

# IBM SPSS Complex Samples 19



*Note:* Before using this information and the product it supports, read the general information under Notices on p. 267.

This document contains proprietary information of SPSS Inc, an IBM Company. It is provided under a license agreement and is protected by copyright law. The information contained in this publication does not include any product warranties, and any statements provided in this manual should not be interpreted as such.

When you send information to IBM or SPSS, you grant IBM and SPSS a nonexclusive right to use or distribute the information in any way it believes appropriate without incurring any obligation to you.

© **Copyright SPSS Inc. 1989, 2010.**

---

# Preface

IBM® SPSS® Statistics is a comprehensive system for analyzing data. The Complex Samples optional add-on module provides the additional analytic techniques described in this manual. The Complex Samples add-on module must be used with the SPSS Statistics Core system and is completely integrated into that system.

## ***About SPSS Inc., an IBM Company***

SPSS Inc., an IBM Company, is a leading global provider of predictive analytic software and solutions. The company's complete portfolio of products — data collection, statistics, modeling and deployment — captures people's attitudes and opinions, predicts outcomes of future customer interactions, and then acts on these insights by embedding analytics into business processes. SPSS Inc. solutions address interconnected business objectives across an entire organization by focusing on the convergence of analytics, IT architecture, and business processes. Commercial, government, and academic customers worldwide rely on SPSS Inc. technology as a competitive advantage in attracting, retaining, and growing customers, while reducing fraud and mitigating risk. SPSS Inc. was acquired by IBM in October 2009. For more information, visit <http://www.spss.com>.

## ***Technical support***

Technical support is available to maintenance customers. Customers may contact Technical Support for assistance in using SPSS Inc. products or for installation help for one of the supported hardware environments. To reach Technical Support, see the SPSS Inc. web site at <http://support.spss.com> or find your local office via the web site at <http://support.spss.com/default.asp?refpage=contactus.asp>. Be prepared to identify yourself, your organization, and your support agreement when requesting assistance.

## ***Customer Service***

If you have any questions concerning your shipment or account, contact your local office, listed on the Web site at <http://www.spss.com/worldwide>. Please have your serial number ready for identification.

## ***Training Seminars***

SPSS Inc. provides both public and onsite training seminars. All seminars feature hands-on workshops. Seminars will be offered in major cities on a regular basis. For more information on these seminars, contact your local office, listed on the Web site at <http://www.spss.com/worldwide>.

## ***Additional Publications***

The *SPSS Statistics: Guide to Data Analysis*, *SPSS Statistics: Statistical Procedures Companion*, and *SPSS Statistics: Advanced Statistical Procedures Companion*, written by Marija Norušis and published by Prentice Hall, are available as suggested supplemental material. These publications cover statistical procedures in the SPSS Statistics Base module, Advanced Statistics module and Regression module. Whether you are just getting starting in data analysis or are ready for advanced applications, these books will help you make best use of the capabilities found within the IBM® SPSS® Statistics offering. For additional information including publication contents and sample chapters, please see the author's website: <http://www.norusis.com>

---

# Contents

## **Part I: User's Guide**

### **1 Introduction to Complex Samples Procedures 1**

Properties of Complex Samples .....	1
Usage of Complex Samples Procedures .....	2
Plan Files .....	2
Further Readings .....	3

### **2 Sampling from a Complex Design 4**

Creating a New Sample Plan .....	4
Sampling Wizard: Design Variables .....	6
Tree Controls for Navigating the Sampling Wizard .....	7
Sampling Wizard: Sampling Method .....	8
Sampling Wizard: Sample Size .....	10
Define Unequal Sizes .....	11
Sampling Wizard: Output Variables .....	12
Sampling Wizard: Plan Summary .....	13
Sampling Wizard: Draw Sample Selection Options .....	14
Sampling Wizard: Draw Sample Output Files .....	15
Sampling Wizard: Finish .....	16
Modifying an Existing Sample Plan .....	16
Sampling Wizard: Plan Summary .....	17
Running an Existing Sample Plan .....	18
CSPLAN and CSSELECT Commands Additional Features .....	18

### **3 Preparing a Complex Sample for Analysis 19**

Creating a New Analysis Plan .....	20
Analysis Preparation Wizard: Design Variables .....	20
Tree Controls for Navigating the Analysis Wizard .....	21

Analysis Preparation Wizard: Estimation Method . . . . .	22
Analysis Preparation Wizard: Size . . . . .	23
Define Unequal Sizes. . . . .	24
Analysis Preparation Wizard: Plan Summary . . . . .	25
Analysis Preparation Wizard: Finish . . . . .	26
Modifying an Existing Analysis Plan . . . . .	26
Analysis Preparation Wizard: Plan Summary . . . . .	27
<b>4   Complex Samples Plan</b>	<b>28</b>
<b>5   Complex Samples Frequencies</b>	<b>29</b>
Complex Samples Frequencies Statistics . . . . .	30
Complex Samples Missing Values. . . . .	31
Complex Samples Options . . . . .	32
<b>6   Complex Samples Descriptives</b>	<b>33</b>
Complex Samples Descriptives Statistics . . . . .	34
Complex Samples Descriptives Missing Values. . . . .	35
Complex Samples Options . . . . .	36
<b>7   Complex Samples Crosstabs</b>	<b>37</b>
Complex Samples Crosstabs Statistics . . . . .	39
Complex Samples Missing Values. . . . .	40
Complex Samples Options . . . . .	41
<b>8   Complex Samples Ratios</b>	<b>42</b>
Complex Samples Ratios Statistics . . . . .	43
Complex Samples Ratios Missing Values . . . . .	44
Complex Samples Options . . . . .	44

**9 Complex Samples General Linear Model 45**

Complex Samples General Linear Model Statistics . . . . . 48  
Complex Samples Hypothesis Tests . . . . . 49  
Complex Samples General Linear Model Estimated Means . . . . . 50  
Complex Samples General Linear Model Save . . . . . 51  
Complex Samples General Linear Model Options . . . . . 52  
CSGLM Command Additional Features . . . . . 53

**10 Complex Samples Logistic Regression 54**

Complex Samples Logistic Regression Reference Category . . . . . 55  
Complex Samples Logistic Regression Model . . . . . 56  
Complex Samples Logistic Regression Statistics . . . . . 57  
Complex Samples Hypothesis Tests . . . . . 59  
Complex Samples Logistic Regression Odds Ratios . . . . . 60  
Complex Samples Logistic Regression Save . . . . . 61  
Complex Samples Logistic Regression Options . . . . . 62  
CSLOGISTIC Command Additional Features . . . . . 63

**11 Complex Samples Ordinal Regression 64**

Complex Samples Ordinal Regression Response Probabilities . . . . . 66  
Complex Samples Ordinal Regression Model . . . . . 66  
Complex Samples Ordinal Regression Statistics . . . . . 68  
Complex Samples Hypothesis Tests . . . . . 69  
Complex Samples Ordinal Regression Odds Ratios . . . . . 70  
Complex Samples Ordinal Regression Save . . . . . 71  
Complex Samples Ordinal Regression Options . . . . . 72  
CSORDINAL Command Additional Features . . . . . 73

**12 Complex Samples Cox Regression 74**

Define Event . . . . . 77

Predictors .....	78
Define Time-Dependent Predictor .....	79
Subgroups .....	80
Model .....	81
Statistics .....	82
Plots .....	84
Hypothesis Tests .....	85
Save .....	86
Export .....	88
Options .....	90
CSCOXREG Command Additional Features .....	91

## ***Part II: Examples***

### ***13 Complex Samples Sampling Wizard 93***

Obtaining a Sample from a Full Sampling Frame .....	93
Using the Wizard .....	93
Plan Summary .....	103
Sampling Summary .....	103
Sample Results .....	104
Obtaining a Sample from a Partial Sampling Frame .....	105
Using the Wizard to Sample from the First Partial Frame .....	105
Sample Results .....	118
Using the Wizard to Sample from the Second Partial Frame .....	118
Sample Results .....	123
Sampling with Probability Proportional to Size (PPS) .....	123
Using the Wizard .....	123
Plan Summary .....	135
Sampling Summary .....	135
Sample Results .....	137
Related Procedures .....	139

### ***14 Complex Samples Analysis Preparation Wizard 140***

Using the Complex Samples Analysis Preparation Wizard to Ready NHIS Public Data .....	140
Using the Wizard .....	140
Summary .....	143

Preparing for Analysis When Sampling Weights Are Not in the Data File . . . . .	143
Computing Inclusion Probabilities and Sampling Weights . . . . .	143
Using the Wizard . . . . .	146
Summary . . . . .	154
Related Procedures . . . . .	154
<b>15 Complex Samples Frequencies</b>	<b>155</b>
Using Complex Samples Frequencies to Analyze Nutritional Supplement Usage . . . . .	155
Running the Analysis . . . . .	155
Frequency Table . . . . .	158
Frequency by Subpopulation . . . . .	158
Summary . . . . .	159
Related Procedures . . . . .	159
<b>16 Complex Samples Descriptives</b>	<b>160</b>
Using Complex Samples Descriptives to Analyze Activity Levels . . . . .	160
Running the Analysis . . . . .	160
Univariate Statistics . . . . .	163
Univariate Statistics by Subpopulation . . . . .	163
Summary . . . . .	164
Related Procedures . . . . .	164
<b>17 Complex Samples Crosstabs</b>	<b>165</b>
Using Complex Samples Crosstabs to Measure the Relative Risk of an Event . . . . .	165
Running the Analysis . . . . .	165
Crosstabulation . . . . .	168
Risk Estimate . . . . .	169
Risk Estimate by Subpopulation . . . . .	170
Summary . . . . .	170
Related Procedures . . . . .	170

**18 Complex Samples Ratios 171**

Using Complex Samples Ratios to Aid Property Value Assessment . . . . . 171  
    Running the Analysis . . . . . 171  
    Ratios . . . . . 174  
    Pivoted Ratios Table . . . . . 174  
    Summary . . . . . 175  
Related Procedures . . . . . 175

**19 Complex Samples General Linear Model 176**

Using Complex Samples General Linear Model to Fit a Two-Factor ANOVA . . . . . 176  
    Running the Analysis . . . . . 176  
    Model Summary . . . . . 181  
    Tests of Model Effects . . . . . 181  
    Parameter Estimates . . . . . 182  
    Estimated Marginal Means . . . . . 183  
    Summary . . . . . 185  
Related Procedures . . . . . 185

**20 Complex Samples Logistic Regression 186**

Using Complex Samples Logistic Regression to Assess Credit Risk . . . . . 186  
    Running the Analysis . . . . . 186  
    Pseudo R-Squares . . . . . 190  
    Classification . . . . . 191  
    Tests of Model Effects . . . . . 191  
    Parameter Estimates . . . . . 192  
    Odds Ratios . . . . . 193  
    Summary . . . . . 194  
Related Procedures . . . . . 194

**21 Complex Samples Ordinal Regression 195**

Using Complex Samples Ordinal Regression to Analyze Survey Results . . . . . 195  
    Running the Analysis . . . . . 195  
    Pseudo R-Squares . . . . . 200  
    Tests of Model Effects . . . . . 200

Parameter Estimates . . . . .	201
Classification. . . . .	202
Odds Ratios . . . . .	203
Generalized Cumulative Model . . . . .	204
Dropping Non-Significant Predictors . . . . .	205
Warnings. . . . .	207
Comparing Models . . . . .	208
Summary . . . . .	209
Related Procedures . . . . .	209

## **22 Complex Samples Cox Regression 210**

Using a Time-Dependent Predictor in Complex Samples Cox Regression. . . . .	210
Preparing the Data . . . . .	210
Running the Analysis . . . . .	216
Sample Design Information . . . . .	221
Tests of Model Effects . . . . .	222
Test of Proportional Hazards . . . . .	222
Adding a Time-Dependent Predictor . . . . .	222
Multiple Cases per Subject in Complex Samples Cox Regression . . . . .	226
Preparing the Data for Analysis . . . . .	227
Creating a Simple Random Sampling Analysis Plan . . . . .	242
Running the Analysis . . . . .	246
Sample Design Information . . . . .	254
Tests of Model Effects . . . . .	255
Parameter Estimates . . . . .	255
Pattern Values . . . . .	256
Log-Minus-Log Plot . . . . .	257
Summary . . . . .	257

***Appendices***

***A Sample Files*** ***258***

***B Notices*** ***267***

***Bibliography*** ***269***

***Index*** ***271***

***Part I:  
User's Guide***



---

# ***Introduction to Complex Samples Procedures***

An inherent assumption of analytical procedures in traditional software packages is that the observations in a data file represent a simple random sample from the population of interest. This assumption is untenable for an increasing number of companies and researchers who find it both cost-effective and convenient to obtain samples in a more structured way.

The Complex Samples option allows you to select a sample according to a complex design and incorporate the design specifications into the data analysis, thus ensuring that your results are valid.

## ***Properties of Complex Samples***

A complex sample can differ from a simple random sample in many ways. In a simple random sample, individual sampling units are selected at random with equal probability and without replacement (WOR) directly from the entire population. By contrast, a given complex sample can have some or all of the following features:

**Stratification.** Stratified sampling involves selecting samples independently within non-overlapping subgroups of the population, or strata. For example, strata may be socioeconomic groups, job categories, age groups, or ethnic groups. With stratification, you can ensure adequate sample sizes for subgroups of interest, improve the precision of overall estimates, and use different sampling methods from stratum to stratum.

**Clustering.** Cluster sampling involves the selection of groups of sampling units, or clusters. For example, clusters may be schools, hospitals, or geographical areas, and sampling units may be students, patients, or citizens. Clustering is common in multistage designs and area (geographic) samples.

**Multiple stages.** In multistage sampling, you select a first-stage sample based on clusters. Then you create a second-stage sample by drawing subsamples from the selected clusters. If the second-stage sample is based on subclusters, you can then add a third stage to the sample. For example, in the first stage of a survey, a sample of cities could be drawn. Then, from the selected cities, households could be sampled. Finally, from the selected households, individuals could be polled. The Sampling and Analysis Preparation wizards allow you to specify three stages in a design.

**Nonrandom sampling.** When selection at random is difficult to obtain, units can be sampled systematically (at a fixed interval) or sequentially.

**Unequal selection probabilities.** When sampling clusters that contain unequal numbers of units, you can use probability-proportional-to-size (PPS) sampling to make a cluster's selection probability equal to the proportion of units it contains. PPS sampling can also use more general weighting schemes to select units.

**Unrestricted sampling.** Unrestricted sampling selects units with replacement (WR). Thus, an individual unit can be selected for the sample more than once.

**Sampling weights.** Sampling weights are automatically computed while drawing a complex sample and ideally correspond to the "frequency" that each sampling unit represents in the target population. Therefore, the sum of the weights over the sample should estimate the population size. Complex Samples analysis procedures require sampling weights in order to properly analyze a complex sample. Note that these weights should be used entirely within the Complex Samples option and should not be used with other analytical procedures via the Weight Cases procedure, which treats weights as case replications.

## ***Usage of Complex Samples Procedures***

Your usage of Complex Samples procedures depends on your particular needs. The primary types of users are those who:

- Plan and carry out surveys according to complex designs, possibly analyzing the sample later. The primary tool for surveyors is the [Sampling Wizard](#).
- Analyze sample data files previously obtained according to complex designs. Before using the Complex Samples analysis procedures, you may need to use the [Analysis Preparation Wizard](#).

Regardless of which type of user you are, you need to supply design information to Complex Samples procedures. This information is stored in a **plan file** for easy reuse.

### ***Plan Files***

A plan file contains complex sample specifications. There are two types of plan files:

**Sampling plan.** The specifications given in the Sampling Wizard define a sample design that is used to draw a complex sample. The sampling plan file contains those specifications. The sampling plan file also contains a default analysis plan that uses estimation methods suitable for the specified sample design.

**Analysis plan.** This plan file contains information needed by Complex Samples analysis procedures to properly compute variance estimates for a complex sample. The plan includes the sample structure, estimation methods for each stage, and references to required variables, such as sample weights. The Analysis Preparation Wizard allows you to create and edit analysis plans.

There are several advantages to saving your specifications in a plan file, including:

- A surveyor can specify the first stage of a multistage sampling plan and draw first-stage units now, collect information on sampling units for the second stage, and then modify the sampling plan to include the second stage.

- An analyst who doesn't have access to the sampling plan file can specify an analysis plan and refer to that plan from each Complex Samples analysis procedure.
- A designer of large-scale public use samples can publish the sampling plan file, which simplifies the instructions for analysts and avoids the need for each analyst to specify his or her own analysis plans.

## ***Further Readings***

For more information on sampling techniques, see the following texts:

Cochran, W. G. 1977. *Sampling Techniques*, 3rd ed. New York: John Wiley and Sons.

Kish, L. 1965. *Survey Sampling*. New York: John Wiley and Sons.

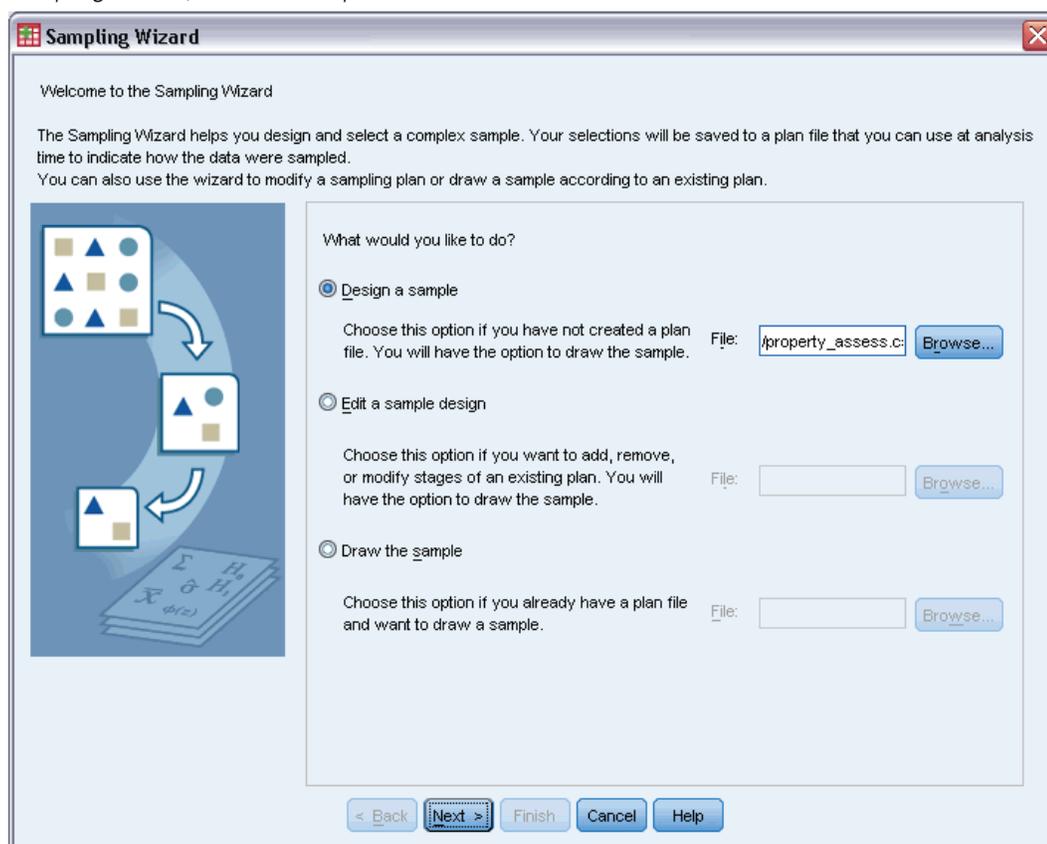
Kish, L. 1987. *Statistical Design for Research*. New York: John Wiley and Sons.

Murthy, M. N. 1967. *Sampling Theory and Methods*. Calcutta, India: Statistical Publishing Society.

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

# Sampling from a Complex Design

Figure 2-1  
Sampling Wizard, Welcome step



The Sampling Wizard guides you through the steps for creating, modifying, or executing a sampling plan file. Before using the Wizard, you should have a well-defined target population, a list of sampling units, and an appropriate sample design in mind.

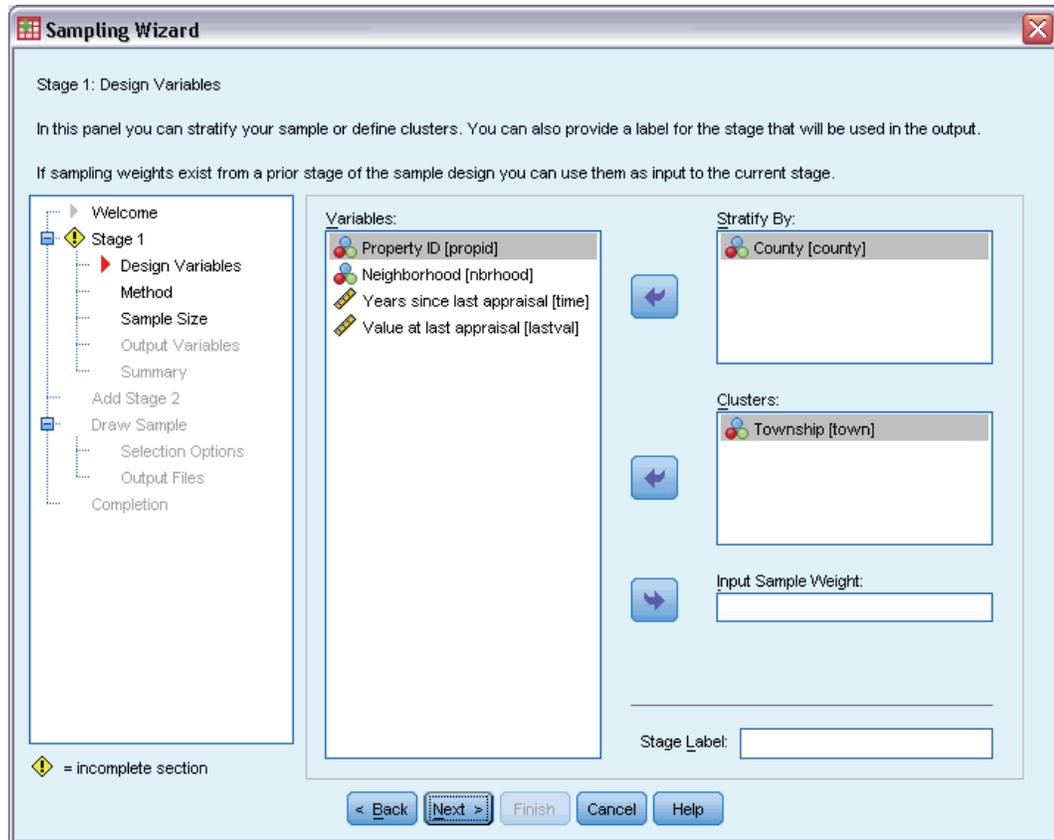
## Creating a New Sample Plan

- ▶ From the menus choose:  
Analyze > Complex Samples > Select a Sample...
- ▶ Select Design a sample and choose a plan filename to save the sample plan.

- ▶ Click **Next** to continue through the Wizard.
- ▶ Optionally, in the **Design Variables** step, you can define strata, clusters, and input sample weights. After you define these, click **Next**.
- ▶ Optionally, in the **Sampling Method** step, you can choose a method for selecting items.  
If you select **PPS Brewer** or **PPS Murthy**, you can click **Finish** to draw the sample. Otherwise, click **Next** and then:
  - ▶ In the **Sample Size** step, specify the number or proportion of units to sample.
  - ▶ You can now click **Finish** to draw the sample.  
Optionally, in further steps you can:
    - Choose output variables to save.
    - Add a second or third stage to the design.
    - Set various selection options, including which stages to draw samples from, the random number seed, and whether to treat user-missing values as valid values of design variables.
    - Choose where to save output data.
    - Paste your selections as command syntax.

## Sampling Wizard: Design Variables

Figure 2-2  
Sampling Wizard, Design Variables step



This step allows you to select stratification and clustering variables and to define input sample weights. You can also specify a label for the stage.

**Stratify By.** The cross-classification of stratification variables defines distinct subpopulations, or strata. Separate samples are obtained for each stratum. To improve the precision of your estimates, units within strata should be as homogeneous as possible for the characteristics of interest.

**Clusters.** Cluster variables define groups of observational units, or clusters. Clusters are useful when directly sampling observational units from the population is expensive or impossible; instead, you can sample clusters from the population and then sample observational units from the selected clusters. However, the use of clusters can introduce correlations among sampling units, resulting in a loss of precision. To minimize this effect, units within clusters should be as heterogeneous as possible for the characteristics of interest. You must define at least one cluster variable in order to plan a multistage design. Clusters are also necessary in the use of several different sampling methods. [For more information, see the topic Sampling Wizard: Sampling Method on p. 8.](#)

**Input Sample Weight.** If the current sample design is part of a larger sample design, you may have sample weights from a previous stage of the larger design. You can specify a numeric variable containing these weights in the first stage of the current design. Sample weights are computed automatically for subsequent stages of the current design.

**Stage Label.** You can specify an optional string label for each stage. This is used in the output to help identify stagewise information.

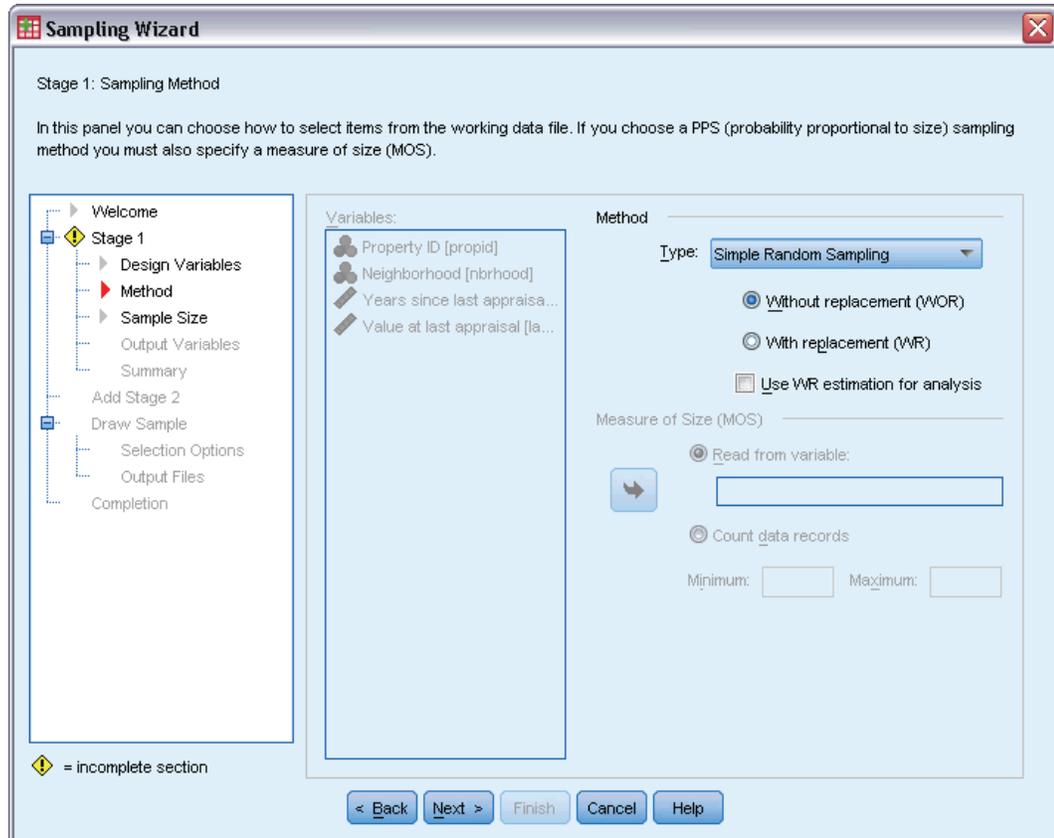
*Note:* The source variable list has the same content across steps of the Wizard. In other words, variables removed from the source list in a particular step are removed from the list in all steps. Variables returned to the source list appear in the list in all steps.

### ***Tree Controls for Navigating the Sampling Wizard***

On the left side of each step in the Sampling Wizard is an outline of all the steps. You can navigate the Wizard by clicking on the name of an enabled step in the outline. Steps are enabled as long as all previous steps are valid—that is, if each previous step has been given the minimum required specifications for that step. See the Help for individual steps for more information on why a given step may be invalid.

## Sampling Wizard: Sampling Method

Figure 2-3  
Sampling Wizard, Sampling Method step



This step allows you to specify how to select cases from the active dataset.

**Method.** Controls in this group are used to choose a selection method. Some sampling types allow you to choose whether to sample with replacement (WR) or without replacement (WOR). See the type descriptions for more information. Note that some probability-proportional-to-size (PPS) types are available only when clusters have been defined and that all PPS types are available only in the first stage of a design. Moreover, WR methods are available only in the last stage of a design.

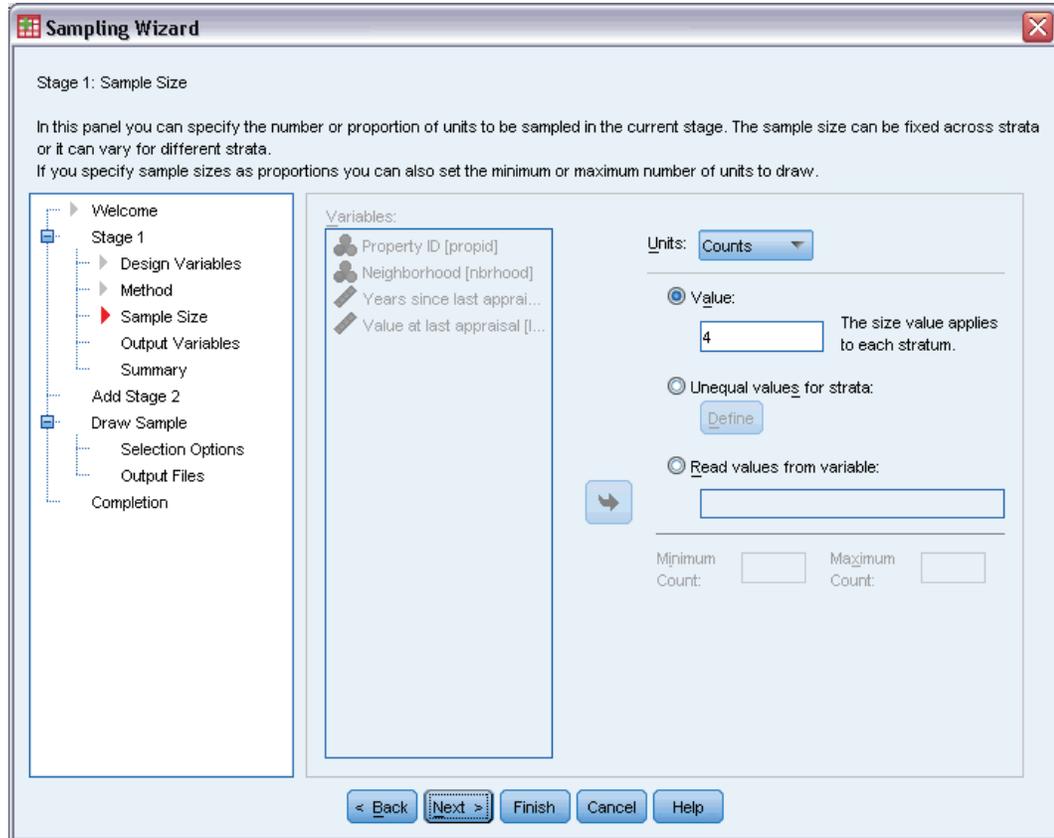
- **Simple Random Sampling.** Units are selected with equal probability. They can be selected with or without replacement.
- **Simple Systematic.** Units are selected at a fixed interval throughout the sampling frame (or strata, if they have been specified) and extracted without replacement. A randomly selected unit within the first interval is chosen as the starting point.
- **Simple Sequential.** Units are selected sequentially with equal probability and without replacement.
- **PPS.** This is a first-stage method that selects units at random with probability proportional to size. Any units can be selected with replacement; only clusters can be sampled without replacement.

- **PPS Systematic.** This is a first-stage method that systematically selects units with probability proportional to size. They are selected without replacement.
- **PPS Sequential.** This is a first-stage method that sequentially selects units with probability proportional to cluster size and without replacement.
- **PPS Brewer.** This is a first-stage method that selects two clusters from each stratum with probability proportional to cluster size and without replacement. A cluster variable must be specified to use this method.
- **PPS Murthy.** This is a first-stage method that selects two clusters from each stratum with probability proportional to cluster size and without replacement. A cluster variable must be specified to use this method.
- **PPS Sampford.** This is a first-stage method that selects more than two clusters from each stratum with probability proportional to cluster size and without replacement. It is an extension of Brewer's method. A cluster variable must be specified to use this method.
- **Use WR estimation for analysis.** By default, an estimation method is specified in the plan file that is consistent with the selected sampling method. This allows you to use with-replacement estimation even if the sampling method implies WOR estimation. This option is available only in stage 1.

**Measure of Size (MOS).** If a PPS method is selected, you must specify a measure of size that defines the size of each unit. These sizes can be explicitly defined in a variable or they can be computed from the data. Optionally, you can set lower and upper bounds on the MOS, overriding any values found in the MOS variable or computed from the data. These options are available only in stage 1.

## Sampling Wizard: Sample Size

Figure 2-4  
Sampling Wizard, Sample Size step



This step allows you to specify the number or proportion of units to sample within the current stage. The sample size can be fixed or it can vary across strata. For the purpose of specifying sample size, clusters chosen in previous stages can be used to define strata.

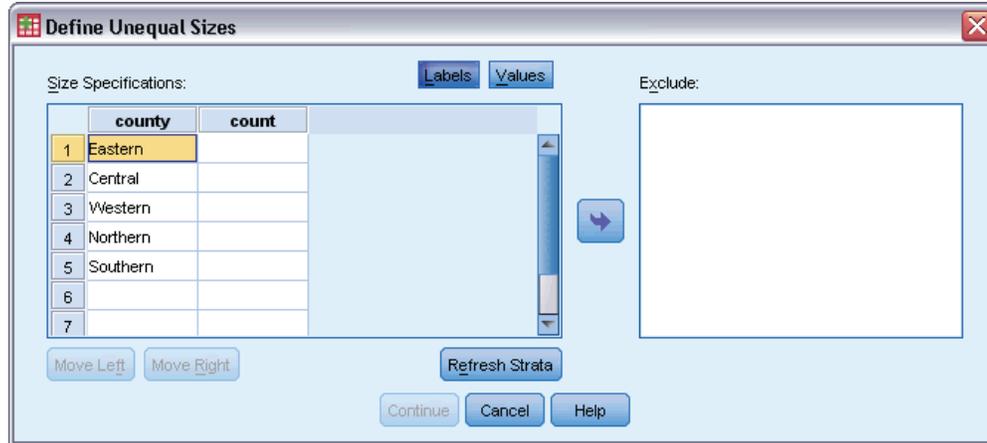
**Units.** You can specify an exact sample size or a proportion of units to sample.

- **Value.** A single value is applied to all strata. If Counts is selected as the unit metric, you should enter a positive integer. If Proportions is selected, you should enter a non-negative value. Unless sampling with replacement, proportion values should also be no greater than 1.
- **Unequal values for strata.** Allows you to enter size values on a per-stratum basis via the Define Unequal Sizes dialog box.
- **Read values from variable.** Allows you to select a numeric variable that contains size values for strata.

If Proportions is selected, you have the option to set lower and upper bounds on the number of units sampled.

## Define Unequal Sizes

Figure 2-5  
Define Unequal Sizes dialog box



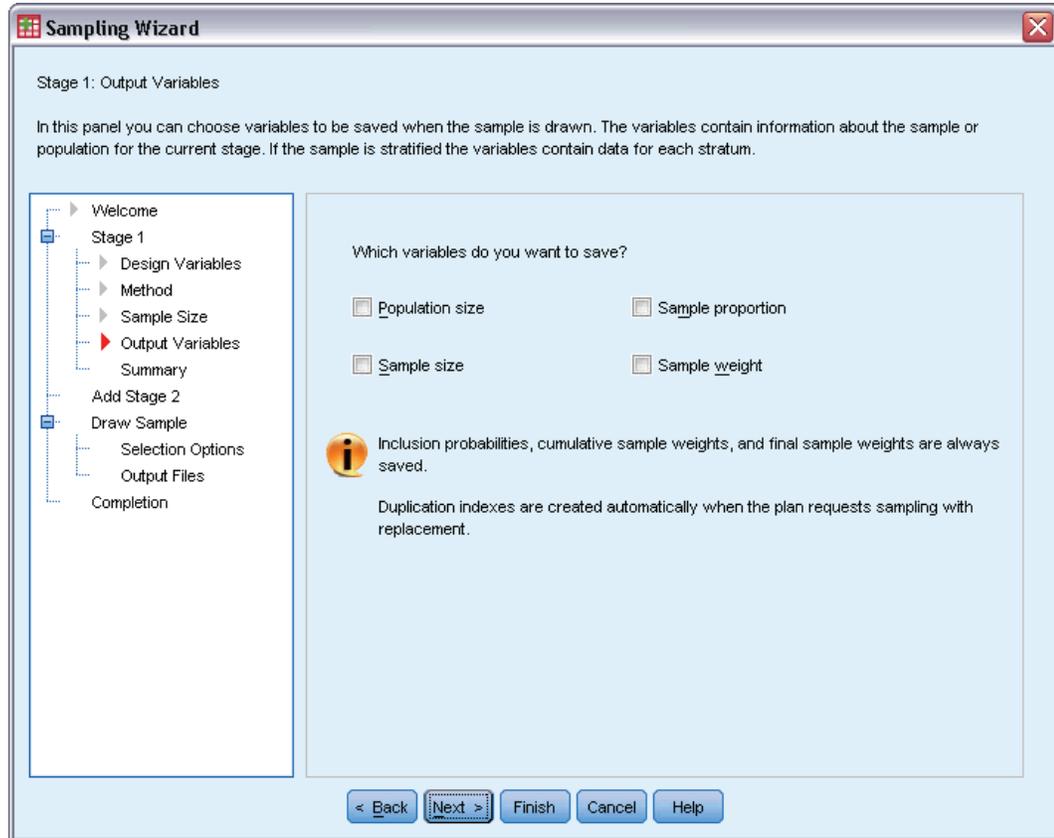
The Define Unequal Sizes dialog box allows you to enter sizes on a per-stratum basis.

**Size Specifications grid.** The grid displays the cross-classifications of up to five strata or cluster variables—one stratum/cluster combination per row. Eligible grid variables include all stratification variables from the current and previous stages and all cluster variables from previous stages. Variables can be reordered within the grid or moved to the Exclude list. Enter sizes in the rightmost column. Click Labels or Values to toggle the display of value labels and data values for stratification and cluster variables in the grid cells. Cells that contain unlabeled values always show values. Click Refresh Strata to repopulate the grid with each combination of labeled data values for variables in the grid.

**Exclude.** To specify sizes for a subset of stratum/cluster combinations, move one or more variables to the Exclude list. These variables are not used to define sample sizes.

## Sampling Wizard: Output Variables

Figure 2-6  
Sampling Wizard, Output Variables step



This step allows you to choose variables to save when the sample is drawn.

**Population size.** The estimated number of units in the population for a given stage. The rootname for the saved variable is *PopulationSize\_.*

**Sample proportion.** The sampling rate at a given stage. The rootname for the saved variable is *SamplingRate\_.*

**Sample size.** The number of units drawn at a given stage. The rootname for the saved variable is *SampleSize\_.*

**Sample weight.** The inverse of the inclusion probabilities. The rootname for the saved variable is *SampleWeight\_.*

Some stagewise variables are generated automatically. These include:

**Inclusion probabilities.** The proportion of units drawn at a given stage. The rootname for the saved variable is *InclusionProbability\_.*

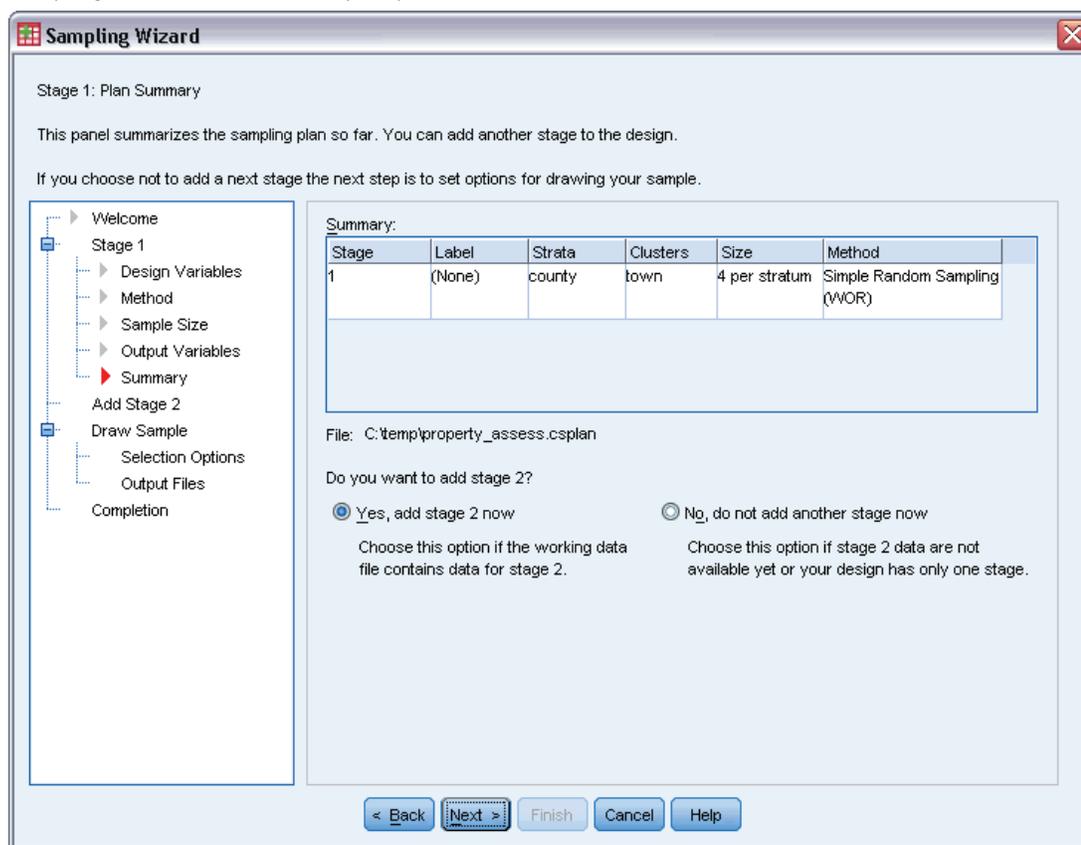
**Cumulative weight.** The cumulative sample weight over stages previous to and including the current one. The rootname for the saved variable is *SampleWeightCumulative\_.*

**Index.** Identifies units selected multiple times within a given stage. The rootname for the saved variable is *Index\_*.

*Note:* Saved variable rootnames include an integer suffix that reflects the stage number—for example, *PopulationSize\_1\_* for the saved population size for stage 1.

## Sampling Wizard: Plan Summary

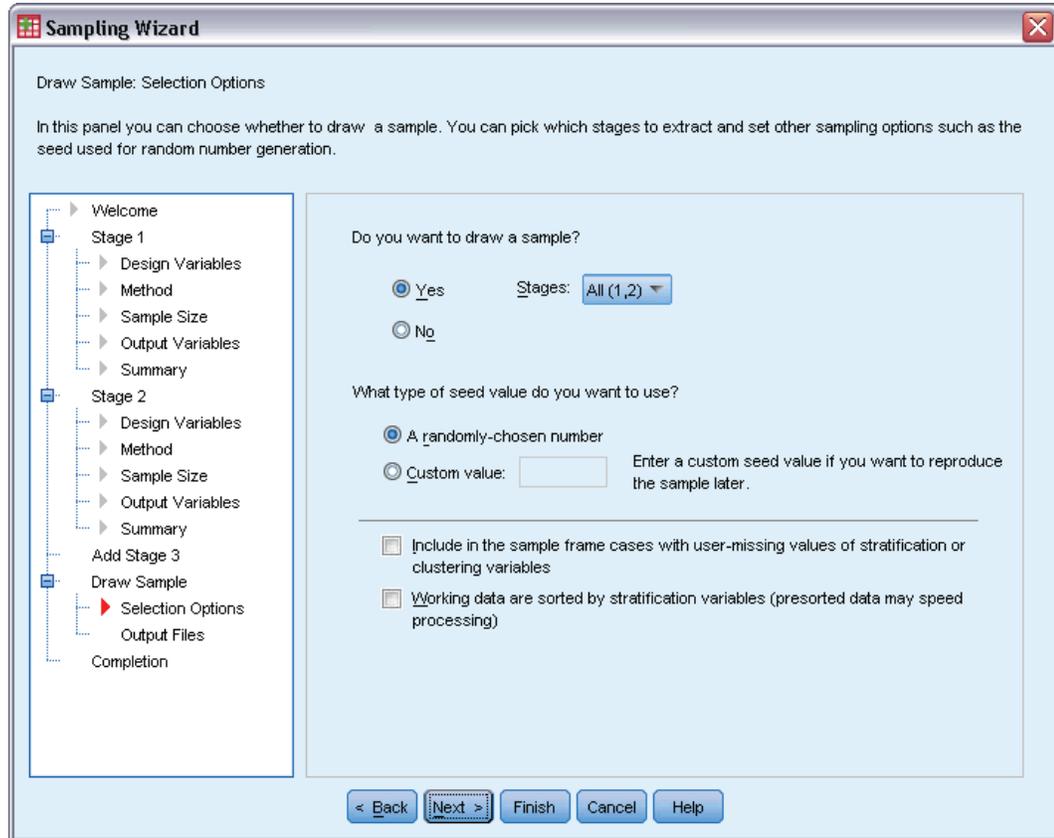
Figure 2-7  
Sampling Wizard, Plan Summary step



This is the last step within each stage, providing a summary of the sample design specifications through the current stage. From here, you can either proceed to the next stage (creating it, if necessary) or set options for drawing the sample.

## Sampling Wizard: Draw Sample Selection Options

Figure 2-8  
Sampling Wizard, Draw Sample Selection Options step



This step allows you to choose whether to draw a sample. You can also control other sampling options, such as the random seed and missing-value handling.

**Draw sample.** In addition to choosing whether to draw a sample, you can also choose to execute part of the sampling design. Stages must be drawn in order—that is, stage 2 cannot be drawn unless stage 1 is also drawn. When editing or executing a plan, you cannot resample locked stages.

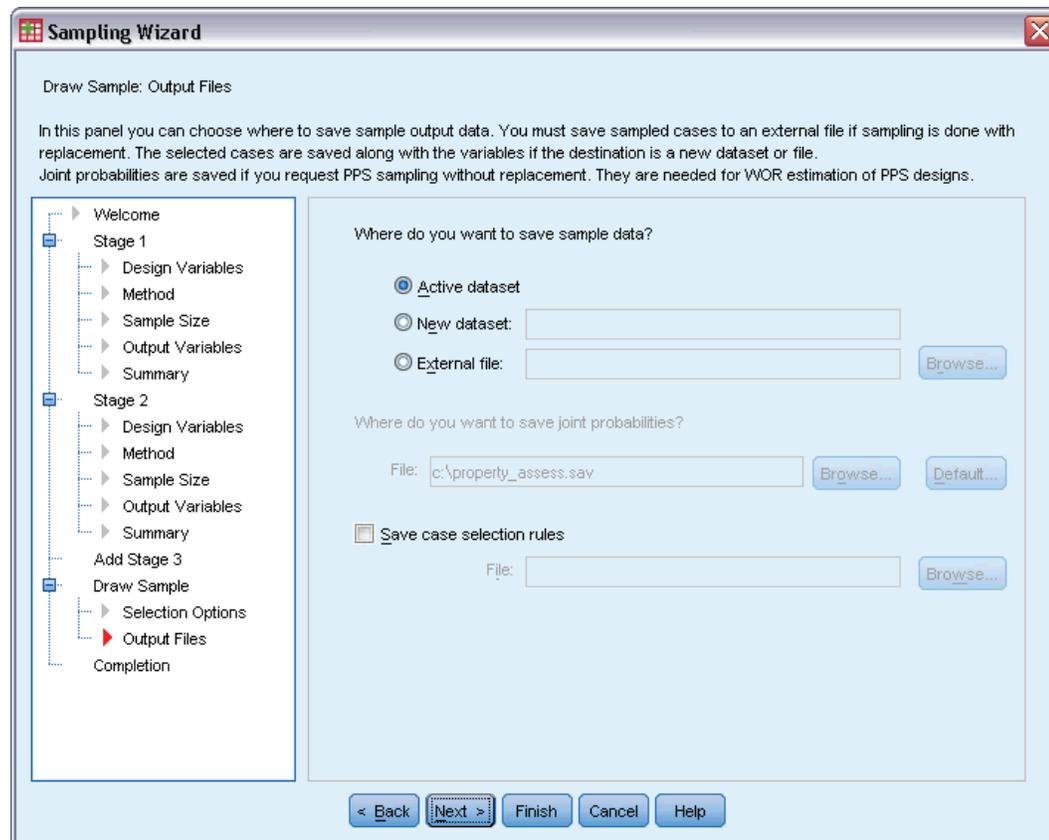
**Seed.** This allows you to choose a seed value for random number generation.

**Include user-missing values.** This determines whether user-missing values are valid. If so, user-missing values are treated as a separate category.

**Data already sorted.** If your sample frame is presorted by the values of the stratification variables, this option allows you to speed the selection process.

## Sampling Wizard: Draw Sample Output Files

Figure 2-9  
Sampling Wizard, Draw Sample Output Files step



This step allows you to choose where to direct sampled cases, weight variables, joint probabilities, and case selection rules.

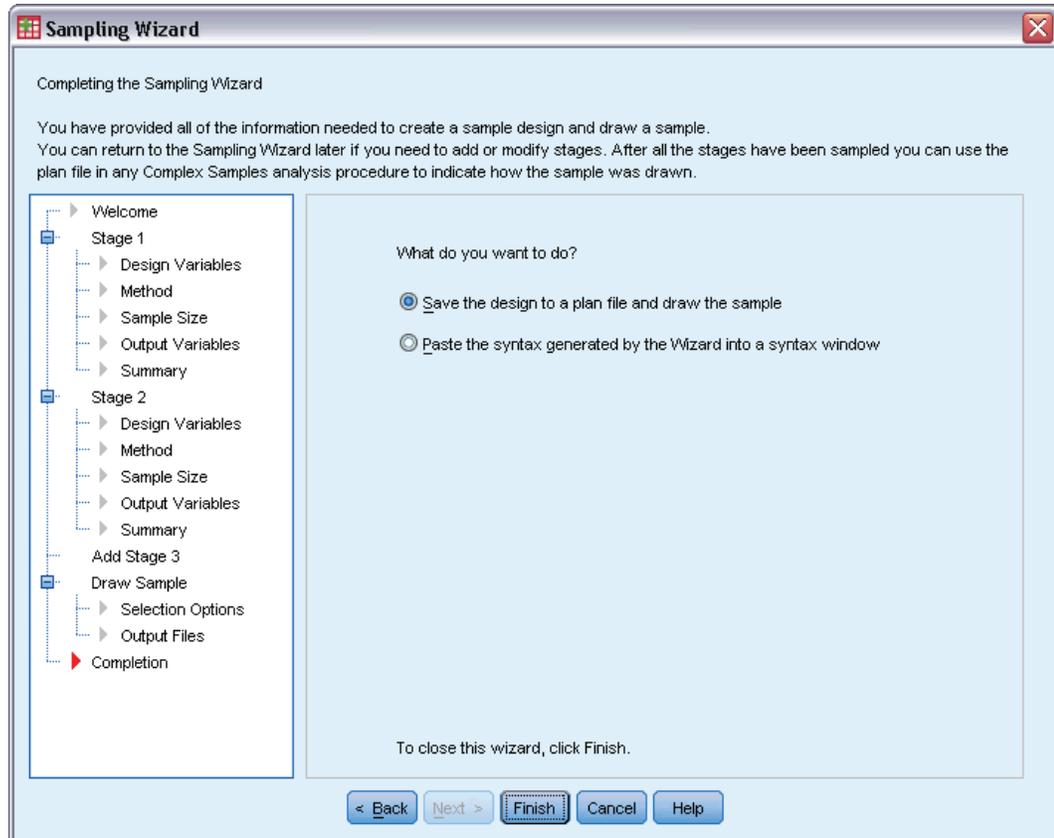
**Sample data.** These options let you determine where sample output is written. It can be added to the active dataset, written to a new dataset, or saved to an external IBM® SPSS® Statistics data file. Datasets are available during the current session but are not available in subsequent sessions unless you explicitly save them as data files. Dataset names must adhere to variable naming rules. If an external file or new dataset is specified, the sampling output variables and variables in the active dataset for the selected cases are written.

**Joint probabilities.** These options let you determine where joint probabilities are written. They are saved to an external SPSS Statistics data file. Joint probabilities are produced if the PPS WOR, PPS Brewer, PPS Sampford, or PPS Murthy method is selected and WR estimation is not specified.

**Case selection rules.** If you are constructing your sample one stage at a time, you may want to save the case selection rules to a text file. They are useful for constructing the subframe for subsequent stages.

## Sampling Wizard: Finish

Figure 2-10  
Sampling Wizard, Finish step



This is the final step. You can save the plan file and draw the sample now or paste your selections into a syntax window.

When making changes to stages in the existing plan file, you can save the edited plan to a new file or overwrite the existing file. When adding stages without making changes to existing stages, the Wizard automatically overwrites the existing plan file. If you want to save the plan to a new file, select Paste the syntax generated by the Wizard into a syntax window and change the filename in the syntax commands.

## Modifying an Existing Sample Plan

- ▶ From the menus choose:  
Analyze > Complex Samples > Select a Sample...
- ▶ Select Edit a sample design and choose a plan file to edit.
- ▶ Click Next to continue through the Wizard.

- ▶ Review the sampling plan in the Plan Summary step, and then click Next.

Subsequent steps are largely the same as for a new design. See the Help for individual steps for more information.

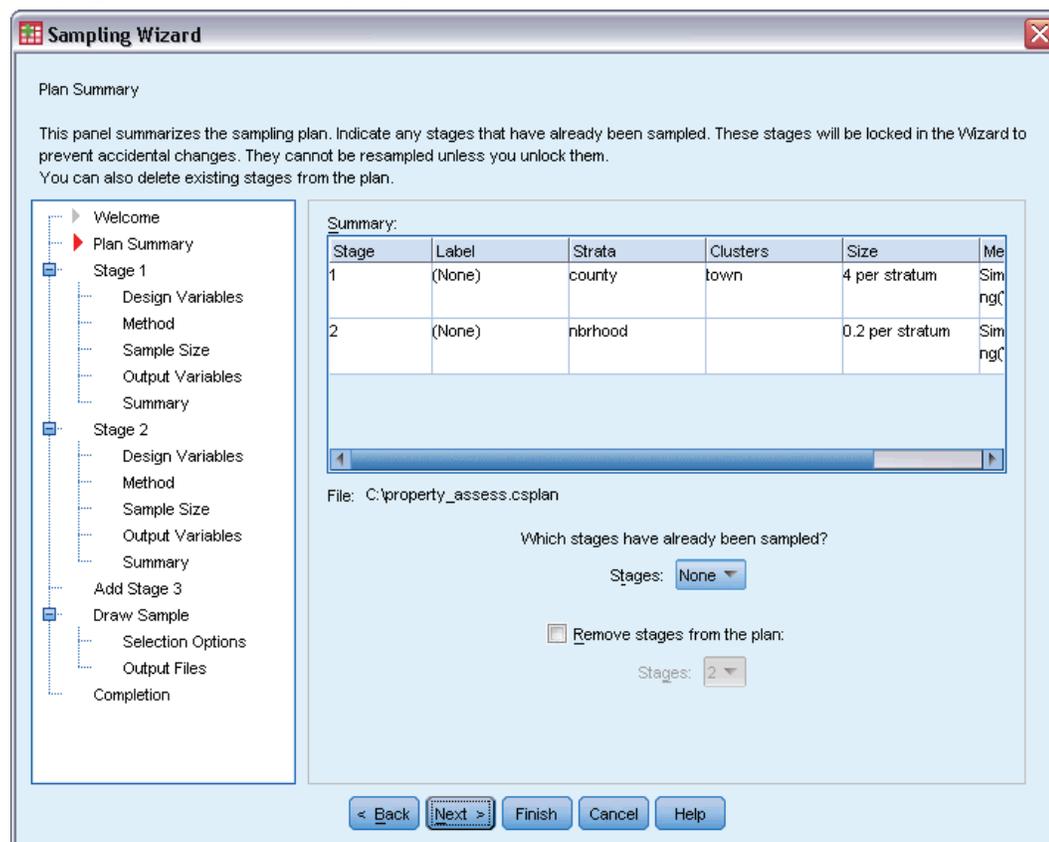
- ▶ Navigate to the Finish step, and specify a new name for the edited plan file or choose to overwrite the existing plan file.

Optionally, you can:

- Specify stages that have already been sampled.
- Remove stages from the plan.

## Sampling Wizard: Plan Summary

Figure 2-11  
Sampling Wizard, Plan Summary step



This step allows you to review the sampling plan and indicate stages that have already been sampled. If editing a plan, you can also remove stages from the plan.

**Previously sampled stages.** If an extended sampling frame is not available, you will have to execute a multistage sampling design one stage at a time. Select which stages have already been sampled from the drop-down list. Any stages that have been executed are locked; they are not available in the Draw Sample Selection Options step, and they cannot be altered when editing a plan.

**Remove stages.** You can remove stages 2 and 3 from a multistage design.

## ***Running an Existing Sample Plan***

- ▶ From the menus choose:  
Analyze > Complex Samples > Select a Sample...
- ▶ Select Draw a sample and choose a plan file to run.
- ▶ Click Next to continue through the Wizard.
- ▶ Review the sampling plan in the Plan Summary step, and then click Next.
- ▶ The individual steps containing stage information are skipped when executing a sample plan. You can now go on to the Finish step at any time.

Optionally, you can specify stages that have already been sampled.

## ***CSPLAN and CSSELECT Commands Additional Features***

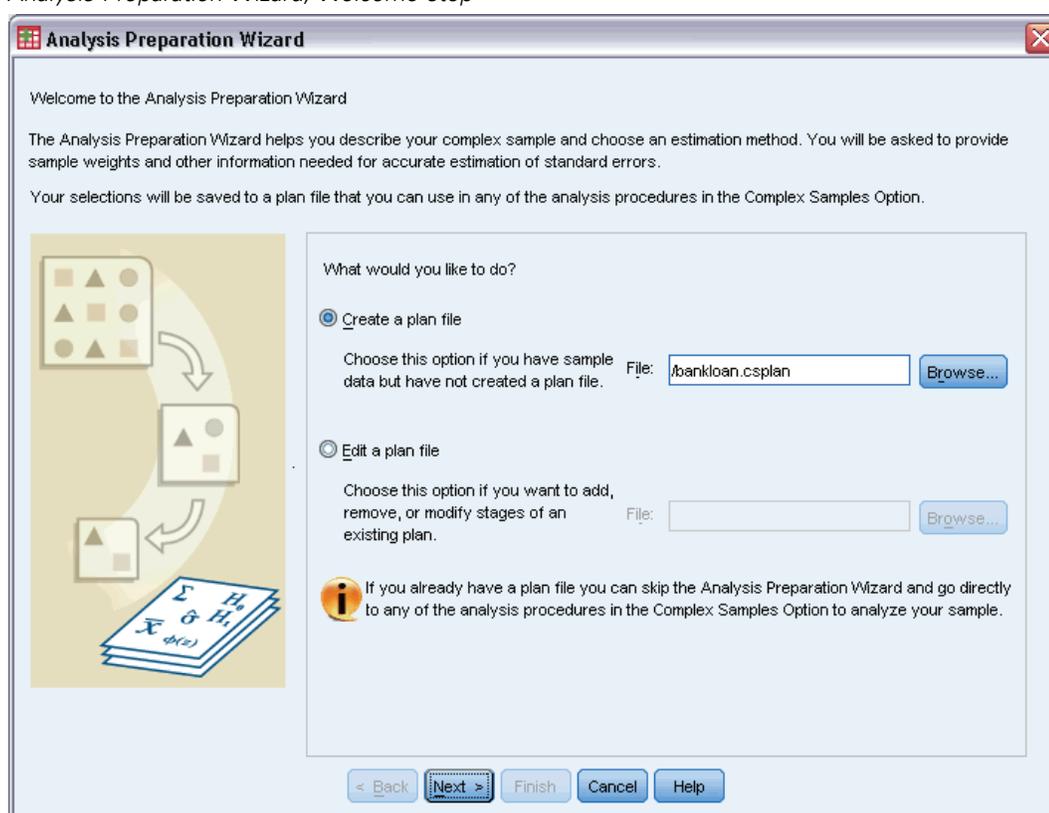
The command syntax language also allows you to:

- Specify custom names for output variables.
- Control the output in the Viewer. For example, you can suppress the stagewise summary of the plan that is displayed if a sample is designed or modified, suppress the summary of the distribution of sampled cases by strata that is shown if the sample design is executed, and request a case processing summary.
- Choose a subset of variables in the active dataset to write to an external sample file or to a different dataset.

See the *Command Syntax Reference* for complete syntax information.

# Preparing a Complex Sample for Analysis

Figure 3-1  
*Analysis Preparation Wizard, Welcome step*



The Analysis Preparation Wizard guides you through the steps for creating or modifying an analysis plan for use with the various Complex Samples analysis procedures. Before using the Wizard, you should have a sample drawn according to a complex design.

Creating a new plan is most useful when you do not have access to the sampling plan file used to draw the sample (recall that the sampling plan contains a default analysis plan). If you do have access to the sampling plan file used to draw the sample, you can use the default analysis plan contained in the sampling plan file or override the default analysis specifications and save your changes to a new file.

## Creating a New Analysis Plan

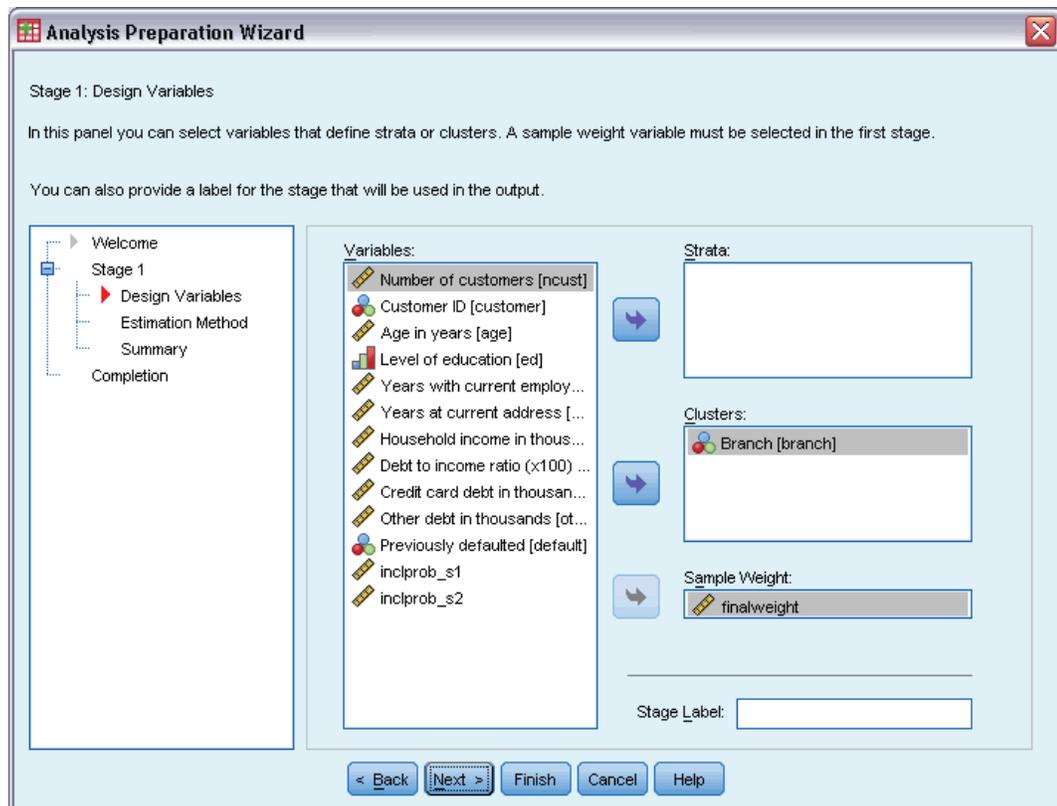
- ▶ From the menus choose:  
Analyze > Complex Samples > Prepare for Analysis...
- ▶ Select Create a plan file, and choose a plan filename to which you will save the analysis plan.
- ▶ Click Next to continue through the Wizard.
- ▶ Specify the variable containing sample weights in the Design Variables step, optionally defining strata and clusters.
- ▶ You can now click Finish to save the plan.

Optionally, in further steps you can:

- Select the method for estimating standard errors in the Estimation Method step.
- Specify the number of units sampled or the inclusion probability per unit in the Size step.
- Add a second or third stage to the design.
- Paste your selections as command syntax.

## Analysis Preparation Wizard: Design Variables

Figure 3-2  
Analysis Preparation Wizard, Design Variables step



This step allows you to identify the stratification and clustering variables and define sample weights. You can also provide a label for the stage.

**Strata.** The cross-classification of stratification variables defines distinct subpopulations, or strata. Your total sample represents the combination of independent samples from each stratum.

**Clusters.** Cluster variables define groups of observational units, or clusters. Samples drawn in multiple stages select clusters in the earlier stages and then subsample units from the selected clusters. When analyzing a data file obtained by sampling clusters with replacement, you should include the duplication index as a cluster variable.

**Sample Weight.** You must provide sample weights in the first stage. Sample weights are computed automatically for subsequent stages of the current design.

**Stage Label.** You can specify an optional string label for each stage. This is used in the output to help identify stagewise information.

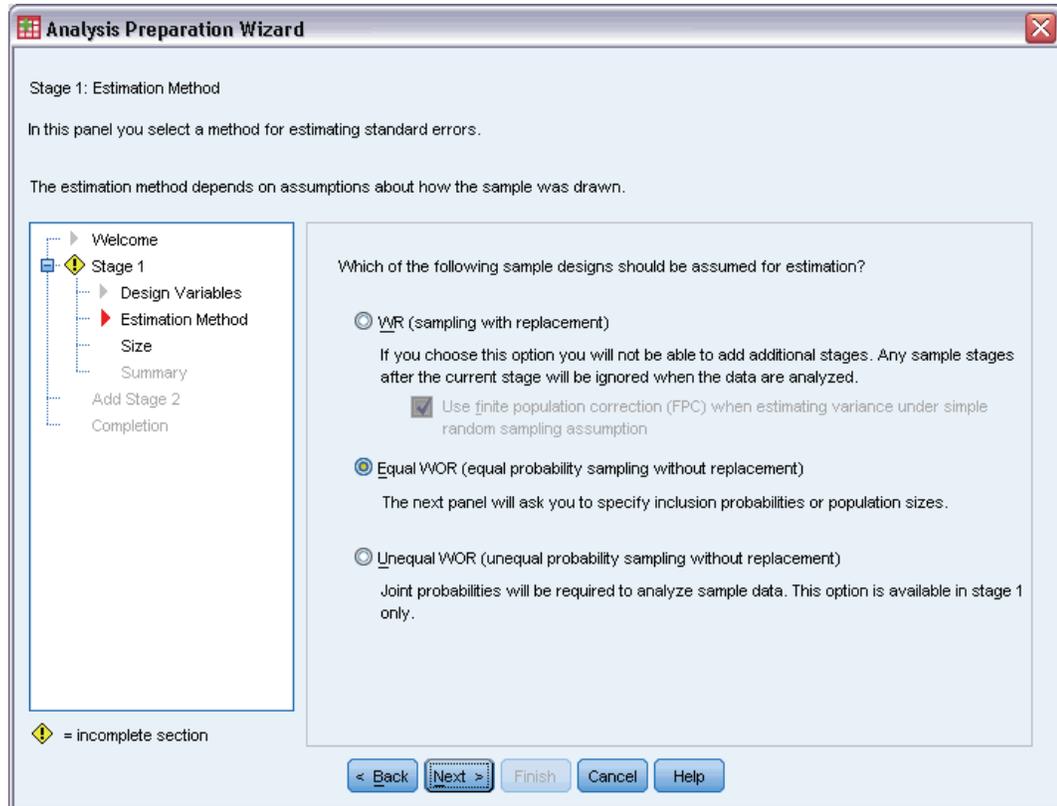
*Note:* The source variable list has the same contents across steps of the Wizard. In other words, variables removed from the source list in a particular step are removed from the list in all steps. Variables returned to the source list show up in all steps.

### ***Tree Controls for Navigating the Analysis Wizard***

At the left side of each step of the Analysis Wizard is an outline of all the steps. You can navigate the Wizard by clicking on the name of an enabled step in the outline. Steps are enabled as long as all previous steps are valid—that is, as long as each previous step has been given the minimum required specifications for that step. For more information on why a given step may be invalid, see the Help for individual steps.

## Analysis Preparation Wizard: Estimation Method

Figure 3-3  
Analysis Preparation Wizard, Estimation Method step



This step allows you to specify an estimation method for the stage.

**WR (sampling with replacement).** WR estimation does not include a correction for sampling from a finite population (FPC) when estimating the variance under the complex sampling design. You can choose to include or exclude the FPC when estimating the variance under simple random sampling (SRS).

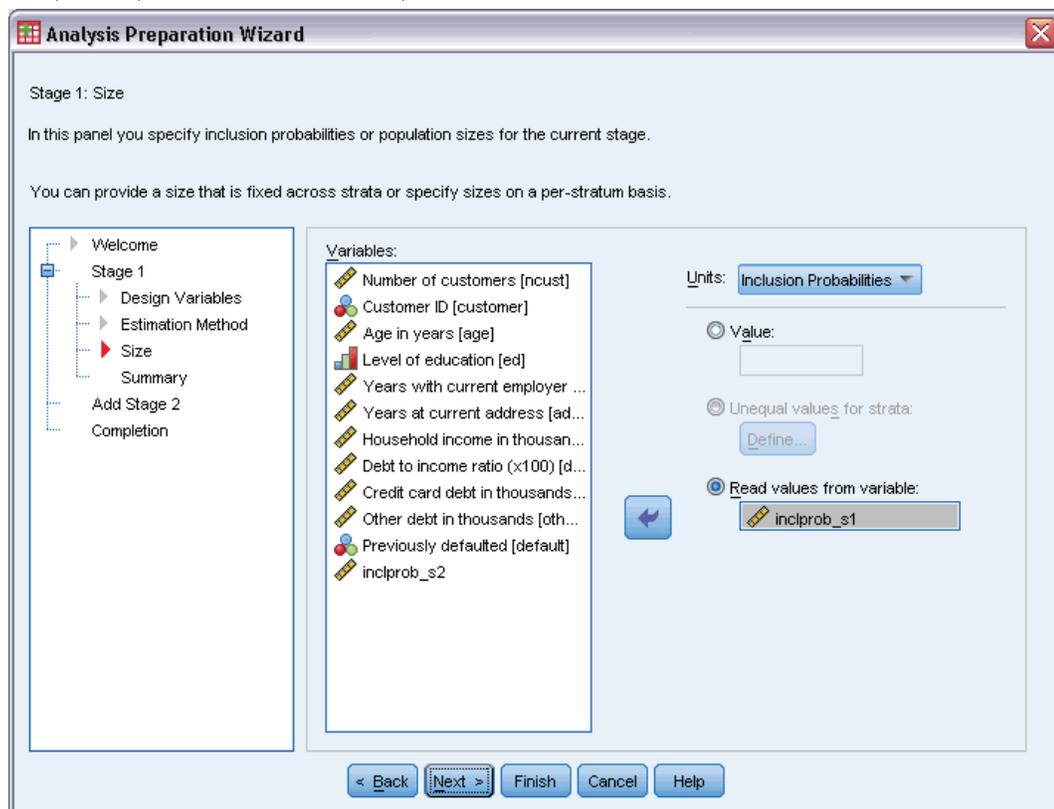
Choosing not to include the FPC for SRS variance estimation is recommended when the analysis weights have been scaled so that they do not add up to the population size. The SRS variance estimate is used in computing statistics like the design effect. WR estimation can be specified only in the final stage of a design; the Wizard will not allow you to add another stage if you select WR estimation.

**Equal WOR (equal probability sampling without replacement).** Equal WOR estimation includes the finite population correction and assumes that units are sampled with equal probability. Equal WOR can be specified in any stage of a design.

**Unequal WOR (unequal probability sampling without replacement).** In addition to using the finite population correction, Unequal WOR accounts for sampling units (usually clusters) selected with unequal probability. This estimation method is available only in the first stage.

## Analysis Preparation Wizard: Size

Figure 3-4  
Analysis Preparation Wizard, Size step



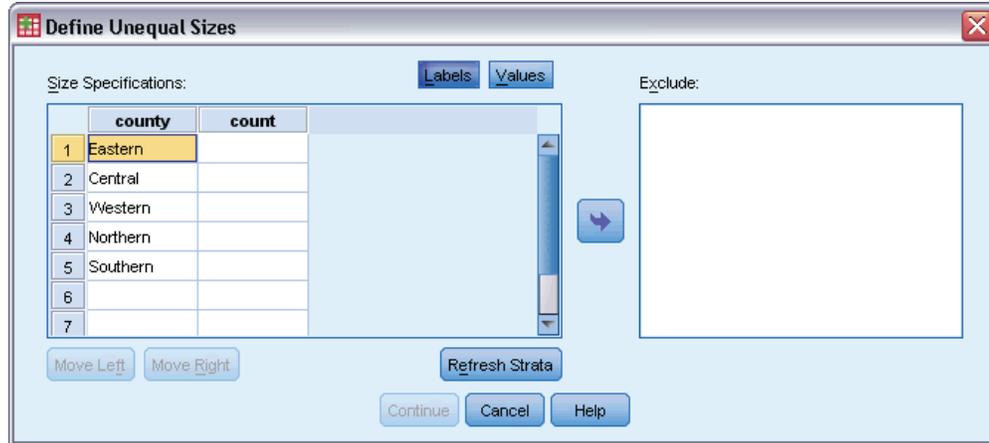
This step is used to specify inclusion probabilities or population sizes for the current stage. Sizes can be fixed or can vary across strata. For the purpose of specifying sizes, clusters specified in previous stages can be used to define strata. Note that this step is necessary only when Equal WOR is chosen as the Estimation Method.

