

# SAS/STAT® 9.2 User's Guide The SURVEYREG Procedure (Book Excerpt)



This document is an individual chapter from SAS/STAT® 9.2 User's Guide.

The correct bibliographic citation for the complete manual is as follows: SAS Institute Inc. 2008. SAS/STAT® 9.2 User's Guide. Cary, NC: SAS Institute Inc.

Copyright © 2008, SAS Institute Inc., Cary, NC, USA

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

**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.

**U.S. Government Restricted Rights Notice**: Use, duplication, or disclosure of this software and related documentation by the U.S. government is subject to the Agreement with SAS Institute and the restrictions set forth in FAR 52.227-19, Commercial Computer Software-Restricted Rights (June 1987).

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

1st electronic book, March 2008 2nd electronic book, February 2009

SAS<sup>®</sup> Publishing provides a complete selection of books and electronic products to help customers use SAS software to its fullest potential. For more information about our e-books, e-learning products, CDs, and hard-copy books, visit the SAS Publishing Web site at **support.sas.com/publishing** or call 1-800-727-3228.

 $SAS^{\textcircled{@}}$  and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries. @ indicates USA registration.

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

# Chapter 86

# The SURVEYREG Procedure

Overview: SURVEYREG	Procedure	652
Getting Started: SURVEY	REG Procedure	652
Simple Random Sar	npling	65
Stratified Sampling		65
Output Data Sets .		65
Syntax: SURVEYREG Pro	ocedure	65
PROC SURVEYRE	G Statement	65
BY Statement		65
<b>CLASS Statement</b>		65
CLUSTER Statemen	nt	65
CONTRAST Statem	nent	65
DOMAIN Statemen	t	65
ESTIMATE Stateme	ent	65
MODEL Statement		65
OUTPUT Statement	t	65
REPWEIGHTS Stat	rement	65
STRATA Statement		65
WEIGHT Statement	t	65
Details: SURVEYREG Pro	ocedure	65
Missing Values		65
Survey Design Infor	rmation	65
Specification	on of Population Totals and Sampling Rates	65
Primary Sa	mpling Units (PSUs)	65
Computational Deta	ils	65
Notation .		65
Regression	Coefficients	65
Variance Es	stimation	65
Hadamard 1	Matrix	65
Degrees of	Freedom	65
Testing Effe	ects	65
Design Effe	ect	65
Stratum Co	llapse	65
Sampling R	tate of the Pooled Stratum from Collapse	65

Computational Resources         6565           Output Data Sets         6566           OUT= Data Set Created by the OUTPUT Statement         6566           Replicate Weights Output Data Set         6567           Jackknife Coefficients Output Data Set         6567           Displayed Output         6568           Design Summary         6568           Design Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6571           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6572           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Coefficients of Estimate         6572           Analysis of Contrast         6572           Coefficients of Estimate         6572           Analysis of Estimate         6572           Analysis of Estimate         6572           Hadamard Matrix         6573 <th>Domain Analysis</th> <th>6565</th>	Domain Analysis	6565
OUT= Data Set Created by the OUTPUT Statement         6566           Replicate Weights Output Data Set         6567           Jackknife Coefficients Output Data Set         6567           Displayed Output         6568           Data Summary         6568           Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6571           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           Examples SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6577 <t< td=""><td>Computational Resources</td><td>6565</td></t<>	Computational Resources	6565
Replicate Weights Output Data Set         6567           Jackknife Coefficients Output Data Set         6567           Displayed Output         6568           Data Summary         6568           Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6570           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Estimate         6572           Analysis of Estimate         6572           Analysis of Estimable Functions         6572           Analysis of Estimable Functions         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6574           Example 8	Output Data Sets	6566
Jackknife Coefficients Output Data Set         6568           Displayed Output         6568           Data Summary         6568           Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6570           Class Level Information         6570           X'X Matrix         6570           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6572           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6574           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.6: Stratum Collapse         6590           <	OUT= Data Set Created by the OUTPUT Statement	6566
Displayed Output         6568           Data Summary         6568           Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6571           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Coefficients of Estimate         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6574           Example 86.4: Stratified Sampling         6580           Example 86.5: Regression Estimator for Stratifie	Replicate Weights Output Data Set	6567
Data Summary         6568           Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6571           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6577           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.6: Stratum Collapse         6590           Exa	Jackknife Coefficients Output Data Set	6567
Design Summary         6568           Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6571           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.4: Stratified Sampling         6581           Example 86.6: Stratum Collapse         6590           Example 86.7: Domain Analysis         6595           <	Displayed Output	6568
Domain Summary         6569           Fit Statistics         6569           Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6570           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6577           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.4: Stratified Sampling         6581           Example 86.6: Stratum Collapse         6590           Example 86.8: Variance Estimate Using the Jackknife Met	Data Summary	6568
Fit Statistics       6569         Variance Estimation       6569         Stratum Information       6570         Class Level Information       6570         X'X Matrix       6570         Inverse Matrix of X'X       6571         ANOVA for Dependent Variable       6571         Tests of Model Effects       6571         Estimated Regression Coefficients       6571         Covariance of Estimated Regression Coefficients       6572         Coefficients of Contrast       6572         Analysis of Contrasts       6572         Coefficients of Estimate       6572         Analysis of Estimable Functions       6572         Hadamard Matrix       6573         ODS Table Names       6573         Examples: SURVEYREG Procedure       6574         Example 86.1: Simple Random Sampling       6574         Example 86.2: Simple Random Cluster Sampling       6577         Example 86.3: Regression Estimator for Simple Random Sample       6580         Example 86.4: Stratified Sampling       6581         Example 86.6: Stratum Collapse       6590         Example 86.7: Domain Analysis       6595         Example 86.8: Variance Estimate Using the Jackknife Method       6598	Design Summary	6568
Variance Estimation         6569           Stratum Information         6570           Class Level Information         6570           X'X Matrix         6570           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.2: Simple Random Cluster Sampling         6574           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.5: Regression Estimator for Stratified Sample         6581           Example 86.6: Stratum Collapse         6590           Example 86.8: Variance Estimate Using the Jackknife Method         6598	Domain Summary	6569
Stratum Information         6570           Class Level Information         6570           X'X Matrix         6570           Inverse Matrix of X'X         6571           ANOVA for Dependent Variable         6571           Tests of Model Effects         6571           Estimated Regression Coefficients         6571           Covariance of Estimated Regression Coefficients         6572           Coefficients of Contrast         6572           Analysis of Contrasts         6572           Coefficients of Estimate         6572           Analysis of Estimable Functions         6572           Hadamard Matrix         6573           ODS Table Names         6573           Examples: SURVEYREG Procedure         6574           Example 86.1: Simple Random Sampling         6574           Example 86.3: Regression Estimator for Simple Random Sample         6580           Example 86.4: Stratified Sampling         6581           Example 86.5: Regression Estimator for Stratified Sample         6587           Example 86.6: Stratum Collapse         6590           Example 86.8: Variance Estimate Using the Jackknife Method         6598	Fit Statistics	6569
Class Level Information       6570         X'X Matrix       6570         Inverse Matrix of X'X       6571         ANOVA for Dependent Variable       6571         Tests of Model Effects       6571         Estimated Regression Coefficients       6571         Covariance of Estimated Regression Coefficients       6572         Coefficients of Contrast       6572         Analysis of Contrasts       6572         Coefficients of Estimate       6572         Analysis of Estimable Functions       6572         Hadamard Matrix       6573         ODS Table Names       6573         Examples: SURVEYREG Procedure       6574         Example 86.1: Simple Random Sampling       6574         Example 86.2: Simple Random Cluster Sampling       6577         Example 86.3: Regression Estimator for Simple Random Sample       6580         Example 86.4: Stratified Sampling       6581         Example 86.5: Regression Estimator for Stratified Sample       6581         Example 86.6: Stratum Collapse       6590         Example 86.7: Domain Analysis       6595         Example 86.8: Variance Estimate Using the Jackknife Method       6598	Variance Estimation	6569
X'X Matrix       6570         Inverse Matrix of X'X       6571         ANOVA for Dependent Variable       6571         Tests of Model Effects       6571         Estimated Regression Coefficients       6571         Covariance of Estimated Regression Coefficients       6572         Coefficients of Contrast       6572         Analysis of Contrasts       6572         Coefficients of Estimate       6572         Analysis of Estimable Functions       6572         Hadamard Matrix       6573         ODS Table Names       6573         Examples: SURVEYREG Procedure       6574         Example 86.1: Simple Random Sampling       6574         Example 86.2: Simple Random Cluster Sampling       6577         Example 86.3: Regression Estimator for Simple Random Sample       6580         Example 86.4: Stratified Sampling       6581         Example 86.5: Regression Estimator for Stratified Sample       6587         Example 86.6: Stratum Collapse       6590         Example 86.7: Domain Analysis       6595         Example 86.8: Variance Estimate Using the Jackknife Method       6598	Stratum Information	6570
Inverse Matrix of X'X ANOVA for Dependent Variable 6571 Tests of Model Effects 6571 Estimated Regression Coefficients 6571 Covariance of Estimated Regression Coefficients 6572 Coefficients of Contrast 6572 Analysis of Contrast 6572 Coefficients of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimable Functions 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.5: Regression Estimator for Stratified Sample 6581 Example 86.6: Stratum Collapse Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	Class Level Information	6570
ANOVA for Dependent Variable 6571 Tests of Model Effects 6571 Estimated Regression Coefficients 6571 Covariance of Estimated Regression Coefficients 6572 Coefficients of Contrast 6572 Analysis of Contrast 6572 Analysis of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimable Functions 6572 Hadamard Matrix 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.4: Stratified Sampling 6581 Example 86.5: Regression Estimator for Stratified Sample 6587 Example 86.6: Stratum Collapse 6590 Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	X'X Matrix	6570
Tests of Model Effects 6571 Estimated Regression Coefficients 6571 Covariance of Estimated Regression Coefficients 6572 Coefficients of Contrast 6572 Analysis of Contrast 6572 Coefficients of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimable Functions 6572 Hadamard Matrix 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.4: Stratified Sampling 6581 Example 86.5: Regression Estimator for Stratified Sample 6587 Example 86.6: Stratum Collapse 6590 Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	Inverse Matrix of X'X	6571
Estimated Regression Coefficients 6571 Covariance of Estimated Regression Coefficients 6572 Coefficients of Contrast 6572 Analysis of Contrasts 6572 Coefficients of Estimate 6572 Analysis of Estimate 6572 Analysis of Estimable Functions 6572 Hadamard Matrix 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.4: Stratified Sampling 6581 Example 86.5: Regression Estimator for Stratified Sample 6587 Example 86.6: Stratum Collapse 6590 Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	ANOVA for Dependent Variable	6571
Covariance of Estimated Regression Coefficients6572Coefficients of Contrast6572Analysis of Contrasts6572Coefficients of Estimate6572Analysis of Estimable Functions6572Hadamard Matrix6573ODS Table Names6573Examples: SURVEYREG Procedure6574Example 86.1: Simple Random Sampling6574Example 86.2: Simple Random Cluster Sampling6577Example 86.3: Regression Estimator for Simple Random Sample6580Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Tests of Model Effects	6571
Coefficients of Contrast6572Analysis of Contrasts6572Coefficients of Estimate6572Analysis of Estimable Functions6572Hadamard Matrix6573ODS Table Names6573Examples: SURVEYREG Procedure6574Example 86.1: Simple Random Sampling6574Example 86.2: Simple Random Cluster Sampling6577Example 86.3: Regression Estimator for Simple Random Sample6580Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Estimated Regression Coefficients	6571
Analysis of Contrasts 6572 Coefficients of Estimate 6572 Analysis of Estimable Functions 6572 Hadamard Matrix 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.4: Stratified Sampling 6581 Example 86.5: Regression Estimator for Stratified Sample 6587 Example 86.6: Stratum Collapse 6590 Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	Covariance of Estimated Regression Coefficients	6572
Coefficients of Estimate 6572 Analysis of Estimable Functions 6572 Hadamard Matrix 6573 ODS Table Names 6573 Examples: SURVEYREG Procedure 6574 Example 86.1: Simple Random Sampling 6574 Example 86.2: Simple Random Cluster Sampling 6577 Example 86.3: Regression Estimator for Simple Random Sample 6580 Example 86.4: Stratified Sampling 6581 Example 86.5: Regression Estimator for Stratified Sample 6587 Example 86.6: Stratum Collapse 6590 Example 86.7: Domain Analysis 6595 Example 86.8: Variance Estimate Using the Jackknife Method 6598	Coefficients of Contrast	6572
Analysis of Estimable Functions  Hadamard Matrix  6573  ODS Table Names  6573  Examples: SURVEYREG Procedure  6574  Example 86.1: Simple Random Sampling  6574  Example 86.2: Simple Random Cluster Sampling  6577  Example 86.3: Regression Estimator for Simple Random Sample  6580  Example 86.4: Stratified Sampling  6581  Example 86.5: Regression Estimator for Stratified Sample  6587  Example 86.6: Stratum Collapse  6590  Example 86.7: Domain Analysis  6598  Example 86.8: Variance Estimate Using the Jackknife Method  6598	Analysis of Contrasts	6572
Hadamard Matrix	Coefficients of Estimate	6572
ODS Table Names 6573  Examples: SURVEYREG Procedure 6574  Example 86.1: Simple Random Sampling 6574  Example 86.2: Simple Random Cluster Sampling 6577  Example 86.3: Regression Estimator for Simple Random Sample 6580  Example 86.4: Stratified Sampling 6581  Example 86.5: Regression Estimator for Stratified Sample 6587  Example 86.6: Stratum Collapse 6590  Example 86.7: Domain Analysis 6595  Example 86.8: Variance Estimate Using the Jackknife Method 6598	Analysis of Estimable Functions	6572
Examples: SURVEYREG Procedure6574Example 86.1: Simple Random Sampling6574Example 86.2: Simple Random Cluster Sampling6577Example 86.3: Regression Estimator for Simple Random Sample6580Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Hadamard Matrix	6573
Example 86.1: Simple Random Sampling	ODS Table Names	6573
Example 86.2: Simple Random Cluster Sampling6577Example 86.3: Regression Estimator for Simple Random Sample6580Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Examples: SURVEYREG Procedure	6574
Example 86.3: Regression Estimator for Simple Random Sample6580Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Example 86.1: Simple Random Sampling	6574
Example 86.4: Stratified Sampling6581Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Example 86.2: Simple Random Cluster Sampling	6577
Example 86.5: Regression Estimator for Stratified Sample6587Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Example 86.3: Regression Estimator for Simple Random Sample	6580
Example 86.6: Stratum Collapse6590Example 86.7: Domain Analysis6595Example 86.8: Variance Estimate Using the Jackknife Method6598	Example 86.4: Stratified Sampling	6581
Example 86.7: Domain Analysis	Example 86.5: Regression Estimator for Stratified Sample	6587
Example 86.8: Variance Estimate Using the Jackknife Method 6598	Example 86.6: Stratum Collapse	6590
	Example 86.7: Domain Analysis	6595
References	Example 86.8: Variance Estimate Using the Jackknife Method	6598
	References	6602

#### **Overview: SURVEYREG Procedure**

The SURVEYREG procedure performs regression analysis for sample survey data. This procedure can handle complex survey sample designs, including designs with stratification, clustering, and unequal weighting. The procedure fits linear models for survey data and computes regression coefficients and their variance-covariance matrix. The procedure also provides significance tests for the model effects and for any specified estimable linear functions of the model parameters. Using the regression model, the procedure can compute predicted values for the sample survey data.

PROC SURVEYREG computes the regression coefficient estimators by generalized least squares estimation using elementwise regression. The procedure assumes that the regression coefficients are the same across strata and primary sampling units (PSUs). To estimate the variance-covariance matrix for the regression coefficients, PROC SURVEYREG uses either the Taylor series (linearization) method or replication (resampling) methods to estimate sampling errors of estimators based on complex sample designs. For details see Woodruff (1971); Fuller (1975); Särndal, Swensson, and Wretman (1992); Wolter (1985); Rust (1985); Dippo, Fay, and Morganstein (1984); Rao and Shao (1999); Rao, Wu, and Yue (1992); and Rao and Shao (1996).

# **Getting Started: SURVEYREG Procedure**

This section demonstrates how you can use PROC SURVEYREG to perform a regression analysis for sample survey data. For a complete description of the usage of PROC SURVEYREG, see the section "Syntax: SURVEYREG Procedure" on page 6534. The section "Examples: SURVEYREG Procedure" on page 6574 provides more detailed examples that illustrate the applications of PROC SURVEYREG.

# **Simple Random Sampling**

Suppose that, in a junior high school, there are a total of 4,000 students in grades 7, 8, and 9. You want to know how household income and the number of children in a household affect students' average weekly spending for ice cream.

In order to answer this question, you draw a sample by using simple random sampling from the student population in the junior high school. You randomly select 40 students and ask them their average weekly expenditure for ice cream, their household income, and the number of children in their household. The answers from the 40 students are saved as the SAS data set IceCream:

```
data IceCream;
   input Grade Spending Income Kids @@;
   datalines;
7
       39
           2
                7
                                   12
                                       47
                    7
                       38
                           1
                                8
                                           1
9
   10
       47
                7
                    1
                                7
                                   10
                                       43
                                            2
7
       44
                8
                   20
                           3
    3
           4
                       60
                                8
                                   19
                                       57
                                            4
7
    2
       35
           2
                7
                    2
                       36
                           1
                                9
                                   15
                                       51
                                            1
8
   16
       53
           1
                7
                    6
                       37
                           4
                                7
                                    6
                                       41
                                            2
7
       39
           2
                  15
                       50
                                   17
                                       57
                                            3
    6
                           4
           2
8
   14
       46
                9
                    8
                       41
                           2
                                9
                                    8
                                       41
                                           1
9
       47
           3
                7
                    3
                       39
                           3
                                7
    7
                                   12
                                       50
                                            2
7
    4
       43
           4
                9 14
                       46
                           3
                                8
                                   18
                                       58
                                            4
9
    9
       44
           3
                7
                   2
                       37
                           1
                                7
                                    1 37
                                           2
7
           2
                7
                           2
                                    8
                                       41
                                           2
    4
       44
                   11
                       42
                                9
                                           3
8
   10
       42
           2
                8 13
                       46
                           1
                                7
                                    2
                                       40
                                    2 36
9
       45
           1
                9 11
                       45
                                7
    6
7
    9
       46
           1
```

In the data set IceCream, the variable Grade indicates a student's grade. The variable Spending contains the dollar amount of each student's average weekly spending for ice cream. The variable Income specifies the household income, in thousands of dollars. The variable Kids indicates how many children are in a student's family.

The following PROC SURVEYREG statements request a regression analysis:

```
title1 'Ice Cream Spending Analysis';
title2 'Simple Random Sample Design';
proc surveyreg data=IceCream total=4000;
   class Kids;
   model Spending = Income Kids / solution;
run;
```

The PROC SURVEYREG statement invokes the procedure. The TOTAL=4000 option specifies the total in the population from which the sample is drawn. The CLASS statement requests that the procedure use the variable Kids as a classification variable in the analysis. The MODEL statement describes the linear model that you want to fit, with Spending as the dependent variable and Income and Kids as the independent variables. The SOLUTION option in the MODEL statement requests that the procedure output the regression coefficient estimates.

Figure 86.1 displays the summary of the data, the summary of the fit, and the levels of the classification variable Kids. The "Fit Statistics" table displays the denominator degrees of freedom, which are used in F tests and t tests in the regression analysis.

Figure 86.1 Summary of Data

Ice Cream Spending Analysis Simple Random Sample Design The SURVEYREG Procedure Regression Analysis for Dependent Variable Spending Data Summary Number of Observations 40
Mean of Spending 8.75000
Sum of Spending 350.00000 Number of Observations 40 Fit Statistics 0.8132 R-square Root MSE 2.4506 Denominator DF 39 Class Level Information Class Variable Levels Values Kids 1 2 3 4

Figure 86.2 displays the tests for model effects. The effect Income is significant in the linear regression model, while the effect Kids is not significant at the 5% level.

