	Test Statistic	Numerator Degrees of Freedom	Denominator Degrees of Freedom
NONE	Chi-square	Q	r	$\infty$
V	Customized F	vQ/rd	r	V
DESIGN	Unadjusted F	Q/r	r	d
DESIGN (V)	Unadjusted F	Q/r	r	V
PARMADJ	Adjusted F	(d-r+1)Q/rd	r	<i>d</i> – <i>r</i> +1
PARMADJ (V)	Adjusted F	(v-r+1)Q/rv	r	<i>v</i> – <i>r</i> +1
DESIGNADJ	Adjusted F	Q/r	r	d

#### **Confidence Intervals**

By default, the SURVEYPHREG procedure computes *t* confidence limits for the estimated regression coefficients. Alternatively, you can specify DF=NONE in the MODEL statement to request standard normal confidence intervals. The *t* confidence interval for a linear combination  $l'\beta$  of the regression coefficients is computed as

$$\left(\mathbf{l}'\hat{\boldsymbol{\beta}} \pm t_{df,\alpha/2}\sqrt{\mathbf{l}'\hat{\mathbf{V}}(\hat{\boldsymbol{\beta}})\mathbf{l}}\right)$$

where  $t_{df,\alpha/2}$  is the 100(1 –  $\alpha/2$ ) percentile point of the *t* distribution with *df* degrees of freedom. See the section "Degrees of Freedom" on page 9386 for more information about *df*. If you use the DF=NONE option in the MODEL statement, then the procedure uses the 100(1 –  $\alpha/2$ ) percentile point of the standard normal distribution.

### **Hazard Ratios**

The hazard ratio for a quantitative effect with regression coefficient  $\beta_j = \mathbf{e}'_j \boldsymbol{\beta}$  is defined as  $\exp(\beta_j)$ , where  $\mathbf{e}_j$  denotes the *j*th unit vector. In general, a log-hazard ratio can be written as  $\mathbf{l}' \boldsymbol{\beta}$ , a linear combination of the regression coefficients, and the hazard ratio  $\exp(\mathbf{l}' \boldsymbol{\beta})$  is obtained by replacing  $e_j$  with **l**.

The confidence intervals for hazard ratios are obtained by exponentiating the confidence limits of the corresponding linear combination. Thus, the  $100(1 - \alpha)$  confidence limits are

$$\exp\left(\mathbf{e}_{j}^{\prime}\hat{\boldsymbol{\beta}}\pm t_{df,\alpha/2}\sqrt{\mathbf{e}_{j}^{\prime}\hat{\mathbf{V}}(\hat{\boldsymbol{\beta}})\mathbf{e}_{j}}\right)$$

where  $t_{df,\alpha/2}$  is the 100(1 –  $\alpha/2$ ) percentile point of the *t* distribution with *df* degrees of freedom. See the section "Degrees of Freedom" on page 9386 for more information about *df*. If you use the DF=NONE option in the MODEL statement, then the procedure uses the 100(1 –  $\alpha/2$ ) percentile point of the standard normal distribution.

#### **Residuals**

This section describes the computation of residuals (RESMART, RESDEV, RESSCH, and RESSCO in the OUTPUT statement). See the section "Notation and Estimation" on page 9369 for definition of notation that is used in this section. The residuals are calculated based on the TIES= option in the MODEL statement.

### TIES=BRESLOW

This is the default option. Let

$$S^{(r)}(\boldsymbol{\beta}, t) = \sum_{A} w_{hij} y_{hij}(t) \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) \mathbf{Z}_{hij}^{\bigotimes r}(t)$$
$$\bar{\mathbf{Z}}(\boldsymbol{\beta}, t) = \frac{S^{(1)}(\boldsymbol{\beta}, t)}{S^{(0)}(\boldsymbol{\beta}, t)}$$

where r = 0, 1; and A be the set of indices in the selected sample.

Further let

$$d\Lambda_0(\boldsymbol{\beta}, t) = \sum_A \frac{w_{hij} dn_{hij}(t)}{S^{(0)}(\boldsymbol{\beta}, t)}$$
  
$$dM_{hij}(\boldsymbol{\beta}, t) = dn_{hij}(t) - y_{hij}(t) \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) d\Lambda_0(\boldsymbol{\beta}, t)$$

The martingale residual at t is defined as

$$\hat{M}_{hij}(t) = \int_0^t dM_{hij}(\hat{\boldsymbol{\beta}}, \tau)$$
  
=  $n_{hij}(t) - \int_0^t y_{hij}(\tau) \exp\left(\hat{\boldsymbol{\beta}}' \mathbf{Z}_{hij}(\tau)\right) d\Lambda_0(\hat{\boldsymbol{\beta}}, \tau)$ 

Here  $\hat{M}_{hij}(t)$  estimates the difference over (0, t] between the observed number of events for the (h, i, j) observation unit and a conditional expected number of events. The quantity  $\hat{M}_{hij} \equiv \hat{M}_{hij}(\infty)$  is referred to as the martingale residual for the (h, i, j) observation unit. For the Cox model with no time-dependent explanatory variables, the martingale residual for the (h, i, j) unit with observation time  $t_{(h,i,j)}$  and event status  $\Delta_{(h,i,j)}$  is

$$\hat{M}_{(h,i,j)} = \Delta_{(h,i,j)} - e^{\hat{\boldsymbol{\beta}}' \mathbf{Z}_{(h,i,j)}} \int_0^{t_{(h,i,j)}} d\Lambda_0(\hat{\boldsymbol{\beta}},s)$$

The deviance residual  $D_{hij}$  for the (h, i, j) observation unit is a transformation of the corresponding martingale residuals,

$$D_{hij} = \operatorname{sign}(\hat{M}_{hij}) \sqrt{2 \left[ -\hat{M}_{hij} - n_{hij}(\infty) \log \left( \frac{n_{hij}(\infty) - \hat{M}_{hij}}{n_{hij}(\infty)} \right) \right]}$$

The square root shrinks large negative martingale residuals, while the logarithmic transformation expands martingale residuals that are close to unity. As such, the deviance residuals are more symmetrically distributed around zero than the martingale residuals. For the Cox model, the deviance residual reduces to the form

$$D_{hij} = \operatorname{sign}(\hat{M}_{hij}) \sqrt{2[-\hat{M}_{hij} - \Delta_{hij} \log(\Delta_{hij} - \hat{M}_{hij})]}$$

The Schoenfeld (1982) residual vector is calculated on a per-event-time basis. At the *k*th event time  $t_{hij,k}$  of the (h, i, j) observation unit, the Schoenfeld residual

$$\hat{\mathbf{U}}_{hij}(t_{hij,k}) = \mathbf{Z}_{hij}(t_{hij,k}) - \bar{\mathbf{Z}}(\hat{\boldsymbol{\beta}}, t_{hij,k})$$

is the difference between the observed covariate vector for the (h, i, j) observation unit and the average of the covariate vectors over the risk set at  $t_{hij,k}$ . Under the proportional hazards assumption, the Schoenfeld residuals have the sample path of a random walk; therefore, they are useful in assessing time trend or lack of proportionality.

The score process for the (h, i, j) subject at time t is

$$\mathbf{L}_{hij}(\boldsymbol{\beta},t) = \int_0^t [\mathbf{Z}_{hij}(\tau) - \bar{\mathbf{Z}}(\boldsymbol{\beta},\tau)] dM_{hij}(\boldsymbol{\beta},\tau)$$

The vector  $\hat{\mathbf{L}}_{hij} \equiv \mathbf{L}_{hij}(\hat{\boldsymbol{\beta}}, \infty)$  is the score residual for the (h, i, j) observation unit.

The score residuals are a decomposition of the first partial derivative of the log likelihood. They are useful in assessing the influence of each subject on individual parameter estimates. They also play an important role in the computation of the variance estimators.

#### **TIES=EFRON**

For TIES=EFRON, the preceding computation is modified to comply with the Efron partial likelihood. For a given uncensored time *t*, let  $\delta_{hij}(t) = 1$  if *t* is an event time for the (h, i, j) observation, and 0 otherwise. Let  $d(t) = \sum_{hij \in A} \delta_{hij}(t)$ , which is the number of observation units that have an event at *t*. For  $1 \le l \le d(t)$ , let

$$S^{(r)}(\boldsymbol{\beta}, l, t) = \sum_{A} w_{hij} y_{hij}(t) \left\{ 1 - \frac{l-1}{d(t)} \delta_{hij}(t) \right\} \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) \mathbf{Z}_{hij}^{\bigotimes r}(t)$$
  

$$\bar{\mathbf{Z}}(\boldsymbol{\beta}, l, t) = \frac{S^{(1)}(\boldsymbol{\beta}, l, t)}{S^{(0)}(\boldsymbol{\beta}, l, t)}$$
  

$$d\Lambda_0(\boldsymbol{\beta}, l, t) = \sum_{A} \frac{w_{hij} dn_{hij}(t)}{S^{(0)}(\boldsymbol{\beta}, l, t)}$$
  

$$dM_{hij}(\boldsymbol{\beta}, l, t) = dn_{hij}(t) - y_{hij}(t) \left(1 - \delta_{hij}(t) \frac{l-1}{d(t)}\right) \exp\left(\boldsymbol{\beta}' \mathbf{Z}_{hij}(t)\right) d\Lambda_0(\boldsymbol{\beta}, l, t)$$

where r = 0, 1, and A are the set of indices in the selected sample.

The martingale residual at t for the (h, i, j) observation unit is defined as

$$\hat{M}_{hij}(t) = \int_{0}^{t} \frac{1}{d(\tau)} \sum_{l=1}^{d(\tau)} dM_{hij}(\hat{\beta}, l, \tau) = n_{hij}(t) - \int_{0}^{t} \frac{1}{d(\tau)} \sum_{l=1}^{d(\tau)} y_{hij}(\tau) \left(1 - \delta_{hij}(\tau) \frac{l-1}{d(\tau)}\right) \exp\left(\hat{\beta}' \mathbf{Z}_{hij}(\tau)\right) d\Lambda_{0}(\hat{\beta}, l, \tau)$$

Deviance residuals are computed by using the same transform on the corresponding martingale residuals as in TIES=BRESLOW.