**Units.** You can specify exact population sizes or the probabilities with which units were sampled.

- **Value.** A single value is applied to all strata. If Population Sizes is selected as the unit metric, you should enter a non-negative integer. If Inclusion Probabilities is selected, you should enter a value between 0 and 1, inclusive.
- **Unequal values for strata.** Allows you to enter size values on a per-stratum basis via the Define Unequal Sizes dialog box.
- **Read values from variable.** Allows you to select a numeric variable that contains size values for strata.

## Define Unequal Sizes

Figure 3-5  
Define Unequal Sizes dialog box



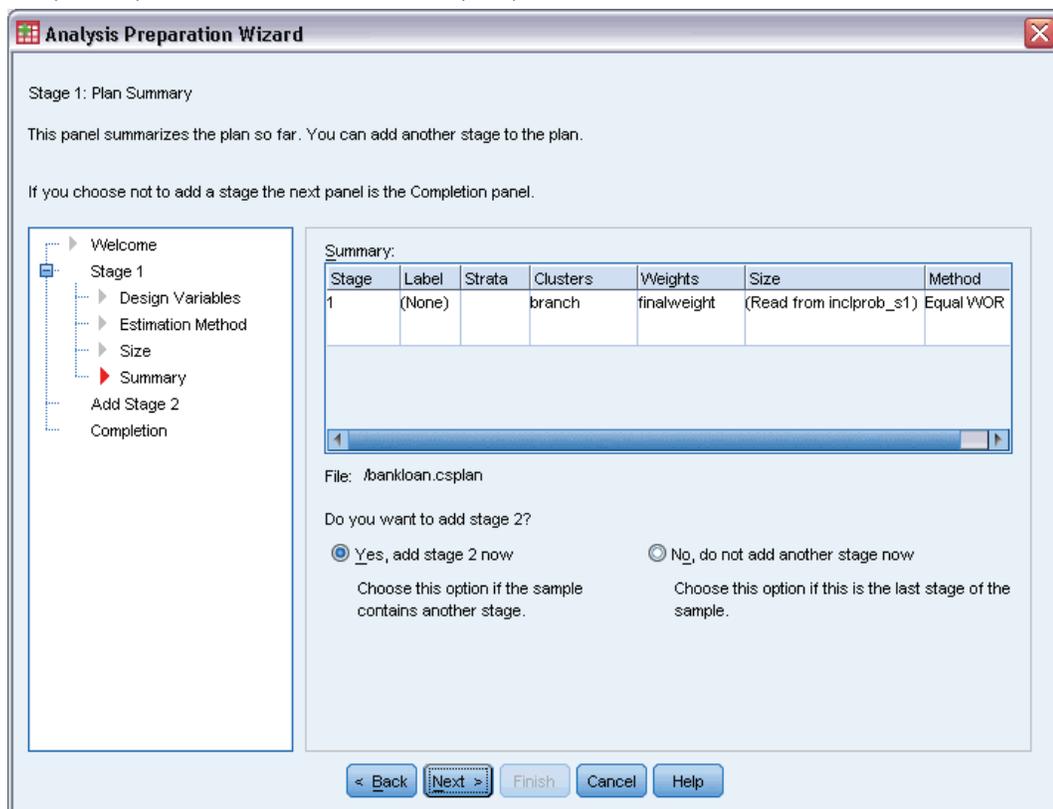
The Define Unequal Sizes dialog box allows you to enter sizes on a per-stratum basis.

**Size Specifications grid.** The grid displays the cross-classifications of up to five strata or cluster variables—one stratum/cluster combination per row. Eligible grid variables include all stratification variables from the current and previous stages and all cluster variables from previous stages. Variables can be reordered within the grid or moved to the Exclude list. Enter sizes in the rightmost column. Click Labels or Values to toggle the display of value labels and data values for stratification and cluster variables in the grid cells. Cells that contain unlabeled values always show values. Click Refresh Strata to repopulate the grid with each combination of labeled data values for variables in the grid.

**Exclude.** To specify sizes for a subset of stratum/cluster combinations, move one or more variables to the Exclude list. These variables are not used to define sample sizes.

## Analysis Preparation Wizard: Plan Summary

Figure 3-6  
Analysis Preparation Wizard, Plan Summary step



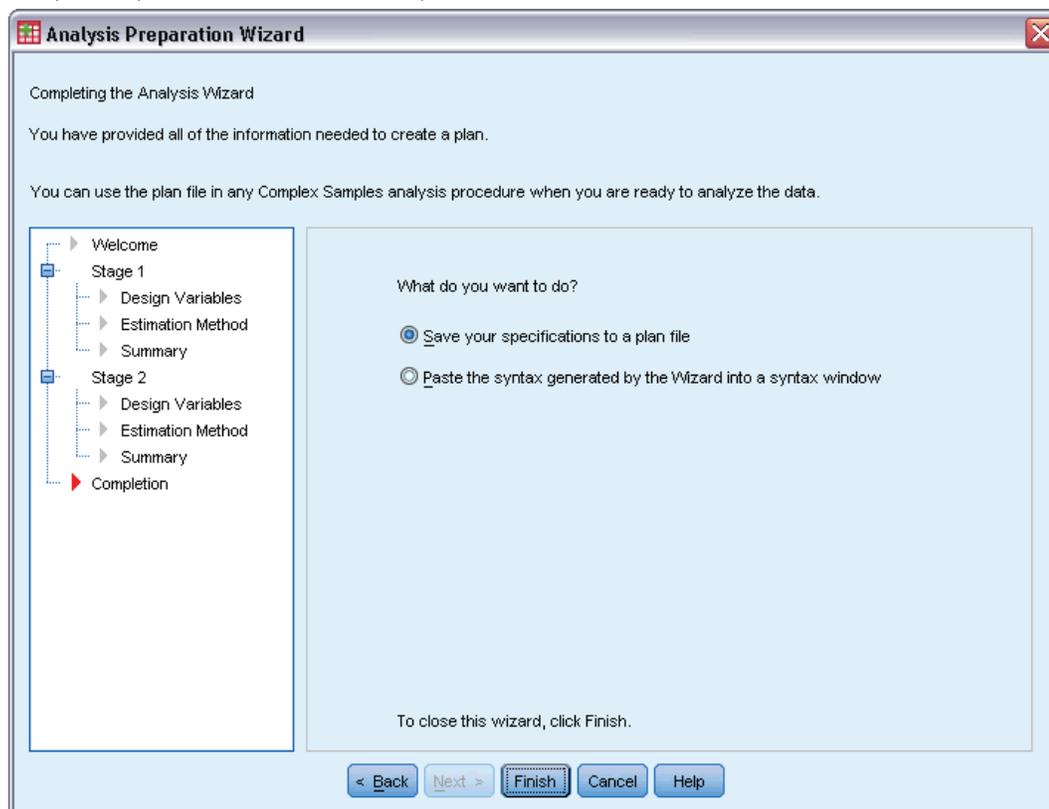
This is the last step within each stage, providing a summary of the analysis design specifications through the current stage. From here, you can either proceed to the next stage (creating it if necessary) or save the analysis specifications.

If you cannot add another stage, it is likely because:

- No cluster variable was specified in the Design Variables step.
- You selected WR estimation in the Estimation Method step.
- This is the third stage of the analysis, and the Wizard supports a maximum of three stages.

## Analysis Preparation Wizard: Finish

Figure 3-7  
Analysis Preparation Wizard, Finish step



This is the final step. You can save the plan file now or paste your selections to a syntax window.

When making changes to stages in the existing plan file, you can save the edited plan to a new file or overwrite the existing file. When adding stages without making changes to existing stages, the Wizard automatically overwrites the existing plan file. If you want to save the plan to a new file, choose to Paste the syntax generated by the Wizard into a syntax window and change the filename in the syntax commands.

## Modifying an Existing Analysis Plan

- ▶ From the menus choose:  
Analyze > Complex Samples > Prepare for Analysis...
- ▶ Select Edit a plan file, and choose a plan filename to which you will save the analysis plan.
- ▶ Click Next to continue through the Wizard.

- ▶ Review the analysis plan in the Plan Summary step, and then click Next.

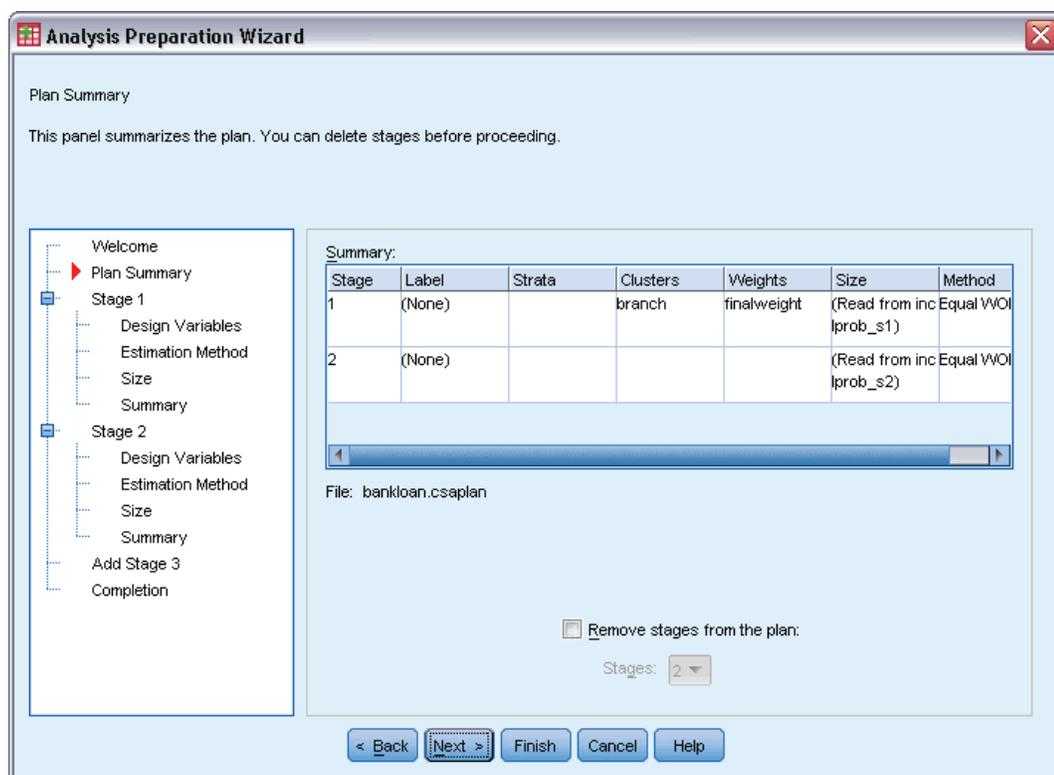
Subsequent steps are largely the same as for a new design. For more information, see the Help for individual steps.

- ▶ Navigate to the Finish step, and specify a new name for the edited plan file, or choose to overwrite the existing plan file.

Optionally, you can remove stages from the plan.

## Analysis Preparation Wizard: Plan Summary

Figure 3-8  
Analysis Preparation Wizard, Plan Summary step



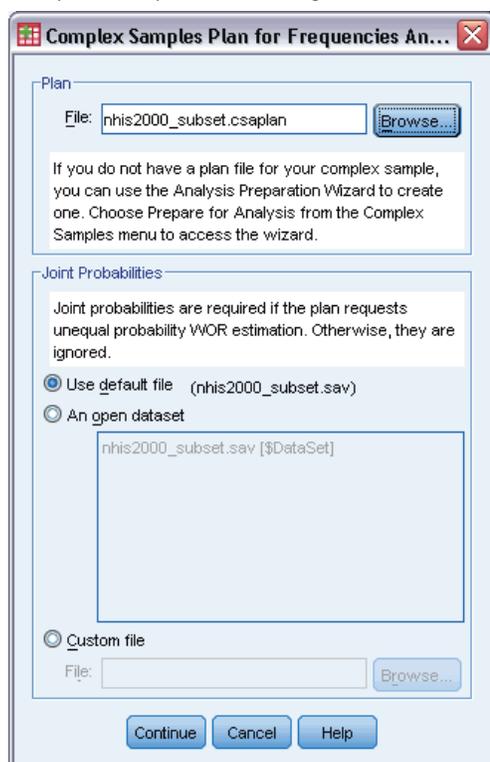
This step allows you to review the analysis plan and remove stages from the plan.

**Remove Stages.** You can remove stages 2 and 3 from a multistage design. Since a plan must have at least one stage, you can edit but not remove stage 1 from the design.

# Complex Samples Plan

Complex Samples analysis procedures require analysis specifications from an analysis or sample plan file in order to provide valid results.

Figure 4-1  
*Complex Samples Plan dialog box*



**Plan.** Specify the path of an analysis or sample plan file.

**Joint Probabilities.** In order to use Unequal WOR estimation for clusters drawn using a PPS WOR method, you need to specify a separate file or an open dataset containing the joint probabilities. This file or dataset is created by the Sampling Wizard during sampling.

# *Complex Samples Frequencies*

The Complex Samples Frequencies procedure produces frequency tables for selected variables and displays univariate statistics. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

**Example.** Using the Complex Samples Frequencies procedure, you can obtain univariate tabular statistics for vitamin usage among U.S. citizens, based on the results of the National Health Interview Survey (NHIS) and with an appropriate analysis plan for this public-use data.

**Statistics.** The procedure produces estimates of cell population sizes and table percentages, plus standard errors, confidence intervals, coefficients of variation, design effects, square roots of design effects, cumulative values, and unweighted counts for each estimate. Additionally, chi-square and likelihood-ratio statistics are computed for the test of equal cell proportions.

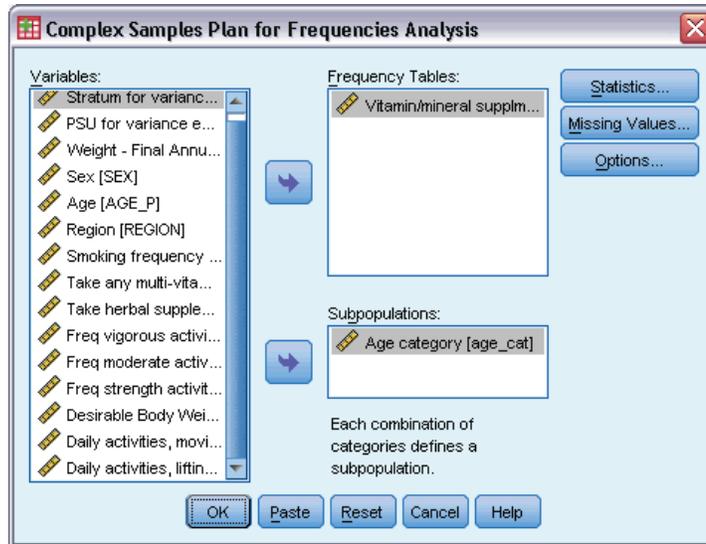
**Data.** Variables for which frequency tables are produced should be categorical. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

## ***Obtaining Complex Samples Frequencies***

- ▶ From the menus choose:  
Analyze > Complex Samples > Frequencies...
- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 5-1  
Frequencies dialog box

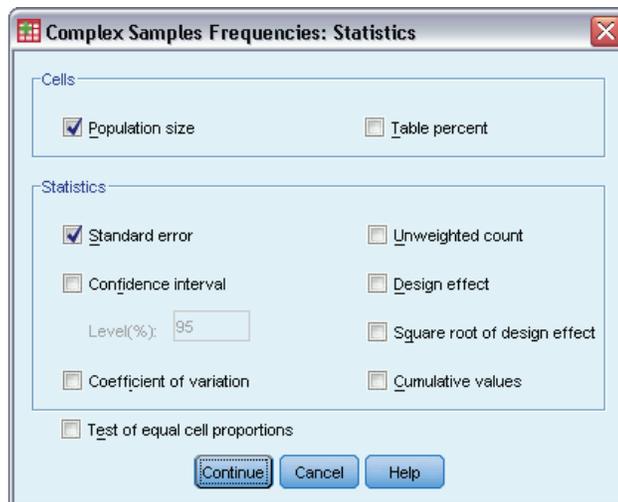


- Select at least one frequency variable.

Optionally, you can specify variables to define subpopulations. Statistics are computed separately for each subpopulation.

## Complex Samples Frequencies Statistics

Figure 5-2  
Frequencies Statistics dialog box



**Cells.** This group allows you to request estimates of the cell population sizes and table percentages.

**Statistics.** This group produces statistics associated with the population size or table percentage.

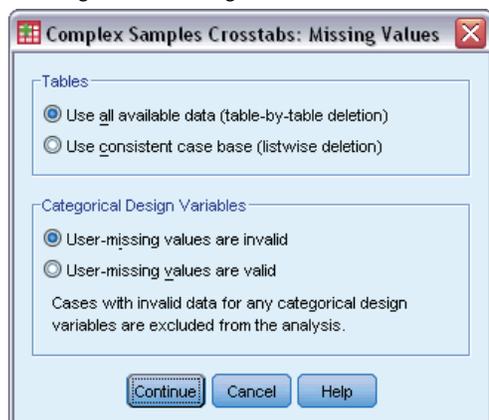
- **Standard error.** The standard error of the estimate.

- **Confidence interval.** A confidence interval for the estimate, using the specified level.
- **Coefficient of variation.** The ratio of the standard error of the estimate to the estimate.
- **Unweighted count.** The number of units used to compute the estimate.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Cumulative values.** The cumulative estimate through each value of the variable.

**Test of equal cell proportions.** This produces chi-square and likelihood-ratio tests of the hypothesis that the categories of a variable have equal frequencies. Separate tests are performed for each variable.

## Complex Samples Missing Values

Figure 5-3  
Missing Values dialog box



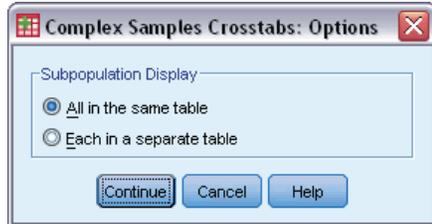
**Tables.** This group determines which cases are used in the analysis.

- **Use all available data.** Missing values are determined on a table-by-table basis. Thus, the cases used to compute statistics may vary across frequency or crosstabulation tables.
- **Use consistent case base.** Missing values are determined across all variables. Thus, the cases used to compute statistics are consistent across tables.

**Categorical Design Variables.** This group determines whether user-missing values are valid or invalid.

## Complex Samples Options

Figure 5-4  
Options dialog box



**Subpopulation Display.** You can choose to have subpopulations displayed in the same table or in separate tables.

# *Complex Samples Descriptives*

The Complex Samples Descriptives procedure displays univariate summary statistics for several variables. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

**Example.** Using the Complex Samples Descriptives procedure, you can obtain univariate descriptive statistics for the activity levels of U.S. citizens, based on the results of the National Health Interview Survey (NHIS) and with an appropriate analysis plan for this public-use data.

**Statistics.** The procedure produces means and sums, plus  $t$  tests, standard errors, confidence intervals, coefficients of variation, unweighted counts, population sizes, design effects, and square roots of design effects for each estimate.

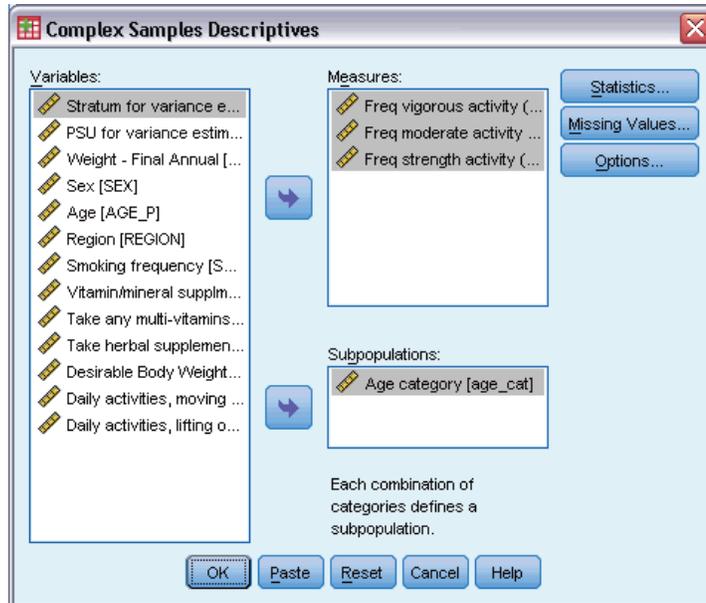
**Data.** Measures should be scale variables. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

## ***Obtaining Complex Samples Descriptives***

- ▶ From the menus choose:  
Analyze > Complex Samples > Descriptives...
- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 6-1  
Descriptives dialog box

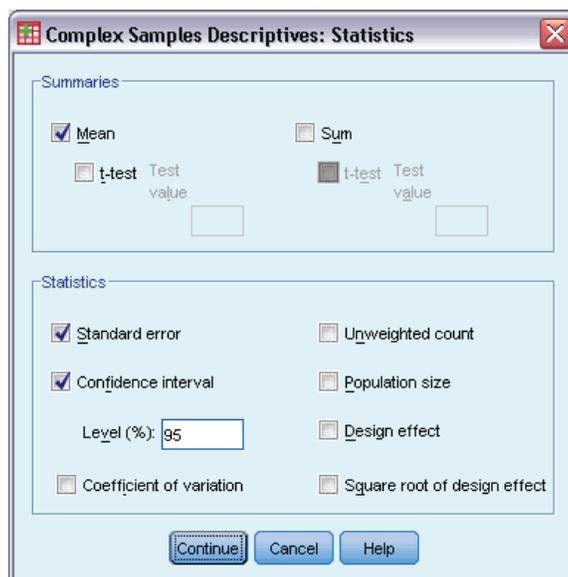


- Select at least one measure variable.

Optionally, you can specify variables to define subpopulations. Statistics are computed separately for each subpopulation.

## Complex Samples Descriptives Statistics

Figure 6-2  
Descriptives Statistics dialog box



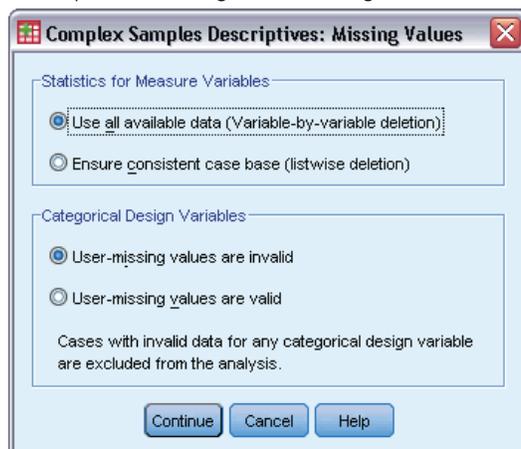
**Summaries.** This group allows you to request estimates of the means and sums of the measure variables. Additionally, you can request  $t$  tests of the estimates against a specified value.

**Statistics.** This group produces statistics associated with the mean or sum.

- **Standard error.** The standard error of the estimate.
- **Confidence interval.** A confidence interval for the estimate, using the specified level.
- **Coefficient of variation.** The ratio of the standard error of the estimate to the estimate.
- **Unweighted count.** The number of units used to compute the estimate.
- **Population size.** The estimated number of units in the population.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.

## Complex Samples Descriptives Missing Values

Figure 6-3  
Descriptives Missing Values dialog box



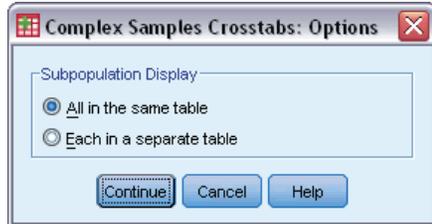
**Statistics for Measure Variables.** This group determines which cases are used in the analysis.

- **Use all available data.** Missing values are determined on a variable-by-variable basis, thus the cases used to compute statistics may vary across measure variables.
- **Ensure consistent case base.** Missing values are determined across all variables, thus the cases used to compute statistics are consistent.

**Categorical Design Variables.** This group determines whether user-missing values are valid or invalid.

## Complex Samples Options

Figure 6-4  
Options dialog box



**Subpopulation Display.** You can choose to have subpopulations displayed in the same table or in separate tables.

# *Complex Samples Crosstabs*

The Complex Samples Crosstabs procedure produces crosstabulation tables for pairs of selected variables and displays two-way statistics. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

**Example.** Using the Complex Samples Crosstabs procedure, you can obtain cross-classification statistics for smoking frequency by vitamin usage of U.S. citizens, based on the results of the National Health Interview Survey (NHIS) and with an appropriate analysis plan for this public-use data.

**Statistics.** The procedure produces estimates of cell population sizes and row, column, and table percentages, plus standard errors, confidence intervals, coefficients of variation, expected values, design effects, square roots of design effects, residuals, adjusted residuals, and unweighted counts for each estimate. The odds ratio, relative risk, and risk difference are computed for 2-by-2 tables. Additionally, Pearson and likelihood-ratio statistics are computed for the test of independence of the row and column variables.

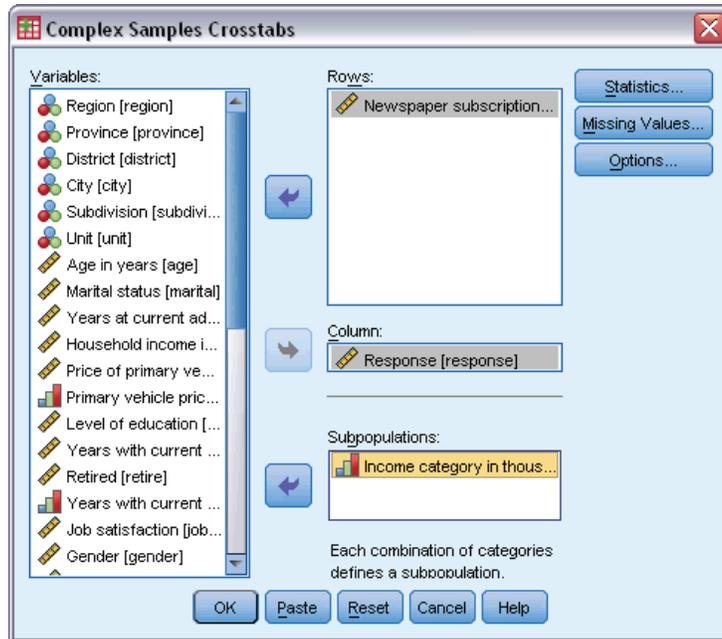
**Data.** Row and column variables should be categorical. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

## ***Obtaining Complex Samples Crosstabs***

- ▶ From the menus choose:  
Analyze > Complex Samples > Crosstabs...
- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 7-1  
Crosstabs dialog box

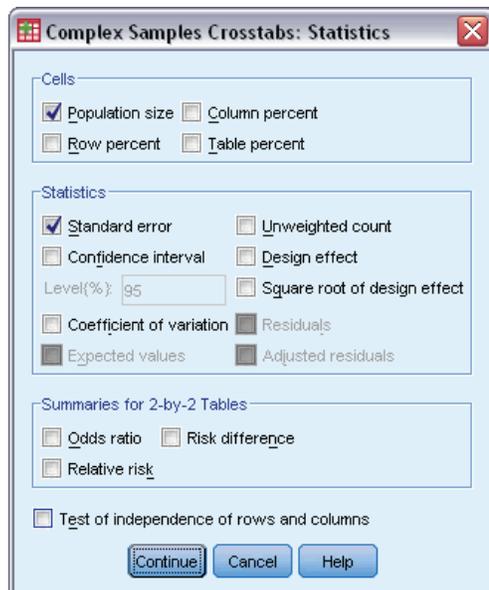


- Select at least one row variable and one column variable.

Optionally, you can specify variables to define subpopulations. Statistics are computed separately for each subpopulation.

## Complex Samples Crosstabs Statistics

Figure 7-2  
Crosstabs Statistics dialog box



**Cells.** This group allows you to request estimates of the cell population size and row, column, and table percentages.

**Statistics.** This group produces statistics associated with the population size and row, column, and table percentages.

- **Standard error.** The standard error of the estimate.
- **Confidence interval.** A confidence interval for the estimate, using the specified level.
- **Coefficient of variation.** The ratio of the standard error of the estimate to the estimate.
- **Expected values.** The expected value of the estimate, under the hypothesis of independence of the row and column variable.
- **Unweighted count.** The number of units used to compute the estimate.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Residuals.** The expected value is the number of cases that you would expect in the cell if there were no relationship between the two variables. A positive residual indicates that there are more cases in the cell than there would be if the row and column variables were independent.
- **Adjusted residuals.** The residual for a cell (observed minus expected value) divided by an estimate of its standard error. The resulting standardized residual is expressed in standard deviation units above or below the mean.

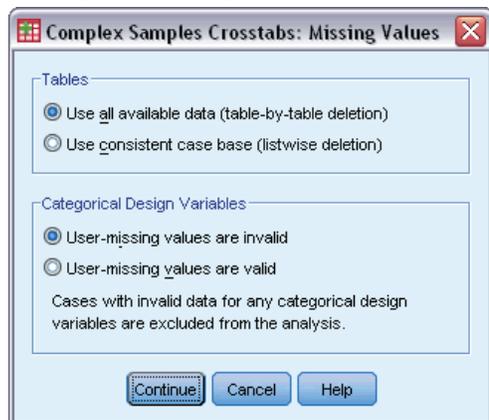
**Summaries for 2-by-2 Tables.** This group produces statistics for tables in which the row and column variable each have two categories. Each is a measure of the strength of the association between the presence of a factor and the occurrence of an event.

- **Odds ratio.** The odds ratio can be used as an estimate of relative risk when the occurrence of the factor is rare.
- **Relative risk.** The ratio of the risk of an event in the presence of the factor to the risk of the event in the absence of the factor.
- **Risk difference.** The difference between the risk of an event in the presence of the factor and the risk of the event in the absence of the factor.

**Test of independence of rows and columns.** This produces chi-square and likelihood-ratio tests of the hypothesis that a row and column variable are independent. Separate tests are performed for each pair of variables.

## Complex Samples Missing Values

Figure 7-3  
Missing Values dialog box



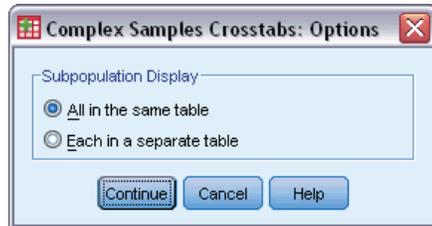
**Tables.** This group determines which cases are used in the analysis.

- **Use all available data.** Missing values are determined on a table-by-table basis. Thus, the cases used to compute statistics may vary across frequency or crosstabulation tables.
- **Use consistent case base.** Missing values are determined across all variables. Thus, the cases used to compute statistics are consistent across tables.

**Categorical Design Variables.** This group determines whether user-missing values are valid or invalid.

## Complex Samples Options

Figure 7-4  
Options dialog box



**Subpopulation Display.** You can choose to have subpopulations displayed in the same table or in separate tables.

# *Complex Samples Ratios*

The Complex Samples Ratios procedure displays univariate summary statistics for ratios of variables. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

**Example.** Using the Complex Samples Ratios procedure, you can obtain descriptive statistics for the ratio of current property value to last assessed value, based on the results of a statewide survey carried out according to a complex design and with an appropriate analysis plan for the data.

**Statistics.** The procedure produces ratio estimates, *t* tests, standard errors, confidence intervals, coefficients of variation, unweighted counts, population sizes, design effects, and square roots of design effects.

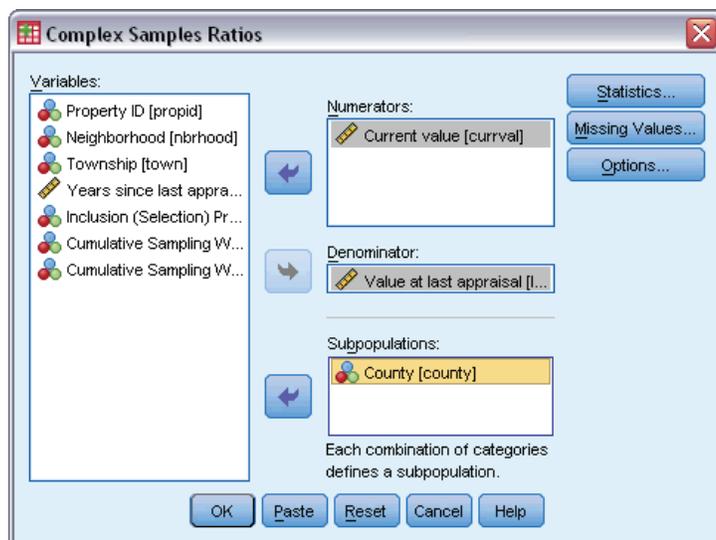
**Data.** Numerators and denominators should be positive-valued scale variables. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

## ***Obtaining Complex Samples Ratios***

- ▶ From the menus choose:  
Analyze > Complex Samples > Ratios...
- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

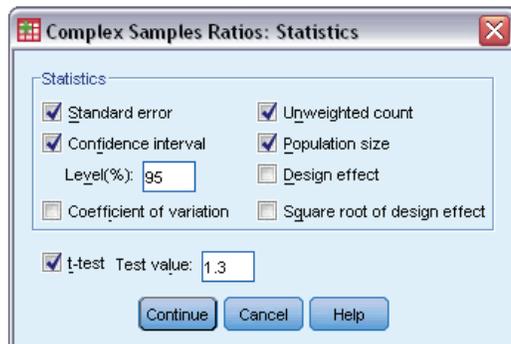
Figure 8-1  
Ratios dialog box



- ▶ Select at least one numerator variable and denominator variable.  
Optionally, you can specify variables to define subgroups for which statistics are produced.

## Complex Samples Ratios Statistics

Figure 8-2  
Ratios Statistics dialog box



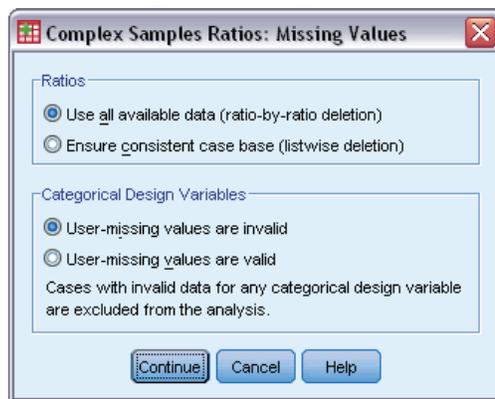
**Statistics.** This group produces statistics associated with the ratio estimate.

- **Standard error.** The standard error of the estimate.
- **Confidence interval.** A confidence interval for the estimate, using the specified level.
- **Coefficient of variation.** The ratio of the standard error of the estimate to the estimate.
- **Unweighted count.** The number of units used to compute the estimate.
- **Population size.** The estimated number of units in the population.

- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
  - **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- T test.** You can request  $t$  tests of the estimates against a specified value.

## Complex Samples Ratios Missing Values

Figure 8-3  
Ratios Missing Values dialog box



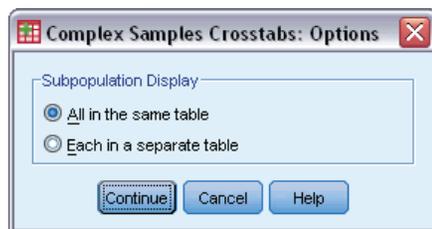
**Ratios.** This group determines which cases are used in the analysis.

- **Use all available data.** Missing values are determined on a ratio-by-ratio basis. Thus, the cases used to compute statistics may vary across numerator-denominator pairs.
- **Ensure consistent case base.** Missing values are determined across all variables. Thus, the cases used to compute statistics are consistent.

**Categorical Design Variables.** This group determines whether user-missing values are valid or invalid.

## Complex Samples Options

Figure 8-4  
Options dialog box



**Subpopulation Display.** You can choose to have subpopulations displayed in the same table or in separate tables.

# Complex Samples General Linear Model

The Complex Samples General Linear Model (CSGLM) procedure performs linear regression analysis, as well as analysis of variance and covariance, for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

**Example.** A grocery store chain surveyed a set of customers concerning their purchasing habits, according to a complex design. Given the survey results and how much each customer spent in the previous month, the store wants to see if the frequency with which customers shop is related to the amount they spend in a month, controlling for the gender of the customer and incorporating the sampling design.

**Statistics.** The procedure produces estimates, standard errors, confidence intervals,  $t$  tests, design effects, and square roots of design effects for model parameters, as well as the correlations and covariances between parameter estimates. Measures of model fit and descriptive statistics for the dependent and independent variables are also available. Additionally, you can request estimated marginal means for levels of model factors and factor interactions.

**Data.** The dependent variable is quantitative. Factors are categorical. Covariates are quantitative variables that are related to the dependent variable. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

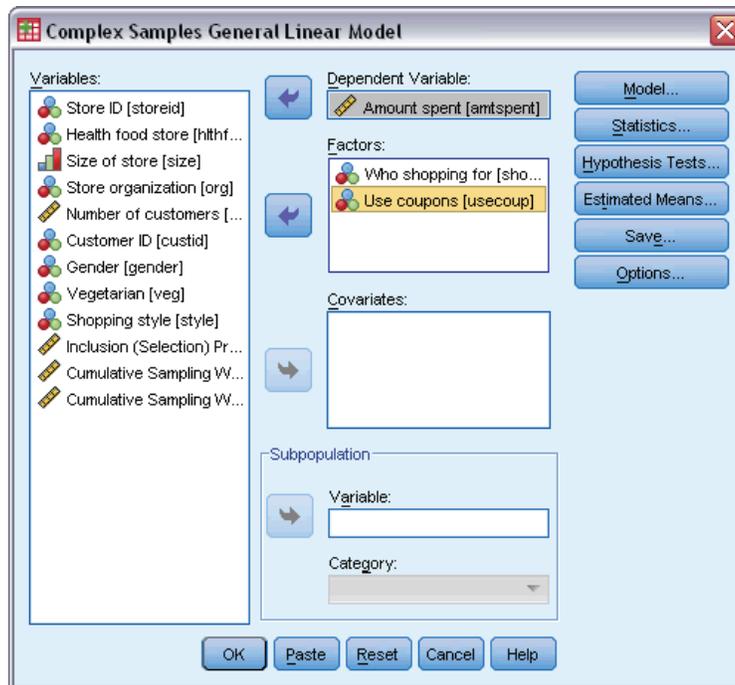
## **Obtaining a Complex Samples General Linear Model**

From the menus choose:

Analyze > Complex Samples > General Linear Model...

- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 9-1  
General Linear Model dialog box

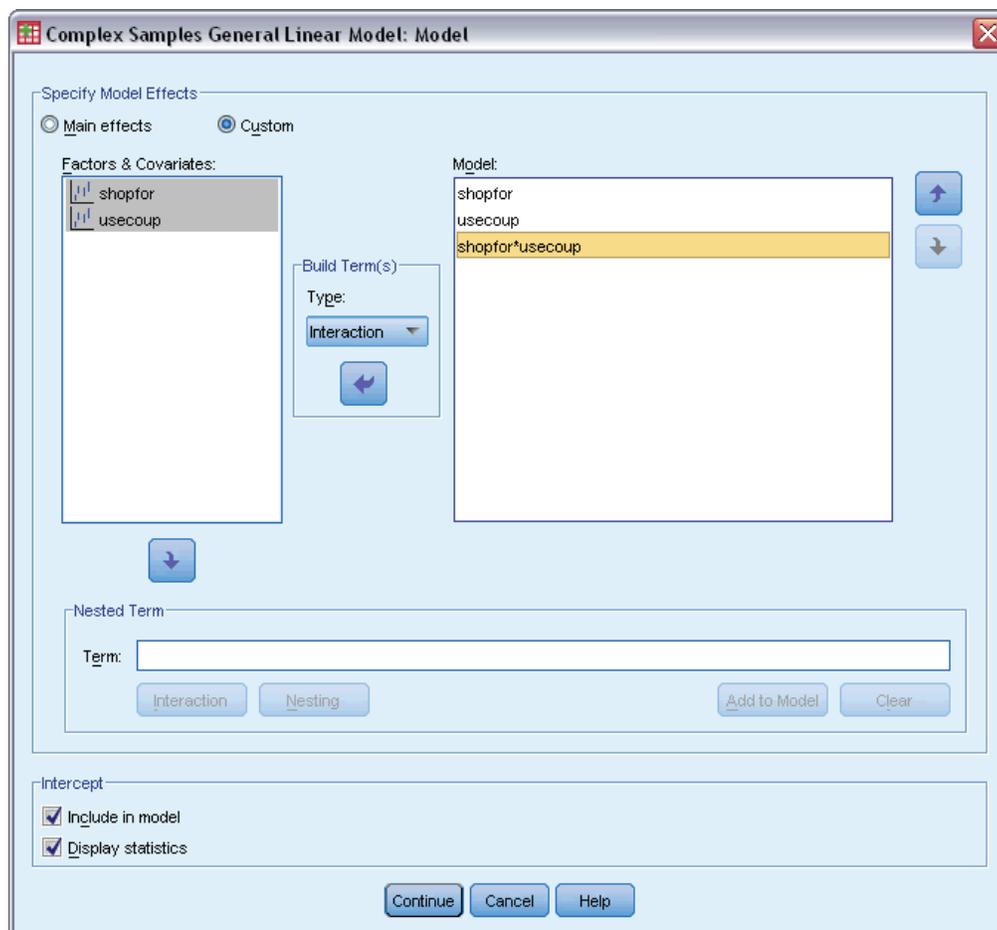


- ▶ Select a dependent variable.

Optionally, you can:

- Select variables for factors and covariates, as appropriate for your data.
- Specify a variable to define a subpopulation. The analysis is performed only for the selected category of the subpopulation variable.

Figure 9-2  
Model dialog box



**Specify Model Effects.** By default, the procedure builds a main-effects model using the factors and covariates specified in the main dialog box. Alternatively, you can build a custom model that includes interaction effects and nested terms.

### ***Non-Nested Terms***

For the selected factors and covariates:

**Interaction.** Creates the highest-level interaction term for all selected variables.

**Main effects.** Creates a main-effects term for each variable selected.

**All 2-way.** Creates all possible two-way interactions of the selected variables.

**All 3-way.** Creates all possible three-way interactions of the selected variables.

**All 4-way.** Creates all possible four-way interactions of the selected variables.

**All 5-way.** Creates all possible five-way interactions of the selected variables.

### ***Nested Terms***

You can build nested terms for your model in this procedure. Nested terms are useful for modeling the effect of a factor or covariate whose values do not interact with the levels of another factor. For example, a grocery store chain may follow the spending habits of its customers at several store locations. Since each customer frequents only one of these locations, the *Customer* effect can be said to be **nested within** the *Store location* effect.

Additionally, you can include interaction effects, such as polynomial terms involving the same covariate, or add multiple levels of nesting to the nested term.

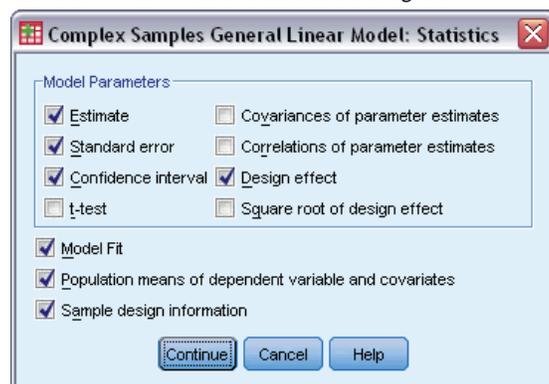
**Limitations.** Nested terms have the following restrictions:

- All factors within an interaction must be unique. Thus, if  $A$  is a factor, then specifying  $A*A$  is invalid.
- All factors within a nested effect must be unique. Thus, if  $A$  is a factor, then specifying  $A(A)$  is invalid.
- No effect can be nested within a covariate. Thus, if  $A$  is a factor and  $X$  is a covariate, then specifying  $A(X)$  is invalid.

**Intercept.** The intercept is usually included in the model. If you can assume the data pass through the origin, you can exclude the intercept. Even if you include the intercept in the model, you can choose to suppress statistics related to it.

## ***Complex Samples General Linear Model Statistics***

Figure 9-3  
General Linear Model Statistics dialog box



**Model Parameters.** This group allows you to control the display of statistics related to the model parameters.

- **Estimate.** Displays estimates of the coefficients.
- **Standard error.** Displays the standard error for each coefficient estimate.
- **Confidence interval.** Displays a confidence interval for each coefficient estimate. The confidence level for the interval is set in the Options dialog box.
- **T test.** Displays a  $t$  test of each coefficient estimate. The null hypothesis for each test is that the value of the coefficient is 0.

- **Covariances of parameter estimates.** Displays an estimate of the covariance matrix for the model coefficients.
- **Correlations of parameter estimates.** Displays an estimate of the correlation matrix for the model coefficients.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.

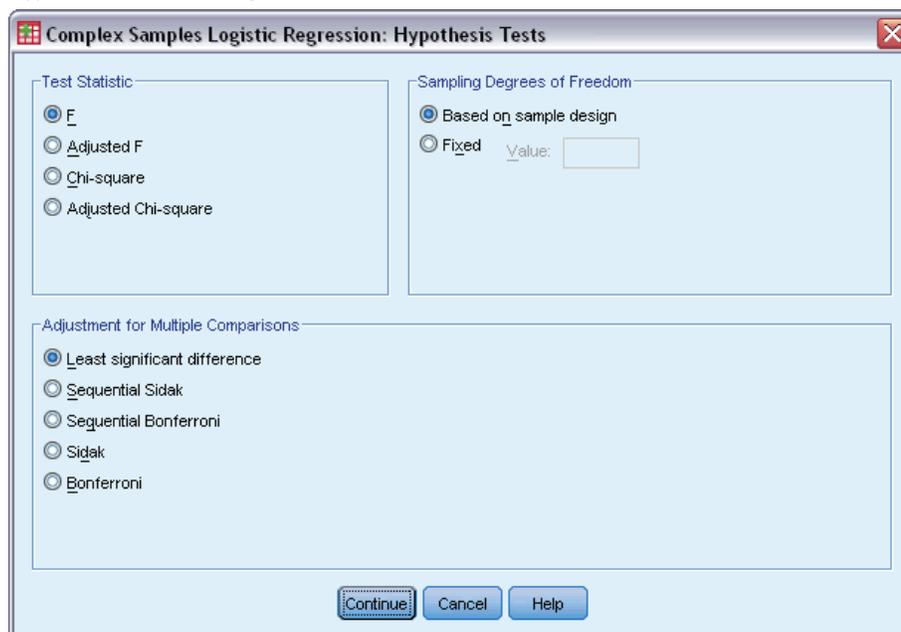
**Model fit.** Displays  $R^2$  and root mean squared error statistics.

**Population means of dependent variable and covariates.** Displays summary information about the dependent variable, covariates, and factors.

**Sample design information.** Displays summary information about the sample, including the unweighted count and the population size.

## Complex Samples Hypothesis Tests

Figure 9-4  
Hypothesis Tests dialog box



**Test Statistic.** This group allows you to select the type of statistic used for testing hypotheses. You can choose between  $F$ , adjusted  $F$ , chi-square, and adjusted chi-square.

**Sampling Degrees of Freedom.** This group gives you control over the sampling design degrees of freedom used to compute  $p$  values for all test statistics. If based on the sampling design, the value is the difference between the number of primary sampling units and the number of strata in the

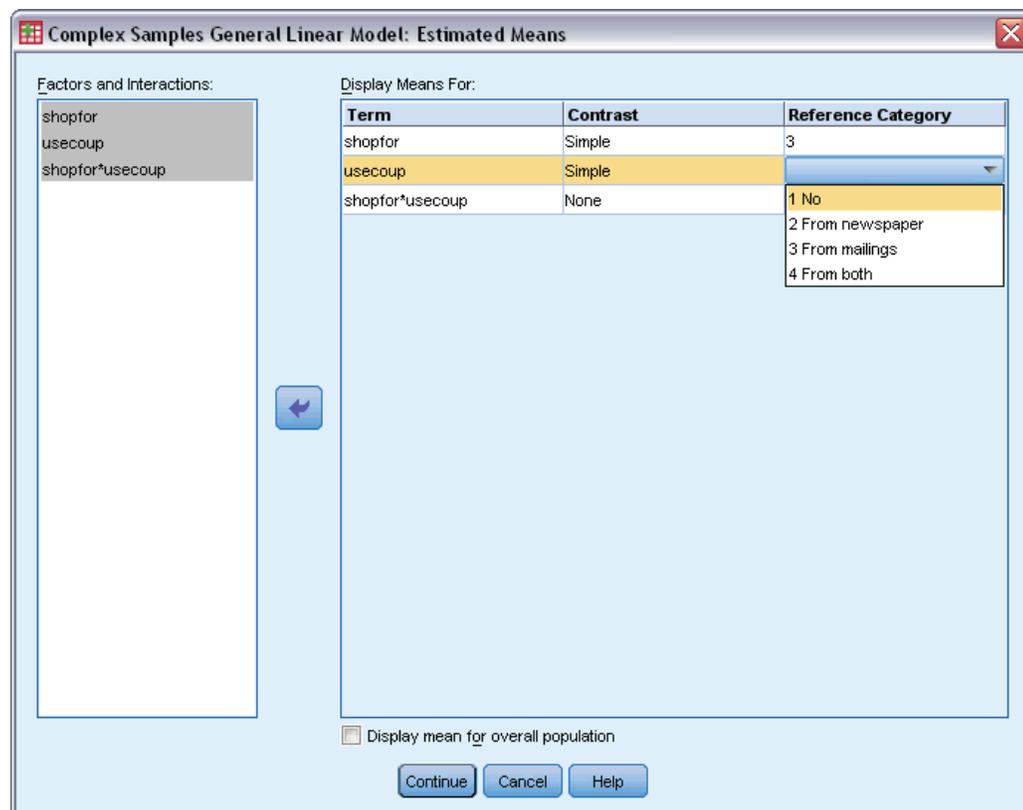
first stage of sampling. Alternatively, you can set a custom degrees of freedom by specifying a positive integer.

**Adjustment for Multiple Comparisons.** When performing hypothesis tests with multiple contrasts, the overall significance level can be adjusted from the significance levels for the included contrasts. This group allows you to choose the adjustment method.

- **Least significant difference.** This method does not control the overall probability of rejecting the hypotheses that some linear contrasts are different from the null hypothesis values.
- **Sequential Sidak.** This is a sequentially step-down rejective Sidak procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sequential Bonferroni.** This is a sequentially step-down rejective Bonferroni procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sidak.** This method provides tighter bounds than the Bonferroni approach.
- **Bonferroni.** This method adjusts the observed significance level for the fact that multiple contrasts are being tested.

## Complex Samples General Linear Model Estimated Means

Figure 9-5  
General Linear Model Estimated Means dialog box



The Estimated Means dialog box allows you to display the model-estimated marginal means for levels of factors and factor interactions specified in the Model subdialog box. You can also request that the overall population mean be displayed.

**Term.** Estimated means are computed for the selected factors and factor interactions.

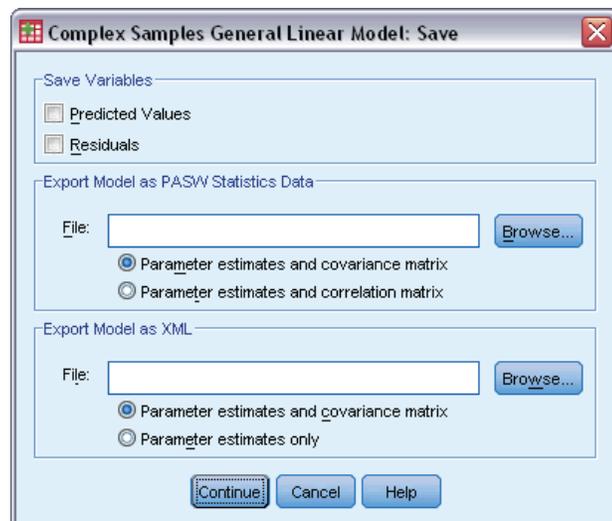
**Contrast.** The contrast determines how hypothesis tests are set up to compare the estimated means.

- **Simple.** Compares the mean of each level to the mean of a specified level. This type of contrast is useful when there is a control group.
- **Deviation.** Compares the mean of each level (except a reference category) to the mean of all of the levels (grand mean). The levels of the factor can be in any order.
- **Difference.** Compares the mean of each level (except the first) to the mean of previous levels. They are sometimes called reverse Helmert contrasts.
- **Helmert.** Compares the mean of each level of the factor (except the last) to the mean of subsequent levels.
- **Repeated.** Compares the mean of each level (except the last) to the mean of the subsequent level.
- **Polynomial.** Compares the linear effect, quadratic effect, cubic effect, and so on. The first degree of freedom contains the linear effect across all categories; the second degree of freedom, the quadratic effect; and so on. These contrasts are often used to estimate polynomial trends.

**Reference Category.** The simple and deviation contrasts require a reference category or factor level against which the others are compared.

## Complex Samples General Linear Model Save

Figure 9-6  
General Linear Model Save dialog box



**Save Variables.** This group allows you to save the model predicted values and residuals as new variables in the working file.

**Export model as SPSS Statistics data.** Writes a dataset in IBM® SPSS® Statistics format containing the parameter correlation or covariance matrix with parameter estimates, standard errors, significance values, and degrees of freedom. The order of variables in the matrix file is as follows.

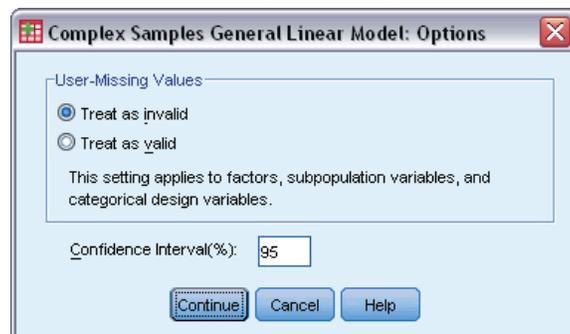
- **rowtype\_.** Takes values (and value labels), COV (Covariances), CORR (Correlations), EST (Parameter estimates), SE (Standard errors), SIG (Significance levels), and DF (Sampling design degrees of freedom). There is a separate case with row type COV (or CORR) for each model parameter, plus a separate case for each of the other row types.
- **varname\_.** Takes values P1, P2, ..., corresponding to an ordered list of all model parameters, for row types COV or CORR, with value labels corresponding to the parameter strings shown in the parameter estimates table. The cells are blank for other row types.
- **P1, P2, ...** These variables correspond to an ordered list of all model parameters, with variable labels corresponding to the parameter strings shown in the parameter estimates table, and take values according to the row type. For redundant parameters, all covariances are set to zero; correlations are set to the system-missing value; all parameter estimates are set at zero; and all standard errors, significance levels, and residual degrees of freedom are set to the system-missing value.

*Note:* This file is not immediately usable for further analyses in other procedures that read a matrix file unless those procedures accept all the row types exported here.

**Export Model as XML.** Saves the parameter estimates and the parameter covariance matrix, if selected, in XML (PMML) format. You can use this model file to apply the model information to other data files for scoring purposes.

## Complex Samples General Linear Model Options

Figure 9-7  
General Linear Model Options dialog box



**User-Missing Values.** All design variables, as well as the dependent variable and any covariates, must have valid data. Cases with invalid data for any of these variables are deleted from the analysis. These controls allow you to decide whether user-missing values are treated as valid among the strata, cluster, subpopulation, and factor variables.

**Confidence Interval.** This is the confidence interval level for coefficient estimates and estimated marginal means. Specify a value greater than or equal to 50 and less than 100.

## ***CSGLM Command Additional Features***

The command syntax language also allows you to:

- Specify custom tests of effects versus a linear combination of effects or a value (using the `CUSTOM` subcommand).
- Fix covariates at values other than their means when computing estimated marginal means (using the `EMMEANS` subcommand).
- Specify a metric for polynomial contrasts (using the `EMMEANS` subcommand).
- Specify a tolerance value for checking singularity (using the `CRITERIA` subcommand).
- Create user-specified names for saved variables (using the `SAVE` subcommand).
- Produce a general estimable function table (using the `PRINT` subcommand).

See the *Command Syntax Reference* for complete syntax information.

# ***Complex Samples Logistic Regression***

The Complex Samples Logistic Regression procedure performs logistic regression analysis on a binary or multinomial dependent variable for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

**Example.** A loan officer has collected past records of customers given loans at several different branches, according to a complex design. While incorporating the sample design, the officer wants to see if the probability with which a customer defaults is related to age, employment history, and amount of credit debt.

**Statistics.** The procedure produces estimates, exponentiated estimates, standard errors, confidence intervals,  $t$  tests, design effects, and square roots of design effects for model parameters, as well as the correlations and covariances between parameter estimates. Pseudo  $R^2$  statistics, classification tables, and descriptive statistics for the dependent and independent variables are also available.

**Data.** The dependent variable is categorical. Factors are categorical. Covariates are quantitative variables that are related to the dependent variable. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

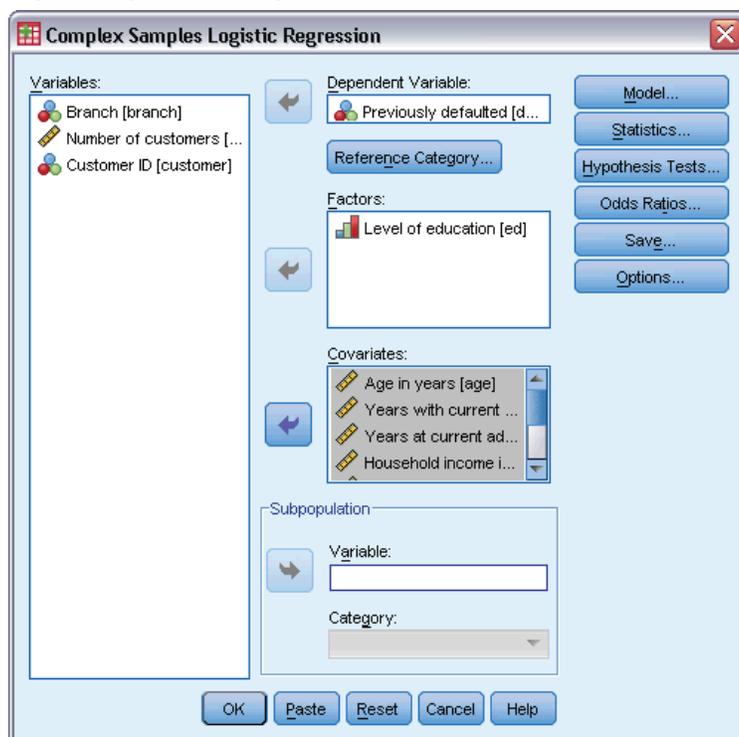
## ***Obtaining Complex Samples Logistic Regression***

From the menus choose:

Analyze > Complex Samples > Logistic Regression...

- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 10-1  
Logistic Regression dialog box



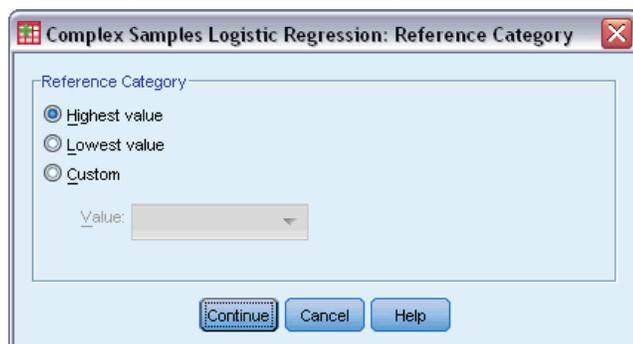
- ▶ Select a dependent variable.

Optionally, you can:

- Select variables for factors and covariates, as appropriate for your data.
- Specify a variable to define a subpopulation. The analysis is performed only for the selected category of the subpopulation variable.

## ***Complex Samples Logistic Regression Reference Category***

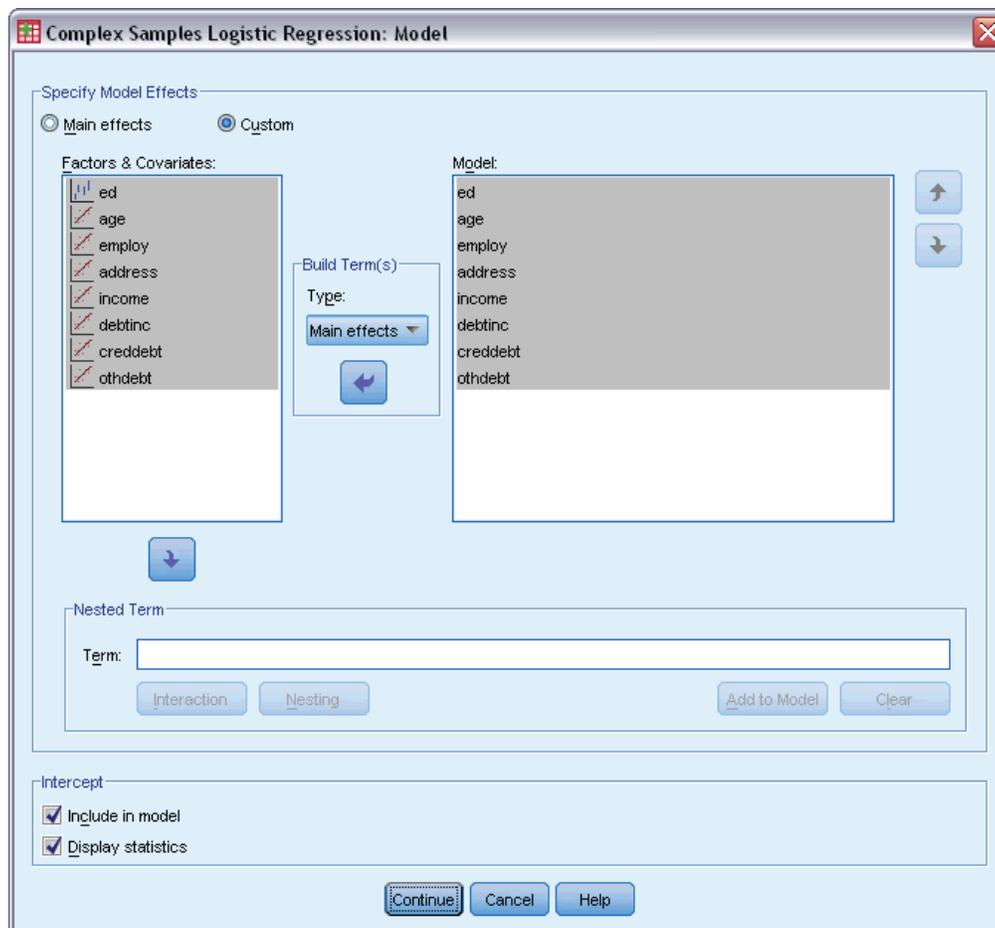
Figure 10-2  
Logistic Regression Reference Category dialog box



By default, the Complex Samples Logistic Regression procedure makes the highest-valued category the reference category. This dialog box allows you to specify the highest value, the lowest value, or a custom category as the reference category.

## Complex Samples Logistic Regression Model

Figure 10-3  
Logistic Regression Model dialog box



**Specify Model Effects.** By default, the procedure builds a main-effects model using the factors and covariates specified in the main dialog box. Alternatively, you can build a custom model that includes interaction effects and nested terms.

### Non-Nested Terms

For the selected factors and covariates:

**Interaction.** Creates the highest-level interaction term for all selected variables.

**Main effects.** Creates a main-effects term for each variable selected.

**All 2-way.** Creates all possible two-way interactions of the selected variables.

**All 3-way.** Creates all possible three-way interactions of the selected variables.

**All 4-way.** Creates all possible four-way interactions of the selected variables.

**All 5-way.** Creates all possible five-way interactions of the selected variables.

### ***Nested Terms***

You can build nested terms for your model in this procedure. Nested terms are useful for modeling the effect of a factor or covariate whose values do not interact with the levels of another factor. For example, a grocery store chain may follow the spending habits of its customers at several store locations. Since each customer frequents only one of these locations, the *Customer* effect can be said to be **nested within** the *Store location* effect.

Additionally, you can include interaction effects, such as polynomial terms involving the same covariate, or add multiple levels of nesting to the nested term.

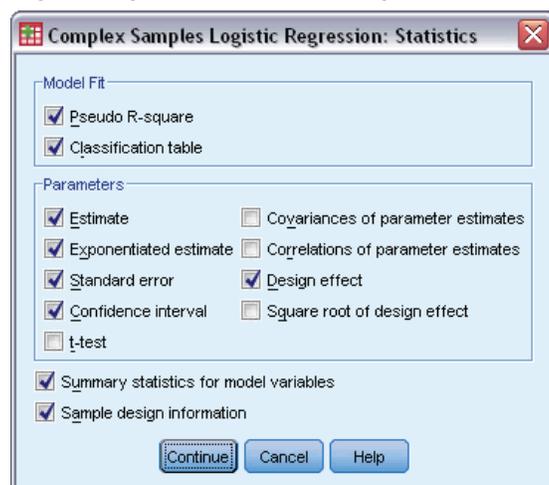
**Limitations.** Nested terms have the following restrictions:

- All factors within an interaction must be unique. Thus, if  $A$  is a factor, then specifying  $A*A$  is invalid.
- All factors within a nested effect must be unique. Thus, if  $A$  is a factor, then specifying  $A(A)$  is invalid.
- No effect can be nested within a covariate. Thus, if  $A$  is a factor and  $X$  is a covariate, then specifying  $A(X)$  is invalid.

**Intercept.** The intercept is usually included in the model. If you can assume the data pass through the origin, you can exclude the intercept. Even if you include the intercept in the model, you can choose to suppress statistics related to it.

## ***Complex Samples Logistic Regression Statistics***

Figure 10-4  
*Logistic Regression Statistics dialog box*



**Model Fit.** Controls the display of statistics that measure the overall model performance.

- **Pseudo R-square.** The  $R^2$  statistic from linear regression does not have an exact counterpart among logistic regression models. There are, instead, multiple measures that attempt to mimic the properties of the  $R^2$  statistic.
- **Classification table.** Displays the tabulated cross-classifications of the observed category by the model-predicted category on the dependent variable.

**Parameters.** This group allows you to control the display of statistics related to the model parameters.

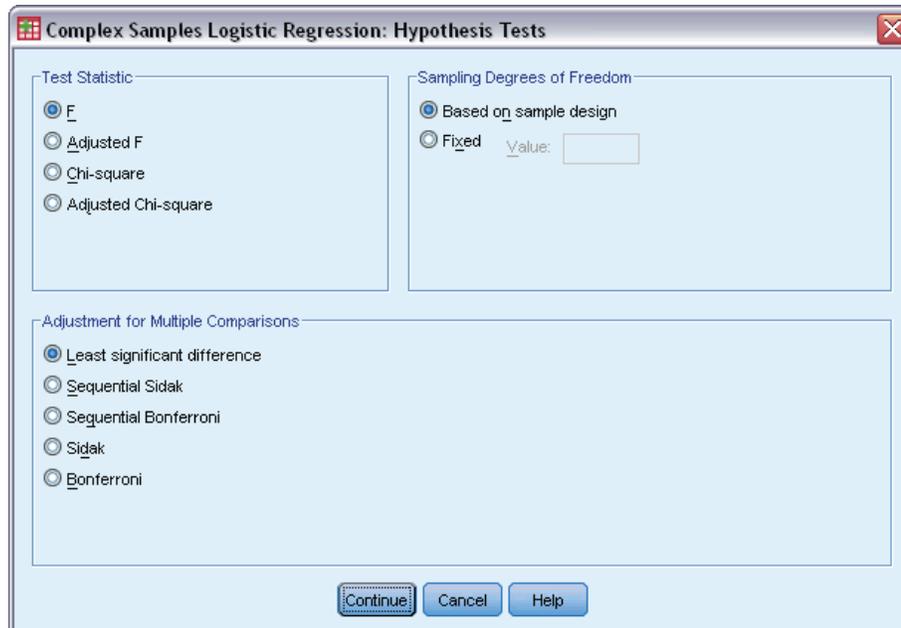
- **Estimate.** Displays estimates of the coefficients.
- **Exponentiated estimate.** Displays the base of the natural logarithm raised to the power of the estimates of the coefficients. While the estimate has nice properties for statistical testing, the exponentiated estimate, or  $\exp(B)$ , is easier to interpret.
- **Standard error.** Displays the standard error for each coefficient estimate.
- **Confidence interval.** Displays a confidence interval for each coefficient estimate. The confidence level for the interval is set in the Options dialog box.
- **T test.** Displays a  $t$  test of each coefficient estimate. The null hypothesis for each test is that the value of the coefficient is 0.
- **Covariances of parameter estimates.** Displays an estimate of the covariance matrix for the model coefficients.
- **Correlations of parameter estimates.** Displays an estimate of the correlation matrix for the model coefficients.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.

**Summary statistics for model variables.** Displays summary information about the dependent variable, covariates, and factors.

**Sample design information.** Displays summary information about the sample, including the unweighted count and the population size.

## Complex Samples Hypothesis Tests

Figure 10-5  
Hypothesis Tests dialog box



**Test Statistic.** This group allows you to select the type of statistic used for testing hypotheses. You can choose between  $F$ , adjusted  $F$ , chi-square, and adjusted chi-square.

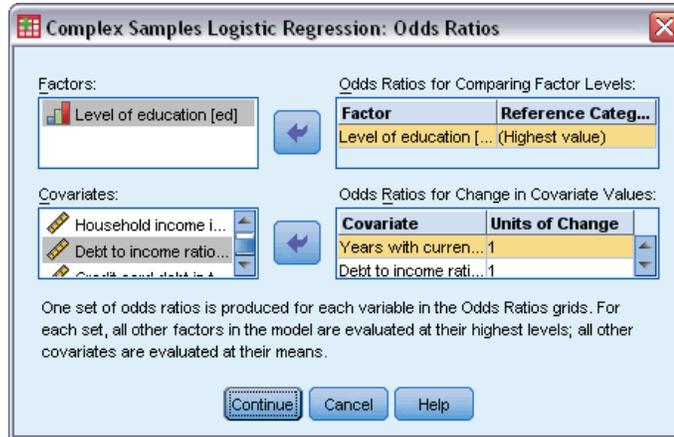
**Sampling Degrees of Freedom.** This group gives you control over the sampling design degrees of freedom used to compute  $p$  values for all test statistics. If based on the sampling design, the value is the difference between the number of primary sampling units and the number of strata in the first stage of sampling. Alternatively, you can set a custom degrees of freedom by specifying a positive integer.

**Adjustment for Multiple Comparisons.** When performing hypothesis tests with multiple contrasts, the overall significance level can be adjusted from the significance levels for the included contrasts. This group allows you to choose the adjustment method.

- **Least significant difference.** This method does not control the overall probability of rejecting the hypotheses that some linear contrasts are different from the null hypothesis values.
- **Sequential Sidak.** This is a sequentially step-down rejective Sidak procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sequential Bonferroni.** This is a sequentially step-down rejective Bonferroni procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sidak.** This method provides tighter bounds than the Bonferroni approach.
- **Bonferroni.** This method adjusts the observed significance level for the fact that multiple contrasts are being tested.

## Complex Samples Logistic Regression Odds Ratios

Figure 10-6  
Logistic Regression Odds Ratios dialog box



The Odds Ratios dialog box allows you to display the model-estimated odds ratios for specified factors and covariates. A separate set of odds ratios is computed for each category of the dependent variable except the reference category.

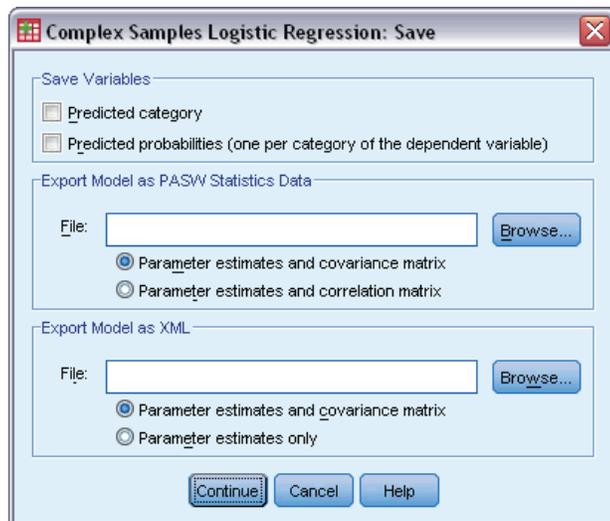
**Factors.** For each selected factor, displays the ratio of the odds at each category of the factor to the odds at the specified reference category.

**Covariates.** For each selected covariate, displays the ratio of the odds at the covariate's mean value plus the specified units of change to the odds at the mean.

When computing odds ratios for a factor or covariate, the procedure fixes all other factors at their highest levels and all other covariates at their means. If a factor or covariate interacts with other predictors in the model, then the odds ratios depend not only on the change in the specified variable but also on the values of the variables with which it interacts. If a specified covariate interacts with itself in the model (for example,  $age*age$ ), then the odds ratios depend on both the change in the covariate and the value of the covariate.

## Complex Samples Logistic Regression Save

Figure 10-7  
Logistic Regression Save dialog box



**Save Variables.** This group allows you to save the model-predicted category and predicted probabilities as new variables in the active dataset.

**Export model as SPSS Statistics data.** Writes a dataset in IBM® SPSS® Statistics format containing the parameter correlation or covariance matrix with parameter estimates, standard errors, significance values, and degrees of freedom. The order of variables in the matrix file is as follows.