Figure 86.2 Testing Effects in the Regression

5	Tests of Mod	del Effects	
Effect	Num DF	F Value	Pr > F
Model	4	119.15	<.0001
Intercept	1	153.32	<.0001
Income	1	324.45	<.0001
Kids	3	0.92	0.4385
NOTE: The denominator	degrees of	freedom for	the F tests is 39.

The regression coefficient estimates and their standard errors and associated *t* tests are displayed in Figure 86.3.

Figure 86.3 Regression Coefficients

		Estimated Re	gression Coeff	icients			
			Standard				
	Parameter	Estimate	Error	t Value	Pr >  t		
	Intercept	-26.084677	2.46720403	-10.57	<.0001		
	Income	0.775330	0.04304415	18.01	<.0001		
	Kids 1	0.897655	1.12352876	0.80	0.4292		
	Kids 2	1.494032	1.24705263	1.20	0.2381		
	Kids 3	-0.513181	1.33454891	-0.38	0.7027		
	Kids 4	0.000000	0.0000000				
NOTE: The denominator degrees of freedom for the t tests is 39.  Matrix X'X is singular and a generalized inverse was used to solve the							
Ma							

#### **Stratified Sampling**

Suppose that the previous student sample is actually drawn using a stratified sample design. The strata are grades in the junior high school: 7, 8, and 9. Within strata, simple random samples are selected. Table 86.1 provides the number of students in each grade.

Table 86.1 Students in Grades

Grade	Number of Students
7	1,824
8	1,025
9	1,151
Total	4,000

In order to analyze this sample by using PROC SURVEYREG, you need to input the stratification information by creating a SAS data set with the information in Table 86.1. The following SAS statements create such a data set called StudentTotals:

```
data StudentTotals;
   input Grade _TOTAL_;
   datalines;
7 1824
8 1025
9 1151
;
```

The variable Grade is the stratification variable, and the variable \_TOTAL\_ contains the total numbers of students in each stratum in the survey population. PROC SURVEYREG requires you to use the keyword \_TOTAL\_ as the name of the variable that contains the population total information.

In a stratified sample design, when the sampling rates in the strata are unequal, you need to use sampling weights to reflect this information. For this example, the appropriate sampling weights are the reciprocals of the probabilities of selection. You can use the following DATA step to create the sampling weights:

```
data IceCream;
  set IceCream;
  if Grade=7 then Prob=20/1824;
  if Grade=8 then Prob=9/1025;
  if Grade=9 then Prob=11/1151;
  Weight=1/Prob;
```

If you use PROC SURVEYSELECT to select your sample, PROC SURVEYSELECT creates these sampling weights for you.

The following statements demonstrate how you can fit a linear model while incorporating the sample design information (stratification):

```
title1 'Ice Cream Spending Analysis';
title2 'Stratified Simple Random Sample Design';
proc surveyreg data=IceCream total=StudentTotals;
    strata Grade /list;
    class Kids;
    model Spending = Income Kids / solution;
    weight Weight;
run:
```

Comparing these statements to those in the section "Simple Random Sampling" on page 6527, you can see how the TOTAL=StudentTotals option replaces the previous TOTAL=4000 option.

The STRATA statement specifies the stratification variable Grade. The LIST option in the STRATA statement requests that the stratification information be included in the output. The WEIGHT statement specifies the weight variable.

Figure 86.4 summarizes the data information, the sample design information, and the fit information. Note that, due to the stratification, the denominator degrees of freedom for F tests and t tests are 37, which is different from the analysis in Figure 86.1.

Figure 86.4 Summary of the Regression

```
Ice Cream Spending Analysis
Stratified Simple Random Sample Design

The SURVEYREG Procedure

Regression Analysis for Dependent Variable Spending

Data Summary

Number of Observations 40
Sum of Weights 4000.0
Weighted Mean of Spending 9.14130
Weighted Sum of Spending 36565.2
```

Figure 86.4 continued

Design Summ	ary
Number of Strata	3
Fit Statist	ics
R-square	0.8219
Root MSE	2.4185
Denominator DF	37

For each stratum, Figure 86.5 displays the value of identifying variables, the number of observations (sample size), the total population size, and the calculated sampling rate or fraction.

Figure 86.5 Stratification and Classification Information

	Stratum	n Informat	ion	
Stratum			Population	Sampling
Index	Grade N	1 Obs	Total	Rate
1	7	20	1824	1.10%
2	8	9	1025	0.88%
3	9	11	1151	0.96%
	Class Lev	vel Inform	ation	
	Class			
	Variable	Levels	Values	
	Kids	4	1 2 3 4	

Figure 86.6 displays the tests for the significance of model effects under the stratified sample design. The Income effect is strongly significant, while the Kids effect is not significant at the 5% level.

Figure 86.6 Testing Effects

	Tests of Mod	lel Effects	
Effect	Num DF	F Value	Pr > F
Model	4	124.85	<.0001
Interd	cept 1	150.95	<.0001
Income	e 1	326.89	<.0001
Kids	3	0.99	0.4081
Kids	3	0.99	

The regression coefficient estimates for the stratified sample, along with their standard errors and associated *t* tests, are displayed in Figure 86.7.

Figure 86.7 Regression Coefficients

		Standard		
Parameter	Estimate	Error	t Value	Pr >  t
Intercept	-26.086882	2.44108058	-10.69	<.0001
Income	0.776699	0.04295904	18.08	<.0001
Kids 1	0.888631	1.07000634	0.83	0.4116
Kids 2	1.545726	1.20815863	1.28	0.2087
Kids 3	-0.526817	1.32748011	-0.40	0.6938
Kids 4	0.00000	0.0000000		

You can request other statistics and tests by using PROC SURVEYREG. You can also analyze data from a more complex sample design. The remainder of this chapter provides more detailed information.

#### **Output Data Sets**

You can use the OUTPUT statement to create a new SAS data set that contains the estimated linear predictors and their standard error estimates, the residuals from the linear regression, and the confidence limits for the predictors. See the section "OUTPUT Statement" on page 6548 for more details.

You can use the Output Delivery System (ODS) to create SAS data sets that capture the outputs from PROC SURVEYREG. For more information about ODS, see Chapter 20, "Using the Output Delivery System."

For example, to save the "ParameterEstimates" table (Figure 86.7) in the previous section in an output data set, you use the ODS OUTPUT statement as follows:

```
title1 'Ice Cream Spending Analysis';
title2 'Stratified Simple Random Sample Design';
proc surveyreg data=IceCream total=StudentTotals;
    strata Grade /list;
    class Kids;
    model Spending = Income Kids / solution;
    weight Weight;
    ods output ParameterEstimates = MyParmEst;
    run;
```

The statement

```
ods output ParameterEstimates = MyParmEst;
```

requests that the "ParameterEstimates" table that appears in Figure 86.7 be placed into a SAS data set MyParmEst.

The PRINT procedure displays observations of the data set MyParmEst:

```
proc print data=MyParmEst;
run;
```

Figure 86.8 displays the observations in the data set MyParmEst.

Figure 86.8 The Data Set MyParmEst

			m Spending Ana	-	_	
		Stratified Si	mple Random San	mpre besign	.1	
Obs	Parameter	Estimate	StdErr	DenDF	tValue	Probt
1	Intercept	-26.086882	2.44108058	37	-10.69	<.0001
2	Income	0.776699	0.04295904	37	18.08	<.0001
3	Kids 1	0.888631	1.07000634	37	0.83	0.4116
4	Kids 2	1.545726	1.20815863	37	1.28	0.2087
5	Kids 3	-0.526817	1.32748011	37	-0.40	0.6938
6	Kids 4	0.00000	0.0000000	37		

The section "ODS Table Names" on page 6573 gives the complete list of the tables produced by PROC SURVEYREG.

# **Syntax: SURVEYREG Procedure**

The following statements are available in PROC SURVEYREG:

```
PROC SURVEYREG < options>;
BY variables;
CLASS variables;
CLUSTER variables;
CONTRAST 'label' effect values < ... effect values > </ options>;
DOMAIN variables < variable* variable * variable* variable* variable ... >;
ESTIMATE 'label' effect values < ... effect values > </ options>;
MODEL dependent = < effects > </ options>;
OUTPUT < keyword < = variable-name > ... keyword < = variable-name >> </ option>;
REPWEIGHTS variables < / options>;
STRATA variables < / options>;
WEIGHT variable;
```

The PROC SURVEYREG and MODEL statements are required. If your model contains classification effects, you must list the classification variables in a CLASS statement, and the CLASS

statement must precede the MODEL statement. If you use a CONTRAST statement or an ESTI-MATE statement, the MODEL statement must precede the CONTRAST or ESTIMATE statement.

The CLASS, CLUSTER, STRATA, CONTRAST, and ESTIMATE statements can appear multiple times. You should use only one MODEL statement and one WEIGHT statement.

#### **PROC SURVEYREG Statement**

#### PROC SURVEYREG < options>;

The PROC SURVEYREG statement invokes the procedure. It optionally names the input data sets and specifies the variance estimation method.

You can specify the following options in the PROC SURVEYREG statement.

#### ALPHA= $\alpha$

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

#### DATA=SAS-data-set

specifies the SAS data set to be analyzed by PROC SURVEYREG. 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 more information, see the section "Missing Values" on page 6552.

#### **NOMCAR**

requests that the procedure treat missing values in the variance computation as *not missing completely at random* (NOMCAR) for Taylor series variance estimation. When you specify the NOMCAR option, PROC SURVEYREG computes variance estimates by analyzing the nonmissing values as a domain or subpopulation, where the entire population includes both nonmissing and missing domains. See the section "Missing Values" on page 6552 for more details.

By default, PROC SURVEYREG completely excludes an observation from analysis if that observation has a missing value, unless you specify the MISSING option. Note that the NOMCAR option has no effect on a classification variable when you specify the MISSING option, which treats missing values as a valid nonmissing level.

The NOMCAR option applies only to Taylor series variance estimation. The replication methods, which you request with the VARMETHOD=BRR and VARMETHOD=JACKKNIFE options, do not use the NOMCAR option.

#### ORDER=DATA | FORMATTED | FREQ | INTERNAL

specifies the sorting 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 might be useful when you use the CONTRAST statement. When the default 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 SURVEYREG interprets values of the ORDER= option.

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

Note that the specified order also determines the ordering for levels of STRATA and CLUS-TER variables.

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

For more information about sorting 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 as a nonnegative *value*, or specifies an input data set that contains the stratum sampling rates. The procedure uses this information to compute a finite population correction for Taylor series variance estimation. The procedure does not use the RATE= option for BRR or jackknife variance estimation, which you request with the VARMETHOD=BRR or VARMETHOD=JACKKNIFE option.

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

For a nonstratified sample design, or for a stratified sample design with the same sampling rate in all strata, you should specify a nonnegative *value* for the RATE= option. If your design is stratified with different sampling rates in the strata, then you should name a SAS data set that contains the stratification variables and the sampling rates. See the section "Specification of Population Totals and Sampling Rates" on page 6553 for more details.

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 SURVEYREG converts that number to a proportion. The procedure treats the value 1 as 100%, and not the percentage form 1%.

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= and RATE= options.

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

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

specifies the total number of primary sampling units in the study population as a positive *value*, or specifies an input data set that contains the stratum population totals. The procedure uses this information to compute a finite population correction for Taylor series variance estimation. The procedure does not use the TOTAL= option for BRR or jackknife variance estimation, which you request with the VARMETHOD=BRR or VARMETHOD=JACKKNIFE option.

For a nonstratified sample design, or for a stratified sample design with the same population total in all strata, you should specify a positive *value* for the TOTAL= option. If your sample design is stratified with different population totals in the strata, then you should name a SAS data set that contains the stratification variables and the population totals. See the section "Specification of Population Totals and Sampling Rates" on page 6553 for more details.

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= and RATE= options.

#### **TRUNCATE**

specifies that class levels should be determined using no more than the first 16 characters of the formatted values of the CLASS, STRATA, and CLUSTER variables. When formatted values are longer than 16 characters, you can use this option in order to revert to the levels as determined in releases before SAS 9.

# VARMETHOD=BRR < (method-options) > VARMETHOD=JACKKNIFE | JK < (method-options) > VARMETHOD=TAYLOR

specifies the variance estimation method. VARMETHOD=TAYLOR requests the Taylor series method, which is the default if you do not specify the VARMETHOD= option or the REPWEIGHTS statement. VARMETHOD=BRR requests variance estimation by balanced repeated replication (BRR), and VARMETHOD=JACKKNIFE requests variance estimation by the delete-1 jackknife method.

For VARMETHOD=BRR and VARMETHOD=JACKKNIFE you can specify *method-options* in parentheses. Table 86.2 summarizes the available *method-options*.

Keyword	Method	(Method-Options)
BRR	Balanced repeated replication	FAY <=value> HADAMARD   H=SAS-data-set OUTWEIGHTS=SAS-data-set PRINTH REPS=number
JACKKNIFE   JK	Jackknife	OUTJKCOEFS= <i>SAS-data-set</i> OUTWEIGHTS= <i>SAS-data-set</i>
TAYLOR	Taylor series	

**Table 86.2** Variance Estimation Method-Options

*Method-options* must be enclosed in parentheses following the method keyword. For example:

#### varmethod=BRR(reps=60 outweights=myReplicateWeights)

The following values are available for the VARMETHOD= option:

BRR < (method-options) > requests balanced repeated replication (BRR) variance estimation. The BRR method requires a stratified sample design with two primary sampling units (PSUs) per stratum. See the section "Balanced Repeated Replication (BRR) Method" on page 6557 for more information.

You can specify the following *method-options* in parentheses following VARMETHOD=BRR.

#### FAY <=value>

requests Fay's method, a modification of the BRR method, for variance estimation. See the section "Fay's BRR Method" on page 6558 for more information.

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 SURVEYREG generates an appropriate Hadamard matrix for replicate construction. See the sections "Balanced Repeated Replication (BRR) Method" on page 6557 and "Hadamard Matrix" on page 6560 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=*SAS-data-set* 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,  $\mathbf{A'A} = R\mathbf{I}$ , where  $\mathbf{A}$  is the Hadamard matrix of dimension R and  $\mathbf{I}$  is an identity matrix. PROC SURVEYREG 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, the procedure 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 taken to be 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 option to display the Hadamard matrix that the procedure uses to construct replicates for BRR.

#### **OUTWEIGHTS=**SAS-data-set

names a SAS data set that contains replicate weights. See the section "Balanced Repeated Replication (BRR) Method" on page 6557 for information about replicate weights. See the section "Replicate Weights Output Data Set" on page 6567 for more details about the contents of the OUTWEIGHTS= data set.

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

#### **PRINTH**

displays the Hadamard matrix.

When you provide your own Hadamard matrix with the HADAMARD= method-option, only the rows and columns of the Hadamard matrix that are used by the procedure are displayed. See the sections "Balanced Repeated Replication (BRR) Method" on page 6557 and "Hadamard Matrix" on page 6560 for details.

The PRINTH method-option is not available when you provide replicate weights with the REPWEIGHTS statement because the procedure 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 with the HADAMARD= method-option, the number of replicates should be greater than the number of strata and should be a multiple of 4. See the section "Balanced Repeated Replication (BRR) Method" on page 6557 for more information. 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 used to be larger than the REPS= value that you specify.

If you provide a Hadamard matrix with the HADAMARD= methodoption, the value of REPS= must not be less 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 the 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 jack-knife method. See the section "Jackknife Method" on page 6559 for details. If you provide replicate weights with a REPWEIGHTS statement, VARMETHOD=JACKKNIFE is the default variance estimation method.

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

#### **OUTWEIGHTS=**SAS-data-set

names a SAS data set that contains replicate weights. See the section "Jackknife Method" on page 6559 for information about replicate weights. See the section "Replicate Weights Output Data Set" on page 6567 for more details about the contents of the OUT-WEIGHTS= data set.

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

#### **OUTJKCOEFS=**SAS-data-set

names a SAS data set that contains jackknife coefficients. See the section "Jackknife Coefficients Output Data Set" on page 6567 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 and if there is no REPWEIGHTS statement. See the section "Taylor Series (Linearization)" on page 6556 for more information.

#### **BY Statement**

#### BY variables;

You can specify a BY statement with PROC SURVEYREG to obtain separate analyses on observations in groups defined by the BY variables.

Note that using a BY statement provides completely separate analyses of the BY groups. It does not provide a statistically valid subpopulation or domain analysis, where the total number of units in the subpopulation is not known with certainty. For more information about subpopulation analysis for sample survey data, see Cochran (1977).

When a BY statement appears, the procedure expects the input data sets to be sorted in order of the BY variables. If you specify more than one BY statement, the procedure uses only the latest BY statement and ignores any previous ones.

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 BY statement option NOTSORTED or DESCENDING. 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.

For more information about the BY statement, see SAS Language Reference: Concepts. For more information about the DATASETS procedure, see the Base SAS Procedures Guide.

#### **CLASS Statement**

#### **CLASS** variables;

The CLASS statement specifies the classification variables to be used in the model. Typical classification variables are TREATMENT, GENDER, RACE, GROUP, and REPLICATION. If you specify the CLASS statement, it must appear before the MODEL statement.

Classification variables can be either character or numeric. Class levels are determined from the formatted values of the CLASS variables. Thus, you can use formats to group values into levels.

See the FORMAT procedure in the *Base SAS Procedures Guide* and the FORMAT statement and SAS formats in *SAS Language Reference: Dictionary* for more information.

By default, class levels are determined from the entire formatted values of the CLASS variables. Note that this represents a slight change from previous releases in the way in which class levels are determined. Prior to SAS 9, class levels were determined by using no more than the first 16 characters of the formatted values. If you want to revert to this previous behavior, you can use the TRUNCATE option in the PROC SURVEYREG statement.

You can use multiple CLASS statements to specify classification variables.

#### **CLUSTER Statement**

#### **CLUSTER** variables;

The CLUSTER statement specifies variables that identify clusters in a clustered sample design. 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 identify only the first-stage clusters, or primary sampling units (PSUs), in the CLUSTER statement.

If you provide replicate weights for BRR or jackknife variance estimation with the 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. The formatted values of the CLUSTER variables determine the CLUSTER variable levels. Thus, you can use formats to group values into levels. See the FORMAT procedure in the *Base SAS Procedures Guide* and the FORMAT statement and SAS formats in *SAS Language Reference: Dictionary* for more information.

By default, clusters are determined from the entire formatted values of the CLUSTER variables. Note that this represents a slight change from previous releases in the way in which clusters are determined. Prior to SAS 9, clusters were determined by using no more than the first 16 characters of the formatted values. If you want to revert to this previous behavior, you can use the TRUNCATE option in the PROC SURVEYREG statement.

You can use multiple CLUSTER statements to specify cluster variables. The procedure uses variables from all CLUSTER statements to create clusters.

#### **CONTRAST Statement**

**CONTRAST** 'label' effect values </ options>;

**CONTRAST** 'label' effect values < . . . effect values > < / options > ;

The CONTRAST statement provides custom hypothesis tests for linear combinations of the regression parameters  $H_0$ :  $L\beta = 0$ , where L is the vector or matrix you specify and  $\beta$  is the vector of regression parameters. Thus, to use this feature, you must be familiar with the details of the model parameterization used by PROC SURVEYREG. For information about the parameterization, see the section "GLM Parameterization of Classification Variables and Effects" on page 369 in Chapter 18, "Shared Concepts and Topics."