The Schoenfeld residual vector for the (h, i, j) observation unit at event time  $t_{hij,k}$  is

$$\hat{\mathbf{U}}_{hij}(t_{hij,k}) = \mathbf{Z}_{hij}(t_{hij,k}) - \frac{1}{d(t_{hij,k})} \sum_{l=1}^{d(t_{hij,k})} \bar{\mathbf{Z}}(\hat{\boldsymbol{\beta}}, l, t_{hij,k})$$

The score process for the (h, i, j) observation unit at time t is

$$\mathbf{L}_{hij}(\boldsymbol{\beta},t) = \int_0^t \frac{1}{d(\tau)} \sum_{l=1}^{d(\tau)} \left( \mathbf{Z}_{hij}(\tau) - \bar{\mathbf{Z}}(\boldsymbol{\beta},l,\tau) \right) dM_{hij}(\boldsymbol{\beta},l,\tau)$$

# **Output Data Sets**

You can use the Output Delivery System to create a SAS data set from any piece of PROC SURVEYPHREG output. See the section "ODS Table Names" on page 9398 for more information. PROC SURVEYPHREG also provides an output data set to store observation-level statistics, an output data set to store the replicate weights for BRR or jackknife variance estimation, and an output data set to store the jackknife coefficients for jackknife variance estimation.

# **OUT= Data Set for the OUTPUT statement**

The OUTPUT statement can be used to store observation-level statistics, such as the predicted values and their standard errors, the (weighted) number of observation units at risk, martingale residuals, Schoenfeld residuals, score residuals, and deviance residuals. See the section "Residuals" on page 9391 for details about how these statistics are calculated.

# **Replicate Weights Output Data Set**

If you specify the OUTWEIGHTS= *method-option* for VARMETHOD=BRR or JACKKNIFE, PROC SURVEYPHREG stores the replicate weights in an output data set. The OUTWEIGHTS= output data set contains all observations that are used in the analysis or all valid observations in the DATA= input data set. See the section "Missing Values" on page 9377 for details about valid observations.

The OUTWEIGHTS= data set contains the following variables:

- all variables in the DATA= input data set
- RepWt\_1, RepWt\_2, ..., RepWt\_R, which are the replicate weight variables, where R is the total number of replicates in the analysis

Each replicate weight variable contains the replicate weights for the corresponding replicate. Replicate weights equal zero for those observations not included in the replicate.

After the procedure creates replicate weights for a particular input data set and survey design, you can use the OUTWEIGHTS= *method-option* to store these replicate weights and then use them again in subsequent analyses, either in PROC SURVEYPHREG or in the other survey procedures. You can use the REPWEIGHTS statement to provide replicate weights for the procedure.

# Jackknife Coefficients Output Data Set

If you specify the OUTJKCOEFS= *method-option* for VARMETHOD=JACKKNIFE, PROC SURVEYPHREG stores the jackknife coefficients in an output data set. The OUTJKCOEFS= output data set contains one observation for each replicate. The OUTJKCOEFS= data set contains the following variables:

- Replicate, which is the replicate number for the jackknife coefficient
- JKCoefficient, which is the jackknife coefficient for the replicate
- DonorStratum, which is the stratum of the PSU that was deleted to construct the replicate, if you specify a STRATA statement

After the procedure creates jackknife coefficients for a particular input data set and survey design, you can use the OUTJKCOEFS= *method-option* to store these coefficients and then use them again in subsequent analyses, either in PROC SURVEYPHREG or in the other survey procedures. You can use the JKCOEFS= option in the REPWEIGHTS statement to provide jackknife coefficients for the procedure.

# **Displayed Output**

If you use the NOPRINT option in the PROC SURVEYPHREG statement, the procedure does not display any output. Otherwise, PROC SURVEYPHREG displays results of the analysis in a collection of tables.

### **Model Information**

The "Model Information" table displays the two-level name of the input data set, the name and label of the failure time variable, the name and label of the censoring variable and the values that indicate censored times, the model, the name and label of the FREQ variable, the name and label of the WEIGHT variable, the name and label of the STRATA variables, the name and label of the CLUSTER variables, and the method of handling ties in the failure time for the Cox model. The ODS name of the "Model Information" table is ModelInfo.

### **Number of Observations**

The "Number of Observations" table displays the number of observations that are read and used, the sum of frequencies read and used, the sum of weights read and used, and the weighted sum of frequencies that are read and used in the analysis. The ODS name of the "Number of Observations" table is NObs.

### Summary of the Number of Event and Censored Values

The "Summary of the Number of Event and Censored Values" table displays the number of events and censored values. The ODS name of the "Summary of the Number of Event and Censored Values" table is CensoredSummary.

### Summary of the Weighted Number of Event and Censored Values

The "Summary of the Weighted Number of Event and Censored Values" table displays the weighted number of events and censored values. The ODS name of the "Summary of the Weighted Number of Event and Censored Values" table is WeightedCensoredSummary.

# **Class Level Information**

The "Class Level Information" table is displayed when there are CLASS variables in the model. The table lists the categories of every CLASS variable that is used in the model and the corresponding design variable values. The ODS name of the "Class Level Information" table is ClassLevelInfo.

# **Design Summary Table**

The "Design Summary" table provides information about the sample design. The table displays the total number of strata that are read and used, and the total number of clusters read and used. The table is displayed only if you specify a STRATA or CLUSTER statement. The ODS name of the "Design Summary" table is DesignSummary.

# **Stratum Information Table**

If you specify the LIST option in the STRATA statement, PROC SURVEYPHREG displays a "Stratum Information" table. The ODS name of the "Stratum Information Table" is StrataInfo. This table provides the following information for each stratum:

- Stratum Index, which is a sequential stratum identification number
- STRATA variables, which list the levels of STRATA variables for the stratum
- Number of Observations, which is the number of observations used in the stratum
- Population Total for the stratum, if you specify the TOTAL= option
- Sampling Rate for the stratum, if you specify the TOTAL= or RATE= option. If you specify the TOTAL= option, the sampling rate is based on the number of valid observations in the stratum.
- Number of Clusters, which is the number of clusters in the stratum, if you specify a CLUSTER statement

# **Convergence Status**

The "Convergence Status" table displays the convergence status of the optimization routine. The procedure displays this table only when you specify the NLOPTIONS statement. The ODS name of the "Convergence Status" table is ConvergenceStatus.

# **Model Fit Statistics**

The "Model Fit Statistics" table displays the values of  $-2 \log$  likelihood and the AIC for the null model and the fitted model. The ODS name of the "Model Fit Statistics" table is FitStatistics.

# Testing Global Null Hypothesis: BETA=0

The "Testing Global Null Hypothesis: BETA=0" table displays results of the likelihood ratio test and the Wald test for testing the hypothesis that all parameters are zero. The ODS name of the "Testing Global Null Hypothesis: BETA=0" table is GlobalTests.

# Analysis of Maximum Likelihood Estimates

The "Analysis of Maximum Likelihood Estimates" table displays the denominator degrees of freedom, which is computed as described in the section "Degrees of Freedom" on page 9386; the maximum likelihood estimate of the parameter; the estimated standard error, computed as the square root of the corresponding diagonal element of the estimated covariance matrix; the *t* statistic, computed as the parameter estimate divided by the standard error; the *p*-value of the *t* statistic with respect to a *t* distribution with denominator degrees of freedom; and the hazard ratio estimate. The *t* confidence limits for the parameter estimates and estimated hazard ratios are displayed if you specify the CLPARM or RISKLIMITS option in the MODEL statement. You can specify the DF=NONE option in the MODEL statement to request *p*-values and confidence intervals from a standard normal distribution. If you specify the VARRATIO=ALL | MODEL | IND option in the MODEL statement, then the variance ratios for model or independence (or both) are displayed. If you specify the SERATIO=ALL | MODEL | IND option in the MODEL statement, then the standard error ratios for model or independence (or both) are displayed.

The ODS name of the "Analysis of Maximum Likelihood Estimates" table is ParameterEstimates.

### **Covariance Matrix**

The "Covariance Matrix" table is displayed if you specify the COVB option in the MODEL statement. The table contains the estimated covariance matrix for the parameter estimates. The ODS name of the "Covariance Matrix" table is CovB.

### **Hessian Matrix**

The "Hessian Matrix" table is displayed if you specify the HESS option in the MODEL statement. The table contains the Hessian matrix that is evaluated at the estimated regression parameters. The ODS name of the "Hessian Matrix" table is Hessian.

### **Inverse Hessian Matrix**

The "Inverse Hessian Matrix" table is displayed if you specify the INVHESS option in the MODEL statement. The table contains the inverse of the Hessian matrix evaluated at the estimated regression parameters. The ODS name of the "Inverse Hessian Matrix" table is InvHessian.

# **Variance Estimation Table**

The "Variance Estimation" table provides the following information:

- Method, which is the variance estimation method—Taylor Series, Balanced Repeated Replication, or Jackknife
- Replicate Weights input data set name, if you provide replicate weights with a REPWEIGHTS statement
- Number of Replicates, for VARMETHOD=BRR or VARMETHOD=JACKKNIFE
- Hadamard Data Set name, if you specify the HADAMARD= method-option for VARMETHOD=BRR
- Fay Coefficient, if you specify the FAY method-option for VARMETHOD=BRR
- Missing Values Included (MISSING), if you specify the MISSING option
- Missing Values Included (NOMCAR), if you specify the NOMCAR option
- Missing Values Excluded, if you have missing values and you do not specify the NOMCAR option

The ODS name of the "Variance Estimation" table is VarianceEstimation.

### **Hadamard Matrix**

If you specify the PRINTH *method-option* for VARMETHOD=BRR, PROC SURVEYPHREG displays the Hadamard matrix that is used to construct replicates for BRR variance estimation. If you provide a Hadamard matrix with the HADAMARD= *method-option* for VARMETHOD=BRR but the procedure does not use the entire matrix, the procedure displays only the rows and columns that are actually used to construct replicates. The ODS name of the "Hadamard Matrix" table is HadamardMatrix.