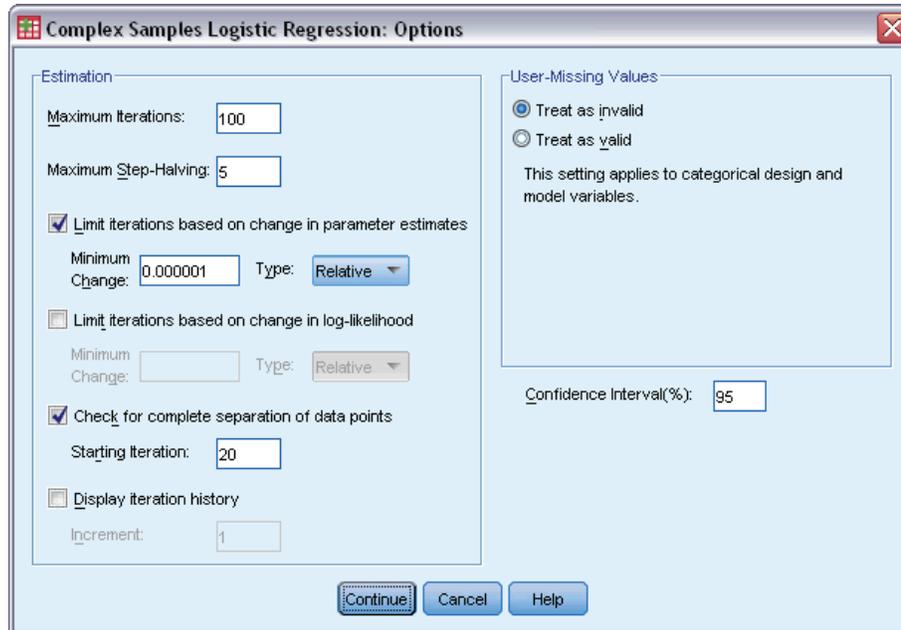
- **rowtype\_.** Takes values (and value labels), COV (Covariances), CORR (Correlations), EST (Parameter estimates), SE (Standard errors), SIG (Significance levels), and DF (Sampling design degrees of freedom). There is a separate case with row type COV (or CORR) for each model parameter, plus a separate case for each of the other row types.
- **varname\_.** Takes values P1, P2, ..., corresponding to an ordered list of all model parameters, for row types COV or CORR, with value labels corresponding to the parameter strings shown in the parameter estimates table. The cells are blank for other row types.
- **P1, P2, ...** These variables correspond to an ordered list of all model parameters, with variable labels corresponding to the parameter strings shown in the parameter estimates table, and take values according to the row type. For redundant parameters, all covariances are set to zero; correlations are set to the system-missing value; all parameter estimates are set at zero; and all standard errors, significance levels, and residual degrees of freedom are set to the system-missing value.

*Note:* This file is not immediately usable for further analyses in other procedures that read a matrix file unless those procedures accept all the row types exported here.

**Export Model as XML.** Saves the parameter estimates and the parameter covariance matrix, if selected, in XML (PMML) format. You can use this model file to apply the model information to other data files for scoring purposes.

## Complex Samples Logistic Regression Options

Figure 10-8  
Logistic Regression Options dialog box



**Estimation.** This group gives you control of various criteria used in the model estimation.

- **Maximum Iterations.** The maximum number of iterations the algorithm will execute. Specify a non-negative integer.
- **Maximum Step-Halving.** At each iteration, the step size is reduced by a factor of 0.5 until the log-likelihood increases or maximum step-halving is reached. Specify a positive integer.
- **Limit iterations based on change in parameter estimates.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the parameter estimates is less than the value specified, which must be non-negative.
- **Limit iterations based on change in log-likelihood.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the log-likelihood function is less than the value specified, which must be non-negative.
- **Check for complete separation of data points.** When selected, the algorithm performs tests to ensure that the parameter estimates have unique values. Separation occurs when the procedure can produce a model that correctly classifies every case.
- **Display iteration history.** Displays parameter estimates and statistics at every  $n$  iterations beginning with the 0<sup>th</sup> iteration (the initial estimates). If you choose to print the iteration history, the last iteration is always printed regardless of the value of  $n$ .

**User-Missing Values.** All design variables, as well as the dependent variable and any covariates, must have valid data. Cases with invalid data for any of these variables are deleted from the analysis. These controls allow you to decide whether user-missing values are treated as valid among the strata, cluster, subpopulation, and factor variables.

---

**Confidence Interval.** This is the confidence interval level for coefficient estimates, exponentiated coefficient estimates, and odds ratios. Specify a value greater than or equal to 50 and less than 100.

## ***CSLOGISTIC Command Additional Features***

The command syntax language also allows you to:

- Specify custom tests of effects versus a linear combination of effects or a value (using the `CUSTOM` subcommand).
- Fix values of other model variables when computing odds ratios for factors and covariates (using the `ODDSRATIOS` subcommand).
- Specify a tolerance value for checking singularity (using the `CRITERIA` subcommand).
- Create user-specified names for saved variables (using the `SAVE` subcommand).
- Produce a general estimable function table (using the `PRINT` subcommand).

See the *Command Syntax Reference* for complete syntax information.

# ***Complex Samples Ordinal Regression***

The Complex Samples Ordinal Regression procedure performs regression analysis on a binary or ordinal dependent variable for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

**Example.** Representatives considering a bill before the legislature are interested in whether there is public support for the bill and how support for the bill is related to voter demographics. Pollsters design and conduct interviews according to a complex sampling design. Using Complex Samples Ordinal Regression, you can fit a model for the level of support for the bill based upon voter demographics.

**Data.** The dependent variable is ordinal. Factors are categorical. Covariates are quantitative variables that are related to the dependent variable. Subpopulation variables can be string or numeric but should be categorical.

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

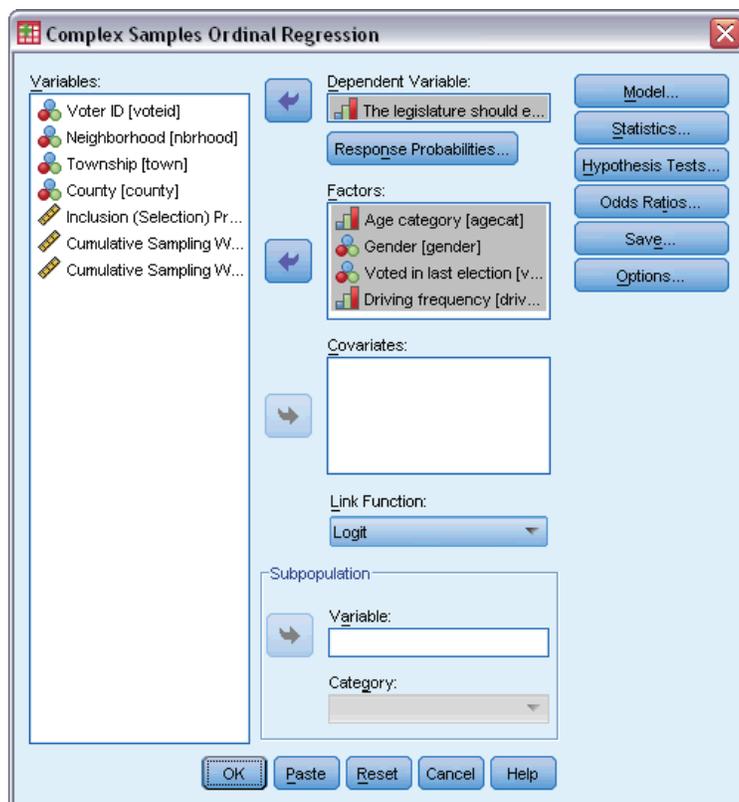
## ***Obtaining Complex Samples Ordinal Regression***

From the menus choose:

Analyze > Complex Samples > Ordinal Regression...

- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

Figure 11-1  
Ordinal Regression dialog box



- ▶ Select a dependent variable.

Optionally, you can:

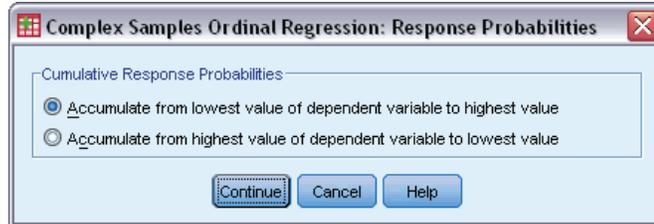
- Select variables for factors and covariates, as appropriate for your data.
- Specify a variable to define a subpopulation. The analysis is performed only for the selected category of the subpopulation variable, although variances are still properly estimated based on the entire dataset.
- Select a link function.

**Link function.** The link function is a transformation of the cumulative probabilities that allows estimation of the model. Five link functions are available, summarized in the following table.

Function	Form	Typical application
Logit	$\log(\xi / (1-\xi))$	Evenly distributed categories
Complementary log-log	$\log(-\log(1-\xi))$	Higher categories more probable
Negative log-log	$-\log(-\log(\xi))$	Lower categories more probable
Probit	$\Phi^{-1}(\xi)$	Latent variable is normally distributed
Cauchit (inverse Cauchy)	$\tan(\pi(\xi-0.5))$	Latent variable has many extreme values

## Complex Samples Ordinal Regression Response Probabilities

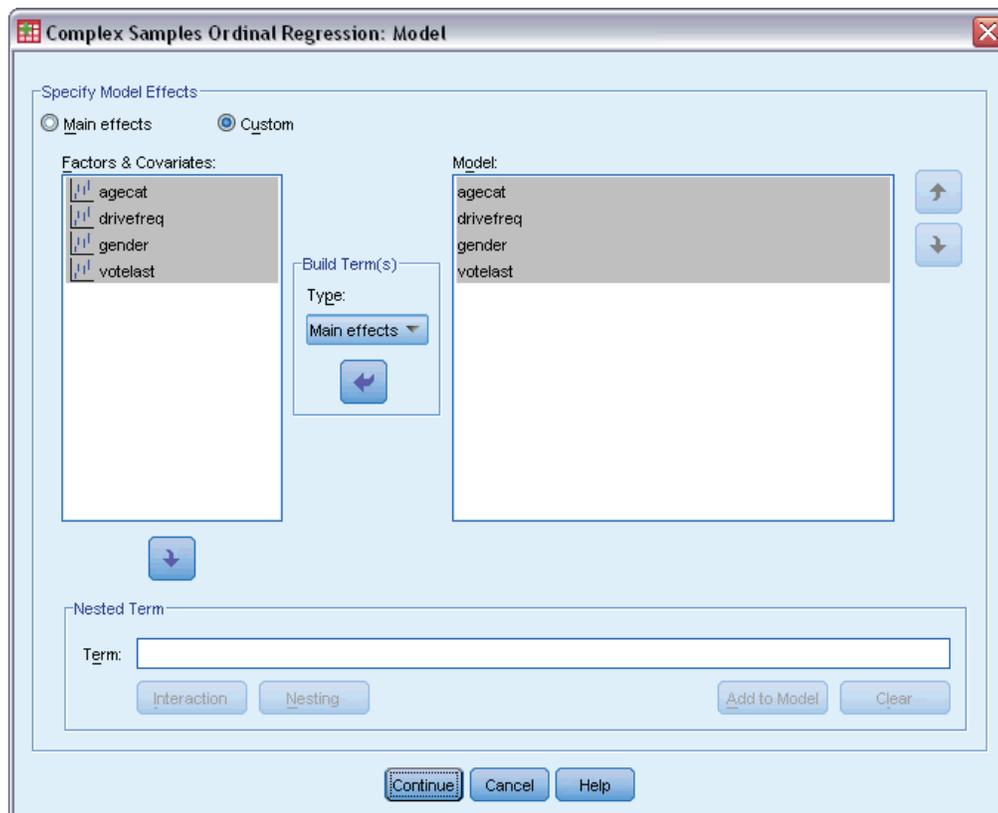
Figure 11-2  
Ordinal Regression Response Probabilities dialog box



The Response Probabilities dialog box allows you to specify whether the cumulative probability of a response (that is, the probability of belonging up to and including a particular category of the dependent variable) increases with increasing or decreasing values of the dependent variable.

## Complex Samples Ordinal Regression Model

Figure 11-3  
Ordinal Regression Model dialog box



**Specify Model Effects.** By default, the procedure builds a main-effects model using the factors and covariates specified in the main dialog box. Alternatively, you can build a custom model that includes interaction effects and nested terms.

### ***Non-Nested Terms***

For the selected factors and covariates:

**Interaction.** Creates the highest-level interaction term for all selected variables.

**Main effects.** Creates a main-effects term for each variable selected.

**All 2-way.** Creates all possible two-way interactions of the selected variables.

**All 3-way.** Creates all possible three-way interactions of the selected variables.

**All 4-way.** Creates all possible four-way interactions of the selected variables.

**All 5-way.** Creates all possible five-way interactions of the selected variables.

### ***Nested Terms***

You can build nested terms for your model in this procedure. Nested terms are useful for modeling the effect of a factor or covariate whose values do not interact with the levels of another factor.

For example, a grocery store chain may follow the spending habits of its customers at several store locations. Since each customer frequents only one of these locations, the *Customer* effect can be said to be **nested within** the *Store location* effect.

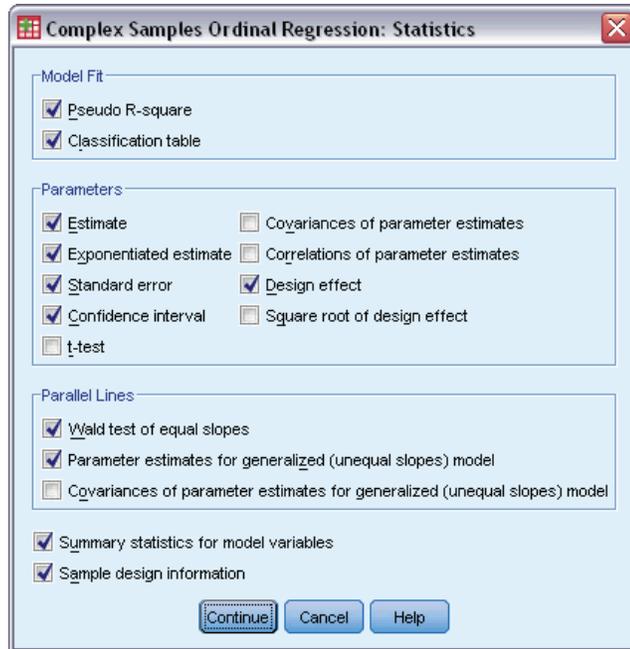
Additionally, you can include interaction effects, such as polynomial terms involving the same covariate, or add multiple levels of nesting to the nested term.

**Limitations.** Nested terms have the following restrictions:

- All factors within an interaction must be unique. Thus, if  $A$  is a factor, then specifying  $A*A$  is invalid.
- All factors within a nested effect must be unique. Thus, if  $A$  is a factor, then specifying  $A(A)$  is invalid.
- No effect can be nested within a covariate. Thus, if  $A$  is a factor and  $X$  is a covariate, then specifying  $A(X)$  is invalid.

## Complex Samples Ordinal Regression Statistics

Figure 11-4  
Ordinal Regression Statistics dialog box



**Model Fit.** Controls the display of statistics that measure the overall model performance.

- **Pseudo R-square.** The  $R^2$  statistic from linear regression does not have an exact counterpart among ordinal regression models. There are, instead, multiple measures that attempt to mimic the properties of the  $R^2$  statistic.
- **Classification table.** Displays the tabulated cross-classifications of the observed category by the model-predicted category on the dependent variable.

**Parameters.** This group allows you to control the display of statistics related to the model parameters.

- **Estimate.** Displays estimates of the coefficients.
- **Exponentiated estimate.** Displays the base of the natural logarithm raised to the power of the estimates of the coefficients. While the estimate has nice properties for statistical testing, the exponentiated estimate, or  $\exp(B)$ , is easier to interpret.
- **Standard error.** Displays the standard error for each coefficient estimate.
- **Confidence interval.** Displays a confidence interval for each coefficient estimate. The confidence level for the interval is set in the Options dialog box.
- **T test.** Displays a  $t$  test of each coefficient estimate. The null hypothesis for each test is that the value of the coefficient is 0.
- **Covariances of parameter estimates.** Displays an estimate of the covariance matrix for the model coefficients.
- **Correlations of parameter estimates.** Displays an estimate of the correlation matrix for the model coefficients.

- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure, expressed in units comparable to those of the standard error, of the effect of specifying a complex design, where values further from 1 indicate greater effects.

**Parallel Lines.** This group allows you to request statistics associated with a model with nonparallel lines where a separate regression line is fitted for each response category (except the last).

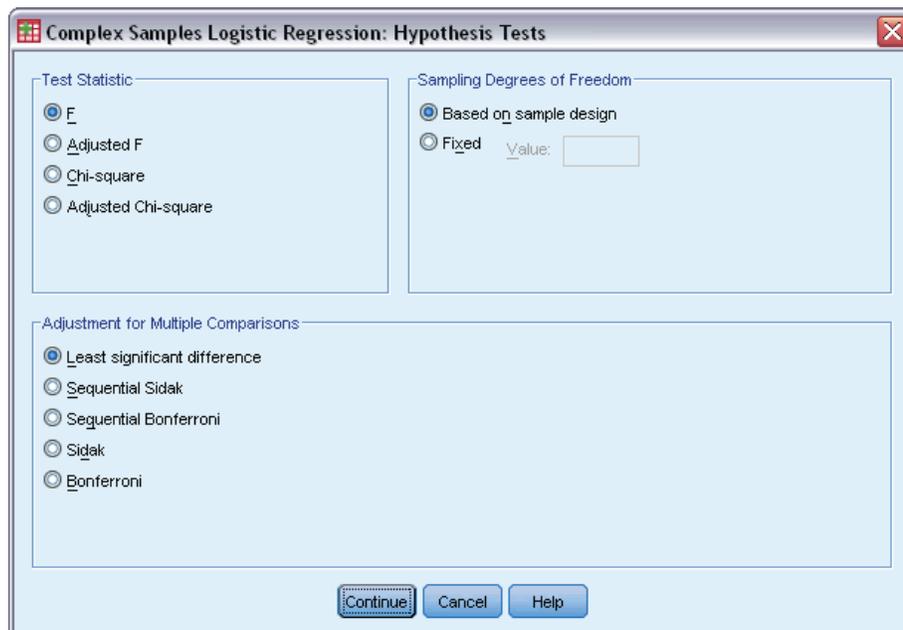
- **Wald test.** Produces a test of the null hypothesis that regression parameters are equal for all cumulative responses. The model with nonparallel lines is estimated and the Wald test of equal parameters is applied.
- **Parameter estimates.** Displays estimates of the coefficients and standard errors for the model with nonparallel lines.
- **Covariances of parameter estimates.** Displays an estimate of the covariance matrix for the coefficients of the model with nonparallel lines.

**Summary statistics for model variables.** Displays summary information about the dependent variable, covariates, and factors.

**Sample design information.** Displays summary information about the sample, including the unweighted count and the population size.

## Complex Samples Hypothesis Tests

Figure 11-5  
Hypothesis Tests dialog box



**Test Statistic.** This group allows you to select the type of statistic used for testing hypotheses. You can choose between  $F$ , adjusted  $F$ , chi-square, and adjusted chi-square.

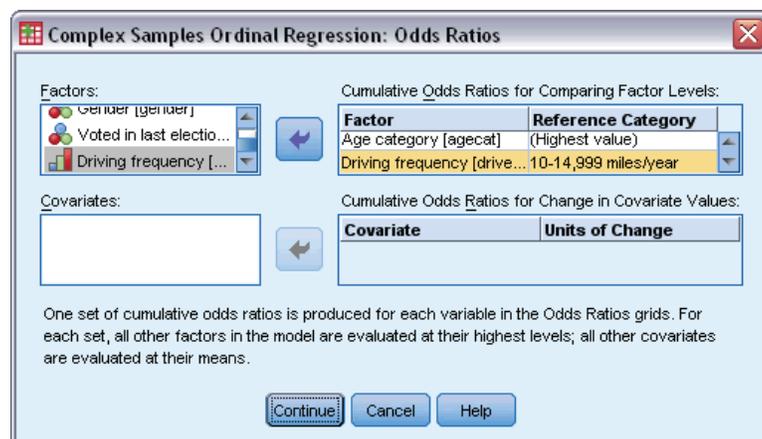
**Sampling Degrees of Freedom.** This group gives you control over the sampling design degrees of freedom used to compute  $p$  values for all test statistics. If based on the sampling design, the value is the difference between the number of primary sampling units and the number of strata in the first stage of sampling. Alternatively, you can set a custom degrees of freedom by specifying a positive integer.

**Adjustment for Multiple Comparisons.** When performing hypothesis tests with multiple contrasts, the overall significance level can be adjusted from the significance levels for the included contrasts. This group allows you to choose the adjustment method.

- **Least significant difference.** This method does not control the overall probability of rejecting the hypotheses that some linear contrasts are different from the null hypothesis values.
- **Sequential Sidak.** This is a sequentially step-down rejective Sidak procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sequential Bonferroni.** This is a sequentially step-down rejective Bonferroni procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sidak.** This method provides tighter bounds than the Bonferroni approach.
- **Bonferroni.** This method adjusts the observed significance level for the fact that multiple contrasts are being tested.

## Complex Samples Ordinal Regression Odds Ratios

Figure 11-6  
Ordinal Regression Odds Ratios dialog box



The Odds Ratios dialog box allows you to display the model-estimated cumulative odds ratios for specified factors and covariates. This feature is only available for models using the Logit link function. A single cumulative odds ratio is computed for all categories of the dependent variable except the last; the proportional odds model postulates that they are all equal.

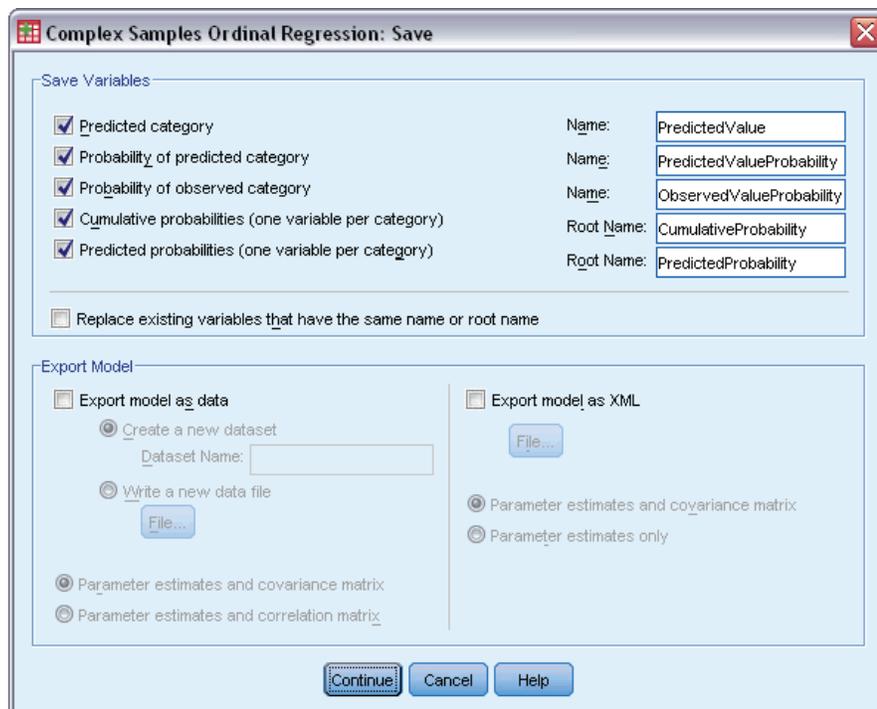
**Factors.** For each selected factor, displays the ratio of the cumulative odds at each category of the factor to the odds at the specified reference category.

**Covariates.** For each selected covariate, displays the ratio of the cumulative odds at the covariate's mean value plus the specified units of change to the odds at the mean.

When computing odds ratios for a factor or covariate, the procedure fixes all other factors at their highest levels and all other covariates at their means. If a factor or covariate interacts with other predictors in the model, then the odds ratios depend not only on the change in the specified variable but also on the values of the variables with which it interacts. If a specified covariate interacts with itself in the model (for example,  $age*age$ ), then the odds ratios depend on both the change in the covariate and the value of the covariate.

## Complex Samples Ordinal Regression Save

Figure 11-7  
Ordinal Regression Save dialog box



**Save Variables.** This group allows you to save the model-predicted category, probability of predicted category, probability of observed category, cumulative probabilities, and predicted probabilities as new variables in the active dataset.

**Export model as SPSS Statistics data.** Writes a dataset in IBM® SPSS® Statistics format containing the parameter correlation or covariance matrix with parameter estimates, standard errors, significance values, and degrees of freedom. The order of variables in the matrix file is as follows.

- **rowtype\_.** Takes values (and value labels), COV (Covariances), CORR (Correlations), EST (Parameter estimates), SE (Standard errors), SIG (Significance levels), and DF (Sampling design degrees of freedom). There is a separate case with row type COV (or CORR) for each model parameter, plus a separate case for each of the other row types.
- **varname\_.** Takes values P1, P2, ..., corresponding to an ordered list of all model parameters, for row types COV or CORR, with value labels corresponding to the parameter strings shown in the parameter estimates table. The cells are blank for other row types.
- **P1, P2, ...** These variables correspond to an ordered list of all model parameters, with variable labels corresponding to the parameter strings shown in the parameter estimates table, and take values according to the row type. For redundant parameters, all covariances are set to zero; correlations are set to the system-missing value; all parameter estimates are set at zero; and all standard errors, significance levels, and residual degrees of freedom are set to the system-missing value.

*Note:* This file is not immediately usable for further analyses in other procedures that read a matrix file unless those procedures accept all the row types exported here.

**Export model as XML.** Saves the parameter estimates and the parameter covariance matrix, if selected, in XML (PMML) format. You can use this model file to apply the model information to other data files for scoring purposes.

## Complex Samples Ordinal Regression Options

Figure 11-8  
Ordinal Regression Options dialog box

**Complex Samples Ordinal Regression: Options**

**Estimation Method**

Newton-Raphson

Fisher scoring

Fisher scoring then Newton-Raphson

Maximum Number of Iteration Before Switching:

**Estimation Criteria**

Maximum Iterations:

Maximum Step-Halving:

Limit iterations based on change in parameter estimates

Minimum Change:  Type: **Relative**

Limit iterations based on change in log-likelihood

Minimum Change:  Type: **Relative**

Check for complete separation of data points

Starting Iteration:

Display iteration history

Increment:

**User-Missing Values**

Treat as invalid

Treat as valid

This setting applies to categorical design and model variables.

Confidence Interval(%):

**Continue** **Cancel** **Help**

**Estimation Method.** You can select a parameter estimation method; choose between Newton-Raphson, Fisher scoring, or a hybrid method in which Fisher scoring iterations are performed before switching to the Newton-Raphson method. If convergence is achieved during the Fisher scoring phase of the hybrid method before the maximum number of Fisher iterations is reached, the algorithm continues with the Newton-Raphson method.

**Estimation.** This group gives you control of various criteria used in the model estimation.

- **Maximum Iterations.** The maximum number of iterations the algorithm will execute. Specify a non-negative integer.
- **Maximum Step-Halving.** At each iteration, the step size is reduced by a factor of 0.5 until the log-likelihood increases or maximum step-halving is reached. Specify a positive integer.
- **Limit iterations based on change in parameter estimates.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the parameter estimates is less than the value specified, which must be non-negative.
- **Limit iterations based on change in log-likelihood.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the log-likelihood function is less than the value specified, which must be non-negative.
- **Check for complete separation of data points.** When selected, the algorithm performs tests to ensure that the parameter estimates have unique values. Separation occurs when the procedure can produce a model that correctly classifies every case.
- **Display iteration history.** Displays parameter estimates and statistics at every  $n$  iterations beginning with the 0<sup>th</sup> iteration (the initial estimates). If you choose to print the iteration history, the last iteration is always printed regardless of the value of  $n$ .

**User-Missing Values.** Scale design variables, as well as the dependent variable and any covariates, should have valid data. Cases with invalid data for any of these variables are deleted from the analysis. These controls allow you to decide whether user-missing values are treated as valid among the strata, cluster, subpopulation, and factor variables.

**Confidence Interval.** This is the confidence interval level for coefficient estimates, exponentiated coefficient estimates, and odds ratios. Specify a value greater than or equal to 50 and less than 100.

## ***CSORDINAL Command Additional Features***

The command syntax language also allows you to:

- Specify custom tests of effects versus a linear combination of effects or a value (using the `CUSTOM` subcommand).
- Fix values of other model variables at values other than their means when computing cumulative odds ratios for factors and covariates (using the `ODDSRATIOS` subcommand).
- Use unlabeled values as custom reference categories for factors when odds ratios are requested (using the `ODDSRATIOS` subcommand).
- Specify a tolerance value for checking singularity (using the `CRITERIA` subcommand).
- Produce a general estimable function table (using the `PRINT` subcommand).
- Save more than 25 probability variables (using the `SAVE` subcommand).

See the *Command Syntax Reference* for complete syntax information.

# ***Complex Samples Cox Regression***

The Complex Samples Cox Regression procedure performs survival analysis for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

**Examples.** A government law enforcement agency is concerned about recidivism rates in their area of jurisdiction. One of the measures of recidivism is the time until second arrest for offenders. The agency would like to model time to rearrest using Cox Regression but are worried the proportional hazards assumption is invalid across age categories.

Medical researchers are investigating survival times for patients exiting a rehabilitation program post-ischemic stroke. There is the potential for multiple cases per subject, since patient histories change as the occurrence of significant nondeath events are noted and the times of these events recorded. The sample is also left-truncated in the sense that the observed survival times are “inflated” by the length of rehabilitation, because while the onset of risk starts at the time of the ischemic stroke, only patients who survive past the rehabilitation program are in the sample.

**Survival Time.** The procedure applies Cox regression to analysis of survival times—that is, the length of time before the occurrence of an event. There are two ways to specify the survival time, depending upon the start time of the interval:

- **Time=0.** Commonly, you will have complete information on the start of the interval for each subject and will simply have a variable containing end times (or create a single variable with end times from Date & Time variables; see below).
- **Varies by subject.** This is appropriate when you have **left-truncation**, also called **delayed entry**; for example, if you are analyzing survival times for patients exiting a rehabilitation program post-stroke, you might consider that their onset of risk starts at the time of the stroke. However, if your sample only includes patients who have survived the rehabilitation program, then your sample is left-truncated in the sense that the observed survival times are “inflated” by the length of rehabilitation. You can account for this by specifying the time at which they exited rehabilitation as the time of entry into the study.

**Date & Time Variables.** Date & Time variables cannot be used to directly define the start and end of the interval; if you have Date & Time variables, you should use them to create variables containing survival times. If there is no left-truncation, simply create a variable containing end times based upon the difference between the date of entry into the study and the observation date. If there is left-truncation, create a variable containing start times, based upon the difference between the date of the start of the study and the date of entry, and a variable containing end times, based upon the difference between the date of the start of the study and the date of observation.

**Event Status.** You need a variable that records whether the subject experienced the event of interest within the interval. Subjects for whom the event has not occurred are right-censored.

**Subject Identifier.** You can easily incorporate piecewise-constant, time-dependent predictors by splitting the observations for a single subject across multiple cases. For example, if you are analyzing survival times for patients post-stroke, variables representing their medical history should be useful as predictors. Over time, they may experience major medical events that alter their medical history. The following table shows how to structure such a dataset: *Patient ID* is the subject identifier, *End time* defines the observed intervals, *Status* records major medical events, and *Prior history of heart attack* and *Prior history of hemorrhaging* are piecewise-constant, time-dependent predictors.

<i>Patient ID</i>	<i>End time</i>	<i>Status</i>	<i>Prior history of heart attack</i>	<i>Prior history of hemorrhaging</i>
1	5	Heart Attack	No	No
1	7	Hemorrhaging	Yes	No
1	8	Died	Yes	Yes
2	24	Died	No	No
3	8	Heart Attack	No	No
3	15	Died	Yes	No

**Assumptions.** The cases in the data file represent a sample from a complex design that should be analyzed according to the specifications in the file selected in the [Complex Samples Plan dialog box](#).

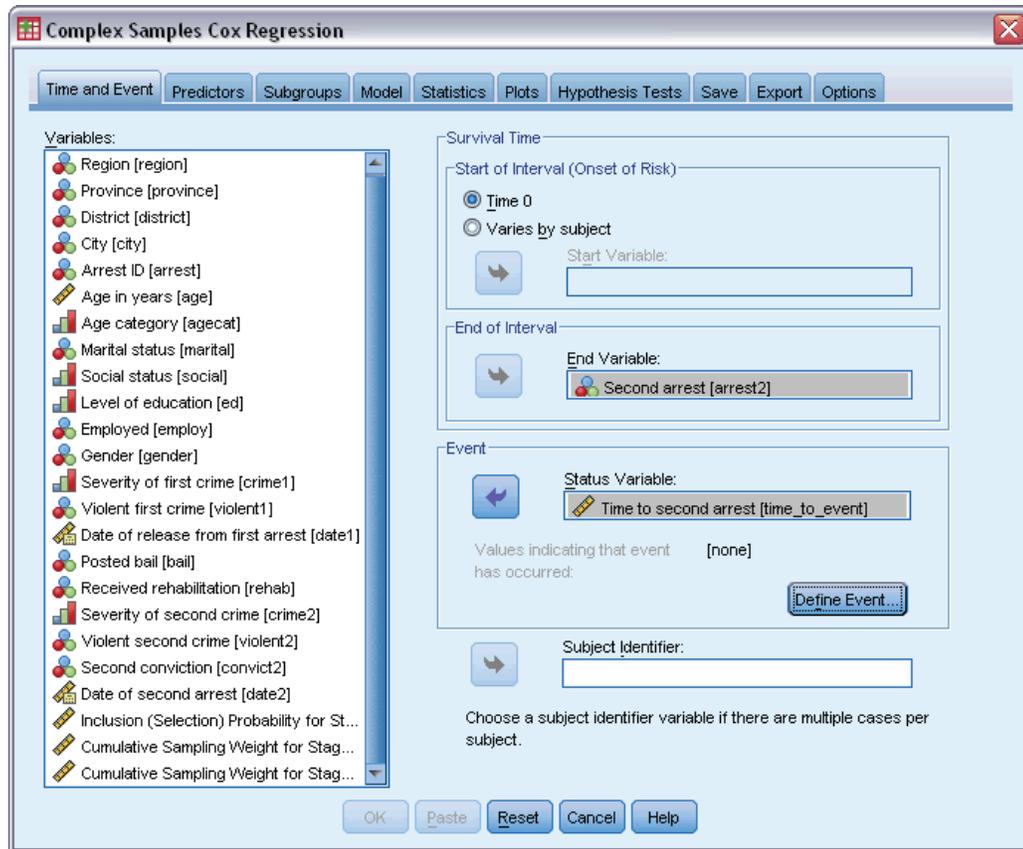
Typically, Cox regression models assume proportional hazards—that is, the ratio of hazards from one case to another should not vary over time. If this assumption does not hold, you may need to add time-dependent predictors to the model.

**Kaplan-Meier Analysis.** If you do not select any predictors (or do not enter any selected predictors into the model) and choose the product limit method for computing the baseline survival curve on the Options tab, the procedure performs a Kaplan-Meier type of survival analysis.

#### **To Obtain Complex Samples Cox Regression**

- ▶ From the menus choose:  
Analyze > Complex Samples > Cox Regression...
- ▶ Select a plan file. Optionally, select a custom joint probabilities file.
- ▶ Click Continue.

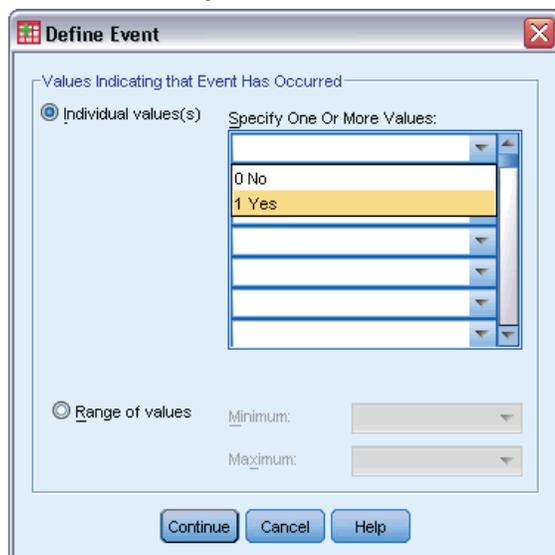
Figure 12-1  
Cox Regression dialog box, Time and Event tab



- ▶ Specify the survival time by selecting the entry and exit times from the study.
  - ▶ Select an event status variable.
  - ▶ Click **Define Event** and define at least one event value.
- Optionally, you can select a subject identifier.

## Define Event

Figure 12-2  
Define Event dialog box

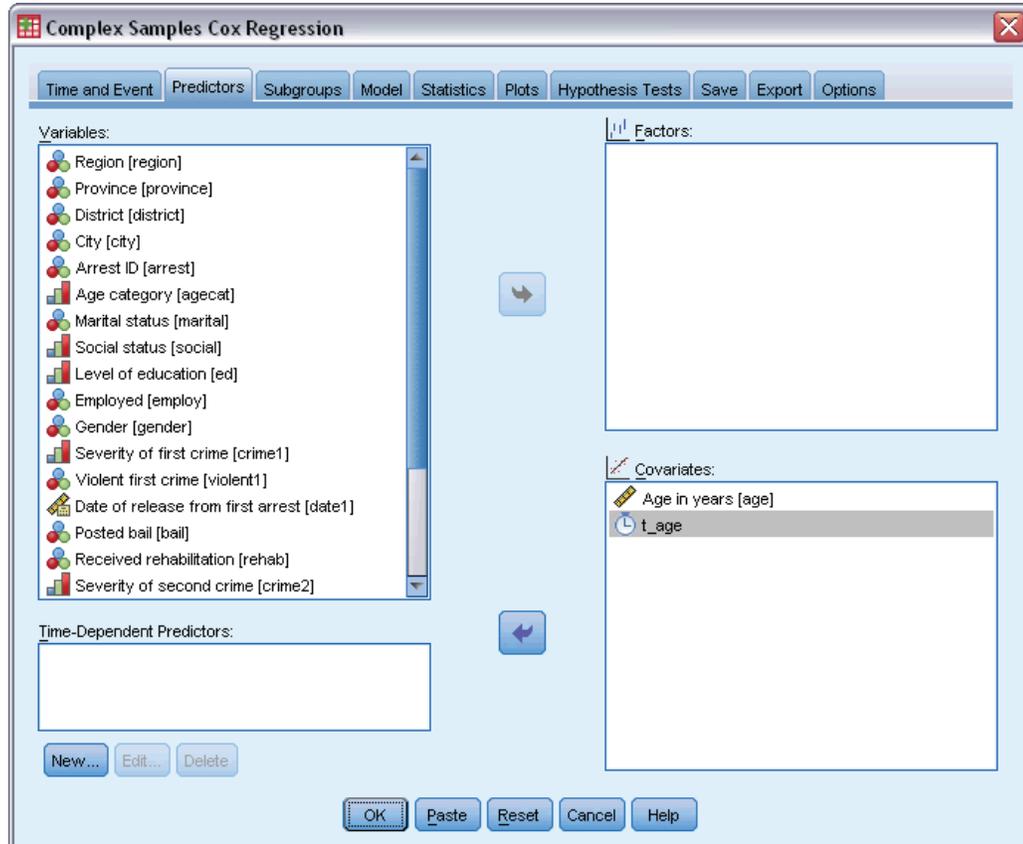


Specify the values that indicate a terminal event has occurred.

- **Individual value(s).** Specify one or more values by entering them into the grid or selecting them from a list of values with defined value labels.
- **Range of values.** Specify a range of values by entering the minimum and maximum values or selecting values from a list with defined value labels.

## Predictors

Figure 12-3  
Cox Regression dialog box, Predictors tab



The Predictors tab allows you to specify the factors and covariates used to build model effects.

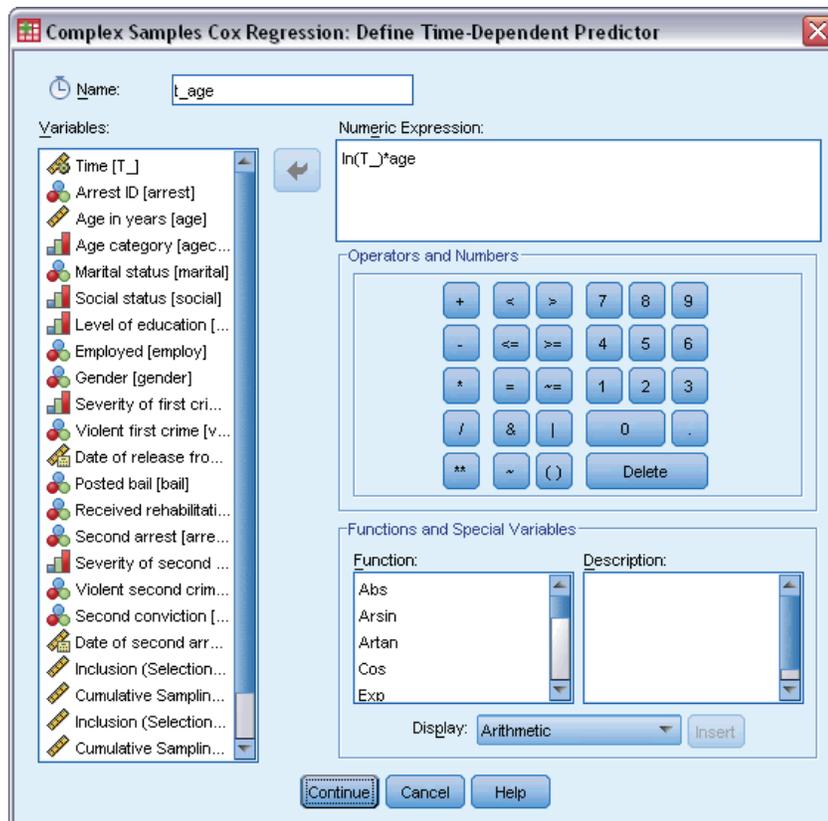
**Factors.** Factors are categorical predictors; they can be numeric or string.

**Covariates.** Covariates are scale predictors; they must be numeric.

**Time-Dependent Predictors.** There are certain situations in which the proportional hazards assumption does not hold. That is, hazard ratios change across time; the values of one (or more) of your predictors are different at different time points. In such cases, you need to specify time-dependent predictors. [For more information, see the topic Define Time-Dependent Predictor on p. 79.](#) Time-dependent predictors can be selected as factors or covariates.

## Define Time-Dependent Predictor

Figure 12-4  
Cox Regression Define Time-Dependent Predictor dialog box



The Define Time-Dependent Predictor dialog box allows you to create a predictor that is dependent upon the built-in time variable,  $T_$ . You can use this variable to define time-dependent covariates in two general ways:

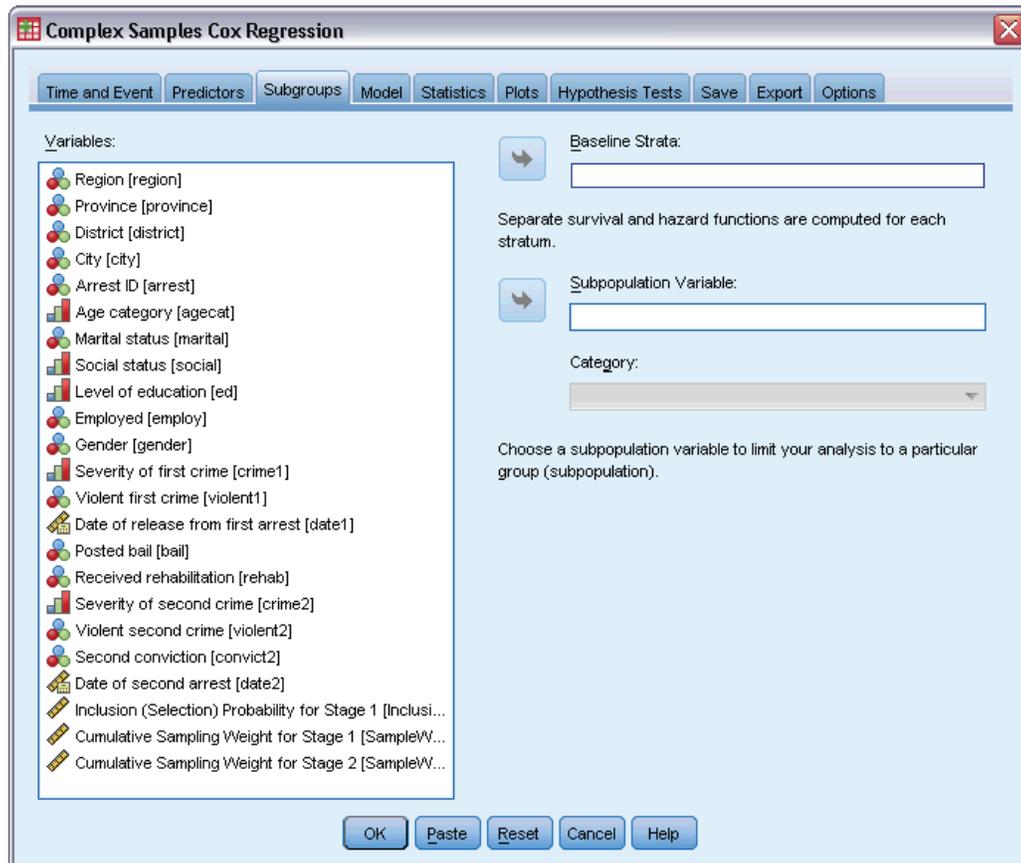
- If you want to estimate an extended Cox regression model that allows nonproportional hazards, you can do so by defining your time-dependent predictor as a function of the time variable  $T_$  and the covariate in question. A common example would be the simple product of the time variable and the predictor, but more complex functions can be specified as well.
- Some variables may have different values at different time periods but aren't systematically related to time. In such cases, you need to define a **segmented time-dependent predictor**, which can be done using logical expressions. Logical expressions take the value 1 if true and 0 if false. Using a series of logical expressions, you can create your time-dependent predictor from a set of measurements. For example, if you have blood pressure measured once a week for the four weeks of your study (identified as  $BP1$  to  $BP4$ ), you can define your time-dependent predictor as  $(T_ < 1) * BP1 + (T_ \geq 1 \ \& \ T_ < 2) * BP2 + (T_ \geq 2 \ \& \ T_ < 3) * BP3 + (T_ \geq 3 \ \& \ T_ < 4) * BP4$ . Notice that exactly one of the terms in parentheses will be equal to 1 for any given case and the rest will all equal 0. In other words, this function means that if time is less than one week, use  $BP1$ ; if it is more than one week but less than two weeks, use  $BP2$ ; and so on.

*Note:* If your segmented, time-dependent predictor is constant within segments, as in the blood pressure example given above, it may be easier for you to specify the piecewise-constant, time-dependent predictor by splitting subjects across multiple cases. See the discussion on Subject Identifiers in [Complex Samples Cox Regression](#) on p. 74 for more information.

In the Define Time-Dependent Predictor dialog box, you can use the function-building controls to build the expression for the time-dependent covariate, or you can enter it directly in the Numeric Expression text area. Note that string constants must be enclosed in quotation marks or apostrophes, and numeric constants must be typed in American format, with the dot as the decimal delimiter. The resulting variable is given the name you specify and should be included as a factor or covariate on the Predictors tab.

## Subgroups

Figure 12-5  
Cox Regression dialog box, Subgroups tab

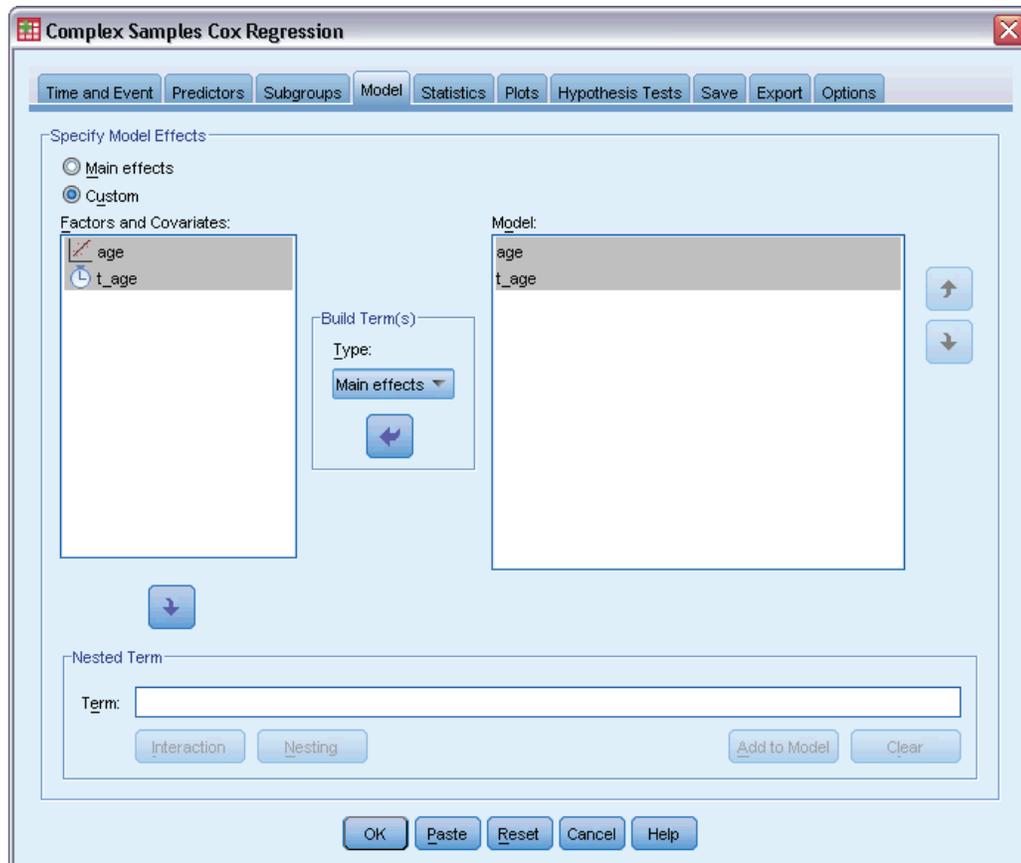


**Baseline Strata.** A separate baseline hazard and survival function is computed for each value of this variable, while a single set of model coefficients is estimated across strata.

**Subpopulation Variable.** Specify a variable to define a subpopulation. The analysis is performed only for the selected category of the subpopulation variable.

## Model

Figure 12-6  
Cox Regression dialog box, Model tab



**Specify Model Effects.** By default, the procedure builds a main-effects model using the factors and covariates specified in the main dialog box. Alternatively, you can build a custom model that includes interaction effects and nested terms.

### **Non-Nested Terms**

For the selected factors and covariates:

**Interaction.** Creates the highest-level interaction term for all selected variables.

**Main effects.** Creates a main-effects term for each variable selected.

**All 2-way.** Creates all possible two-way interactions of the selected variables.

**All 3-way.** Creates all possible three-way interactions of the selected variables.

**All 4-way.** Creates all possible four-way interactions of the selected variables.

**All 5-way.** Creates all possible five-way interactions of the selected variables.

### Nested Terms

You can build nested terms for your model in this procedure. Nested terms are useful for modeling the effect of a factor or covariate whose values do not interact with the levels of another factor. For example, a grocery store chain may follow the spending habits of its customers at several store locations. Since each customer frequents only one of these locations, the *Customer* effect can be said to be **nested within** the *Store location* effect.

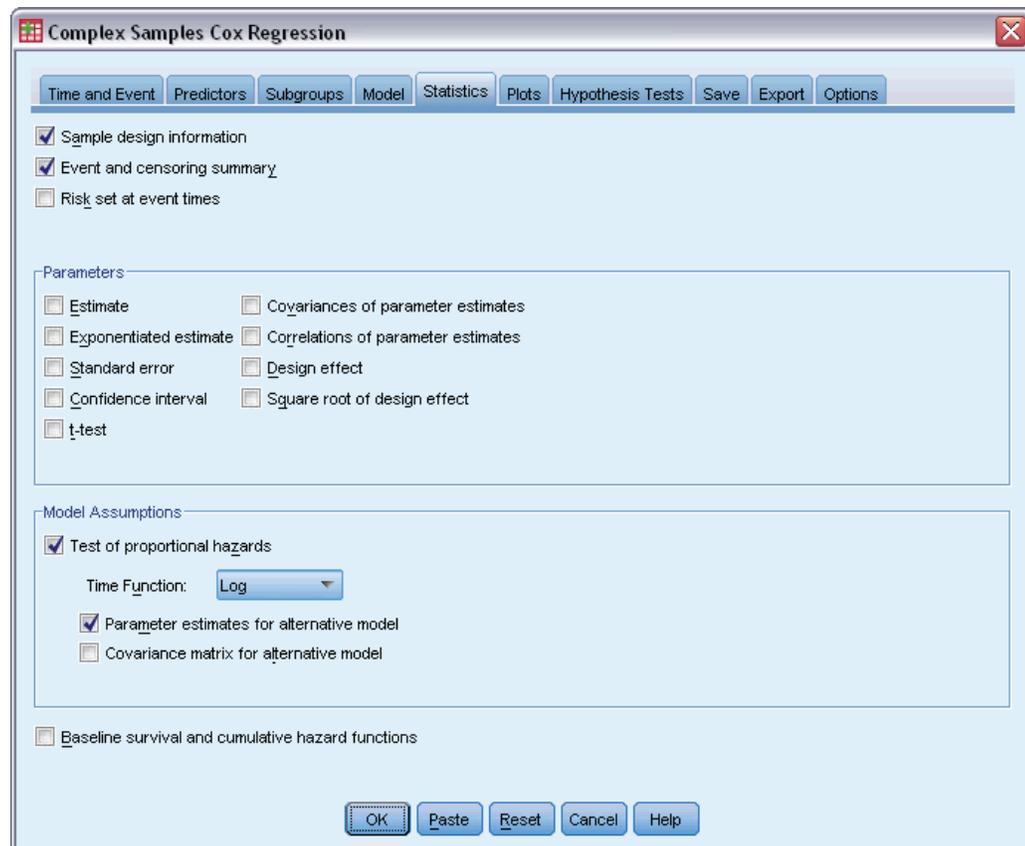
Additionally, you can include interaction effects, such as polynomial terms involving the same covariate, or add multiple levels of nesting to the nested term.

**Limitations.** Nested terms have the following restrictions:

- All factors within an interaction must be unique. Thus, if  $A$  is a factor, then specifying  $A*A$  is invalid.
- All factors within a nested effect must be unique. Thus, if  $A$  is a factor, then specifying  $A(A)$  is invalid.
- No effect can be nested within a covariate. Thus, if  $A$  is a factor and  $X$  is a covariate, then specifying  $A(X)$  is invalid.

## Statistics

Figure 12-7  
Cox Regression dialog box, Statistics tab



**Sample design information.** Displays summary information about the sample, including the unweighted count and the population size.

**Event and censoring summary.** Displays summary information about the number and percentage of censored cases.

**Risk set at event times.** Displays number of events and number at risk for each event time in each baseline stratum.

**Parameters.** This group allows you to control the display of statistics related to the model parameters.

- **Estimate.** Displays estimates of the coefficients.
- **Exponentiated estimate.** Displays the base of the natural logarithm raised to the power of the estimates of the coefficients. While the estimate has nice properties for statistical testing, the exponentiated estimate, or  $\exp(B)$ , is easier to interpret.
- **Standard error.** Displays the standard error for each coefficient estimate.
- **Confidence interval.** Displays a confidence interval for each coefficient estimate. The confidence level for the interval is set in the Options dialog box.
- **t-test.** Displays a  $t$  test of each coefficient estimate. The null hypothesis for each test is that the value of the coefficient is 0.
- **Covariances of parameter estimates.** Displays an estimate of the covariance matrix for the model coefficients.
- **Correlations of parameter estimates.** Displays an estimate of the correlation matrix for the model coefficients.
- **Design effect.** The ratio of the variance of the estimate to the variance obtained by assuming that the sample is a simple random sample. This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.
- **Square root of design effect.** This is a measure of the effect of specifying a complex design, where values further from 1 indicate greater effects.

**Model Assumptions.** This group allows you to produce a test of the proportional hazards assumption. The test compares the fitted model to an alternative model that includes time-dependent predictors  $x*_TF$  for each predictor  $x$ , where  $_TF$  is the specified time function.

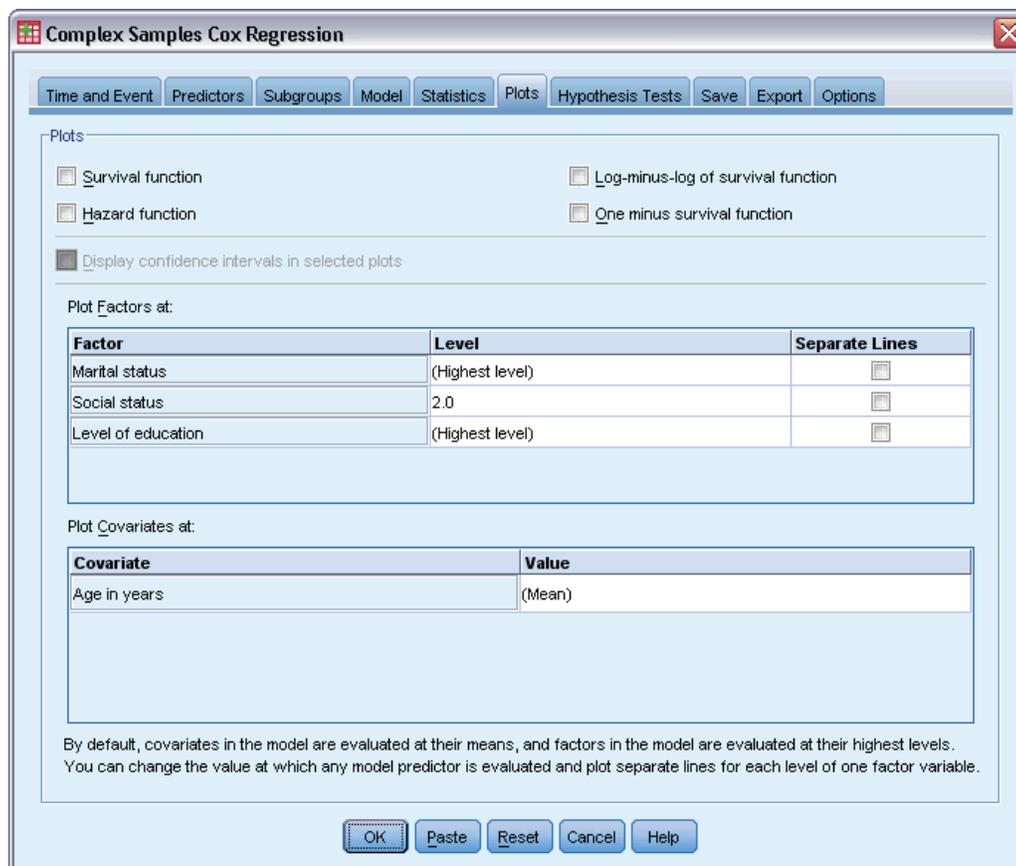
- **Time Function.** Specifies the form of  $_TF$  for the alternative model. For the **identity** function,  $_TF=T_$ . For the **log** function,  $_TF=\log(T_)$ . For **Kaplan-Meier**,  $_TF=1-S_{KM}(T_)$ , where  $S_{KM}(\cdot)$  is the Kaplan-Meier estimate of the survival function. For **rank**,  $_TF$  is the rank-order of  $T_$  among the observed end times.
- **Parameter estimates for alternative model.** Displays the estimate, standard error, and confidence interval for each parameter in the alternative model.
- **Covariance matrix for alternative model.** Displays the matrix of estimated covariances between parameters in the alternative model.

**Baseline survival and cumulative hazard functions.** Displays the baseline survival function and baseline cumulative hazards function along with their standard errors.

*Note:* If time-dependent predictors defined on the Predictors tab are included in the model, this option is not available.

## Plots

Figure 12-8  
Cox Regression dialog box, Plots tab



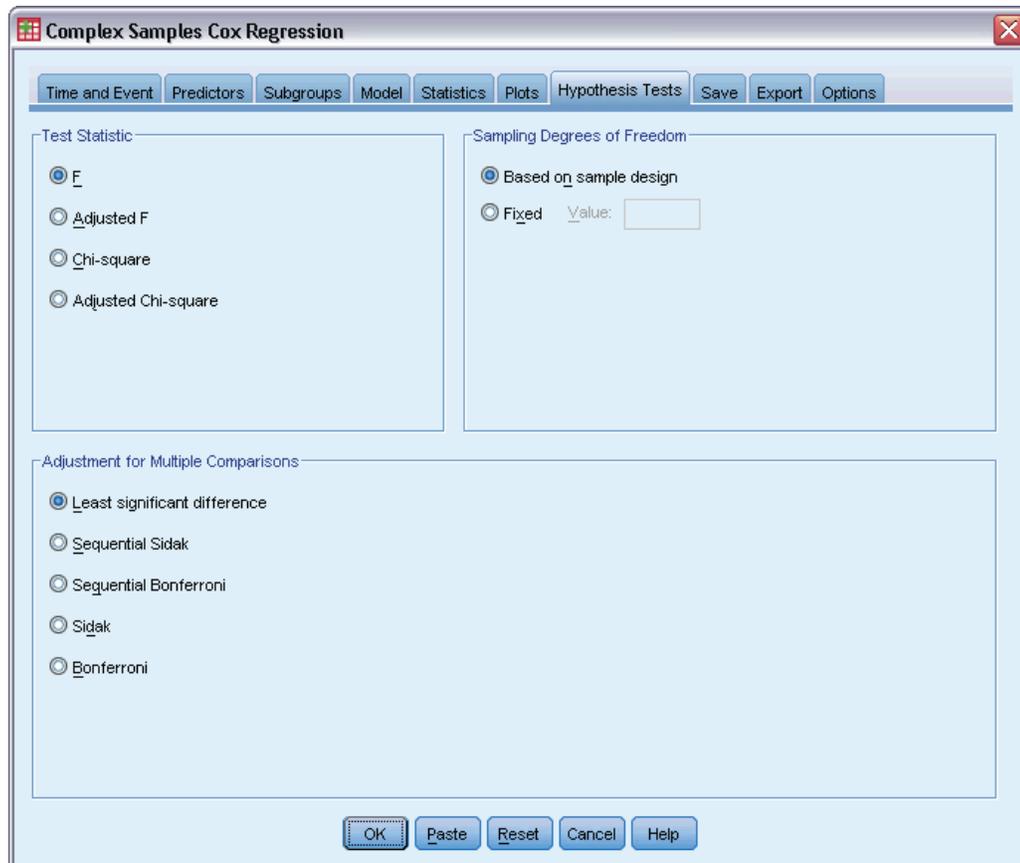
The Plots tab allows you to request plots of the hazard function, survival function, log-minus-log of the survival function, and one minus the survival function. You can also choose to plot confidence intervals along the specified functions; the confidence level is set on the Options tab.

**Predictor patterns.** You can specify a pattern of predictor values to be used for the requested plots and the exported survival file on the Export tab. Note that these options are not available if time-dependent predictors defined on the Predictors tab are included in the model.

- **Plot Factors at.** By default, each factor is evaluated at its highest level. Enter or select a different level if desired. Alternatively, you can choose to plot separate lines for each level of a single factor by selecting the check box for that factor.
- **Plot Covariates at.** Each covariate is evaluated at its mean. Enter or select a different value if desired.

## Hypothesis Tests

Figure 12-9  
Cox Regression dialog box, Hypothesis Tests tab



**Test Statistic.** This group allows you to select the type of statistic used for testing hypotheses. You can choose between  $F$ , adjusted  $F$ , chi-square, and adjusted chi-square.

**Sampling Degrees of Freedom.** This group gives you control over the sampling design degrees of freedom used to compute  $p$  values for all test statistics. If based on the sampling design, the value is the difference between the number of primary sampling units and the number of strata in the first stage of sampling. Alternatively, you can set a custom degrees of freedom by specifying a positive integer.

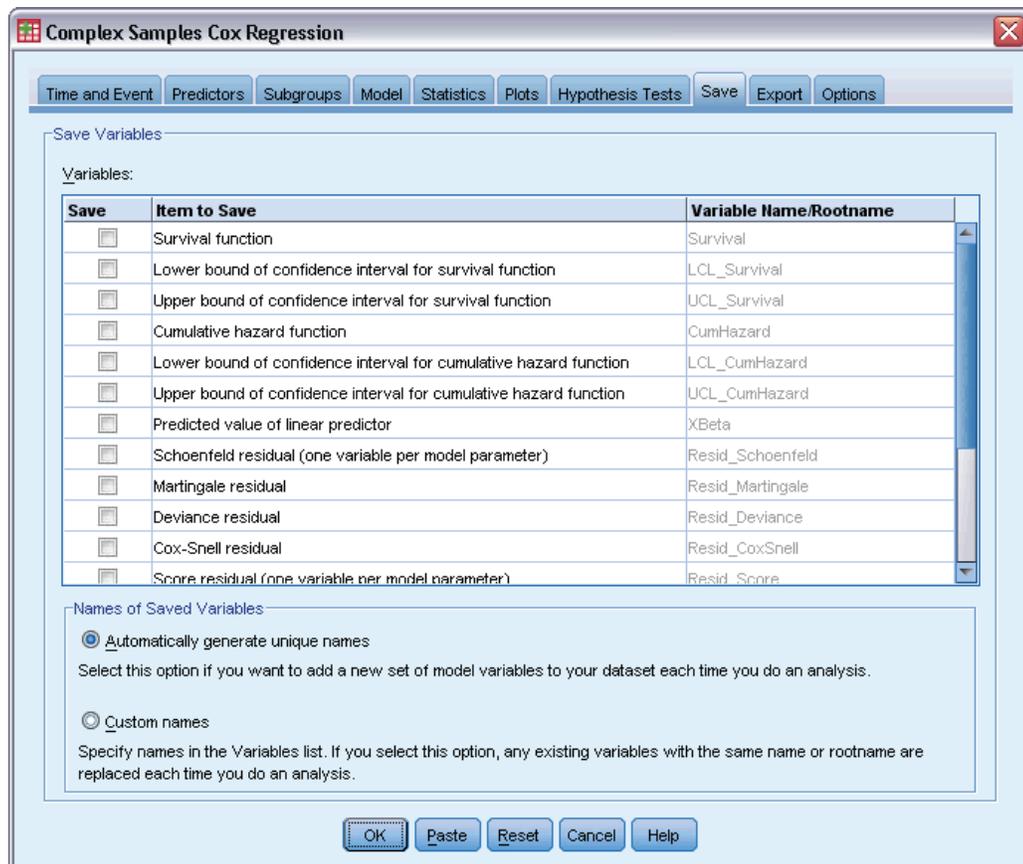
**Adjustment for Multiple Comparisons.** When performing hypothesis tests with multiple contrasts, the overall significance level can be adjusted from the significance levels for the included contrasts. This group allows you to choose the adjustment method.

- **Least significant difference.** This method does not control the overall probability of rejecting the hypotheses that some linear contrasts are different from the null hypothesis values.
- **Sequential Sidak.** This is a sequentially step-down rejective Sidak procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.

- **Sequential Bonferroni.** This is a sequentially step-down rejective Bonferroni procedure that is much less conservative in terms of rejecting individual hypotheses but maintains the same overall significance level.
- **Sidak.** This method provides tighter bounds than the Bonferroni approach.
- **Bonferroni.** This method adjusts the observed significance level for the fact that multiple contrasts are being tested.

## Save

Figure 12-10  
Cox Regression dialog box, Save tab



**Save Variables.** This group allows you to save model-related variables to the active dataset for further use in diagnostics and reporting of results. Note that none of these are available when time-dependent predictors are included in the model.

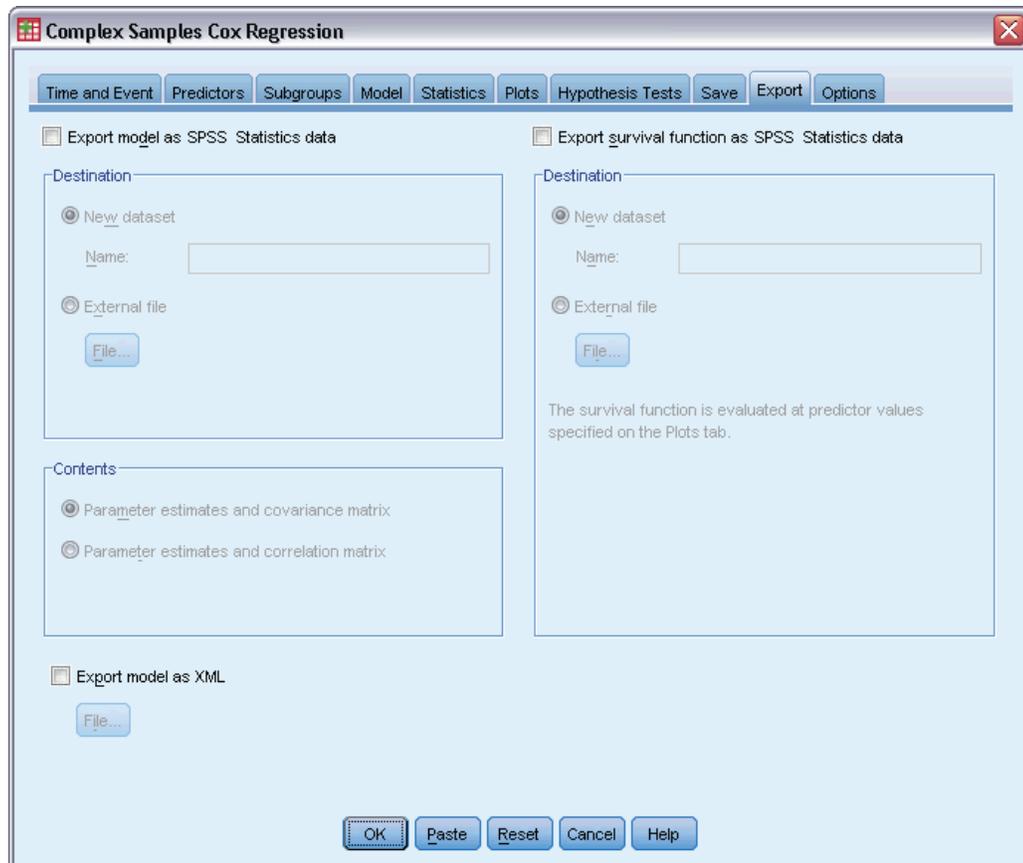
- **Survival function.** Saves the probability of survival (the value of the survival function) at the observed time and predictor values for each case.
- **Lower bound of confidence interval for survival function.** Saves the lower bound of the confidence interval for the survival function at the observed time and predictor values for each case.

- **Upper bound of confidence interval for survival function.** Saves the upper bound of the confidence interval for the survival function at the observed time and predictor values for each case.
- **Cumulative hazard function.** Saves the cumulative hazard, or  $-\ln(\text{survival})$ , at the observed time and predictor values for each case.
- **Lower bound of confidence interval for cumulative hazard function.** Saves the lower bound of the confidence interval for the cumulative hazard function at the observed time and predictor values for each case.
- **Upper bound of confidence interval for cumulative hazard function.** Saves the upper bound of the confidence interval for the cumulative hazard function at the observed time and predictor values for each case.
- **Predicted value of linear predictor.** Saves the linear combination of reference value corrected predictors times regression coefficients. The linear predictor is the ratio of the hazard function to the baseline hazard. Under the proportional hazards model, this value is constant across time.
- **Schoenfeld residual.** For each uncensored case and each nonredundant parameter in the model, the Schoenfeld residual is the difference between the observed value of the predictor associated with the model parameter and the expected value of the predictor for cases in the risk set at the observed event time. Schoenfeld residuals can be used to help assess the proportional hazards assumption; for example, for a predictor  $x$ , plots of the Schoenfeld residuals for the time-dependent predictor  $x \cdot \ln(T_)$  versus time should show a horizontal line at 0 if proportional hazards holds. A separate variable is saved for each nonredundant parameter in the model. Schoenfeld residuals are only computed for uncensored cases.
- **Martingale residual.** For each case, the martingale residual is the difference between the observed censoring (0 if censored, 1 if not) and the expectation of an event during the observation time.
- **Deviance residual.** Deviance residuals are martingale residuals “adjusted” to appear more symmetrical about 0. Plots of deviance residuals against predictors should reveal no patterns.
- **Cox-Snell residual.** For each case, the Cox-Snell residual is the expectation of an event during the observation time, or the observed censoring minus the martingale residual.
- **Score residual.** For each case and each nonredundant parameter in the model, the score residual is the contribution of the case to the first derivative of the pseudo-likelihood. A separate variable is saved for each nonredundant parameter in the model.
- **DFBeta residual.** For each case and each nonredundant parameter in the model, the DFBeta residual approximates the change in the value of the parameter estimate when the case is removed from the model. Cases with relatively large DFBeta residuals may be exerting undue influence on the analysis. A separate variable is saved for each nonredundant parameter in the model.
- **Aggregated residuals.** When multiple cases represent a single subject, the aggregated residual for a subject is simply the sum of the corresponding case residuals over all cases belonging to the same subject. For Schoenfeld’s residual, the aggregated version is the same as that of the non-aggregated version because Schoenfeld’s residual is only defined for uncensored cases. These residuals are only available when a subject identifier is specified on the Time and Event tab.

**Names of Saved Variables.** Automatic name generation ensures that you keep all your work. Custom names allow you to discard/replace results from previous runs without first deleting the saved variables in the Data Editor.

## Export

Figure 12-11  
Cox Regression dialog box, Export tab



**Export model as SPSS Statistics data.** Writes a dataset in IBM® SPSS® Statistics format containing the parameter correlation or covariance matrix with parameter estimates, standard errors, significance values, and degrees of freedom. The order of variables in the matrix file is as follows.

- **rowtype\_.** Takes values (and value labels), COV (Covariances), CORR (Correlations), EST (Parameter estimates), SE (Standard errors), SIG (Significance levels), and DF (Sampling design degrees of freedom). There is a separate case with row type COV (or CORR) for each model parameter, plus a separate case for each of the other row types.

- **varname\_.** Takes values P1, P2, ..., corresponding to an ordered list of all model parameters, for row types COV or CORR, with value labels corresponding to the parameter strings shown in the parameter estimates table. The cells are blank for other row types.
- **P1, P2, ...** These variables correspond to an ordered list of all model parameters, with variable labels corresponding to the parameter strings shown in the parameter estimates table, and take values according to the row type. For redundant parameters, all covariances are set to zero; correlations are set to the system-missing value; all parameter estimates are set at zero; and all standard errors, significance levels, and residual degrees of freedom are set to the system-missing value.