Each term in the MODEL statement, called an *effect*, is a variable or a combination of variables. You can specify an effect with a variable name or a special notation by using variable names and operators. For more details about how to specify an effect, see the section "Specification of Effects" on page 2486 in Chapter 39, "The GLM Procedure."

For each CONTRAST statement, PROC SURVEYREG computes Wald's F test. The procedure displays this value with the degrees of freedom, and identifies it with the contrast label. The numerator degrees of freedom for Wald's F test equal rank( $\mathbf{L}$ ). The denominator degrees of freedom equal the number of clusters (or the number of observations if there is no CLUSTER statement) minus the number of strata. Alternatively, you can use the DF= option in the MODEL statement to specify the denominator degrees of freedom.

You can specify any number of CONTRAST statements, but they must appear after the MODEL statement.

In the CONTRAST statement,

label identifies the contrast in the output. A label is required for every contrast

specified. Labels must be enclosed in single quotes.

effect identifies an effect that appears in the MODEL statement. You can use the

INTERCEPT keyword as an effect when an intercept is fitted in the model. You do not need to include all effects that are in the MODEL statement.

values are constants that are elements of L associated with the effect.

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

Ε

displays the entire coefficient L vector or matrix.

#### **NOFILL**

requests no filling in higher-order effects. When you specify only certain portions of **L**, by default PROC SURVEYREG constructs the remaining elements from the context. (For more information, see the section "Specification of ESTIMATE Expressions" on page 2507 in Chapter 39, "The GLM Procedure.")

When you specify the NOFILL option, PROC SURVEYREG does not construct the remaining portions and treats the vector or matrix **L** as it is defined in the CONTRAST statement.

#### SINGULAR=value

tunes the estimability checking. If v is a vector, define ABS(v) to be the largest absolute value of the elements of v. For a row vector  $\mathbf{l}$  of the matrix  $\mathbf{L}$ , define

$$c = \begin{cases} ABS(\mathbf{I}) & \text{if } ABS(\mathbf{I}) > 0\\ 1 & \text{otherwise} \end{cases}$$

If ABS( $\mathbf{l}-\mathbf{l}\mathbf{H}$ ) is greater than c\*value, then  $\mathbf{l}\boldsymbol{\beta}$  is declared nonestimable. Here,  $\mathbf{H}$  is the matrix  $(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}$ . The *value* must be between 0 and 1; the default is  $10^{-4}$ .

As stated previously, the CONTRAST statement enables you to perform hypothesis tests  $H_0$ :  $L\beta = 0$ .

If the L matrix contains more than one contrast, then you can separate the rows of the L matrix with commas.

For example, for the model

```
proc surveyreg;
  class A B;
  model Y=A B;
run;
```

with A at 5 levels and B at 2 levels, the parameter vector is

$$(\mu \alpha_1 \alpha_2 \alpha_3 \alpha_4 \alpha_5 \beta_1 \beta_2)$$

To test the hypothesis that the pooled A linear and A quadratic effect is zero, you can use the following L matrix:

$$\mathbf{L} = \begin{bmatrix} 0 & -2 & -1 & 0 & 1 & 2 & 0 & 0 \\ 0 & 2 & -1 & -2 & -1 & 2 & 0 & 0 \end{bmatrix}$$

The corresponding CONTRAST statement is

```
contrast 'A Linear & Quadratic'
a -2 -1 0 1 2,
a 2 -1 -2 -1 2;
```

#### **DOMAIN Statement**

**DOMAIN** *variables* < *variable\* variable\* variab* 

The DOMAIN statement requests analysis for subpopulations, or domains, 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 be unrelated to the sample design. Therefore, the sample sizes for the domains are random variables. In order to incorporate this variability into the variance estimation, you should use a DOMAIN statement.

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.

You should use the DOMAIN statement on the entire data set to perform the domain analysis. Creating a new data set from a single domain and analyzing that with SURVEYREG 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 defined by the domain variables.

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

#### **ESTIMATE Statement**

**ESTIMATE** 'label' effect values < / options > ;

**ESTIMATE** 'label' effect values < ... effect values > < / options > ;

You can use an ESTIMATE statement to estimate a linear function of the regression parameters by multiplying a row vector  $\mathbf{l}$  by the parameter estimate vector  $\hat{\boldsymbol{\beta}}$ .

Each term in the MODEL statement, called an *effect*, is a variable or a combination of variables. You can specify an effect with a variable name or with a special notation by using variable names and operators.

For more details about how to specify an effect, see the section "Specification of Effects" on page 2486 in Chapter 39, "The GLM Procedure."

PROC SURVEYREG checks the linear function for estimability. (See the SINGULAR= option).

The procedure displays the estimate  $\mathbf{l}\hat{\boldsymbol{\beta}}$  along with its standard error and t test. If you specify the CLPARM option in the MODEL statement, PROC SURVEYREG also displays confidence limits for the linear function. By default, the degrees of freedom for the t test equal the number of clusters (or the number of observations if there is no CLUSTER statement) minus the number of strata. Alternatively, you can specify the degrees of freedom with the DF= option in the MODEL statement.

You can specify any number of ESTIMATE statements, but they must appear after the MODEL statement.

In the ESTIMATE statement,

identifies the linear function I in the output. A label is required for every function specified. Labels must be enclosed in single quotes.

effect identifies an effect that appears in the MODEL statement. You can use the INTERCEPT keyword as an effect when an intercept is fitted in the model. You do not need to include all effects that are in the MODEL statement.

values are constants that are elements of the vector **l** associated with the effect. For example, the following statement forms an estimate that is the difference between the parameters estimated for the first and second levels of the CLASS variable A:

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

#### **DIVISOR**=value

specifies a value by which to divide all coefficients so that fractional coefficients can be entered as integers.

For example, you can use

```
estimate '1/3(A1+A2) - 2/3A3' a 1 1 -2 / divisor=3; instead of estimate '1/3(A1+A2) - 2/3A3' a 0.33333 0.33333 -0.66667;
```

Ε

displays the entire coefficient vector **l**.

#### **NOFILL**

requests no filling in higher-order effects. When you specify only certain portions of the vector **l**, by default PROC SURVEYREG constructs the remaining elements from the context. (See the section "Specification of ESTIMATE Expressions" on page 2507 in Chapter 39, "The GLM Procedure.") When you specify the NOFILL option, PROC SURVEYREG does not construct the remaining portions and treats the vector **l** as it is defined in the ESTIMATE statement.

#### SINGULAR=value

tunes the estimability checking. If  $\mathbf{v}$  is a vector, define ABS( $\mathbf{v}$ ) to be the largest absolute value of the elements of  $\mathbf{v}$ . For a row vector  $\mathbf{l}$ , define

$$c = \begin{cases} ABS(\mathbf{I}) & \text{if } ABS(\mathbf{I}) > 0\\ 1 & \text{otherwise} \end{cases}$$

If ABS( $\mathbf{l}-\mathbf{l}\mathbf{H}$ ) is greater than c\*value, then  $\mathbf{l}\boldsymbol{\beta}$  is declared nonestimable. Here,  $\mathbf{H}$  is the matrix  $(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}$ . The *value* must be between 0 and 1; the default is  $10^{-4}$ .

#### **MODEL Statement**

**MODEL** dependent = < effects > </ options > ;

The MODEL statement specifies the dependent (response) variable and the independent (regressor) variables or effects. The dependent variable must be numeric. Each term in a MODEL statement, called an *effect*, is a variable or a combination of variables. You can specify an effect with a variable name or with special notation by using variable names and operators. For more information about how to specify an effect, see the section "Specification of Effects" on page 2486 in Chapter 39, "The GLM Procedure."

Only one MODEL statement is allowed for each PROC SURVEYREG statement. If you specify more than one MODEL statement, the procedure uses the first model and ignores the rest.

You can specify the following options in the MODEL statement after a slash (/).

#### **ADJRSQ**

requests the procedure compute the adjusted multiple R-square.

#### **ANOVA**

requests the ANOVA table be produced in the output. By default, the ANOVA table is not printed in the output.

#### **CLPARM**

requests confidence limits for the parameter estimates. The SURVEYREG procedure determines the confidence coefficient by using the ALPHA= option, which by default equals 0.05 and produces 95% confidence bounds. The CLPARM option also requests confidence limits for all the estimable linear functions of regression parameters in the ESTIMATE statements.

Note that when there is a CLASS statement, you need to use the SOLUTION option with the CLPARM option to obtain the parameter estimates and their confidence limits.

#### COVB

displays the estimated covariance matrix of the estimated regression estimates.

#### **DEFF**

displays design effects for the regression coefficient estimates.

#### DF=value

specifies the denominator degrees of freedom for the *F* tests and the degrees of freedom for the *t* tests. The default is the number of clusters (or the number of observations if there is no CLUSTER statement) minus the number of actual strata. The number of actual strata equals the number of strata in the data before collapsing minus the number of strata collapsed plus 1. See the section "Stratum Collapse" on page 6564 for details about "collapsing of strata."

#### I | INVERSE

displays the inverse or the generalized inverse of the X'X matrix. When there is a WEIGHT variable, the procedure displays the inverse or the generalized inverse of the X'WX matrix, where W is the diagonal matrix constructed from WEIGHT variable values.

#### **NOINT**

omits the intercept from the model.

#### **PARMLABEL**

displays the labels of the parameters in the "Estimated Regression Coefficients" table.

#### SINGULAR=value

tunes the estimability checking. If v is a vector, define ABS(v) to be the largest absolute value of the elements of v. For a row vector l of the matrix L, define

$$c = \begin{cases} ABS(\mathbf{l}) & \text{if } ABS(\mathbf{l}) > 0\\ 1 & \text{otherwise} \end{cases}$$

If ABS( $\mathbf{l}-\mathbf{l}\mathbf{H}$ ) is greater than c\*value, then  $\mathbf{l}\boldsymbol{\beta}$  is declared nonestimable. Here,  $\mathbf{H}$  is the matrix  $(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{X}$ . The *value* must be between 0 and 1; the default is  $10^{-4}$ .

#### **SOLUTION**

displays a solution to the normal equations, which are the parameter estimates. The SOLU-TION option is useful only when you use a CLASS statement. If you do not specify a CLASS statement, PROC SURVEYREG displays parameter estimates by default. But if you specify a CLASS statement, PROC SURVEYREG does not display parameter estimates unless you also specify the SOLUTION option.

#### **VADJUST=DF | NONE**

specifies whether to use degrees of freedom adjustment (n-1)/(n-p) in the computation of the matrix **G** for the variance estimation. If you do not specify the VADJUST= option, by default, PROC SURVEYREG uses the degrees-of-freedom adjustment that is equivalent to the VARADJ=DF option. If you do not want to use this variance adjustment, you can specify the VADJUST=NONE option.

#### X | XPX

displays the X'X matrix, or the X'WX matrix when there is a WEIGHT variable, where W is the diagonal matrix constructed from WEIGHT variable values. The X option also displays the crossproducts vector X'y or X'Wy.

#### **OUTPUT Statement**

```
OUTPUT < OUT=SAS-data-set> < keyword <=variable-name> ... keyword <=variable-name> ... keyword <=variable-name> ...
```

The OUTPUT statement creates a new SAS data set that contains all the variables in the input data set and, optionally, the estimated linear predictors and their standard error estimates, the residuals from the linear regression, and the confidence limits for the predictors.

You can specify the following options in the OUTPUT statement.

#### OUT=SAS-data-set

gives the name of the new output data set. By default, the procedure uses the DATA*n* convention to name the new data set.

#### keyword < =variable-name >

specifies the statistics to include in the output data set and names the new variables that contain the statistics. You can specify a keyword for each desired statistic (see the following list of keywords). Optionally, you can name a statistic by providing a variable name followed an equal sign to contain the statistic. For example,

#### output out=myOutDataSet p=myPredictor;

creates a SAS data set myOutDataSet that contains the predicted values in the variable myPredictor.

The keywords allowed and the statistics they represent are as follows:

LCLM | L

lower bound of a  $100(1-\alpha)\%$  confidence interval for the expected value (mean) of the predicted value. The  $\alpha$  level is equal to the value of the AL-PHA= option in the OUTPUT statement or, if this option is not specified, to the ALPHA= option in the PROC SURVEYREG statement. If neither of these options is set, then  $\alpha=0.05$  by default, resulting in the lower bound for a 95% confidence interval. If no variable name is given for this keyword, the default variable name is \_LCLM\_.

PREDICTED | PRED | P predicted values. If no variable name is given for this keyword, the default variable name is PREDICTED .

RESIDUAL | R residuals, calculated as ACTUAL - PREDICTED. If no variable name is given for this keyword, the default variable name is RESIDUAL .

STDP | STD standard error of the mean predicted value. If no variable name is given for this keyword, the default variable name is STD.

UCLM | U

upper bound of a  $100(1-\alpha)\%$  confidence interval for the expected value (mean) of the predicted value. The  $\alpha$  level is equal to the value of the AL-PHA= option in the OUTPUT statement or, if this option is not specified, to the ALPHA= option in the PROC SURVEYREG statement. If neither of these options is set, then  $\alpha=0.05$  by default, resulting in the upper bound for a 95% confidence interval. If no variable name is given for this keyword, the default variable name is UCLM.

The following option is available in the OUTPUT statement and is specified after a slash(/):

#### ALPHA= $\alpha$

specifies the level of significance  $\alpha$  for  $100(1-\alpha)\%$  confidence intervals. By default,  $\alpha$  is equal to the value of the ALPHA= option in the PROC SURVEYREG statement or 0.05 if that option is not specified. You can use values between 0 and 1.

#### **REPWEIGHTS Statement**

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

The REPWEIGHTS statement names variables that provide replicate weights for BRR or jackknife variance estimation, which you request with the VARMETHOD=BRR or VARMETHOD=JACKKNIFE option in the PROC SURVEYREG statement. If you do not provide replicate weights for these methods by using a REPWEIGHTS statement, then the procedure constructs replicate weights for the analysis. See the sections "Balanced Repeated Replication (BRR) Method" on page 6557 and "Jackknife Method" on page 6559 for information about replicate weights.

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 SURVEYREG statement, the procedure uses VARMETHOD=JACKKNIFE by default.

If you specify a REPWEIGHTS statement but do not include a WEIGHT statement, the procedure 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 (/):

#### **DF**=df

specifies the denominator degrees of freedom for various tests in the regression analysis. The value of *df* must be a positive number. See the section "Degrees of Freedom" on page 6560 for details. By default, the denominator degrees of freedom equals the number of REPWEIGHTS variables.

#### JKCOEFS=value

specifies a jackknife coefficient for VARMETHOD=JACKKNIFE. The coefficient *value* must be a nonnegative number less than one. See the section "Jackknife Method" on page 6559 for details about jackknife coefficients.

You can use this option to specify a single value of the jackknife coefficient, which the procedure uses for all replicates. To specify different coefficients for different replicates, use the JKCOEFS=*values* or JKCOEFS=*SAS-data-set* option.

#### 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 less than one. 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.

See the section "Jackknife Method" on page 6559 for details about jackknife coefficients.

To specify different coefficients for different replicates, you can also use the JKCOEFS=SAS-data-set option. To specify a single jackknife coefficient for all replicates, use the JKCOEFS=value option.

#### 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 less than one. The observations in the JKCOEFS= data set should correspond to the replicates that are identified by the REP-WEIGHTS 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 REP-WEIGHTS 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 6559 for details about jackknife coefficients.

To specify different coefficients for different replicates, you can also use the JKCOEFS=*values* option. To specify a single jackknife coefficient for all replicates, use the JKCOEFS=*value* option.

#### STRATA Statement

#### STRATA variables < /options > ;

The STRATA statement specifies variables that form the strata in a stratified sample design. The combinations of categories of STRATA variables define the strata in the sample.

If your sample design has stratification at multiple stages, you should identify only the first-stage strata in the STRATA statement. See the section "Specification of Population Totals and Sampling Rates" on page 6553 for more information.

If you provide replicate weights for BRR or jackknife variance estimation with the 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. By default, strata are determined from the entire formatted values of the STRATA variables. Note that this represents a slight change from previous releases in the way in which strata are determined. Prior to SAS 9, strata were determined by using no more than the first 16 characters of the formatted values. If you want to revert to this previous behavior, you can use the TRUNCATE option in the PROC SURVEYREG statement.

Thus, you can use formats to group values into levels. See the FORMAT procedure in the *Base SAS Procedures Guide* and the FORMAT statement and SAS formats in *SAS Language Reference: Dictionary* for more information.

You can use multiple STRATA statements to specify stratum variables.

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

#### LIST

displays a "Stratum Information" table, which includes values of the STRATA variables, and the number of observations, number of clusters, population total, and sampling rate for each stratum. This table also displays stratum collapse information.

#### **NOCOLLAPSE**

prevents the procedure from collapsing, or combining, strata that have only one sampling unit for the Taylor series variance estimation. By default, the procedure collapses strata that contain only one sampling unit for the Taylor series method. See the section "Stratum Collapse" on page 6564 for details.

#### **WEIGHT Statement**

#### **WEIGHT** | **WGT** *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 6552 for more information. If you specify more than one WEIGHT statement, the procedure uses only the first WEIGHT statement and ignores the rest.

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

If you do not specify a WEIGHT statement or a REPWEIGHTS statement, PROC SURVEYREG assigns all observations a weight of one.

# **Details: SURVEYREG Procedure**

# **Missing Values**

If you have missing values in your survey data for any reason, such as nonresponse, this can compromise the quality of your survey results. If the respondents are different from the nonrespondents with regard to a survey effect or outcome, then survey estimates might be biased and cannot 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/or 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 SURVEYREG. See Cochran (1977), Kalton and Kaspyzyk (1986), and Brick and Kalton (1996) for more information.

If an observation has a missing value or a nonpositive value for the WEIGHT variable, then that observation is excluded from the analysis.

An observation is also excluded if it has a missing value for design variables such as STRATA variables, CLUSTER variables, and DOMAIN variables, unless missing values are regarded as a legitimate categorical level for these variables, as specified by the MISSING option.

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. This treatment is based on the assumption that the missing values are missing completely at random (MCAR). However, this assumption sometimes is not true. For example, evidence from other surveys might suggest that observations with missing values are systematically different from observations without missing values. If you believe that missing values are not missing completely at random, then you can specify the NOMCAR option to include these observations with missing values in the dependent variable and the independent variables in the variance estimation.

Whether or not the NOMCAR option is used, observations with missing or invalid values for WEIGHT, STRATA, CLUSTER, or DOMAIN variables are always excluded, unless the MISSING option is also specified.

When you specify the NOMCAR option, the procedure treats observations with and without missing values for variables in the regression model as two different domains, and it performs a domain analysis in the domain of nonmissing observations.

If you use a REPWEIGHTS statement, all REPWEIGHTS variables must contain nonmissing values.

# **Survey Design Information**

#### **Specification of Population Totals and Sampling Rates**

To include a finite population correction (*fpc*) in Taylor series variance estimation, you can input either the sampling rate or the population total by using the RATE= or TOTAL= option in the PROC SURVEYREG statement. (You cannot specify both of these options in the same PROC SURVEYREG statement.) The RATE= and TOTAL= options apply only to Taylor series variance estimation. The procedure does not use a finite population correction for 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, it is appropriate to ignore this correction. 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. See the section "Primary Sampling Units (PSUs)" on page 6554 for more details.

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. 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. Or 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, then 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 SURVEYREG converts that number to a proportion. The procedure treats the value 1 as 100%, and not the percentage form 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.

#### **Primary Sampling Units (PSUs)**

When you have clusters, or primary sampling units (PSUs), in your sample design, the procedure estimates variance from the variation among PSUs when the Taylor series variance method is used. See the section "Variance Estimation" on page 6556 for more information.

BRR or jackknife variance estimation methods draw multiple replicates (or subsamples) from the full sample by following a specific resampling scheme. These subsamples are constructed by deleting PSUs from the full sample.