### Maximum Likelihood Estimates for Replicate Samples

If you specify the DETAILS *method-option* for VARMETHOD=BRR or the DETAILS *method-option* for VARMETHOD=JACKKNIFE, PROC SURVEYPHREG displays the "Maximum Likelihood Estimates for Replicate Samples" table. The Replicate Number column displays the replication number, the Replicate Weight column displays the name of the replicate weight variable, and the Status column displays the convergence status. The replicate number for the full sample is set to 0. If you do not specify replicate weights, then PROC SURVEYPHRG uses default names to identify the replicate weights. For more information, see "Replicate Weights Output Data Set" on page 9394. The convergence status is 1 if the maximum likelihood estimates are not available for a replicate sample and 0 otherwise. If the maximum likelihood estimates are not available for a replicate sample, then the parameter estimates are set to missing for that replicate. The ODS name of the "Maximum Likelihood Estimates for Replicate Samples" table is RepEstimates.

# **ODS Table Names**

PROC SURVEYPHREG assigns a name to each table it creates. You can use this name to refer to the table when using the Output Delivery System (ODS) to select tables and create output data sets. For more information about ODS, see Chapter 20, "Using the Output Delivery System." Table 115.8 lists the table names, along with the corresponding analysis options.

ODS Table Name	Description	Statement / Option
CensoredSummary	Summary of event and censored	Default
	observations	
ClassLevelInfo	CLASS variable levels	CLASS
ConvergenceStatus	Convergence status	NLOPTIONS / PALL
CovB	Covariance of parameter estimates	MODEL / COVB
DesignSummary	Design summary	STRATA or CLUSTER
FitStatistics	Model fit statistics	Default
GlobalTests	Tests of the global null	Default
	hypothesis	
Hadamard	Hadamard matrix	PROC /
		VARMETHOD=BRR(PRINTH)
Hessian	Observed Hessian matrix	MODEL / HESSIAN
InvHessian	Inverse Hessian matrix	MODEL / INVHESS
IterHist	Iteration history	NLOPTIONS / PHISTORY
ModelInfo	Model information	Default
NObs	Number of observations	Default
ParameterEstimates	Maximum likelihood estimates	Default
ParameterEstimatesStart	Initial parameter values	NLOPTIONS / PALL
RepEstimates	Maximum likelihood estimates for	PROC /
_	replicate samples	VARMETHOD=BRR(DETAILS)
		or
		VARMETHOD=JACKKNIFE(DETAILS
StrataInfo	Stratum information	STRATA / LIST

Table 115.8 ODS Tables Produced by PROC SURVEYPHREG

ODS Table Name	Description	Statement / Option
VarianceEstimation WeightedCensoredSummary	Variance estimation Summary of weighted number of event and censored observations	Default WEIGHT

Table 115.8 continued

# **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."

Before you create graphs, ODS Graphics must be enabled (for example, by specifying the ODS GRAPH-ICS ON statement). For more information about enabling and disabling ODS Graphics, see the section "Enabling and Disabling ODS Graphics" on page 607 in Chapter 21, "Statistical Graphics Using ODS."

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" on page 606 in Chapter 21, "Statistical Graphics Using ODS."

When ODS Graphics is enabled, the ESTIMATE, LSMEANS, LSMESTIMATE, and SLICE statements can produce plots that are associated with their analyses. For information about these plots, see the corresponding sections of Chapter 19, "Shared Concepts and Topics."

# **Examples: SURVEYPHREG Procedure**

# Example 115.1: Analysis of Clustered Data

When experimental units are naturally or artificially clustered, failure times of experimental units within a cluster are correlated. Lee, Wei, and Amato (1992) estimate the regression parameters in the Cox model by maximizing a partial likelihood function under an independent working correlation assumption and estimate the variance of the estimated regression coefficients by using a robust sandwich variance estimator that accounts for the intracluster dependence.

The Diabetic Retinopathy Study (DRS) is a randomized, controlled clinical trial of more than 1,700 patients across 15 medical centers. One objective of this study was to determine if photocoagulation treatment delays the occurrence of blindness. One eye of each patient was randomly assigned to treatment and the other eye to control. For more information about the data set and a similar analysis, see Example 86.11 in Chapter 86, "The PHREG Procedure."

Each patient is a cluster that contributes two observations to the input data set, one for each eye. The following variables are available:

• ID, patient's identification

- Time, failure time
- Status, event indicator (0=censored, and 1=uncensored)
- Treatment, treatment received (1=laser photocoagulation, and 0=otherwise)
- DiabeticType, type of diabetes (0=juvenile onset with age of onset at 20 or under, and 1= adult onset with age of onset over 20)

The following DATA step creates the data set Blind, which represents 197 diabetic patients from the DRS:

```
data Blind;
  input ID Time Status DiabeticType Treatment @@;
  datalines;
  5 46.23 0 1 1
                  5 46.23 0 1 0
                                  14 42.50 0 0 1
                                                  14 31.30 1 0 0
 16 42.27 0 0 1 16 42.27 0 0 0
                                  25 20.60 0 0 1
                                                  25 20.60 0 0 0
 29 38.77 0 0 1 29 0.30 1 0 0 46 65.23 0 0 1
                                                  46 54.27 1 0 0
  49 63.50 0 0 1 49 10.80 1 0 0 56 23.17 0 0 1 56 23.17 0 0 0
 61 1.47 0 0 1 61 1.47 0 0 0 71 58.07 0 1 1 71 13.83 1 1 0
 100 46.43 1 1 1 100 48.53 0 1 0 112 44.40 0 1 1 112 7.90 1 1 0
 120 39.57 0 1 1 120 39.57 0 1 0 127 30.83 1 1 1 127 38.57 1 1 0
 133 66.27 0 1 1 133 14.10 1 1 0 150 20.17 1 0 1 150 6.90 1 0 0
 167 58.43 0 1 1 167 41.40 1 1 0 176 58.20 0 0 1 176 58.20 0 0 0
 185 57.43 0 1 1 185 57.43 0 1 0 190 56.03 0 0 1 190 56.03 0 0 0
   ... more lines ...
1705 8.00 0 0 1 1705 8.00 0 0 0 1717 51.60 0 1 1 1717 42.33 1 1 0
1727 49.97 0 1 1 1727 2.90 1 1 0 1746 45.90 0 0 1 1746 1.43 1 0 0
1749 41.93 0 1 1 1749 41.93 0 1 0
```

The following statements request a proportional hazards regression of Time on Treatment, DiabeticType, and the Treatment  $\times$  DiabeticType interaction, with Status as the censoring indicator. The CLUSTER statement indicates the observations that came from the same patient.

```
proc surveyphreg data=Blind;
  model Time*Status(0) = Treatment DiabeticType Treatment*DiabeticType;
  cluster id;
run;
```

Output 115.1.1 displays some summary information. There are 394 observations and 197 patients (clusters). Almost 61% of the observations are censored. The *p*-values for the null model are less than 0.0001 for both the likelihood ratio test and the Wald test (Output 115.1.2), indicating that the survival time is highly dependent on Treatment and DiabeticType. In this example, the likelihood ratio statistic has an approximate chi-square distribution with 3 degrees of freedom, and the Wald statistic has an approximate *F* distribution with 3 numerator degrees of freedom and 194 denominator degrees of freedom. The denominator degrees of freedom are calculated as the number of clusters (197) minus the number of estimable parameters (3).

Output 115.1.1	Summary	Information
----------------	---------	-------------

#### The SURVEYPHREG Procedure

Numb	er of Obsei	vations I	Read	394
Numb	er of Obsei	vations	Jsed	394
_	Design	Summary	/	
Ν	umber of C	lusters	197	
C	e			
Sumr	and Censo	Number pred Valu	of Ev es	ent
Total	and Censo Event Cer	Number ored Valu	of Ev es Perc Censo	ent ent red
Total 394	and Censo Event Cer 155	Number ored Valu nsored C 239	of Ev es Perc Censo	ent ent red
Total 394	and Censo Event Cer 155	Number ored Valu nsored C 239	of Ev es Perc Censo 60	ent red
Total 394	and Censo Event Cer 155 Variance	Number ored Valu nsored C 239 Estimatic	of Ev es Perc Censo 60	ent ent red

Output 115.1.2 Global Test Results

Testing Global Null Hypothesis: BETA=0					
	Test				
Test	Statistic	Num DF	Den DF	p-Value	
Likelihood Ratio	28.4556	3	Infty	<.0001	
Wald	11.4455	3	194	<.0001	

Output 115.1.3 displays parameter estimates, standard errors, *t* statistics, denominator degrees of freedom, *p*-values, and hazard ratios. In this example data set, Treatment and Treatment × DiabeticType interaction are significant with *p*-values 0.023 and 0.006, respectively. Since the model contains Treatment × DiabeticType interaction, the exponential of the estimated regression coefficient is not the hazard ratio. Use the ESTIMATE statement to calculate the hazard ratios.

|--|

Analysis of Maximum Likelihood Estimates									
	Standard								
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio			
Treatment	196	-0.424672	0.185438	-2.29	0.0231	0.654			
DiabeticType	196	0.340841	0.196076	1.74	0.0837	1.406			
Treatment*DiabeticTy	196	-0.845665	0.304303	-2.78	0.0060	0.429			

# **Example 115.2: Stratification, Clustering, and Unequal Weights**

This example uses a data set from the National Health and Nutrition Examination Survey I (NHANES I) Epidemiologic Followup Study (NHEFS). The NHEFS is a national longitudinal survey that is conducted by the National Center for Health Statistics, the National Institute on Aging, and some other agencies of the Public Health Service in the United States. Some important objectives of this survey are to determine the relationships between clinical, nutritional, and behavioral factors; to determine mortality and hospital utilizations; and to

monitor changes in risk factors for the initial cohort that represents the NHANES I population. A cohort of size 14,407, which includes all persons 25 to 74 years old who completed a medical examination at NHANES I in 1971–1975, was selected for the NHEFS. Personal interviews were conducted for every selected unit during the first wave of data collection from the year 1982 to 1984. Follow-up studies were conducted in 1986, 1987, and 1992. In the year 1986, only nondeceased persons 55 to 74 years old (as reported in the base year survey) were interviewed. The 1987 and 1992 NHEFS contain the entire nondeceased NHEFS cohort. Vital and tracing status data, interview data, health care facility stay data, and mortality data for all four waves are available for public use. See http://www.cdc.gov/nchs/nhanes/nhefs/nhefs.htm for more information about the survey and the data sets.