*Note:* This file is not immediately usable for further analyses in other procedures that read a matrix file unless those procedures accept all the row types exported here.

**Export survival function as SPSS Statistics data.** Writes a dataset in SPSS Statistics format containing the survival function; standard error of the survival function; upper and lower bounds of the confidence interval of the survival function; and the cumulative hazards function for each failure or event time, evaluated at the baseline and at the predictor patterns specified on the Plot tab. The order of variables in the matrix file is as follows.

- **Baseline strata variable.** Separate survival tables are produced for each value of the strata variable.
- **Survival time variable.** The event time; a separate case is created for each unique event time.
- **Sur\_0, LCL\_Sur\_0, UCL\_Sur\_0.** Baseline survival function and the upper and lower bounds of its confidence interval.
- **Sur\_R, LCL\_Sur\_R, UCL\_Sur\_R.** Survival function evaluated at the “reference” pattern (see the pattern values table in the output) and the upper and lower bounds of its confidence interval.
- **Sur\_##, LCL\_Sur\_##, UCL\_Sur\_##, ...** Survival function evaluated at each of the predictor patterns specified on the Plots tab and the upper and lower bounds of their confidence intervals. See the pattern values table in the output to match patterns with the number ##.
- **Haz\_0, LCL\_Haz\_0, UCL\_Haz\_0.** Baseline cumulative hazard function and the upper and lower bounds of its confidence interval.
- **Haz\_R, LCL\_Haz\_R, UCL\_Haz\_R.** Cumulative hazard function evaluated at the “reference” pattern (see the pattern values table in the output) and the upper and lower bounds of its confidence interval.
- **Haz\_##, LCL\_Haz\_##, UCL\_Haz\_##, ...** Cumulative hazard function evaluated at each of the predictor patterns specified on the Plots tab and the upper and lower bounds of their confidence intervals. See the pattern values table in the output to match patterns with the number ##.

**Export model as XML.** Saves all information needed to predict the survival function, including parameter estimates and the baseline survival function, in XML (PMML) format. You can use this model file to apply the model information to other data files for scoring purposes.

## Options

Figure 12-12  
Cox Regression dialog box, Options tab

**Estimation.** These controls specify criteria for estimation of regression coefficients.

- **Maximum Iterations.** The maximum number of iterations the algorithm will execute. Specify a non-negative integer.
- **Maximum Step-Halving.** At each iteration, the step size is reduced by a factor of 0.5 until the log-likelihood increases or maximum step-halving is reached. Specify a positive integer.
- **Limit iterations based on change in parameter estimates.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the parameter estimates is less than the value specified, which must be positive.
- **Limit iterations based on change in log-likelihood.** When selected, the algorithm stops after an iteration in which the absolute or relative change in the log-likelihood function is less than the value specified, which must be positive.
- **Display iteration history.** Displays the iteration history for the parameter estimates and pseudo log-likelihood and prints the last evaluation of the change in parameter estimates and pseudo log-likelihood. The iteration history table prints every  $n$  iterations beginning with the  $0$ th

iteration (the initial estimates), where  $n$  is the value of the increment. If the iteration history is requested, then the last iteration is always displayed regardless of  $n$ .

- **Tie breaking method for parameter estimation.** When there are tied observed failure times, one of these methods is used to break the ties. The Efron method is more computationally expensive.

**Survival Functions.** These controls specify criteria for computations involving the survival function.

- **Method for estimating baseline survival functions.** The **Breslow** (or Nelson-Aalan or empirical) method estimates the baseline cumulative hazard by a nondecreasing step function with steps at the observed failure times, then computes the baseline survival by the relation  $\text{survival} = \exp(-\text{cumulative hazard})$ . The **Efron** method is more computationally expensive and reduces to the Breslow method when there are no ties. The **product limit** method estimates the baseline survival by a non-increasing right continuous function; when there are no predictors in the model, this method reduces to Kaplan-Meier estimation.
- **Confidence intervals of survival functions.** The confidence interval can be calculated in three ways: in original units, via a log transformation, or a log-minus-log transformation. Only the log-minus-log transformation guarantees that the bounds of the confidence interval will lie between 0 and 1, but the log transformation generally seems to perform “best.”

**User Missing Values.** All variables must have valid values for a case to be included in the analysis. These controls allow you to decide whether user-missing values are treated as valid among categorical models (including factors, event, strata, and subpopulation variables) and sampling design variables.

**Confidence interval(%).** This is the confidence interval level used for coefficient estimates, exponentiated coefficient estimates, survival function estimates, and cumulative hazard function estimates. Specify a value greater than or equal to 0, and less than 100.

## ***CSCOXREG Command Additional Features***

The command language also allows you to:

- Perform custom hypothesis tests (using the `CUSTOM` subcommand and `/PRINT LMATRIX`).
- Tolerance specification (using `/CRITERIA SINGULAR`).
- General estimable function table (using `/PRINT GEF`).
- Multiple predictor patterns (using multiple `PATTERN` subcommands).
- Maximum number of saved variables when a rootname is specified (using the `SAVE` subcommand). The dialog honors the `CSCOXREG` default of 25 variables.

See the *Command Syntax Reference* for complete syntax information.

# ***Part II: Examples***

# ***Complex Samples Sampling Wizard***

The Sampling Wizard guides you through the steps for creating, modifying, or executing a sampling plan file. Before using the wizard, you should have a well-defined target population, a list of sampling units, and an appropriate sample design in mind.

## ***Obtaining a Sample from a Full Sampling Frame***

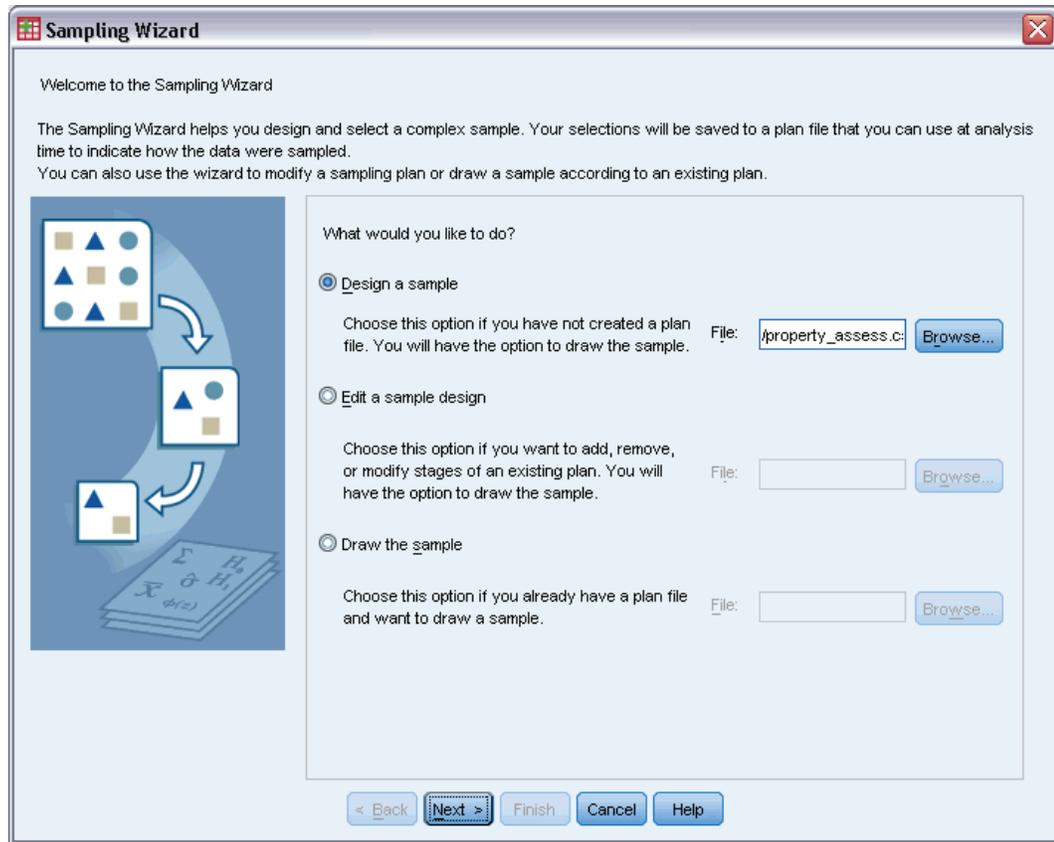
A state agency is charged with ensuring fair property taxes from county to county. Taxes are based on the appraised value of the property, so the agency wants to survey a sample of properties by county to be sure that each county's records are equally up to date. However, resources for obtaining current appraisals are limited, so it's important that what is available is used wisely. The agency decides to employ complex sampling methodology to select a sample of properties.

A listing of properties is collected in *property\_assess\_cs.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use the Complex Samples Sampling Wizard to select a sample.

## ***Using the Wizard***

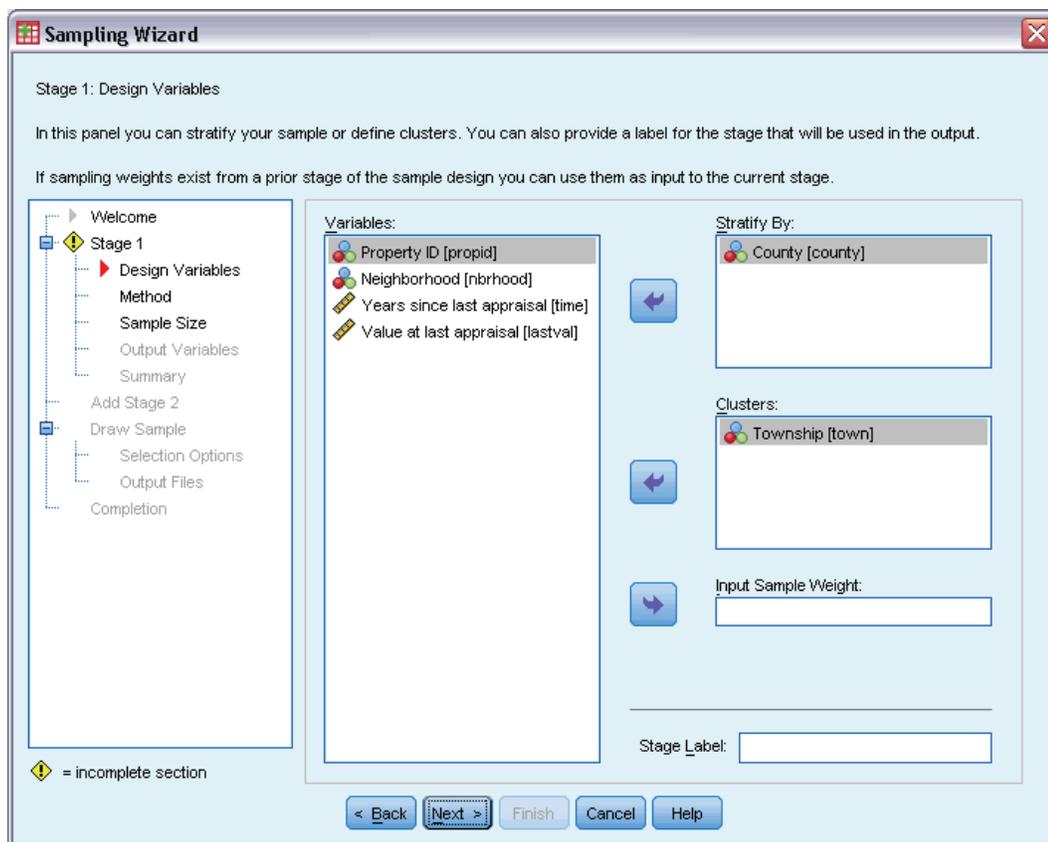
- ▶ To run the Complex Samples Sampling Wizard, from the menus choose:  
Analyze > Complex Samples > Select a Sample...

Figure 13-1  
Sampling Wizard, Welcome step



- ▶ Select Design a sample, browse to where you want to save the file, and type property\_assess.csplan as the name of the plan file.
- ▶ Click Next.

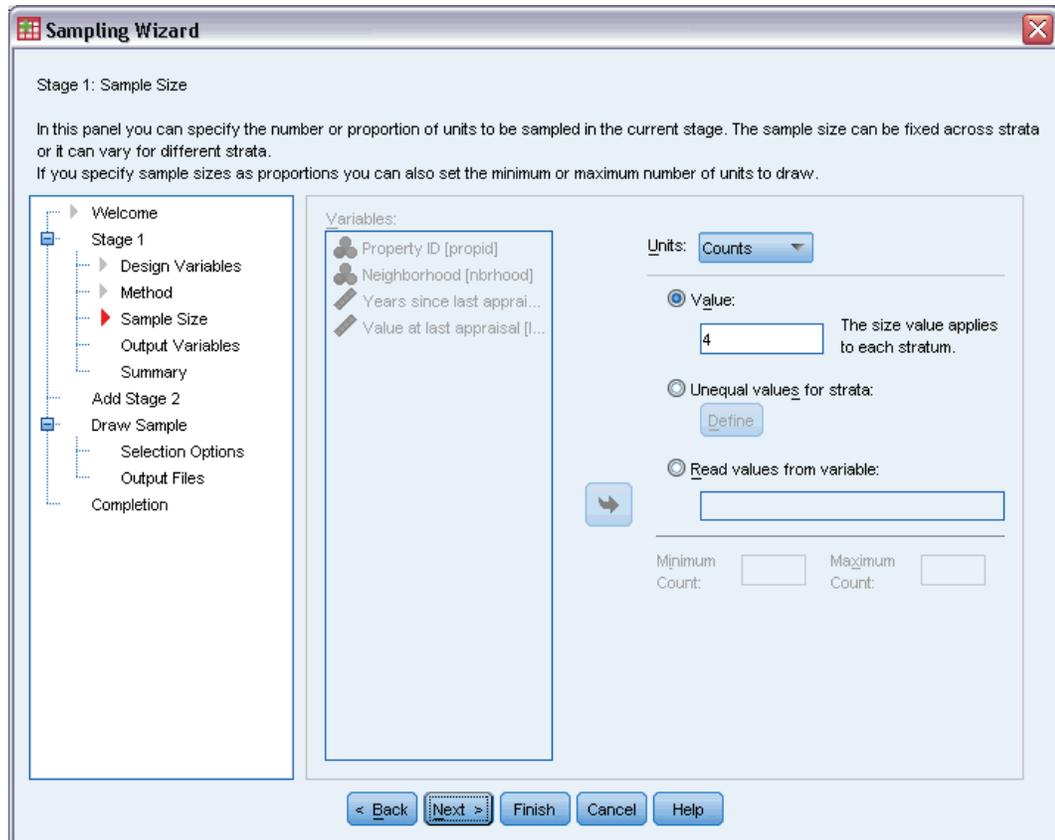
Figure 13-2  
Sampling Wizard, Design Variables step (stage 1)



- ▶ Select *County* as a stratification variable.
- ▶ Select *Township* as a cluster variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

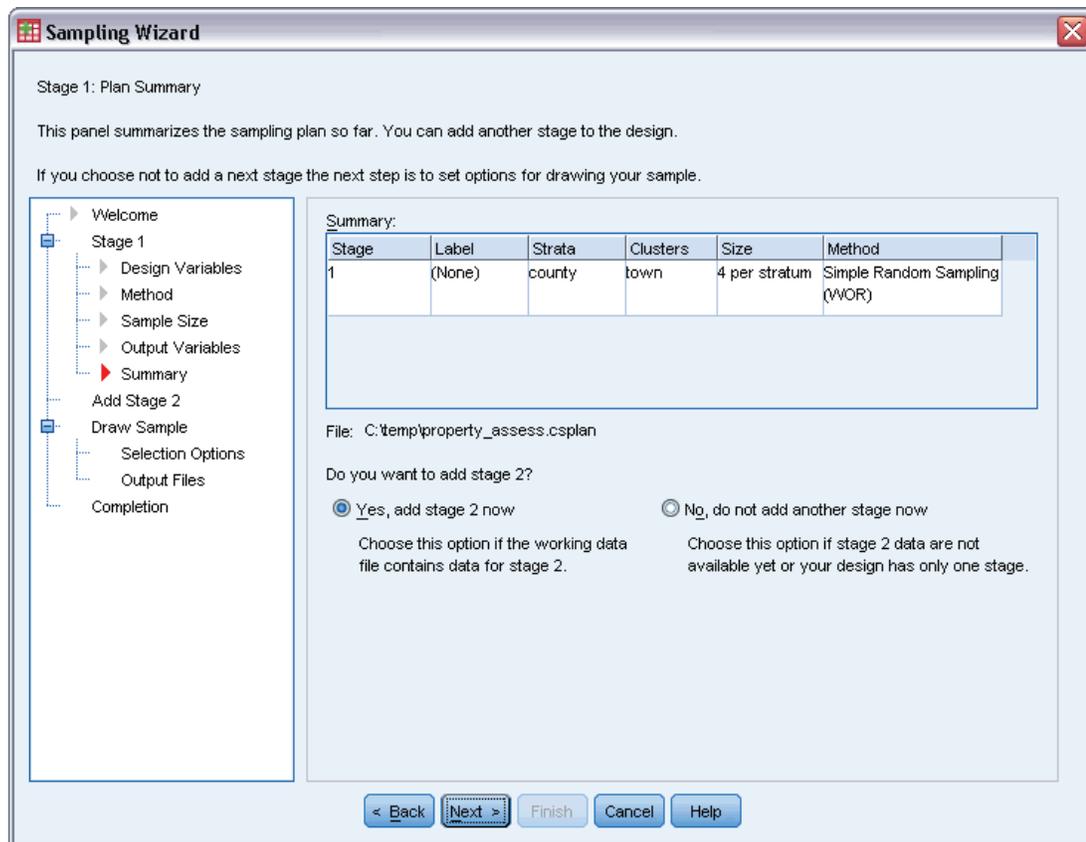
This design structure means that independent samples are drawn for each county. In this stage, townships are drawn as the primary sampling unit using the default method, simple random sampling.

Figure 13-3  
Sampling Wizard, Sample Size step (stage 1)



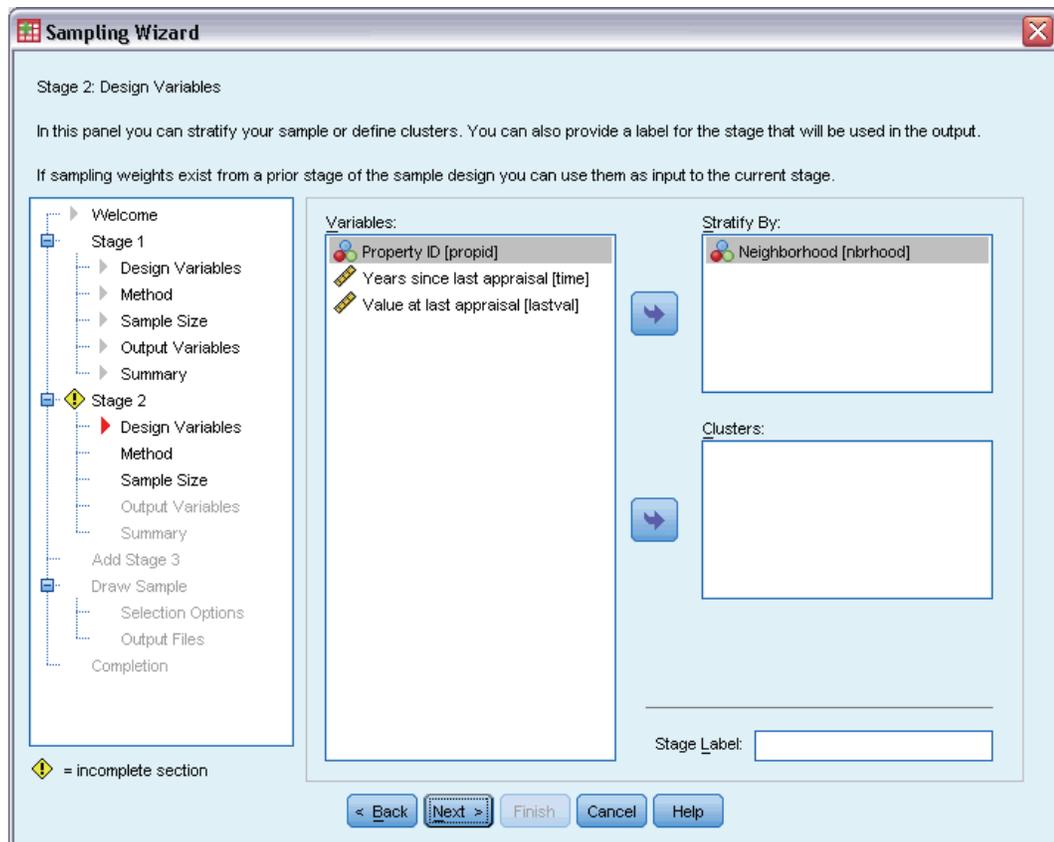
- ▶ Select Counts from the Units drop-down list.
- ▶ Type 4 as the value for the number of units to select in this stage.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-4  
Sampling Wizard, Plan Summary step (stage 1)



- ▶ Select Yes, add stage 2 now.
- ▶ Click Next.

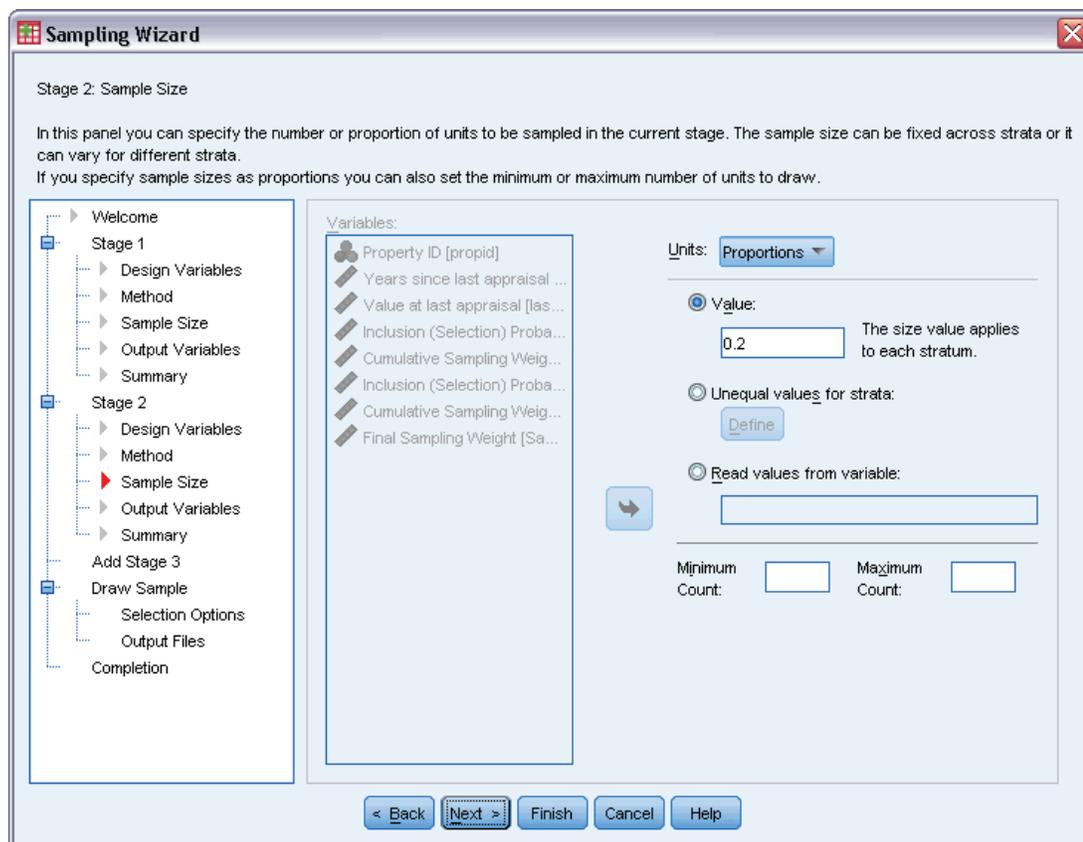
Figure 13-5  
Sampling Wizard, Design Variables step (stage 2)



- ▶ Select *Neighborhood* as a stratification variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

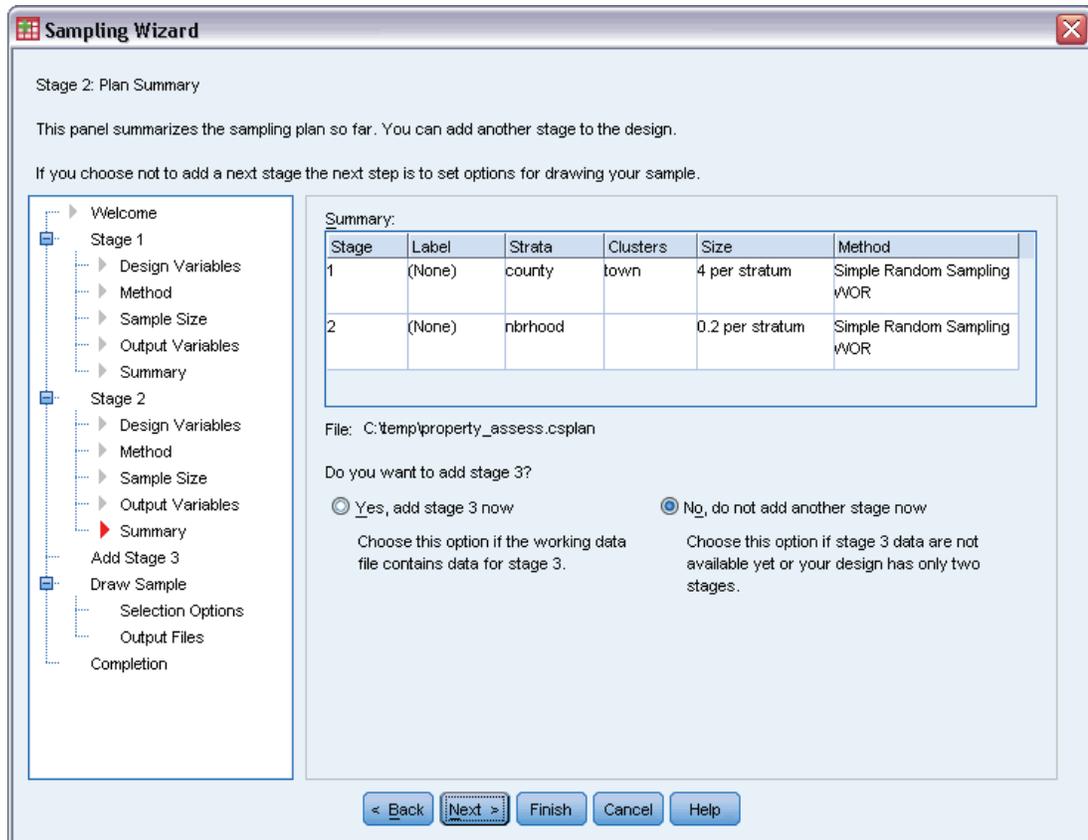
This design structure means that independent samples are drawn for each neighborhood of the townships drawn in stage 1. In this stage, properties are drawn as the primary sampling unit using simple random sampling.

Figure 13-6  
Sampling Wizard, Sample Size step (stage 2)



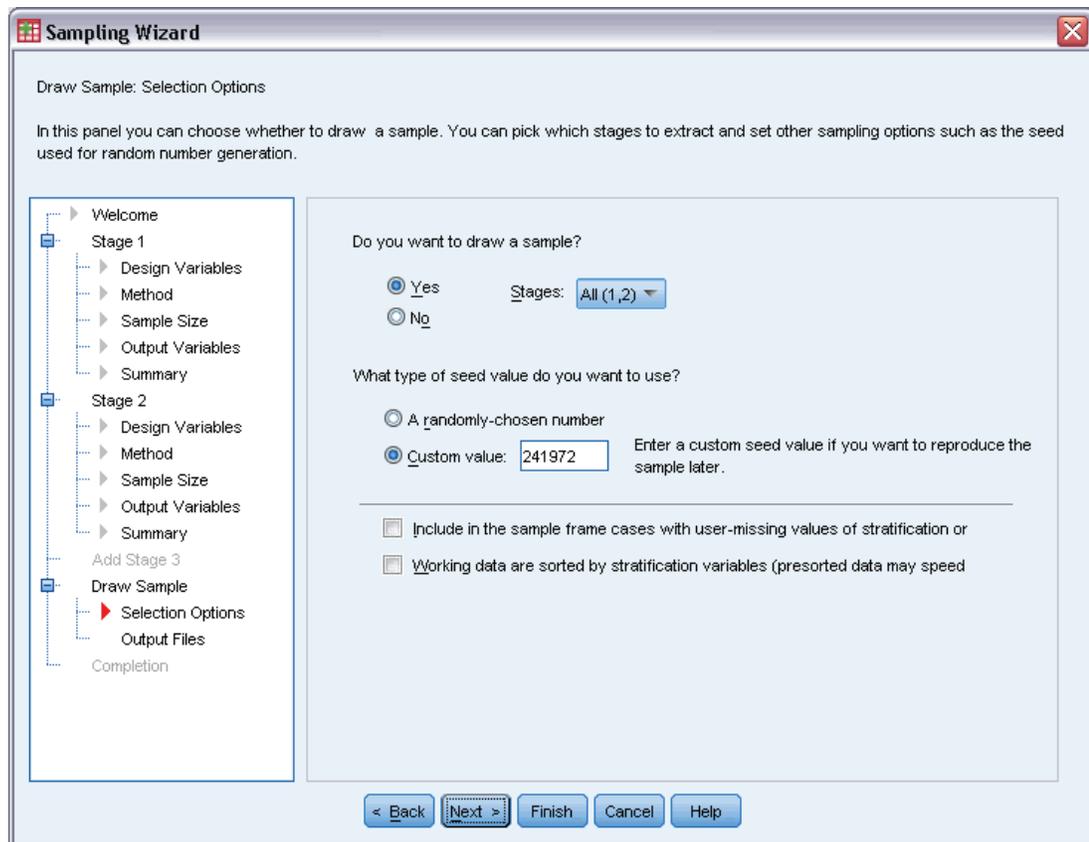
- ▶ Select Proportions from the Units drop-down list.
- ▶ Type 0.2 as the value of the proportion of units to sample from each stratum.
- ▶ Click Next, and then click Next in the Output Variables step.

**Figure 13-7**  
*Sampling Wizard, Plan Summary step (stage 2)*



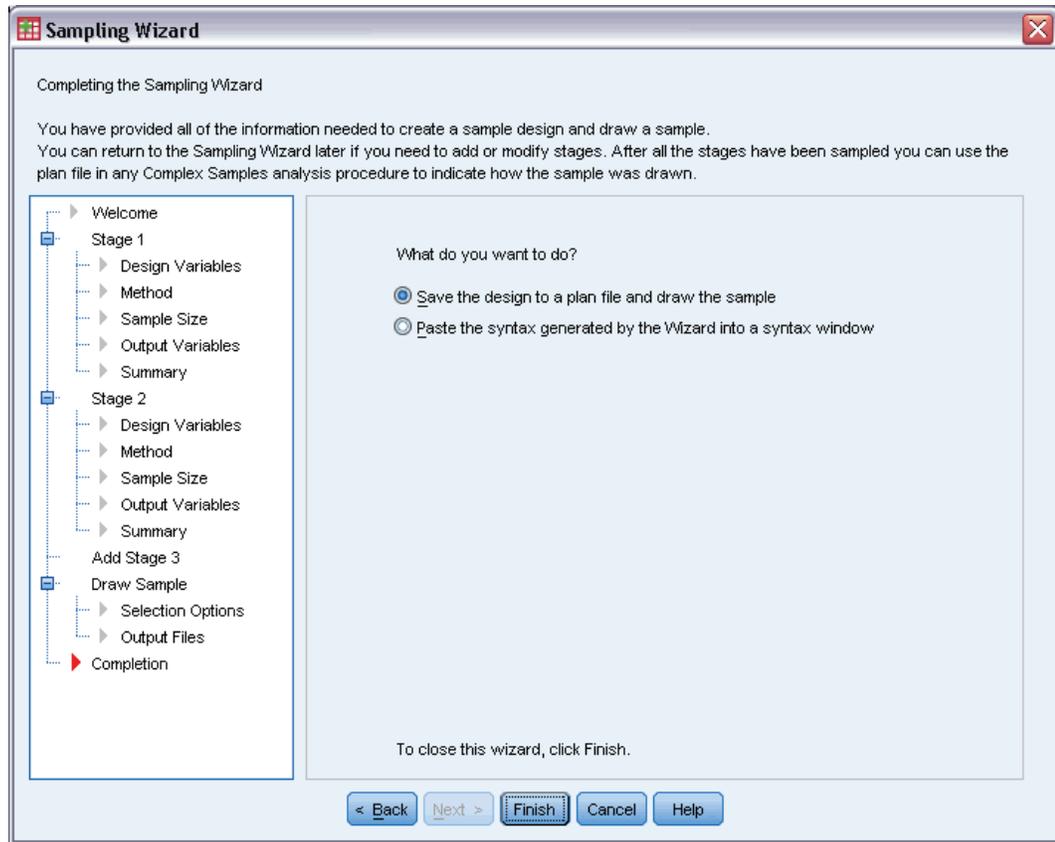
- Look over the sampling design, and then click Next.

Figure 13-8  
Sampling Wizard, Draw Sample, Selection Options step



- ▶ Select Custom value for the type of random seed to use, and type 241972 as the value. Using a custom value allows you to replicate the results of this example exactly.
- ▶ Click Next, and then click Next in the Draw Sample Output Files step.

Figure 13-9  
Sampling Wizard, Finish step



- ▶ Click Finish.

These selections produce the sampling plan file *property\_assess.csplan* and draw a sample according to that plan.

### Plan Summary

Figure 13-10  
Plan summary

			Stage 1	Stage 2
Design Variables	Stratification	1	County	Neighborhood
	Cluster	1	Township	
Sample Information	Selection Method		Simple random sampling without replacement	Simple random sampling without replacement
	Number of Units Sampled		4	
	Variables Created or Modified	Stagewise Inclusion (Selection) Probability	Inclusion Probability_1_	Inclusion Probability_2_
		Stagewise Cumulative Sample Weight	Sample Weight Cumulative_1_	Sample Weight Cumulative_2_
	Proportion of Units Sampled			.2
Analysis Information	Estimator Assumption		Equal probability sampling without replacement	Equal probability sampling without replacement
	Inclusion Probability		Obtained from variable Inclusion Probability_1_	Obtained from variable Inclusion Probability_2_

Plan File: c:\property\_assess.csplan  
Weight Variable: SampleWeight\_Final\_

The summary table reviews your sampling plan and is useful for making sure that the plan represents your intentions.

### Sampling Summary

Figure 13-11  
Stage summary

County	Number of Units Sampled		Proportion of Units Sampled	
	Requested	Actual	Requested	Actual
Eastern	4	4	44.4%	44.4%
Central	4	4	57.1%	57.1%
Western	4	4	25.0%	25.0%
Northern	4	4	44.4%	44.4%
Southern	4	4	50.0%	50.0%

Plan File: c:\property\_assess.csplan

This summary table reviews the first stage of sampling and is useful for checking that the sampling went according to plan. Four townships were sampled from each county, as requested.

Figure 13-12  
Stage summary

County	Township	Neighborhood	Number of Units Sampled		Proportion of Units Sampled	
			Requested	Actual	Requested	Actual
Eastern	2	8	4	4	20.0%	19.0%
		9	14	14	20.0%	20.6%
		10	7	7	20.0%	18.9%
		11	14	14	20.0%	20.0%
	6	36	13	13	20.0%	20.3%
		37	14	14	20.0%	20.6%
		38	13	13	20.0%	20.6%
	7	43	12	12	20.0%	20.7%
		44	11	11	20.0%	19.6%
		45	11	11	20.0%	20.8%
		46	13	13	20.0%	20.0%
	9	57	13	13	20.0%	20.6%
		58	5	5	20.0%	18.5%
		59	11	11	20.0%	19.3%
60		13	13	20.0%	19.4%	
Central	22	148	9	9	20.0%	19.6%
		149	8	8	20.0%	20.0%

This summary table (the top part of which is shown here) reviews the second stage of sampling. It is also useful for checking that the sampling went according to plan. Approximately 20% of the properties were sampled from each neighborhood from each township sampled in the first stage, as requested.

## Sample Results

Figure 13-13  
Data Editor with sample results

	propid	nbrhood	town	county	time	lastval	InclusionPr obability_1	SampleWei ghtCumulat ive_1	InclusionPr obability_2	SampleWei ghtCumulat ive_2	SampleWei ght_Final_
273	577.0	8	2	1	4	181.70	.	.	.	.	.
274	578.0	8	2	1	5	189.60	.	.	.	.	.
275	579.0	8	2	1	4	200.10	.	.	.	.	.
276	580.0	8	2	1	5	211.50	.	.	.	.	.
277	581.0	8	2	1	4	181.50	.	.	.	.	.
278	641.0	9	2	1	7	192.40	.	.	.	.	.
279	642.0	9	2	1	6	236.70	.44	2.25	.21	10.93	10.93
280	643.0	9	2	1	6	150.40	.44	2.25	.21	10.93	10.93
281	644.0	9	2	1	8	204.80	.	.	.	.	.
282	645.0	9	2	1	6	225.40	.	.	.	.	.
283	646.0	9	2	1	7	180.80	.44	2.25	.21	10.93	10.93
284	647.0	9	2	1	5	176.90	.	.	.	.	.

You can see the sampling results in the Data Editor. Five new variables were saved to the working file, representing the inclusion probabilities and cumulative sampling weights for each stage, plus the final sampling weights.

- Cases with values for these variables were selected to the sample.
- Cases with system-missing values for the variables were not selected.

The agency will now use its resources to collect current valuations for the properties selected in the sample. Once those valuations are available, you can process the sample with Complex Samples analysis procedures, using the sampling plan *property\_assess.csplan* to provide the sampling specifications.

## ***Obtaining a Sample from a Partial Sampling Frame***

A company is interested in compiling and selling a database of high-quality survey information. The survey sample should be representative but efficiently carried out, so complex sampling methods are used. The full sampling design calls for the following structure:

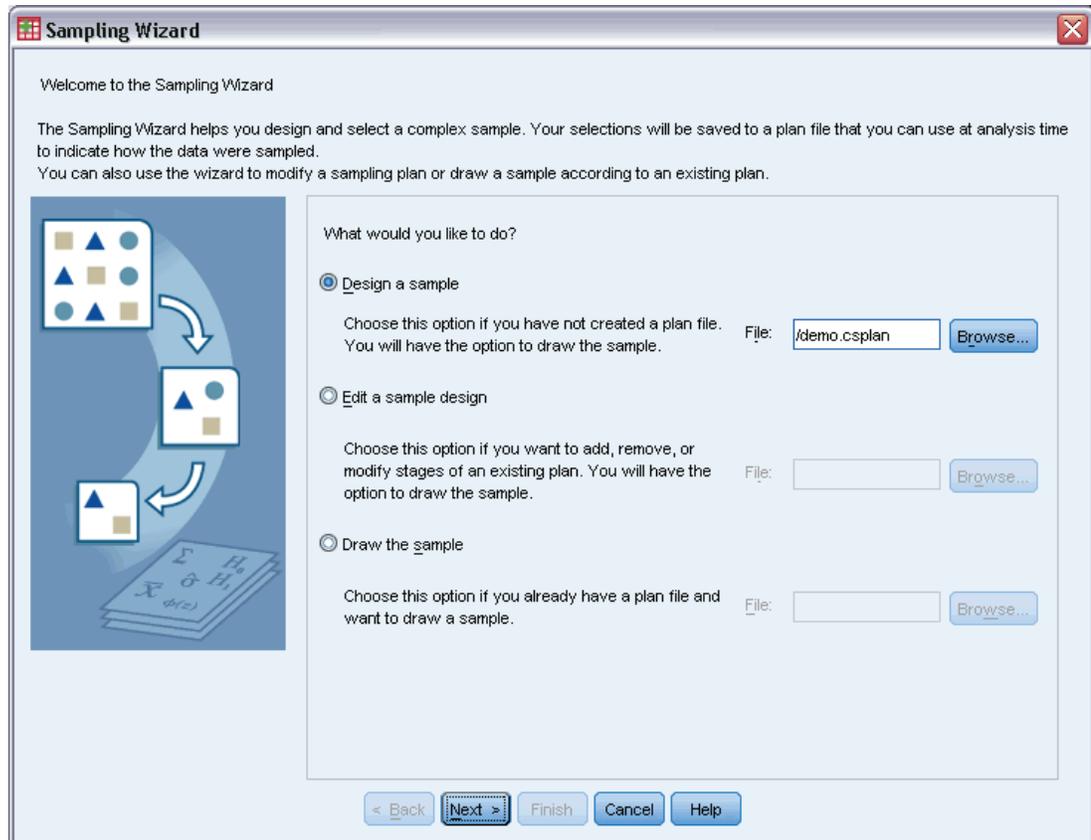
Stage	Strata	Clusters
1	Region	Province
2	District	City
3	Subdivision	

In the third stage, households are the primary sampling unit, and selected households will be surveyed. However, since information is easily available only to the city level, the company plans to execute the first two stages of the design now and then collect information on the numbers of subdivisions and households from the sampled cities. The available information to the city level is collected in *demo\_cs\_1.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Note that this file contains a variable *Subdivision* that contains all 1's. This is a placeholder for the “true” variable, whose values will be collected after the first two stages of the design are executed, that allows you to specify the full three-stage sampling design now. Use the Complex Samples Sampling Wizard to specify the full complex sampling design, and then draw the first two stages.

## ***Using the Wizard to Sample from the First Partial Frame***

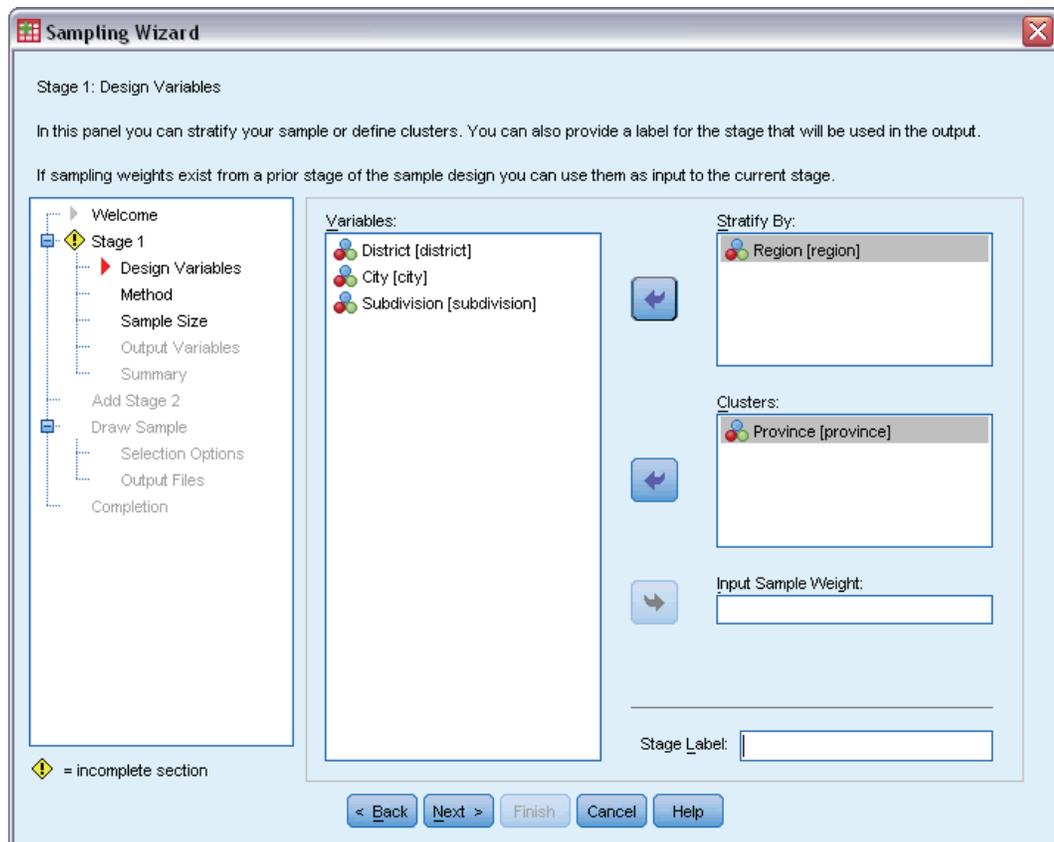
- ▶ To run the Complex Samples Sampling Wizard, from the menus choose:  
Analyze > Complex Samples > Select a Sample...

Figure 13-14  
Sampling Wizard, Welcome step



- ▶ Select Design a sample, browse to where you want to save the file, and type demo.csplan as the name of the plan file.
- ▶ Click Next.

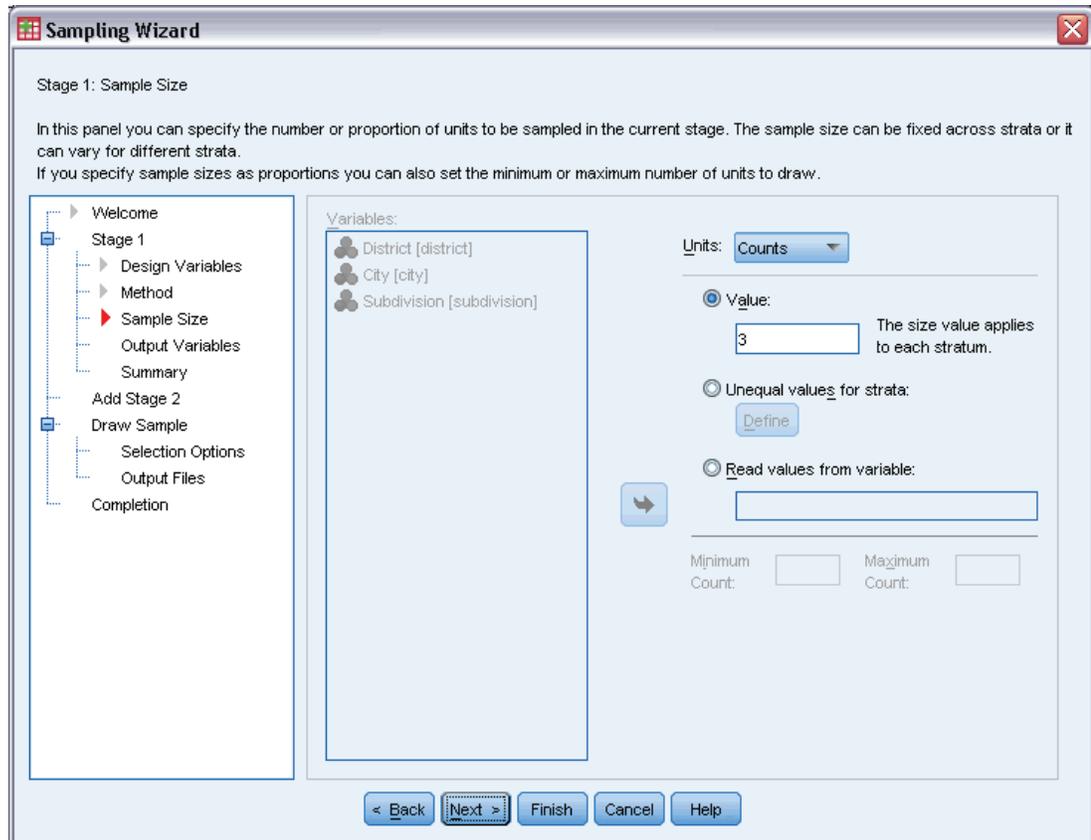
Figure 13-15  
Sampling Wizard, Design Variables step (stage 1)



- ▶ Select *Region* as a stratification variable.
- ▶ Select *Province* as a cluster variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

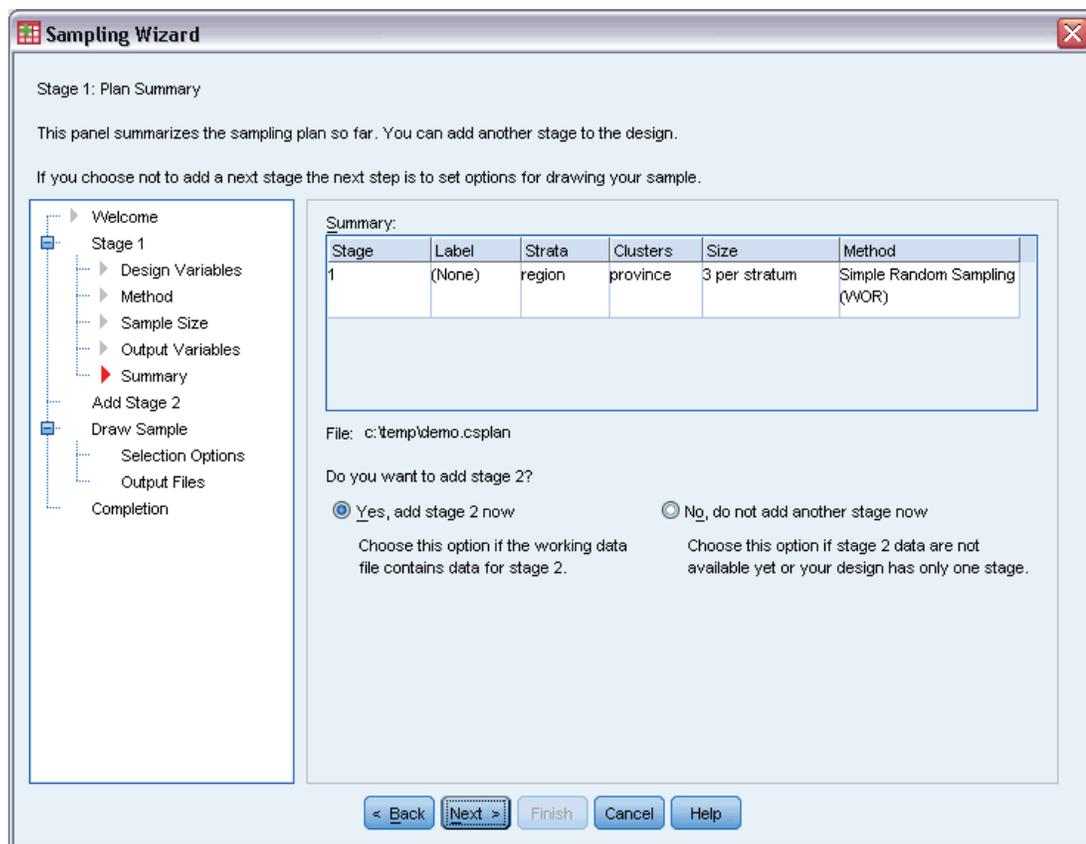
This design structure means that independent samples are drawn for each region. In this stage, provinces are drawn as the primary sampling unit using the default method, simple random sampling.

Figure 13-16  
Sampling Wizard, Sample Size step (stage 1)



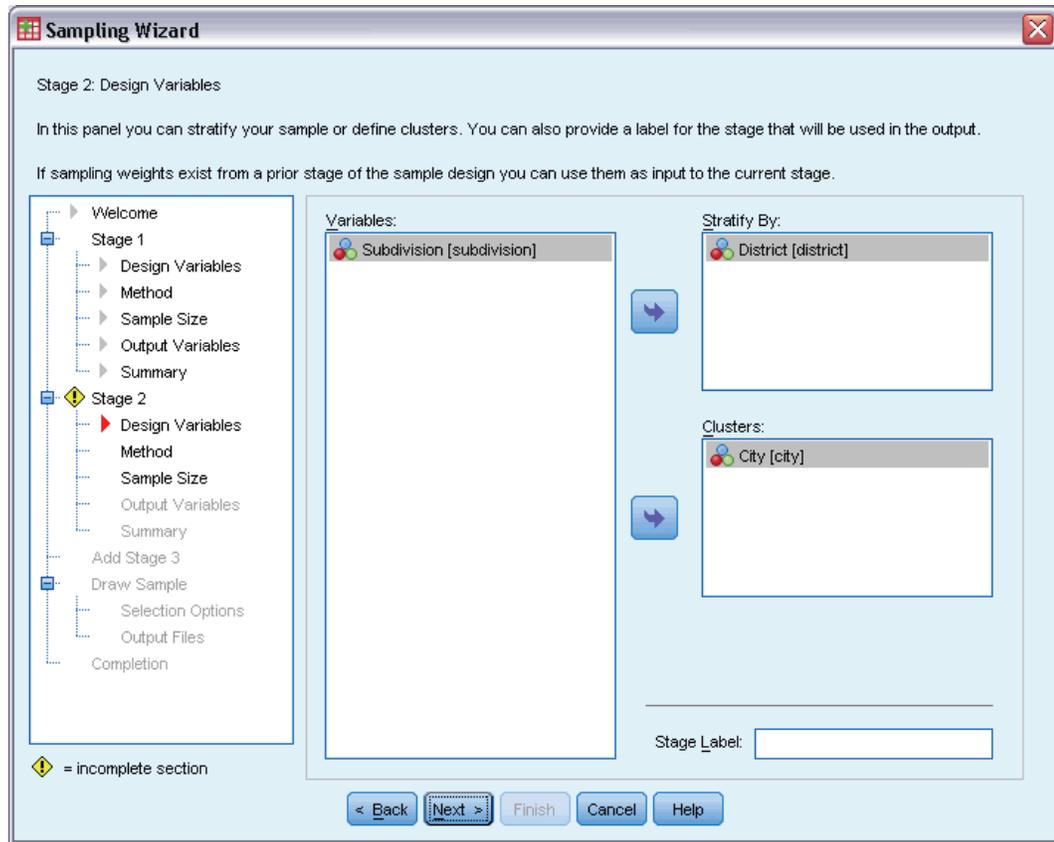
- ▶ Select Counts from the Units drop-down list.
- ▶ Type 3 as the value for the number of units to select in this stage.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-17  
Sampling Wizard, Plan Summary step (stage 1)



- ▶ Select Yes, add stage 2 now.
- ▶ Click Next.

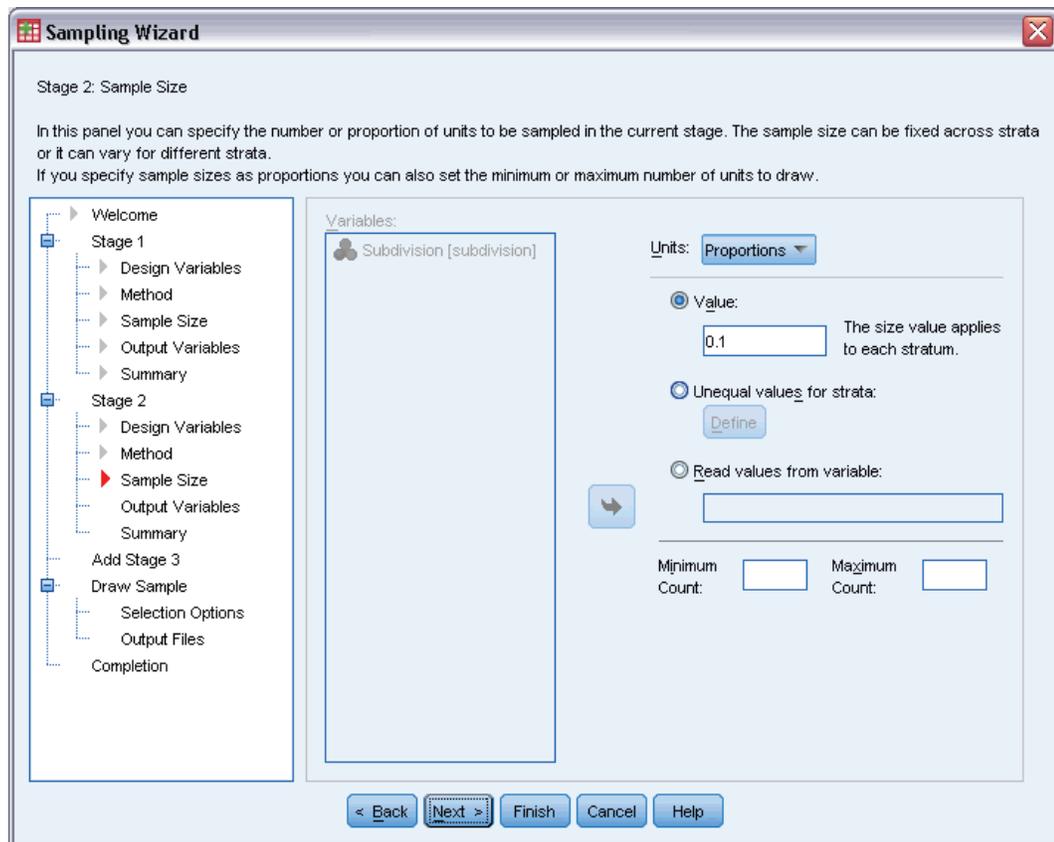
Figure 13-18  
Sampling Wizard, Design Variables step (stage 2)



- ▶ Select *District* as a stratification variable.
- ▶ Select *City* as a cluster variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

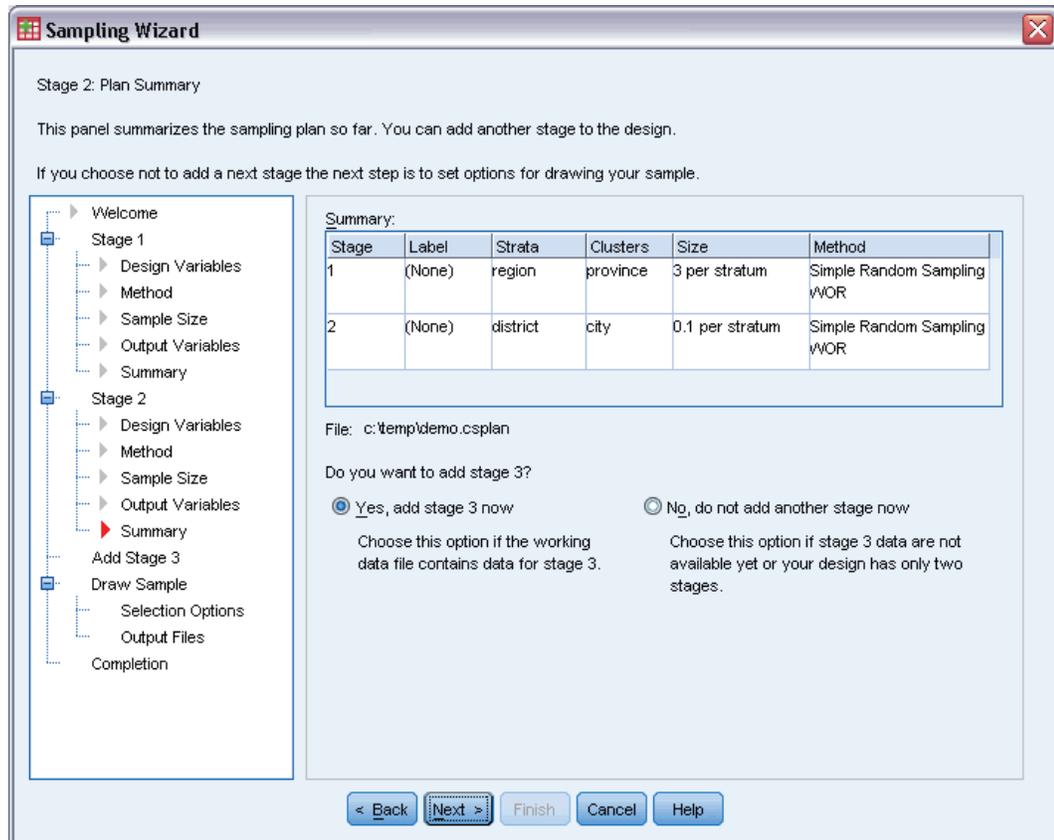
This design structure means that independent samples are drawn for each district. In this stage, cities are drawn as the primary sampling unit using the default method, simple random sampling.

Figure 13-19  
Sampling Wizard, Sample Size step (stage 2)



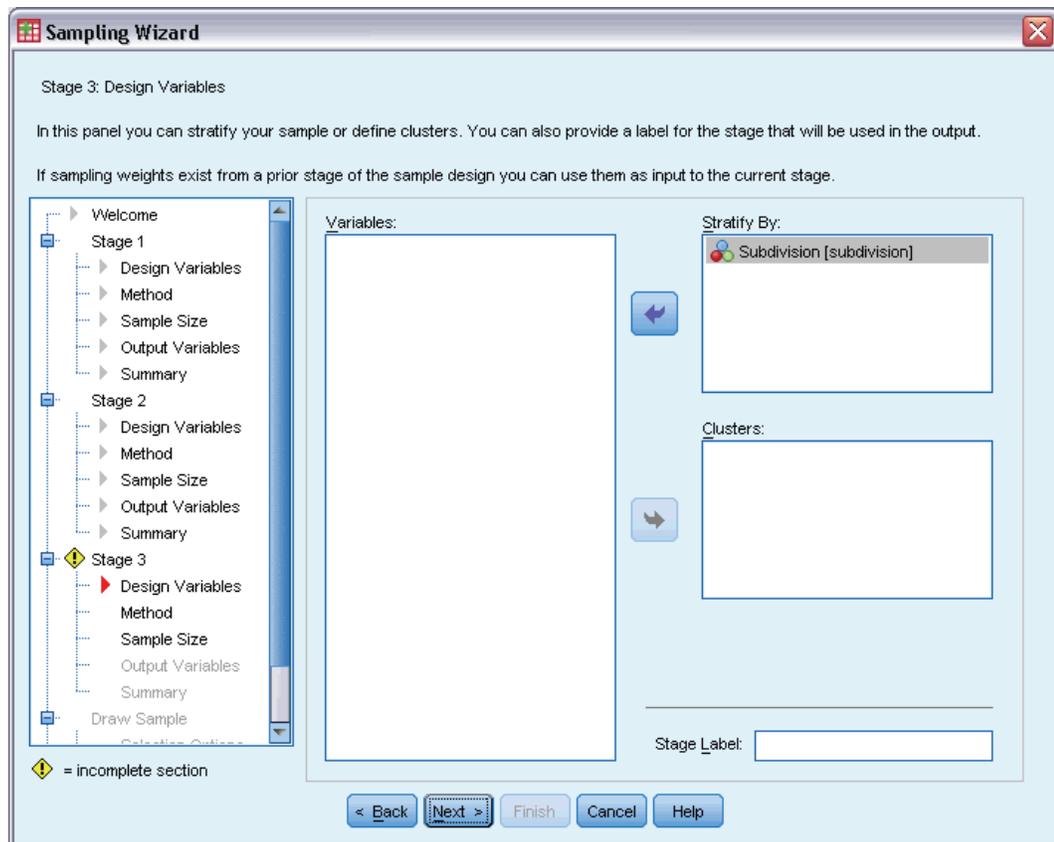
- ▶ Select Proportions from the Units drop-down list.
- ▶ Type 0.1 as the value of the proportion of units to sample from each strata.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-20  
Sampling Wizard, Plan Summary step (stage 2)



- ▶ Select Yes, add stage 3 now.
- ▶ Click Next.

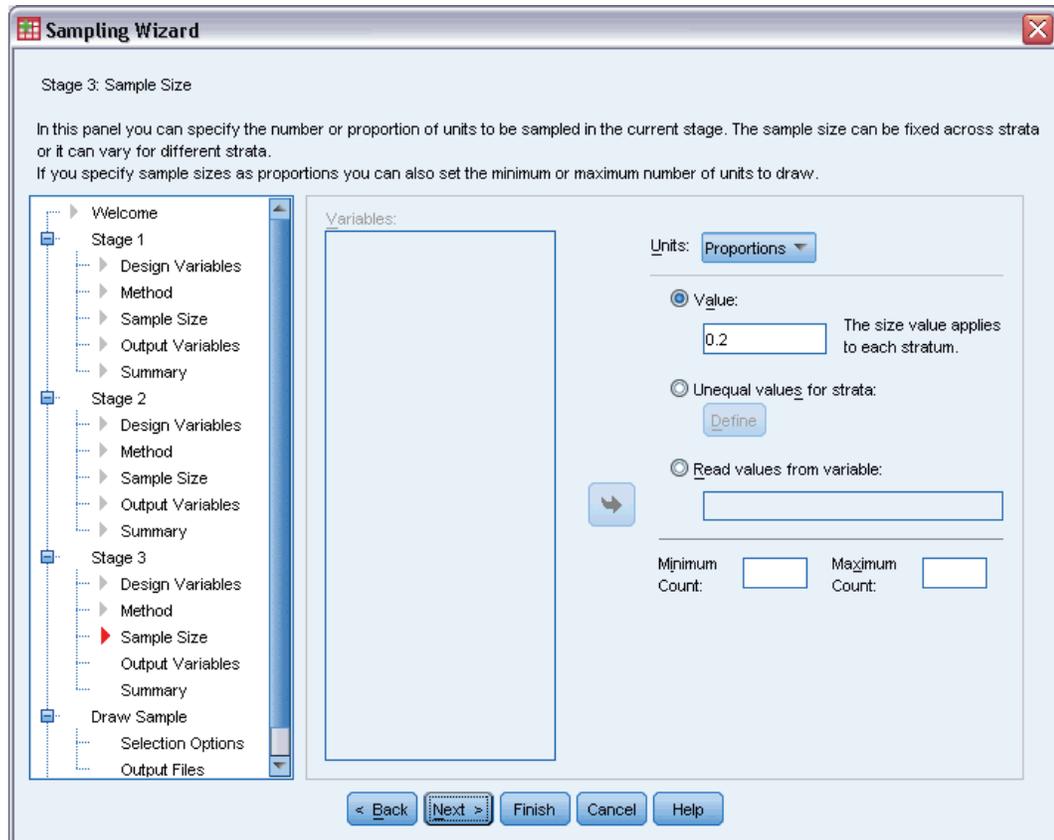
Figure 13-21  
Sampling Wizard, Design Variables step (stage 3)



- ▶ Select *Subdivision* as a stratification variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

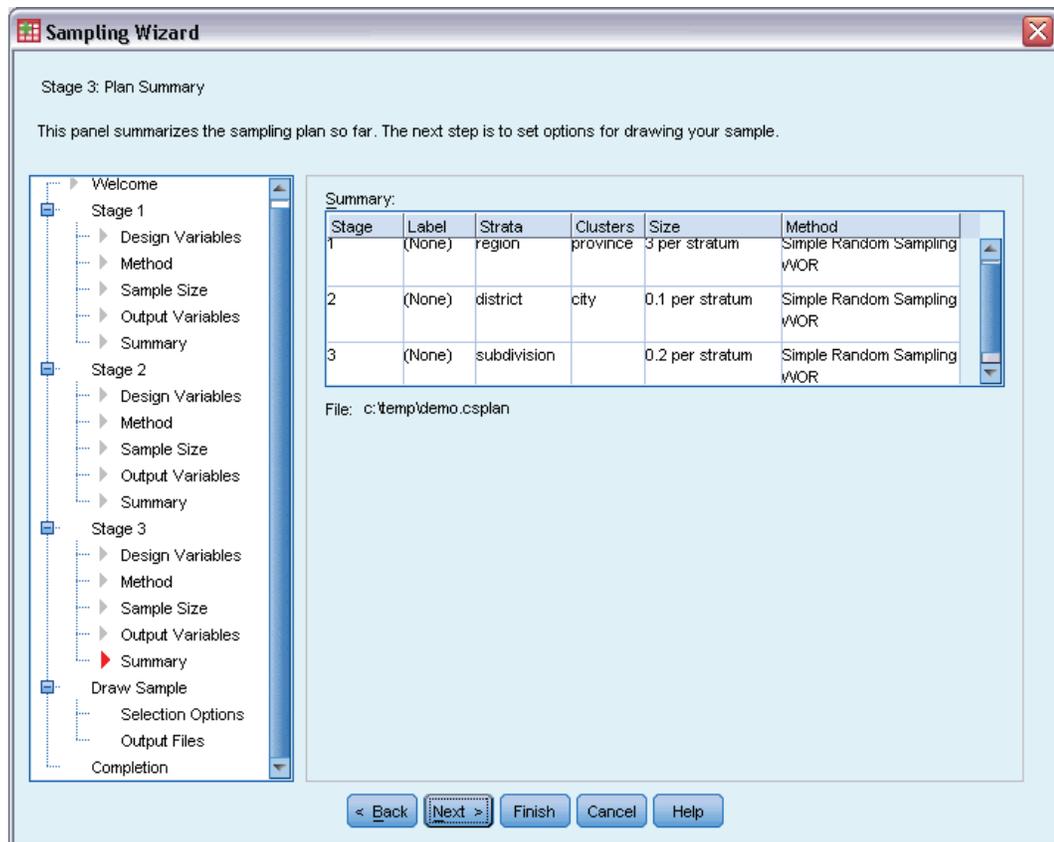
This design structure means that independent samples are drawn for each subdivision. In this stage, household units are drawn as the primary sampling unit using the default method, simple random sampling.

Figure 13-22  
Sampling Wizard, Sample Size step (stage 3)



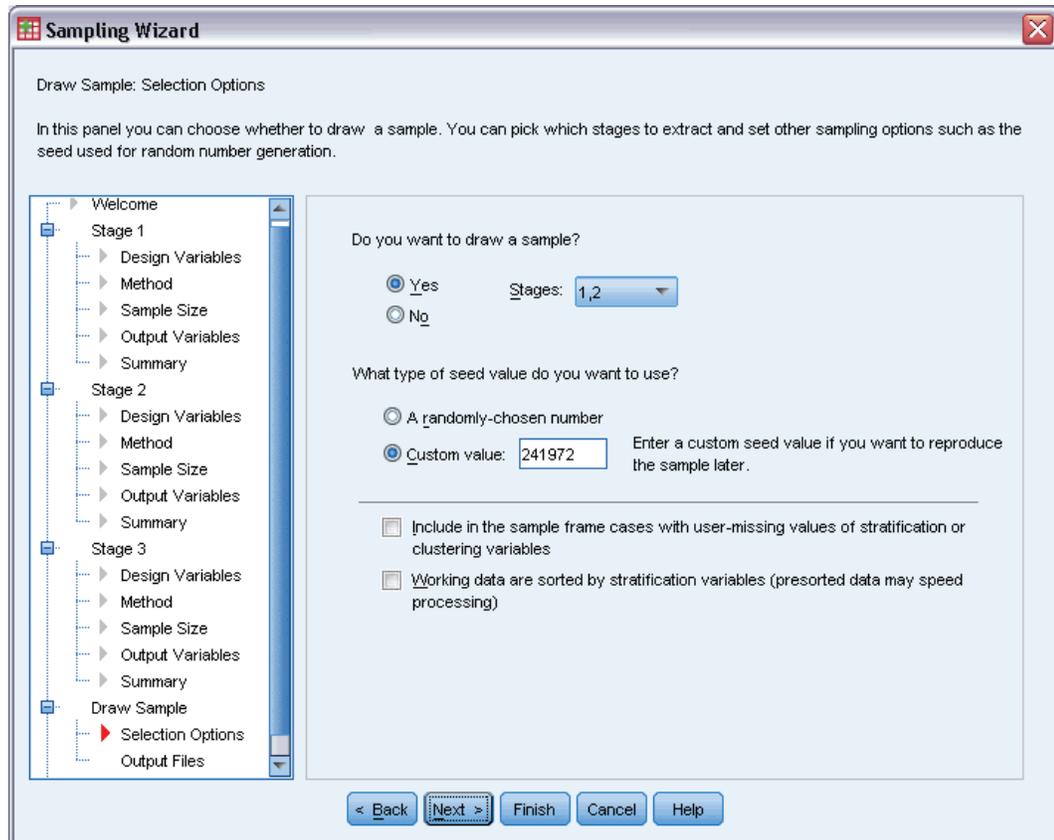
- ▶ Select Proportions from the Units drop-down list.
- ▶ Type 0.2 as the value for the proportion of units to select in this stage.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-23  
Sampling Wizard, Plan Summary step (stage 3)



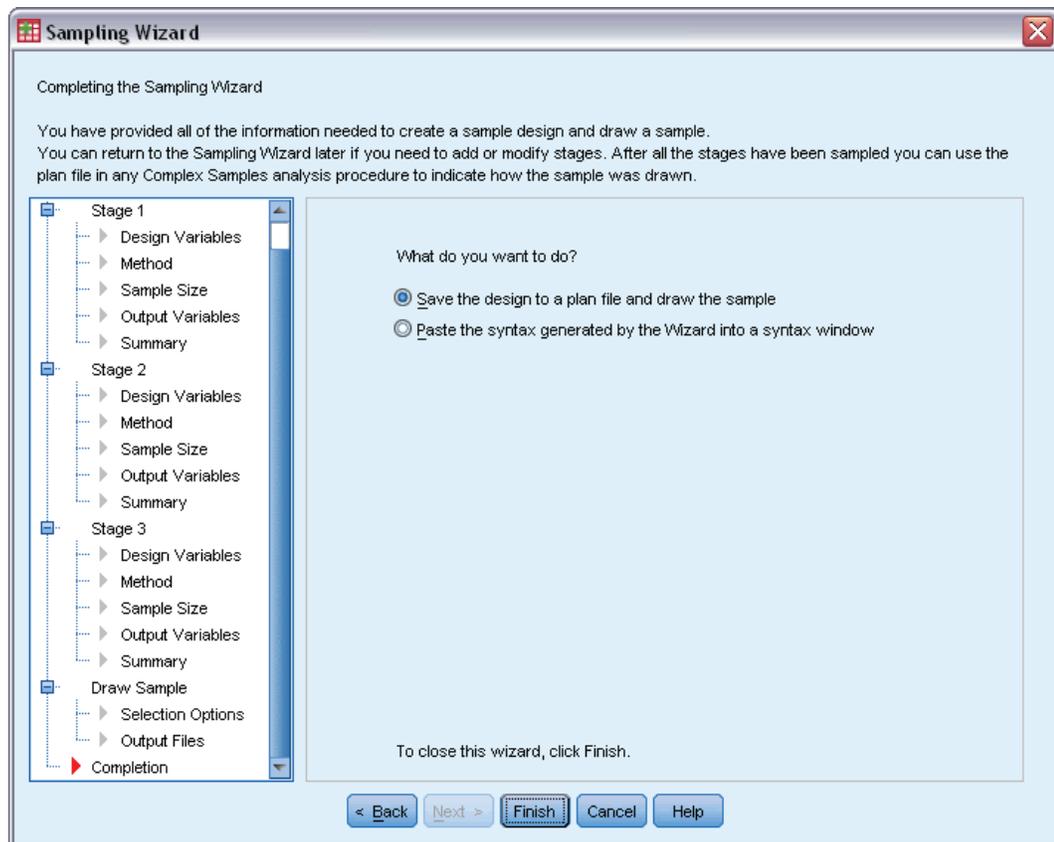
- Look over the sampling design, and then click Next.