If you use a REPWEIGHTS statement to provide replicate weights for BRR or jackknife 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 SURVEYREG treats each observation as a PSU.

#### **Computational Details**

#### **Notation**

For a stratified clustered sample design, observations are represented by an  $n \times (p+2)$  matrix

$$(\mathbf{w}, \mathbf{y}, \mathbf{X}) = (w_{hij}, y_{hij}, \mathbf{x}_{hij})$$

where

- w denotes the sampling weight vector
- y denotes the dependent variable
- **X** denotes the  $n \times p$  design matrix. (When an effect contains only classification variables, the columns of **X** that correspond this effect contain only 0s and 1s; no reparameterization is made.)
- h = 1, 2, ..., 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 (including an intercept if the INTERCEPT effect is included in the MODEL statement)
- $n = \sum_{h=1}^{H} \sum_{i=1}^{n_h} m_{hi}$  is the total number of observations in the sample

Also,  $f_h$  denotes the sampling rate for stratum h. You can use the TOTAL= or RATE= option to input population totals or sampling rates. See the section "Specification of Population Totals and Sampling Rates" on page 6553 for details. If you input stratum totals, PROC SURVEYREG computes  $f_h$  as the ratio of the stratum sample size to the stratum total. If you input stratum sampling rates, PROC SURVEYREG uses these values directly for  $f_h$ . If you do not specify the TOTAL= or RATE= option, then the procedure assumes that the stratum sampling rates  $f_h$  are negligible, and a finite population correction is not used when computing variances.

#### **Regression Coefficients**

PROC SURVEYREG solves the normal equations  $X'WX\beta = X'Wy$  by using a modified sweep routine that produces a generalized (g2) inverse  $(X'WX)^-$  and a solution (Pringle and Rayner 1971)

$$\hat{\beta} = (X'WX)^{-}X'Wy$$

where **W** is the diagonal matrix constructed from WEIGHT variable values.

For models with class variables, there are more design matrix columns than there are degrees of freedom (df) for the effect. Thus, there are linear dependencies among the columns. In this case, the parameters are not estimable; there is an infinite number of least squares solutions. PROC

SURVEYREG uses a generalized (g2) inverse to obtain values for the estimates. The solution values are not displayed unless you specify the SOLUTION option in the MODEL statement. The solution has the characteristic that estimates are zero whenever the design column for that parameter is a linear combination of previous columns. (In strict terms, the solution values should not be called estimates.) With this full parameterization, hypothesis tests are constructed to test linear functions of the parameters that are estimable.

#### **Variance Estimation**

PROC SURVEYREG uses the Taylor series method or replication (resampling) methods to estimate sampling errors of estimators based on complex sample designs (Woodruff (1971); Fuller (1975); Fuller, Kennedy, Schnell, Sullivan, and Park (1989); Särndal, Swensson, and Wretman (1992); Wolter (1985); Rust (1985); Dippo, Fay, and Morganstein (1984); Rao and Shao (1999); Rao, Wu, and Yue (1992); and Rao and Shao (1996)). You can use the VARMETHOD= option to specify a variance estimation method to use. By default, the Taylor series method is used. 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. 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 most commonly used method to estimate the covariance matrix of the regression coefficients for complex survey data. It is the default variance estimation method used by PROC SURVEYREG.

Use the notation described in the section "Notation" on page 6555 to denote the residuals from the linear regression as

$$\mathbf{r} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$$

with  $r_{hij}$  as its elements. Let the  $p \times p$  matrix **G** be defined as

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

where

$$\mathbf{e}_{hij} = w_{hij}r_{hij}\mathbf{x}_{hij}$$

$$\mathbf{e}_{hi} = \sum_{j=1}^{m_{hi}} \mathbf{e}_{hij}$$

$$\bar{\mathbf{e}}_{h..} = \frac{1}{n_h} \sum_{i=1}^{n_h} \mathbf{e}_{hi}.$$

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

$$\widehat{\mathbf{V}}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}'\mathbf{W}\mathbf{X})^{-}\mathbf{G}(\mathbf{X}'\mathbf{W}\mathbf{X})^{-}$$

The factor (n-1)/(n-p) in the computation of the matrix **G** reduces the small sample bias associated with using the estimated function to calculate deviations (Hidiroglou, Fuller, and Hickman 1980). 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 use this adjustment in the variance estimation. If you do not want to use this multiplier in variance estimation, you can specify the VADJUST=NONE option in the MODEL statement to suppress this factor.

#### 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. 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=number option. If a  $number \times number$  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 the 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 rth (r = 1, 2, ..., R) replicate is drawn from the full sample according to the rth row of the Hadamard matrix as follows:

- If the (r, h)th 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)th element of the Hadamard matrix is -1, then the second PSU of stratum h is included in the rth 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 (1985) and Lohr (1999).

By default, an appropriate Hadamard matrix is generated automatically to create the replicates. You can request that the Hadamard matrix be displayed 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.

Let  $\hat{\beta}$  be the estimated regression coefficients from the full sample for  $\beta$ , and let  $\hat{\beta}_r$  be the estimated regression coefficient from the rth replicate by using replicate weights. PROC SURVEYREG estimates the covariance matrix of  $\hat{\beta}$  by

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

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

You can use the VARMETHOD=BRR(OUTWEIGHTS=) method-option to save the replicate weights into a SAS data set.

#### Fay's BRR Method

Fay's method is a modification of the BRR method, and it requires a stratified sample design with two primary sampling units (PSUs) per stratum. The total number of replicates R is the smallest multiple of 4 that is greater than the total number of strata H. However, if you prefer a larger number of replicates, you can specify the REPS= option.

For each replicate, Fay's method uses a Fay coefficient  $0 \le \epsilon < 1$  to impose a perturbation of the original weights in the full sample that is gentler than using only half samples, as in the traditional BRR method. The Fay coefficient  $0 \le \epsilon < 1$  can be optionally set by the FAY =  $\epsilon$  method-option. By default,  $\epsilon = 0.5$  if only the FAY method-option is used without specifying a value for  $\epsilon$  (Judkins 1990; Rao and Shao 1999). When  $\epsilon = 0$ , Fay's method becomes the traditional BRR method. For more details, see Dippo, Fay, and Morganstein (1984), Fay (1984), Fay (1989), and Judkins (1990).

Let H be the number of strata. Replicates are constructed by using the first H columns of the  $R \times R$  Hadamard matrix, where R is the number of replicates, R > H. The rth (r = 1, 2, ..., R) replicate is created from the full sample according to the rth row of the Hadamard matrix as follows:

- If the (r, h)th element of the Hadamard matrix is 1, then the full sample weight of the first PSU in stratum h is multiplied by  $\epsilon$  and that of the second PSU is multiplied by  $2-\epsilon$  to obtain the rth replicate weights.
- If the (r, h)th element of the Hadamard matrix is -1, then the full sample weight of the first PSU in stratum h is multiplied by  $2-\epsilon$  and that of the second PSU is multiplied by  $\epsilon$  to obtain the rth replicate weights.

You can use the VARMETHOD=BRR(OUTWEIGHTS=) method-option to save the replicate weights into a SAS data set.

By default, an appropriate Hadamard matrix is generated automatically to create the replicates. You can request that the Hadamard matrix be displayed 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.

Let  $\hat{\beta}$  be the estimated regression coefficients from the full sample for  $\beta$ . Let  $\hat{\beta}_r$  be the estimated regression coefficient obtained from the rth replicate by using replicate weights. PROC SUR-VEYREG 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.

#### Jackknife Method

The jackknife method of variance estimation deletes one PSU at a time from the full sample to create replicates. The total number of replicates R is the same as the total number of PSUs. In each replicate, the sample weights of the remaining PSUs are modified by the *jackknife coefficient*  $\alpha_r$ . The modified weights are called replicate weights.

The jackknife coefficient and replicate weights are described as follows.

**Without Stratification** If there is no stratification in the sample design (no STRATA statement), the jackknife coefficients  $\alpha_r$  are the same for all replicates:

$$\alpha_r = \frac{R-1}{R}$$
 where  $r = 1, 2, ..., R$ 

Denote the original weight in the full sample for the jth member of the ith PSU as  $w_{ij}$ . If the ith PSU is included in the rth replicate (r=1,2,...,R), then the corresponding replicate weight for the jth member of the ith PSU is defined as

$$w_{ij}^{(r)} = w_{ij}/\alpha_r$$

**With Stratification** If the sample design involves stratification, each stratum must have at least two PSUs to use the jackknife method.

Let stratum  $\tilde{h}_r$  be the stratum from which a PSU is deleted for the rth replicate. Stratum  $\tilde{h}_r$  is called the *donor stratum*. Let  $n_{\tilde{h}_r}$  be the total number of PSUs in the donor stratum  $\tilde{h}_r$ . The jackknife coefficients are defined as

$$\alpha_r = \frac{n_{\tilde{h}_r} - 1}{n_{\tilde{h}_r}}$$
 where  $r = 1, 2, ..., R$ 

Denote the original weight in the full sample for the jth member of the ith PSU as  $w_{ij}$ . If the ith PSU is included in the rth replicate (r = 1, 2, ..., R), then the corresponding replicate weight for the jth member of the ith PSU is defined as

$$w_{ij}^{(r)} = \begin{cases} w_{ij} & \text{if } i \text{th PSU is not in the donor stratum } \tilde{h}_r \\ w_{ij}/\alpha_r & \text{if } i \text{th PSU is in the donor stratum } \tilde{h}_r \end{cases}$$

You can use the VARMETHOD=JACKKNIFE(OUTJKCOEFS=) method-option to save the jack-knife coefficients into a SAS data set and use the VARMETHOD=JACKKNIFE(OUTWEIGHTS=) method-option to save the replicate weights into 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.

Let  $\hat{\beta}$  be the estimated regression coefficients from the full sample for  $\beta$ . Let  $\hat{\beta}_r$  be the estimated regression coefficient obtained from the rth replicate by using replicate weights. PROC SUR-VEYREG 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.

#### **Hadamard Matrix**

A Hadamard matrix **H** is a square matrix whose elements are either 1 or -1 such that

$$\mathbf{H}\mathbf{H}' = k\mathbf{I}$$

where k is the dimension of **H** and **I** is the identity matrix of order k. The order k is necessarily 1, 2, or a positive integer that is a multiple of 4.

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

### **Degrees of Freedom**

PROC SURVEYREG produces tests for the significance of model effects, regression parameters, estimable functions specified in the ESTIMATE statement, and contrasts specified in the CONTRAST

statement. It computes all these tests taking into account the sample design. The degrees of freedom for these tests differ from the degrees of freedom for the ANOVA table, which does not consider the sample design.

## **Denominator Degrees of Freedom**

The denominator df refers to the denominator degrees of freedom for F tests and to the degrees of freedom for t tests in the analysis.

For the Taylor series method, the denominator df equals the number of clusters minus the actual number of strata. If there are no clusters, the denominator df equals the number of observations minus the actual number of strata. The actual number of strata equals the following:

- one, if there is no STRATA statement
- the number of strata in the sample, if there is a STRATA statement but the procedure does not collapse any strata
- the number of strata in the sample after collapsing, if there is a STRATA statement and the procedure collapses strata that have only one sampling unit

Alternatively, you can specify your own denominator *df* by using the DF= option in the MODEL statement.

For the BRR method (including Fay's method) without a REPWEIGHTS statement, the denominator *df* equals the number of strata.

For the jackknife method without a REPWEIGHTS statement, the denominator df is equal to the number of replicates minus the *actual number of strata*.

When there is a REPWEIGHTS statement, the denominator *df* equals the number of REP-WEIGHTS variables, unless you specify an alternative in the DF= option in a REPWEIGHTS statement.

### **Numerator Degrees of Freedom**

The numerator df refers to the numerator degrees of freedom for the Wald F statistic associated with an effect or with a contrast. The procedure computes the Wald F statistic for an effect as a Type III test; that is, the test has the following properties:

- The hypothesis for an effect does not involve parameters of other effects except for containing effects (which it must involve to be estimable).
- The hypotheses to be tested are invariant to the ordering of effects in the model.

See the section "Testing Effects" on page 6562 for more information. The numerator df for the Wald F statistic for a contrast is the rank of the  $\mathbf{L}$  matrix that defines the contrast.

## **Testing Effects**

For each effect in the model, PROC SURVEYREG computes an L matrix such that every element of  $L\beta$  is estimable; the L matrix has the maximum possible rank associated with the effect. To test the effect, the procedure uses the Wald F statistic for the hypothesis  $H_0$ :  $L\beta = 0$ . The Wald F statistic equals

$$F_{\text{Wald}} = \frac{(\mathbf{L}\hat{\boldsymbol{\beta}})'(\mathbf{L}'\widehat{\mathbf{V}}\mathbf{L})^{-1}(\mathbf{L}\hat{\boldsymbol{\beta}})}{\text{rank}(\mathbf{L})}$$

with numerator degrees of freedom equal to rank(**L**) and denominator degrees of freedom equal to the number of clusters minus the number of strata (unless you have specified the denominator degrees of freedom with the DF= option in the MODEL statement; see the section "Denominator Degrees of Freedom" on page 6561). It is possible that the **L** matrix cannot be constructed for an effect, in which case that effect is not testable. For more information about how the matrix **L** is constructed, see the discussion in Chapter 15, "The Four Types of Estimable Functions."

## **Analysis of Variance (ANOVA)**

PROC SURVEYREG produces an analysis of variance table for the model specified in the MODEL statement. This table is identical to the one produced by the GLM procedure for the model. PROC SURVEYREG computes ANOVA table entries by using the sampling weights, but not the sample design information about stratification and clustering.

The degrees of freedom (*df*) displayed in the ANOVA table are the same as those in the ANOVA table produced by PROC GLM. The Total DF is the total degrees of freedom used to obtain the regression coefficient estimates. The Total DF equals the total number of observations minus 1 if the model includes an intercept. If the model does not include an intercept, the Total DF equals the total number of observations. The Model DF equals the degrees of freedom for the effects in the MODEL statement, not including the intercept. The Error DF equals the Total DF minus the Model DF.

#### Multiple R-Square

PROC SURVEYREG computes a multiple R-square for the weighted regression as

$$R^2 = 1 - \frac{SS_{error}}{SS_{total}}$$

where SS<sub>error</sub> is the error sum of squares in the ANOVA table

$$SS_{error} = \mathbf{r}'\mathbf{W}\mathbf{r}$$

and  $SS_{total}$  is the total sum of squares

$$SS_{total} = \begin{cases} \mathbf{y'Wy} & \text{if no intercept} \\ \mathbf{y'Wy} - \left(\sum_{h=1}^{H} \sum_{i=1}^{n_h} \sum_{j=1}^{m_{hi}} w_{hij} y_{hij}\right)^2 / w... & \text{otherwise} \end{cases}$$

where w... is the sum of the sampling weights over all observations.

## **Adjusted R-Square**

If you specify the option ADJRSQ in the MODEL statement, PROC SURVEYREG computes an multiple R-square adjusted as the weighted regression as

ADJRSQ = 
$$\begin{cases} 1 - \frac{n(1 - R^2)}{n - p} & \text{if no intercept} \\ 1 - \frac{(n - 1)(1 - R^2)}{n - p} & \text{otherwise} \end{cases}$$

where  $R^2$  is the multiple R-square.

## **Root Mean Square Errors**

PROC SURVEYREG computes the square root of mean square errors as

$$\sqrt{\text{MSE}} = \sqrt{n \text{ SS}_{error} / (n - p) w...}$$

where w... is the sum of the sampling weights over all observations.

## **Design Effect**

If you specify the DEFF option in the MODEL statement, PROC SURVEYREG calculates the design effects for the regression coefficients. The design effect of an estimate is the ratio of the actual variance to the variance computed under the assumption of simple random sampling:

$$DEFF = \frac{variance\ under\ the\ sample\ design}{variance\ under\ simple\ random\ sampling}$$

See Kish (1965, p. 258) for more details. PROC SURVEYREG computes the numerator as described in the section "Variance Estimation" on page 6556. And the denominator is computed under the assumption that the sample design is simple random sampling, with no stratification and no clustering.

To compute the variance under the assumption of simple random sampling, PROC SURVEYREG calculates the sampling rate as follows. If you specify both sampling weights and sampling rates (or population totals) for the analysis, then the sampling rate under simple random sampling is calculated as

$$f_{SRS} = n / w...$$

where n is the sample size and w... (the sum of the weights over all observations) estimates the population size. If the sum of the weights is less than the sample size,  $f_{SRS}$  is set to zero. If you specify sampling rates for the analysis but not sampling weights, then PROC SURVEYREG computes the sampling rate under simple random sampling as the average of the stratum sampling rates:

$$f_{SRS} = \frac{1}{H} \sum_{h=1}^{H} f_h$$

If you do not specify sampling rates (or population totals) for the analysis, then the sampling rate under simple random sampling is assumed to be zero:

$$f_{\rm SRS} = 0$$

## **Stratum Collapse**

If there is only one sampling unit in a stratum, then PROC SURVEYREG cannot estimate the variance for this stratum for the Taylor series method. To estimate stratum variances, by default the procedure collapses, or combines, those strata that contain only one sampling unit. If you specify the NOCOLLAPSE option in the STRATA statement, PROC SURVEYREG does not collapse strata and uses a variance estimate of zero for any stratum that contains only one sampling unit.

Note that stratum collapse only applies to Taylor series variance estimation (the default method, also specified by VARMETHOD=TAYLOR). The procedure does not collapse strata for BRR or jackknife variance estimation, which you request with the VARMETHOD=BRR or VARMETHOD=JACKKNIFE option.

If you do not specify the NOCOLLAPSE option for the Taylor series method, PROC SURVEYREG collapses strata according to the following rules. If there are multiple strata that contain only one sampling unit each, then the procedure collapses, or combines, all these strata into a new pooled stratum. If there is only one stratum with a single sampling unit, then PROC SURVEYREG collapses that stratum with the preceding stratum, where strata are ordered by the STRATA variable values. If the stratum with one sampling unit is the first stratum, then the procedure combines it with the following stratum.

If you specify stratum sampling rates by using the RATE=SAS-data-set option, PROC SUR-VEYREG computes the sampling rate for the new pooled stratum as the weighted average of the sampling rates for the collapsed strata. See the section "Computational Details" on page 6555 for details. If the specified sampling rate equals 0 for any of the collapsed strata, then the pooled stratum is assigned a sampling rate of 0. If you specify stratum totals by using the TOTAL=SAS-data-set option, PROC SURVEYREG combines the totals for the collapsed strata to compute the sampling rate for the new pooled stratum.

## Sampling Rate of the Pooled Stratum from Collapse

Assuming that PROC SURVEYREG collapses single-unit strata  $h_1, h_2, \dots, h_c$  into the pooled stratum, the procedure calculates the sampling rate for the pooled stratum as

$$f_{\text{Pooled Stratum}} = \begin{cases} 0 & \text{if any of } f_{h_l} = 0 \text{ where } l = 1, 2, \dots, c \\ \left(\sum_{l=1}^{c} n_{h_l} f_{h_l}^{-1}\right)^{-1} \sum_{l=1}^{c} n_{h_l} & \text{otherwise} \end{cases}$$

#### **Contrasts**

You can use the CONTRAST statement to perform custom hypothesis tests. If the hypothesis is testable in the univariate case, the Wald F statistic for  $H_0: \mathbf{L}\boldsymbol{\beta} = 0$  is computed as

$$F_{\text{Wald}} = \frac{(\mathbf{L}_{\text{Full}} \hat{\boldsymbol{\beta}})' (\mathbf{L}_{\text{Full}}' \widehat{\mathbf{V}} \mathbf{L}_{\text{Full}})^{-1} (\mathbf{L}_{\text{Full}} \hat{\boldsymbol{\beta}})}{\text{rank}(\mathbf{L})}$$

where **L** is the contrast vector or matrix you specify,  $\hat{\beta}$  is the vector of regression parameters,  $\hat{\beta} = (X'WX)^-X'WY$ ,  $\hat{V}$  is the estimated covariance matrix of  $\hat{\beta}$ , rank(**L**) is the rank of **L**, and **L**<sub>Full</sub> is a matrix such that

- L<sub>Full</sub> has the same number of columns as L
- L<sub>Full</sub> has full row rank
- $\bullet$  the rank of  $L_{Full}$  equals the rank of the L matrix
- all rows of L<sub>Full</sub> are estimable functions
- the Wald F statistic computed using the L<sub>Full</sub> matrix is equivalent to the Wald F statistic
  computed by using the L matrix with any row deleted that is a linear combination of previous
  rows

If L is a full-rank matrix and all rows of L are estimable functions, then  $L_{Full}$  is the same as L. It is possible that  $L_{Full}$  matrix cannot be constructed for contrasts in a CONTRAST statement, in which case the contrasts are not testable.