For illustration purposes, 1,018 observations from the 1987 NHEFS public use interview data are used to create the data set cancer. The observations are obtained from 10 strata that contain 596 PSUs. The sum of observation weights for these selected units is over 19 million. Observation weights range from 359 to 129,359 with a mean of 18,747.69 and a median of 11,414. Several observation weights have large values; therefore it is reasonable to rescale the observation weights to facilitate the optimization routine. Different scaling techniques are proposed in the literature. For example, Binder (1992) uses scaled weights such that the sum of weights over the sampled units is one. Without loss of generality, the analysis weights in this example are obtained by dividing each observation weight by a large number (130,000). Because of this rescaling, you must be careful interpreting some results from PROC SURVEYPHREG.

The following variables are used in this example:

- ObsNo, unit identification
- Strata, stratum identification
- PSU, identification for primary sampling units
- ObservationWt, sampling weight associated with each unit
- AnalysisWt, obtained from the sampling weights by dividing each ObservationWt by 130,000
- Smoke, smoking status (-1 = not applicable, 1 = never smoked, 2 = current or former smoker in 1982-1984 follow-up, and 3 = current or former smoker in 1987 follow-up)
- Age, the event-time variable, defined as follows:
  - age of the subject when the first cancer was reported for subjects with reported cancer
  - age of the subject at death for deceased subjects without reported cancer
  - age of the subject as reported in 1987 follow-up (this value is used for nondeceased subjects who never reported cancer)
  - age of the subject for the entry year 1971–1975 survey if the subject has cancer (or is deceased) but the date of incident is not reported
- Cancer, cancer indicator (1 = cancer reported, 0 = cancer not reported)
- BodyWeight, body weight of the subject as reported in the 1987 follow-up, or an imputed body weight based on the subject's age in the entry year 1971–1975 survey

The following SAS statements create the data set cancer. Note that BodyWeight for a few observations (8%) is imputed based on Age by using a deterministic regression imputation model (Särndal and Lundström (2005, chapter 12)). The imputed values are treated as observed values in this example. In other words, this example treats the data set cancer as the observed data set.

```
data cancer;
  input ObsNo Strata PSU AnalysisWt ObservationWt Smoke
        Age Cancer BodyWeight;
  datalines;
    3 002 0.02927
                      3805
                                     175
  1
                          2 53 1
                      6107
    3
       002 0.04698
                            2
                               77 0
  2
                                     175
  3
     3
       039 0.02283
                      2968
                           2
                               50
                                  0
                                     160
  Δ
    3
       084 0.23414
                      30438 2 52 0 145
    3
       007 0.03908
                      5081 1 80 0 127
  5
    3
       009 0.02993
  6
                      3891
                            1 62 0 180
       009 0.02754
                      3580 2
                               50 0
  7
     3
                                    157
                      2968 2 56 0 142
  8
    3 022 0.02283
  9
     3 050 0.18268
                      23748 2 60 0 140
  ... more lines ...
1016 4 002
            0.02068
                      2689 2 40 0
                                    120
                    45888
1017 4
       092
             0.35298
                            2 52 0
                                     166
1018 4 035
             0.03344
                      4347 -1 58 0
                                    156
:
```

Suppose you want to study the occurrence of cancer for the base year survey population and its relation to smoking status and body weight. The following statements request a proportional hazards regression of Age on BodyWeight and Smoke with Cancer as the censor indicator. The STRATA, CLUSTER, and WEIGHT statements identify the variance strata, PSUs, and analysis weights respectively. The CLASS statement specifies that Smoke is a categorical variable, and the MODEL statement provides information about the analysis model. The TIES= option in the MODEL statement requests the Efron likelihood to handle tied events. If you do not specify the TIES= option in the MODEL statement is used to display the iteration history table. The ESTIMATE statement computes a contrast between subjects who are reported as current (or former) smokers and the others. The EXP option in the ESTIMATE statement requests that the linear contrast be estimated in the exponential scale, which is the hazard ratio. The TEST statement requests the Type 3 test for each effect that is specified in the MODEL statement.

```
proc surveyphreg data = cancer;
  strata strata;
  cluster psu;
  weight analysiswt;
  class smoke;
  model age*cancer(0) = bodyweight smoke / ties = efron;
  nloptions phistory;
  estimate smoke 0.5 0.5 -0.5 -0.5 / exp;
  test ;
run;
```

Some summary statistics are shown in Output 115.2.1. The "Model Information" table contains information about the model such as the names for the dependent and censoring variables, and the likelihood. The "Number of Observations" table displays the number of observations and the sum of weights. A total of 1,018 observations are read from the cancer data set, but one observation is not used in the analysis because it has a zero sampling weight. The sum of weights is 146.81, which gives an estimated population size of 19,085,105 (=  $146.8085 \times 130,000$ ). Note that the estimated population size would be 19,085,151 if you use the sampling weights (ObservationWt) instead of the analysis weights (AnalysisWt). The difference is due to the rounding errors in AnalysisWt. For simplicity, analysis weights are rounded at the fifth decimal place. The

"Design Summary" table shows that there are 596 PSUs and 10 strata. From the censored summary tables, 11.7% subjects in the sample have reported cancer and an estimated 11.6% subjects in the study population have cancer. The "Variance Estimation" table shows that the Taylor series linearization variance estimation method is used and the observation units with missing values are excluded from the analysis. Note that the only missing unit in this data set has a zero sampling weight and hence it is not included in the analysis.

Output 115.2.1 Model Information, Data Summary, Design Summary, and Information about Variance Estimation

Model Info	rmation	
Data Set	WORK.	CANCER
Dependent Variable	Age	
Censoring Variable	Cancer	
Censoring Value(s)	0	
Weight Variable	Analysis	sWt
Stratum Variable	Strata	
Cluster Variable	PSU	
Ties Handling	EFRON	
Number of Observatio	ns Read	1018
Sum of Weights Read	113 0300	146 8085
Sum of Weights Used		146 8085
Sum of Weights OSca		140.0005
Design Su	mmary	
Number of Stra	ata	10
Number of Clu	sters 5	596
Summary of the Nu and Censore	umber of d Values	f Event s
Total Event Cone	F F	Percent
1017 119	898	88.30
Summary of the Weig Event and Cens	ghted Nu ored Val	umber of ues
Table Frank C		Percent
	ensored	Censored
	/u /unn	88.41
140.0005 17.01105 12	23.7300	
Variance Es	timation	
Variance Es Method	timation Taylor Se	eries

#### The SURVEYPHREG Procedure

The "Iteration History" table in Output 115.2.2 shows that the procedure converged after four iterations. The "Objective Function" column contains the value of the likelihood after every iteration. The "Objective Function Change" column measures the change in the objective function between iterations; however, this is not the monitored convergence criterion. The SURVEYPHREG procedure monitors several features simultaneously to determine whether to stop an optimization.

Maximum Likelihood Iteration History								
		Function	Active	Objective	Objective Function	Max Abs Gradient		Ratio Between Actual and Predicted
Iteration	Restarts	Calls	Constraints	Function	Change	Element	Ridge	Change
1	0	4	0	-63.34004	1.6501	21.9620	0	0.916
2	0	6	0	-63.29819	0.0418	0.2005	0	1.052
3	0	8	0	-63.29776	0.000430	0.00293	0	1.012
4	0	10	0	-63.29776	1.528E-7	1.102E-6	0	1.000

#### Output 115.2.2 Iteration History

Estimates for proportional hazards regression coefficients and their standard errors are shown in Output 115.2.3. The categorical variable Smoke has four levels, and GLM parameterization is used by PROC SURVEYPHREG. You can use the PARAM= option in the CLASS statement to specify other types of parameterizations. The estimated regression coefficient for BodyWeight is 0.012 with a standard error of 0.003. The degrees of freedom for the *t* test are equal to the number of PSUs (596) minus the number of strata (10). The "Estimates" table displays the estimated contrast and the corresponding hypothesis test. The estimated value for the contrast is -0.75. The estimated hazard for the nonsmokers is 0.47 times the estimated hazard for the current or former smokers. In this example data set, the contrast of interest is not significant at 0.05 levels. The "Type III Tests of Model Effects" table displays the Type 3 analysis. The effect variable Smoke has four levels. The F Value for Smoke is 1.49 with three numerator degrees of freedom and 584 denominator degrees of freedom.

#### Output 115.2.3 Parameter Estimates

Analysis of Maximum Likelihood Estimates									
Parameter	DF	Estimate	Standard Error	t Value	Pr >  t	Hazard Ratio			
BodyWeight	586	0.011920	0.003150	3.78	0.0002	1.012			
Smoke -1	586	-1.174048	0.738358	-1.59	0.1124	0.309			
Smoke 1	586	-1.006515	0.577955	-1.74	0.0821	0.365			
Smoke 2	586	-0.674183	0.557587	-1.21	0.2271	0.510			
Smoke 3	586	0				1.000			

	Num	Den		
Effect	DF	DF	F Value	Pr > F
BodyWeight	1	586	14.32	0.0002
Smoke	3	584	1.49	0.2162

			Estir	nate		
		Standard				
Label	Estimate	Error	DF	t Value	Pr >  t	Exponentiated
Row 1	-0.7532	0.3864	586	-1.95	0.0518	0.4709

# Example 115.3: Domain Analysis

This example uses a data set from the NHANES I Epidemiologic Followup Study (NHEFS); see Example 115.2 for more information about the NHEFS.

For illustration purposes, 1,891 observations from the 1992 NHEFS vital and tracing status data set are used to estimate the regression coefficients of a proportional hazards model. The observations are obtained from 22 strata; each stratum contains either two or three primary sampling units. The sum of observation weights for these selected units is almost 103 million. Observation weights range from 1,498 to 470,154 with a mean of 54,457.11 and a median of 45,246. The following variables are used in this example. Although this example uses the observation weights directly, Binder (1992) suggests that a scaled version of the observation weights would be useful to improve the performance of the optimization routine.