Figure 13-24  
Sampling Wizard, Draw Sample Selection Options step



- ▶ Select 1, 2 as the stages to sample now.
- ▶ Select Custom value for the type of random seed to use, and type 241972 as the value.  
Using a custom value allows you to replicate the results of this example exactly.
- ▶ Click Next, and then click Next in the Draw Sample Output Files step.

Figure 13-25  
Sampling Wizard, Finish step



- ▶ Click Finish.

These selections produce the sampling plan file *demo.csplan* and draw a sample according to the first two stages of that plan.

## Sample Results

Figure 13-26  
Data Editor with sample results

	region	province	district	city	InclusionPr obability_1_	SampleWei ghtCumulat ve_1	InclusionPr obability_2_	SampleWei ghtCumulat ve_2	SampleWei ght_Final_
295	1	2	10	295	.	.	.	.	.
296	1	2	10	296	.	.	.	.	.
297	1	2	10	297	.	.	.	.	.
298	1	2	10	298	.20	5.00	.10	50.00	50.00
299	1	2	10	299	.	.	.	.	.
300	1	2	10	300	.20	5.00	.10	50.00	50.00
301	1	2	11	301	.	.	.	.	.
302	1	2	11	302	.	.	.	.	.
303	1	2	11	303	.	.	.	.	.
304	1	2	11	304	.	.	.	.	.
305	1	2	11	305	.	.	.	.	.
306	1	2	11	306	.	.	.	.	.
307	1	2	11	307	.20	5.00	.10	50.00	50.00
308	1	2	11	308	.	.	.	.	.

You can see the sampling results in the Data Editor. Five new variables were saved to the working file, representing the inclusion probabilities and cumulative sampling weights for each stage, plus the “final” sampling weights for the first two stages.

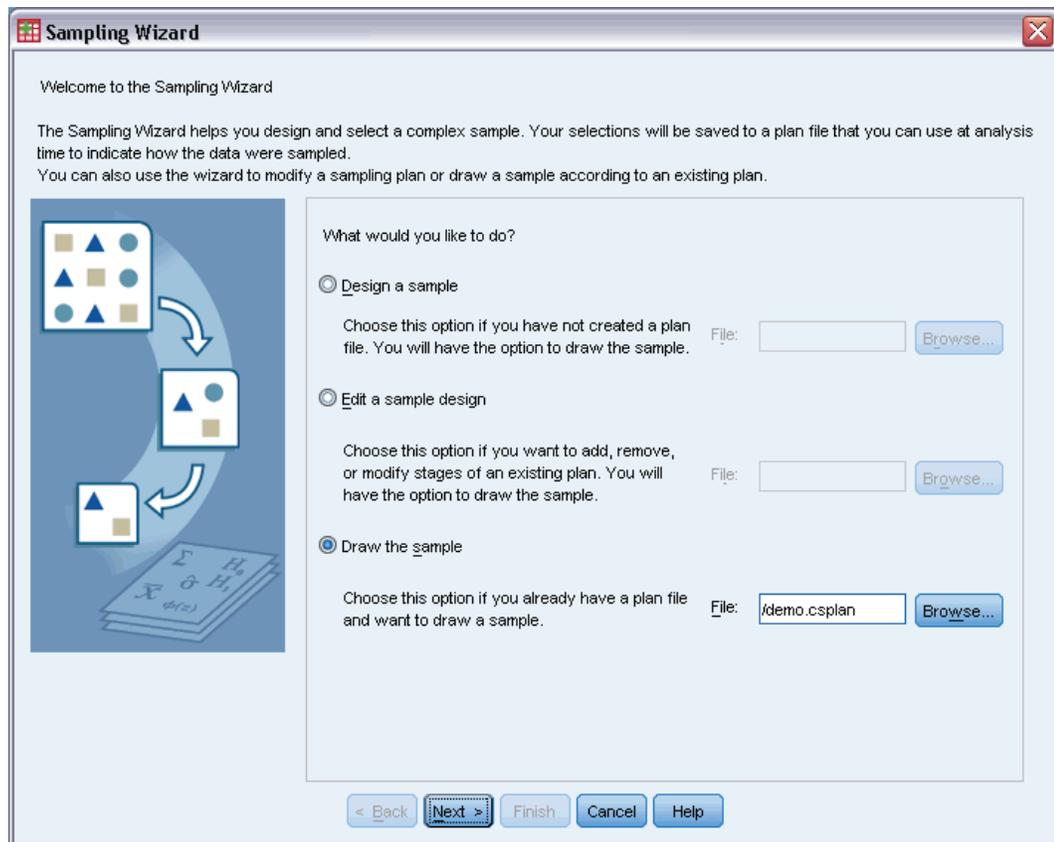
- Cities with values for these variables were selected to the sample.
- Cities with system-missing values for the variables were not selected.

For each city selected, the company acquired subdivision and household unit information and placed it in *demo\_cs\_2.sav*. Use this file and the Sampling Wizard to sample the third stage of this design.

### Using the Wizard to Sample from the Second Partial Frame

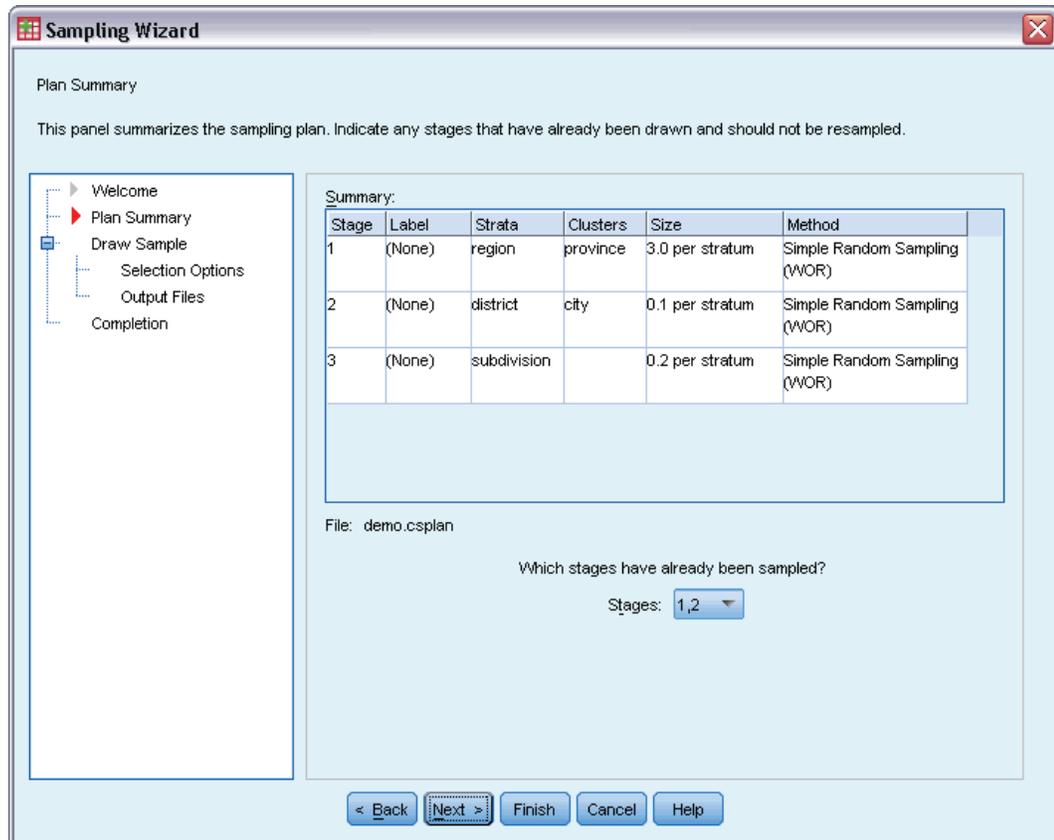
- To run the Complex Samples Sampling Wizard, from the menus choose:  
Analyze > Complex Samples > Select a Sample...

Figure 13-27  
Sampling Wizard, Welcome step



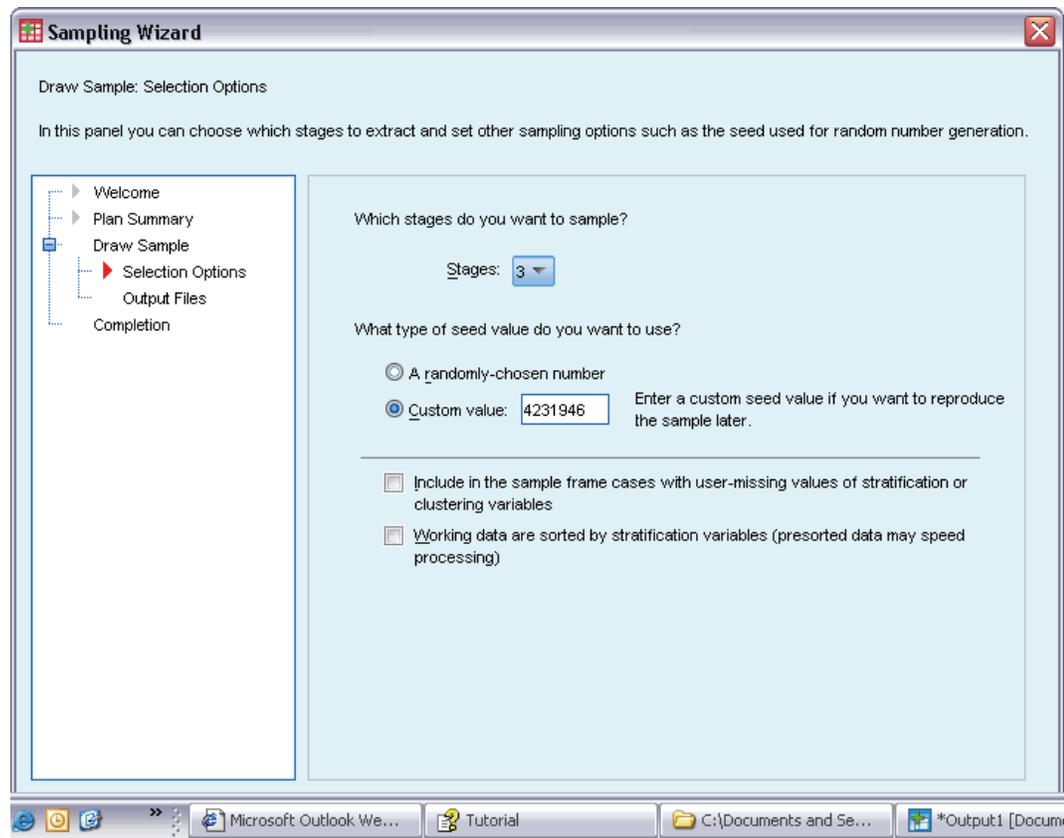
- ▶ Select Draw a sample, browse to where you saved the plan file, and select the demo.csplan plan file that you created.
- ▶ Click Next.

Figure 13-28  
Sampling Wizard, Plan Summary step (stage 3)



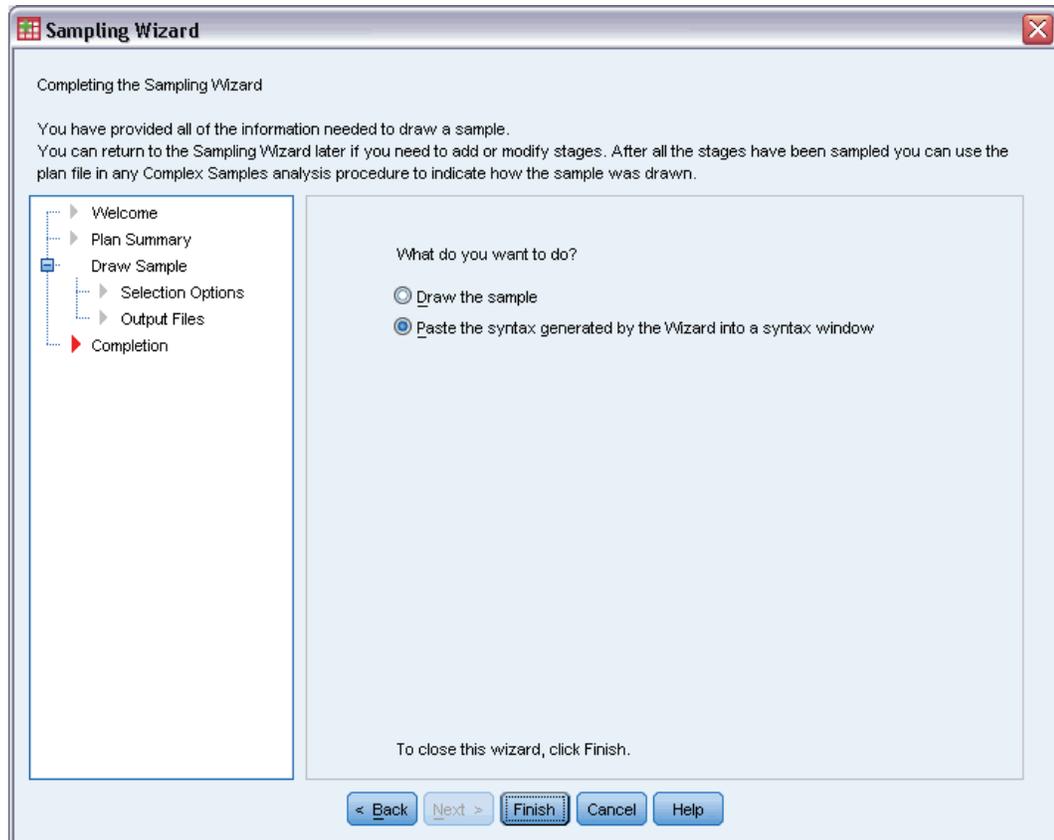
- ▶ Select 1, 2 as stages already sampled.
- ▶ Click Next.

Figure 13-29  
Sampling Wizard, Draw Sample Selection Options step



- ▶ Select Custom value for the type of random seed to use and type 4231946 as the value.
- ▶ Click Next, and then click Next in the Draw Sample Output Files step.

Figure 13-30  
Sampling Wizard, Finish step



- ▶ Select Paste the syntax generated by the Wizard into a syntax window.
- ▶ Click Finish.

The following syntax is generated:

```
* Sampling Wizard.
CSSELECT
/PLAN FILE='demo.csplan'
/CRITERIA STAGES = 3 SEED = 4231946
/CLASSMISSING EXCLUDE
/DATA RENAMEVARS
/PRINT SELECTION.
```

Printing the sampling summary in this case produces a cumbersome table that causes problems in the Output Viewer. To turn off display of the sampling summary, replace SELECTION with CPS in the PRINT subcommand. Then run the syntax within the syntax window.

These selections draw a sample according to the third stage of the *demo.csplan* sampling plan.

## Sample Results

Figure 13-31  
Data Editor with sample results

	city	subdivision	unit	InclusionPr obability_1	SampleWei ghtCumulat ve_1	InclusionPr obability_2	SampleWei ghtCumulat ve_2	InclusionPr obability_3	SampleWei ghtCumulat ve_3	SampleWei ght_Final_
14	190	946	94514	.20	5.00	.10	50.00	.	.	.
15	190	946	94515	.20	5.00	.10	50.00	.	.	.
16	190	946	94516	.20	5.00	.10	50.00	.20	244.44	244.44
17	190	946	94517	.20	5.00	.10	50.00	.	.	.
18	190	946	94518	.20	5.00	.10	50.00	.	.	.
19	190	946	94519	.20	5.00	.10	50.00	.	.	.
20	190	946	94520	.20	5.00	.10	50.00	.	.	.
21	190	946	94521	.20	5.00	.10	50.00	.	.	.
22	190	946	94522	.20	5.00	.10	50.00	.	.	.
23	190	946	94523	.20	5.00	.10	50.00	.	.	.
24	190	946	94524	.20	5.00	.10	50.00	.20	244.44	244.44
25	190	946	94525	.20	5.00	.10	50.00	.	.	.
26	190	946	94526	.20	5.00	.10	50.00	.	.	.
27	190	946	94527	.20	5.00	.10	50.00	.	.	.
28	190	946	94528	.20	5.00	.10	50.00	.	.	.
29	190	946	94529	.20	5.00	.10	50.00	.20	244.44	244.44
30	190	946	94530	.20	5.00	.10	50.00	.	.	.

You can see the sampling results in the Data Editor. Three new variables were saved to the working file, representing the inclusion probabilities and cumulative sampling weights for the third stage, plus the final sampling weights. These new weights take into account the weights computed during the sampling of the first two stages.

- Units with values for these variables were selected to the sample.
- Units with system-missing values for these variables were not selected.

The company will now use its resources to obtain survey information for the housing units selected in the sample. Once the surveys are collected, you can process the sample with Complex Samples analysis procedures, using the sampling plan *demo.csplan* to provide the sampling specifications.

## Sampling with Probability Proportional to Size (PPS)

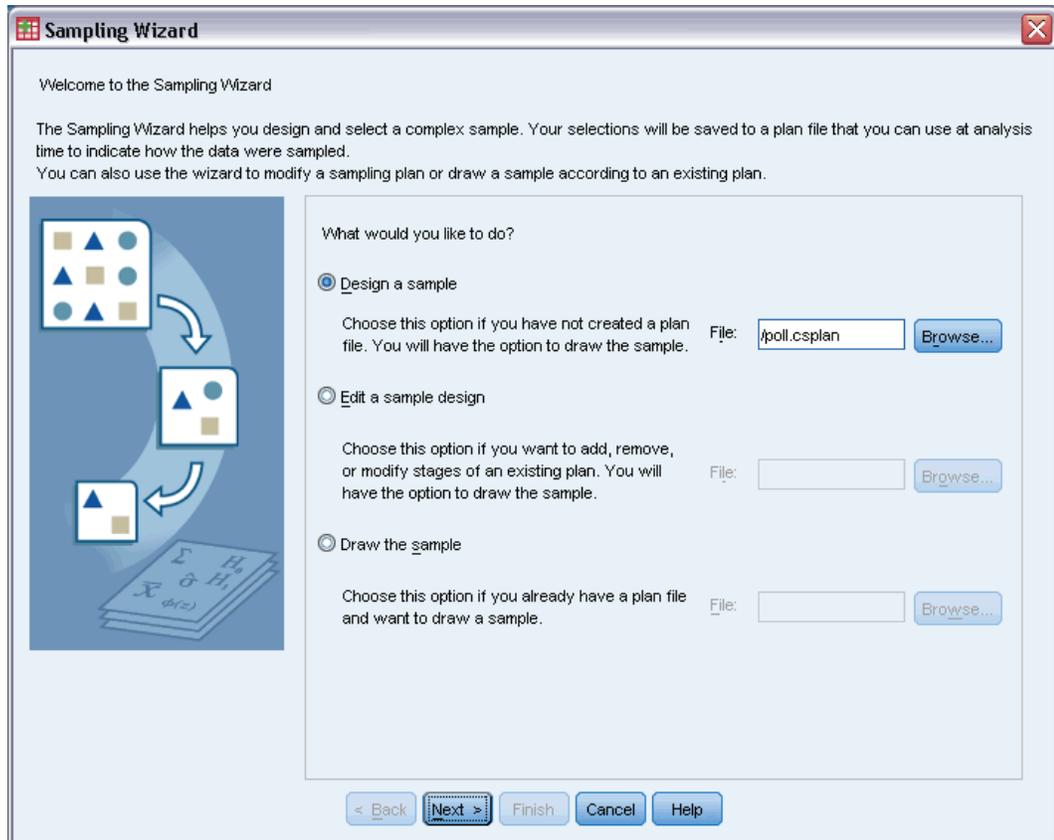
Representatives considering a bill before the legislature are interested in whether there is public support for the bill and how support for the bill is related to voter demographics. Pollsters design and conduct interviews according to a complex sampling design.

A list of registered voters is collected in *poll\_cs.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use the Complex Samples Sampling Wizard to select a sample for further analysis.

## Using the Wizard

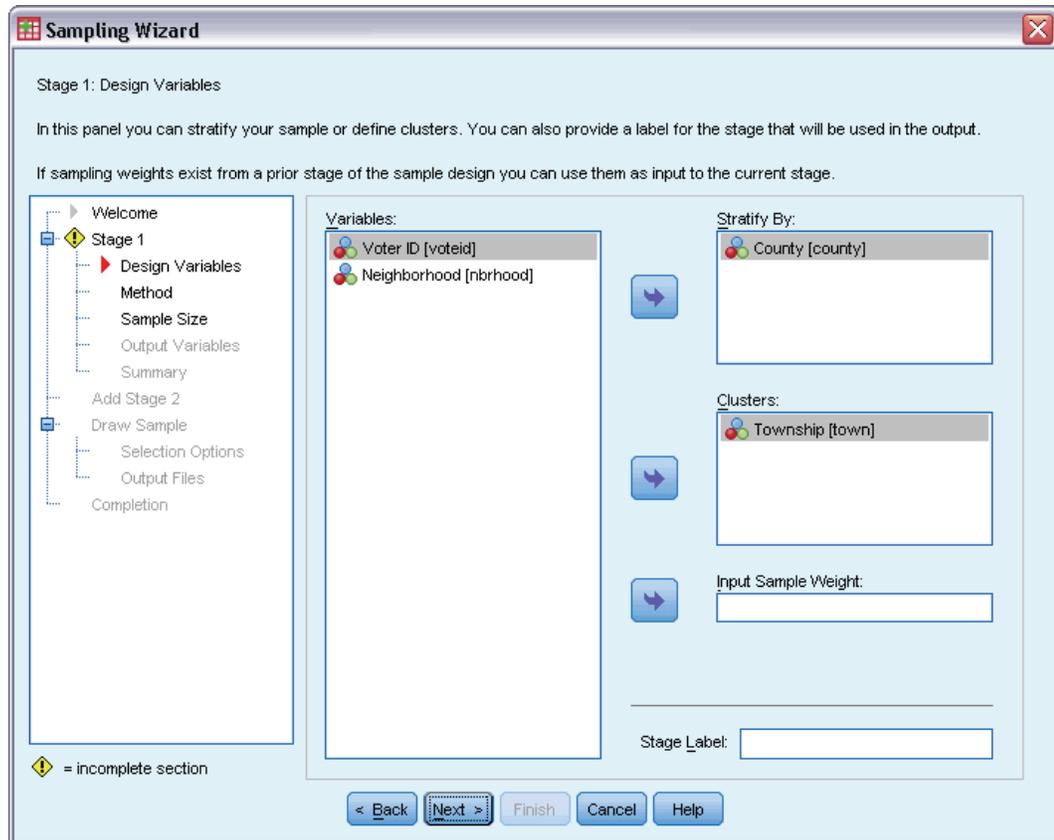
- To run the Complex Samples Sampling Wizard, from the menus choose:  
Analyze > Complex Samples > Select a Sample...

Figure 13-32  
Sampling Wizard, Welcome step



- ▶ Select Design a sample, browse to where you want to save the file, and type poll.csplan as the name of the plan file.
- ▶ Click Next.

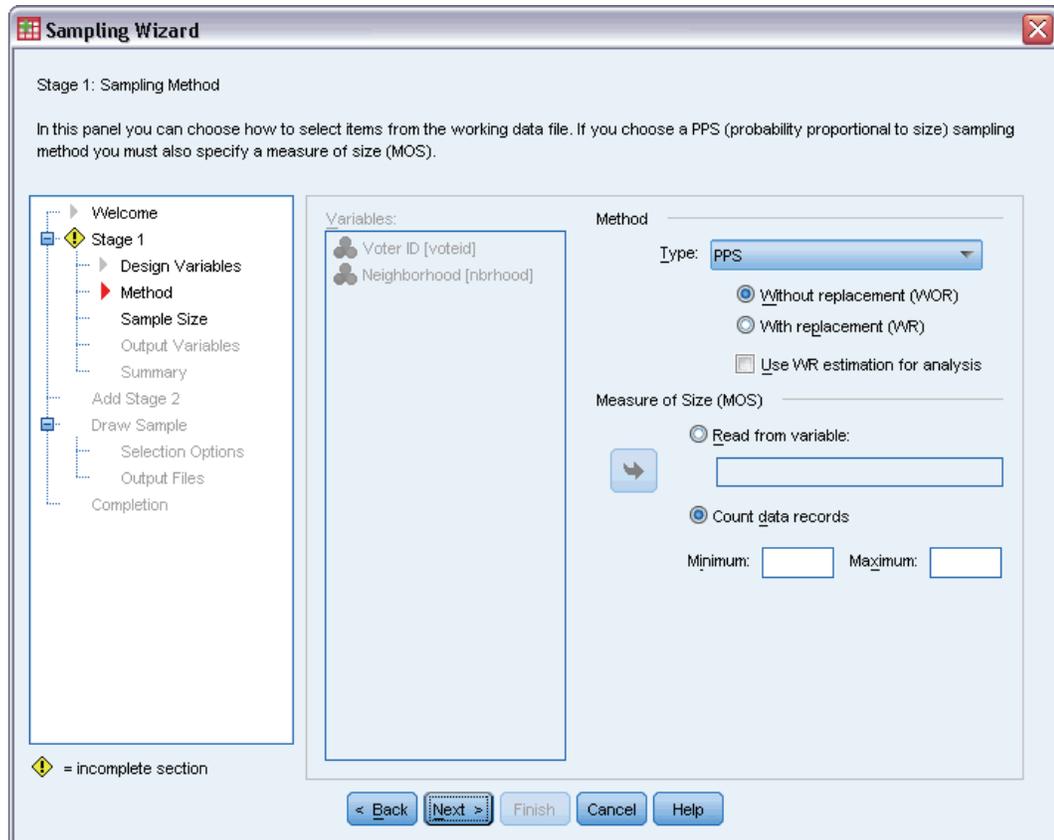
Figure 13-33  
Sampling Wizard, Design Variables step (stage 1)



- ▶ Select *County* as a stratification variable.
- ▶ Select *Township* as a cluster variable.
- ▶ Click Next.

This design structure means that independent samples are drawn for each county. In this stage, townships are drawn as the primary sampling unit.

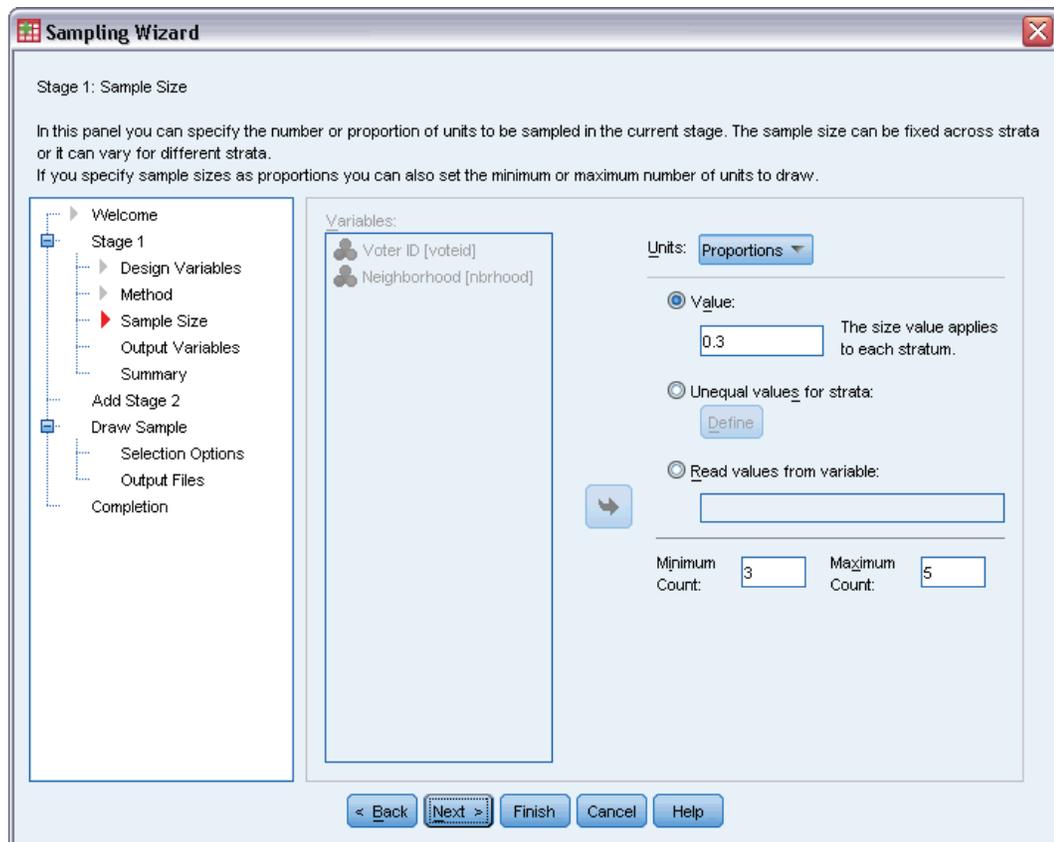
**Figure 13-34**  
*Sampling Wizard, Sampling Method step (stage 1)*



- ▶ Select PPS as the sampling method.
- ▶ Select Count data records as the measure of size.
- ▶ Click Next.

Within each county, townships are drawn without replacement with probability proportional to the number of records for each township. Using a PPS method generates joint sampling probabilities for the townships; you will specify where to save these values in the Output Files step.

Figure 13-35  
Sampling Wizard, Sample Size step (stage 1)

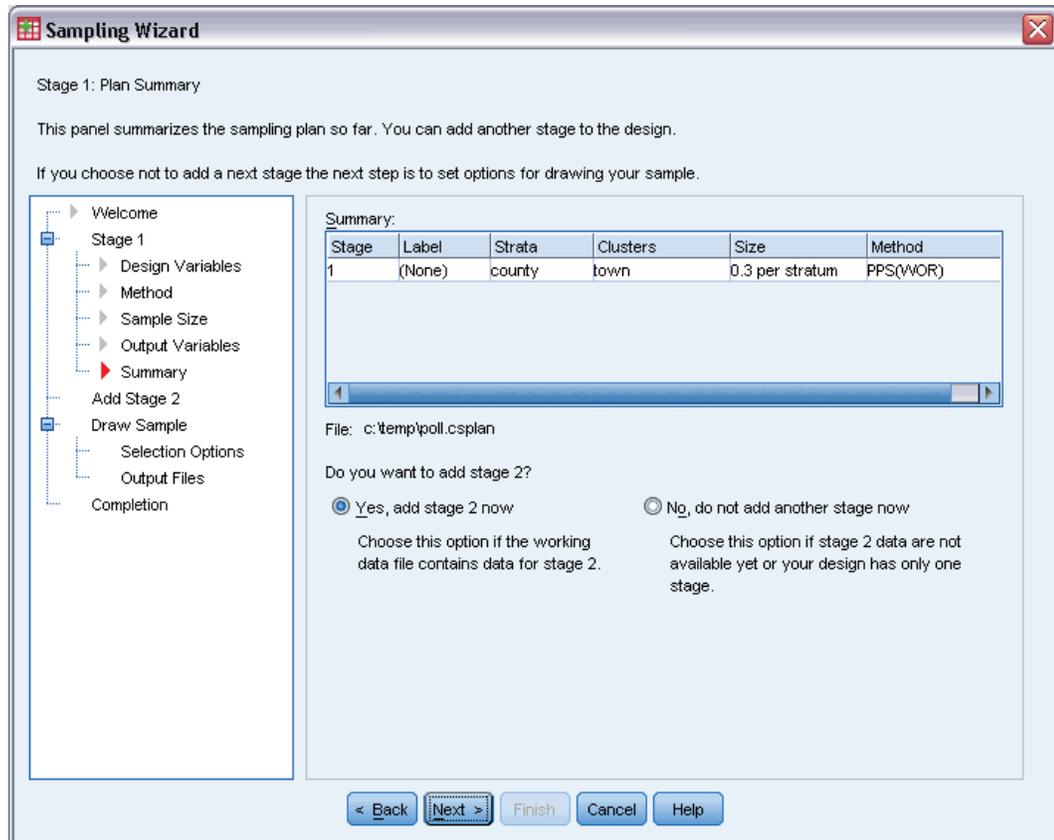


- ▶ Select Proportions from the Units drop-down list.
- ▶ Type 0.3 as the value for the proportion of townships to select per county in this stage.

Legislators from the Western county point out that there are fewer townships in their county than in others. In order to ensure adequate representation, they would like to establish a minimum of 3 townships sampled from each county.

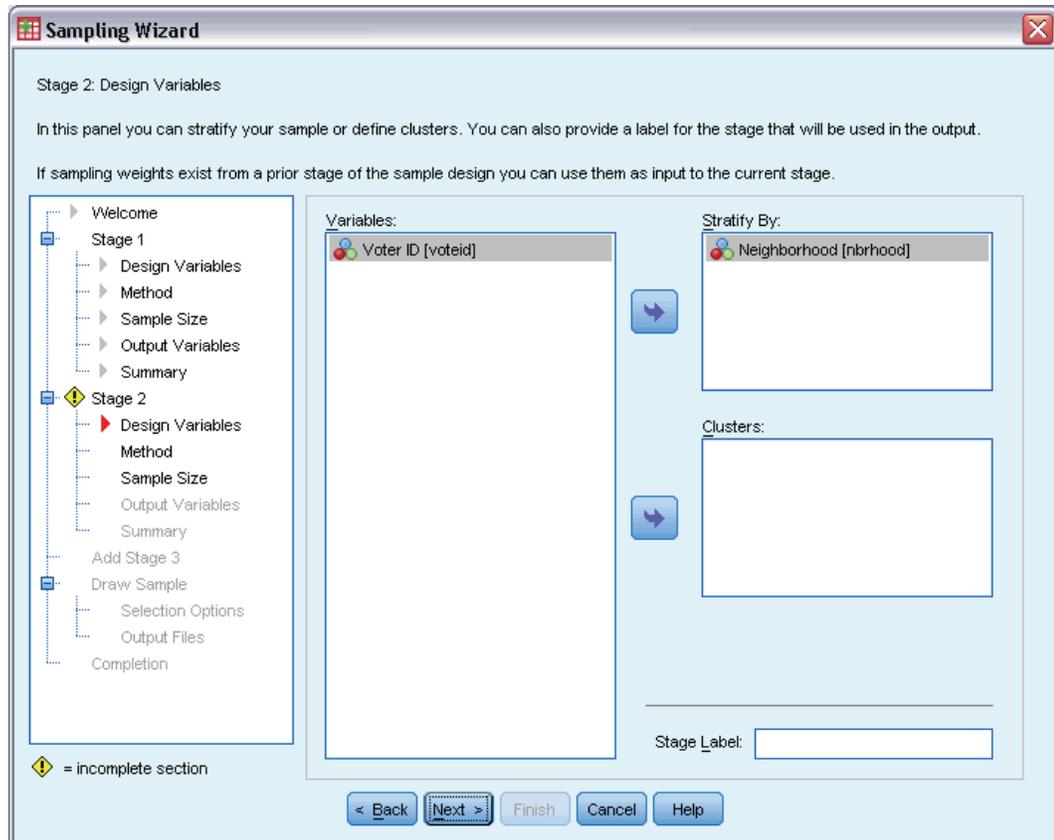
- ▶ Type 3 as the minimum number of townships to select and 5 as the maximum.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-36  
Sampling Wizard, Plan Summary step (stage 1)



- ▶ Select Yes, add stage 2 now.
- ▶ Click Next.

Figure 13-37  
Sampling Wizard, Design Variables step (stage 2)



- ▶ Select *Neighborhood* as a stratification variable.
- ▶ Click Next, and then click Next in the Sampling Method step.

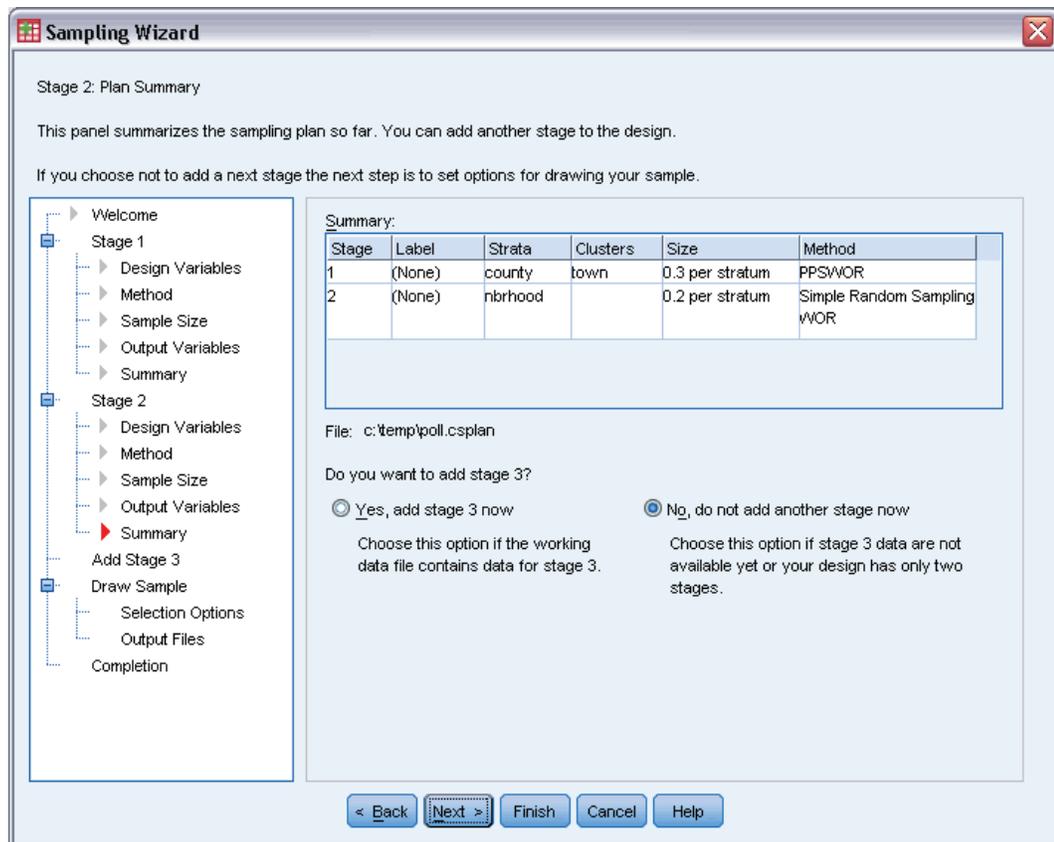
This design structure means that independent samples are drawn for each neighborhood of the townships drawn in stage 1. In this stage, voters are drawn as the primary sampling unit using simple random sampling without replacement.

Figure 13-38  
Sampling Wizard, Sample Size step (stage 2)



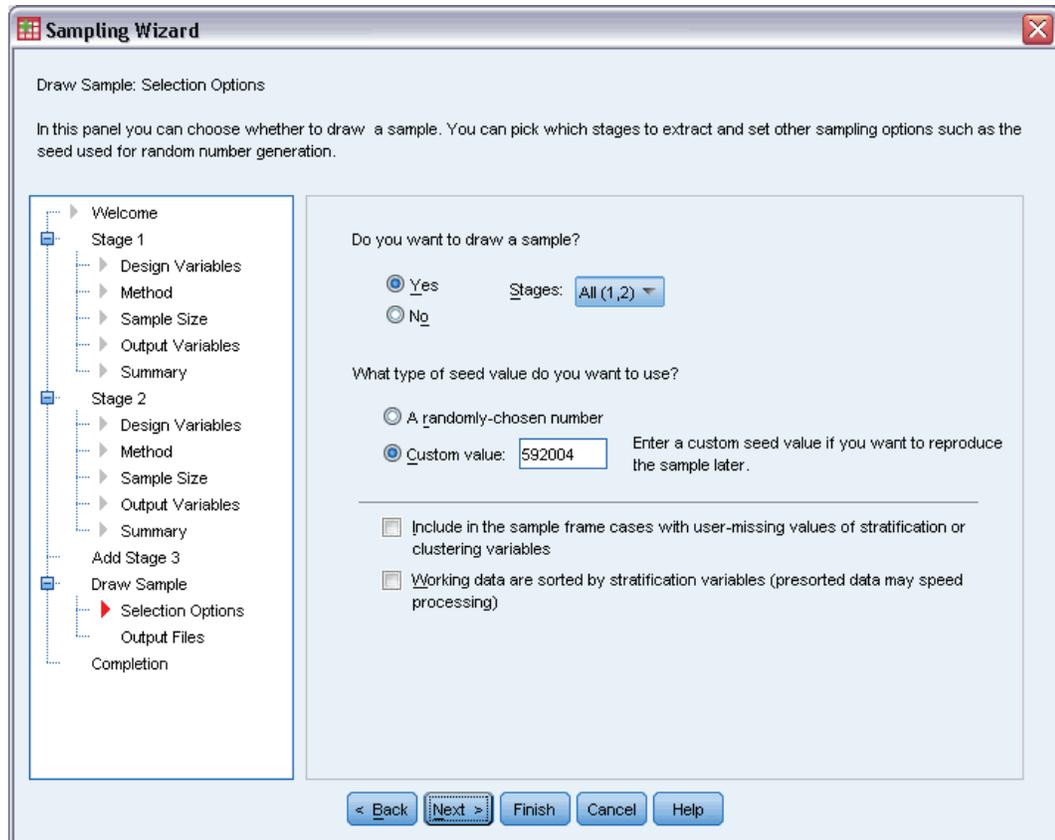
- ▶ Select Proportions from the Units drop-down list.
- ▶ Type 0.2 as the value of the proportion of units to sample from each strata.
- ▶ Click Next, and then click Next in the Output Variables step.

Figure 13-39  
Sampling Wizard, Plan Summary step (stage 2)



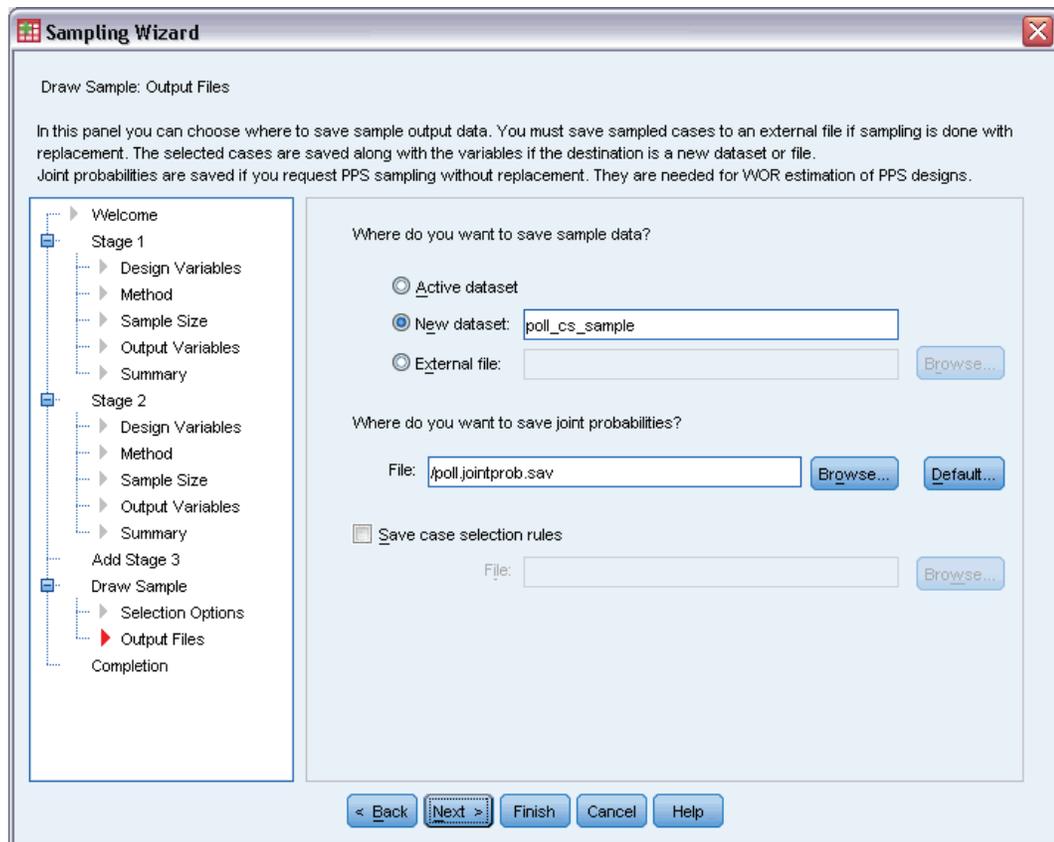
- Look over the sampling design, and then click Next.

Figure 13-40  
Sampling Wizard, Draw Sample Selection Options step



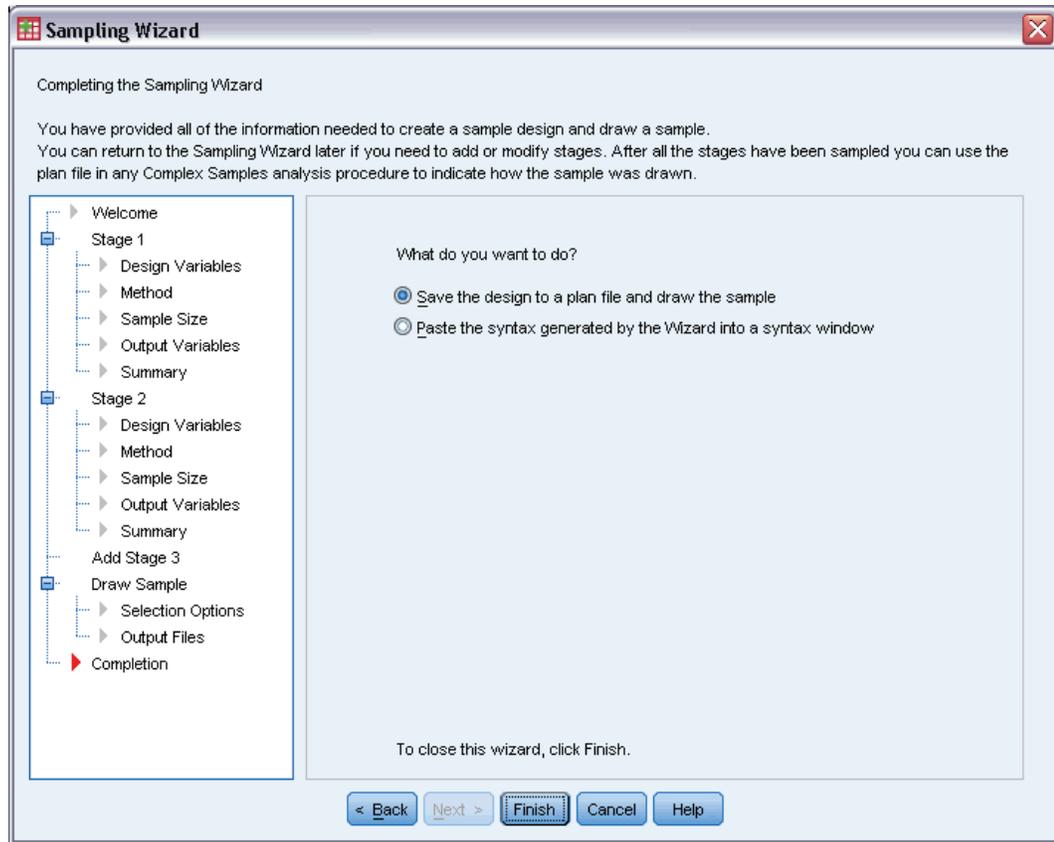
- ▶ Select Custom value for the type of random seed to use, and type 592004 as the value. Using a custom value allows you to replicate the results of this example exactly.
- ▶ Click Next.

Figure 13-41  
Sampling Wizard, Draw Sample Selection Options step



- ▶ Choose to save the sample to a new dataset, and type poll\_cs\_sample as the name of the dataset.
- ▶ Browse to where you want to save the joint probabilities and type poll\_jointprob.sav as the name of the joint probabilities file.
- ▶ Click Next.

Figure 13-42  
Sampling Wizard, Finish step



- Click Finish.

These selections produce the sampling plan file *poll.csplan* and draw a sample according to that plan, save the sample results to the new dataset *poll\_cs\_sample*, and save the joint probabilities file to the external data file *poll\_jointprob.sav*.

## Plan Summary

Figure 13-43  
Plan summary

			Stage 1	Stage 2
Design	Stratification	1	County	Neighborhood
Variables	Cluster	1	Township	
Sample Information	Selection Method		PPS sampling without replacement	Simple random sampling without replacement
	Measure of Size		Obtained from data	
	Proportion of Units Sampled		.3	.2
	Minimum Number of Units Sampled		3	
	Maximum Number of Units Sampled		5	
	Variables Created or Modified	Stagewise Inclusion (Selection) Probability Stagewise Cumulative Sample Weight	Inclusion Probability_1_ Sample/Weight Cumulative_1_	Inclusion Probability_2_ Sample/Weight Cumulative_2_
Analysis Information	Estimator Assumption		Unequal probability sampling without replacement (using joint inclusion probabilities)	Equal probability sampling without replacement
	Inclusion Probability		Obtained from variable Inclusion Probability_1_	Obtained from variable Inclusion Probability_2_

Plan File: c:\poll.csplan  
Weight Variable: SampleWeight\_Final\_

The summary table reviews your sampling plan and is useful for making sure that the plan represents your intentions.

## Sampling Summary

Figure 13-44  
Stage summary

County	Number of Units Sampled		Proportion of Units Sampled	
	Requested	Actual	Requested	Actual
Eastern	4	4	30.0%	30.8%
Central	4	4	30.0%	30.8%
Western	3	3	30.0%	50.0%
Northern	5	5	30.0%	33.3%
Southern	3	3	30.0%	50.0%

Plan File: c:\poll.csplan

This summary table reviews the first stage of sampling and is useful for checking that the sampling went according to plan. Recall that you requested a 30% sample of townships by county; the actual proportions sampled are close to 30%, except in the Western and Southern counties. This is because these counties each have only six townships, and you also specified that a minimum of three townships should be selected per county.

**Figure 13-45**  
*Stage summary*

County	Township	Neighborhood	Number of Units Sampled		Proportion of Units Sampled		
			Requested	Actual	Requested	Actual	
Eastern	9	1	49	49	20.0%	19.9%	
		2	143	143	20.0%	20.0%	
		3	113	113	20.0%	20.0%	
		4	77	77	20.0%	20.0%	
		5	139	139	20.0%	20.0%	
		6	120	120	20.0%	20.0%	
	10	1	149	149	20.0%	20.1%	
		2	117	117	20.0%	20.0%	
		3	116	116	20.0%	20.0%	
		4	69	69	20.0%	19.9%	
	11	1	65	65	20.0%	19.9%	
		2	72	72	20.0%	19.9%	
		3	109	109	20.0%	20.0%	
		4	140	140	20.0%	20.0%	
		5	42	42	20.0%	19.8%	
		6	142	142	20.0%	20.0%	
	12	1	145	145	20.0%	20.1%	
		2	69	69	20.0%	20.1%	
		3	98	98	20.0%	20.1%	
		4	134	134	20.0%	20.0%	
		5	114	114	20.0%	20.0%	
		6	137	137	20.0%	19.9%	
	Central	2	1	119	119	20.0%	20.1%
			2	153	153	20.0%	19.9%
3			101	101	20.0%	20.0%	
4			52	52	20.0%	19.8%	
5			144	144	20.0%	20.0%	

Plan File: c:\poll.csplan

This summary table (the top part of which is shown here) reviews the second stage of sampling. It is also useful for checking that the sampling went according to plan. Approximately 20% of the voters were sampled from each neighborhood from each township sampled in the first stage, as requested.

## Sample Results

Figure 13-46  
Data Editor with sample results

	voteid	nbrhood	town	county	InclusionPr obability_1_	SampleWei ghtCumulat ve_1	InclusionPr obability_2_	SampleWei ghtCumulat ve_2	SampleWei ght_Final_
376	368	4	9	1	.44	2.26	.20	11.28	11.28
377	369	4	9	1	.44	2.26	.20	11.28	11.28
378	374	4	9	1	.44	2.26	.20	11.28	11.28
379	376	4	9	1	.44	2.26	.20	11.28	11.28
380	379	4	9	1	.44	2.26	.20	11.28	11.28
381	380	4	9	1	.44	2.26	.20	11.28	11.28
382	382	4	9	1	.44	2.26	.20	11.28	11.28
383	13	5	9	1	.44	2.26	.20	11.26	11.26
384	18	5	9	1	.44	2.26	.20	11.26	11.26
385	23	5	9	1	.44	2.26	.20	11.26	11.26
386	38	5	9	1	.44	2.26	.20	11.26	11.26
387	39	5	9	1	.44	2.26	.20	11.26	11.26
388	40	5	9	1	.44	2.26	.20	11.26	11.26
389	41	5	9	1	.44	2.26	.20	11.26	11.26
390	43	5	9	1	.44	2.26	.20	11.26	11.26

You can see the sampling results in the newly created dataset. Five new variables were saved to the working file, representing the inclusion probabilities and cumulative sampling weights for each stage, plus the final sampling weights. Voters who were not selected to the sample are excluded from this dataset.

The final sampling weights are identical for voters within the same neighborhood because they are selected according to a simple random sampling method within neighborhoods. However, they are different across neighborhoods within the same township because the sampled proportions are not exactly 20% in all neighborhoods.

Figure 13-47  
Data Editor with sample results

	voteid	nbrhood	town	county	InclusionPr obability_1	SampleWei ghtCumulat ve_1	InclusionPr obability_2	SampleWei ghtCumulat ve_2	SampleWei ght_Final_
635	577	6	9	1	.44	2.26	.20	11.30	11.30
636	578	6	9	1	.44	2.26	.20	11.30	11.30
637	582	6	9	1	.44	2.26	.20	11.30	11.30
638	590	6	9	1	.44	2.26	.20	11.30	11.30
639	594	6	9	1	.44	2.26	.20	11.30	11.30
640	597	6	9	1	.44	2.26	.20	11.30	11.30
641	600	6	9	1	.44	2.26	.20	11.30	11.30
642	4	1	10	1	.31	3.21	.20	16.00	16.00
643	5	1	10	1	.31	3.21	.20	16.00	16.00
644	9	1	10	1	.31	3.21	.20	16.00	16.00
645	10	1	10	1	.31	3.21	.20	16.00	16.00
646	12	1	10	1	.31	3.21	.20	16.00	16.00
647	16	1	10	1	.31	3.21	.20	16.00	16.00
648	17	1	10	1	.31	3.21	.20	16.00	16.00
649	19	1	10	1	.31	3.21	.20	16.00	16.00

Unlike voters in the second stage, the first-stage sampling weights are not identical for townships within the same county because they are selected with probability proportional to size.

Figure 13-48  
Joint probabilities file

	county	town	Unit_No_	Joint_Prob _1_	Joint_Prob _2_	Joint_Prob _3_	Joint_Prob _4_	Joint_Prob _5_
1	1	10	1	.31	.10	.11	.12	.
2	1	11	2	.10	.39	.15	.16	.
3	1	9	3	.11	.15	.44	.21	.
4	1	12	4	.12	.16	.21	.48	.
5	2	12	1	.22	.04	.07	.08	.
6	2	6	2	.04	.23	.07	.08	.
7	2	7	3	.07	.07	.41	.19	.
8	2	2	4	.08	.08	.19	.45	.
9	3	5	1	.58	.31	.32	.	.
10	3	3	2	.31	.61	.36	.	.
11	3	4	3	.32	.36	.63	.	.
12	4	14	1	.26	.06	.06	.07	.09
13	4	8	2	.06	.29	.07	.08	.10
14	4	4	3	.06	.07	.29	.08	.10
15	4	2	4	.07	.08	.08	.33	.12
16	4	13	5	.09	.10	.10	.12	.43
17	5	3	1	.74	.25	.27	.	.
18	5	6	2	.25	.41	.13	.	.
19	5	4	3	.27	.13	.43	.	.

The file *poll\_jointprob.sav* contains first-stage joint probabilities for selected townships within counties. *County* is a first-stage stratification variable, and *Township* is a cluster variable. Combinations of these variables identify all first-stage PSUs uniquely. *Unit\_No\_* labels PSUs within each stratum and is used to match up with *Joint\_Prob\_1\_*, *Joint\_Prob\_2\_*, *Joint\_Prob\_3\_*, *Joint\_Prob\_4\_*, and *Joint\_Prob\_5\_*. The first two strata each have 4 PSUs; therefore, the joint

inclusion probability matrices are 4×4 for these strata, and the *Joint\_Prob\_5\_* column is left empty for these rows. Similarly, strata 3 and 5 have 3×3 joint inclusion probability matrices, and stratum 4 has a 5×5 joint inclusion probability matrix.

The need for a joint probabilities file is seen by perusing the values of the joint inclusion probability matrices. When the sampling method is not a PPS WOR method, the selection of a PSU is independent of the selection of another PSU, and their joint inclusion probability is simply the product of their inclusion probabilities. In contrast, the joint inclusion probability for Townships 9 and 10 of County 1 is approximately 0.11 (see the first case of *Joint\_Prob\_3\_* or the third case of *Joint\_Prob\_1\_*), or less than the product of their individual inclusion probabilities (the product of the first case of *Joint\_Prob\_1\_* and the third case of *Joint\_Prob\_3\_* is  $0.31 \times 0.44 = 0.1364$ ).

The pollsters will now conduct interviews for the selected sample. Once the results are available, you can process the sample with Complex Samples analysis procedures, using the sampling plan *poll.csplan* to provide the sampling specifications and *poll\_jointprob.sav* to provide the needed joint inclusion probabilities.

## **Related Procedures**

The Complex Samples Sampling Wizard procedure is a useful tool for creating a sampling plan file and drawing a sample.

- To ready a sample for analysis when you do not have access to the sampling plan file, use the [Analysis Preparation Wizard](#).

# ***Complex Samples Analysis Preparation Wizard***

The Analysis Preparation Wizard guides you through the steps for creating or modifying an analysis plan for use with the various Complex Samples analysis procedures. It is most useful when you do not have access to the sampling plan file used to draw the sample.

## ***Using the Complex Samples Analysis Preparation Wizard to Ready NHIS Public Data***

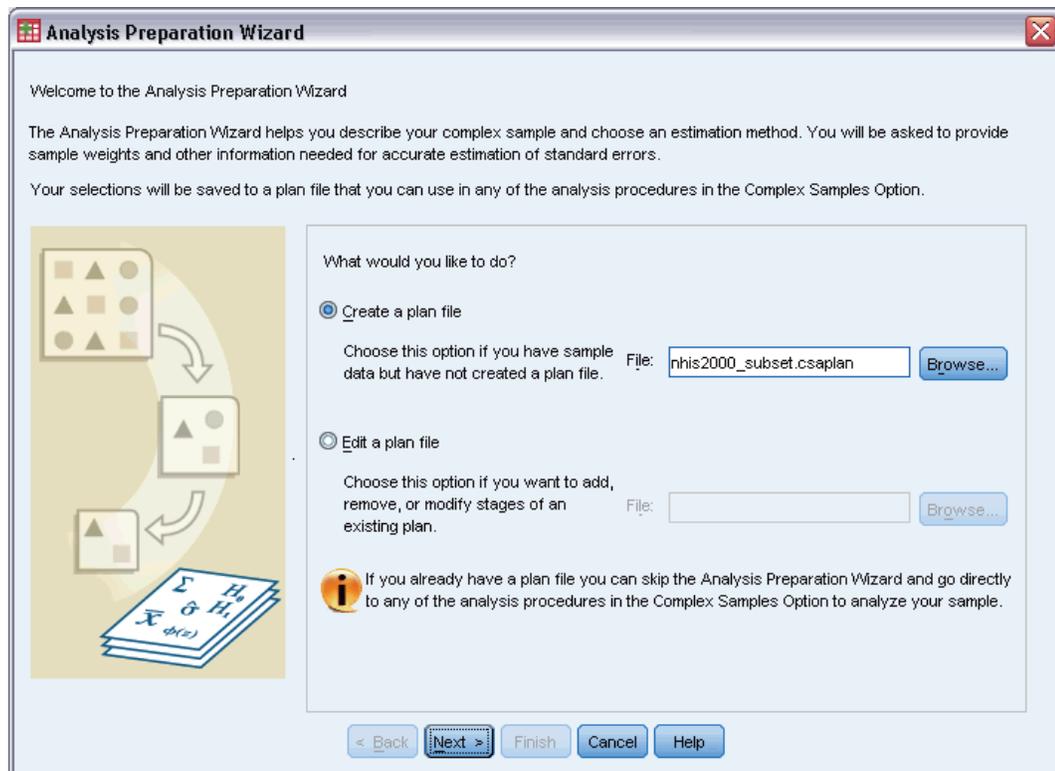
The National Health Interview Survey (NHIS) is a large, population-based survey of the U.S. civilian population. Interviews are carried out face-to-face in a nationally representative sample of households. Demographic information and observations about health behavior and status are obtained for members of each household.

A subset of the 2000 survey is collected in *nhis2000\_subset.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use the Complex Samples Analysis Preparation Wizard to create an analysis plan for this data file so that it can be processed by Complex Samples analysis procedures.

### ***Using the Wizard***

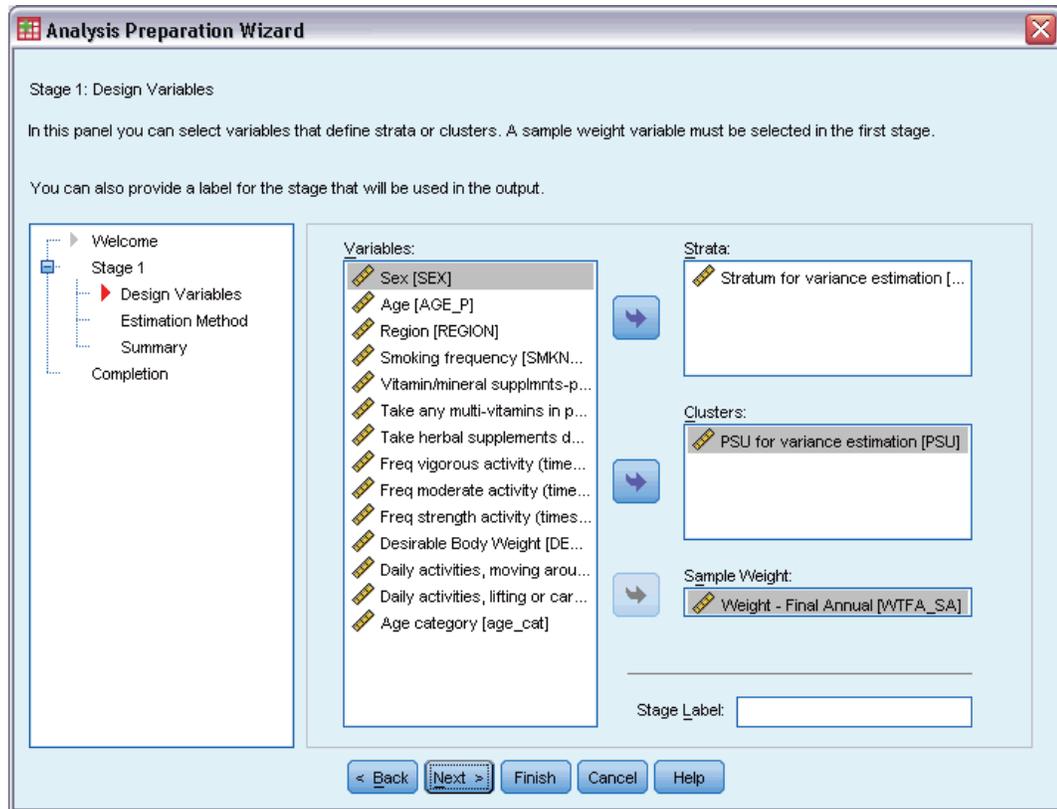
- ▶ To prepare a sample using the Complex Samples Analysis Preparation Wizard, from the menus choose:  
Analyze > Complex Samples > Prepare for Analysis...

Figure 14-1  
Analysis Preparation Wizard, Welcome step



- ▶ Browse to where you want to save the plan file and type `nhis2000_subset.csaplan` as the name for the analysis plan file.
- ▶ Click Next.

Figure 14-2  
Analysis Preparation Wizard, Design Variables step (stage 1)



The data are obtained using a complex multistage sample. However, for end users, the original NHIS design variables were transformed to a simplified set of design and weight variables whose results approximate those of the original design structures.

- ▶ Select *Stratum for variance estimation* as a strata variable.
- ▶ Select *PSU for variance estimation* as a cluster variable.
- ▶ Select *Weight - Final Annual* as the sample weight variable.
- ▶ Click Finish.

## Summary

Figure 14-3  
Summary

			Stage 1
Design Variables	Stratification	1	Stratum for variance estimation
	Cluster	1	PSU for variance estimation
Analysis Information	Estimator Assumption		Sampling with replacement

Plan File: c:\nhis2000\_subset.csaplan  
Weight Variable: Weight - Final Annual  
SRS Estimator: Sampling without replacement

The summary table reviews your analysis plan. The plan consists of one stage with a design of one stratification variable and one cluster variable. With-replacement (WR) estimation is used, and the plan is saved to *c:\nhis2000\_subset.csaplan*. You can now use this plan file to process *nhis2000\_subset.sav* with Complex Samples analysis procedures.

## Preparing for Analysis When Sampling Weights Are Not in the Data File

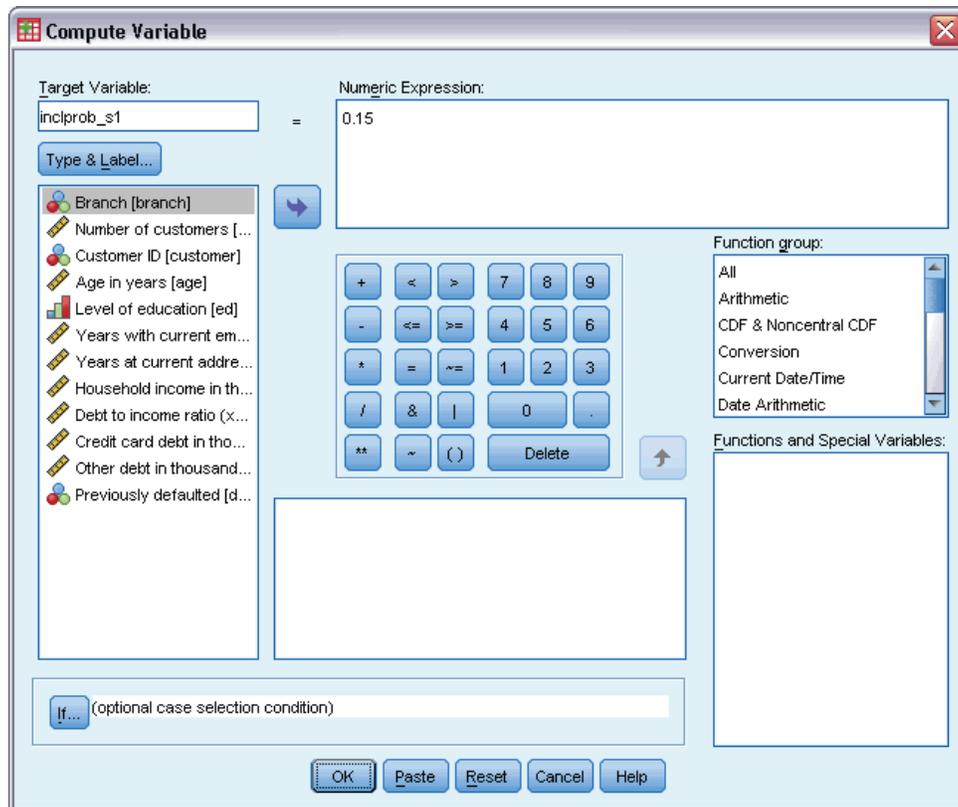
A loan officer has a collection of customer records, taken according to a complex design; however, the sampling weights are not included in the file. This information is contained in *bankloan\_cs\_noweights.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Starting with what she knows about the sampling design, the officer wants to use the Complex Samples Analysis Preparation Wizard to create an analysis plan for this data file so that it can be processed by Complex Samples analysis procedures.

The loan officer knows that the records were selected in two stages, with 15 out of 100 bank branches selected with equal probability and without replacement in the first stage. One hundred customers were then selected from each of those banks with equal probability and without replacement in the second stage, and information on the number of customers at each bank is included in the data file. The first step to creating an analysis plan is to compute the stagewise inclusion probabilities and final sampling weights.

### Computing Inclusion Probabilities and Sampling Weights

- ▶ To compute the inclusion probabilities for the first stage, from the menus choose:  
Transform > Compute Variable...

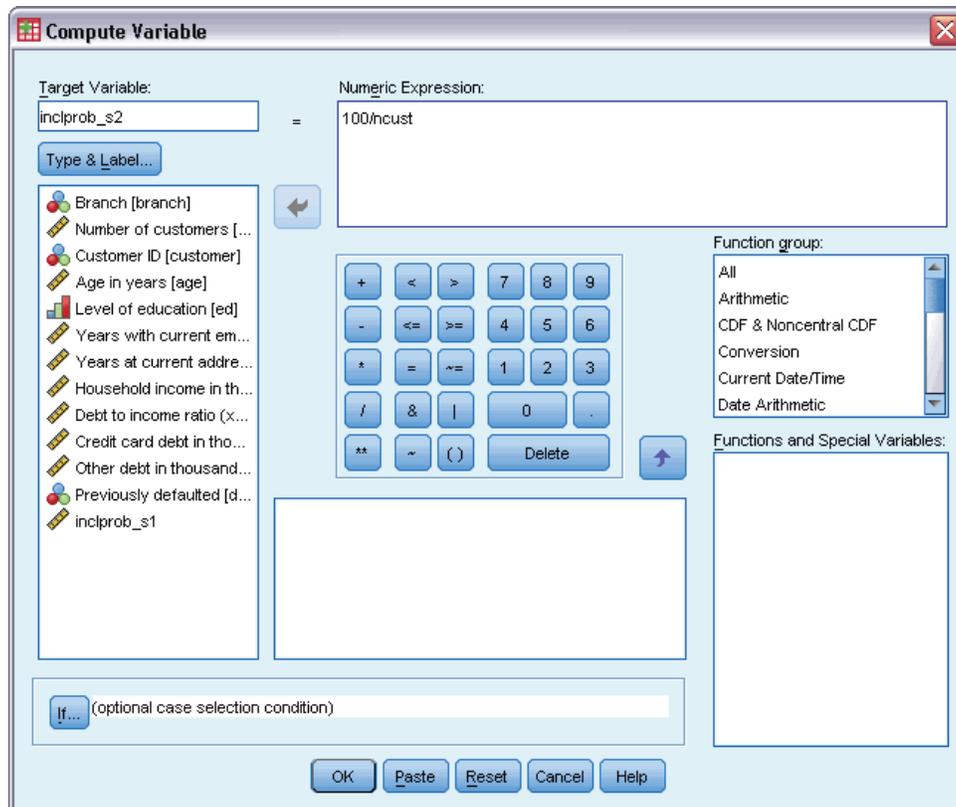
Figure 14-4  
Compute Variable dialog box



Fifteen out of one hundred bank branches were selected without replacement in the first stage; thus, the probability that a given bank was selected is  $15/100 = 0.15$ .

- ▶ Type `inclprob_s1` as the target variable.
- ▶ Type `0.15` as the numeric expression.
- ▶ Click OK.

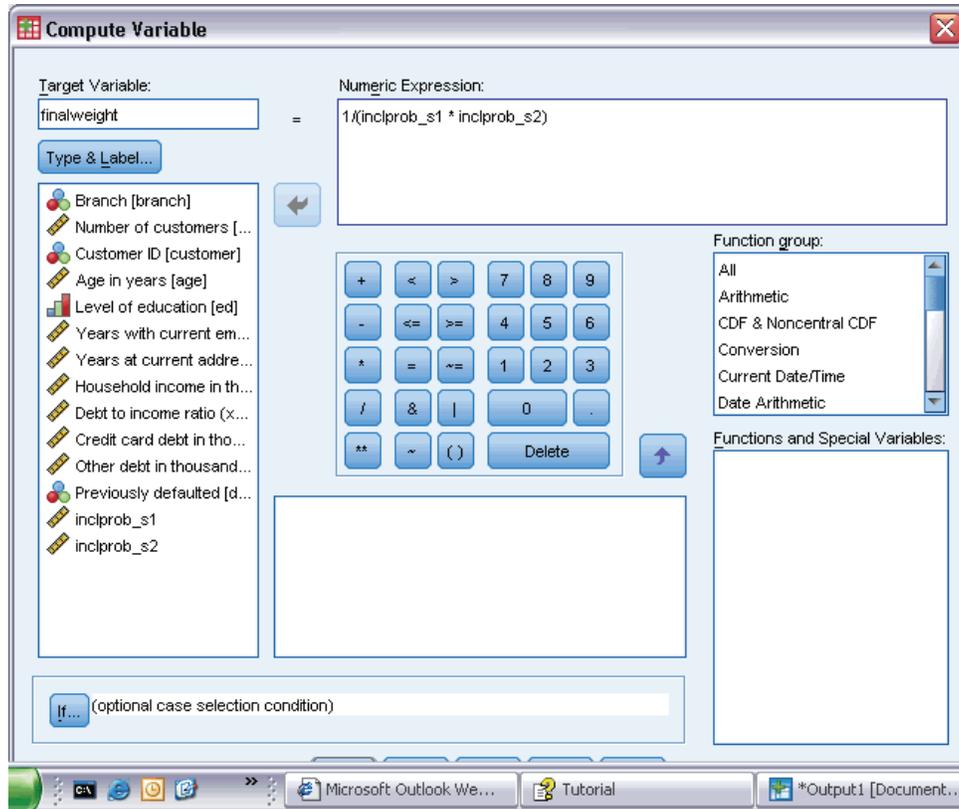
Figure 14-5  
Compute Variable dialog box



One hundred customers were selected from each branch in the second stage; thus, the stage 2 inclusion probability for a given customer at a given bank is  $100/\text{the number of customers at that bank}$ .