#### **Domain Analysis**

A DOMAIN statement requests that the procedure perform regression analysis for each domain.

For a domain D, let  $I_D$  be the corresponding indicator variable:

$$I_D(h, i, j) = \begin{cases} 1 & \text{if observation } (h, i, j) \text{ belongs to domain } D \\ 0 & \text{otherwise} \end{cases}$$

Let

$$v_{hij} = w_{hij} I_D(h, i, j) = \begin{cases} w_{hij} & \text{if observation } (h, i, j) \text{ belongs to domain } D \\ 0 & \text{otherwise} \end{cases}$$

The regression in domain D uses v as the weight variable.

### **Computational Resources**

Due to the nature of survey data analysis, the SURVEYREG procedure requires more memory compared to analysis that uses the same regression model in the GLM procedure. For details about

the amount of memory related to the modeling, see the section "Computational Resources" on page 2542 in Chapter 39, "The GLM Procedure."

The memory needed in the SURVEYREG procedure to handle the survey design is described as follows.

Let

- *H* be the total number of strata
- $n_c$  be the total clusters in your sample across all H strata if it is a cluster sampling
- p be the total number of parameters in the model

The memory needed (in bytes) is

$$48H + 8pH + 4p(p+1)H$$

For a cluster sampling, the additional memory needed (in bytes) is

$$48H + 8pH + 4p(p+1)H + 4p(p+1)n_c + 16n_c$$

Other small amounts of additional memory are also used by the SURVEYREG procedure. However, when you have large number of clusters or strata, or large number of parameters in your model, the memory described previously dominates the total memory usage in the procedure.

# **Output Data Sets**

You can use the Output Delivery System (ODS) to create a SAS data set from any piece of PROC SURVEYREG output. See the section "ODS Table Names" on page 6573 for more information. For a more detailed description of using ODS, see Chapter 20, "Using the Output Delivery System."

PROC SURVEYREG also provides an OUTPUT statement to create a data set that contains estimated linear predictors and their standard error estimates, the residuals from the linear regression, and the confidence limits for the predictors.

If you use BRR or jackknife variance estimation, PROC SURVEYREG provides an output data set that stores the replicate weights and an output data set that stores the jackknife coefficients for jackknife variance estimation.

## **OUT= Data Set Created by the OUTPUT Statement**

The OUTPUT statement produces an output data set that contains the following:

- all original data from the SAS data set input to PROC SURVEYREG
- the new variables corresponding to the diagnostic measures specified with statistics keywords in the OUTPUT statement (PREDICTED=, RESIDUAL=, and so on)

When any independent variable in the analysis (including all classification variables) is missing for an observation, then all new variables that correspond to diagnostic measures are missing for the observation in the output data set.

When a dependent variable in the analysis is missing for an observation, then the residual variable that corresponds to R is also missing in the output data set. However, the variables corresponding to LCLM, P, STDP, and UCLM are not missing.

## **Replicate Weights Output Data Set**

If you specify the OUTWEIGHTS= method-option for BRR or jackknife method in the VARMETHOD= option, PROC SURVEYREG stores the replicate weights in an output data set. The OUTWEIGHTS= output data set contains all observations used in the analysis or all valid observations in the DATA= input data set. A valid observation is an observation that has a positive value of the WEIGHT variable. Valid observations must also have nonmissing values of the STRATA and CLUSTER variables, unless you specify the MISSING option.

The OUTWEIGHTS= data set contains the following variables:

- all variables in the DATA= input data set
- RepWt\_1, RepWt\_2, ..., RepWt\_n, which are the replicate weight variables

where n 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 SURVEYREG 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 SUR-VEYREG 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

• 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 SURVEYREG 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**

The SURVEYREG procedure produces the following output.

## **Data Summary**

By default, PROC SURVEYREG displays the following information in the "Data Summary" table:

- Number of Observations, which is the total number of observations used in the analysis, excluding observations with missing values
- Sum of Weights, if you specify a WEIGHT statement
- Mean of the dependent variable in the MODEL statement, or Weighted Mean if you specify a WEIGHT statement
- Sum of the dependent variable in the MODEL statement, or Weighted Sum if you specify a WEIGHT statement

## **Design Summary**

When you specify a CLUSTER statement or a STRATA statement, the procedure displays a "Design Summary" table, which provides the following sample design information:

- Number of Strata, if you specify a STRATA statement
- Number of Strata Collapsed, if the procedure collapses strata
- Number of Clusters, if you specify a CLUSTER statement
- Overall Sampling Rate used to calculate the design effect, if you specify the DEFF option in the MODEL statement

## **Domain Summary**

By default, PROC SURVEYREG displays the following information in the "Domain Summary" table:

- Number of Observations, which is the total number of observations used in the analysis
- total number of observations in the current domain
- total number of observations not in the current domain
- Sum of Weights for the observations in the current domain, if you specify a WEIGHT statement

### **Fit Statistics**

By default, PROC SURVEYREG displays the following regression statistics in the "Fit Statistics" table:

- R-square for the regression
- Root MSE, which is the square root of the mean square error
- Denominator DF, which is the denominator degrees of freedom for the *F* tests and also the degrees of freedom for the *t* tests produced by the procedure

#### **Variance Estimation**

If the variance method is not Taylor series (see the section "Variance Estimation" on page 6556) or if the NOMCAR option is used, by default, PROC SURVEYREG displays the following variance estimation specifications the "Variance Estimation" table:

- Method, which is the variance estimation method
- Number of Replicates, which is the number of replicates if you specify the VARMETHOD=BRR or VARMETHOD=JACKKNIFE option
- Hadamard Data Set, which is the name of the SAS data set for the HADAMARD matrix if you specify the VARMETHOD=BRR(HADAMARD=) method-option in the VARMETHOD=BRR option
- Fay Coefficient, which is the value of the FAY coefficient if you specify the VARMETHOD=BRR(FAY) method-option
- Replicate Weights Dataset, which is the name of the SAS data set that contains the replicate weights
- Missing Levels, which indicates whether a missing category for categorical variables is created depending on whether you specify the MISSING option

• Missing Values, which indicates whether observations with missing values are included in the analysis, depending on whether you specify the NOMCAR option

#### Stratum Information

When you specify the LIST option in the STRATA statement, PROC SURVEYREG displays a "Stratum Information" table, which provides the following information for each stratum:

- Stratum Index, which is a sequential stratum identification number
- STRATA variable(s), which lists the levels of STRATA variables for the stratum
- Population Total, if you specify the TOTAL= option
- Sampling Rate, if you specify the TOTAL= option or the RATE= option. If you specify the TOTAL= option, the sampling rate is based on the number of nonmissing observations in the stratum.
- N Obs, which is the number of observations
- number of Clusters, if you specify a CLUSTER statement
- Collapsed, which has the value 'Yes' if the stratum is collapsed with another stratum before analysis

If PROC SURVEYREG collapses strata, the "Stratum Information" table also displays stratum information for the new, collapsed stratum. The new stratum has a Stratum Index of 0 and is labeled 'Pooled.'

#### **Class Level Information**

If you use a CLASS statement to name classification variables, PROC SURVEYREG displays a "Class Level Information" table. This table contains the following information for each classification variable:

- Class Variable, which lists each CLASS variable name
- Levels, which is the number of values or levels of the classification variable
- Values, which lists the values of the classification variable. The values are separated by a white space character; therefore, to avoid confusion, you should not include a white space character within a classification variable value.

#### X'X Matrix

If you specify the XPX option in the MODEL statement, PROC SURVEYREG displays the X'X matrix. When there is a WEIGHT variable, the procedure displays the X'WX matrix. This option also displays the crossproducts vector X'y or X'Wy, where y is the response vector (dependent variable).

#### **Inverse Matrix of X'X**

If you specify the INV option in the MODEL statement, PROC SURVEYREG displays the inverse or the generalized inverse of the  $\mathbf{X}'\mathbf{X}$  matrix. When there is a WEIGHT variable, the procedure displays the inverse or the generalized inverse of the  $\mathbf{X}'\mathbf{W}\mathbf{X}$  matrix.

## **ANOVA for Dependent Variable**

If you specify the ANOVA option in the model statement, PROC SURVEYREG displays an analysis of variance table for the dependent variable. This table is identical to the ANOVA table displayed by the GLM procedure.

#### **Tests of Model Effects**

By default, PROC SURVEYREG displays a "Tests of Model Effects" table, which provides Wald's *F* test for each effect in the model. The table contains the following information for each effect:

- Effect, which is the effect name
- Num DF, which is the numerator degrees of freedom for Wald's F test
- F Value, which is Wald's F statistic
- Pr > F, which is the significance probability corresponding to the F Value

A footnote displays the denominator degrees of freedom, which is the same for all effects.

## **Estimated Regression Coefficients**

PROC SURVEYREG displays the "Estimated Regression Coefficients" table by default when there is no CLASS statement. Also, the procedure displays this table when you specify a CLASS statement and also specify the SOLUTIONS option in the MODEL statement. This table contains the following information for each regression parameter:

- Parameter, which identifies the effect or regressor variable
- Estimate, which is the estimate of the regression coefficient
- Standard Error, which is the standard error of the estimate
- t Value, which is the t statistic for testing  $H_0$ : Parameter = 0
- Pr > | t |, which is the two-sided significance probability corresponding to the t Value

## **Covariance of Estimated Regression Coefficients**

When you specify the COVB option in the MODEL statement, PROC SURVEYREG displays the "Covariance of Estimated Regression Coefficients" matrix.

#### **Coefficients of Contrast**

When you specify the E option in a CONTRAST statement, PROC SURVEYREG displays a "Coefficients of Contrast" table for the contrast. You can use this table to check the coefficients you specified in the CONTRAST statement. Also, this table gives a note for a nonestimable contrast.

## **Analysis of Contrasts**

If you specify a CONTRAST statement, PROC SURVEYREG produces an "Analysis of Contrasts" table, which displays Wald's *F* test for the contrast. If you use more than one CONTRAST statement, the procedure displays all results in the same table. The "Analysis of Contrasts" table contains the following information for each contrast:

- Contrast, which is the label of the contrast
- Num DF, which is the numerator degrees of freedom for Wald's F test
- F Value, which is Wald's F statistic for testing  $H_0$ : Contrast = 0
- Pr > F, which is the significance probability corresponding to the F Value

#### **Coefficients of Estimate**

When you specify the E option in an ESTIMATE statement, PROC SURVEYREG displays a "Coefficients of Estimate" table for the linear function of the regression parameters in the ESTIMATE statement. You can use this table to check the coefficients you specified in the ESTIMATE statement. Also, this table gives a note for a nonestimable function.

## **Analysis of Estimable Functions**

If you specify an ESTIMATE statement, PROC SURVEYREG checks the function for estimability. If the function is estimable, PROC SURVEYREG produces an "Analysis of Estimable Functions" table, which displays the estimate and the corresponding *t* test. If you use more than one ESTIMATE statement, the procedure displays all results in the same table. The table contains the following information for each estimable function:

- Parameter, which is the label of the function
- Estimate, which is the estimate of the estimable linear function

- Standard Error, which is the standard error of the estimate
- t Value, which is the t statistic for testing  $H_0$ : Estimable Function = 0
- Pr > | t |, which is the two-sided significance probability corresponding to the t Value

## **Hadamard Matrix**

If you specify the VARMETHOD=BRR(PRINTH) method-option in the PROC statement, PROC SURVEYREG displays the Hadamard matrix.

When you provide a Hadamard matrix with the VARMETHOD=BRR(HADAMARD=) method-option, the procedure displays only used rows and columns of the Hadamard matrix.

## **ODS Table Names**

PROC SURVEYREG assigns a name to each table it creates. You can use these names to reference the table when by using the Output Delivery System (ODS) to select tables and create output data sets. These names are listed in Table 86.3. For more information about ODS, see Chapter 20, "Using the Output Delivery System."

Table 86.3 ODS Tables Produced by PROC SURVEYREG

ODS Table Name	Description	Statement	Option
ANOVA	ANOVA for dependent variable	MODEL	ANOVA
ClassVarInfo	Class level information	CLASS	default
ContrastCoef	Coefficients of contrast	CONTRAST	E
Contrasts	Analysis of contrasts	CONTRAST	default
CovB	Covariance of estimated regression coefficients	MODEL	COVB
DataSummary	Data summary	PROC	default
DesignSummary	Design summary	STRATA   CLUSTER	default
DomainSummary	Domain summary	DOMAIN	default
Effects	Tests of model effects	MODEL	defect
EstimateCoef	Coefficients of estimate	ESTIMATE	E
Estimates	Analysis of estimable functions	ESTIMATE	default
FitStatistics	Fit statistics	MODEL	default
HadamardMatrix	Hadamard matrix	PROC	PRINTH
InvXPX	Inverse matrix of $X'X$	MODEL	INV
ParameterEstimates	Estimated regression coefficients	MODEL	SOLUTION
StrataInfo	Stratum information	STRATA	LIST

Table 86.3 (continued)

ODS Table Name	Description	Statement	Option
VarianceEstimation	Variance estimation	PROC	default
XPX	X'X matrix	MODEL	XPX

By referring to the names of such tables, you can use the ODS OUTPUT statement to place one or more of these tables in output data sets.

For example, the following statements create an output data set MyStrata, which contains the "StrataInfo" table, an output data set MyParmEst, which contains the "ParameterEstimates" table, and an output data set Cov, which contains the "CovB" table for the ice cream study discussed in the section "Stratified Sampling" on page 6530:

Note that the option CovB is specified in the MODEL statement in order to produce the covariance matrix table.

# **Examples: SURVEYREG Procedure**

# **Example 86.1: Simple Random Sampling**

This example investigates the relationship between the labor force participation rate (LFPR) of women in 1968 and 1972 in large cities in the United States. A simple random sample of 19 cities is drawn from a total of 200 cities. For each selected city, the LFPRs are recorded and saved in a SAS data set Labor. In the following DATA step, LFPR in 1972 is contained in the variable LFPR1972, and the LFPR in 1968 is identified by the variable LFPR1968:

```
data Labor;
   input City $ 1-16 LFPR1972 LFPR1968;
   datalines;
             . 45
New York
                    . 42
Los Angeles
             .50
                     . 50
Chicago
              . 52
                     . 52
Philadelphia .45
                     . 45
                    . 43
Detroit
              .46
San Francisco .55
                    . 55
              . 60
                     . 45
Boston
Pittsburgh
              .49
                     . 34
             . 35
St. Louis
                    . 45
Connecticut .55
                     . 54
Washington D.C. .52
                     . 42
Cincinnati .53
                    . 51
Baltimore
              . 57
                    .49
Newark
              . 53
                     . 54
Minn/St. Paul
              .59
                     .50
Buffalo
              . 64
                     . 58
Houston
              .50
                     .49
Patterson
              .57
                      .56
              . 64
Dallas
                      . 63
```

Assume that the LFPRs in 1968 and 1972 have a linear relationship, as shown in the following model:

```
LFPR1972 = \beta_0 + \beta_1 * LFPR1968 + error
```

You can use PROC SURVEYREG to obtain the estimated regression coefficients and estimated standard errors of the regression coefficients. The following statements perform the regression analysis:

```
title 'Study of Labor Force Participation Rates of Women';
proc surveyreg data=Labor total=200;
   model LFPR1972 = LFPR1968;
run;
```

Here, the TOTAL=200 option specifies the finite population total from which the simple random sample of 19 cities is drawn. You can specify the same information by using the sampling rate option RATE=0.095 (19/200=.095).

Output 86.1.1 summarizes the data information and the fit information.

Output 86.1.1 Summary of Regression Using Simple Random Sampling

Study of Labor Force Participation Rates of Women

The SURVEYREG Procedure

Regression Analysis for Dependent Variable LFPR1972

Data Summary

Number of Observations 19
Mean of LFPR1972 0.52684
Sum of LFPR1972 10.01000

Fit Statistics

R-square 0.3970
Root MSE 0.05657
Denominator DF 18

Output 86.1.2 presents the significance tests for the model effects and estimated regression coefficients. The F tests and t tests for the effects in the model are also presented in these tables.

Output 86.1.2 Regression Coefficient Estimates

•	Effect	Num DF	F Value	Pr > F	
1	Model	1	13.84	0.0016	
	Intercept	1	4.63	0.0452	
:	LFPR1968	1	13.84	0.0016	
NOTE: The de	nominator de Estimated	_		for the F to	ests is 18.
NOTE: The de		Regressi			ests is 18.
		Regressi St	on Coeffi andard		
Parameter	Estimated	Regressi St e	on Coeffi andard Error	icients t Value	Pr >  t

From the regression performed by PROC SURVEYREG, you obtain a positive estimated slope for the linear relationship between the LFPR in 1968 and the LFPR in 1972. The regression coefficients are all significant at the 5% level. The effects Intercept and LFPR1968 are significant in the model at the 5% level. In this example, the F test for the overall model without intercept is the same as the effect LFPR1968.

## **Example 86.2: Simple Random Cluster Sampling**

This example illustrates the use of regression analysis in a simple random cluster sample design. The data are from Särndal, Swensson, and Wretman (1992, p. 652).

A total of 284 Swedish municipalities are grouped into 50 clusters of neighboring municipalities. Five clusters with a total of 32 municipalities are randomly selected. The results from the regression analysis in which clusters are used in the sample design are compared to the results of a regression analysis that ignores the clusters. The linear relationship between the population in 1975 and in 1985 is investigated.

The 32 selected municipalities in the sample are saved in the data set Municipalities:

```
data Municipalities;
   input Municipality Cluster Population85 Population75;
   datalines;
       37
205
              5
                    5
206
       37
            11
                  11
207
       37
            13
                  13
208
              8
                   8
       37
209
       37
            17
                  19
  6
        2
            16
                  15
  7
        2
            70
                  62
  8
        2
            66
                  54
  9
        2
            12
                  12
        2
 10
            60
                  50
 94
       17
             7
                   7
 95
       17
            16
                  16
 96
       17
            13
                  11
 97
       17
            12
                  11
 98
       17
            70
                  67
 99
       17
            20
                  20
100
       17
            31
                  28
            49
101
       17
                  48
276
       50
              6
                   7
277
       50
              9
                  10
278
       50
            24
                  26
279
       50
            10
            67
280
       50
                  64
281
       50
            39
                  35
282
       50
            29
                  27
283
       50
            10
                    9
284
            27
       50
                  31
       10
              7
 52
                    6
 53
       10
              9
                    8
 54
       10
            28
                  27
 55
       10
            12
                  11
 56
       10
           107
                 108
```

The variable Municipality identifies the municipalities in the sample; the variable Cluster indicates the

cluster to which a municipality belongs; and the variables Population85 and Population75 contain the municipality populations in 1985 and in 1975 (in thousands), respectively. A regression analysis is performed by PROC SURVEYREG with a CLUSTER statement:

```
title1 'Regression Analysis for Swedish Municipalities';
title2 'Cluster Simple Random Sampling';
proc surveyreg data=Municipalities total=50;
   cluster Cluster;
   model Population85=Population75;
run;
```

The TOTAL=50 option specifies the total number of clusters in the sampling frame.