The following variables are created in the data set mortality:

- ID, unit identification
- VARSTRATA, stratum identification
- VARPSU, identification for primary sampling units
- SWEIGHT, sampling weight associated with each unit
- AGE, the subject's reported age at the 1992 interview if the subject was alive at that time; otherwise, the subject's age at death
- VITALSTATUS, vital status of subject in 1992 (1 = alive, 3 = dead, 4 = unknown, 5 = traced alive with direct subject contact, 6 = traced alive without direct subject contact)
- POVARIND, indicator for poverty area where subject's household was located at NHANES I (1971– 1975) exam, (1 = poverty area, 2 = non-poverty area)
- GENDER, (1 = male, 2 = female)

```
data mortality;
```

;

```
input ID VARSTRATA VARPSU SWEIGHT AGE VITALSTATUS POVARIND GENDER;
datalines;
  1 03 1 13312
                    66 1
                            1
                               1
  2
     03 1
           7941
                    71 3
                            1
                               2
                     . 4
                               1
  3 03 1 16048
                           1
  4
     03
         3
            9298
                    58 3
                            1
                               1
  5
     03 2 15336
                    56 3
                           1
                               2
     03 1 14744
                    63 1
                           1
                               1
  6
  7
     03 2 83729
                    70 1
                           2
                               2
                    57 1
                            2
                               1
  8
     03
         3 106492
     03 3 78083
                            2
  9
                    81 3
                               2
 10
     03 3 55957
                    79 3
                            2
                               1
... more lines ...
1890 13 1 88939 59
                      1
                            1
                         2
1891 13 1 59218 75 1
                         2
                            2
```
Suppose you want to estimate the hazard function for mortality time after adjusting for the poverty area indicator in the base year survey population. The following SAS statements request a proportional hazards regression of age (AGE) on poverty indicator (POVARIND):

```
proc surveyphreg data = mortality nomcar;
    class povarind;
    strata varstrata;
    cluster varpsu;
    weight sweight;
    model age*vitalstatus(1 4 5 6) = povarind;
    domain gender;
run;
```

Subjects with VITALSTATUS 1, 4, 5, or 6 are considered alive. The CLASS statement specifies that POVARIND is a categorical variable, the WEIGHT statement identifies the sampling weights, the STRATA statement identifies variance strata, and the CLUSTER statement identifies variance PSUs. The DOMAIN statement requests three separate analyses: for the overall data set, the male subpopulation, and the female subpopulation respectively. There are 223 observation units with missing values on age. All the units with missing age have vital status 1, 4, 5, or 6. Therefore, these subjects are considered to be alive in the current survey year 1992. Age for every observation unit in the base year survey was known from 1971–1975 NHANES I. One reasonable approach is to determine the age of these 223 units based on their age from the NHANES I data set. However, for illustration purposes, this example does not include the observation units with missing age when estimating the regression coefficients. Instead, an analysis of just the set of respondents is requested by specifying the NOMCAR option in the PROC SURVEYPHREG statement. This option uses a variance estimator that accounts for the random size of the set of respondents.

Output 115.3.1 shows summary statistics for the overall analysis. A total of 1,891 observations are read from the input DATA= data set, but only 1,668 observations are used in the analysis. The remaining 223 observations have missing values in the variable age. The respondent data set represents almost 89.5 million units in the population. There are 22 strata and 55 clusters. Although only 57% observation units in the sample are alive, an estimated 69% observation units in the population are alive. This difference is reasonable because selection probabilities for observation units are not the same. If you do not use the sampling weights, then your sample-based estimators might be biased for the corresponding finite population quantities. The "Variance Estimation" table indicates that the NOMCAR option is used for variance estimation.

Output 115.3.1 Summary Statistics for the Entire Population

## The SURVEYPHREG Procedure

Number of Observations Read	1891
Number of Observations Used	1668
Sum of Weights Read	1.0298E8
Sum of Weights Used	89439590

Design Summary					
Number of Strata	22				
Number of Clusters	55				

Out	put 1	15.3	<b>8.1</b> cc	ontii	nuec	1		
Summary of the Number of Event and Censored Values								
Total	Event	Cer	sored	P Cer	erce 1sore	nt ed		
1668	717		951		57.0	)1		
Summary of the Weighted Number of Event and Censored Values								
Ev	ent and	l Cei	nsored	Val	ues	1 01		
Ev	ent and Ev	d Cei	Censo	Val red	ues Pe Cen	rcent		
Ev Total 89439590	ent and Ev 27650	d Cei rent 348	Censo 61789	Val red	ues Pe Cen	ercent sored 69.08		
Ev Total 89439590	ent and Ev 27650 Varia	ent 348	Censo 61789	Val red 242	Pe Cen	ercent sored 69.08		
Ev Total 89439590	ent and Ev 27650 Variat	ent 348	Censo 61789 Estimat	Val red 242	Pe Cen	ercent sored 69.08		

Output 115.3.2 displays the estimated regression coefficients and their standard errors. Poverty index has two levels, and only one level is estimable. By default, PROC SURVEYPHREG estimates the first level (POVARIND 1) and assigns a zero value for the second level. The estimated regression coefficient is 0.385 with a standard error of 0.078. The estimated hazard for the poverty areas is 1.47 times higher than the estimated hazard for the non-poverty areas. The degrees of freedom are equal to the number of PSUs (55) minus the number of strata (22).

Output 115.3.2	Inference	for the	Entire	Population
----------------	-----------	---------	--------	------------

Analysis of Maximum Likelihood Estimates								
Standard						Hazard		
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio		
POVARIND 1	33	0.384961	0.077586	4.96	<.0001	1.470		
POVARIND 2	33	0				1.000		

Output 115.3.3 shows that 813 observation units in the sample are male, and they account for over 42.6 million males in the base year survey population. Approximately half of these observation units in the sample are censored, and an estimated 64.5% observation units are censored for the male subpopulation.

Output 115.3.3 Summary Statistics for the Male Subpopulation

### The SURVEYPHREG Procedure

#### Domain Analysis for domain GENDER=1

Number of Observations Read	1891
Number of Observations Used	813
Sum of Weights Read	48887067
Sum of Weights Used	42629905

Summary of the Number of Event and Censored Values								
	Total	Event	Cei	nsored	P Cer	ercent nsored		
	813	404		409		50.31		
Summary of the Weighted Number of Event and Censored Values								
	Total	Εv	ent	Censo	red	Perc Censo	ent red	
426	529905	15126	321	27503	584	64	1.52	

Output 115.3.3 continued

Output 115.3.4 shows that the estimated regression coefficient for POVARIND 1 is 0.425 with a standard error of 0.157. The estimated hazard for the males in the poverty areas is 1.53 times higher than the estimated hazard for the males in the non-poverty areas. The degrees of freedom for the *t* significant test for the male subpopulation are equal to the total number of PSUs (55) minus the total number of strata (22).

Output 115.3.4 Inference for the Male Subpopulation

Analysis of Maximum Likelihood Estimates									
Standard									
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio			
POVARIND 1	33	0.424922	0.156583	2.71	0.0105	1.529			
POVARIND 2	33	0				1.000			

Output 115.3.5 displays some summary statistics for the female subpopulation. There are 855 observation units for females in the sample, and they represent over 46.8 million females in the base year survey population. Although 63.4% females in the sample are alive, an estimated 73.2% females in the subpopulation are alive.

Output 115.3.5 Summary Statistics for the Female Subpopulation

#### The SURVEYPHREG Procedure

Domain Analysis for domain GENDER=2

Nu	mber o	of Obse	rvat	tions R	ead	1891		
Nu	Number of Observations Used 855							
Su	Sum of Weights Read 5409160							
Su	m of W	/eights	Use	ed		46809685		
	Summary of the Number of Event and Censored Values							
					Р	ercent		
			~		~ .			
	Total	Event	Cer	nsored	Cer	nsored		
	Total 855	<b>Event</b> 313	Cer	n <b>sored</b> 542	Cer	63.39		
	Total 855 Summa Ev	Event 313 ary of th ent and	Cei le W	nsored 542 /eightee nsored	Cer d Nu Val	63.39 mber of ues		
	Total 855 Summa Ev Total	Event 313 ary of th ent and Ev	Cei le W l Ce	150red 542 Veightee nsored Censo	Cer d Nu Val	63.39 mber of ues Percent Censored		

Output 115.3.6 shows that the estimated proportional hazards regression coefficients for POVARIND for the females subpopulation (0.435) is higher than the estimated proportional hazards regression coefficients for POVARIND for the males subpopulation. The estimated hazard for the females in the poverty areas is 1.54 times higher than the estimated hazard for the females in the non-poverty areas. The degrees of freedom for the *t* significant test for the female subpopulation are equal to the total number of PSUs (55) minus the total number of strata (22).

Output 115.3.6 Inference for the Female Subpopulation

Analysis of Maximum Likelihood Estimates								
Standard						Hazard		
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio		
POVARIND 1	33	0.434579	0.115766	3.75	0.0007	1.544		
POVARIND 2	33	0				1.000		

## **Example 115.4: Variance Estimation by Using Replicate Weights**

Consider the data set LibrarySurvey from "Getting Started: SURVEYPHREG Procedure" on page 9339. The selected sample contains 100 transactions from ten branch libraries. A set of replicate weights and jackknife coefficients are created by randomly assigning observation units in disjoint groups of nearly equal size within each stratum. A total of 46 different groups are created. The data set LibraryRepWeights is similar to the data set LibrarySurvey except that it also contains replicate weights repwt\_1 to repwt\_46. Each column of replicate weights is obtained by deleting one group of observations and adjusting the sampling weights for the other groups in that stratum (Rust 1985).