- ▶ Recall the Compute Variable dialog box.
- ▶ Type `inclprob_s2` as the target variable.
- ▶ Type `100/ncust` as the numeric expression.
- ▶ Click OK.

Figure 14-6  
Compute Variable dialog box



Now that you have the inclusion probabilities for each stage, it's easy to compute the final sampling weights.

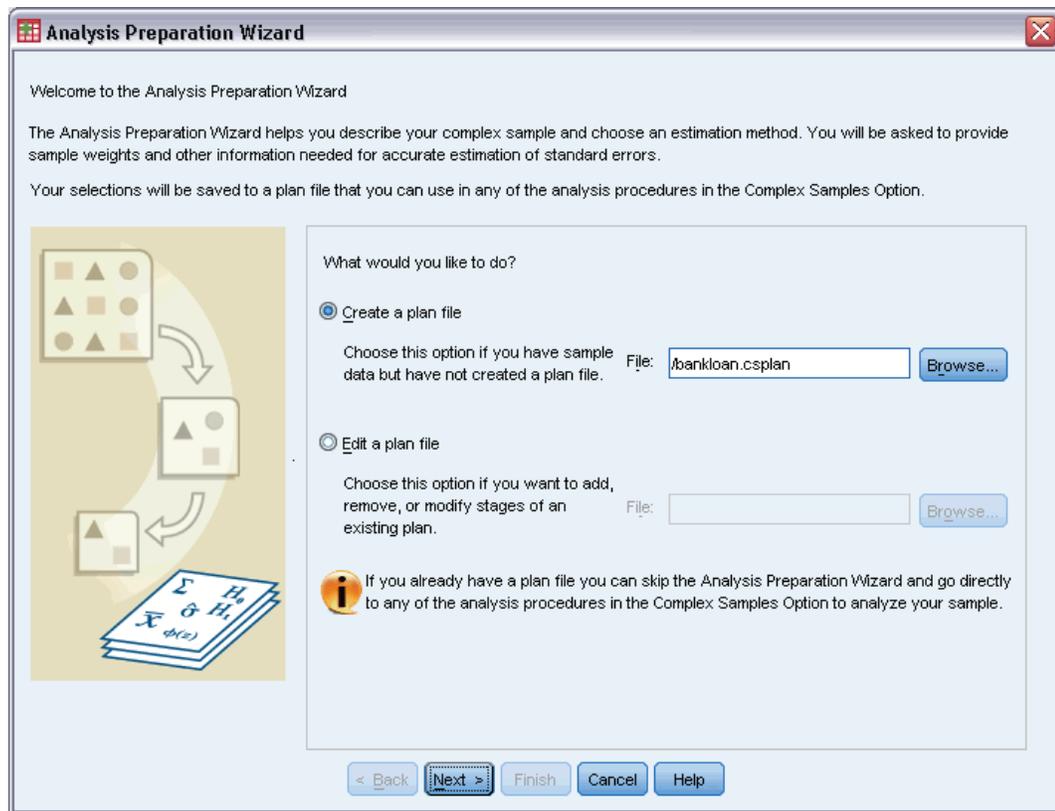
- ▶ Recall the Compute Variable dialog box.
- ▶ Type finalweight as the target variable.
- ▶ Type  $1/(inclprob\_s1 * inclprob\_s2)$  as the numeric expression.
- ▶ Click OK.

You are now ready to create the analysis plan.

### Using the Wizard

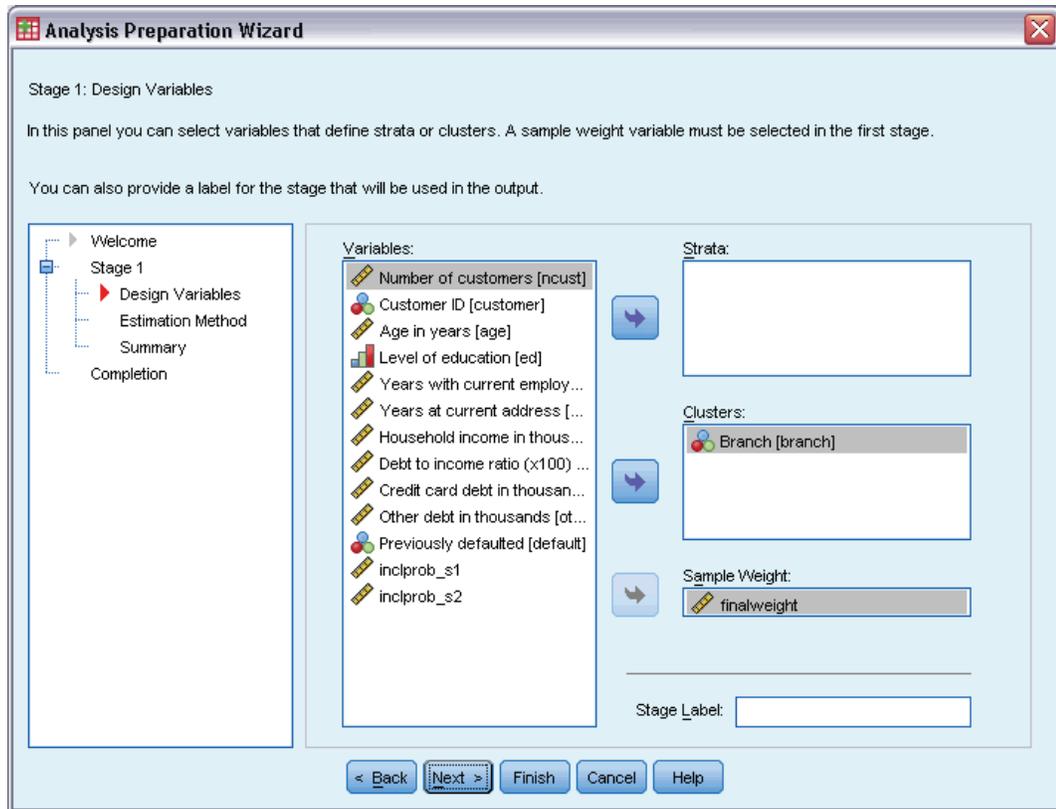
- ▶ To prepare a sample using the Complex Samples Analysis Preparation Wizard, from the menus choose:  
Analyze > Complex Samples > Prepare for Analysis...

Figure 14-7  
Analysis Preparation Wizard, Welcome step



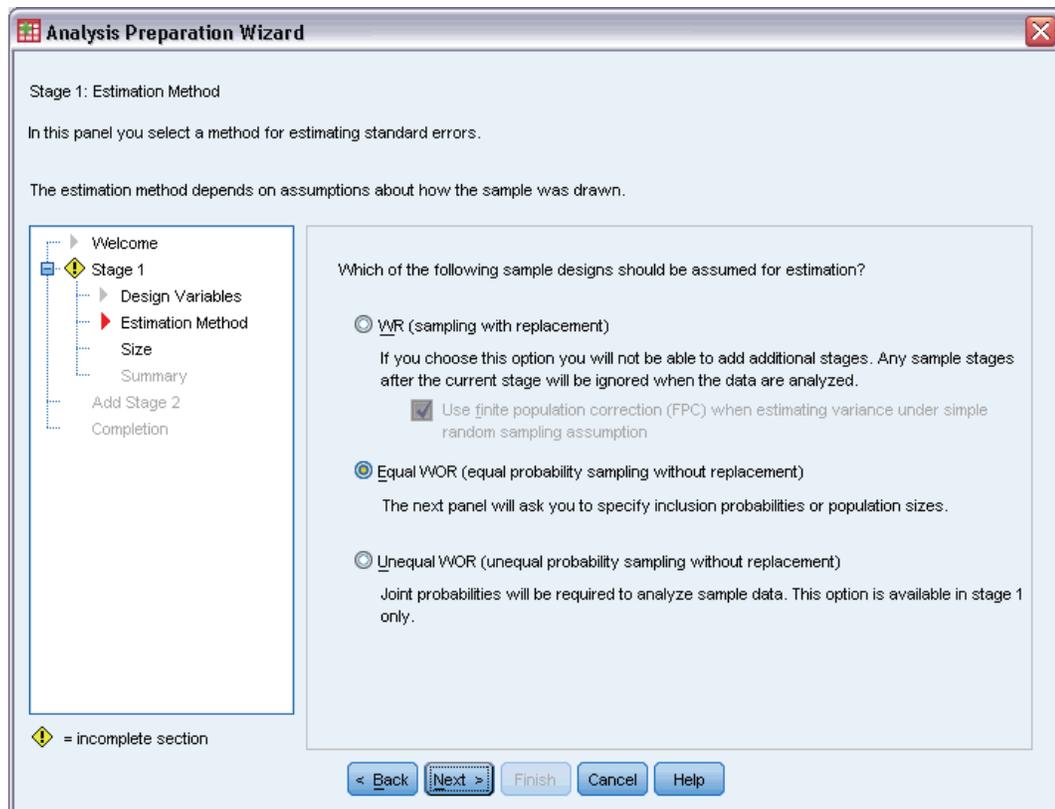
- ▶ Browse to where you want to save the plan file and type bankloan.csplan as the name for the analysis plan file.
- ▶ Click Next.

Figure 14-8  
Analysis Preparation Wizard, Design Variables step (stage 1)



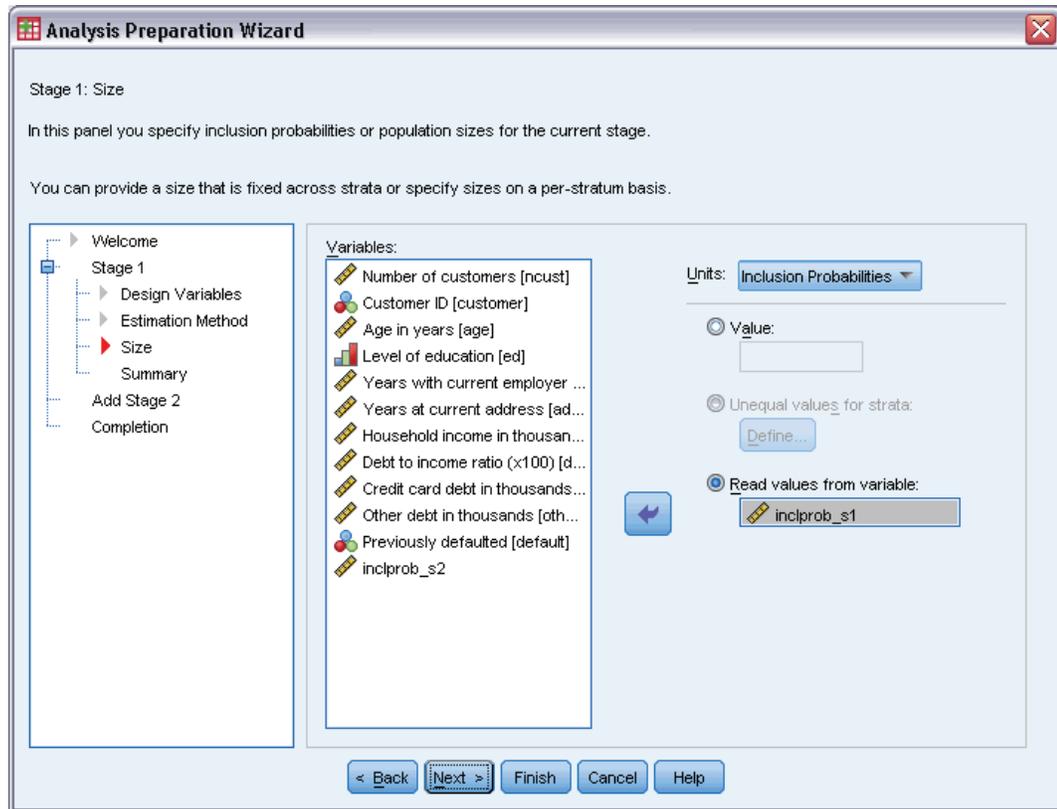
- ▶ Select *Branch* as a cluster variable.
- ▶ Select *finalweight* as the sample weight variable.
- ▶ Click Next.

Figure 14-9  
Analysis Preparation Wizard, Estimation Method step (stage 1)



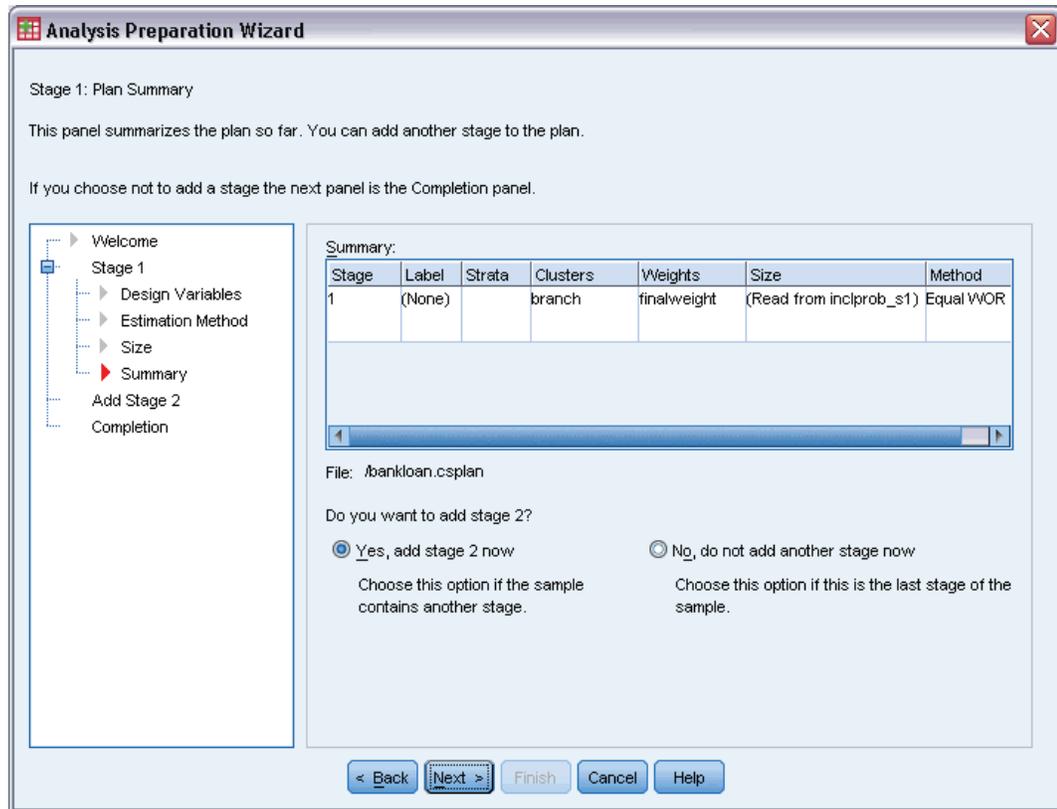
- ▶ Select Equal WOR as the first-stage estimation method.
- ▶ Click Next.

Figure 14-10  
Analysis Preparation Wizard, Size step (stage 1)



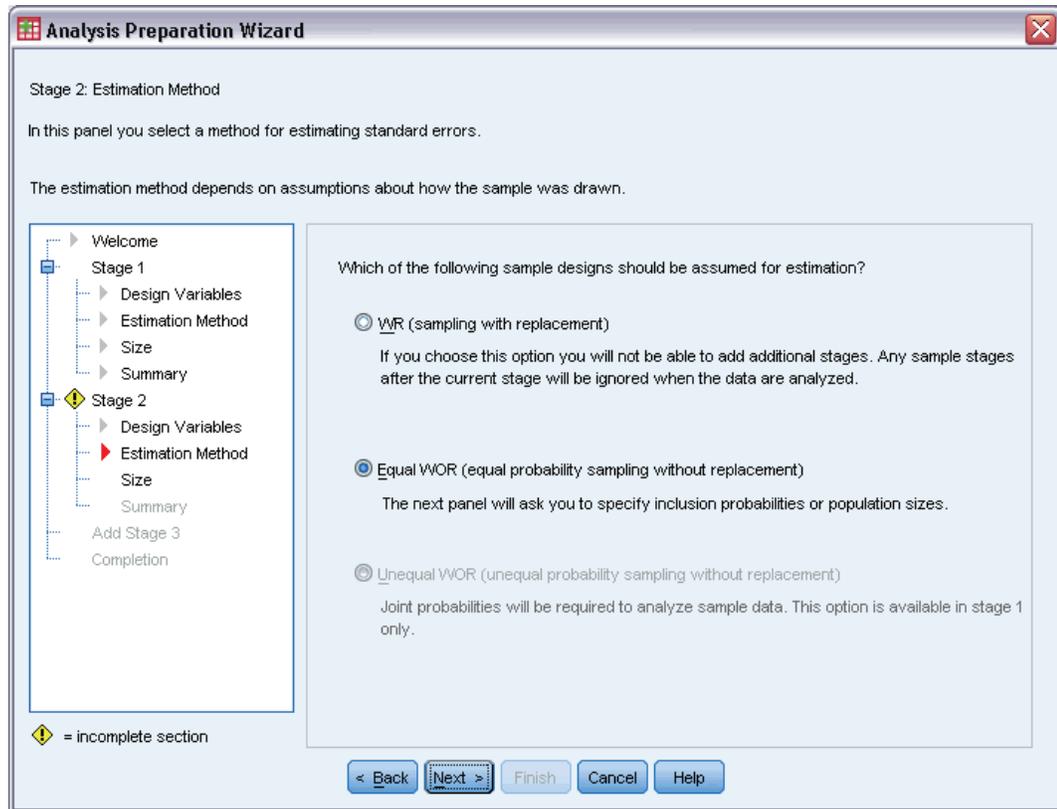
- ▶ Select Read values from variable and select *inclprob\_s1* as the variable containing the first-stage inclusion probabilities.
- ▶ Click Next.

Figure 14-11  
Analysis Preparation Wizard, Plan Summary step (stage 1)



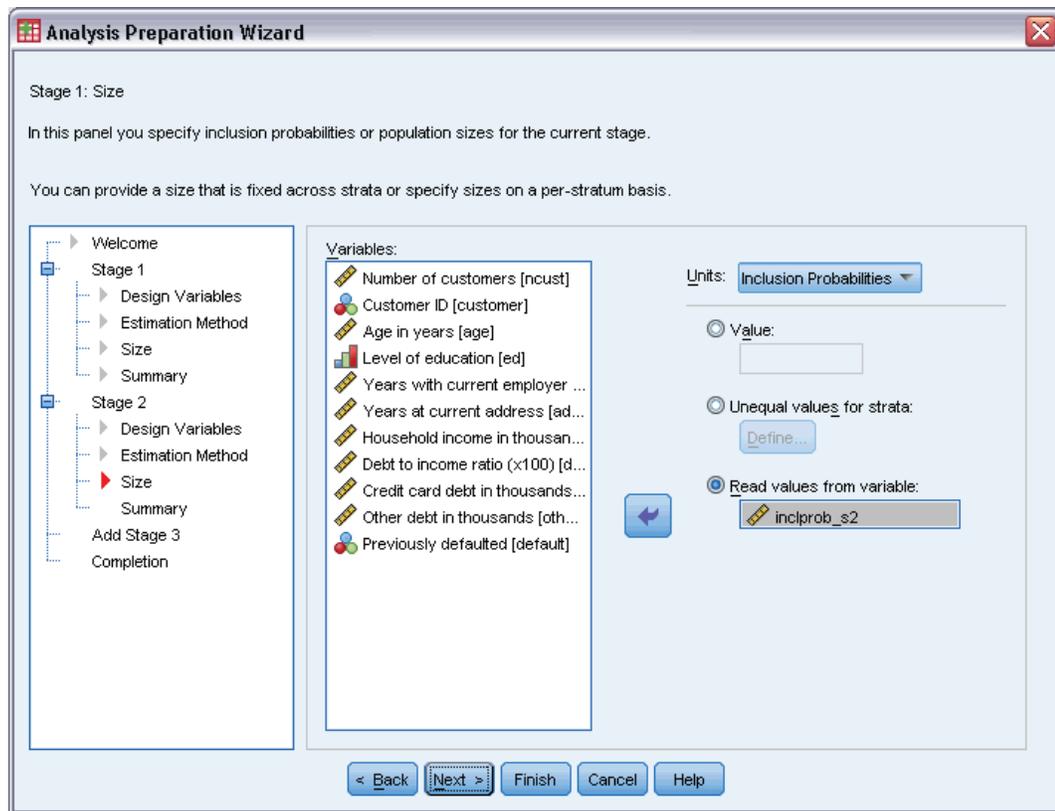
- ▶ Select Yes, add stage 2 now.
- ▶ Click Next, and then click Next in the Design Variables step.

Figure 14-12  
Analysis Preparation Wizard, Estimation Method step (stage 2)



- ▶ Select Equal WOR as the second-stage estimation method.
- ▶ Click Next.

Figure 14-13  
Analysis Preparation Wizard, Size step (stage 2)



- ▶ Select Read values from variable and select *inclprob\_s2* as the variable containing the second-stage inclusion probabilities.
- ▶ Click Finish.

## Summary

Figure 14-14  
Summary table

			Stage 1	Stage 2
Design Variables	Cluster	1	Branch	
Analysis Information	Estimator Assumption		Equal probability sampling without replacement	Equal probability sampling without replacement
	Inclusion Probability		Obtained from variable inclprob_s1	Obtained from variable inclprob_s2

Plan File: c:\bankloan.csaplan  
Weight Variable: finalweight  
SRS Estimator: Sampling without replacement

The summary table reviews your analysis plan. The plan consists of two stages with a design of one cluster variable. Equal probability without replacement (WOR) estimation is used, and the plan is saved to *c:\bankloan.csaplan*. You can now use this plan file to process *bankloan\_noweights.sav* (with the inclusion probabilities and sampling weights you've computed) with Complex Samples analysis procedures.

## Related Procedures

The Complex Samples Analysis Preparation Wizard procedure is a useful tool for readying a sample for analysis when you do not have access to the sampling plan file.

- To create a sampling plan file and draw a sample, use the [Sampling Wizard](#).

# ***Complex Samples Frequencies***

The Complex Samples Frequencies procedure produces frequency tables for selected variables and displays univariate statistics. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

## ***Using Complex Samples Frequencies to Analyze Nutritional Supplement Usage***

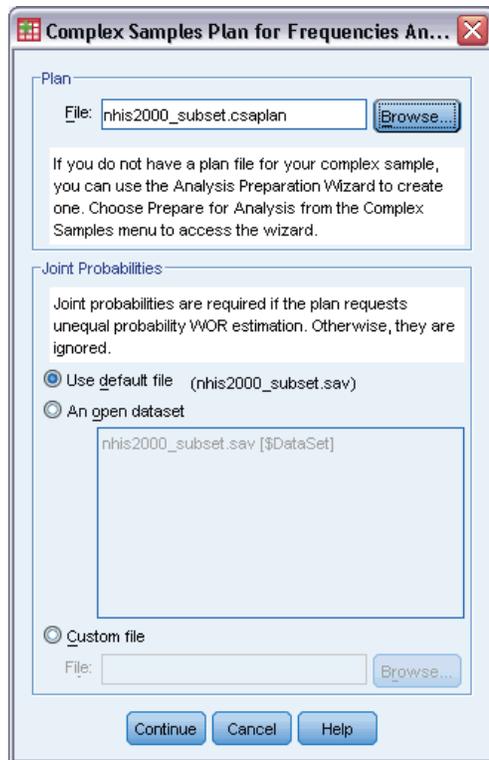
A researcher wants to study the use of nutritional supplements among U.S. citizens, using the results of the National Health Interview Survey (NHIS) and a previously created analysis plan. For more information, see the topic [Using the Complex Samples Analysis Preparation Wizard to Ready NHIS Public Data in Chapter 14 on p. 140](#).

A subset of the 2000 survey is collected in *nhis2000\_subset.sav*. The analysis plan is stored in *nhis2000\_subset.csaplan*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use Complex Samples Frequencies to produce statistics for nutritional supplement usage.

### ***Running the Analysis***

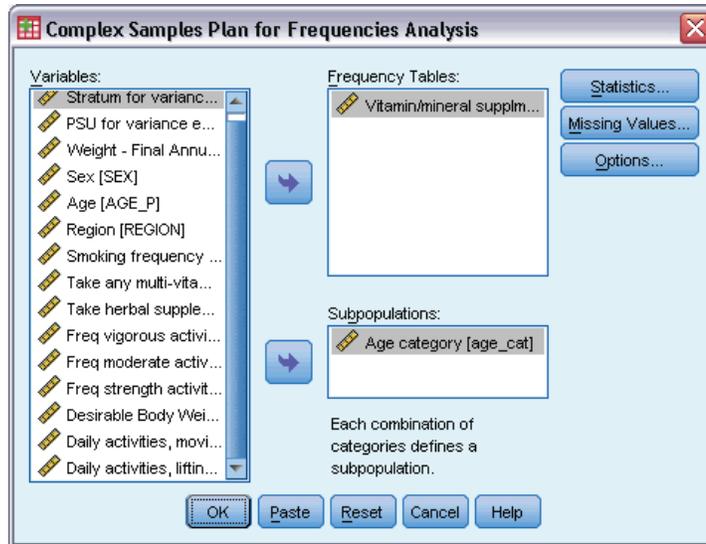
- ▶ To run a Complex Samples Frequencies analysis, from the menus choose:  
Analyze > Complex Samples > Frequencies...

Figure 15-1  
Complex Samples Plan dialog box



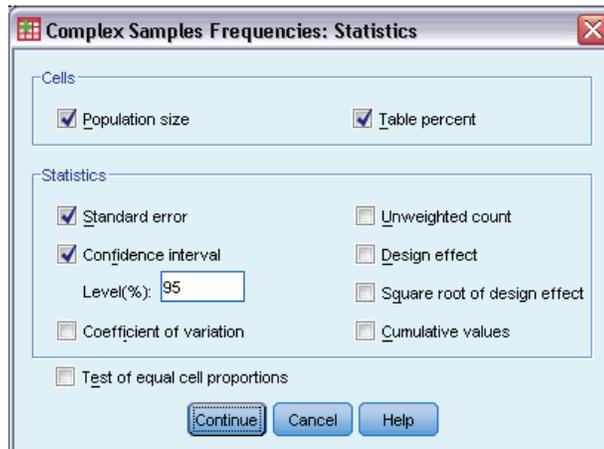
- ▶ Browse to and select *nhis2000\_subset.csaplan*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#).
- ▶ Click Continue.

Figure 15-2  
Frequencies dialog box



- ▶ Select *Vitamin/mineral supplmnts-past 12 m* as a frequency variable.
- ▶ Select *Age category* as a subpopulation variable.
- ▶ Click Statistics.

Figure 15-3  
Frequencies Statistics dialog box



- ▶ Select Table percent in the Cells group.
- ▶ Select Confidence interval in the Statistics group.
- ▶ Click Continue.
- ▶ Click OK in the Frequencies dialog box.

## Frequency Table

Figure 15-4  
Frequency table for variable/situation

		Estimate	Standard Error	95% Confidence Interval	
				Lower	Upper
Population Size	Yes	102767095	1185126.709	100435967	105098223
	No	90794234	1094401.949	88641560	92946908
	Total	193561329	1789098.713	190042196	197080462
% of Total	Yes	53.1%	.4%	52.4%	53.8%
	No	46.9%	.4%	46.2%	47.6%
	Total	100.0%	.0%	100.0%	100.0%

Each selected statistic is computed for each selected cell measure. The first column contains estimates of the number and percentage of the population that do or do not take vitamin/mineral supplements. The confidence intervals are non-overlapping; thus, you can conclude that, overall, more Americans take vitamin/mineral supplements than not.

## Frequency by Subpopulation

Figure 15-5  
Frequency table by subpopulation

Age category			Estimate	Standard Error	95% Confidence Interval	
					Lower	Upper
18-24	Population Size	Yes	10018312	350602.352	9328681.9	10707942
		No	15472368	499182.391	14490483	16454253
		Total	25490680	680732.812	24151688	26829672
	% of Total	Yes	39.3%	1.0%	37.4%	41.2%
		No	60.7%	1.0%	58.8%	62.6%
		Total	100.0%	.0%	100.0%	100.0%
25-44	Population Size	Yes	39163840	660855.719	37863946	40463734
		No	39503150	645934.187	38232606	40773694
		Total	78666990	961114.325	76776491	80557489
	% of Total	Yes	49.8%	.6%	48.7%	50.9%
		No	50.2%	.6%	49.1%	51.3%
		Total	100.0%	.0%	100.0%	100.0%
45-64	Population Size	Yes	34154952	598603.728	32977507	35332397
		No	24005512	497723.833	23026496	24984528
		Total	58160464	814680.415	56557999	59762929
	% of Total	Yes	58.7%	.6%	57.5%	60.0%
		No	41.3%	.6%	40.0%	42.5%
		Total	100.0%	.0%	100.0%	100.0%
65+	Population Size	Yes	19429991	439459.793	18565580	20294402
		No	11813204	314238.078	11195102	12431306
		Total	31243195	587623.439	30087348	32399042
	% of Total	Yes	62.2%	.7%	60.7%	63.6%
		No	37.8%	.7%	36.4%	39.3%
		Total	100.0%	.0%	100.0%	100.0%

When computing statistics by subpopulation, each selected statistic is computed for each selected cell measure by value of *Age category*. The first column contains estimates of the number and percentage of the population of each category that do or do not take vitamin/mineral supplements. The confidence intervals for the table percentages are all non-overlapping; thus, you can conclude that the use of vitamin/mineral supplements increases with age.

## **Summary**

Using the Complex Samples Frequencies procedure, you have obtained statistics for the use of nutritional supplements among U.S. citizens.

- Overall, more Americans take vitamin/mineral supplements than not.
- When broken down by age category, greater proportions of Americans take vitamin/mineral supplements with increasing age.

## **Related Procedures**

The Complex Samples Frequencies procedure is a useful tool for obtaining univariate descriptive statistics of categorical variables for observations obtained via a complex sampling design.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to set analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples Crosstabs](#) procedure provides descriptive statistics for the crosstabulation of categorical variables.
- The [Complex Samples Descriptives](#) procedure provides univariate descriptive statistics for scale variables.

# ***Complex Samples Descriptives***

The Complex Samples Descriptives procedure displays univariate summary statistics for several variables. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

## ***Using Complex Samples Descriptives to Analyze Activity Levels***

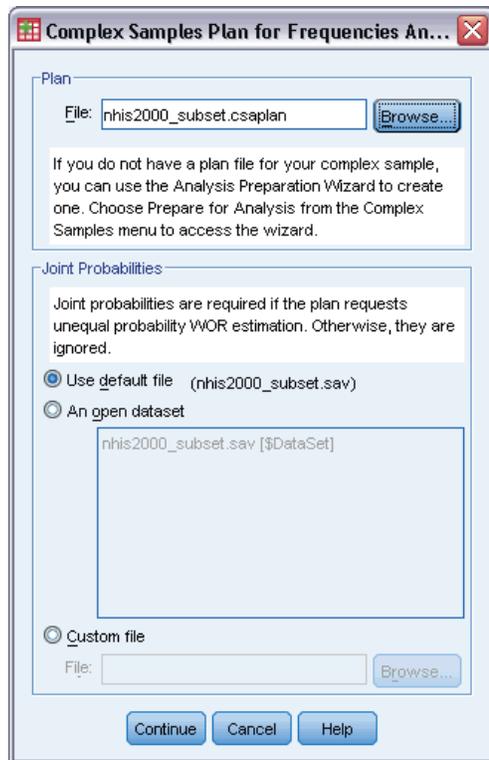
A researcher wants to study the activity levels of U.S. citizens, using the results of the National Health Interview Survey (NHIS) and a previously created analysis plan. [For more information, see the topic Using the Complex Samples Analysis Preparation Wizard to Ready NHIS Public Data in Chapter 14 on p. 140.](#)

A subset of the 2000 survey is collected in *nhis2000\_subset.sav*. The analysis plan is stored in *nhis2000\_subset.csaplan*. [For more information, see the topic Sample Files in Appendix A in IBM SPSS Complex Samples 19.](#) Use Complex Samples Descriptives to produce univariate descriptive statistics for activity levels.

### ***Running the Analysis***

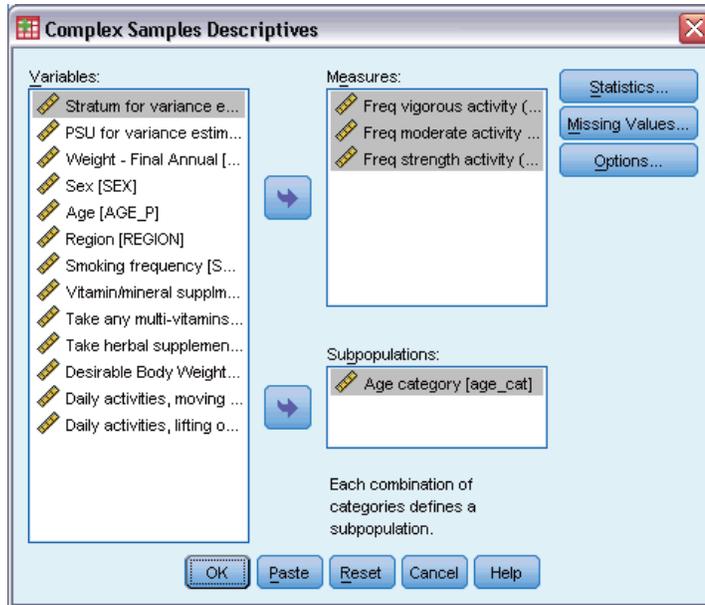
- ▶ To run a Complex Samples Descriptives analysis, from the menus choose:  
Analyze > Complex Samples > Descriptives...

Figure 16-1  
Complex Samples Plan dialog box



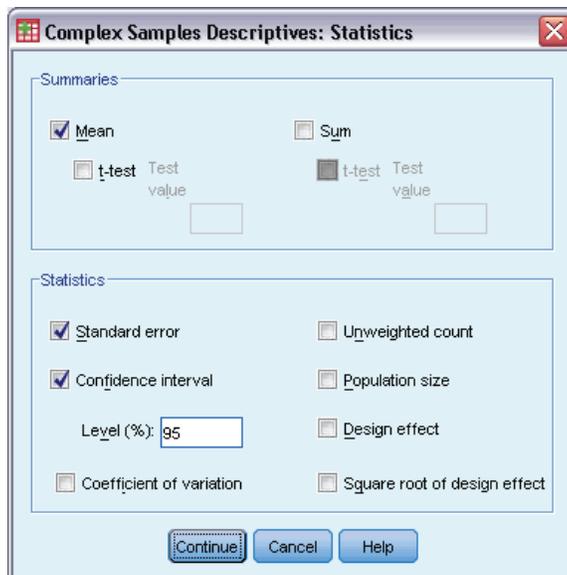
- ▶ Browse to and select *nhis2000\_subset.csaplan*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#).
- ▶ Click Continue.

Figure 16-2  
Descriptives dialog box



- ▶ Select *Freq vigorous activity (times per wk)* through *Freq strength activity (times per wk)* as measure variables.
- ▶ Select *Age category* as a subpopulation variable.
- ▶ Click Statistics.

Figure 16-3  
Descriptives Statistics dialog box



- ▶ Select Confidence interval in the Statistics group.

- ▶ Click Continue.
- ▶ Click OK in the Complex Samples Descriptives dialog box.

## Univariate Statistics

Figure 16-4  
Univariate statistics

		Estimate	Standard Error	95% Confidence Interval	
Mean				Lower	Upper
	Freq vigorous activity (times per wk)	3.73	.033	3.66	3.79
	Freq moderate activity (times per wk)	4.90	.041	4.82	4.98
	Freq strength activity (times per wk)	3.52	.042	3.43	3.60

Each selected statistic is computed for each measure variable. The first column contains estimates of the average number of times per week that a person engages in a particular type of activity. The confidence intervals for the means are non-overlapping. Thus, you can conclude that, overall, Americans engage in a strength activity less often than vigorous activity, and they engage in vigorous activity less often than moderate activity.

## Univariate Statistics by Subpopulation

Figure 16-5  
Univariate statistics by subpopulation

Age category			Estimate	Standard Error	95% Confidence Interval	
	Mean				Lower	Upper
18-24	Mean	Freq vigorous activity (times per wk)	3.92	.087	3.75	4.09
		Freq moderate activity (times per wk)	5.18	.137	4.91	5.45
		Freq strength activity (times per wk)	3.45	.085	3.28	3.62
25-44	Mean	Freq vigorous activity (times per wk)	3.55	.048	3.46	3.65
		Freq moderate activity (times per wk)	4.73	.056	4.62	4.84
		Freq strength activity (times per wk)	3.28	.052	3.18	3.38
45-64	Mean	Freq vigorous activity (times per wk)	3.79	.063	3.66	3.91
		Freq moderate activity (times per wk)	4.88	.070	4.74	5.02
		Freq strength activity (times per wk)	3.65	.092	3.47	3.84
65+	Mean	Freq vigorous activity (times per wk)	4.18	.111	3.96	4.39
		Freq moderate activity (times per wk)	5.22	.084	5.06	5.39
		Freq strength activity (times per wk)	4.66	.155	4.36	4.97

Each selected statistic is computed for each measure variable by values of *Age category*. The first column contains estimates of the average number of times per week that people of each category engage in a particular type of activity. The confidence intervals for the means allow you to make some interesting conclusions.

- In terms of vigorous and moderate activities, 25–44-year-olds are less active than those 18–24 and 45–64, and 45–64-year-olds are less active than those 65 or older.
- In terms of strength activity, 25–44-year-olds are less active than those 45–64, and 18–24 and 45–64-year-olds are less active than those 65 or older.

## **Summary**

Using the Complex Samples Descriptives procedure, you have obtained statistics for the activity levels of U.S. citizens.

- Overall, Americans spend varying amounts of time at different types of activities.
- When broken down by age, it roughly appears that post-collegiate Americans are initially less active than they were while in school but become more conscientious about exercising as they age.

## **Related Procedures**

The Complex Samples Descriptives procedure is a useful tool for obtaining univariate descriptive statistics of scale measures for observations obtained via a complex sampling design.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to set analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples Ratios](#) procedure provides descriptive statistics for ratios of scale measures.
- The [Complex Samples Frequencies](#) procedure provides univariate descriptive statistics of categorical variables.

---

# ***Complex Samples Crosstabs***

The Complex Samples Crosstabs procedure produces crosstabulation tables for pairs of selected variables and displays two-way statistics. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

## ***Using Complex Samples Crosstabs to Measure the Relative Risk of an Event***

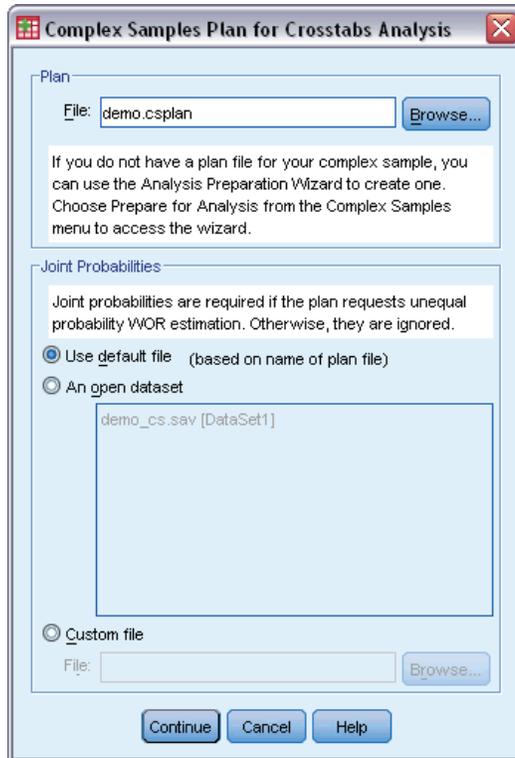
A company that sells magazine subscriptions traditionally sends monthly mailings to a purchased database of names. The response rate is typically low, so you need to find a way to better target prospective customers. One suggestion is to focus mailings on people with newspaper subscriptions, on the assumption that people who read newspapers are more likely to subscribe to magazines.

Use the Complex Samples Crosstabs procedure to test this theory by constructing a two-by-two table of *Newspaper subscription* by *Response* and computing the relative risk that a person with a newspaper subscription will respond to the mailing. This information is collected in *demo\_cs.sav* and should be analyzed using the sampling plan file *demo.csplan*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#).

### ***Running the Analysis***

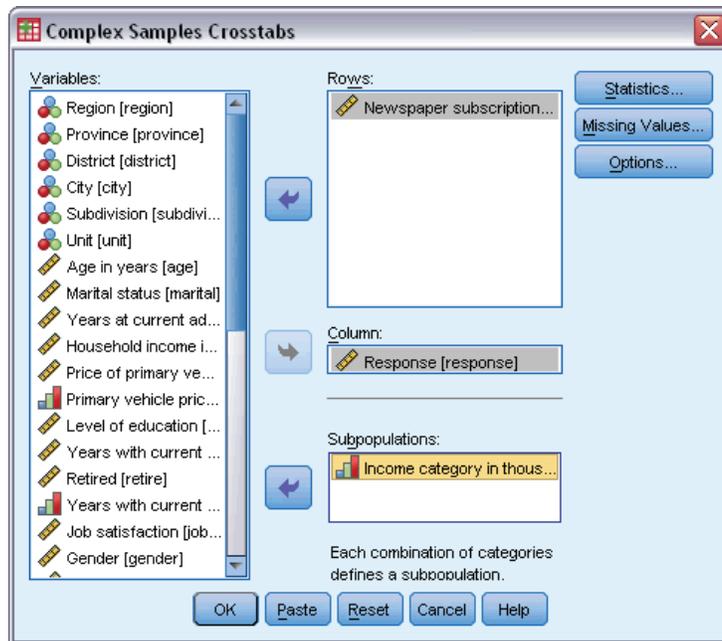
- ▶ To run a Complex Samples Crosstabs analysis, from the menus choose:  
Analyze > Complex Samples > Crosstabs...

Figure 17-1  
Complex Samples Plan dialog box



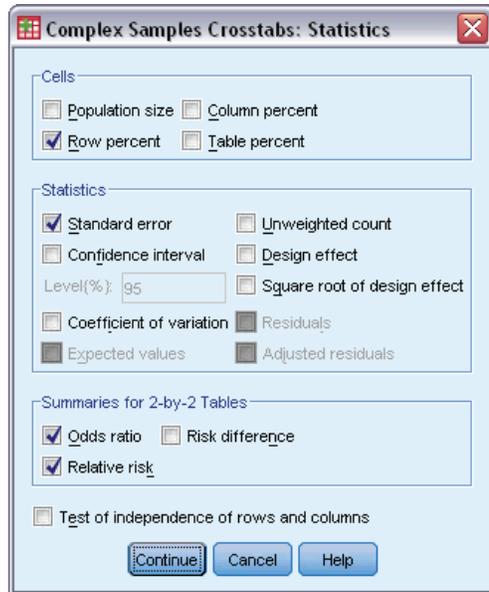
- ▶ Browse to and select *demo.csplan*. For more information, see the topic Sample Files in Appendix A in *IBM SPSS Complex Samples 19*.
- ▶ Click Continue.

Figure 17-2  
Crosstabs dialog box



- ▶ Select *Newspaper subscription* as a row variable.
- ▶ Select *Response* as a column variable.
- ▶ There is also some interest in seeing the results broken down by income categories, so select *Income category in thousands* as a subpopulation variable.
- ▶ Click Statistics.

Figure 17-3  
Crosstabs Statistics dialog box



- ▶ Deselect Population size and select Row percent in the Cells group.
- ▶ Select Odds ratio and Relative risk in the Summaries for 2-by-2 Tables group.
- ▶ Click Continue.
- ▶ Click OK in the Complex Samples Crosstabs dialog box.

These selections produce a crosstabulation table and risk estimate for *Newspaper subscription* by *Response*. Separate tables with results split by *Income category in thousands* are also created.

## Crosstabulation

Figure 17-4  
Crosstabulation for newspaper subscription by response

Newspaper subscription			Response		
			Yes	No	Total
Yes	% within Newspaper subscription	Estimate	17.2%	82.8%	100.0%
		Standard Error	1.0%	1.0%	.0%
No	% within Newspaper subscription	Estimate	10.3%	89.7%	100.0%
		Standard Error	.7%	.7%	.0%
Total	% within Newspaper subscription	Estimate	12.8%	87.2%	100.0%
		Standard Error	.7%	.7%	.0%

The crosstabulation shows that, overall, few people responded to the mailing. However, a greater proportion of newspaper subscribers responded.

## Risk Estimate

Figure 17-5  
Risk estimate for newspaper subscription by response

		Estimate
Newspaper subscription * Response	Odds Ratio	1.812
	Relative Risk	1.673
		.923

Statistics are computed only for 2-by-2 tables with all cells observed.

The relative risk is a ratio of event probabilities. The relative risk of a response to the mailing is the ratio of the probability that a newspaper subscriber responds to the probability that a nonsubscriber responds. Thus, the estimate of the relative risk is simply  $17.2\%/10.3\% = 1.673$ . Likewise, the relative risk of nonresponse is the ratio of the probability that a subscriber does not respond to the probability that a nonsubscriber does not respond. Your estimate of this relative risk is 0.923. Given these results, you can estimate that a newspaper subscriber is 1.673 times as likely to respond to the mailing as a nonsubscriber, or 0.923 times as likely as a nonsubscriber not to respond.

The odds ratio is a ratio of event odds. The odds of an event is the ratio of the probability that the event occurs to the probability that the event does not occur. Thus, the estimate of the odds that a newspaper subscriber responds to the mailing is  $17.2\%/82.8\% = 0.208$ . Likewise, the estimate of the odds that a nonsubscriber responds is  $10.3\%/89.7\% = 0.115$ . The estimate of the odds ratio is therefore  $0.208/0.115 = 1.812$  (note there is some rounding error in the intervening steps). The odds ratio is also the ratio of the relative risk of responding to the relative risk of not responding, or  $1.673/0.923 = 1.812$ .

### Odds Ratio versus Relative Risk

Since it is a ratio of ratios, the odds ratio is very difficult to interpret. The relative risk is easier to interpret, so the odds ratio alone is not very helpful. However, there are certain commonly occurring situations in which the estimate of the relative risk is not very good, and the odds ratio can be used to approximate the relative risk of the event of interest. The odds ratio should be used as an approximation of the relative risk of the event of interest when both of the following conditions are met:

- The probability of the event of interest is small ( $< 0.1$ ). This condition guarantees that the odds ratio will make a good approximation to the relative risk. In this example, the event of interest is a response to the mailing.
- The design of the study is case control. This condition signals that the usual estimate of the relative risk will likely not be good. A case-control study is retrospective, most often used when the event of interest is unlikely or when the design of a prospective experiment is impractical or unethical.

Neither condition is met in this example, since the overall proportion of respondents was 12.8% and the design of the study was not case control, so it's safer to report 1.673 as the relative risk, rather than the value of the odds ratio.

## Risk Estimate by Subpopulation

Figure 17-6

Risk estimate for newspaper subscription by response, controlling for income category

Income category			Estimate
Under \$25	Newspaper subscription * Response	Odds Ratio	2.712
		Relative Risk	2.241
		For cohort Response = Yes For cohort Response = No	.826
\$25 - \$49	Newspaper subscription * Response	Odds Ratio	1.794
		Relative Risk	1.645
		For cohort Response = Yes For cohort Response = No	.917
\$50 - \$74	Newspaper subscription * Response	Odds Ratio	1.168
		Relative Risk	1.152
		For cohort Response = Yes For cohort Response = No	.986
\$75+	Newspaper subscription * Response	Odds Ratio	1.242
		Relative Risk	1.227
		For cohort Response = Yes For cohort Response = No	.988

Statistics are computed only for 2-by-2 tables with all cells observed.

Relative risk estimates are computed separately for each income category. Note that the relative risk of a positive response for newspaper subscribers appears to gradually decrease with increasing income, which indicates that you may be able to further target the mailings.

## Summary

Using Complex Samples Crosstabs risk estimates, you found that you can increase your response rate to direct mailings by targeting newspaper subscribers. Further, you found some evidence that the risk estimates may not be constant across *Income category*, so you may be able to increase your response rate even more by targeting lower-income newspaper subscribers.

## Related Procedures

The Complex Samples Crosstabs procedure is a useful tool for obtaining descriptive statistics of the crosstabulation of categorical variables for observations obtained via a complex sampling design.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to set analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples Frequencies](#) procedure provides univariate descriptive statistics of categorical variables.

# ***Complex Samples Ratios***

The Complex Samples Ratios procedure displays univariate summary statistics for ratios of variables. Optionally, you can request statistics by subgroups, defined by one or more categorical variables.

## ***Using Complex Samples Ratios to Aid Property Value Assessment***

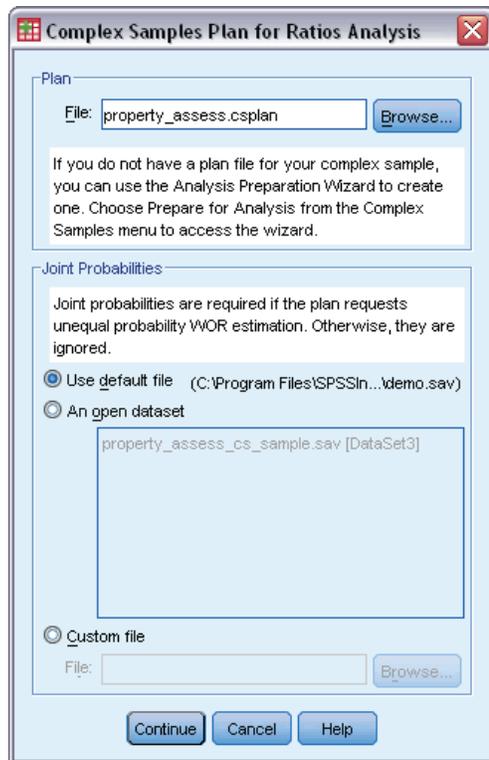
A state agency is charged with ensuring that property taxes are fairly assessed from county to county. Taxes are based on the appraised value of the property, so the agency wants to track property values across counties to be sure that each county's records are equally up-to-date. Since resources for obtaining current appraisals are limited, the agency chose to employ complex sampling methodology to select properties.

The sample of properties selected and their current appraisal information is collected in *property\_assess\_cs\_sample.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use Complex Samples Ratios to assess the change in property values across the five counties since the last appraisal.

### ***Running the Analysis***

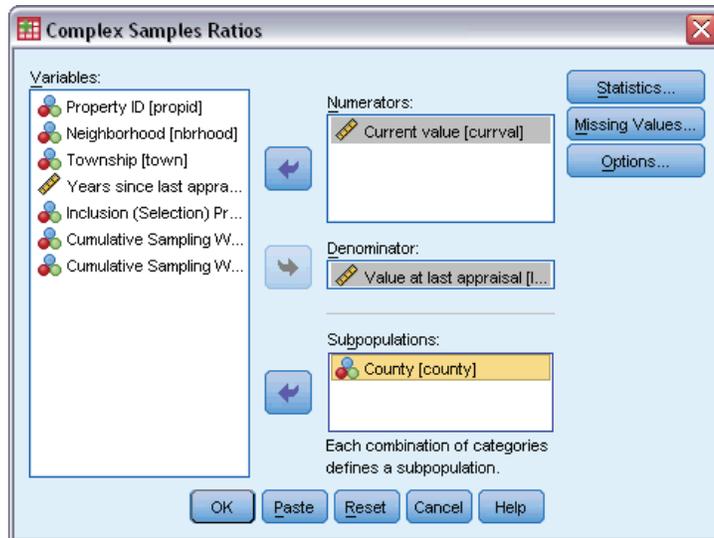
- ▶ To run a Complex Samples Ratios analysis, from the menus choose:  
Analyze > Complex Samples > Ratios...

Figure 18-1  
Complex Samples Plan dialog box



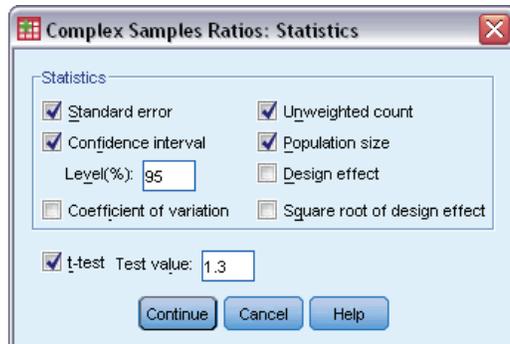
- ▶ Browse to and select *property\_assess.csplan*. For more information, see the topic *Sample Files* in *Appendix A* in *IBM SPSS Complex Samples 19*.
- ▶ Click *Continue*.

Figure 18-2  
Ratios dialog box



- ▶ Select *Current value* as a numerator variable.
- ▶ Select *Value at last appraisal* as the denominator variable.
- ▶ Select *County* as a subpopulation variable.
- ▶ Click Statistics.

Figure 18-3  
Ratios Statistics dialog box



- ▶ Select Confidence interval, Unweighted count, and Population size in the Statistics group.
- ▶ Select t-test and enter 1.3 as the test value.
- ▶ Click Continue.
- ▶ Click OK in the Complex Samples Ratios dialog box.

## Ratios

Figure 18-4  
Ratios table

County	Numerator	Denominator	Ratio Estimate	Standard Error	95% Confidence Interval		Te
					Lower	Upper	
Eastern	Current value	Value at last appraisal	1.381	.068	1.236	1.525	
Central	Current value	Value at last appraisal	1.364	.064	1.227	1.502	
Western	Current value	Value at last appraisal	1.524	.053	1.410	1.638	
Northern	Current value	Value at last appraisal	1.277	.032	1.208	1.346	
Southern	Current value	Value at last appraisal	1.195	.029	1.134	1.256	

The default display of the table is very wide, so you will need to pivot it for a better view.

### Pivoting the Ratios Table

- ▶ Double-click the table to activate it.
- ▶ From the Viewer menus choose:  
Pivot > Pivoting Trays
- ▶ Drag *Numerator* and then *Denominator* from the row to the layer.
- ▶ Drag *County* from the row to the column.
- ▶ Drag *Statistics* from the column to the row.
- ▶ Close the pivoting trays window.

### Pivoted Ratios Table

Figure 18-5  
Pivoted ratios table

Numerator: Current value  
Denominator: Value at last appraisal

	County					
	Eastern	Central	Western	Northern	Southern	
Ratio Estimate	1.381	1.364	1.524	1.277	1.195	
Standard Error	.068	.064	.053	.032	.029	
95% Confidence Interval	Lower	1.236	1.227	1.410	1.208	1.134
	Upper	1.525	1.502	1.638	1.346	1.256
Hypothesis Test	Test Value	1.3	1.3	1.3	1.3	1.3
	t	1.191	.997	4.201	-.702	-3.646
	df	15	15	15	15	15
	Sig.	.252	.334	.001	.493	.002
Unweighted Count	168	179	202	205	220	

The ratios table is now pivoted so that statistics are easier to compare across counties.

- The ratio estimates range from a low of 1.195 in the Southern county to a high of 1.524 in the Western county.
- There is also quite a bit of variability in the standard errors, which range from a low of 0.029 in the Southern county to 0.068 in the Eastern county.

- Some of the confidence intervals do not overlap; thus, you can conclude that the ratios for the Western county are higher than the ratios for the Northern and Southern counties.
- Finally, as a more objective measure, note that the significance values of the  $t$  tests for the Western and Southern counties are less than 0.05. Thus, you can conclude that the ratio for the Western county is greater than 1.3 and the ratio for the Southern county is less than 1.3.

## Summary

Using the Complex Samples Ratios procedure, you have obtained various statistics for the ratios of *Current value* to *Value at last appraisal*. The results suggest that there may be certain inequities in the assessment of property taxes from county to county, namely:

- The ratios for the Western county are high, indicating that their records are not as up-to-date as other counties with respect to the appreciation of property values. Property taxes are probably too low in this county.
- The ratios for the Southern county are low, indicating that their records are more up-to-date than the other counties with respect to the appreciation of property values. Property taxes are probably too high in this county.
- The ratios for the Southern county are lower than those of the Western county but are still within the objective goal of 1.3.

Resources used to track property values in the Southern county will be reassigned to the Western county to bring these counties' ratios in line with the others and with the goal of 1.3.

## Related Procedures

The Complex Samples Ratios procedure is a useful tool for obtaining univariate descriptive statistics of the ratio of scale measures for observations obtained via a complex sampling design.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to set analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples Descriptives](#) procedure provides descriptive statistics for scale variables.

# ***Complex Samples General Linear Model***

The Complex Samples General Linear Model (CSGLM) procedure performs linear regression analysis, as well as analysis of variance and covariance, for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

## ***Using Complex Samples General Linear Model to Fit a Two-Factor ANOVA***

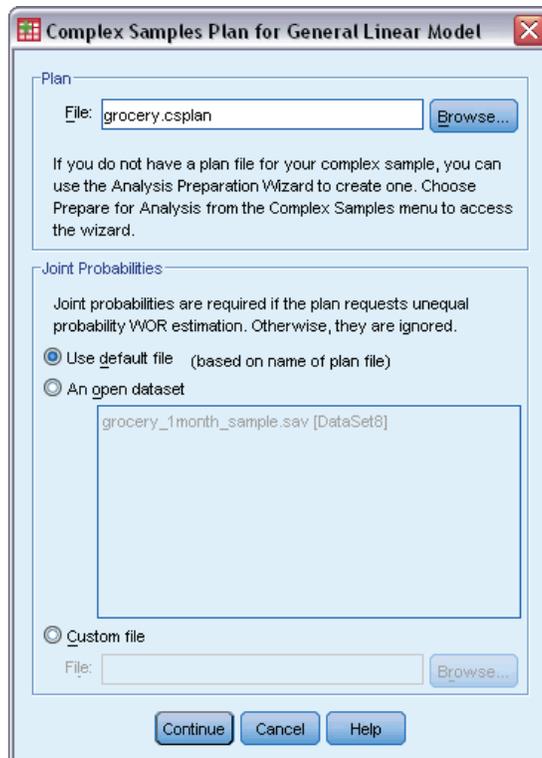
A grocery store chain surveyed a set of customers concerning their purchasing habits, according to a complex design. Given the survey results and how much each customer spent in the previous month, the store wants to see if the frequency with which customers shop is related to the amount they spend in a month, controlling for the gender of the customer and incorporating the sampling design.

This information is collected in *grocery\_1month\_sample.sav*. For more information, see the [topic Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use the Complex Samples General Linear Model procedure to perform a two-factor (or two-way) ANOVA on the amounts spent.

### ***Running the Analysis***

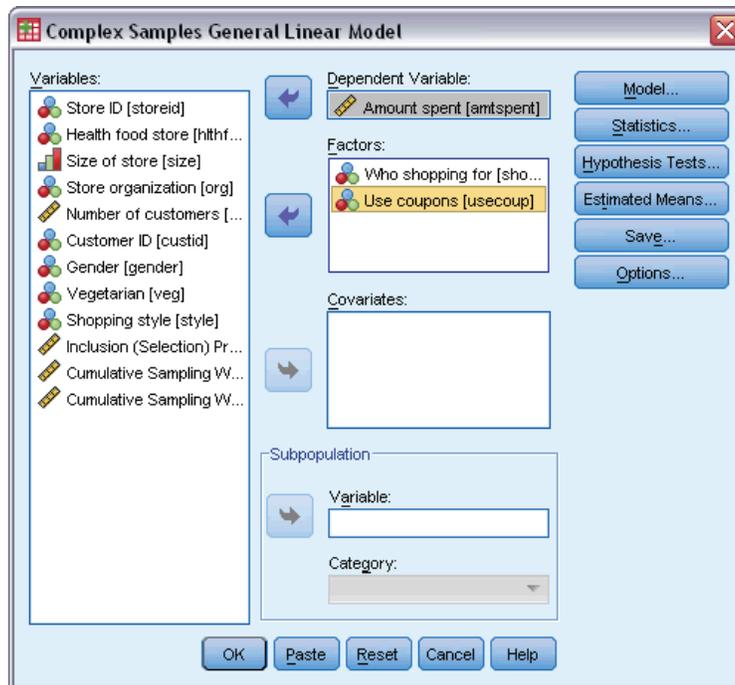
- ▶ To run a Complex Samples General Linear Model analysis, from the menus choose:  
Analyze > Complex Samples > General Linear Model...

Figure 19-1  
Complex Samples Plan dialog box



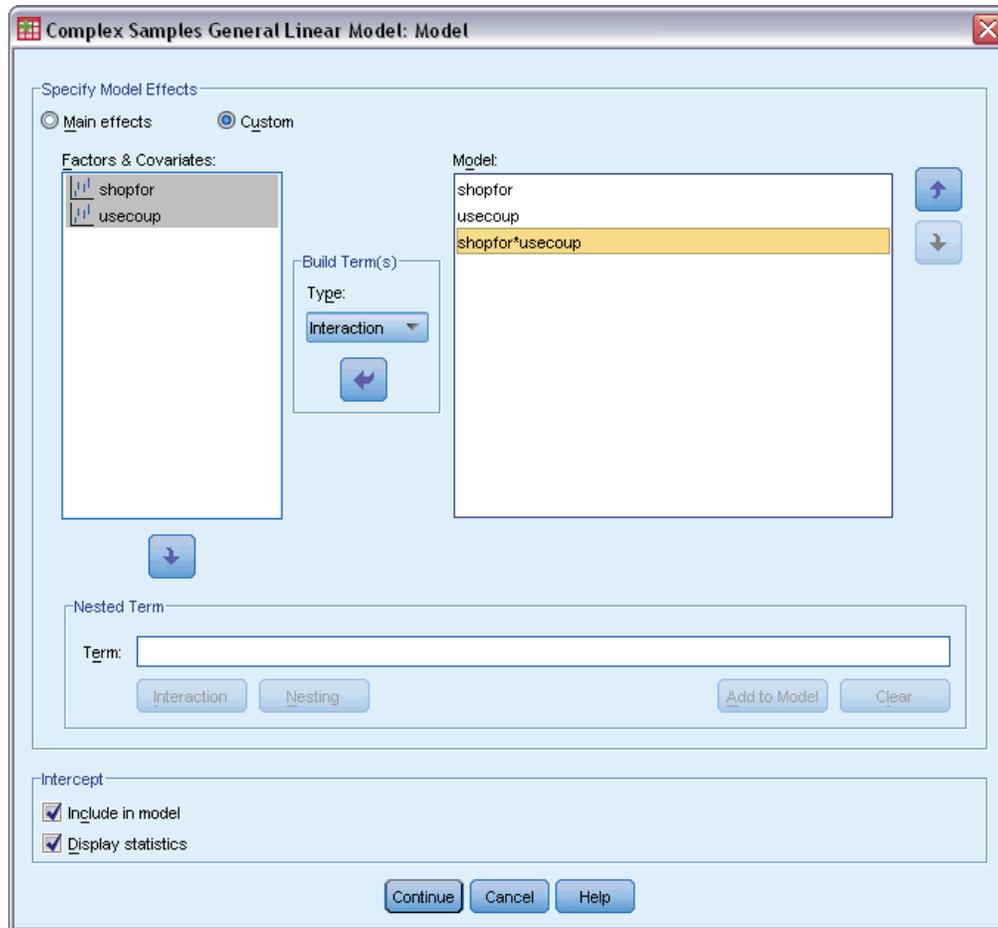
- ▶ Browse to and select *grocery.csplan*. For more information, see the topic *Sample Files* in *Appendix A* in *IBM SPSS Complex Samples 19*.
- ▶ Click *Continue*.

Figure 19-2  
General Linear Model dialog box



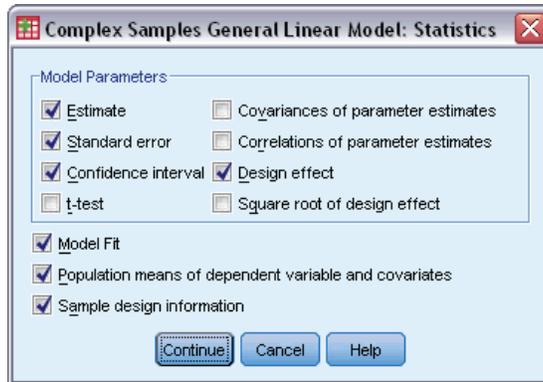
- ▶ Select *Amount spent* as the dependent variable.
- ▶ Select *Who shopping for* and *Use coupons* as factors.
- ▶ Click Model.

Figure 19-3  
Model dialog box



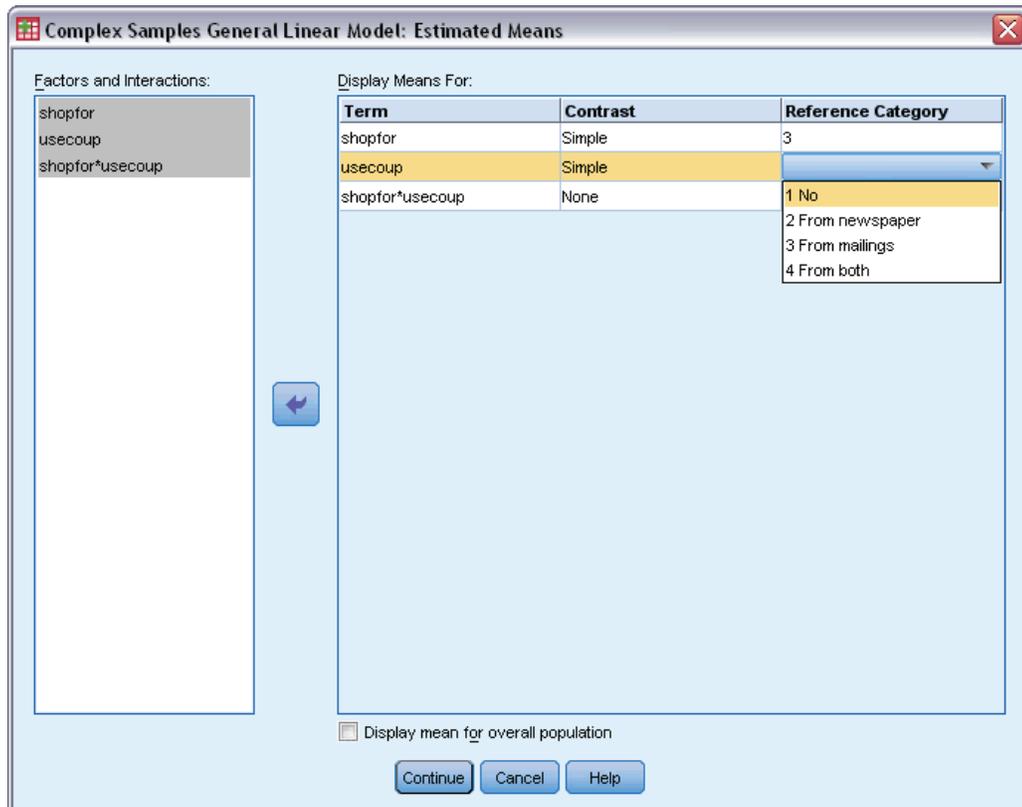
- ▶ Choose to build a Custom model.
- ▶ Select Main effects as the type of term to build and select *shopfor* and *usecoup* as model terms.
- ▶ Select Interaction as the type of term to build and add the *shopfor\*usecoup* interaction as a model term.
- ▶ Click Continue.
- ▶ Click Statistics in the General Linear Model dialog box.

Figure 19-4  
General Linear Model Statistics dialog box



- ▶ Select Estimate, Standard error, Confidence interval, and Design effect in the Model Parameters group.
- ▶ Click Continue.
- ▶ Click Estimated Means in the General Linear Model dialog box.

Figure 19-5  
General Linear Model Estimated Means dialog box



- ▶ Choose to display means for *shopfor*, *usecoup*, and the *shopfor\*usecoup* interaction.

- ▶ Select a Simple contrast and 3 Self and family as the reference category for *shopfor*. Note that, once selected, the category appears as “3” in the dialog box.
- ▶ Select a Simple contrast and 1 No as the reference category for *usecoup*.
- ▶ Click Continue.
- ▶ Click OK in the General Linear Model dialog box.

## Model Summary

Figure 19-6  
*R-square statistic*

R Square	.601
----------	------

a. Model: Amount spent = (Intercept) +  
shopfor + usecoup + shopfor \* usecoup

*R*-square, the coefficient of determination, is a measure of the strength of the model fit. It shows that about 60% of the variation in *Amount spent* is explained by the model, which gives you good explanatory ability. You may still want to add other predictors to the model to further improve the fit.

## Tests of Model Effects

Figure 19-7  
*Tests of between-subjects effects*

Source	df1	df2	Wald F	Sig.
(Corrected Model)	11.000	3.000	127.231	.001
(Intercept)	1.000	13.000	6321.597	.000
shopfor	2.000	12.000	643.593	.000
usecoup	3.000	11.000	87.453	.000
shopfor * usecoup	6.000	8.000	10.688	.002

a. Model: Amount spent = (Intercept) + shopfor + usecoup + shopfor \*  
usecoup

Each term in the model, plus the model as a whole, is tested for whether the value of its effect equals 0. Terms with significance values of less than 0.05 have some discernible effect. Thus, all model terms contribute to the model.

## Parameter Estimates

Figure 19-8  
Parameter estimates

Parameter	Estimate	Std. Error	95% Confidence Interval		Design Effect
			Lower	Upper	
(Intercept)	518.249	11.731	492.905	543.592	1.387
[shopfor=1]	-174.757	10.762	-198.0	-151.51	.950
[shopfor=2]	-129.443	11.455	-154.2	-104.70	.925
[shopfor=3]	.000 <sup>a</sup>	.	.	.	.
[usecoup=1]	-140.838	10.180	-162.8	-118.85	.649
[usecoup=2]	-63.026	13.195	-91.531	-34.520	.940
[usecoup=3]	-31.375	9.726	-52.387	-10.363	.564
[usecoup=4]	.000 <sup>a</sup>	.	.	.	.
[shopfor=1] * [usecoup=1]	41.693	11.170	17.562	65.824	.606
[shopfor=1] * [usecoup=2]	44.505	18.068	5.471	83.539	1.413
[shopfor=1] * [usecoup=3]	9.204	11.057	-14.684	33.092	.594
[shopfor=1] * [usecoup=4]	.000 <sup>a</sup>	.	.	.	.
[shopfor=2] * [usecoup=1]	89.211	10.967	65.518	112.903	.533
[shopfor=2] * [usecoup=2]	54.267	14.949	21.972	86.562	.836
[shopfor=2] * [usecoup=3]	17.884	13.753	-11.828	47.595	.797
[shopfor=2] * [usecoup=4]	.000 <sup>a</sup>	.	.	.	.
[shopfor=3] * [usecoup=1]	.000 <sup>a</sup>	.	.	.	.
[shopfor=3] * [usecoup=2]	.000 <sup>a</sup>	.	.	.	.
[shopfor=3] * [usecoup=3]	.000 <sup>a</sup>	.	.	.	.
[shopfor=3] * [usecoup=4]	.000 <sup>a</sup>	.	.	.	.

a. Set to zero because this parameter is redundant.

b. Model: Amount spent = (Intercept) + shopfor + usecoup + shopfor \* usecoup

The parameter estimates show the effect of each predictor on *Amount spent*. The value of 518.249 for the intercept term indicates that the grocery chain can expect a shopper with a family who uses coupons from the newspaper and targeted mailings to spend \$518.25, on average. You can tell that the intercept is associated with these factor levels because those are the factor levels whose parameters are redundant.

- The *shopfor* coefficients suggest that among customers who use both mailed coupons and newspaper coupons, those without family tend to spend less than those with spouses, who in turn spend less than those with dependents at home. Since the tests of model effects showed that this term contributes to the model, these differences are not due to chance.
- The *usecoup* coefficients suggest that spending among customers with dependents at home decreases with decreased coupon usage. There is a moderate amount of uncertainty in the estimates, but the confidence intervals do not include 0.
- The interaction coefficients suggest that customers who do not use coupons or only clip from the newspaper and do not have dependents tend to spend more than you would otherwise expect. If any portion of an interaction parameter is redundant, the interaction parameter is redundant.
- The deviation in the values of the design effects from 1 indicate that some of the standard errors computed for these parameter estimates are larger than those you would obtain if you assumed that these observations came from a simple random sample, while others are smaller. It is vitally important to incorporate the sampling design information in your analysis because you might otherwise infer, for example, that the *usecoup*=3 coefficient is not different from 0!

The parameter estimates are useful for quantifying the effect of each model term, but the estimated marginal means tables can make it easier to interpret the model results.

## Estimated Marginal Means

Figure 19-9  
Estimated marginal means by levels of *Who shopping for*

Who shopping for	Mean	Std. Error	95% Confidence Interval	
			Lower	Upper
Self	308.5326	3.94286	300.0145	317.0506
Self and spouse	370.3361	4.87908	359.7955	380.8767
Self and family	459.4392	7.19769	443.8895	474.9888

This table displays the model-estimated marginal means and standard errors of *Amount spent* at the factor levels of *Who shopping for*. This table is useful for exploring the differences between the levels of this factor. In this example, a customer who shops for him- or herself is expected to spend about \$308.53, while a customer with a spouse is expected to spend \$370.34, and a customer with dependents will spend \$459.44. To see whether this represents a real difference or is due to chance variation, look at the test results.

Figure 19-10  
Individual test results for estimated marginal means of gender

Who shopping for Simple Contrast <sup>a</sup>	Contrast Estimate	Hypothesized Value	Difference (Estimate - Hypothesized)	Std. Error	df1	df2	Wald F	Sig.
Level Self vs. Level Self and family	-150.907	.000	-150.907	4.903	1.000	13.00	947.41	.000
Level Self and spouse vs. Level Self and family	-89.103	.000	-89.103	5.903	1.000	13.00	227.84	.000

a. Reference Category = Self and family

The individual tests table displays two simple contrasts in spending.

- The contrast estimate is the difference in spending for the listed levels of *Who shopping for*.
- The hypothesized value of 0.00 represents the belief that there is no difference in spending.
- The Wald *F* statistic, with the displayed degrees of freedom, is used to test whether the difference between a contrast estimate and hypothesized value is due to chance variation.
- Since the significance values are less than 0.05, you can conclude that there are differences in spending.

The values of the contrast estimates are different from the parameter estimates. This is because there is an interaction term containing the *Who shopping for* effect. As a result, the parameter estimate for *shopfor=1* is a simple contrast between the levels *Self* and *Self and Family* at the level *From both* of the variable *Use coupons*. The contrast estimate in this table is averaged over the levels of *Use coupons*.