Output 86.2.1 displays the data and design summary. Since the sample design includes clusters, the procedure displays the total number of clusters in the sample in the "Design Summary" table.

Output 86.2.1 Regression Analysis for Simple Random Cluster Sampling

```
Regression Analysis for Swedish Municipalities
Cluster Simple Random Sampling

The SURVEYREG Procedure

Regression Analysis for Dependent Variable Population85

Data Summary

Number of Observations 32
Mean of Population85 27.50000
Sum of Population85 880.00000

Design Summary

Number of Clusters 5
```

Output 86.2.2 displays the fit statistics and regression coefficient estimates. In the "Estimated Regression Coefficients" table, the estimated slope for the linear relationship is 1.05, which is significant at the 5% level; but the intercept is not significant. This suggests that a regression line crossing the original can be established between populations in 1975 and in 1985.

Output 86.2.2 Regression Analysis for Simple Random Cluster Sampling

Fit Statistics					
	R-square	0.9860			
	Root MSE	3.0488			
	Denominator DF	4			

#### Output 86.2.2 continued

Estimated Regression Coefficients					
		Standard			
Parameter	Estimate	Error	t Value	Pr >  t	
Intercept	-0.0191292	0.89204053	-0.02	0.9839	
Population75	1.0546253	0.05167565	20.41	<.0001	
NOTE: The den	ominator degre	es of freedom	for the t t	ests is 4.	

The CLUSTER statement is necessary in PROC SURVEYREG in order to incorporate the sample design. If you do not specify a CLUSTER statement in the regression analysis, as in the following statements, the standard deviation of the regression coefficients are incorrectly estimated.

```
title1 'Regression Analysis for Swedish Municipalities';
title2 'Simple Random Sampling';
proc surveyreg data=Municipalities total=284;
   model Population85=Population75;
run:
```

The analysis ignores the clusters in the sample, assuming that the sample design is a simple random sampling. Therefore, the TOTAL= option specifies the total number of municipalities, which is 284.

Output 86.2.3 displays the regression results ignoring the clusters. Compared to the results in Output 86.2.1, the regression coefficient estimates are the same. However, without using clusters, the regression coefficients have a smaller variance estimate, as in Output 86.2.3. By using clusters in the analysis, the estimated regression coefficient for effect Population75 is 1.05, with the estimated standard error 0.05, as displayed in Output 86.2.1; without using the clusters, the estimate is 1.05, but with the estimated standard error 0.04, as displayed in Output 86.2.3. To estimated the variance of the regression coefficients correctly, you should include the clustering information in the regression analysis.

Output 86.2.3 Regression Analysis for Simple Random Sampling

Regression Analysis for Swedish Municipalities
Simple Random Sampling

The SURVEYREG Procedure

Regression Analysis for Dependent Variable Population85

Data Summary

Number of Observations 32
Mean of Population85 27.50000
Sum of Population85 880.00000

#### Output 86.2.3 continued

```
Fit Statistics
                                    0.9860
                 R-square
                 Root MSE
                                  3.0488
                 Denominator DF
                                        31
              Estimated Regression Coefficients
                             Standard
                                Error t Value Pr > |t|
Parameter
                Estimate
                                          -0.03
                                                     0.9775
              -0.0191292
                           0.67417606
Intercept
               1.0546253 0.03668414
                                          28.75
                                                     <.0001
Population75
NOTE: The denominator degrees of freedom for the t tests is 31.
```

## **Example 86.3: Regression Estimator for Simple Random Sample**

Using auxiliary information, you can construct the regression estimators to provide more accurate estimates of the population characteristics that are of interest. With ESTIMATE statements in PROC SURVEYREG, you can specify a regression estimator as a linear function of the regression parameters to estimate the population total. This example illustrates this application, by using the data in the previous example.

In this sample, a linear model between the Swedish populations in 1975 and in 1985 is established:

```
Population85 = \alpha + \beta * Population75 + error
```

Assuming that the total population in 1975 is known to be 8200 (in thousands), you can use the ESTIMATE statement to predict the 1985 total population by using the following statements:

```
title1 'Regression Analysis for Swedish Municipalities';
title2 'Estimate Total Population';
proc surveyreg data=Municipalities total=50;
   cluster Cluster;
   model Population85=Population75;
   estimate '1985 population' Intercept 284 Population75 8200;
run;
```

Since each observation in the sample is a municipality and there is a total of 284 municipalities in Sweden, the coefficient for Intercept ( $\alpha$ ) in the ESTIMATE statement is 284 and the coefficient for Population75 ( $\beta$ ) is the total population in 1975 (8.2 million).

Output 86.3.1 displays the regression results and the estimation of the total population. By using the linear model, you can predict the total population in 1985 to be 8.64 million, with a standard error of 0.26 million.

Output 86.3.1 Use the Regression Estimator to Estimate the Population Total

Regression Analysis for Swedish Municipalities
Estimate Total Population

The SURVEYREG Procedure

Regression Analysis for Dependent Variable Population85

Analysis of Estimable Functions

		Standard		
Parameter	Estimate	Error	t Value	Pr >  t
1985 population	8642.49485	258.558613	33.43	<.0001

NOTE: The denominator degrees of freedom for the t tests is 4.

## **Example 86.4: Stratified Sampling**

This example illustrates the use of the SURVEYREG procedure to perform a regression in a stratified sample design. Consider a population of 235 farms producing corn in Nebraska and Iowa. You are interested in the relationship between corn yield (CornYield) and total farm size (FarmArea).

Each state is divided into several regions, and each region is used as a stratum. Within each stratum, a simple random sample with replacement is drawn. A total of 19 farms is selected by using a stratified simple random sample. The sample size and population size within each stratum are displayed in Table 86.4.

Table 86.4 Number of Farms in Each Stratum

			<b>Number of Farms</b>		
Stratum	State	Region	Population	Sample	
1	Iowa	1	100	3	
2		2	50	5	
3		3	15	3	
4	Nebraska	1	30	6	
5		2	40	2	
Total			235	19	

Three models for the data are considered:

• Model I — Common intercept and slope:

Corn Yield =  $\alpha + \beta * Farm Area$ 

• Model II — Common intercept, different slope:

Corn Yield = 
$$\begin{cases} \alpha + \beta_{\text{Iowa}} * \text{Farm Area} & \text{if the farm is in Iowa} \\ \alpha + \beta_{\text{Nebraska}} * \text{Farm Area} & \text{if the farm is in Nebraska} \end{cases}$$

• Model III — Different intercept and different slope:

Corn Yield = 
$$\begin{cases} \alpha_{\text{Iowa}} + \beta_{\text{Iowa}} * \text{Farm Area} & \text{if the farm is in Iowa} \\ \alpha_{\text{Nebraska}} + \beta_{\text{Nebraska}} * \text{Farm Area} & \text{if the farm is in Nebraska} \end{cases}$$

Data from the stratified sample are saved in the SAS data set Farms. In the data set Farms, the variable Weight represents the sampling weight. In the following DATA step, the sampling weights are the reciprocals of selection probabilities:

```
data Farms;
   input State $ Region FarmArea CornYield Weight;
   datalines;
Iowa
       1 100 54 33.333
        1 83 25 33.333
Iowa
Iowa
        1 25 10 33.333
       2 120 83 10.000
Iowa
Iowa 2 50 35 10.000
Iowa 2 110 65 10.000
Iowa
       2 60 35 10.000
Iowa 2 45 20 10.000
Iowa 3 23 5 5.000
Iowa
        3 10
                8 5.000
       3 350 125 5.000
Iowa
Nebraska 1 130 20 5.000
Nebraska 1 245 25 5.000
Nebraska 1 150 33 5.000
Nebraska 1 263 50 5.000
Nebraska 1 320 47 5.000
Nebraska 1 204 25 5.000
Nebraska 2 80 11 20.000
Nebraska 2 48 8 20.000
```

The information about population size in each stratum is saved in the SAS data set StratumTotals:

```
data StratumTotals;
   input State $ Region _TOTAL_;
   datalines;
Iowa     1 100
Iowa     2 50
Iowa     3 15
Nebraska 1 30
Nebraska 2 40
;
```

Using the sample data from the data set Farms and the control information data from the data set StratumTotals, you can fit Model I by using PROC SURVEYREG with the following statements:

```
title1 'Analysis of Farm Area and Corn Yield';
title2 'Model I: Same Intercept and Slope';
proc surveyreg data=Farms total=StratumTotals;
    strata State Region / list;
    model CornYield = FarmArea / covB;
    weight Weight;
run;
```

Output 86.4.1 displays the data summary and stratification information fitting Model I. The sampling rates are automatically computed by the procedure based on the sample sizes and the population totals in strata.

Output 86.4.1 Data Summary and Stratum Information Fitting Model I

•					
	Ana	lysis of Farm Ar	ea and Co	rn Yield	
	Мо	del I: Same Inte	rcept and	Slope	
		The SURVEYREG	Procedur	e	
	Regression A	Analysis for Dep	endent Va	riable CornYiel	d
		Data Su	mmary		
	Numbe	r of Observation	-	19	
		f Weights		234.99900	
	Weight	ted Mean of Corn ted Sum of CornY	Yield	31.56029	
	Weigh	ted Sum of CornY	ield	7416.6	
		Design S	ummary		
	1	Number of Strata	ļ	5	
		Fit Stat	istics		
		R-square	0.3	882	
		Root MSE	20.6	422	
		Denominator DF	•	14	
		Stratum Inf	ormation		
Stratum				Population	Sampling
Index	State	Region	N Obs	Total	Rate
1	Iowa	1	3	100	3.00%
2		2	5	50	10.0%
3		3	3	15	20.0%
4	Nebraska	1	6	30	20.0%
5		2	2	40	5.00%

Output 86.4.2 displays tests of model effects and the estimated regression coefficients.

Output 86.4.2 Estimated Regression Coefficients and the Estimated Covariance Matrix

	Te	sts of Mod	el Effects	3		
Effect Num DF F Value Pr > F						
	Model	1	21.74	0.0004		
	Intercept	1	4.93	0.0433		
	FarmArea	1	21.74	0.0004		
NOTE: The d	enominator d	legrees of	freedom fo	or the F te	sts is 14.	
	Estimate	ed Regressi	on Coeffic	cients		
		St	andard			
Parameter	Estima	ite	Error	t Value	Pr >  t	
Intercept	11.81629	978 5.31	981027	2.22	0.0433	
FarmArea	0.21265	0.04	560949	4.66	0.0004	
NOTE: The d	enominator d	legrees of	freedom fo	or the t te	sts is 14.	
	Cox	variance of	Estimated	Ī		
		ression Co		=		
	1.09	,		-		
		Interc	ept	FarmAre	a	
In	tercept	28.300381	277 -	-0.14647153	8	
	rmArea					

Alternatively, you can assume that the linear relationship between corn yield (CornYield) and farm area (FarmArea) is different among the states (Model II). In order to analyze the data by using this model, you create auxiliary variables FarmAreaNE and FarmArealA to represent farm area in different states:

$$\begin{aligned} & \mathsf{FarmAreaNE} = \left\{ \begin{array}{ll} 0 & \text{if the farm is in Iowa} \\ & \mathsf{FarmArea} & \text{if the farm is in Nebraska} \end{array} \right. \\ & \mathsf{FarmAreaIA} = \left\{ \begin{array}{ll} \mathsf{FarmArea} & \text{if the farm is in Iowa} \\ 0 & \text{if the farm is in Nebraska} \end{array} \right. \end{aligned}$$

The following statements create these variables in a new data set called FarmsByState and use PROC SURVEYREG to fit Model II:

```
data FarmsByState; set Farms;
  if State='Iowa' then do;
    FarmAreaIA=FarmArea; FarmAreaNE=0; end;
  else do;
    FarmAreaIA=0 ; FarmAreaNE=FarmArea; end;
run;
```

The following statements perform the regression by using the new data set FarmsByState. The analysis uses the auxiliary variables FarmArealA and FarmAreaNE as the regressors:

```
title1 'Analysis of Farm Area and Corn Yield';
title2 'Model II: Same Intercept, Different Slopes';
proc SURVEYREG data=FarmsByState total=StratumTotals;
    strata State Region;
    model CornYield = FarmAreaIA FarmAreaNE / covB;
    weight Weight;
run;
```

Output 86.4.3 displays the fit statistics and parameter estimates. The estimated slope parameters for each state are quite different from the estimated slope in Model I. The results from the regression show that Model II fits these data better than Model I.

Output 86.4.3 Regression Results from Fitting Model II

Analysis of Farm Area and Corn Yield Model II: Same Intercept, Different Slopes

The SURVEYREG Procedure

Regression Analysis for Dependent Variable CornYield

Fit Statistics

R-square 0.8158
Root MSE 11.6759
Denominator DF 14

#### Estimated Regression Coefficients

		Standard		
Parameter	Estimate	Error	t Value	Pr >  t
Intercept	4.04234816	3.80934848	1.06	0.3066
FarmAreaIA	0.41696069	0.05971129	6.98	<.0001
FarmAreaNE	0.12851012	0.02495495	5.15	0.0001

NOTE: The denominator degrees of freedom for the t tests is 14.

#### Covariance of Estimated Regression Coefficients

	Intercept	FarmAreaIA	FarmAreaNE
Intercept	14.511135861	-0.118001232	-0.079908772
FarmAreaIA	-0.118001232	0.0035654381	0.0006501109
FarmAreaNE	-0.079908772	0.0006501109	0.0006227496

For Model III, different intercepts are used for the linear relationship in two states. The following statements illustrate the use of the NOINT option in the MODEL statement associated with the CLASS statement to fit Model III:

```
title2 'Model III: Different Intercepts and Slopes';
proc SURVEYREG data=FarmsByState total=StratumTotals;
    strata State Region;
    class State;
    model CornYield = State FarmAreaIA FarmAreaNE / noint covB solution;
    weight Weight;
run;
```

The model statement includes the classification effect State as a regressor. Therefore, the parameter estimates for effect State present the intercepts in two states.

Output 86.4.4 displays the regression results for fitting Model III, including parameter estimates, and covariance matrix of the regression coefficients. The estimated covariance matrix shows a lack of correlation between the regression coefficients from different states. This suggests that Model III might be the best choice for building a model for farm area and corn yield in these two states.

However, some statistics remain the same under different regression models—for example, Weighted Mean of CornYield. These estimators do not rely on the particular model you use.

Output 86.4.4 Regression Results for Fitting Model III

Analysis of Farm Area and Corn Yield	
Model III: Different Intercepts and Slopes	

#### The SURVEYREG Procedure

#### Regression Analysis for Dependent Variable CornYield

#### Fit Statistics

R-square	0.9300
Root MSE	11.9810
Denominator DF	14

#### Estimated Regression Coefficients

		Standard		
Parameter	Estimate	Error	t Value	Pr >  t
State Iowa	5.27797099	5.27170400	1.00	0.3337
State Nebraska	0.65275201	1.70031616	0.38	0.7068
FarmAreaIA	0.40680971	0.06458426	6.30	<.0001
FarmAreaNE	0.14630563	0.01997085	7.33	<.0001

 $\ensuremath{\mathsf{NOTE}}\xspace$  . The denominator degrees of freedom for the t tests is 14.

#### Covariance of Estimated Regression Coefficients

	State Iowa	State Nebraska	FarmAreaIA	FarmAreaNE
State Iowa	27.790863033	0	-0.205517205	0
State Nebraska	0	2.8910750385	0	-0.027354011
FarmAreaIA	-0.205517205	0	0.0041711265	0
FarmAreaNE	0	-0.027354011	0	0.0003988349

## **Example 86.5: Regression Estimator for Stratified Sample**

This example uses the corn yield data from the previous example to illustrate how to construct a regression estimator for a stratified sample design.

As in Example 86.3, by incorporating auxiliary information into a regression estimator, the procedure can produce more accurate estimates of the population characteristics that are of interest. In this example, the sample design is a stratified sample design. The auxiliary information is the total farm areas in regions of each state, as displayed in Table 86.5. You want to estimate the total corn yield by using this information under the three linear models given in Example 86.4.

	<b>Number of Farms</b>						
Stratum	State	Region	Population	Sample	Total Farm Area		
1	Iowa	1	100	3			
2		2	50	5	13,200		
3		3	15	3			
4	Nebraska	1	30	6	8,750		
5		2	40	2			
Total			235	19	21,950		

**Table 86.5** Information for Each Stratum

The regression estimator to estimate the total corn yield under Model I can be obtained by using PROC SURVEYREG with an ESTIMATE statement:

To apply the constraint in each stratum that the weighted total number of farms equals to the total number of farms in the stratum, you can include the strata as an effect in the MODEL statement, effect State\*Region. Thus, the CLASS statement must list the STRATA variables, State and Region, as classification variables. The following ESTIMATE statement specifies the regression estimator, which is a linear function of the regression parameters:

```
estimate 'Estimate of CornYield under Model I'
INTERCEPT 235 FarmArea 21950
State*Region 100 50 15 30 40 /e;
```

This linear function contains the total for each explanatory variable in the model. Because the sampling units are farms in this example, the coefficient for Intercept in the ESTIMATE statement is

the total number of farms (235); the coefficient for FarmArea is the total farm area listed in Table 86.5 (21950); and the coefficients for effect State\*Region are the total number of farms in each strata (as displayed in Table 86.5).

Output 86.5.1 displays the results of the ESTIMATE statement. The regression estimator for the total of CornYield in Iowa and Nebraska is 7464 under Model I, with a standard error of 927.

Output 86.5.1 Regression Estimator for the Total of CornYield under Model I

```
Estimate Corn Yield from Farm Size
Model I: Same Intercept and Slope

The SURVEYREG Procedure

Regression Analysis for Dependent Variable CornYield

Analysis of Estimable Functions

Standard

Parameter

Estimate

Error t Value Pr > |t|

Estimate of CornYield under Model I 7463.52329 926.841541 8.05 <.0001

NOTE: The denominator degrees of freedom for the t tests is 14.
```

Under Model II, a regression estimator for totals can be obtained by using the following statements:

In this model, you also need to include strata as a fixed effect in the MODEL statement. Other regressors are the auxiliary variables FarmArealA and FarmAreaNE (defined in Example 86.4). In the following ESTIMATE statement, the coefficient for Intercept is still the total number of farms; and the coefficients for FarmArealA and FarmAreaNE are the total farm area in Iowa and Nebraska, respectively, as displayed in Table 86.5. The total number of farms in each strata are the coefficients for the strata effect:

```
estimate 'Total of CornYield under Model II'

INTERCEPT 235 FarmAreaIA 13200 FarmAreaNE 8750

State*Region 100 50 15 30 40 /e;
```

Output 86.5.2 displays that the results of the regression estimator for the total of corn yield in two states under Model II is 7580 with a standard error of 859. The regression estimator under Model II has a slightly smaller standard error than under Model I.