The data set LibraryJKCOEF contains the jackknife coefficient for every replicate sample. The variable replicate denotes the replicate number, donorstratum denotes the stratum identification for that replicate, and jkcoefficient denotes the jackknife coefficient for that replicate sample.

```
data LibrarySurvey;
   set LibrarySurvey;
   randomorder = ranuni(12345);
run:
proc sort data = LibrarySurvey out = LibrarySurvey;
   by Branch randomorder;
run:
data LibrarySurvey;
   set LibrarySurvey;
   array nGroup{10} (2 2 2 4 4 4 4 8 8 8);
   GroupPSU = mod(_N_, nGroup{Branch});
   drop randomorder nGroup1 nGroup2 nGroup3 nGroup4
        nGroup5 nGroup6 nGroup7 nGroup8 nGroup9 nGroup10;
run;
proc surveymeans data = LibrarySurvey varmethod = jk
                (outweights = LibraryRepWeights outjkcoefs = LibraryJKCOEF);
   weight SamplingWeight;
   strata Branch;
   cluster GroupPSU;
   var Age;
run;
```

It is not necessary to provide replicate weights to compute jackknife variance estimates using the SUR-VEYPHREG procedure. If you do not specify the replicate weights, then the procedure creates replicate weights for you. For this illustration, assume that LibraryRepWeights and LibraryJKCOEF are the only two data sets available for analysis.

The following SAS statements request a proportional hazards regression of lenBorrow on Age. The variable Returned is the censor indicator, and the value 0 indicates a censored observation. The WEIGHT statement specifies the sampling weight variable, and the REPWEIGHTS statement specifies replicate weight variables RepWt 1 to RepWt 46. The JKCOEFS= option in the REPWEIGHTS statement specifies the jackknife coefficient for each replicate sample. The VARMETHOD= option in the MODEL statement requests the jackknife variance estimation method. A STRATA statement is not required when the REPWEIGHTS statement is specified.

```
proc surveyphreg data = LibraryRepWeights varmethod = jk;
   weight SamplingWeight;
   repweights RepWt_: / jkcoefs = LibraryJKCOEF;
   model lenBorrow*Returned(0) = Age;
run;
```

Output 115.4.1 displays some summary information. The "Number of Observations," "Censored Summary," and "Weighted Censored Summary" tables are exactly the same as in the example discussed in "Getting Started: SURVEYPHREG Procedure" on page 9339. The "Variance Estimation" table displays information about the variance estimation, such as the name of the variance estimation method and the number of replicate samples.

Output 115.4.1 Summary Statistics for Overall Analysis

Number	of Observa	ations Re	ad 100
Number	of Observa	ations Us	<b>ed</b> 100
Sum of V	Veights Re	ad	11616.79
Sum of V	Veights Us	ed	11616.79
Sum	nary of the and Cens	e Number ored Valu	of Event les Percent
Total	Event Ce	nsored (	Censored
100 Summa	90 ary of the V	10 Veighted	10.00 Number of
100 Summa Ev	90 ary of the V rent and Co	10 Veighted ensored V	10.00 Number of /alues Percent
100 Summa Ev Total 11616.79	90 ary of the V rent and Ce Event 10449.22	10 Veighted ensored V Censore 1167.5	10.00 Number of /alues Percent d Censored
100 Summa Ev Total 11616.79	90 ary of the V rent and Co Event 10449.22 Variance	10 Veighted ensored V Censore 1167.5 Estimatio	10.00 Number of /alues Percent d Censored 7 10.05
100 Summa Ev Total 11616.79	90 ary of the V rent and Co Event 10449.22 Variance	10 Veighted ensored V Censore 1167.5 Estimatio	10.00 Number of /alues Percent of Censored 7 10.05
100 Summa Ev Total 11616.79 od cate Weig	90 ary of the V rent and Ce Event 10449.22 Variance hts W0	10 Veighted ensored V Censore 1167.5 Estimatio	10.00 Number of /alues Percent d Censored 7 10.05 Dn Ja XARYREPWE

### The SURVEYPHREG Procedure

Output 115.4.2 shows that the estimated regression coefficient is 0.0616 with a standard error of 0.009. The denominator degrees of freedom (46) for the *t* test is equal to the number of replicates used. Note that the estimated proportional hazards regression coefficient is the same as the estimated proportional hazards regression coefficient is the same as the estimated proportional hazards regression coefficient. This is not surprising because these two examples use the same estimator to estimate the regression coefficients but different estimators to estimate the variance.

Output 115.4.2 Inferences Based on Survey Design for Overall Analysis

Analysis of Maximum Likelihood Estimates								
Standard Ha								
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio		
Age	46	0.061593	0.009159	6.73	<.0001	1.064		

## Example 115.5: A Test of the Proportional Hazards Assumption by Using the Programming Statements

You can use programming statements in PROC SURVEYPHREG to create time-dependent covariates to test the proportional hazards assumption for complex survey data. Consider the data set mortality from Example 115.3. The data set contains 1,891 observations from the 1992 NHANES I Epidemiologic Followup study (NHEFS) vital and tracing status.

Suppose you want to fit a proportional hazards model to this data and construct a test for the proportional hazards assumption on gender. The following statements request a proportional hazards regression of age on gender and x, where the time-dependent covariate x is created using the programing statements. The explanatory variable x assumes the value of the time variable age for the male subgroup. The variable vitalstatus is the censor indicator, and a value of 1, 4, 5, or 6 indicates a censored observation. The WEIGHT statement specifies the sampling weight, and the CLASS statement specifies that gender is a classification variable.

```
proc surveyphreg data = mortality nomcar;
    class gender;
    strata varstrata;
    cluster varpsu;
    weight sweight;
    model age*vitalstatus(1 4 5 6) = gender x;
    x = age*(gender=1);
run;
```

Output 115.5.1 displays some summary information. The "Number of Observations," "Censored Summary," and "Weighted Censored Summary" tables are exactly the same as in the example discussed in "Example 115.3: Domain Analysis" on page 9406.

Output 115.5.1 Data Summary, Censored Summary, and Information about Variance Estimation

Numb	1891								
Numb	1891								
Sum o	1.0298E8								
Sum o	1.0298E8								
Summary of the Number of Event and Censored Values									
Percent Total Event Censored Censored									
	1891 717			62.00					
18	91	/1/	11/4	62.08					
	91 ma Eve	ry of the W ent and Ce	11/4 /eighted Nu nsored Val	imber of ues					
	91 ma Eve	ry of the W ent and Ce Event	/eighted Nu nsored Val Censored	imber of ues Percent Censored					
	91 ma Eve tal	ry of the W ent and Ce Event 27650348	1174 Veighted Nu nsored Val Censored 75328323	imber of ues Percent Censored 73.15					
18 Sum To 1.0298	91 ma Eve tal	ry of the W ent and Ce Event 27650348 Variance	11/4 Veighted Nu nsored Val Censored 75328323 Estimation	umber of ues Percent Censored 73.15					
_18 Sum To 1.0298	91 ma Eve tal E8	ry of the Went and Ce Event 27650348 Variance	11/4 Veighted Nu nsored Val 75328323 Estimation Taylor Se	b2.08 umber of ues Percent Censored 73.15 eries					

#### The SURVEYPHREG Procedure

Output 115.5.2 displays the estimated regression coefficients and their standard errors. The variable gender has two levels, and only one level is estimable. By default, PROC SURVEYPHREG estimates the first level (GENDER 1) and assigns a zero value for the second level. The estimated regression coefficient is 1.61 with a standard error of 0.71. The estimated regression coefficient for x is -0.02 with a standard error of 0.01. The *t* statistic for x is -1.55 with a *p*-value of 0.13 on 33 degrees of freedom. This test suggests that an interaction between the time variable age and gender is not significant. Therefore, there is little evidence of an exponential trend over time in the hazard ratio for gender.

Analysis of Maximum Likelihood Estimates									
Standard									
Parameter	DF	Estimate	Error	t Value	Pr >  t	Ratio			
GENDER 1	33	1.605505	0.709081	2.26	0.0303	4.980			
GENDER 2	33	0				1.000			
x	33	-0.015648	0.010079	-1.55	0.1301	0.984			

## References

- Binder, D. A. (1983). "On the Variances of Asymptotically Normal Estimators from Complex Surveys." *International Statistical Review* 51:279–292.
- Binder, D. A. (1990). "Fitting Cox's Proportional Hazards Models from Survey Data." In *Proceedings of the Survey Research Methods Section*, 342–347. Alexandria, VA: American Statistical Association.

- Binder, D. A. (1992). "Fitting Cox's Proportional Hazards Models from Survey Data." *Biometrika* 79:139–147.
- Boudreau, C., and Lawless, J. F. (2006). "Survival Analysis Based on the Proportional Hazards Model and Survey Data." *Canadian Journal of Statistics* 34:203–216.
- Breslow, N. E. (1974). "Covariance Analysis of Censored Survival Data." Biometrics 30:89–99.
- Brick, J. M., and Kalton, G. (1996). "Handling Missing Data in Survey Research." Statistical Methods in Medical Research 5:215–238.
- Chambers, R. L., and Skinner, C. J. (2003). Analysis of Survey Data. Chichester, UK: John Wiley & Sons.
- Cochran, W. G. (1977). Sampling Techniques. 3rd ed. New York: John Wiley & Sons.
- Cox, D. R. (1972). "Regression Models and Life Tables." *Journal of the Royal Statistical Society, Series B* 20:187–220. With discussion.
- Cox, D. R. (1975). "Partial Likelihood." Biometrika 62:269-276.
- Dippo, C. S., Fay, R. E., and Morganstein, D. H. (1984). "Computing Variances from Complex Samples with Replicate Weights." In *Proceedings of the Survey Research Methods Section*, 489–494. Alexandria, VA: American Statistical Association.
- Efron, B. (1977). "The Efficiency of Cox's Likelihood Function for Censored Data." *Journal of the American Statistical Association* 72:557–565.
- Fay, R. E. (1984). "Some Properties of Estimates of Variance Based on Replication Methods." In Proceedings of the Survey Research Methods Section, 495–500. Alexandria, VA: American Statistical Association.
- Fay, R. E. (1989). "Theory and Application of Replicate Weighting for Variance Calculations." In *Proceedings* of the Survey Research Methods Section, 212–217. Alexandria, VA: American Statistical Association.
- Fuller, W. A. (1975). "Regression Analysis for Sample Survey." Sankhyā, Series C 37:117–132.
- Fuller, W. A. (2009). Sampling Statistics. Hoboken, NJ: John Wiley & Sons.
- Fuller, W. A., Kennedy, W. J., Schnell, D., Sullivan, G., and Park, H. J. (1989). PC CARP. Ames: Iowa State University Statistical Laboratory.
- Godambe, V. P., and Thompson, M. E. (1986). "Parameters of Superpopulation and Survey Population: Their Relationships and Estimation." *International Statistical Review* 54:127–138.
- Harrell, F. E. (1986). "The PHGLM Procedure." In *SUGI Supplemental Library Guide, Version 5 Edition*. Cary, NC: SAS Institute Inc.
- Judkins, D. R. (1990). "Fay's Method for Variance Estimation." Journal of Official Statistics 6:223-239.
- Kalbfleisch, J. D., and Prentice, R. L. (2002). *The Statistical Analysis of Failure Time Data*. 2nd ed. Hoboken, NJ: John Wiley & Sons.
- Kalton, G., and Kasprzyk, D. (1986). "The Treatment of Missing Survey Data." *Survey Methodology* 12:1–16.
- Kish, L. (1965). Survey Sampling. New York: John Wiley & Sons.