**Figure 19-11**  
Overall test results for estimated marginal means of gender

df1	df2	Wald F	Sig.
2.000	12.000	643.593	.000

The overall test table reports the results of a test of all of the contrasts in the individual test table. Its significance value of less than 0.05 confirms that there is a difference in spending among the levels of *Who shopping for*.

**Figure 19-12**  
Estimated marginal means by levels of shopping style

Use coupons	Mean	Std. Error	95% Confidence Interval	
			Lower	Upper
No	319.6455	6.51429	305.5722	333.7188
From newspaper	386.7469	4.32295	377.4077	396.0861
From mailings	394.5028	5.54218	382.5297	406.4760
From both	416.8486	6.51260	402.7790	430.9182

This table displays the model-estimated marginal means and standard errors of *Amount spent* at the factor levels of *Use coupons*. This table is useful for exploring the differences between the levels of this factor. In this example, a customer who does not use coupons is expected to spend about \$319.65, and those who do use coupons are expected to spend considerably more.

**Figure 19-13**  
Individual test results for estimated marginal means of shopping style

Use coupons Simple Contrast <sup>a</sup>	Contrast Estimate	Hypothesized Value	Difference (Estimate - Hypothesized)	Std. Error	df1	df2	Wald F	Sig.
Level From newspaper vs. Level No	67.101	.000	67.101	6.537	1.000	13.000	105.35	.000
Level From mailings vs. Level No	74.857	.000	74.857	5.875	1.000	13.000	162.33	.000
Level From both vs. Level No	97.203	.000	97.203	5.603	1.000	13.000	300.92	.000

a. Reference Category = No

The individual tests table displays three simple contrasts, comparing the spending of customers who do not use coupons to those who do.

Since the significance values of the tests are less than 0.05, you can conclude that customers who use coupons tend to spend more than those who don't.

**Figure 19-14**  
Overall test results for estimated marginal means of shopping style

df1	df2	Wald F	Sig.
3.000	11.000	87.453	.000

The overall test table reports the results of a test of all the contrasts in the individual test table. Its significance value of less than 0.05 confirms that there is a difference in spending among the levels of *Use coupons*. Note that the overall tests for *Use coupons* and *Who shopping for* are equivalent to the tests of model effects because the hypothesized contrast values are equal to 0.

**Figure 19-15**  
*Estimated marginal means by levels of gender by shopping style*

Who shopping for	Use coupons	Mean	Std. Error	95% Confidence Interval	
				Lower	Upper
Self	No	244.3471	6.00949	231.3644	257.3298
	From newspaper	324.9708	5.94134	312.1353	337.8063
	From mailings	321.3207	4.11028	312.4410	330.2005
	From both	343.4916	6.57845	329.2797	357.7034
Self and spouse	No	337.1783	7.12181	321.7925	352.5640
	From newspaper	380.0468	7.91038	362.9574	397.1361
	From mailings	375.3141	6.22468	361.8665	388.7617
	From both	388.8054	7.12101	373.4214	404.1894
Self and family	No	377.4111	11.58215	352.3894	402.4328
	From newspaper	455.2232	6.14420	441.9494	468.4969
	From mailings	486.8736	10.76529	463.6166	510.1306
	From both	518.2488	11.73120	492.9050	543.5925

This table displays the model-estimated marginal means, standard errors, and confidence intervals of *Amount spent* at the factor combinations of *Who shopping for* and *Use coupons*. This table is useful for exploring the interaction effect between these two factors that was found in the tests of model effects.

## Summary

In this example, the estimated marginal means revealed differences in spending between customers at varying levels of *Who shopping for* and *Use coupons*. The tests of model effects confirmed this, as well as the fact that there appears to be a *Who shopping for\*Use coupons* interaction effect. The model summary table revealed that the present model explains somewhat more than half of the variation in the data and could likely be improved by adding more predictors.

## Related Procedures

The Complex Samples General Linear Model procedure is a useful tool for modeling a scale variable when the cases have been drawn according to a complex sampling scheme.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to specify analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples Logistic Regression](#) procedure allows you to model a categorical response.
- The [Complex Samples Ordinal Regression](#) procedure allows you to model an ordinal response.

# ***Complex Samples Logistic Regression***

The Complex Samples Logistic Regression procedure performs logistic regression analysis on a binary or multinomial dependent variable for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

## ***Using Complex Samples Logistic Regression to Assess Credit Risk***

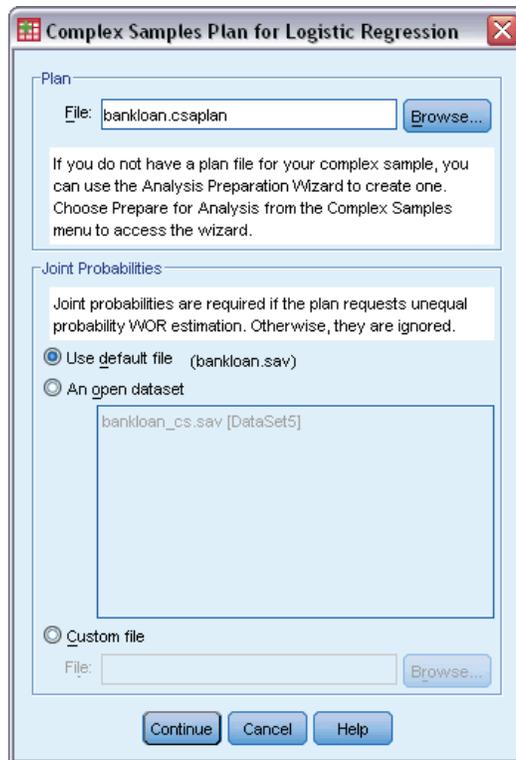
If you are a loan officer at a bank, you want to be able to identify characteristics that are indicative of people who are likely to default on loans and then use those characteristics to identify good and bad credit risks.

Suppose that a loan officer has collected past records of customers given loans at several different branches, according to a complex design. This information is contained in *bankloan\_cs.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). The officer wants to see if the probability with which a customer defaults is related to age, employment history, and amount of credit debt, incorporating the sampling design.

### ***Running the Analysis***

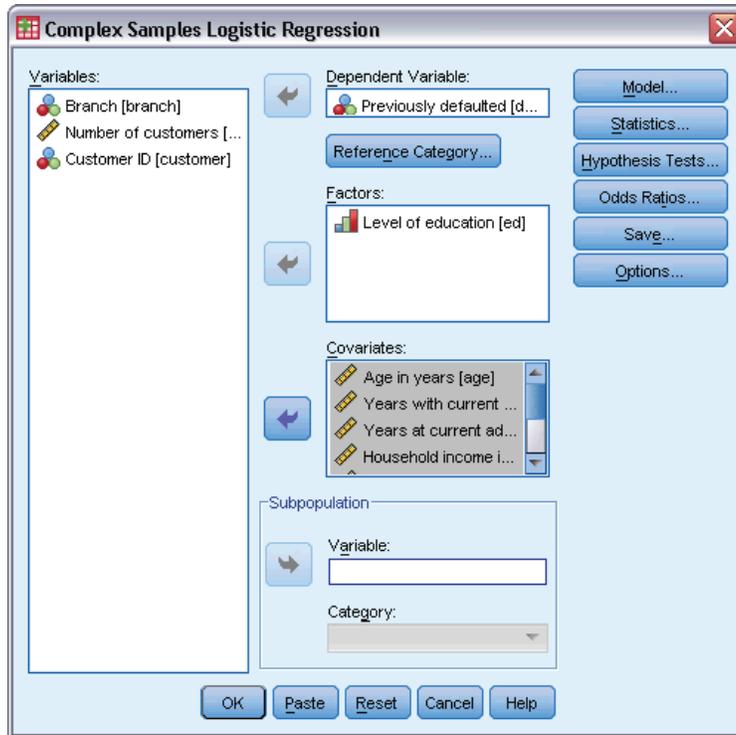
- ▶ To create the logistic regression model, from the menus choose:  
Analyze > Complex Samples > Logistic Regression...

Figure 20-1  
Complex Samples Plan dialog box



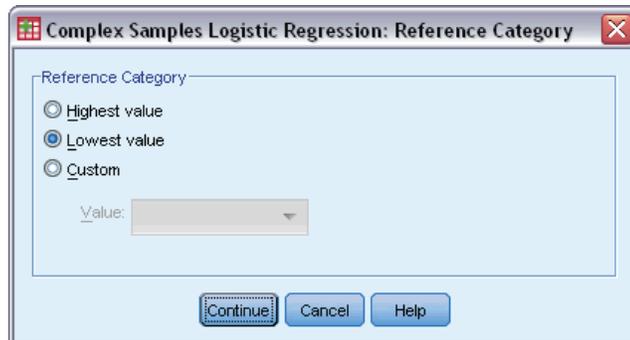
- ▶ Browse to and select *bankloan.csaplan*. For more information, see the topic [Sample Files](#) in [Appendix A in IBM SPSS Complex Samples 19](#).
- ▶ Click Continue.

Figure 20-2  
Logistic Regression dialog box



- ▶ Select *Previously defaulted* as the dependent variable.
- ▶ Select *Level of education* as a factor.
- ▶ Select *Age in years* through *Other debt in thousands* as covariates.
- ▶ Select *Previously defaulted* and click Reference Category.

Figure 20-3  
Logistic Regression Reference Category dialog box

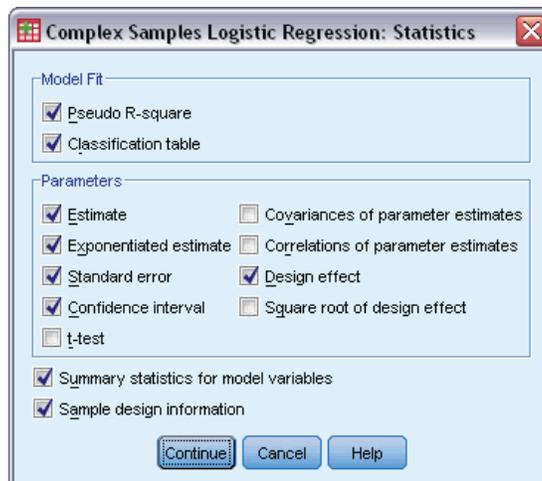


- ▶ Select Lowest value as the reference category.

This sets the “did not default” category as the reference category; thus, the odds ratios reported in the output will have the property that increasing odds ratios correspond to increasing probability of default.

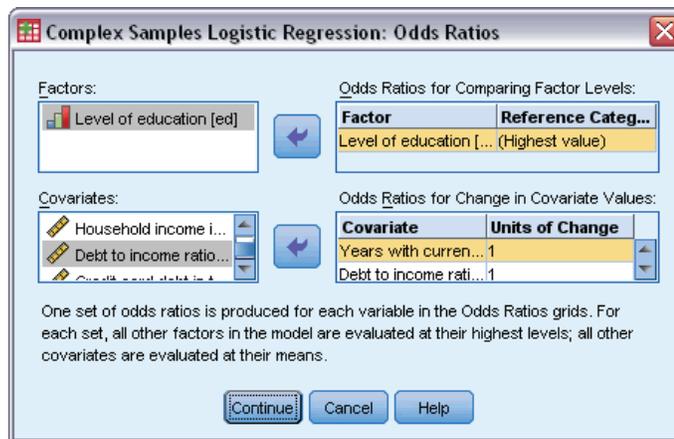
- ▶ Click Continue.
- ▶ Click Statistics in the Logistic Regression dialog box.

Figure 20-4  
Logistic Regression Statistics dialog box



- ▶ Select Classification table in the Model Fit group
- ▶ Select Estimate, Exponentiated estimate, Standard error, Confidence interval, and Design effect in the Parameters group.
- ▶ Click Continue.
- ▶ Click Odds Ratios in the Logistic Regression dialog box.

Figure 20-5  
Logistic Regression Odds Ratios dialog box



- ▶ Choose to create odds ratios for the factor *ed* and the covariates *employ* and *debtinc*.
- ▶ Click Continue.
- ▶ Click OK in the Logistic Regression dialog box.

## Pseudo R-Squares

Figure 20-6  
Pseudo R-square statistics

Cox and Snell	.330
Nagelkerke	.451
McFadden	.304

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

In the linear regression model, the coefficient of determination,  $R^2$ , summarizes the proportion of variance in the dependent variable associated with the predictor (independent) variables, with larger  $R^2$  values indicating that more of the variation is explained by the model, to a maximum of 1. For regression models with a categorical dependent variable, it is not possible to compute a single  $R^2$  statistic that has all of the characteristics of  $R^2$  in the linear regression model, so these approximations are computed instead. The following methods are used to estimate the coefficient of determination.

- Cox and Snell's  $R^2$  (Cox and Snell, 1989) is based on the log likelihood for the model compared to the log likelihood for a baseline model. However, with categorical outcomes, it has a theoretical maximum value of less than 1, even for a "perfect" model.
- Nagelkerke's  $R^2$  (Nagelkerke, 1991) is an adjusted version of the Cox & Snell  $R$ -square that adjusts the scale of the statistic to cover the full range from 0 to 1.
- McFadden's  $R^2$  (McFadden, 1974) is another version, based on the log-likelihood kernels for the intercept-only model and the full estimated model.

What constitutes a “good”  $R^2$  value varies between different areas of application. While these statistics can be suggestive on their own, they are most useful when comparing competing models for the same data. The model with the largest  $R^2$  statistic is “best” according to this measure.

## Classification

Figure 20-7  
Classification table

Observed	Predicted		
	No	Yes	Percent Correct
No	188289.667	31871.267	85.5%
Yes	49970.600	77675.133	60.9%
Overall Percent	68.5%	31.5%	76.5%

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

The classification table shows the practical results of using the logistic regression model. For each case, the predicted response is *Yes* if that case’s model-predicted logit is greater than 0. Cases are weighted by *finalweight*, so that the classification table reports the expected model performance in the population.

- Cells on the diagonal are correct predictions.
- Cells off the diagonal are incorrect predictions.

Based upon the cases used to create the model, you can expect to correctly classify 85.5% of the nondefaulters in the population using this model. Likewise, you can expect to correctly classify 60.9% of the defaulters. Overall, you can expect to classify 76.5% of the cases correctly; however, because this table was constructed with the cases used to create the model, these estimates are likely to be overly optimistic.

## Tests of Model Effects

Figure 20-8  
Tests of between-subjects effects

Source	df1	df2	Wald F	Sig.
(Corrected Model)	11.000	4.000	14.669	.010
(Intercept)	1.000	14.000	5.777	.031
ed	4.000	11.000	1.683	.224
age	1.000	14.000	5.352	.036
employ	1.000	14.000	88.244	.000
address	1.000	14.000	1.123	.307
income	1.000	14.000	.007	.932
debtinc	1.000	14.000	27.632	.000
creddebt	1.000	14.000	33.402	.000
othdebt	1.000	14.000	.709	.414

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

Each term in the model, plus the model as a whole, is tested for whether its effect equals 0. Terms with significance values less than 0.05 have some discernible effect. Thus, *age*, *employ*, *debtinc*, and *creddebt* contribute to the model, while the other main effects do not. In a further analysis of the data, you would probably remove *ed*, *address*, *income*, and *othdebt* from model consideration.

## Parameter Estimates

Figure 20-9  
Parameter estimates

Previously defaulted	Parameter	B	Std. Error	95% Confidence Interval		Design Effect	Exp(B)	95% Confidence Interval for Exp(B)	
				Lower	Upper			Lower	Upper
Yes	(Intercept)	-1.140	.399	-1.995	-.284	.665	.320	.136	.753
	[ed=1]	.720	.340	-.010	1.449	.862	2.054	.990	4.259
	[ed=2]	.684	.371	-.112	1.481	1.247	1.983	.894	4.397
	[ed=3]	.518	.307	-.140	1.177	.813	1.679	.869	3.244
	[ed=4]	.789	.302	.142	1.437	.817	2.202	1.152	4.208
	[ed=5]	.000 <sup>a</sup>	.	.	.	.	1.000	.	.
	age	-.023	.010	-.043	-.002	.418	.978	.958	.998
	employ	-.225	.024	-.277	-.174	1.200	.798	.758	.840
	address	-.028	.026	-.085	.029	.651	.972	.919	1.029
	income	.000	.003	-.007	.006	1.410	1.000	.993	1.006
	debtinc	.095	.018	.056	.134	1.222	1.100	1.058	1.143
	creddebt	.493	.085	.310	.676	1.373	1.637	1.363	1.966
	othdebt	.026	.031	-.041	.094	1.219	1.027	.960	1.098

Dependent Variable: Previously defaulted (reference category = No)

Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

a. Set to zero because this parameter is redundant.

The parameter estimates table summarizes the effect of each predictor. Note that parameter values affect the likelihood of the “did default” category relative to the “did not default” category. Thus, parameters with positive coefficients increase the likelihood of default, while parameters with negative coefficients decrease the likelihood of default.

The meaning of a logistic regression coefficient is not as straightforward as that of a linear regression coefficient. While  $B$  is convenient for testing the model effects,  $Exp(B)$  is easier to interpret.  $Exp(B)$  represents the ratio change in the odds of the event of interest attributable to a one-unit increase in the predictor for predictors that are not part of interaction terms. For example,  $Exp(B)$  for *employ* is equal to 0.798, which means that the odds of default for people who have been with their current employer for two years are 0.798 times the odds of default for those who have been with their current employer for one year, all other things being equal.

The design effects indicate that some of the standard errors computed for these parameter estimates are larger than those you would obtain if you assumed that these observations came from a simple random sample, while others are smaller. It is vitally important to incorporate the sampling design information in your analysis because you might otherwise infer, for example, that the age coefficient is no different from 0!

## Odds Ratios

Figure 20-10  
Odds ratios for level of education

		Previously defaulted	Odds Ratio	95% Confidence Interval	
Level of education				Lower	Upper
Did not complete high school vs. Post-undergraduate degree	Yes		2.054	.990	4.259
High school degree vs.	Yes		1.983	.894	4.397
Some college vs.	Yes		1.679	.869	3.244
College degree vs.	Yes		2.202	1.152	4.208

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

- a. Factors and covariates used in the computation are fixed at the following values:  
Level of education=Post-undergraduate degree; Age in years=34.19; Years with current employer=6.99; Years at current address=6.32; Household income in thousands=60.1581; Debt to income ratio (x100)=9.9341; Credit card debt in thousands=1.9764; Other debt in thousands=3.9164

This table displays the odds ratios of *Previously defaulted* at the factor levels of *Level of education*. The reported values are the ratios of the odds of default for *Did not complete high school* through *College degree*, compared to the odds of default for *Post-undergraduate degree*. Thus, the odds ratio of 2.054 in the first row of the table means that the odds of default for a person who did not complete high school are 2.054 times the odds of default for a person who has a post-undergraduate degree.

Figure 20-11  
Odds ratios for years with current employer

Units of Change		Previously defaulted	Odds Ratio	95% Confidence Interval	
Years with current employer				Lower	Upper
1.000	Yes		.798	.758	.840

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

- a. Factors and covariates used in the computation are fixed at the following values:  
Level of education=Post-undergraduate degree; Age in years=34.19; Years with current employer=6.99; Years at current address=6.32; Household income in thousands=60.1581; Debt to income ratio (x100)=9.9341; Credit card debt in thousands=1.9764; Other debt in thousands=3.9164

This table displays the odds ratio of *Previously defaulted* for a unit change in the covariate *Years with current employer*. The reported value is the ratio of the odds of default for a person with 7.99 years at their current job compared to the odds of default for a person with 6.99 years (the mean).

Figure 20-12  
Odds ratios for debt to income ratio

Units of Change		Previously defaulted	Odds Ratio	95% Confidence Interval	
Debt to income ratio (x100)				Lower	Upper
1.000	Yes		1.100	1.058	1.143

Dependent Variable: Previously defaulted (reference category = No)  
Model: (Intercept), ed, age, employ, address, income, debtinc, creddebt, othdebt

- a. Factors and covariates used in the computation are fixed at the following values:  
Level of education=Post-undergraduate degree; Age in years=34.19; Years with current employer=6.99; Years at current address=6.32; Household income in thousands=60.1581; Debt to income ratio (x100)=9.9341; Credit card debt in thousands=1.9764; Other debt in thousands=3.9164

This table displays the odds ratio of *Previously defaulted* for a unit change in the covariate *Debt to income ratio*. The reported value is the ratio of the odds of default for a person with a debt/income ratio of 10.9341 compared to the odds of default for a person with 9.9341 (the mean).

Note that because none of these predictors are part of interaction terms, the values of the odds ratios reported in these tables are equal to the values of the exponentiated parameter estimates. When a predictor is part of an interaction term, its odds ratio as reported in these tables will also depend on the values of the other predictors that make up the interaction.

## **Summary**

Using the Complex Samples Logistic Regression Procedure, you have constructed a model for predicting the probability that a given customer will default on a loan.

A critical issue for loan officers is the cost of Type I and Type II errors. That is, what is the cost of classifying a defaulter as a nondefaulter (Type I)? What is the cost of classifying a nondefaulter as a defaulter (Type II)? If bad debt is the primary concern, then you want to lower your Type I error and maximize your **sensitivity**. If growing your customer base is the priority, then you want to lower your Type II error and maximize your **specificity**. Usually, both are major concerns, so you have to choose a decision rule for classifying customers that gives the best mix of sensitivity and specificity.

## **Related Procedures**

The Complex Samples Logistic Regression procedure is a useful tool for modeling a categorical variable when the cases have been drawn according to a complex sampling scheme.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to specify analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples General Linear Model](#) procedure allows you to model a scale response.
- The [Complex Samples Ordinal Regression](#) procedure allows you to model an ordinal response.

# ***Complex Samples Ordinal Regression***

The Complex Samples Ordinal Regression procedure creates a predictive model for an ordinal dependent variable for samples drawn by complex sampling methods. Optionally, you can request analyses for a subpopulation.

## ***Using Complex Samples Ordinal Regression to Analyze Survey Results***

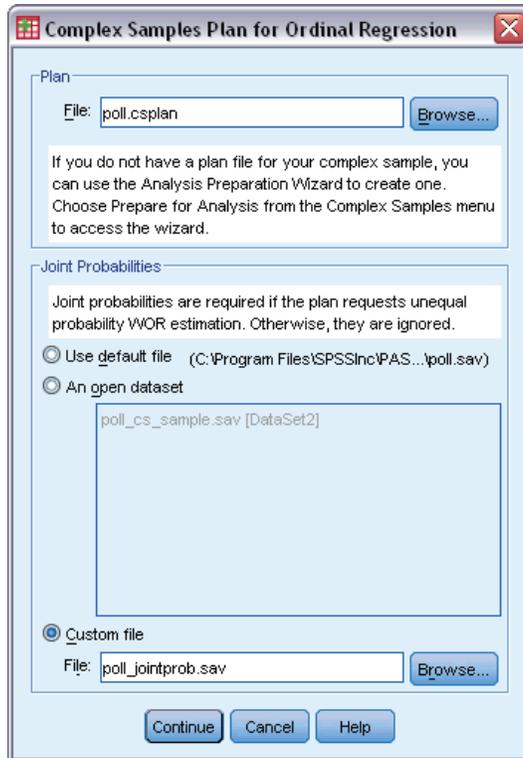
Representatives considering a bill before the legislature are interested in whether there is public support for the bill and how support for the bill is related to voter demographics. Pollsters design and conduct interviews according to a complex sampling design.

The survey results are collected in *poll\_cs\_sample.sav*. The sampling plan used by the pollsters is contained in *poll.csplan*; because it makes use of a probability-proportional-to-size (PPS) method, there is also a file containing the joint selection probabilities (*poll\_jointprob.sav*). [For more information, see the topic Sample Files in Appendix A in IBM SPSS Complex Samples 19.](#) Use Complex Samples Ordinal Regression to fit a model for the level of support for the bill based upon voter demographics.

### ***Running the Analysis***

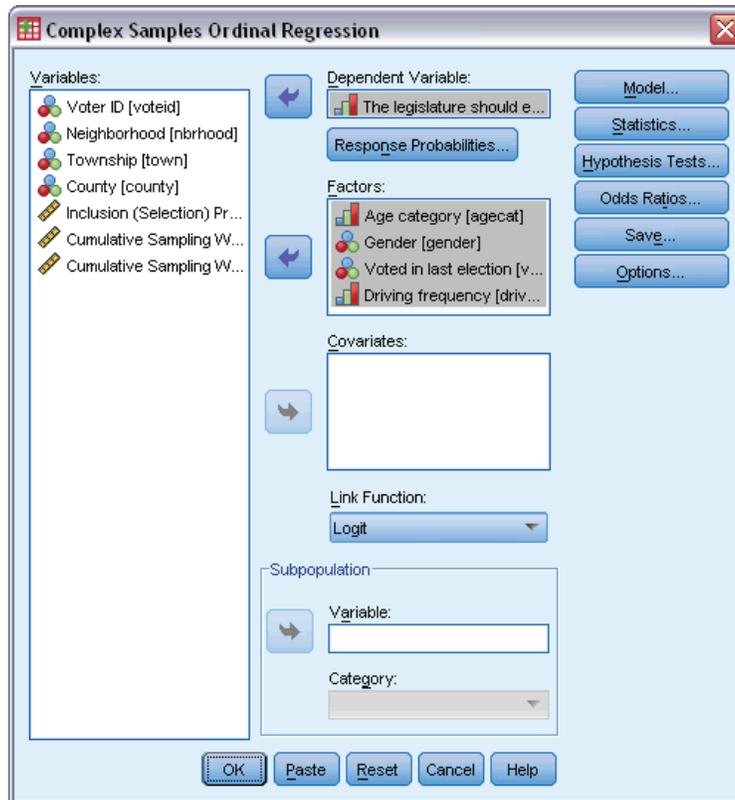
- ▶ To run a Complex Samples Ordinal Regression analysis, from the menus choose:  
Analyze > Complex Samples > Ordinal Regression...

Figure 21-1  
Complex Samples Plan dialog box



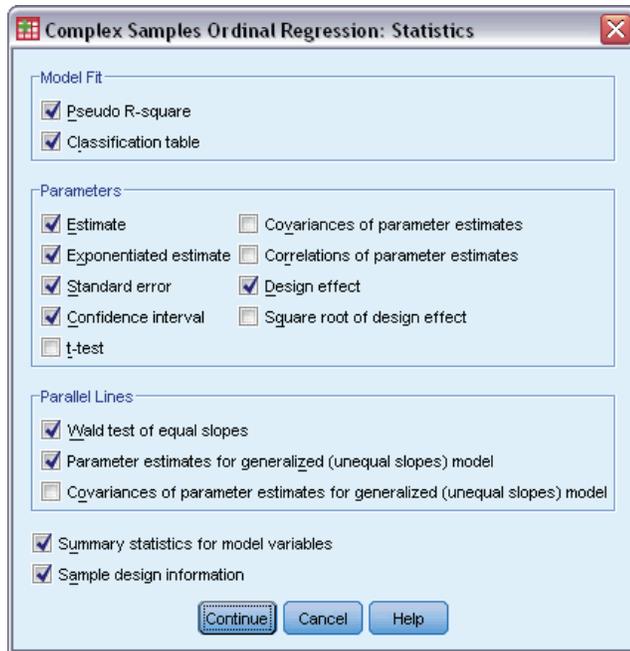
- ▶ Browse to and select *poll.csplan* as the plan file. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#).
- ▶ Select *poll\_jointprob.sav* as the joint probabilities file.
- ▶ Click Continue.

Figure 21-2  
Ordinal Regression dialog box



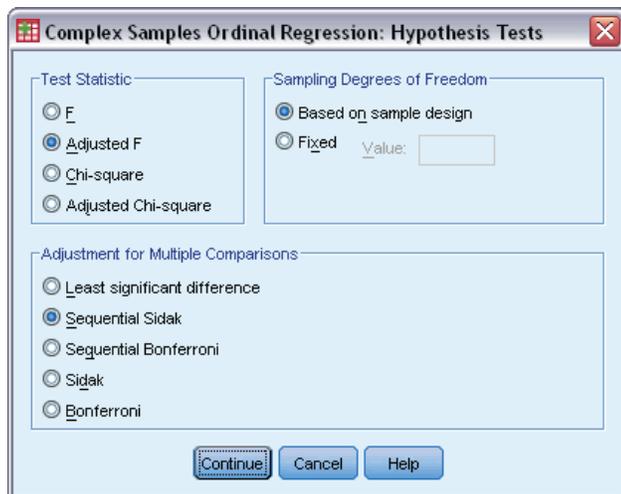
- ▶ Select *The legislature should enact a gas tax* as the dependent variable.
- ▶ Select *Age category* through *Driving frequency* as factors.
- ▶ Click *Statistics*.

Figure 21-3  
Ordinal Regression Statistics dialog box



- ▶ Select Classification table in the Model Fit group.
- ▶ Select Estimate, Exponentiated estimate, Standard error, Confidence interval, and Design effect in the Parameters group.
- ▶ Select Wald test of equal slopes and Parameter estimates for generalized (unequal slopes) model.
- ▶ Click Continue.
- ▶ Click Hypothesis Tests in the Complex Samples Ordinal Regression dialog box.

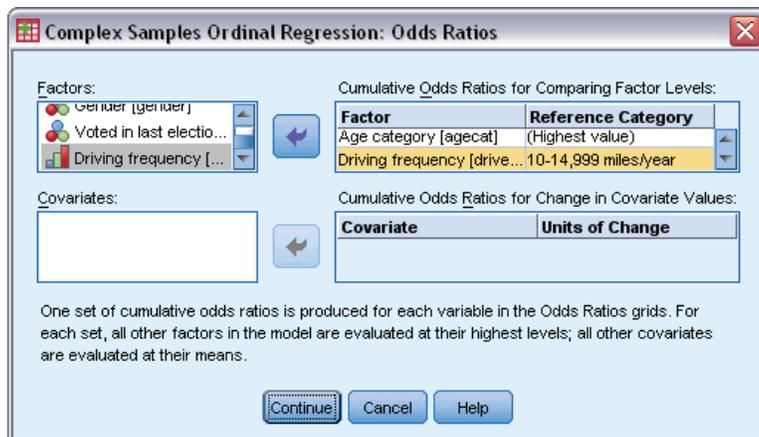
Figure 21-4  
Hypothesis Tests dialog box



Even for a moderate number of predictors and response categories, the Wald  $F$  test statistic can be inestimable for the test of parallel lines.

- ▶ Select Adjusted F in the Test Statistic group.
- ▶ Select Sequential Sidak as the adjustment method for multiple comparisons.
- ▶ Click Continue.
- ▶ Click Odds Ratios in the Complex Samples Ordinal Regression dialog box.

Figure 21-5  
Ordinal Regression Odds Ratios dialog box



- ▶ Choose to produce cumulative odds ratios for *Age category* and *Driving frequency*.
- ▶ Select 10-14,999 miles/year, a more “typical” yearly mileage than the maximum, as the reference category for *Driving frequency*.
- ▶ Click Continue.

- Click OK in the Complex Samples Ordinal Regression dialog box.

## Pseudo R-Squares

Figure 21-6  
Pseudo R-Squares

Cox and Snell	.179
Nagelkerke	.191
McFadden	.071

Dependent Variable: The legislature should enact a gas tax (Ascending)  
Model: (Threshold), agecat, gender, votelast, drivefreq  
Link function: Logit

In the linear regression model, the coefficient of determination,  $R^2$ , summarizes the proportion of variance in the dependent variable associated with the predictor (independent) variables, with larger  $R^2$  values indicating that more of the variation is explained by the model, to a maximum of 1. For regression models with a categorical dependent variable, it is not possible to compute a single  $R^2$  statistic that has all of the characteristics of  $R^2$  in the linear regression model, so these approximations are computed instead. The following methods are used to estimate the coefficient of determination.

- Cox and Snell's  $R^2$  (Cox and Snell, 1989) is based on the log likelihood for the model compared to the log likelihood for a baseline model. However, with categorical outcomes, it has a theoretical maximum value of less than 1, even for a “perfect” model.
- Nagelkerke's  $R^2$  (Nagelkerke, 1991) is an adjusted version of the Cox & Snell  $R$ -square that adjusts the scale of the statistic to cover the full range from 0 to 1.
- McFadden's  $R^2$  (McFadden, 1974) is another version, based on the log-likelihood kernels for the intercept-only model and the full estimated model.

What constitutes a “good”  $R^2$  value varies between different areas of application. While these statistics can be suggestive on their own, they are most useful when comparing competing models for the same data. The model with the largest  $R^2$  statistic is “best” according to this measure.

## Tests of Model Effects

Figure 21-7  
Tests of model effects

Source	df1	df2	Adjusted Wald F	Sig.	Sequential Sidak Sig.
agecat	2.283	31.966	6.215	.004	.003
gender	1.000	14.000	.046	.834	.834
votelast	1.000	14.000	.076	.787	.787
drivefreq	3.785	52.987	228.015	.000	.000

Dependent Variable: The legislature should enact a gas tax (Ascending)  
Model: (Threshold), agecat, gender, votelast, drivefreq  
Link function: Logit

Each term in the model is tested for whether its effect equals 0. Terms with significance values less than 0.05 have some discernable effect. Thus, *agecat* and *drivefreq* contribute to the model, while the other main effects do not. In a further analysis of the data, you would consider removing *gender* and *votelast* from the model.

## Parameter Estimates

The parameter estimates table summarizes the effect of each predictor. While interpretation of the coefficients in this model is difficult due to the nature of the link function, the signs of the coefficients for covariates and relative values of the coefficients for factor levels can give important insights into the effects of the predictors in the model.

- For covariates, positive (negative) coefficients indicate positive (inverse) relationships between predictors and outcome. An increasing value of a covariate with a positive coefficient corresponds to an increasing probability of being in one of the “higher” cumulative outcome categories.
- For factors, a factor level with a greater coefficient indicates a greater probability of being in one of the “higher” cumulative outcome categories. The sign of a coefficient for a factor level is dependent upon that factor level’s effect relative to the reference category.

Figure 21-8  
Parameter estimates

Parameter	B	Std. Error	95% Confidence Interval		Design Effect	Exp(B)	95% Confidence Interval for Exp(B)		
			Lower	Upper			Lower	Upper	
Threshold	[opinion_gastax=1]	-3.343	.104	-3.566	-3.120	1.132	.035	.028	.044
	[opinion_gastax=2]	-1.910	.098	-2.120	-1.700	1.058	.148	.120	.183
	[opinion_gastax=3]	-.674	.090	-.866	-.482	.915	.510	.421	.618
Regression	[agecat=1]	-.324	.079	-.494	-.154	1.793	.723	.610	.858
	[agecat=2]	-.138	.054	-.255	-.022	1.158	.871	.775	.978
	[agecat=3]	-.095	.076	-.257	.068	2.206	.909	.773	1.070
	[agecat=4]	.000 <sup>a</sup>	.	.	.	.	1.000	.	.
	[gender=0]	-.008	.035	-.084	.068	.949	.992	.920	1.071
	[gender=1]	.000 <sup>a</sup>	.	.	.	.	1.000	.	.
	[votelast=0]	-.011	.039	-.095	.073	1.103	.989	.909	1.076
	[votelast=1]	.000 <sup>a</sup>	.	.	.	.	1.000	.	.
	[drivefreq=1]	-3.751	.153	-4.079	-3.423	1.117	.023	.017	.033
	[drivefreq=2]	-3.003	.116	-3.251	-2.755	1.226	.050	.039	.064
	[drivefreq=3]	-2.295	.114	-2.540	-2.050	1.585	.101	.079	.129
	[drivefreq=4]	-1.570	.092	-1.769	-1.372	1.078	.208	.171	.254
[drivefreq=5]	-.812	.089	-1.003	-.621	.941	.444	.367	.537	
[drivefreq=6]	.000 <sup>a</sup>	.	.	.	.	1.000	.	.	

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

a. Set to zero because this parameter is redundant.

You can make the following interpretations based on the parameter estimates:

- Those in lower age categories show greater support for the bill than those in the highest age category.

- Those who drive less frequently show greater support for the bill than those who drive more frequently.
- The coefficients for the variables *gender* and *votelast*, in addition to not being statistically significant, appear to be small compared to other coefficients.

The design effects indicate that some of the standard errors computed for these parameter estimates are larger than those you would obtain if you used a simple random sample, while others are smaller. It is vitally important to incorporate the sampling design information in your analysis because you might otherwise infer, for example, that the coefficient for the third level of *Age category*, [*agecat*=3], is significantly different from 0!

## Classification

Figure 21-9  
Categorical variable information

		Weighted Count	Weighted Percent
The legislature should enact a gas tax <sup>a</sup>	Strongly agree	25132.955	21.3%
	Agree	32261.425	27.3%
	Disagree	29477.417	24.9%
	Strongly disagree	31314.203	26.5%
Age category	18-30	20509.504	17.4%
	31-45	35380.506	29.9%
	46-60	34865.792	29.5%
	>60	27430.198	23.2%
Gender	Male	61424.547	52.0%
	Female	56761.453	48.0%
Voted in last election	No	70607.216	59.7%
	Yes	47578.784	40.3%
Driving frequency	Do not own car	3437.137	2.9%
	<10,000 miles/year	10816.349	9.2%
	10-14,999 miles/year	32539.364	27.5%
	15-19,999 miles/year	39179.814	33.2%
	20-29,999 miles/year	25617.804	21.7%
	>=30,000 miles/year	6595.532	5.6%
Population Size		118186.000	100.0%

a. Dependent variable values are sorted in ascending order.

Given the observed data, the “null” model (that is, one without predictors) would classify all customers into the modal group, *Agree*. Thus, the null model would be correct 27.3% of the time.

Figure 21-10  
Classification table

Observed	Predicted				Percent Correct
	Strongly agree	Agree	Disagree	Strongly disagree	
Strongly agree	7067.567	12130.814	3875.825	2058.750	28.1%
Agree	4271.234	14464.286	7320.767	6205.137	44.8%
Disagree	2024.816	11703.368	7108.487	8640.746	24.1%
Strongly disagree	889.869	8169.109	6946.522	15308.703	48.9%
Overall Percent	12.1%	39.3%	21.4%	27.3%	37.2%

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

The classification table shows the practical results of using the model. For each case, the predicted response is the response category with the highest model-predicted probability. Cases are weighted by *Final Sampling Weight*, so that the classification table reports the expected model performance in the population.

- Cells on the diagonal are correct predictions.
- Cells off the diagonal are incorrect predictions.

The model correctly classifies 9.9% more, or 37.2% of the cases. In particular, the model does considerably better at classifying those who *Agree* or *Strongly disagree*, and slightly worse with those who *Disagree*.

## Odds Ratios

**Cumulative odds** are defined as the ratio of the probability that the dependent variable takes a value less than or equal to a given response category to the probability that it takes a value greater than that response category. The **cumulative odds ratio** is the ratio of cumulative odds for different predictor values and is closely related to the exponentiated parameter estimates. Interestingly, the cumulative odds ratio itself does not depend upon the response category.

Figure 21-11  
Cumulative odds ratios for Age category

	Cumulative Odds Ratio	95% Confidence Interval		Design Effect	Square Root Design Effect
		Lower	Upper		
Age category 18-30 vs. >60	1.383	1.166	1.639	1.793	1.339
31-45 vs. >60	1.148	1.022	1.290	1.158	1.076
46-60 vs. >60	1.100	.935	1.294	2.206	1.485

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

- a. Factors and covariates used in the computation are fixed at the following values: Age category=>60; Gender=Female; Voted in last election=Yes; Driving frequency=>=30,000 miles/year

This table displays cumulative odds ratios for the factor levels of *Age category*. The reported values are the ratios of the cumulative odds for 18–30 through 46–60, compared to the cumulative odds for >60. Thus, the odds ratio of 1.383 in the first row of the table means that the cumulative odds for a person aged 18–30 are 1.383 times the cumulative odds for a person older than 60. Note that because *Age category* is not involved in any interaction terms, the odds ratios are

merely the ratios of the exponentiated parameter estimates. For example, the cumulative odds ratio for 18–30 vs. >60 is  $1.00/0.723 = 1.383$ .

Figure 21-12  
Odds ratios for driving frequency

		Cumulative Odds Ratio	95% Confidence Interval		Design Effect	Square Root Design Effect
			Lower	Upper		
Driving frequency	Do not own car vs. 10-14,999 miles/year	4.288	2.878	6.390	2.345	1.531
	<10,000 miles/year vs. 10-14,999 miles/year	2.030	1.656	2.488	1.838	1.356
	15-19,999 miles/year vs. 10-14,999 miles/year	.484	.430	.546	1.450	1.204
	20-29,999 miles/year vs. 10-14,999 miles/year	.227	.193	.267	2.095	1.448
	>=30,000 miles/year vs. 10-14,999 miles/year	.101	.079	.129	1.585	1.259

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

a. Factors and covariates used in the computation are fixed at the following values: Age category=>60; Gender=Female; Voted in last election=Yes; Driving frequency=>=30,000 miles/year

This table displays the cumulative odds ratios for the factor levels of *Driving frequency*, using *10–14,999 miles/year* as the reference category. Since *Driving frequency* is not involved in any interaction terms, the odds ratios are merely the ratios of the exponentiated parameter estimates. For example, the cumulative odds ratio for *20–29,999 miles/year* vs. *10–14,999 miles/year* is  $0.101/0.444 = 0.227$ .

## Generalized Cumulative Model

Figure 21-13  
Test of parallel lines

df1	df2	Adjusted Wald F	Sig.	Sequential Sidak Sig.
8.769	122.767	1.894	.061	.392

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

The test of parallel lines can help you assess whether the assumption that the parameters are the same for all response categories is reasonable. This test compares the estimated model with one set of coefficients for all categories to a generalized model with a separate set of coefficients for each category.

The Wald  $F$  test is an omnibus test of the contrast matrix for the parallel lines assumption that provides asymptotically correct  $p$  values; for small to mid-sized samples, the adjusted Wald  $F$  statistic performs well. The significance value is near 0.05, suggesting that the generalized model may give an improvement in the model fit; however, the Sequential Sidak adjusted test reports a significance value high enough (0.392) that, overall, there is no clear evidence for rejecting the parallel lines assumption. The Sequential Sidak test starts with individual contrast Wald tests to provide an overall  $p$  value, and these results should be comparable to the omnibus Wald test result. The fact that they are so different in this example is somewhat surprising but could be due to the existence of many contrasts in the test and a relatively small design degrees of freedom.

Figure 21-14  
Parameter estimates for generalized cumulative model (shown in part)

The legislature should enact a gas tax	Parameter	B	Std. Error	95% Confidence Interval	
				Lower	Upper
Strongly agree	(Threshold)	-3.681	.221	-4.155	-3.207
	[agecat=1]	-.320	.096	-.525	-.115
	[agecat=2]	-.075	.071	-.227	.077
	[agecat=3]	-.022	.073	-.180	.135
	[agecat=4]	.000 <sup>a</sup>	.		
	[gender=0]	-.082	.054	-.197	.033
	[gender=1]	.000 <sup>a</sup>	.		
	[votelastr=0]	.008	.052	-.104	.120
	[votelastr=1]	.000 <sup>a</sup>	.		
	[drivfreq=1]	-4.096	.267	-4.669	-3.523
	[drivfreq=2]	-3.367	.237	-3.876	-2.857
	[drivfreq=3]	-2.678	.224	-3.158	-2.199
	[drivfreq=4]	-1.928	.213	-2.384	-1.471
	[drivfreq=5]	-1.015	.252	-1.555	-.476
[drivfreq=6]	.000 <sup>a</sup>	.			
Agree	(Threshold)	-1.963	.153	-2.291	-1.635
	[agecat=1]	-.385	.095	-.587	-.182
	[agecat=2]	-.130	.069	-.279	.018
	[agecat=3]	-.139	.101	-.356	.077
	[agecat=4]	.000 <sup>a</sup>	.		
	[gender=0]	-.004	.040	-.090	.082
	[gender=1]	.000 <sup>a</sup>	.		
	[votelastr=0]	.009	.059	-.117	.135
	[votelastr=1]	.000 <sup>a</sup>	.		
	[drivfreq=1]	-3.867	.318	-4.549	-3.185
	[drivfreq=2]	-3.005	.175	-3.380	-2.630
	[drivfreq=3]	-2.290	.187	-2.691	-1.888
	[drivfreq=4]	-1.633	.166	-1.988	-1.278
	[drivfreq=5]	-.909	.137	-1.204	-.615
[drivfreq=6]	.000 <sup>a</sup>	.			

Moreover, the estimated values of the generalized model coefficients don't appear to differ much from the estimates under the parallel lines assumption.

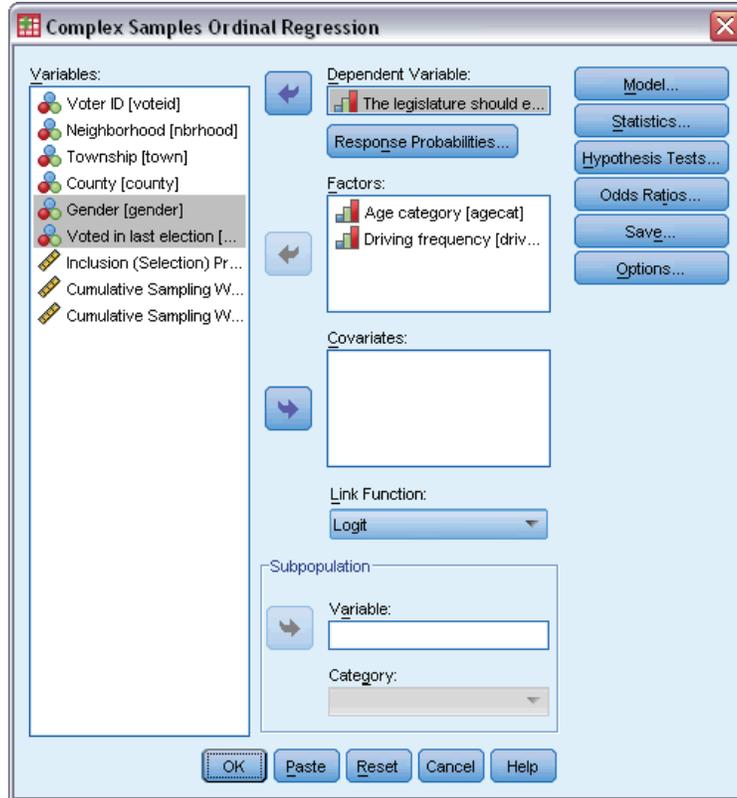
### ***Dropping Non-Significant Predictors***

The tests of model effects showed that the model coefficients for *Gender* and *Voted in last election* are not statistically significantly different from 0.

- To produce a reduced model, recall the Complex Samples Ordinal Regression dialog box.

- ▶ Click Continue in the Plan dialog box.

Figure 21-15  
Ordinal Regression dialog box



- ▶ Deselect *Gender* and *Voted in last election* as factors.
- ▶ Click Options.

Figure 21-16  
Ordinal Regression Options dialog box

- ▶ Select Display iteration history.

The iteration history is useful for diagnosing problems encountered by the estimation algorithm.

- ▶ Click Continue.
- ▶ Click OK in the Complex Samples Ordinal Regression dialog box.

## Warnings

Figure 21-17  
Warnings for reduced model

The log-likelihood value cannot be increased after the maximum number of steps in the step-halving method.

The CSORDINAL procedure continues despite the above warning(s). Subsequent results shown are based on the last iteration. Validity of the model fit is uncertain.

The following message applies to the generalized cumulative model.

The log-likelihood value cannot be increased after the maximum number of steps in the step-halving method.

The warnings note that estimation of the reduced model ended before the parameter estimates reached convergence because the log-likelihood could not be increased with any change, or “step,” in the current values of the parameter estimates.

**Figure 21-18**  
*Warnings for reduced model*

Iteration Number <sup>b</sup>	N Step-halving	Pseudo -2 Log Likelihood	Threshold			Regression								
			[opinion _gastax =1]	[opinion _gastax =2]	[opinion _gastax =3]	[agec at=1]	[agec at=2]	[agec at=3]	[drive freq= 1]	[drive freq= 2]	[drive freq= 3]	[drive freq= 4]	[drive freq= 5]	
0	0	326640.3	-1.309	-.058	1.020	.000	.000	.000	.000	.000	.000	.000	.000	.000
1	0	303567.5	-3.242	-1.881	-.704	-.323	-.137	-.094	-3.841	-2.970	-2.248	-1.563	-.835	
2	0	303336.3	-3.327	-1.897	-.664	-.325	-.139	-.095	-3.740	-2.998	-2.291	-1.568	-.811	
3	0	303335.9	-3.333	-1.900	-.664	-.326	-.139	-.096	-3.750	-3.003	-2.295	-1.570	-.812	
4	0	303335.9	-3.333	-1.900	-.664	-.326	-.139	-.096	-3.750	-3.003	-2.295	-1.570	-.812	
5 <sup>a</sup>	5	303335.9	-3.333	-1.900	-.664	-.326	-.139	-.096	-3.750	-3.003	-2.295	-1.570	-.812	

Redundant parameters are not displayed. Their values are always zero in all iterations.

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, drivefreq

Link function: Logit

- a. The log-likelihood value cannot be increased after the maximum number of steps in the step-halving method.
- b. Newton-Raphson method was used to estimate the parameters.

Looking at the iteration history, the changes in the parameter estimates over the last few iterations are slight enough that you're not terribly concerned about the warning message.

### Comparing Models

**Figure 21-19**  
*Pseudo R-Squares for reduced model*

Cox and Snell	.179
Nagelkerke	.191
McFadden	.071

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, gender, votelast, drivefreq

Link function: Logit

The  $R^2$  values for the reduced model are identical to those for the original model. This is evidence in favor of the reduced model.

**Figure 21-20**  
*Classification table for reduced model*

Observed	Predicted				Percent Correct
	Strongly agree	Agree	Disagree	Strongly disagree	
Strongly agree	7067.567	12823.258	3183.380	2058.750	28.1%
Agree	4271.234	15684.090	6100.963	6205.137	48.6%
Disagree	2024.816	13157.809	5654.047	8640.746	19.2%
Strongly disagree	889.869	9226.578	5889.053	15308.703	48.9%
Overall Percent	12.1%	43.1%	17.6%	27.3%	37.0%

Dependent Variable: The legislature should enact a gas tax (Ascending)

Model: (Threshold), agecat, drivefreq

Link function: Logit

The classification table somewhat complicates matters. The overall classification rate of 37.0% for the reduced model is comparable to the original model, which is evidence in favor of the reduced model. However, the reduced model shifts the predicted response of 3.8% of the voters

from *Disagree* to *Agree*, more than half of whom were observed to respond *Disagree* or *Strongly disagree*. This is a very important distinction that deserves careful consideration before choosing the reduced model.

## Summary

Using the Complex Samples Ordinal Regression Procedure, you have constructed competing models for the level of support for the proposed bill based on voter demographics. The test of parallel lines shows that a generalized cumulative model is not necessary. The tests of model effects suggest that *Gender* and *Voted in last election* could be dropped from the model, and the reduced model performs well in terms of pseudo- $R^2$  and overall classification rate compared to the original model. However, the reduced model misclassifies more voters across the *Agree/Disagree* split, so the legislators prefer to keep the original model for now.

## Related Procedures

The Complex Samples Ordinal Regression procedure is a useful tool for modeling an ordinal variable when the cases have been drawn according to a complex sampling scheme.

- The [Complex Samples Sampling Wizard](#) is used to specify complex sampling design specifications and obtain a sample. The sampling plan file created by the Sampling Wizard contains a default analysis plan and can be specified in the Plan dialog box when you are analyzing the sample obtained according to that plan.
- The [Complex Samples Analysis Preparation Wizard](#) is used to specify analysis specifications for an existing complex sample. The analysis plan file created by the Sampling Wizard can be specified in the Plan dialog box when you are analyzing the sample corresponding to that plan.
- The [Complex Samples General Linear Model](#) procedure allows you to model a scale response.
- The [Complex Samples Logistic Regression](#) procedure allows you to model a categorical response.

# ***Complex Samples Cox Regression***

The Complex Samples Cox Regression procedure performs survival analysis for samples drawn by complex sampling methods.

## ***Using a Time-Dependent Predictor in Complex Samples Cox Regression***

A government law enforcement agency is concerned about recidivism rates in their area of jurisdiction. One of the measures of recidivism is the time until second arrest for offenders. The agency would like to model time to rearrest using Cox Regression on a sample drawn by complex sampling methods, but they are worried the proportional hazards assumption is invalid across age categories.

Persons released from their first arrest during the month of June 2003 were selected from sampled departments, and their case history inspected through the end of June 2006. The sample is collected in *recidivism\_cs\_sample.sav*. The sampling plan used is contained in *recidivism\_cs.csplan*; because it makes use of a probability-proportional-to-size (PPS) method, there is also a file containing the joint selection probabilities (*recidivism\_cs\_jointprob.sav*). [For more information, see the topic Sample Files in Appendix A in IBM SPSS Complex Samples 19.](#) Use Complex Samples Cox Regression to assess the validity of the proportional hazards assumption and fit a model with time-dependent predictors, if appropriate.

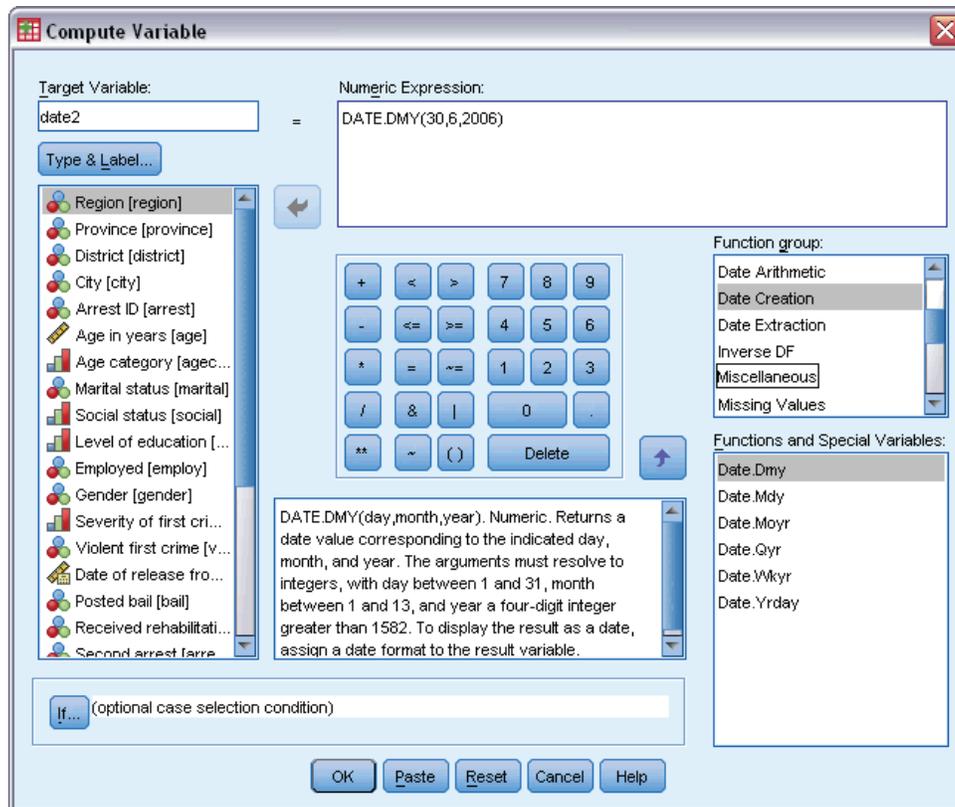
### ***Preparing the Data***

The dataset contains the dates of release from first arrest and second arrest; since Cox regression analyzes survival times, you need to compute the amount of time between these dates.

However, *Date of second arrest [date2]* contains cases with the value 10/03/1582, a missing value for date variables. These are people who have not had a second offense, and we definitely want to include them as right-censored cases in the model. The end of the follow-up period was June 30, 2006, so we are going to recode 10/03/1582 to 06/30/2006.

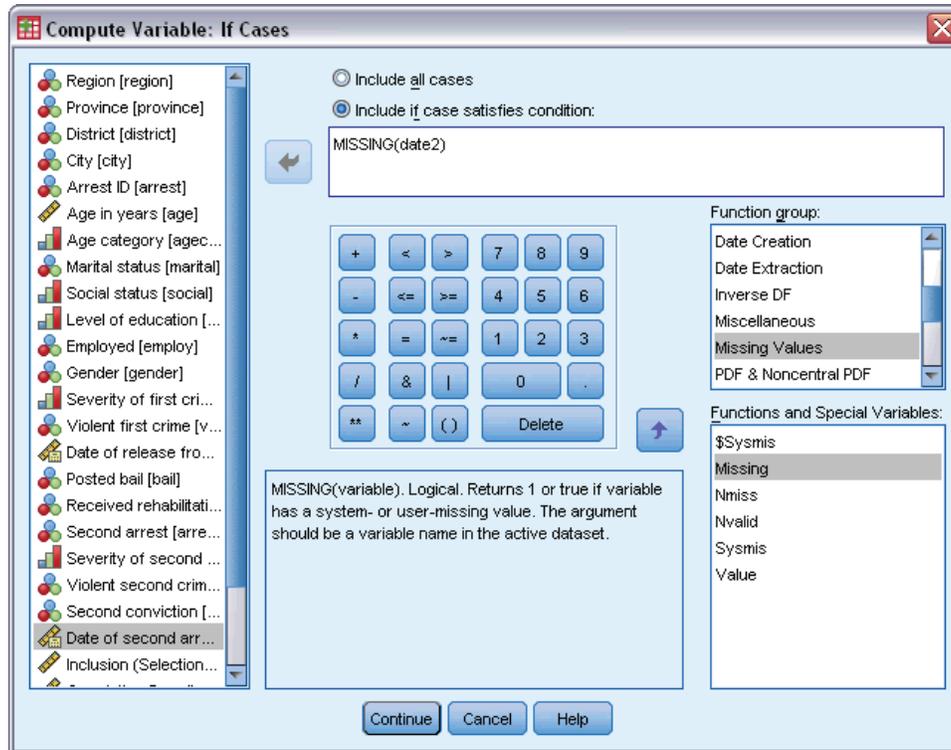
- ▶ To recode these values, from the menus choose:  
Transform > Compute Variable...

Figure 22-1  
Compute Variable dialog box



- ▶ Type date2 as the target variable.
- ▶ Type DATE.DMY(30,6,2006) as the expression.
- ▶ Click If.

Figure 22-2  
 Compute Variable If Cases dialog box



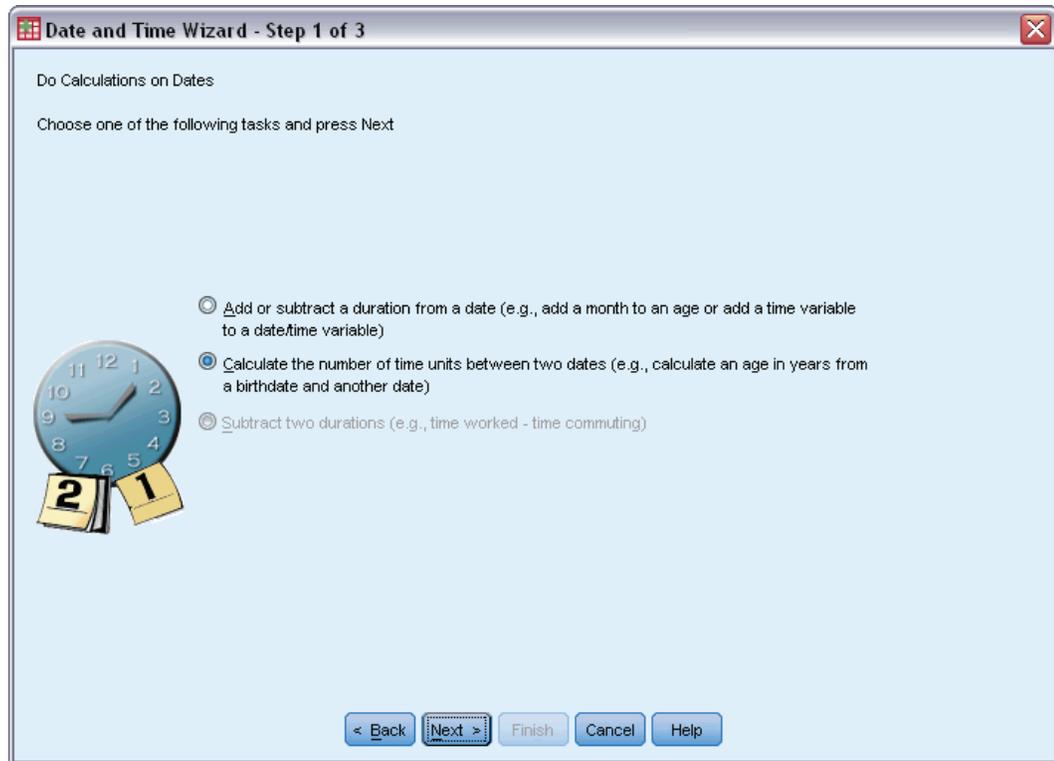
- ▶ Select Include if case satisfies condition.
- ▶ Type MISSING(date2) as the expression.
- ▶ Click Continue.
- ▶ Click OK in the Compute Variable dialog box.
- ▶ Next, to compute the time between first and second arrest, from the menus choose: Transform > Date and Time Wizard...

Figure 22-3  
Date and Time Wizard, Welcome step



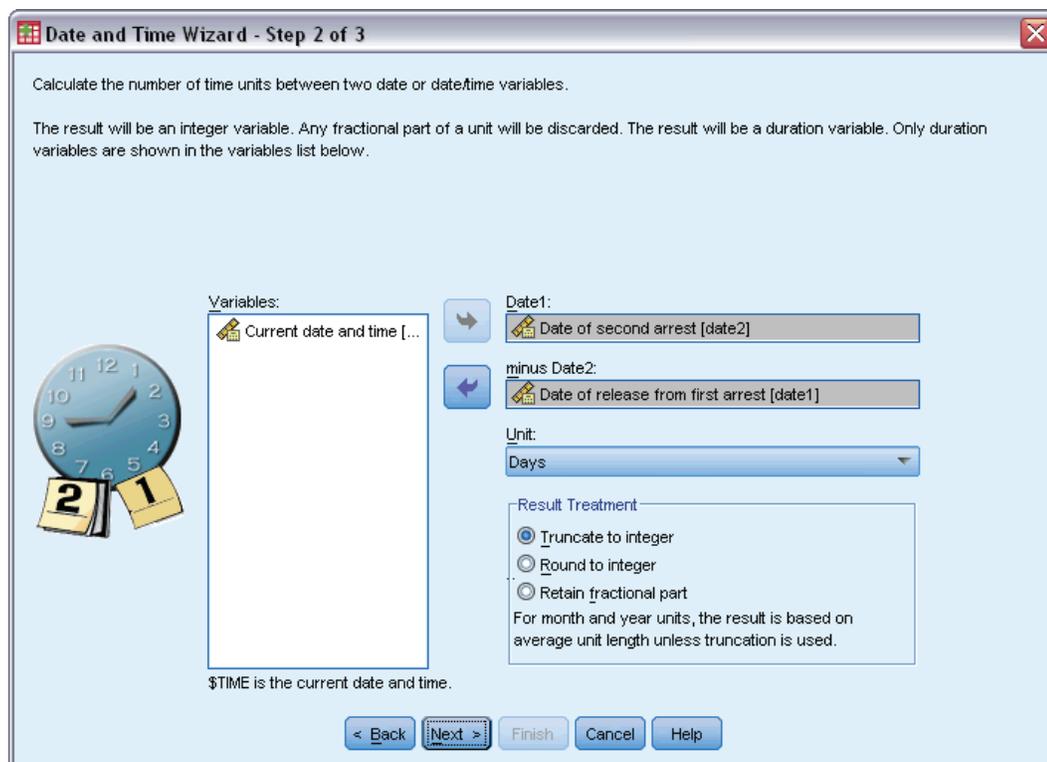
- ▶ Select Calculate with dates and times.
- ▶ Click Next.

Figure 22-4  
*Date and Time Wizard, Do Calculations on Dates step*



- ▶ Select Calculate the number of time units between two dates.
- ▶ Click Next.

Figure 22-5  
Date and Time Wizard, Calculate the number of time units between two dates step



- ▶ Select *Date of second arrest [date2]* as the first date.
- ▶ Select *Date of release from first arrest [date1]* as the date to subtract from the first date.
- ▶ Select Days as the unit.
- ▶ Click Next.

Figure 22-6  
Date and Time Wizard, Calculation step

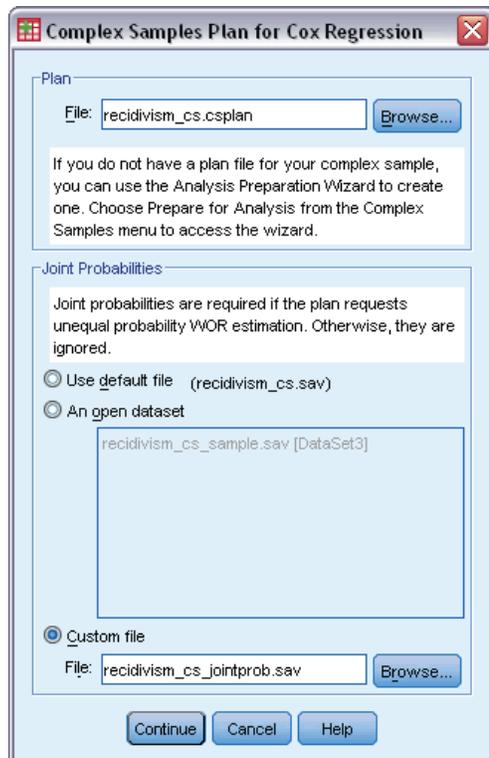


- ▶ Type *time\_to\_event* as the name of the variable representing the time between the two dates.
- ▶ Type *Time to second arrest* as the variable label.
- ▶ Click Finish.

### **Running the Analysis**

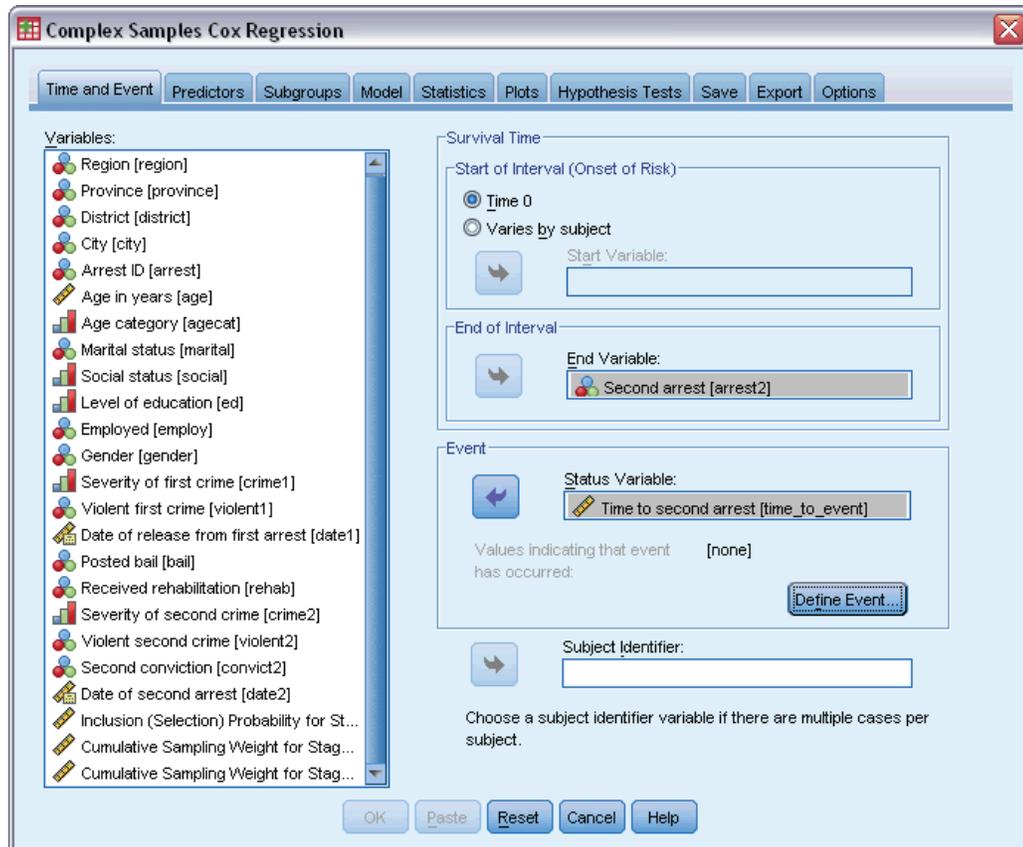
- ▶ To run a Complex Samples Cox Regression analysis, from the menus choose:  
Analyze > Complex Samples > Cox Regression...

Figure 22-7  
Complex Samples Plan for Cox Regression dialog box



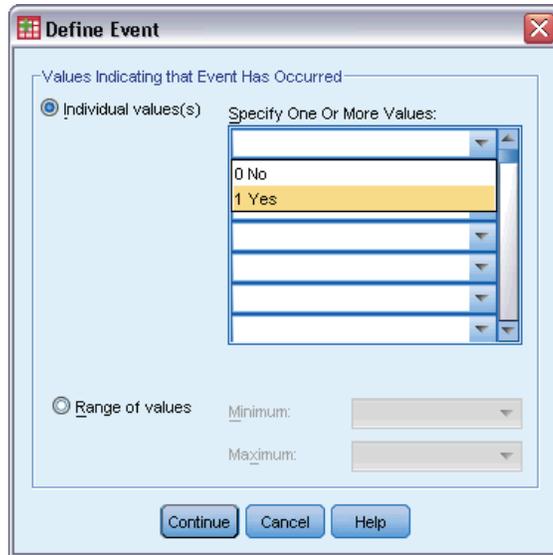
- ▶ Browse to the sample files directory and select *recidivism\_cs.csplan* as the plan file.
- ▶ Select Custom file in the Joint Probabilities group, browse to the sample files directory, and select *recidivism\_cs\_jointprob.sav*.
- ▶ Click Continue.

Figure 22-8  
Cox Regression dialog box, Time and Event tab



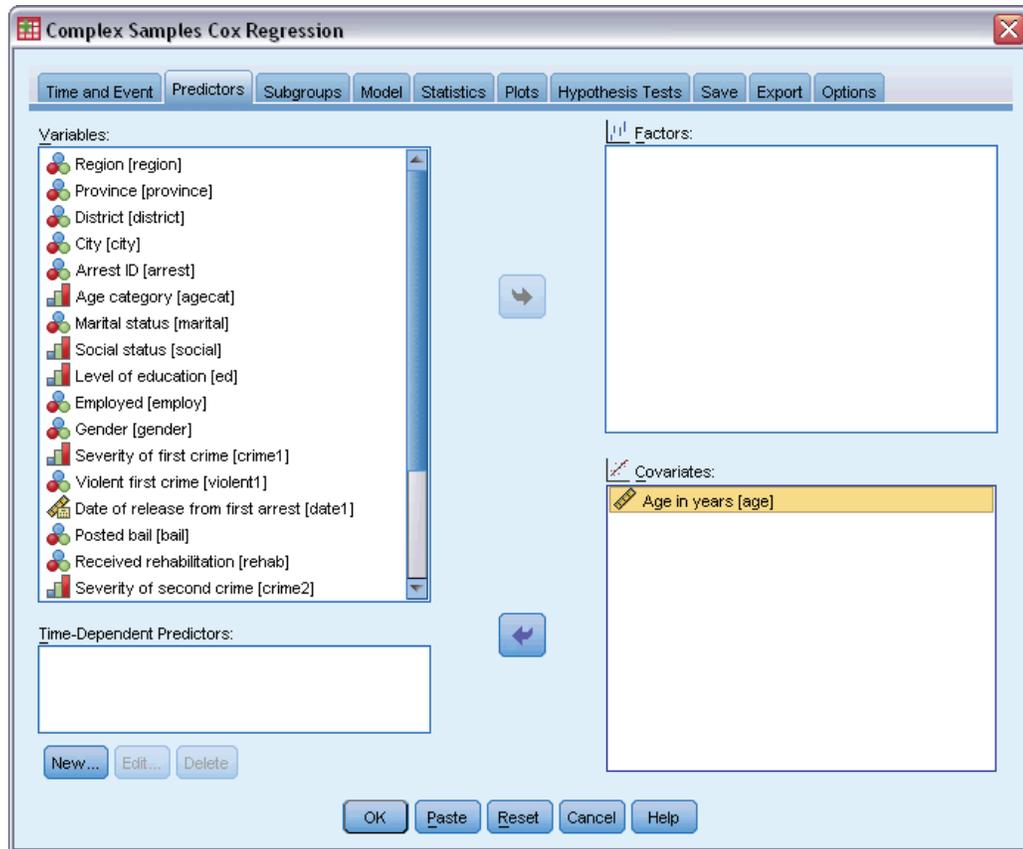
- ▶ Select *Time to second arrest [time\_to\_event]* as the variable defining the end of the interval.
- ▶ Select *Second arrest [arrest2]* as the variable defining whether the event has occurred.
- ▶ Click Define Event.