Output 86.5.2 Regression Estimator for the Total of CornYield under Model II

```
Estimate Corn Yield from Farm Size
Model II: Same Intercept, Different Slopes

The SURVEYREG Procedure

Regression Analysis for Dependent Variable CornYield

Analysis of Estimable Functions

Standard

Parameter

Estimate

Estimate

Error t Value

Pr > |t|

Total of CornYield under Model II 7580.48657 859.180439 8.82 <.0001

NOTE: The denominator degrees of freedom for the t tests is 14.
```

Finally, you can apply Model III to the data and estimate the total corn yield. Under Model III, you can also obtain the regression estimators for the total corn yield for each state. Three ESTIMATE statements are used in the following statements to create the three regression estimators:

```
title1 'Estimate Corn Yield from Farm Size';
title2 'Model III: Different Intercepts and Slopes';
proc SURVEYREG data=FarmsByState total=StratumTotals;
   strata State Region;
  class State Region;
  model CornYield = state FarmAreaIA FarmAreaNE
      State*Region /noint solution;
   weight Weight;
   estimate 'Total CornYield in Iowa under Model III'
             State 165 0 FarmAreaIA 13200 FarmAreaNE 0
             State*region 100 50 15 0 0 /e;
   estimate 'Total CornYield in Nebraska under Model III'
             State 0 70 FarmAreaIA 0 FarmAreaNE 8750
             State * Region 0 0 0 30 40 /e;
   estimate 'Total CornYield in both states under Model III'
             State 165 70 FarmAreaIA 13200 FarmAreaNE 8750
             State*Region 100 50 15 30 40 /e;
run;
```

The fixed effect State is added to the MODEL statement to obtain different intercepts in different states, by using the NOINT option. Among the ESTIMATE statements, the coefficients for explanatory variables are different depending on which regression estimator is estimated. For example, in the ESTIMATE statement

```
estimate 'Total CornYield in Iowa under Model III'

State 165 0 FarmAreaIA 13200 FarmAreaNE 0

State*region 100 50 15 0 0 /e;
```

the coefficients for the effect State are 165 and 0, respectively. This indicates that the total number of farms in Iowa is 165 and the total number of farms in Nebraska is 0, because the estimation is the total corn yield in Iowa only. Similarly, the total numbers of farms in three regions in Iowa are

used for the coefficients of the strata effect State\*Region, as displayed in Table 86.5.

Output 86.5.3 displays the results from the three regression estimators by using Model III. Since the estimations are independent in each state, the total corn yield from both states is equal to the sum of the estimated total of corn yield in Iowa and Nebraska, 6246 + 1334 = 7580. This regression estimator is the same as the one under Model II. The variance of regression estimator of the total corn yield in both states is the sum of variances of regression estimators for total corn yield in each state. Therefore, it is not necessary to use Model III to obtain the regression estimator for the total corn yield unless you need to estimate the total corn yield for each individual state.

Output 86.5.3 Regression Estimator for the Total of CornYield under Model III

Estimate Corn Yield from	Farm Size		
Model III: Different Interce	pts and Slop	es	
The SURVEYREG Proce	dure		
Regression Analysis for Dependent	Variable Co	rnYield	
Analysis of Estimable F	unctions		
		Standard	
Parameter	Estimate	Error	t Value
Total CornYield in Iowa under Model III	6246.10697	851.272372	7.34
Total CornYield in Nebraska under Model III	1334.37961	116.302948	11.47
Total CornYield in both states under Model III	7580.48657	859.180439	8.82
Analysis of Estimable F	unctions		
Parameter		Pr >  t	
Total CornYield in Iowa under Model	III	<.0001	
Total CornYield in Nebraska under M	odel III	<.0001	
Total CornYield in both states unde	r Model III	<.0001	
NOTE: The denominator degrees of freedo	m for the t	tests is 14.	

# **Example 86.6: Stratum Collapse**

In a stratified sample, it is possible that some strata might have only one sampling unit. When this happens, PROC SURVEYREG collapses the strata that contain a single sampling unit into a pooled stratum. For more detailed information about stratum collapse, see the section "Stratum Collapse" on page 6564.

Suppose that you have the following data:

```
data Sample;

input Stratum X Y W;

datalines;

10 0 0 5

10 1 1 5

11 1 1 10

11 1 2 10

12 3 3 16

33 4 4 45

14 6 7 50

12 3 4 16

.
```

The variable Stratum is again the stratification variable, the variable X is the independent variable, and the variable Y is the dependent variable. You want to regress Y on X. In the data set Sample, both Stratum=33 and Stratum=14 contain one observation. By default, PROC SURVEYREG collapses these strata into one pooled stratum in the regression analysis.

To input the finite population correction information, you create the SAS data set StratumTotals:

```
data StratumTotals;
    input Stratum _TOTAL_;
    datalines;

10  10
11  20
12  32
33  40
33  45
14  50
15    .
66  70
;
```

The variable Stratum is the stratification variable, and the variable \_TOTAL\_ contains the stratum totals. The data set StratumTotals contains more strata than the data set Sample. Also in the data set StratumTotals, more than one observation contains the stratum totals for Stratum=33:

33 40 33 45

PROC SURVEYREG allows this type of input. The procedure simply ignores strata that are not present in the data set Sample; for the multiple entries of a stratum, the procedure uses the first observation. In this example, Stratum=33 has the stratum total TOTAL =40.

The following SAS statements perform the regression analysis:

```
title1 'Stratified Sample with Single Sampling Unit in Strata';
title2 'With Stratum Collapse';
proc SURVEYREG data=Sample total=StratumTotals;
    strata Stratum/list;
    model Y=X;
    weight W;
run;
```

Output 86.6.1 shows that there are a total of five strata in the input data set and two strata are collapsed into a pooled stratum. The denominator degrees of freedom is 4, due to the collapse (see the section "Denominator Degrees of Freedom" on page 6561).

Output 86.6.1 Summary of Data and Regression

<u>-</u>	gle Sampling Unit in Strata
With Strat	um Collapse
The SURVEYR	EG Procedure
Regression Analysis f	or Dependent Variable Y
Data	Summary
Number of Observa	
Sum of Weights	157.00000
	Y 4.31210
Weighted Sum of Y	677.00000
Design	Summary
Number of Strata	5
Number of Strata Co	llapsed 2
Fit St	atistics
R-square	0.9564
	0.5111
Denominator	DF 4

Output 86.6.2 displays the stratification information, including stratum collapse. Under the column Collapsed, the fourth stratum (Stratum=14) and the fifth (Stratum=33) are marked as 'Yes,' which indicates that these two strata are collapsed into the pooled stratum (Stratum Index=0). The sampling rate for the pooled stratum is 2% (see the section "Sampling Rate of the Pooled Stratum from Collapse" on page 6564).

Output 86.6.2 Stratification Information

Stratum Information					
Stratum				Population	Sampling
Index	Collapsed	Stratum	N Obs	Total	Rate
1		10	2	10	20.0%
2		11	2	20	10.0%
3		12	2	32	6.25%
4	Yes	14	1	50	2.00%
5	Yes	33	1	40	2.50%
0	Pooled		2	90	2.22%

NOTE: Strata with only one observation are collapsed into the stratum with Stratum Index "0".

Output 86.6.3 displays the parameter estimates and the tests of the significance of the model effects.

Output 86.6.3 Parameter Estimates and Effect Tests

	Effect	Num DF	F Value	Pr > F	
	Model	1	173.01	0.0002	
	Intercept	1	0.00	0.9961	
	X	1	173.01	0.0002	
NOTE: The	denominator d	-	freedom for coeffic		tests is 4
NOTE: The		d Regressi			tests is 4
NOTE: The	Estimate	d Regressi St	on Coeffic	ients	tests is 4
Parameter	Estimate	d Regressi St te	on Coeffic andard Error	ients t Value	Pr >  t

Alternatively, if you prefer not to collapse strata with a single sampling unit, you can specify the NOCOLLAPSE option in the STRATA statement:

```
title1 'Stratified Sample with Single Sampling Unit in Strata';
title2 'Without Stratum Collapse';
proc SURVEYREG data=Sample total=StratumTotals;
    strata Stratum/list nocollapse;
model Y = X;
weight W;
run;
```

Output 86.6.4 does not contain the stratum collapse information displayed in Output 86.6.1, and the denominator degrees of freedom are 3 instead of 4.

Output 86.6.4 Summary of Data and Regression

```
Stratified Sample with Single Sampling Unit in Strata
Without Stratum Collapse

The SURVEYREG Procedure

Regression Analysis for Dependent Variable Y

Data Summary

Number of Observations 8
Sum of Weights 157.00000
Weighted Mean of Y 4.31210
Weighted Sum of Y 677.00000
```

## Output 86.6.4 continued

Design Sum	nmary
Number of Strata	5
Fit Statis	stics
R-square	0.9564
Root MSE	0.5111
Denominator DF	3

In Output 86.6.5, although the fourth stratum and the fifth stratum contain only one observation, no stratum collapse occurs.

Output 86.6.5 Stratification Information

Stratum					
Index	Stratum	N Obs	Population Total	Sampling Rate	
1	10	2	10	20.0%	
2	11	2	20	10.0%	
3	12	2	32	6.25%	
4	14	1	50	2.00%	
5	33	1	40	2.50%	
	1 2 3 4	1 10 2 11 3 12 4 14	1 10 2 2 11 2 3 12 2 4 14 1	1 10 2 10 2 11 2 20 3 12 2 32 4 14 1 50	1 10 2 10 20.0% 2 11 2 20 10.0% 3 12 2 32 6.25% 4 14 1 50 2.00%

As a result of not collapsing strata, the standard error estimates of the parameters, shown in Output 86.6.6, are different from those in Output 86.6.3, as are the tests of the significance of model effects.

Output 86.6.6 Parameter Estimates and Effect Tests

Tests of Model Effects					
Effect	Num DF	F Value	Pr > F		
Model	1	347.27	0.0003		
Intercept	1	0.00	0.9962		
x	1	347.27	0.0003		

#### Output 86.6.6 continued

Estimated Regression Coefficients							
Standard							
Parameter	Estimate	Error	t Value	Pr >  t			
Intercept	0.00179469	0.34302581	0.01	0.9962			
x	1.12598708	0.06042241	18.64	0.0003			
NOTE: The denominator degrees of freedom for the t tests is 3.							

### **Example 86.7: Domain Analysis**

Recall that in the section "Getting Started: SURVEYREG Procedure" on page 6527, you collected a stratified simple random sample from a junior high school to examine how household income and the number of children in a household affect students' average weekly spending for ice cream. You can also use the same sample to estimate the average weekly spending among male and female students, respectively. This is often called domain analysis (subgroup analysis). You can use PROC SURVEYREG to perform domain analysis as in this example. The data set follows:

```
data IceCreamDataDomain;
   input Grade Spending Income Gender$ @@;
   datalines;
7
    7
        39
            М
                 7
                     7
                        38
                             F
                                  8
                                     12
                                          47
                                              F
9
   10
        47
            М
                 7
                     1
                         34
                             М
                                  7
                                     10
                                          43
                                              М
7
    3
        44
            М
                 8
                    20
                         60
                             F
                                  8
                                     19
                                          57
                                              М
7
    2
        35
                 7
                     2
                        36
                                  9
                                     15
                                              F
                             F
                                         51
            М
                 7
                                  7
8
   16
        53
            F
                     6
                         37
                             F
                                      6
                                          41
                                              М
7
                 9
                    15
                                     17
                                              F
    6
        39
                        50
                                  8
                                         57
            М
                             M
8
   14
        46
            М
                 9
                     8
                         41
                                  9
                                      8
                                          41
                             М
9
    7
        47
            F
                 7
                     3
                        39
                             F
                                  7
                                     12
                                         50
                                              М
7
        43
            М
                 9
                    14
                         46
                             F
                                  8
                                     18
                                         58
                                              M
9
                 7
        44
            F
                     2
                        37
                             F
                                  7
                                      1
                                         37
                                              M
7
        44
                    11
                         42
                                         41
            М
                             М
                                              M
8
   10
                                  7
                                      2
        42
            M
                 8
                    13
                         46
                             F
                                          40
                                              F
9
    6
        45
            F
                 9
                         45
                                         36 F
7
        46
            F
data IceCreamDataDomain;
   set IceCreamDataDomain:
   if Grade=7 then Prob=20/1824;
   if Grade=8 then Prob=9/1025;
   if Grade=9 then Prob=11/1151;
   Weight=1/Prob;
```

In the data set IceCreamDataDomain, the variable Grade indicates a student's grade, which is the stratification variable. The variable Spending contains the dollar amount of each student's average weekly spending for ice cream. The variable Income specifies the household income, in thousands

of dollars. The variable Gender indicates a student's gender. The sampling weights are created by using the reciprocals of the probabilities of selection, as follows:

```
data StudentTotals;
   input Grade _TOTAL_;
   datalines;
7 1824
8 1025
9 1151
;
```

In the data set StudentTotals, the variable Grade is the stratification variable, and the variable \_TO-TAL\_ contains the total numbers of students in the strata in the survey population.

The following statements demonstrate how you can estimate the average spending in the subgroup of male students:

```
title1 'Ice Cream Spending Analysis';
title2 'Domain Analysis by Gender';
proc surveyreg data=IceCreamDataDomain total=StudentTotals;
    strata Grade;
    model Spending = Income;
    domain Gender;
run;
```

Output 86.7.1 gives a summary of the domains.

Output 86.7.1 Domain Analysis Summary

```
Ice Cream Spending Analysis
            Domain Analysis by Gender
             The SURVEYREG Procedure
                     Gender=F
 Domain Regression Analysis for Variable Spending
                  Domain Summary
Number of Observations
                                                 40
Number of Observations in Domain
                                                 19
Number of Observations Not in Domain
                                                 21
Mean of Spending
                                           8.94737
Sum of Spending
                                         170.00000
```

### Output 86.7.1 continued

Ice Cream Spending Analysis
Domain Analysis by Gender

The SURVEYREG Procedure

Gender=M

Domain Regression Analysis for Variable Spending

Domain Summary

Number of Observations 40
Number of Observations in Domain 21
Number of Observations Not in Domain 19
Mean of Spending 8.57143
Sum of Spending 180.00000

Output 86.7.2 shows that parameter estimates for the model within each domain.

0.737649

### Output 86.7.2 Parameter Estimates within Domain

Income

Ice Cream Spending Analysis
Domain Analysis by Gender

The SURVEYREG Procedure

Gender=F

Domain Regression Analysis for Variable Spending

Estimated Regression Coefficients

0.04973471

14.83

<.0001

NOTE: The denominator degrees of freedom for the t tests is 37.

#### Output 86.7.2 continued

```
Ice Cream Spending Analysis
                  Domain Analysis by Gender
                   The SURVEYREG Procedure
                          Gender=M
       Domain Regression Analysis for Variable Spending
              Estimated Regression Coefficients
                             Standard
                                Error t Value
 Parameter
                Estimate
                                                   Pr > |t|
 Intercept
              -23.342282 2.11458083
                                         -11.04
                                                     <.0001
                0.730052 0.04587826
 Income
                                         15.91
                                                     <.0001
NOTE: The denominator degrees of freedom for the t tests is 37.
```

### **Example 86.8: Variance Estimate Using the Jackknife Method**

This example uses the stratified sample in the section "Stratified Sampling" on page 6530 to illustrate how to estimate the variances with replication methods.

As shown in the section "Stratified Sampling," the selected sample is saved in the SAS data set lceCream. The variable Grade that indicates a student's grade is the stratification variable. The variable Spending contains the dollar amount of each student's average weekly spending for ice cream. The variable Income specifies the household income, in thousands of dollars. The variable Kids indicates how many children are in a student's family. The variable Weight contains sampling weights.

In this example, we use the jackknife method to estimate the variance, saving the replicate weights generated by the procedure into a SAS data set:

```
title1 'Ice Cream Spending Analysis';
title2 'Use the jackknife method to estimate the variance';
proc surveyreg data=IceCream
   varmethod=JACKKNIFE(outweights=JKWeights);
   strata Grade;
   class Kids;
   model Spending = Income Kids / solution;
   weight Weight;
run;
```

The option VARMETHOD=JACKKNIFE requests the procedure to estimate the variance by using the jackknife method. The option OUTWEIGHTS=JKWeights provides a SAS data set named JKWeights that contains the replicate weights used in the computation.

Output 86.8.1 shows the summary of the data and the variance estimation method. There are a total of 40 replicates generated by the procedure.

Output 86.8.1 Variance Estimation Using the Jackknife Method

Ice Cream Spending Analysis Use the jackknife method to estimate the variance The SURVEYREG Procedure Regression Analysis for Dependent Variable Spending Data Summary Number of Observations 40 Sum of Weights 4000.0 Weighted Mean of Spending 9.14130 Weighted Sum of Spending 36565.2 Design Summary Number of Strata 3 Variance Estimation Method Jackknife Number of Replicates 40

Output 86.8.2 displays the parameter estimates and their standard errors, as well as the tests of model effects that use the jackknife method.

Output 86.8.2 Variance Estimation Using the Jackknife Method

E	Effect	Num DF	F Value	Pr > F
N	Model	4	110.48	<.0001
3	Intercept	1	133.30	<.0001
3	Income	1	289.16	<.0001
F	Kids	3	0.90	0.4525

### Output 86.8.2 continued

Estimated Regression Coefficients						
		Standard				
Parameter	Estimate	Error	t Value	Pr >  t		
Intercept	-26.086882	2.58771182	-10.08	<.0001		
Income	0.776699	0.04567521	17.00	<.0001		
Kids 1	0.888631	1.12799263	0.79	0.4358		
Kids 2	1.545726	1.25598146	1.23	0.2262		
Kids 3	-0.526817	1.42555453	-0.37	0.7138		
Kids 4	0.000000	0.0000000				

NOTE: The denominator degrees of freedom for the t tests is 37. Matrix X'WX is singular and a generalized inverse was used to solve the normal equations. Estimates are not unique.

Output 86.8.3 prints the first 6 observation in the output data set JKWeights, which contains the replicate weights.

The data set JKWeights contains all the variable in the data set IceCream, in addition to the replicate weights variables named RepWt\_1, RepWt\_2, ..., RepWt\_40.

For example, the first observation (student) from stratum Grade=7 is deleted to create the first replicate. Therefore, stratum Grade=7 is the donor stratum for the first replicate, and the corresponding replicate weights are saved in the variable RepWt\_1.

Because the first observation is deleted in the first replicate, RepWt\_1=0 for the first observation. For observations from strata other than the donor stratum Grade=7, their replicate weights remain the same as in the variable Weight, while the rest of the observations in stratum Grade=7 are multiplied by the reciprocal of the corresponding jackknife coefficient, 0.95 for the first replicate.

Output 86.8.3 The Jackknife Replicate Weights for the First 6 Observations

The Jackknife Weights for the First 6 Obs									
Obs	Grade	Spending	Income K	ids Pro	ob We	ight Repl	Wt_1 Rep	Wt_2 Repl	Wt_3 RepWt_4
1	7	7	39	2 0.01	0965 91	.200 0	.000 96	.000 91	.200 91.200
2	7	7	38	1 0.01		.200 96	.000 0		.200 91.200
3	8	12	47	1 0.00	8780 113	.889 113	.889 113	.889 0	.000 113.889
4	9	10	47	4 0.00	9557 104	.636 104	.636 104	.636 104	.636 0.000
5	7	1	34	4 0.01					.200 91.200
6	7	10	43	2 0.01					.200 91.200
						RepWt_	RepWt_	RepWt_	RepWt_
Obs	RepWt_	5 RepWt_6	RepWt_7	RepWt_8	RepWt_9	10	11	12	13
	• –	• –	• –	• -	• –				
1	96.00	96.000	96.000	91.200	91.200	96.000	96.000	91.200	91.200
2	96.00			91.200	91.200	96.000	96.000	91.200	91.200
3		9 113.889							
		6 104.636							
5	0.00			91.200	91.200	96.000	96.000	91.200	91.200
6	96.00				91.200				91.200
					0-1-00				
	RepWt	_ RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_
Obs	14	_ 1.6p1.6_ 15	16	17	18	19	20	21	22
020					0				
1	96.00	96.000	96.000	91.200	91.200	91.200	91.200	91.200	91.200
2	96.00			91.200	91.200	91.200	91.200	91.200	91.200
		9 113.889							
		6 104.636							
5	96.00								91.200
6	96.00				91.200			91.200	91.200
	30.00		30.000	31.200	31.200	31.200	31.200	31.200	31.200
	RepWt	_ RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_
Obs	23	24	25	26	27	28	29	30	31
020								30	32
1	96.00	0 96.000	96.000	91.200	91.200	91.200	96.000	96.000	96.000
2	96.00			91.200	91.200	91.200	96.000	96.000	96.000
		9 113.889							
4		6 104.636							
5	96.00			91.200	91.200	91.200	96.000	96.000	96.000
6	96.00			91.200	91.200	91.200	96.000	96.000	96.000
"	30.00	0 30.000	30.000	JI.200	31.200	31.200	30.000	30.000	30.000
	RepWt	_ RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_	RepWt_
Obs	32	_ 1.cp1.c_ 33	34	35	36	37	38	39	40
223	J <u>.</u>	33	3-1	55	50	5,	30		
1	96.00	0 91.200	91.200	91.200	96.000	91.200	91.200	96.000	96.000
2	96.00					91.200			96.000
		91.200 9 113.889							
		6 115.009							
5	96.00			91.200	96.000		91.200	96.000	96.000
6	96.00							96.000	
0	90.00	0 91.200	91.200	91.200	96.000	91.200	91.200	30.000	96.000

### References

Brick, J. M. and Kalton, G. (1996), "Handling Missing Data in Survey Research," *Statistical Methods in Medical Research*, 5, 215–238.