- Kish, L., and Frankel, M. R. (1974). "Inference from Complex Samples." *Journal of the Royal Statistical Society, Series B* 36:1–37.
- Korn, E. L., and Graubard, B. I. (1999). Analysis of Health Surveys. New York: John Wiley & Sons.
- Lawless, J. F. (2003). *Statistical Model and Methods for Lifetime Data*. 2nd ed. New York: John Wiley & Sons.
- Lee, E. W., Wei, L. J., and Amato, D. A. (1992). "Cox-Type Regression Analysis for Large Numbers of Small Groups of Correlated Failure Time Observations." In *Survival Analysis: State of the Art*, edited by J. P. Klein and P. K. Goel, 237–247. Dordrecht, Netherlands: Kluwer Academic.
- Lin, D. Y. (2000). "On Fitting Cox's Proportional Hazards Models to Survey Data." Biometrika 87:37-47.
- Lin, D. Y., and Wei, L. J. (1989). "The Robust Inference for the Proportional Hazards Model." *Journal of the American Statistical Association* 84:1074–1078.
- Lohr, S. L. (2010). Sampling: Design and Analysis. 2nd ed. Boston: Brooks/Cole.
- Pfeffermann, D. (1993). "The Role of Sampling Weights When Modeling Survey Data." *International Statistical Review* 61:317–337.
- Rao, J. N. K., Scott, A. J., and Skinner, C. J. (1998). "Quasi-score Tests with Survey Data." *Statistica Sinica* 8:1059–1070.
- Rao, J. N. K., and Shao, J. (1996). "On Balanced Half-Sample Variance Estimation in Stratified Random Sampling." *Journal of the American Statistical Association* 91:343–348.
- Rao, J. N. K., and Shao, J. (1999). "Modified Balanced Repeated Replication for Complex Survey Data." *Biometrika* 86:403–415.
- Rao, J. N. K., and Wu, C. F. J. (1988). "Resampling Inference with Complex Survey Data." *Journal of the American Statistical Association* 83:231–241.
- Rust, K. (1985). "Variance Estimation for Complex Estimators in Sample Surveys." *Journal of Official Statistics* 1:381–397.
- Särndal, C.-E., and Lundström, S. (2005). *Estimation in Surveys with Nonresponse*. Chichester, UK: John Wiley & Sons.
- Särndal, C.-E., Swensson, B., and Wretman, J. (1992). *Model Assisted Survey Sampling*. New York: Springer-Verlag.
- Schoenfeld, D. A. (1982). "Partial Residuals for the Proportional Hazards Regression Model." *Biometrika* 69:239–241.
- Therneau, T. M., Grambsch, P. M., and Fleming, T. R. (1990). "Martingale-Based Residuals and Survival Models." *Biometrika* 77:147–160.
- Wolter, K. M. (2007). Introduction to Variance Estimation. 2nd ed. New York: Springer.

## Subject Index

Akaike's information criterion SURVEYPHREG procedure, 9389 alpha level hazard ratio estimates (SURVEYPHREG), 9358 balanced repeated replication SURVEYPHREG procedure, 9381 variance estimation (SURVEYPHREG), 9381 bootstrap SURVEYPHREG procedure, 9381 variance estimation (SURVEYPHREG), 9381 bootstrap variance estimation SURVEYPHREG procedure, 9381 Breslow method likelihood (SURVEYPHREG), 9361 BRR SURVEYPHREG procedure, 9381 **BRR** variance estimation SURVEYPHREG procedure, 9381 censored survival times (SURVEYPHREG), 9369 censored values summary SURVEYPHREG procedure, 9395 censoring variable (SURVEYPHREG), 9363 CLASS variables programming statements (SURVEYPHREG), 9363 **CLUSTER** variables programming statements (SURVEYPHREG), 9363 clustering SURVEYPHREG procedure, 9353, 9375 covariance matrix SURVEYPHREG procedure, 9359 Cox regression analysis semiparametric model (SURVEYPHREG), 9338 degrees of freedom SURVEYPHREG procedure, 9386 deviance residuals SURVEYPHREG procedure, 9363, 9392 domain analysis SURVEYPHREG procedure, 9388 **DOMAIN** variables programming statements (SURVEYPHREG), 9363 domains

donor stratum SURVEYPHREG procedure, 9383 Efron method likelihood (SURVEYPHREG), 9361 event values summary SURVEYPHREG procedure, 9395 Fay coefficient SURVEYPHREG procedure, 9382 Fay's BRR method variance estimation (SURVEYPHREG), 9382 finite population correction SURVEYPHREG procedure, 9345 frequency variable programming statements (SURVEYPHREG), 9363 value (SURVEYPHREG), 9355 global null hypothesis SURVEYPHREG procedure, 9396 Hadamard matrix BRR variance estimation (SURVEYPHREG), 9383 hazard function baseline (SURVEYPHREG), 9369 hazard ratios Wald's confidence limits (SURVEYPHREG), 9360 Hessian matrix SURVEYPHREG procedure, 9360 inverse Hessian matrix SURVEYPHREG procedure, 9360 jackknife SURVEYPHREG procedure, 9383 jackknife coefficients SURVEYPHREG procedure, 9383 jackknife variance estimation SURVEYPHREG procedure, 9383 Lee-Wei-Amato model SURVEYPHREG procedure, 9399 likelihood ratio test SURVEYPHREG procedure, 9389, 9396 linear predictor SURVEYPHREG procedure, 9362, 9363

SURVEYPHREG procedure, 9353

linearization method SURVEYPHREG procedure, 9380 local influence score residuals (SURVEYPHREG), 9363, 9392 martingale residuals SURVEYPHREG procedure, 9363 missing values SURVEYPHREG procedure, 9364, 9377 model fit criteria (SURVEYPHREG), 9389 model information SURVEYPHREG procedure, 9395 number of observations SURVEYPHREG procedure, 9395 number of replicates SURVEYPHREG procedure, 9381-9383, 9385 number of subjects at risk SURVEYPHREG procedure, 9363 options summary ESTIMATE statement, 9354 parameter estimates SURVEYPHREG procedure, 9396 partial likelihood SURVEYPHREG procedure, 9369, 9374 primary sampling units (PSUs) SURVEYPHREG procedure, 9353 programming statements SURVEYPHREG procedure, 9363, 9365 proportional hazards model SURVEYPHREG procedure, 9338 replicate coefficients SURVEYPHREG procedure, 9385 replicate weights SURVEYPHREG procedure, 9365, 9379, 9385 replicate weights variance estimation SURVEYPHREG procedure, 9385 replication methods SURVEYPHREG procedure, 9379 replication variance estimation SURVEYPHREG procedure, 9385 residuals deviance (SURVEYPHREG), 9363, 9392 martingale (SURVEYPHREG), 9363 Schoenfeld (SURVEYPHREG), 9363, 9392, 9393 score (SURVEYPHREG), 9363, 9392 response variable SURVEYPHREG procedure, 9363 risk set SURVEYPHREG procedure, 9374

sample design SURVEYPHREG procedure, 9375 sampling rates SURVEYPHREG procedure, 9345, 9376 sampling weights SURVEYPHREG procedure, 9368, 9376 Schoenfeld residuals SURVEYPHREG procedure, 9363, 9392, 9393 score residuals SURVEYPHREG procedure, 9363, 9392 semiparametric model SURVEYPHREG procedure, 9338 singularity criterion SURVEYPHREG procedure, 9361 standard error SURVEYPHREG procedure, 9362, 9363, 9396 standard error ratio SURVEYPHREG procedure, 9387 STRATA variables programming statements (SURVEYPHREG), 9363 stratification SURVEYPHREG procedure, 9367, 9375 subdomain analysis, see domain analysis subgroup analysis, see domain analysis subpopulation analysis, see domain analysis survey data analysis SURVEYPHREG procedure, 9338 survey sampling data analysis (SURVEYPHREG), 9338 SURVEYPHREG procedure, 9338 Akaike's information criterion, 9389 alpha level, 9358 balanced repeated replication, 9381 bootstrap, 9381 bootstrap variance estimation, 9381 Breslow likelihood, 9361 BRR, 9381 BRR variance estimation, 9381 censored values summary, 9395 clustering, 9353, 9375 continuous time scale, 9361 covariance matrix, 9359 Cox regression analysis, 9338 DATA step statements, 9363 degrees of freedom, 9386 design summary table, 9395 displayed output, 9395 domain analysis, 9388 domain variable, 9353 domains, 9353 donor stratum, 9383 Efron likelihood, 9361 event values summary, 9395

Fay coefficient, 9382 Fay's BRR variance estimation, 9382 finite population correction, 9345 global null hypothesis, 9396 Hadamard matrix (BRR variance estimation), 9383 hazard ratio confidence intervals, 9358, 9360 Hessian matrix, 9360 hypothesis tests and confidence intervals, 9389 inverse Hessian matrix, 9360 jackknife, 9383 jackknife coefficients, 9383 jackknife variance estimation, 9383 Lee-Wei-Amato model, 9399 likelihood ratio test, 9389, 9396 linear predictor, 9362, 9363 linearization method, 9380 missing values, 9364, 9377 model fit statistics, 9389 model information, 9395 number of observations, 9395 number of replicates, 9381-9383, 9385 number of subjects at risk, 9363 ODS graph names, 9399 ODS graphics, 9399 ODS table names, 9398 ordering of effects, 9345 output data sets, 9393 OUTPUT statistics, 9362, 9363 parameter estimates, 9396 parameter estimates confidence intervals, 9359 partial likelihood, 9369, 9374 population totals, 9346, 9376 primary sampling units (PSUs), 9353 programming statements, 9363, 9365 proportional hazards model, 9338 replicate coefficients, 9385 replicate weights, 9365, 9379, 9385 replicate weights variance estimation, 9385 replication methods, 9379 replication variance estimation, 9385 residuals, 9363, 9391-9393 risk set, 9374 sample design, 9375 sampling rates, 9345, 9376 sampling weights, 9368, 9376 singularity criterion, 9361 standard error, 9362, 9363, 9396 standard error ratio, 9387 stratification, 9367, 9375 survival distribution function, 9370 survival times, 9369 survivor function, 9369, 9370 Taylor series linearized variance estimation, 9350

Taylor series variance estimation, 9380 ties, 9361, 9395 time-dependent covariates, 9339, 9363 variance adjustment, 9386 variance estimation, 9379 variance ratio, 9387 Wald test, 9390, 9396 weighting, 9368, 9376 survival distribution function SURVEYPHREG procedure, 9370 survival times SURVEYPHREG procedure, 9369 survivor function definition (SURVEYPHREG), 9370 SURVEYPHREG procedure, 9369 Taylor series linearized variance estimation SURVEYPHREG procedure, 9350 Taylor series variance estimation SURVEYPHREG procedure, 9380 ties SURVEYPHREG procedure, 9361, 9395 time-dependent covariates SURVEYPHREG procedure, 9339, 9363 variance adjustment SURVEYPHREG procedure, 9386 variance estimation bootstrap (SURVEYPHREG), 9381 BRR (SURVEYPHREG), 9381 jackknife (SURVEYPHREG), 9383 replicate weights (SURVEYPHREG), 9385 SURVEYPHREG procedure, 9379 Taylor series (SURVEYPHREG), 9350, 9380 variance ratio SURVEYPHREG procedure, 9387 Wald test SURVEYPHREG procedure, 9390, 9396 WEIGHT variable programming statements (SURVEYPHREG), 9363 weighting SURVEYPHREG procedure, 9368, 9376

# Syntax Index

ALPHA= option MODEL statement (SURVEYPHREG), 9358 BY statement SURVEYPHREG procedure, 9350 CENTER= option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9347 VARMETHOD=JK (PROC SURVEYPHREG statement), 9349 CLASS statement SURVEYPHREG procedure, 9350 **CLPARM** option MODEL statement (SURVEYPHREG), 9359 **CLUSTER** statement SURVEYPHREG procedure, 9353 COVB option MODEL statement (SURVEYPHREG), 9359 DATA= option PROC SURVEYPHREG statement, 9344 **DESCENDING** option CLASS statement (SURVEYPHREG), 9351 **DETAILS** option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9347 VARMETHOD=JK (PROC SURVEYPHREG statement), 9349 DF= option MODEL statement (SURVEYPHREG), 9359 DF=ALLREPS DF= (SURVEYPHREG), 9359 DF=DESIGN DF= (SURVEYPHREG), 9359 DF=DESIGN (value) DF= (SURVEYPHREG), 9359 DF=DESIGNADJ DF= (SURVEYPHREG), 9359 DF=NONE DF= (SURVEYPHREG), 9360 DF=PARMADJ DF= (SURVEYPHREG), 9360 DF=PARMADJ (value) DF= (SURVEYPHREG), 9360 DOMAIN statement SURVEYPHREG procedure, 9353 ESTIMATE statement

SURVEYPHREG procedure, 9354

FAY= option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9347 FREO statement SURVEYPHREG procedure, 9355 HADAMARD= option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9347 **HESS** option MODEL statement (SURVEYPHREG), 9360 **INVHESS** option MODEL statement (SURVEYPHREG), 9360 JKCOEFS= option **REPWEIGHTS** statement (SURVEYPHREG), 9365 keyword= option OUTPUT statement (SURVEYPHREG), 9362 LIST option STRATA statement (SURVEYPHREG), 9368 LSMEANS statement SURVEYPHREG procedure, 9355 LSMESTIMATE statement SURVEYPHREG procedure, 9356 MISSING option CLASS statement (SURVEYPHREG), 9351 PROC SURVEYPHREG statement, 9344 MODEL statement SURVEYPHREG procedure, 9358 NLOPTIONS statement SURVEYPHREG procedure, 9362 NOMCAR option PROC SURVEYPHREG statement, 9344 NOPRINT option PROC SURVEYPHREG statement, 9344 ORDER= option CLASS statement (SURVEYPHREG), 9351 PROC SURVEYPHREG statement, 9345 OUT= option OUTPUT statement (SURVEYPHREG), 9362 OUTJKCOEFS= option

VARMETHOD=JK (PROC SURVEYPHREG statement), 9350 **OUTPUT** statement SURVEYPHREG procedure, 9362 OUTWEIGHTS= option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9348 VARMETHOD=JK (PROC SURVEYPHREG statement), 9349 PARAM= option CLASS statement (SURVEYPHREG), 9351 **PRINTH** option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9348 PROC SURVEYPHREG statement, see SURVEYPHREG procedure RATE= option PROC SURVEYPHREG statement, 9345 REF= option CLASS statement (SURVEYPHREG), 9352 **REPCOEFS**= option **REPWEIGHTS** statement (SURVEYPHREG), 9366 **REPS**= option VARMETHOD=BRR (PROC SURVEYPHREG statement), 9348 **REPWEIGHTS** statement SURVEYPHREG procedure, 9365 **RISKLIMITS**= option MODEL statement (SURVEYPHREG), 9360 SERATIO= option MODEL statement (SURVEYPHREG), 9360 SINGULAR= option MODEL statement (SURVEYPHREG), 9361 SLICE statement SURVEYPHREG procedure, 9367 SRVEYPHREG procedure, PROC SURVEYPHREG statement DATA= option, 9344 MISSING option, 9344 STORE statement SURVEYPHREG procedure, 9367 STRATA statement SURVEYPHREG procedure, 9367 SURVEYPHREG procedure DF=ALLREPS, 9359 DF=DESIGN, 9359 DF=DESIGN (value), 9359 DF=DESIGNADJ, 9359 **DF=NONE**, 9360 DF=PARMADJ, 9360 DF=PARMADJ (value), 9360

NLOPTIONS statement, 9362 SURVEYPHREG procedure, BY statement, 9350 SURVEYPHREG procedure, CLASS statement, 9350 **DESCENDING** option, 9351 MISSING option, 9351 ORDER= option, 9351 PARAM= option, 9351 REF= option, 9352 TRUNCATE option, 9352 SURVEYPHREG procedure, CLUSTER statement, 9353 SURVEYPHREG procedure, DOMAIN statement, 9353 SURVEYPHREG procedure, ESTIMATE statement, 9354 SURVEYPHREG procedure, FREQ statement, 9355 SURVEYPHREG procedure, LSMEANS statement, 9355 SURVEYPHREG procedure, LSMESTIMATE statement, 9356 SURVEYPHREG procedure, MODEL statement, 9358 ALPHA= option, 9358 CLPARM option, 9359 COVB option, 9359 DF= option, 9359 HESS option, 9360 **INVHESS** option, 9360 RISKLIMITS= option, 9360 SERATIO= option, 9360 SINGULAR= option, 9361 TIES= option, 9361 VADJUST= option, 9361 VARRATIO= option, 9361 SURVEYPHREG procedure, NLOPTIONS statement, 9362 SURVEYPHREG procedure, OUTPUT statement, 9362 keyword= option, 9362 OUT= option, 9362 SURVEYPHREG procedure, PROC SURVEYPHREG statement, 9344 CENTER= option (VARMETHOD=BRR), 9347 CENTER= option (VARMETHOD=JK), 9349 DETAILS option (VARMETHOD=BRR), 9347 DETAILS option (VARMETHOD=JK), 9349 FAY= option (VARMETHOD=BRR), 9347 HADAMARD= option (VARMETHOD=BRR), 9347 NOMCAR option, 9344 NOPRINT option, 9344 ORDER= option, 9345 OUTJKCOEFS= option (VARMETHOD=JK), 9350

OUTWEIGHTS= option (VARMETHOD=BRR), 9348 OUTWEIGHTS= option (VARMETHOD=JK), 9349 PRINTH option (VARMETHOD=BRR), 9348 RATE= option, 9345REPS= option (VARMETHOD=BRR), 9348 TOTAL= option, 9346 VARMETHOD= option, 9346 SURVEYPHREG procedure, REPWEIGHTS statement, 9365 JKCOEFS= option, 9365 REPCOEFS= option, 9366 SURVEYPHREG procedure, SLICE statement, 9367 SURVEYPHREG procedure, STORE statement, 9367 SURVEYPHREG procedure, STRATA statement, 9367 LIST option, 9368 SURVEYPHREG procedure, TEST statement, 9368 SURVEYPHREG procedure, WEIGHT statement, 9368 TEST statement SURVEYPHREG procedure, 9368 TIES= option MODEL statement (SURVEYPHREG), 9361 TOTAL= option PROC SURVEYPHREG statement, 9346 **TRUNCATE** option CLASS statement (SURVEYPHREG), 9352 VADJUST= option MODEL statement (SURVEYPHREG), 9361 VARMETHOD= option PROC SURVEYPHREG statement, 9346 VARRATIO= option MODEL statement (SURVEYPHREG), 9361 WEIGHT statement SURVEYPHREG procedure, 9368