Cochran, W. G. (1977), Sampling Techniques, Third Edition, New York: John Wiley & Sons.

Dippo, C. S., Fay, R. E., and Morganstein, D. H. (1984), "Computing Variances from Complex Samples with Replicate Weights," *Proceedings of the Survey Research Methods Section, ASA*, 489–494.

Fay, R. E. (1984), "Some Properties of Estimators of Variance Based on Replication Methods," *Proceedings of the Survey Research Methods Section, ASA*, 495–500.

Fay, R. E. (1989), "Theory and Application of Replicate Weighting for Variance Calculations," *Proceedings of the Survey Research Methods Section, ASA*, 212–217.

Fuller, W. A. (1975), "Regression Analysis for Sample Survey," Sankhyā, 37 (3), Series C, 117–132.

Fuller, W. A., Kennedy, W., Schnell, D., Sullivan, G., and Park, H. J. (1989), *PC CARP*, Ames: Statistical Laboratory, Iowa State University.

Hidiroglou, M. A., Fuller, W. A., and Hickman, R. D. (1980), *SUPER CARP*, Ames: Statistical Laboratory, Iowa State University.

Judkins, D. (1990), "Fay's Method for Variance Estimation," *Journal of Official Statistics*, 6, 223–239.

Kalton, G., and Kaspyzyk, D. (1986), "The Treatment of Missing Survey Data," *Survey Methodology*, 12, 1–16.

Kish, L. (1965), Survey Sampling, New York: John Wiley & Sons.

Lohr, S. L. (1999), Sampling: Design and Analysis, Pacific Grove, CA: Duxbury Press.

Pringle, R. M. and Rayner, A. A. (1971), *Generalized Inverse Matrices with Applications to Statistics*, New York: Hafner Publishing.

Rao, J. N. K., and Shao, J. (1996), "On Balanced Half Sample Variance Estimation in Stratified 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., Wu, C. F. J., and Yue, K. (1992), "Some Recent Work on Resampling Methods for Complex Surveys," *Survey Methodology*, 18, 209–217.

Särndal, C. E., Swensson, B., and Wretman, J. (1992), *Model Assisted Survey Sampling*, New York: Springer-Verlag.

Rust, K. (1985), "Variance Estimation for Complex Estimators in Sample Surveys," *Journal of Official Statistics*, 1, 381–397.

Wolter, K. M. (1985), Introduction to Variance Estimation, New York: Springer-Verlag.

Woodruff, R. S. (1971), "A Simple Method for Approximating the Variance of a Complicated Estimate," *Journal of the American Statistical Association*, 66, 411–414.

## Subject Index

ADJRSQ	SURVEYREG procedure, 6538, 6558
SURVEYREG procedure, 6547	Fay's BRR method
adjusted R-square	variance estimation (SURVEYREG), 6558
SURVEYREG procedure, 6563	finite population correction
alpha level	SURVEYREG procedure, 6536, 6537, 6553
SURVEYREG procedure, 6535, 6549	•
analysis of variance	Hadamard matrix
SURVEYREG procedure, 6562	SURVEYREG procedure, 6538, 6560
ANOVA	
SURVEYREG procedure, 6547, 6562	jackknife
r,,	SURVEYREG procedure, 6559
balanced repeated replication	jackknife coefficients
SURVEYREG procedure, 6557	SURVEYREG procedure, 6559, 6567
variance estimation (SURVEYREG), 6557	jackknife variance estimation
BRR	SURVEYREG procedure, 6559
SURVEYREG procedure, 6557	-
BRR variance estimation	linearization method
SURVEYREG procedure, 6557	SURVEYREG procedure, 6556
SORVETREO procedure, 0557	
classification variables	missing values
SURVEYREG procedure, 6541	SURVEYREG procedure, 6535, 6552
clustering	MSE
SURVEYREG procedure, 6542	SURVEYREG procedure, 6563
computational details	multiple R-square
-	SURVEYREG procedure, 6562
SURVEYREG procedure, 6555	•
computational resources	number of replicates
SURVEYREG procedure, 6565	SURVEYREG procedure, 6540, 6557–6559
confidence level	
SURVEYREG procedure, 6535	ODS table names
confidence limits	SURVEYREG procedure, 6573
SURVEYREG procedure, 6547	output data sets
contrasts	SURVEYREG procedure, 6566
SURVEYREG procedure, 6543, 6565	output jackknife coefficient
	SURVEYREG procedure, 6567
degrees of freedom	output replicate weights
SURVEYREG procedure, 6560	SURVEYREG procedure, 6567
design effects	output table names
SURVEYREG procedure, 6563	SURVEYREG procedure, 6573
domain analysis	BORVETRES procedure, 6575
SURVEYREG procedure, 6565, 6595	pooled stratum
donor stratum	SURVEYREG procedure, 6564
SURVEYREG procedure, 6559	primary sampling units (PSUs)
•	SURVEYREG procedure, 6554
effect testing	BORVETRES procedure, 655 i
SURVEYREG procedure, 6562	regression coefficients
estimable functions	SURVEYREG procedure, 6555
SURVEYREG procedure, 6545	regression estimator
-	SURVEYREG procedure, 6587
Fay coefficient	regression estimators

SURVEYREG procedure, 6580	data summary table, 6568
replicate weights	degrees of freedom, 6560
SURVEYREG procedure, 6550, 6556	design effects, 6563
replication methods	design summary table, 6568
SURVEYREG procedure, 6537, 6556, 6598	domain analysis, 6565, 6595
root MSE	domain summary table, 6569
SURVEYREG procedure, 6563	domain variable, 6544
•	donor stratum, 6559
sampling rates	effect testing, 6562
SURVEYREG procedure, 6536, 6553	estimable functions, 6545
sampling weights	Fay coefficient, 6538, 6558
SURVEYREG procedure, 6550, 6552	Fay's BRR variance estimation, 6558
simple random cluster sampling	finite population correction, 6536, 6537,
SURVEYREG procedure, 6577	6553
simple random sampling	first-stage sampling rate, 6536
SURVEYREG procedure, 6527, 6574	fit statistics table, 6569
singularity level	Hadamard matrix, 6538, 6560, 6573
SURVEYREG procedure, 6544, 6546, 6548	inverse matrix of X'X, 6571
stratification	jackknife, 6559
SURVEYREG procedure, 6551	jackknife coefficients, 6559, 6567
stratified sampling	jackknife variance estimation, 6559
SURVEYREG procedure, 6530, 6581	linearization method, 6556
stratum collapse	missing values, 6535, 6552
SURVEYREG procedure, 6564, 6590	MSE, 6563
subdomain analysis, see also domain analysis	
subgroup analysis, see also domain analysis	multiple R-square, 6562
subpopulation analysis, see also domain analysis	number of replicates, 6540, 6557–6559
survey sampling, see also SURVEYREG	output data sets, 6533, 6566
procedure	output jackknife coefficient, 6567
regression analysis, 6526	output replicate weights, 6567
SURVEYREG procedure, 6526	output table names, 6573
ADJRSQ, 6547	pooled stratum, 6564
adjusted R-square, 6563	population totals, 6537, 6553
	primary sampling units (PSUs), 6554
alpha level, 6535, 6549	regression coefficients, 6555
analysis of contrasts table, 6572	regression coefficients table, 6571
analysis of estimable functions table, 6572	regression estimator, 6587
analysis of variance, 6562	regression estimators, 6580
ANOVA 6547, 6562	replicate weights, 6550, 6556
ANOVA table, 6571	replication methods, 6537, 6556, 6598
balanced repeated replication, 6557	root MSE, 6563
BRR, 6557	sampling rates, 6536, 6553
BRR variance estimation, 6557	sampling weights, 6550, 6552
classification level table, 6570	simple random cluster sampling, 6577
classification variables, 6541	simple random sampling, 6527, 6574
clustering, 6542	singularity level, 6544, 6546, 6548
coefficients of contrast table, 6572	stratification, 6551
coefficients of estimate table, 6572	stratified sampling, 6530, 6581
computational details, 6555	stratum collapse, 6564, 6590
computational resources, 6565	stratum information table, 6570
confidence level, 6535	subpopulation analysis, 6595
confidence limits, 6547	Taylor series variance estimation, 6541,
contrasts, 6543, 6565	6556
covariance of estimated regression	testing effect, 6562
coefficients table, 6572	tests of model effects table, 6571

variance estimation, 6556 variance estimation table, 6569 Wald test, 6562, 6565 weighting, 6550, 6552 X'X matrix, 6570 Taylor series variance estimation SURVEYREG procedure, 6541, 6556 testing effect SURVEYREG procedure, 6562 variance estimation BRR (SURVEYREG), 6557 jackknife (SURVEYREG), 6559 SURVEYREG procedure, 6556 Taylor series (SURVEYREG), 6541, 6556 Wald test SURVEYREG procedure, 6562, 6565 weighting SURVEYREG procedure, 6550, 6552

### Syntax Index

ADJRSQ option	VARMETHOD=BRR (PROC
MODEL statement (SURVEYREG), 6547	SURVEYREG statement), 6538
ALPHA= option	HADAMARD= option
OUTPUT statement (SURVEYREG), 6549	VARMETHOD=BRR (PROC
PROC SURVEYREG statement, 6535	SURVEYREG statement), 6538
ANOVA option	Serit Errice statement), esce
MODEL statement (SURVEYREG), 6547	INVERSE option
into BBB simicano (S CIT / B TITE C), GC T/	MODEL statement (SURVEYREG), 6547
BY statement	
SURVEYREG procedure, 6541	JKCOEFS= option
r r	REPWEIGHTS statement (SURVEYREG)
CLASS statement	6550
SURVEYREG procedure, 6541	
CLPARM option	keyword= option
MODEL statement (SURVEYREG), 6547	OUTPUT statement (SURVEYREG), 6549
CLUSTER statement	
SURVEYREG procedure, 6542	LCLM keyword
CONTRAST statement	OUTPUT statement (SURVEYREG), 6549
SURVEYREG procedure, 6543	LIST option
COVB option	STRATA statement (SURVEYREG), 6552
MODEL statement (SURVEYREG), 6547	) doctor
WODEL statement (SORVETRES), 0547	MISSING option
DATA= option	PROC SURVEYREG statement, 6535
PROC SURVEYREG statement, 6535	MODEL statement
DEFF option	SURVEYREG procedure, 6547
MODEL statement (SURVEYREG), 6547	NT
DF= option	N= option
MODEL statement (SURVEYREG), 6547	PROC SURVEYREG statement, 6537
REPWEIGHTS statement (SURVEYREG),	NOCOLLAPSE option
6550	STRATA statement (SURVEYREG), 6552
	NOFILL option
DIVISOR= option  ESTIMATE statement (SURVEYBEC)	CONTRAST statement (SURVEYREG),
ESTIMATE statement (SURVEYREG),	6543
6546	ESTIMATE statement (SURVEYREG),
DOMAIN statement	6546
SURVEYREG procedure, 6544	NOINT option
Farting	MODEL statement (SURVEYREG), 6548
E option  CONTRACT -total and (SURVEY/REC)	NOMCAR option
CONTRAST statement (SURVEYREG),	PROC SURVEYREG statement, 6535
6543	
ESTIMATE statement (SURVEYREG),	ORDER= option
6546	CLASS statement (SURVEYREG), 6536
ESTIMATE statement	OUT= option
SURVEYREG procedure, 6545	OUTPUT statement (SURVEYREG), 6548
THE STATE OF THE S	OUTJKCOEFS= option
FAY= option	VARMETHOD=JK (PROC SURVEYREG
VARMETHOD=BRR (PROC	statement), 6540
SURVEYREG statement), 6538	OUTPUT statement
II antian	SURVEYREG procedure, 6548
H= option	OUTWEIGHTS= option
	<del>-</del>

VARMETHOD=BRR (PROC	SINGULAR= option, 6544
SURVEYREG statement), 6539	SURVEYREG procedure, DOMAIN statement,
VARMETHOD=JK (PROC SURVEYREG	6544
statement), 6540	SURVEYREG procedure, ESTIMATE statement
	6545
PARMLABEL option	DIVISOR= option, 6546
MODEL statement (SURVEYREG), 6548	E option, 6546
PREDICTED keyword	NOFILL option, 6546
OUTPUT statement (SURVEYREG), 6549	SINGULAR= option, 6546
PRINTH option	SURVEYREG procedure, MODEL statement,
VARMETHOD=BRR (PROC	6547
SURVEYREG statement), 6539	ADJRSQ option, 6547
PROC SURVEYREG statement, see	ANOVA option, 6547
SURVEYREG procedure	CLPARM option, 6547
D	COVB option, 6547
R= option	DEFF option, 6547
PROC SURVEYREG statement, 6536	INVERSE option, 6547
RATE= option	NOINT option, 6548
PROC SURVEYREG statement, 6536	PARMLABEL option, 6548
REPS= option	SINGULAR= option, 6548
VARMETHOD=BRR (PROC	SOLUTION option, 6548
SURVEYREG statement), 6540	VADJUST= option, 6548
REPWEIGHTS statement	XPX option, 6548
SURVEYREG procedure, 6550	SURVEYREG procedure, MODEL statement
RESIDUAL keyword	(SURVEYREG)
OUTPUT statement (SURVEYREG), 6549	DF= option, 6547
CDICITI AD	SURVEYREG procedure, OUTPUT statement,
SINGULAR= option	6548
CONTRAST statement (SURVEYREG),	ALPHA= option, 6549
6544	keyword= option, 6549
ESTIMATE statement (SURVEYREG),	LCLM keyword, 6549
6546	OUT= option, 6548
MODEL statement (SURVEYREG), 6548	PREDICTED keyword, 6549
SOLUTION option	RESIDUAL keyword, 6549
MODEL statement (SURVEYREG), 6548	STD keyword, 6549
STD keyword	STDP keyword, 6549
OUTPUT statement (SURVEYREG), 6549	UCLM keyword, 6549
STDP keyword	SURVEYREG procedure, PROC SURVEYREG
OUTPUT statement (SURVEYREG), 6549	statement, 6535
STRATA statement	
SURVEYREG procedure, 6551	ALPHA= option, 6535
SUBGROUP statement	DATA= option, 6535
SURVEYREG procedure, 6544	FAY= option (VARMETHOD=BRR), 6538
SURVEYREG procedure	H= option (VARMETHOD=BRR), 6538
syntax, 6534	HADAMARD= option
SURVEYREG procedure, BY statement, 6541	(VARMETHOD=BRR), 6538
SURVEYREG procedure, CLASS statement,	MISSING option, 6535
6541	N= option, 6537
ORDER= option, 6536	NOMCAR option, 6535
SURVEYREG procedure, CLUSTER statement,	OUTJKCOEFS= option
6542	(VARMETHOD=JK), 6540
SURVEYREG procedure, CONTRAST	OUTWEIGHTS= option
statement, 6543	(VARMETHOD=BRR), 6539
E option, 6543	OUTWEIGHTS= option
NOFILL option, 6543	(VARMETHOD=JK), 6540

```
PRINTH option (VARMETHOD=BRR),
       6539
    R= option, 6536
    RATE= option, 6536
    REPS= option (VARMETHOD=BRR),
       6540
   TOTAL= option, 6537
   TRUNCATE option, 6537
    VARMETHOD= option, 6537
SURVEYREG procedure, REPWEIGHTS
       statement, 6550
   DF= option, 6550
    JKCOEFS= option, 6550
SURVEYREG procedure, STRATA statement,
        6551
    LIST option, 6552
    NOCOLLAPSE option, 6552
SURVEYREG procedure, WEIGHT statement,
       6552
TOTAL= option
    PROC SURVEYREG statement, 6537
TRUNCATE option
    PROC SURVEYREG statement, 6537
UCLM keyword
    OUTPUT statement (SURVEYREG), 6549
VADJUST= option
   MODEL statement (SURVEYREG), 6548
VARMETHOD= option
   PROC SURVEYREG statement, 6537
WEIGHT statement
    SURVEYREG procedure, 6552
XPX option
    MODEL statement (SURVEYREG), 6548
```

### **Your Turn**

We welcome your feedback.

- If you have comments about this book, please send them to yourturn@sas.com. Include the full title and page numbers (if applicable).
- If you have comments about the software, please send them to suggest@sas.com.

# **SAS®** Publishing Delivers!

Whether you are new to the work force or an experienced professional, you need to distinguish yourself in this rapidly changing and competitive job market. SAS® Publishing provides you with a wide range of resources to help you set yourself apart. Visit us online at support.sas.com/bookstore.

### SAS® Press

Need to learn the basics? Struggling with a programming problem? You'll find the expert answers that you need in example-rich books from SAS Press. Written by experienced SAS professionals from around the world, SAS Press books deliver real-world insights on a broad range of topics for all skill levels.

### support.sas.com/saspress

### **SAS®** Documentation

To successfully implement applications using SAS software, companies in every industry and on every continent all turn to the one source for accurate, timely, and reliable information: SAS documentation. We currently produce the following types of reference documentation to improve your work experience:

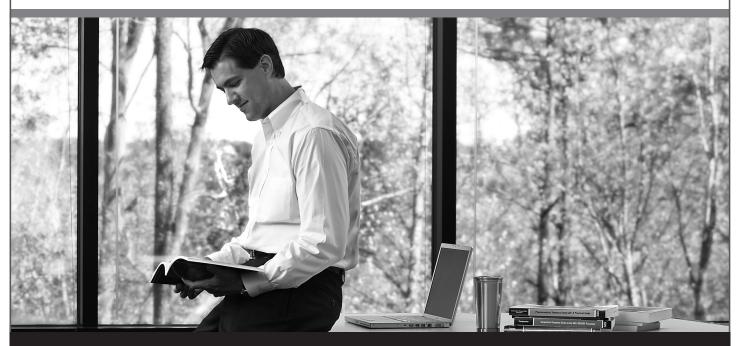
- Online help that is built into the software.
- Tutorials that are integrated into the product.
- Reference documentation delivered in HTML and PDF free on the Web.
- Hard-copy books.

### support.sas.com/publishing

### **SAS®** Publishing News

Subscribe to SAS Publishing News to receive up-to-date information about all new SAS titles, author podcasts, and new Web site features via e-mail. Complete instructions on how to subscribe, as well as access to past issues, are available at our Web site.

support.sas.com/spn



Sas THE POWER TO KNOW.