Figure 22-9  
Define Event dialog box



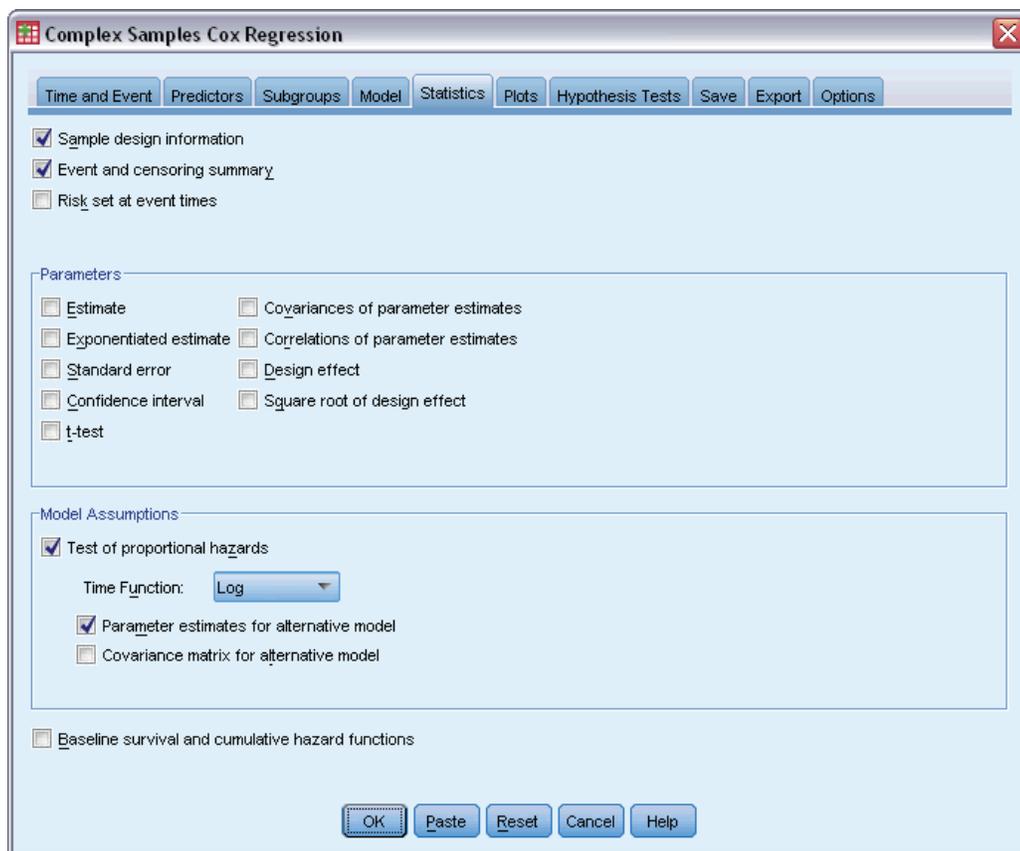
- ▶ Select 1 Yes as the value indicating the event of interest (rearrest) has occurred.
- ▶ Click Continue.
- ▶ Click the Predictors tab.

Figure 22-10  
Cox Regression dialog box, Predictors tab



- ▶ Select *Age in years [age]* as a covariate.
- ▶ Click the Statistics tab.

Figure 22-11  
Cox Regression dialog box, Statistics tab



- ▶ Select Test of proportional hazards and then select Log as the time function in the Model Assumptions group.
- ▶ Select Parameter estimates for alternative model.
- ▶ Click OK.

### Sample Design Information

Figure 22-12  
Sample design information

			N
Unweighted Counts	Valid	Subjects	5687
		Cases	5687
	Invalid Cases		0
	Total Cases		5687
Population Subject Size			307583.898
Stage 1	Strata		4
	Units		20
Sampling Design Degrees of Freedom			16

This table contains information on the sample design pertinent to the estimation of the model.

- There is one case per subject, and all 5,687 cases are used in the analysis.
- The sample represents less than 2% of the entire estimated population.
- The design requested 4 strata and 5 units per strata for a total of 20 units in the first stage of the design. The sampling design degrees of freedom are estimated by  $20-4=16$ .

### Tests of Model Effects

Figure 22-13  
Tests of model effects

Source	df1	df2	Wald F	Sig.
age	1.000	16	504.787	1.580E-13

Survival Time Variable: Time to second arrest  
Event Status Variable: Second arrest = 1.0  
Model: age

In the proportional hazards model, the significance value for the predictor *age* is less than 0.05 and, therefore, appears to contribute to the model.

### Test of Proportional Hazards

Figure 22-14  
Overall test of proportional hazards

df1	df2	Wald F	Sig.
1.000	16.000	29.924	5.136E-5

Survival Time Variable: Time to second arrest  
Event Status Variable: Second arrest = 1.0  
Model: age, age\*\_TF

Figure 22-15  
Parameter estimates for alternative model

Parameter	B	Std. Error	90% Confidence Interval	
			Lower	Upper
age	-.002	.014	-.025	0.02
age*_TF <sup>a</sup>	-.012	.002	-.016	-.009

Survival Time Variable: Time to second arrest  
Event Status Variable: Second arrest = 1.0  
Model: age, age\*\_TF

a. Time function: Log

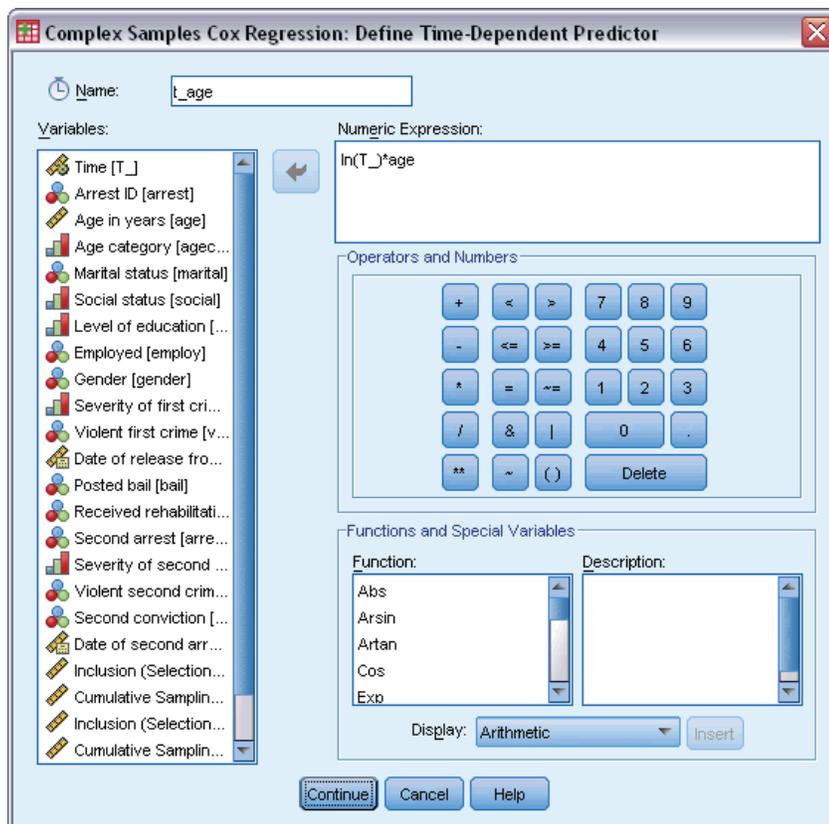
The significance value for the overall test of proportional hazards is less than 0.05, indicating that the proportional hazards assumption is violated. The log time function is used for the alternative model, so it will be easy to replicate this time-dependent predictor.

### Adding a Time-Dependent Predictor

- Recall the Complex Samples Cox Regression dialog box and click the Predictors tab.

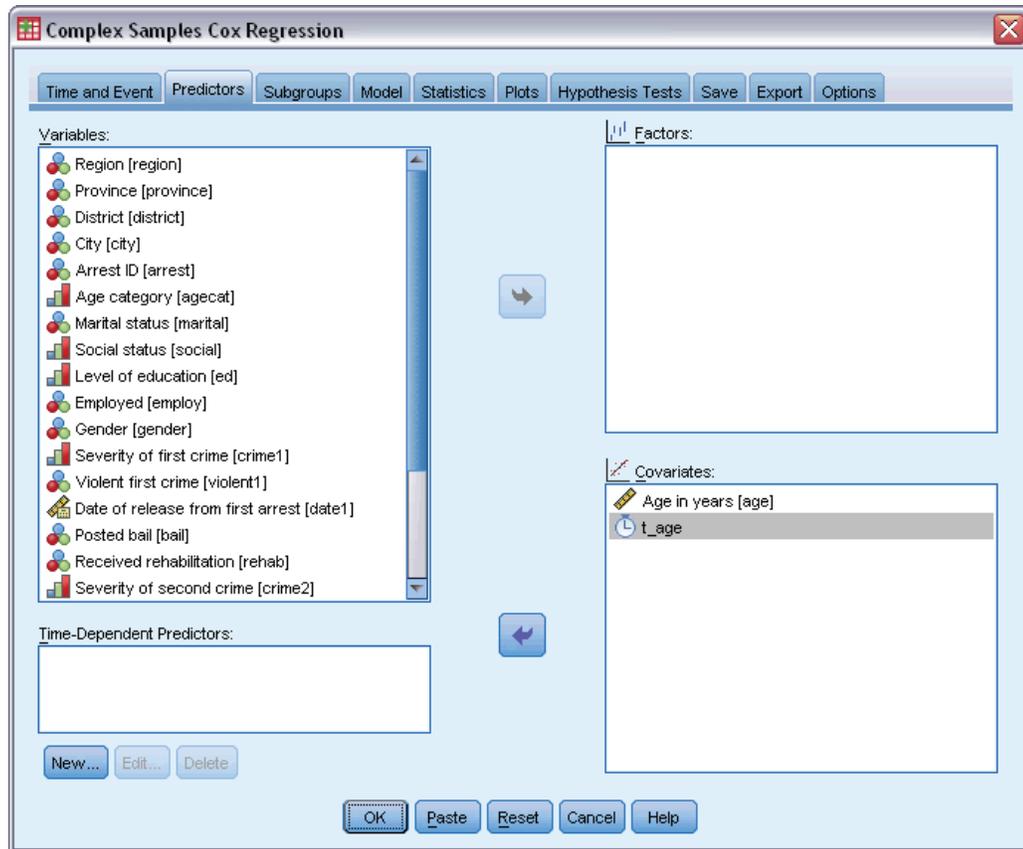
- ▶ Click New.

Figure 22-16  
Cox Regression Define Time-Dependent Predictor dialog box



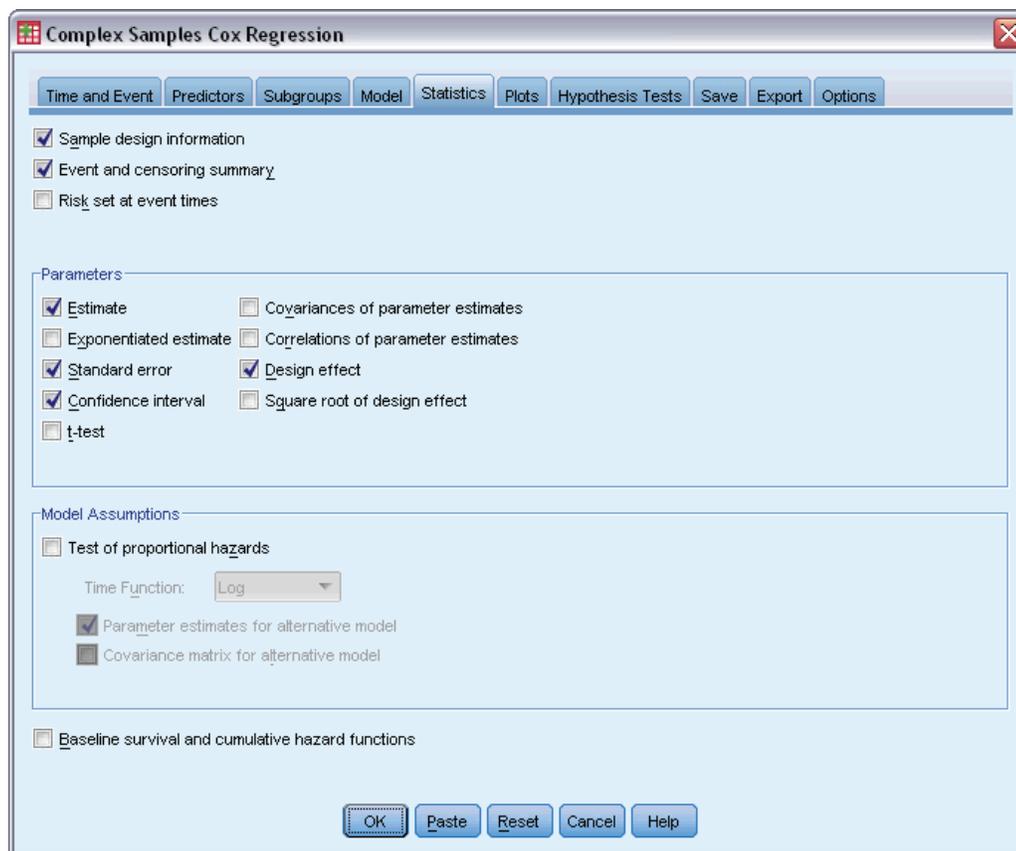
- ▶ Type `t_age` as the name of the time-dependent predictor you want to define.
- ▶ Type `ln(T_)*age` as the numeric expression.
- ▶ Click Continue.

Figure 22-17  
Cox Regression dialog box, Predictors tab



- ▶ Select  $t\_age$  as a covariate.
- ▶ Click the Statistics tab.

Figure 22-18  
Cox Regression dialog box, Predictors tab



- ▶ Select Estimate, Standard error, Confidence interval, and Design effect in the Parameters group.
- ▶ Deselect Test of proportional hazards and Parameter estimates for alternative model in the Model Assumptions group.
- ▶ Click OK.

### Tests of Model Effects

Figure 22-19  
Tests of model effects

Source	df1	df2	Wald F	Sig.
age	1.000	16.000	.015	0.91
t_age	1.000	16.000	29.924	5.136E-5

Survival Time Variable: Time to second arrest  
Event Status Variable: Second arrest = 1  
Model: age, t\_age

With the addition of the time-dependent predictor, the significance value for *age* is 0.91, indicating that its contribution to the model is superseded by that of *t\_age*.

### Parameter Estimates

Figure 22-20  
Parameter estimates

Parameter	B	Std. Error	95% Confidence Interval		Design Effect
			Lower	Upper	
age	-.002	0.01	-.030	.027	.702
t_age	-.012	.002	-.017	-.008	.666

Survival Time Variable: Time to second arrest  
Event Status Variable: Second arrest = 1  
Model: age, t\_age

Looking at the parameter estimates and standard errors, you can see that you have replicated the alternative model from the test of proportional hazards. By explicitly specifying the model, you can request additional parameter statistics and plots. Here we have requested the design effect; the value for *t\_age* of less than 1 indicates that the standard error for *t\_age* is smaller than what you would obtain if you assumed that the dataset was a simple random sample. In this case, the effect of *t\_age* would still be statistically significant, but the confidence intervals would be wider.

## Multiple Cases per Subject in Complex Samples Cox Regression

Researchers investigating survival times for patients exiting a rehabilitation program post-ischemic stroke face a number of challenges.

**Multiple cases per subject.** Variables representing patient medical history should be useful as predictors. Over time, patients may experience major medical events that alter their medical history. In this dataset, the occurrence of myocardial infarction, ischemic stroke, or hemorrhagic stroke is noted and the time of the event recorded. You could create computable time-dependent covariates within the procedure to include this information in the model, but it should be more convenient to use multiple cases per subject. Note that the variables were originally coded so that the patient history is recorded across variables, so you will need to restructure the dataset.

**Left-truncation.** The onset of risk starts at the time of the ischemic stroke. However, the sample only includes patients who have survived the rehabilitation program, thus the sample is left-truncated in the sense that the observed survival times are “inflated” by the length of rehabilitation. You can account for this by specifying the time at which they exited rehabilitation as the time of entry into the study.

**No sampling plan.** The dataset was not collected via a complex sampling plan and is considered to be a simple random sample. You will need to create an analysis plan to use Complex Samples Cox Regression.

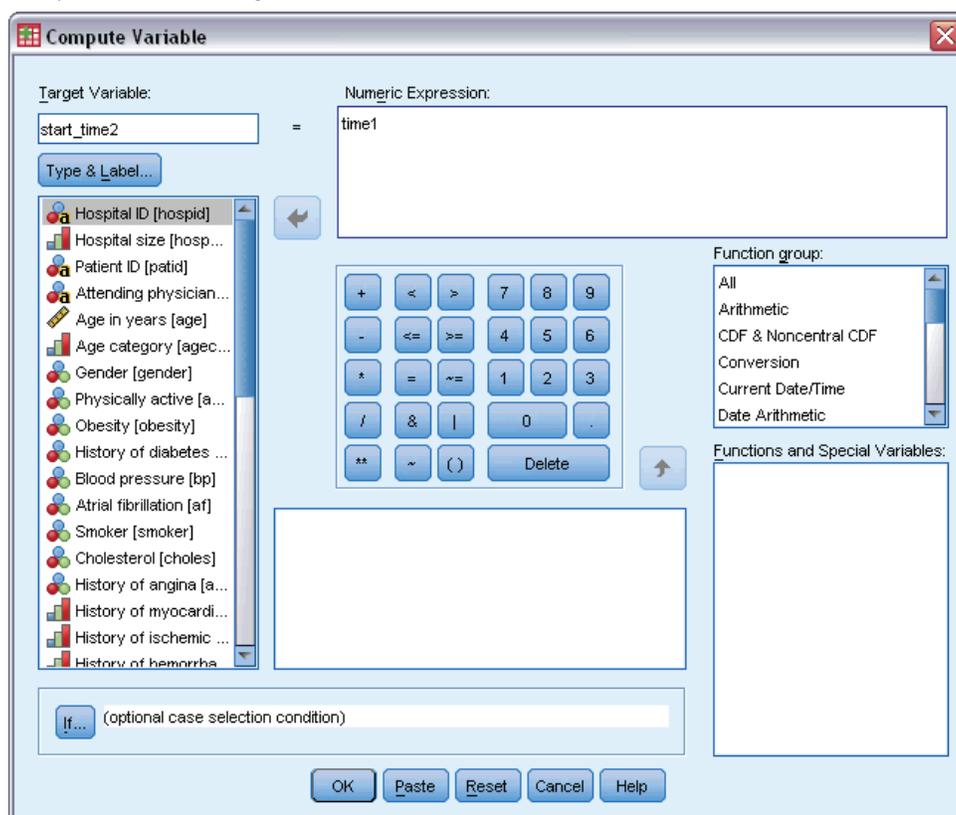
The dataset is collected in *stroke\_survival.sav*. For more information, see the topic [Sample Files in Appendix A in IBM SPSS Complex Samples 19](#). Use the Restructure Data Wizard to prepare the data for analysis, then the Analysis Preparation Wizard to create a simple random sampling plan, and finally Complex Samples Cox Regression to build a model for survival times.

## Preparing the Data for Analysis

Before restructuring the data, you will need to create two ancillary variables to help with the restructuring.

- ▶ To compute a new variable, from the menus choose:  
Transform > Compute Variable...

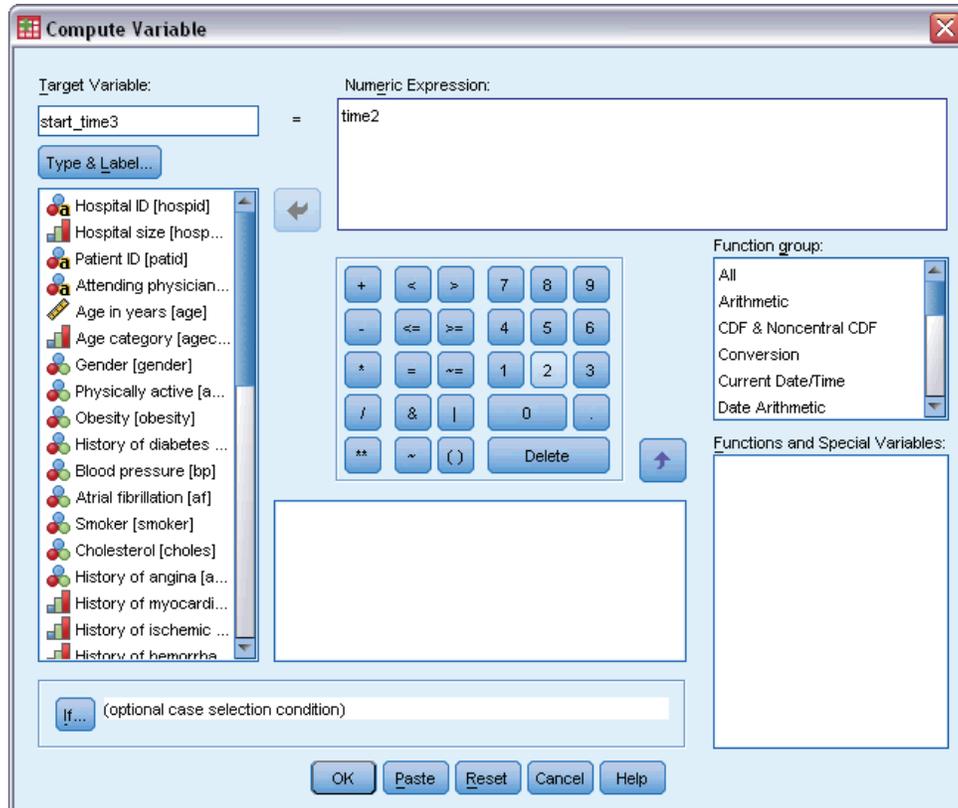
Figure 22-21  
Compute Variable dialog box



- ▶ Type `start_time2` as the target variable.
- ▶ Type `time1` as the numeric expression.
- ▶ Click OK.

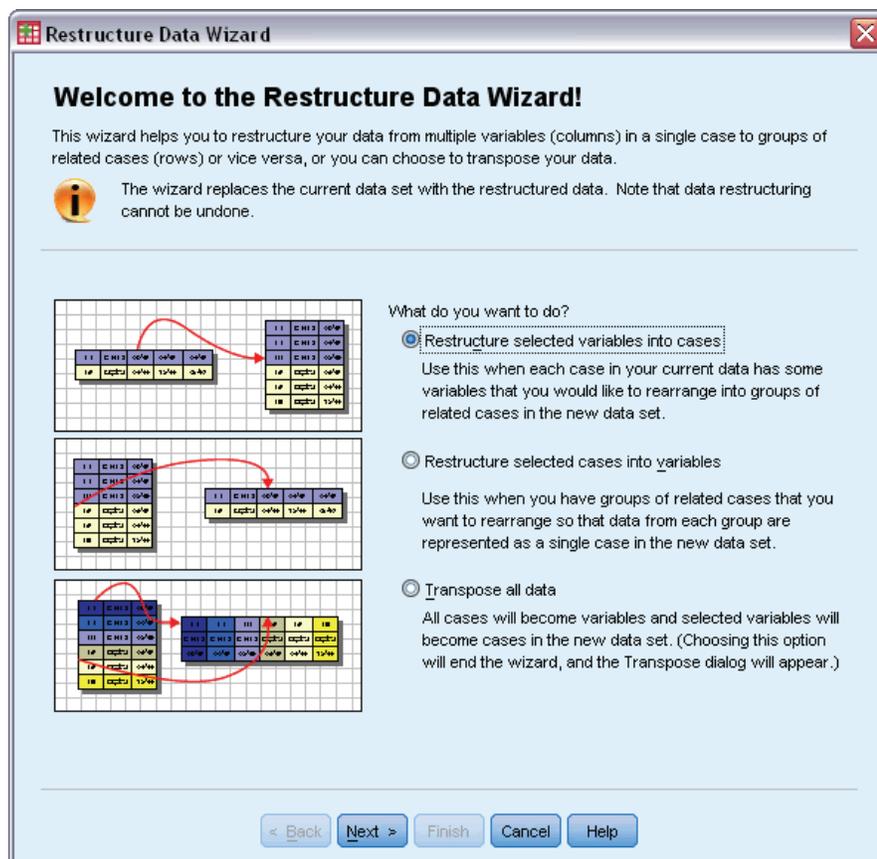
- ▶ Recall the Compute Variable dialog box.

Figure 22-22  
Compute Variable dialog box



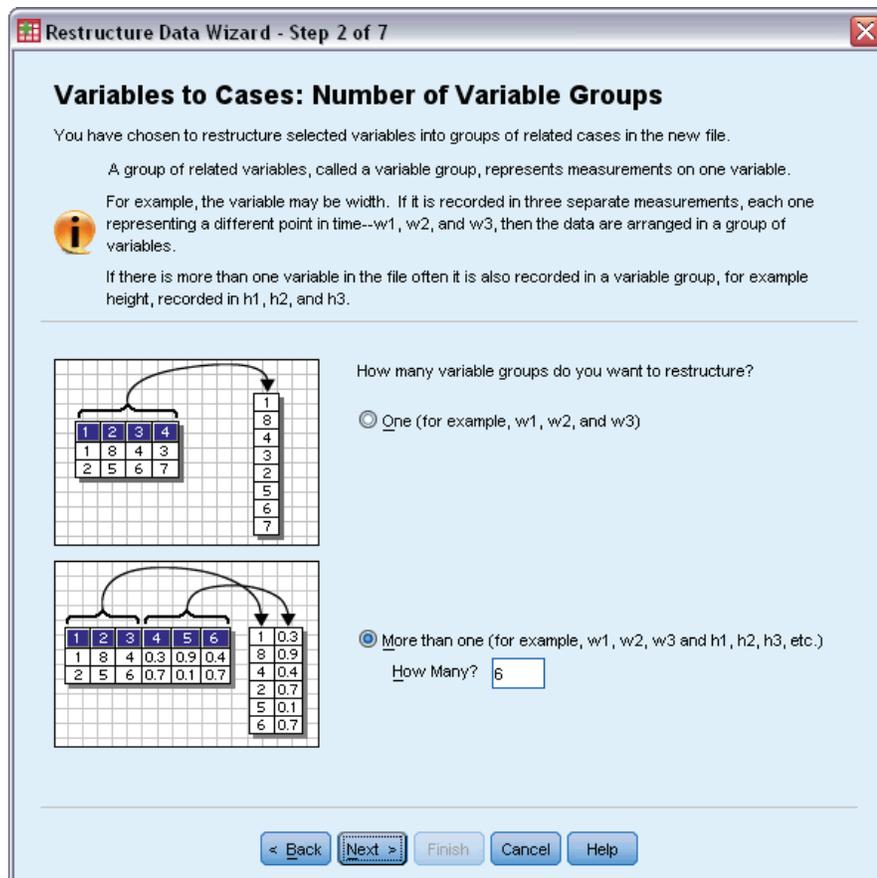
- ▶ Type `start_time3` as the target variable.
- ▶ Type `time2` as the numeric expression.
- ▶ Click OK.
- ▶ To restructure the data from variables to cases, from the menus choose:  
Data > Restructure...

Figure 22-23  
Restructure Data Wizard, Welcome step



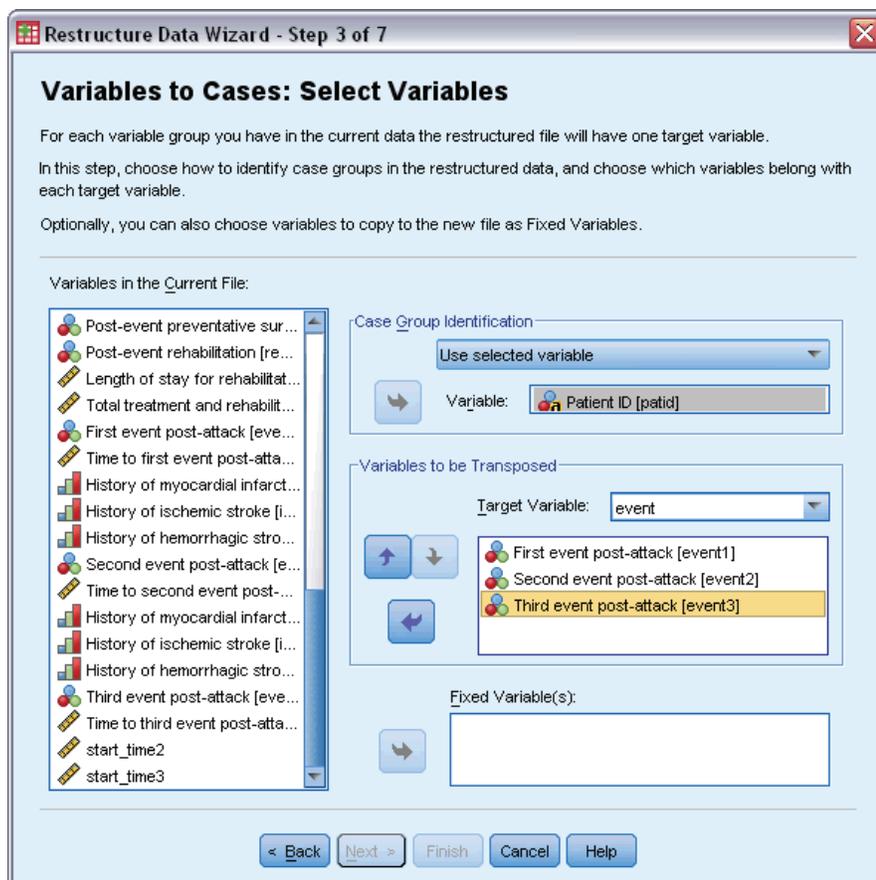
- ▶ Make sure Restructure selected variables into cases is selected.
- ▶ Click Next.

Figure 22-24  
Restructure Data Wizard, Variables to Cases Number of Variable Groups step



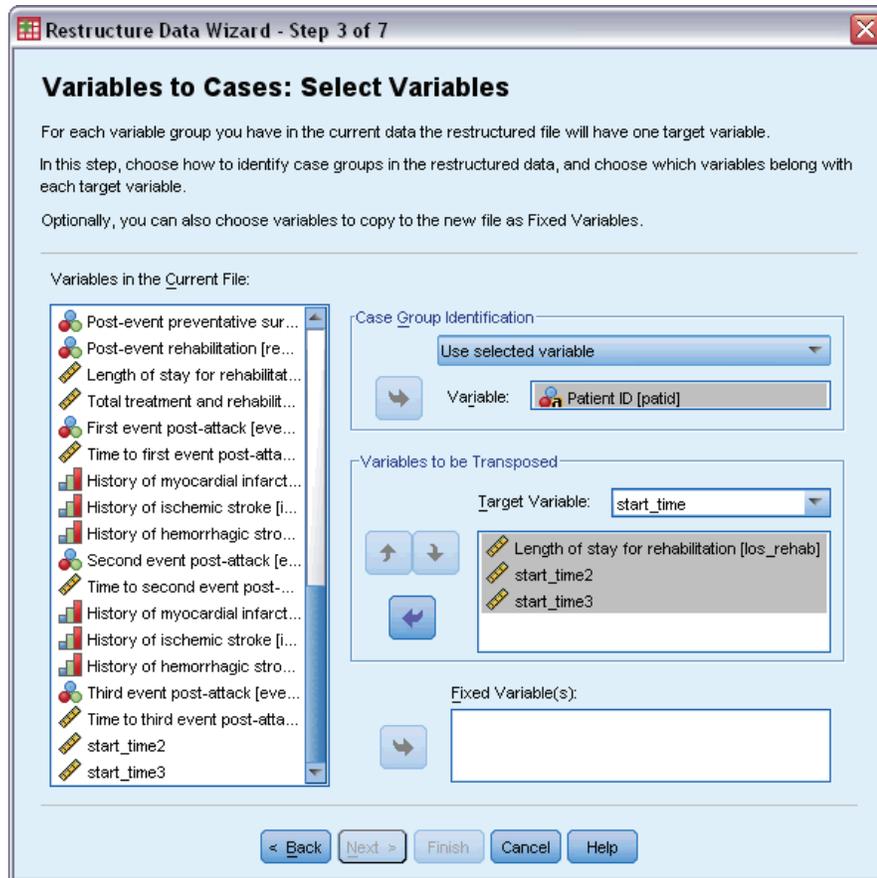
- ▶ Select More than one variable group to restructure.
- ▶ Type 6 as the number of groups.
- ▶ Click Next.

Figure 22-25  
Restructure Data Wizard, Variables to Cases Select Variables step



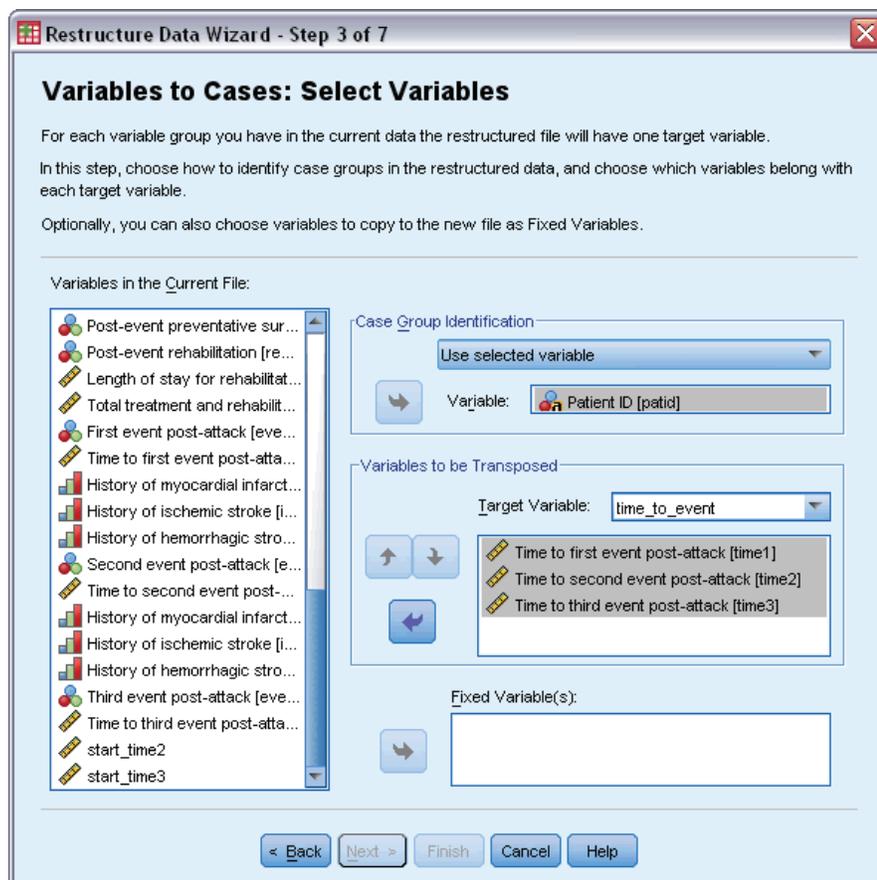
- ▶ In the Case Group Identification group, select Use selected variable and select *Patient ID [patid]* as the subject identifier.
- ▶ Type event as the first target variable.
- ▶ Select *First event post-attack [event1]*, *Second event post-attack [event2]*, and *Third event post-attack [event3]* as variables to be transposed.
- ▶ Select *trans2* from the target variable list.

Figure 22-26  
Restructure Data Wizard, Variables to Cases Select Variables step



- ▶ Type `start_time` as the target variable.
- ▶ Select *Length of stay for rehabilitation [los\_rehab]*, *start\_time2*, and *start\_time3* as variables to be transposed. *Time to first event post-attack [time1]* and *Time to second event post-attack [time2]* will be used to create the end times, and each variable can only appear in one list of variables to be transposed, thus *start\_time2* and *start\_time3* were necessary.
- ▶ Select *trans3* from the target variable list.

Figure 22-27  
Restructure Data Wizard, Variables to Cases Select Variables step



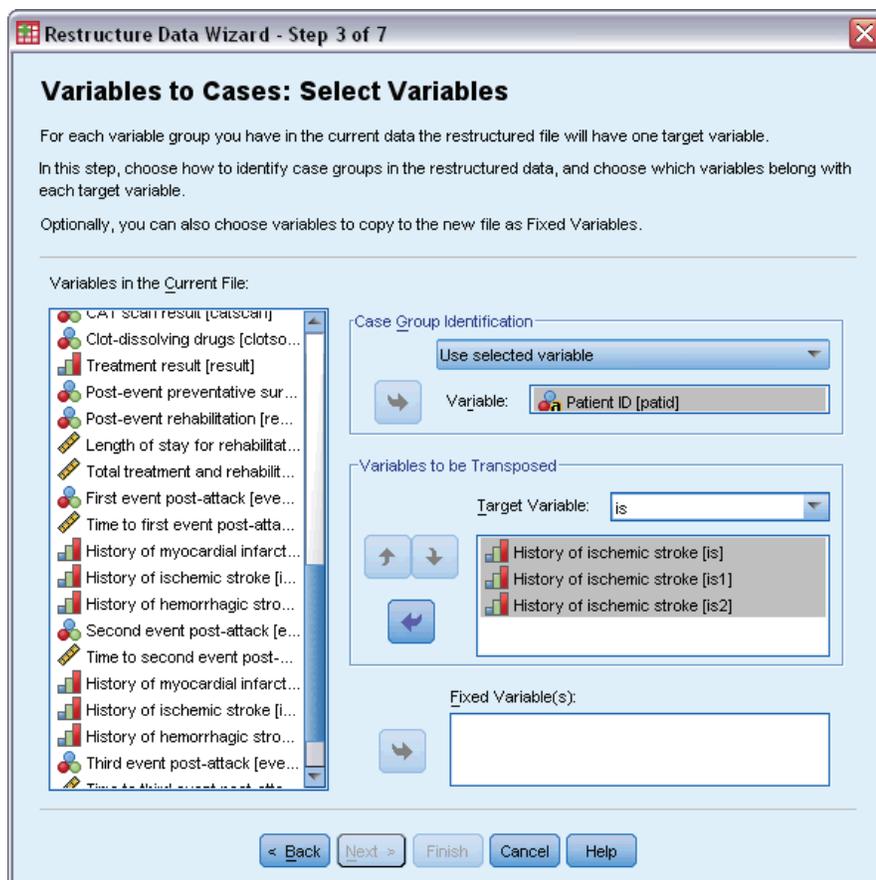
- ▶ Type `time_to_event` as the target variable.
- ▶ Select *Time to first event post-attack [time1]*, *Time to second event post-attack [time2]*, and *Time to third event post-attack [time3]* as variables to be transposed.
- ▶ Select *trans4* from the target variable list.

Figure 22-28  
Restructure Data Wizard, Variables to Cases Select Variables step



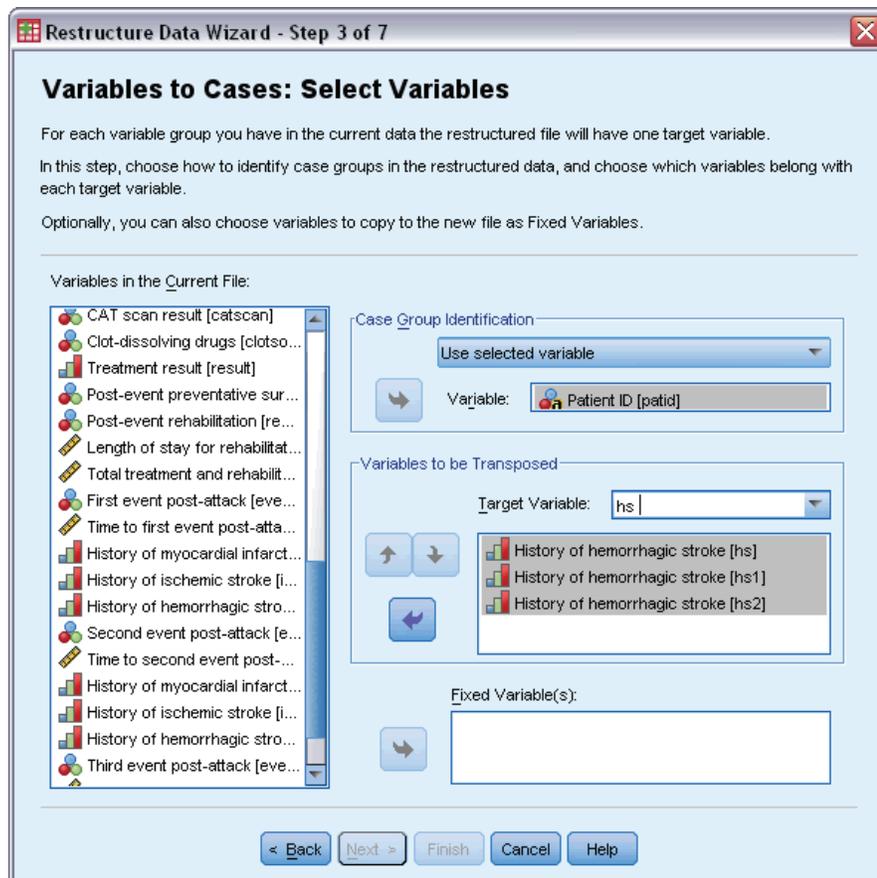
- ▶ Type *mi* as the target variable.
- ▶ Select *History of myocardial infarction [mi]*, *History of myocardial infarction [mi1]*, and *History of myocardial infarction [mi2]* as variables to be transposed.
- ▶ Select *trans5* from the target variable list.

Figure 22-29  
Restructure Data Wizard, Variables to Cases Select Variables step



- ▶ Type *is* as the target variable.
- ▶ Select *History of ischemic stroke [is]*, *History of ischemic stroke [is1]*, and *History of ischemic stroke [is2]* as variables to be transposed.
- ▶ Select *trans6* from the target variable list.

Figure 22-30  
Restructure Data Wizard, Variables to Cases Select Variables step



- ▶ Type *hs* as the target variable.
- ▶ Select *History of hemorrhagic stroke [hs]*, *History of hemorrhagic stroke [hs1]*, and *History of hemorrhagic stroke [hs2]* as variables to be transposed.
- ▶ Click Next, then click Next in the Create Index Variables step.

Figure 22-31  
Restructure Data Wizard, Variables to Cases Create One Index Variable step

**Variables to Cases: Create One Index Variable**

You have chosen to create one index variable. The variable's values can be sequential numbers or the names of variables in a group.

In the table you can specify the name and label for the index variable.

What kind of index values?

Sequential numbers  
Index Values: 1, 2, 3

Variable names  
Index Values: event1, event2, event3

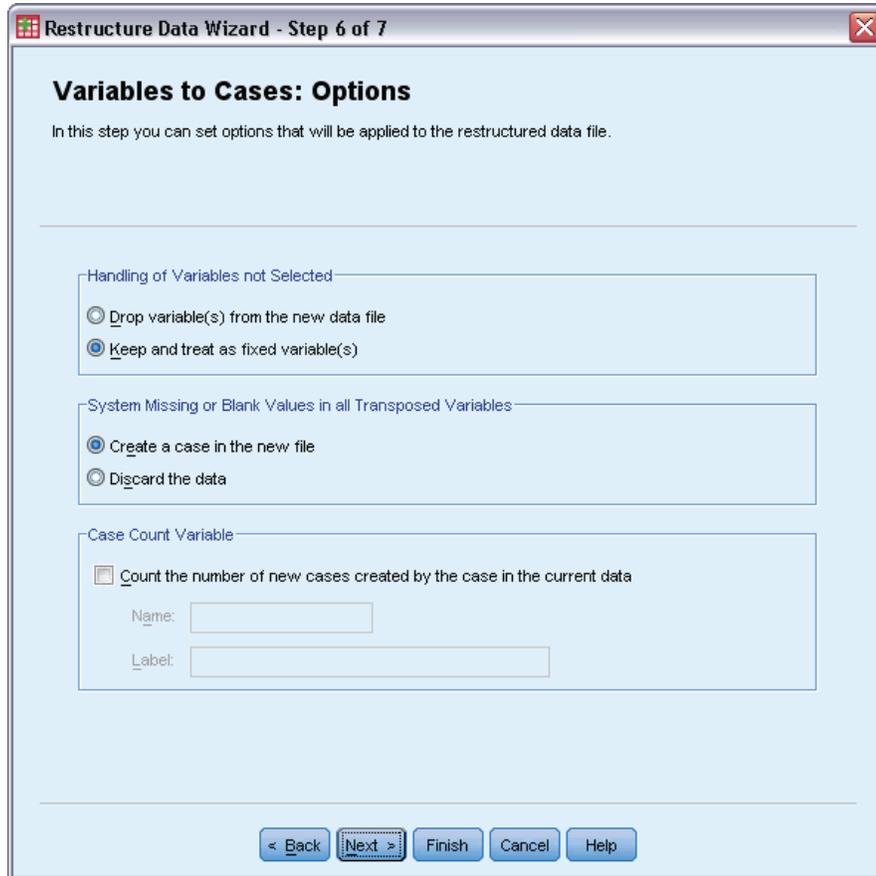
Edit the Index Variable Name and Label:

	Name	Label	Levels	Index Values
1	event_index	Event Index	3	1, 2, 3

< Back Next > Finish Cancel Help

- ▶ Type event\_index as the name of the index variable and type Event index as the variable label.
- ▶ Click Next.

Figure 22-32  
Restructure Data Wizard, Variables to Cases Create One Index Variable step



- ▶ Make sure Keep and treat as fixed variable(s) is selected.
- ▶ Click Finish.

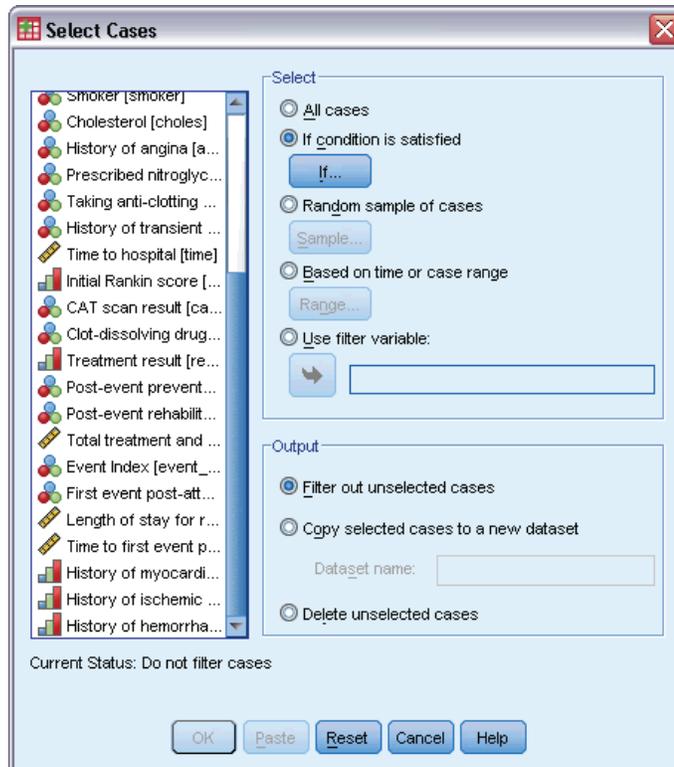
Figure 22-33  
Restructured data

event_index	event	start_time	time_to_event	mi	is	hs
1	0	3	1500	0	1	0
2	-4	1500	-4	-4	-4	-4
3	-4	.	-4	-4	-4	-4
1	1	33	1311	0	1	0
2	4	1311	1325	1	1	0
3	-3	1325	-3	-3	-3	-3
1	4	12	1098	1	1	0
2	-3	1098	-3	-3	-3	-3
3	-3	.	-3	-3	-3	-3
1	4	4	1356	0	1	0
2	-3	1356	-3	-3	-3	-3
3	-3	.	-3	-3	-3	-3

The restructured data contains three cases for every patient; however, many patients experienced fewer than three events, so there are many cases with negative (missing) values for *event*. You can simply filter these from the dataset.

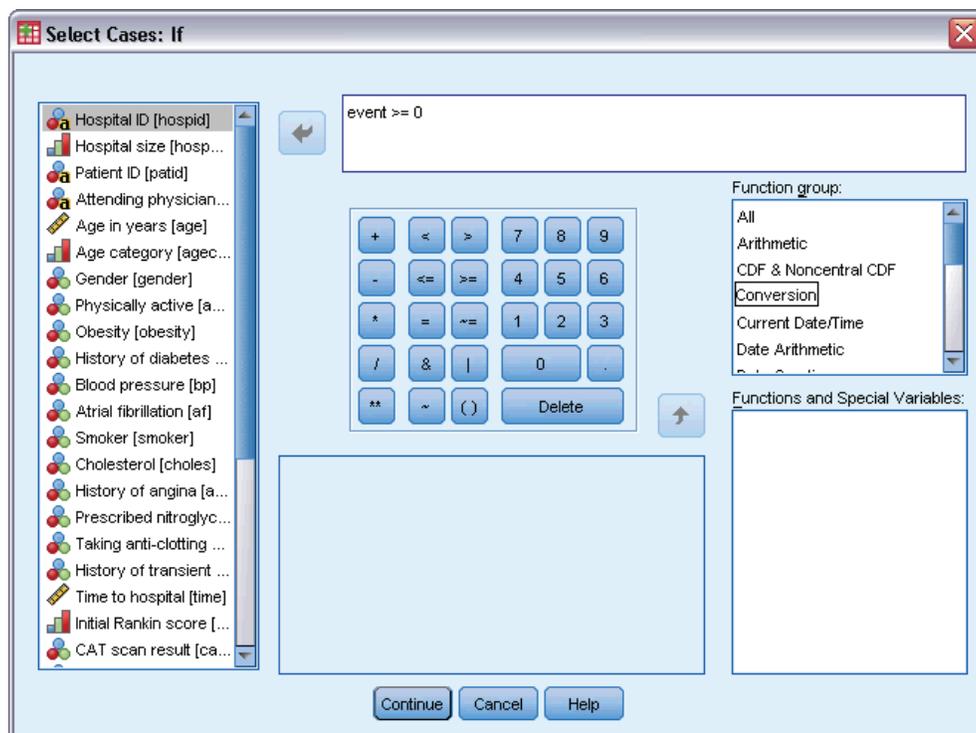
- ▶ To filter these cases, from the menus choose:  
Data > Select Cases...

Figure 22-34  
Select Cases dialog box



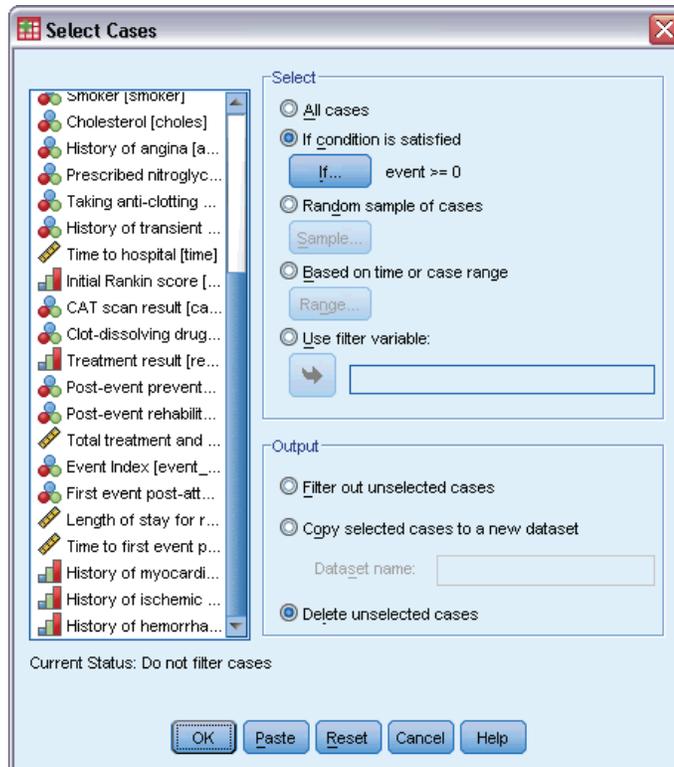
- ▶ Select If condition is satisfied.
- ▶ Click If.

Figure 22-35  
Select Cases If dialog box



- ▶ Type event >= 0 as the conditional expression.
- ▶ Click Continue.

Figure 22-36  
Select Cases dialog box



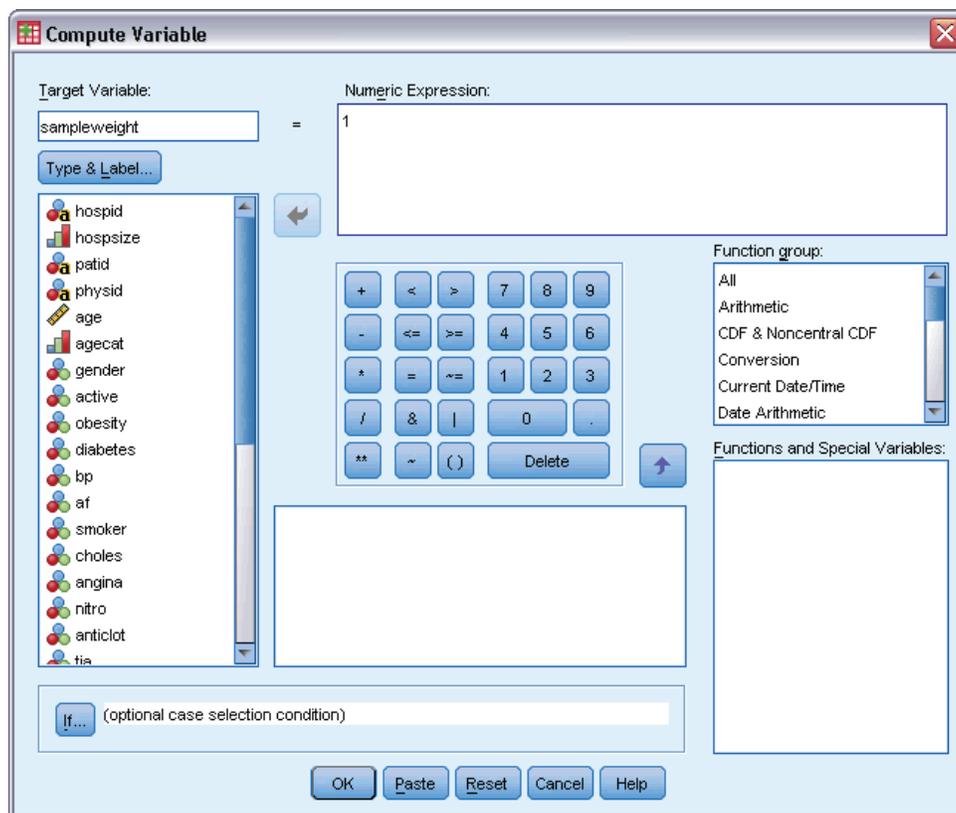
- ▶ Select Delete unselected cases.
- ▶ Click OK.

### ***Creating a Simple Random Sampling Analysis Plan***

Now you are ready to create the simple random sampling analysis plan.

- ▶ First, you need to create a sampling weight variable. From the menus choose: Transform > Compute Variable...

Figure 22-37  
Cox Regression main dialog box



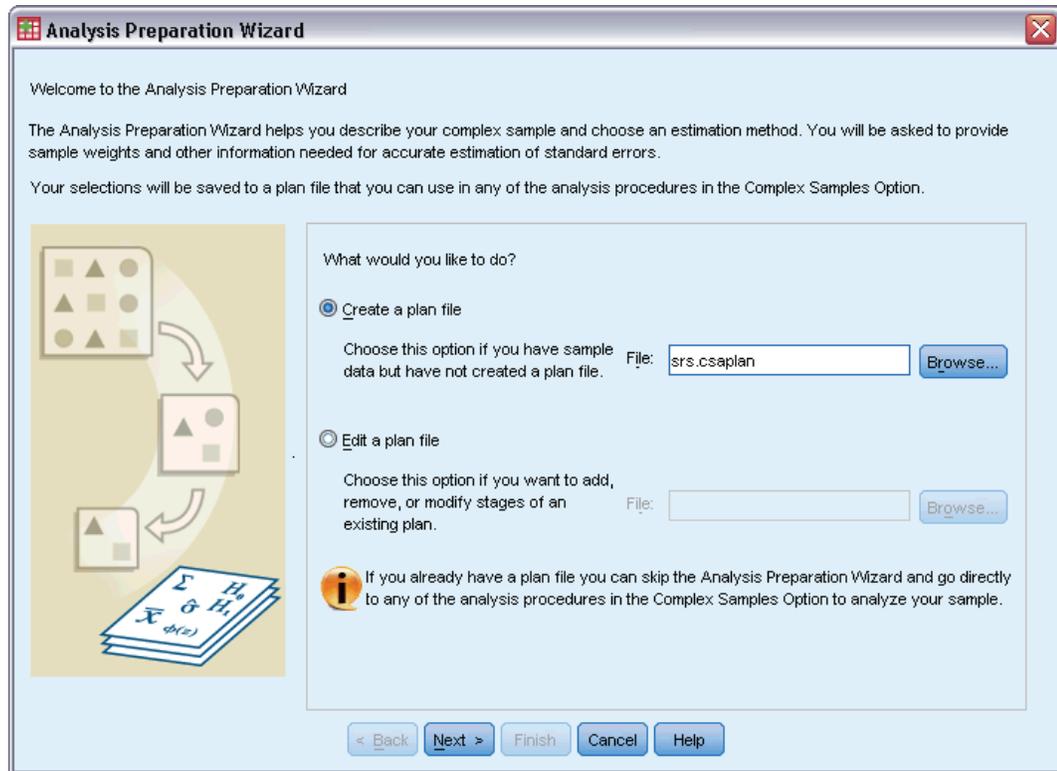
- ▶ Type sampleweight as the target variable.
- ▶ Type 1 as the numeric expression.
- ▶ Click OK.

You are now ready to create the analysis plan.

*Note:* There is an existing plan file, *srs.csaplan*, in the sample files directory that you can use if you want to skip the following instructions and proceed to analysis of the data.

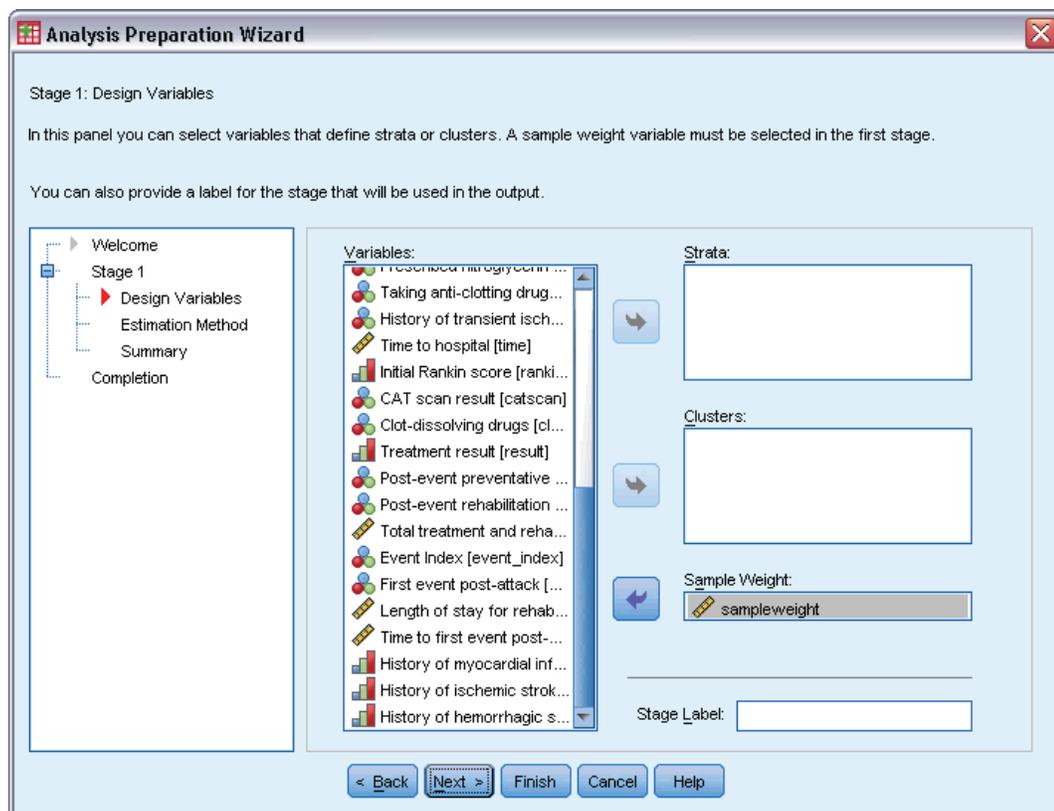
- ▶ To create the analysis plan, from the menus choose:  
Analyze > Complex Samples > Prepare for Analysis...

Figure 22-38  
Analysis Preparation Wizard, Welcome step



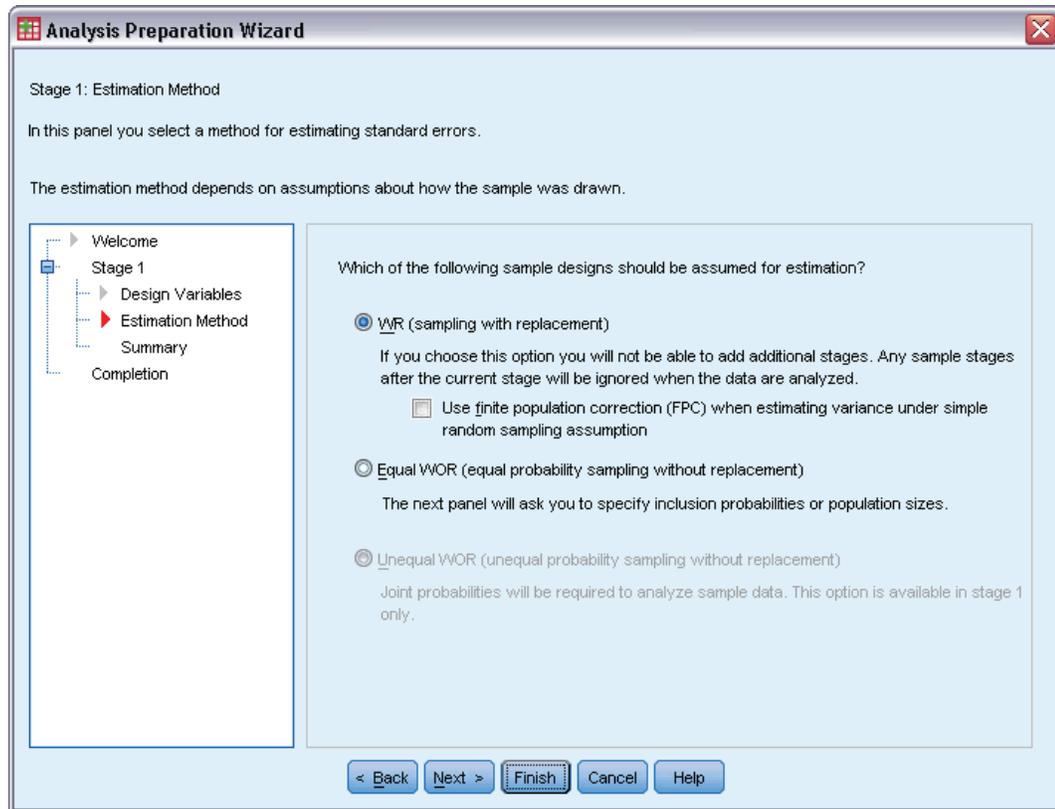
- ▶ Select Create a plan file and type srs.csaplan as the name of the file. Alternatively, browse to the location you want to save it.
- ▶ Click Next.

Figure 22-39  
Analysis Preparation Wizard, Design Variables



- ▶ Select *sampleweight* as the sample weight variable.
- ▶ Click Next.

Figure 22-40  
Analysis Preparation Wizard, Estimation Method



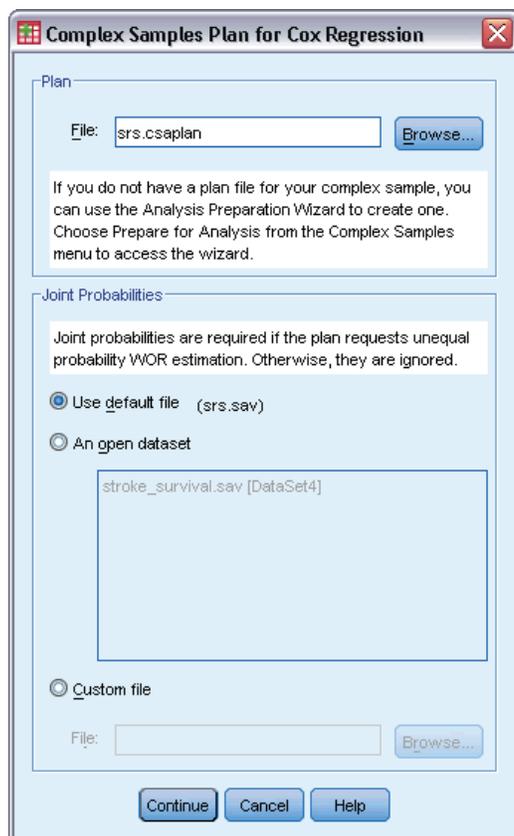
- ▶ Deselect Use finite population correction.
- ▶ Click Finish.

You are now ready to run the analysis.

### **Running the Analysis**

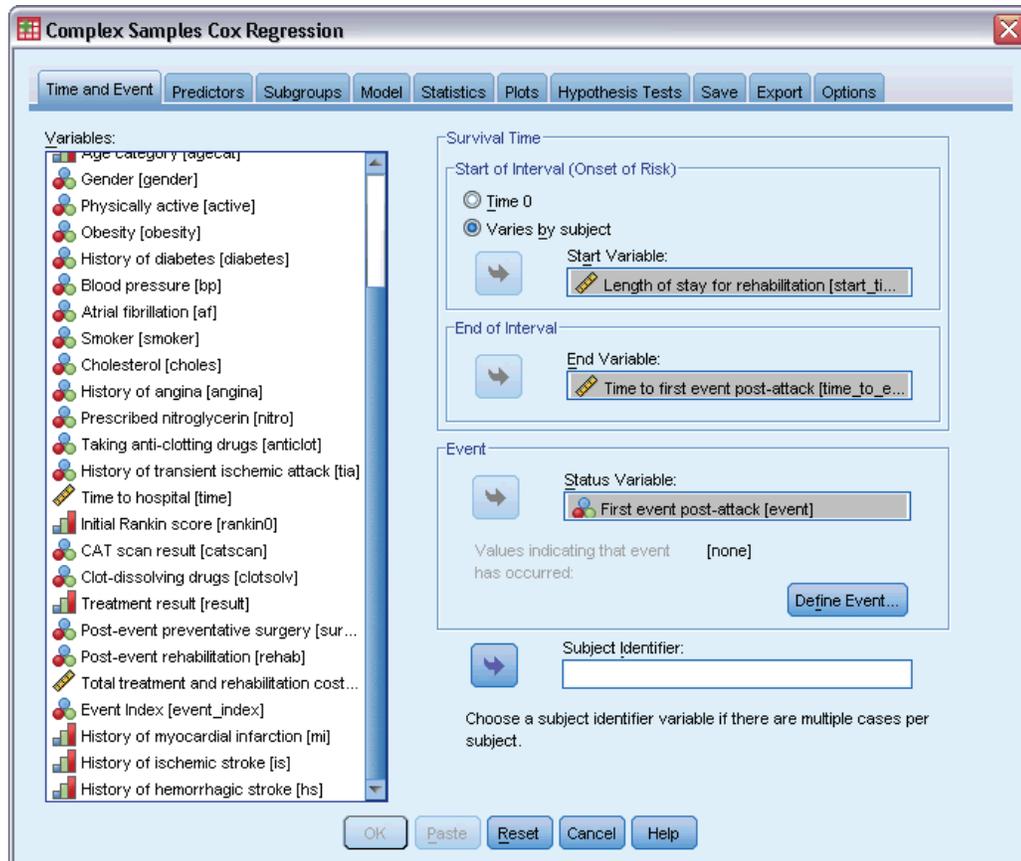
- ▶ To run a Complex Samples Cox Regression analysis, from the menus choose: Analyze > Complex Samples > Cox Regression...

Figure 22-41  
Plan for Cox Regression dialog box



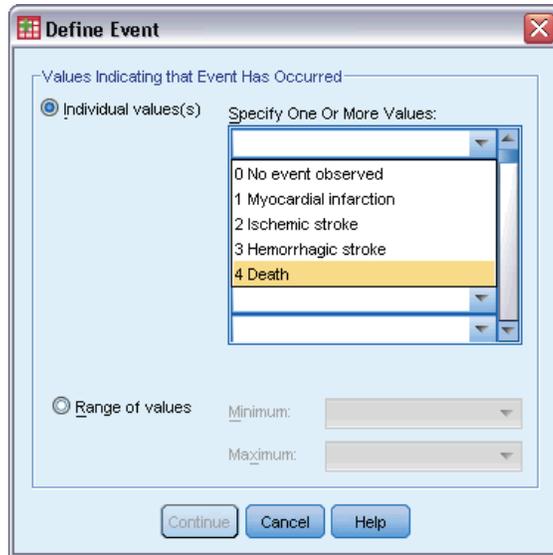
- ▶ Browse to where you saved the simple random sampling analysis plan, or to the sample files directory, and select *srs.csaplan*.
- ▶ Click Continue.

Figure 22-42  
Cox Regression dialog box, Time and Event tab



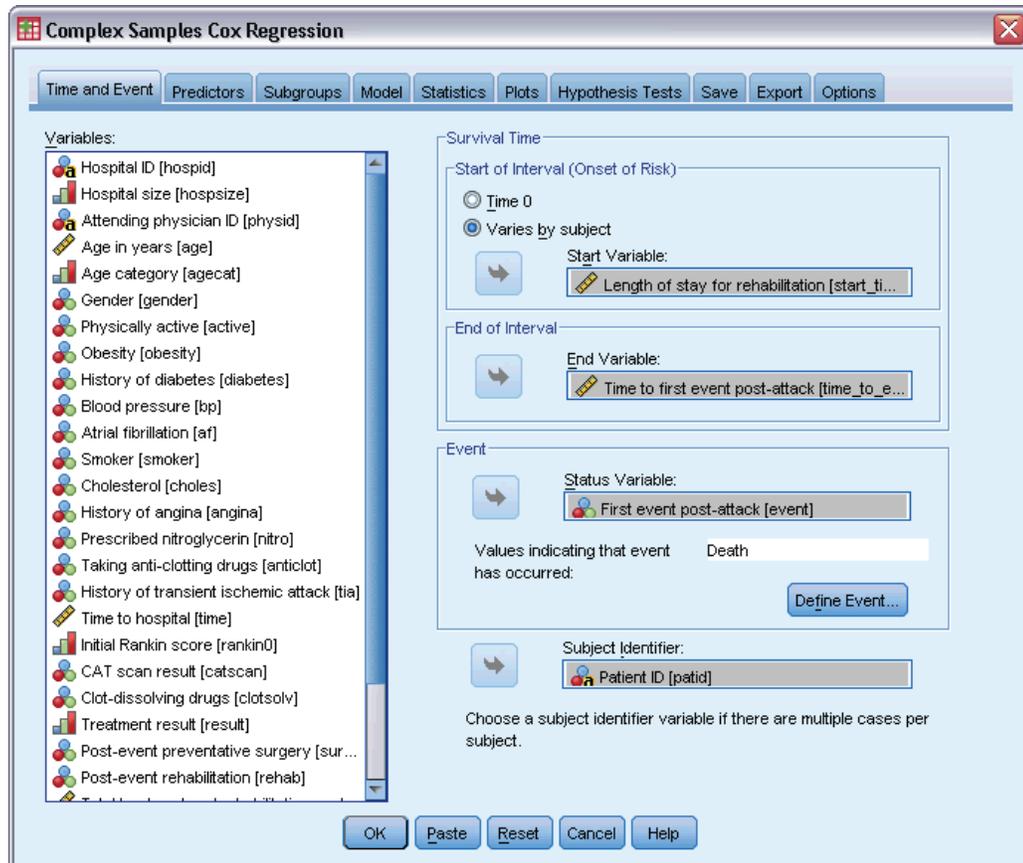
- ▶ Select *Varies by subject* and select *Length of stay for rehabilitation [los\_rehab]* as the start variable. Note that the restructured variable has taken the variable label from the first variable used to construct it, though the label is not necessarily appropriate for the constructed variable.
- ▶ Select *Time to first event post-attack [time\_to\_event]* as the end variable.
- ▶ Select *First event post-attack [event]* as the status variable.
- ▶ Click *Define Event*.

Figure 22-43  
Define Event dialog box



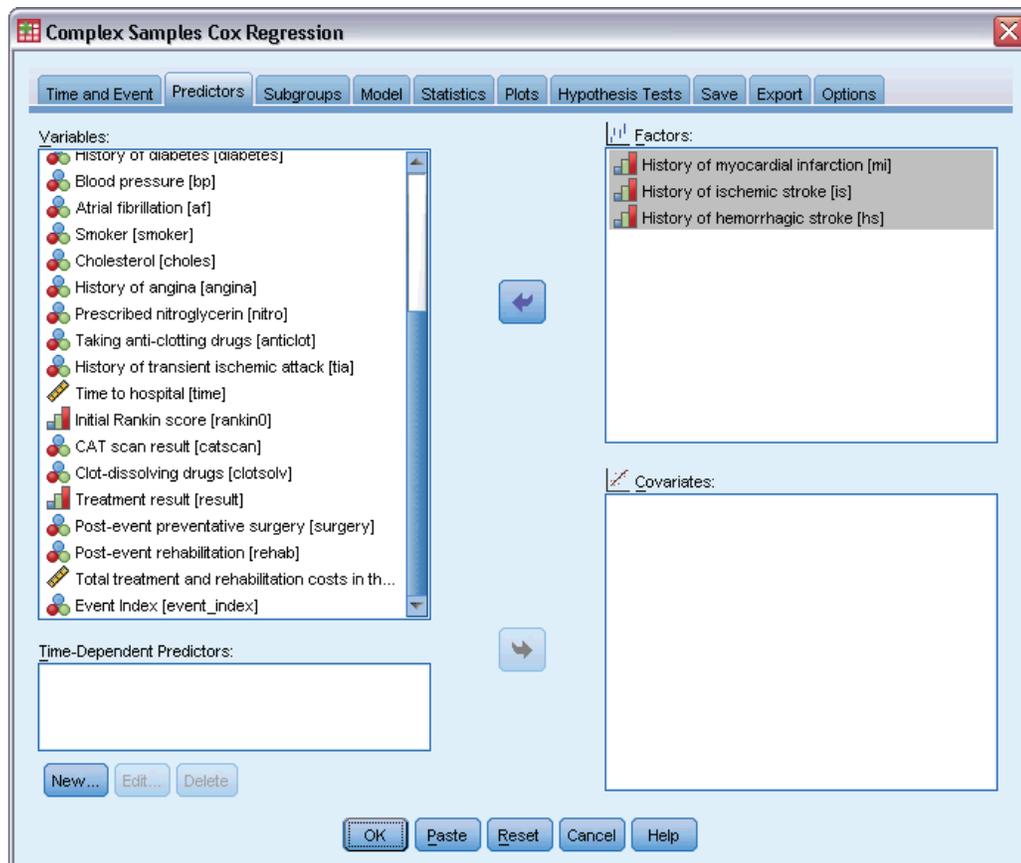
- ▶ Select 4 Death as the value indicating the terminal event has occurred.
- ▶ Click Continue.

Figure 22-44  
Cox Regression dialog box, Time and Event tab



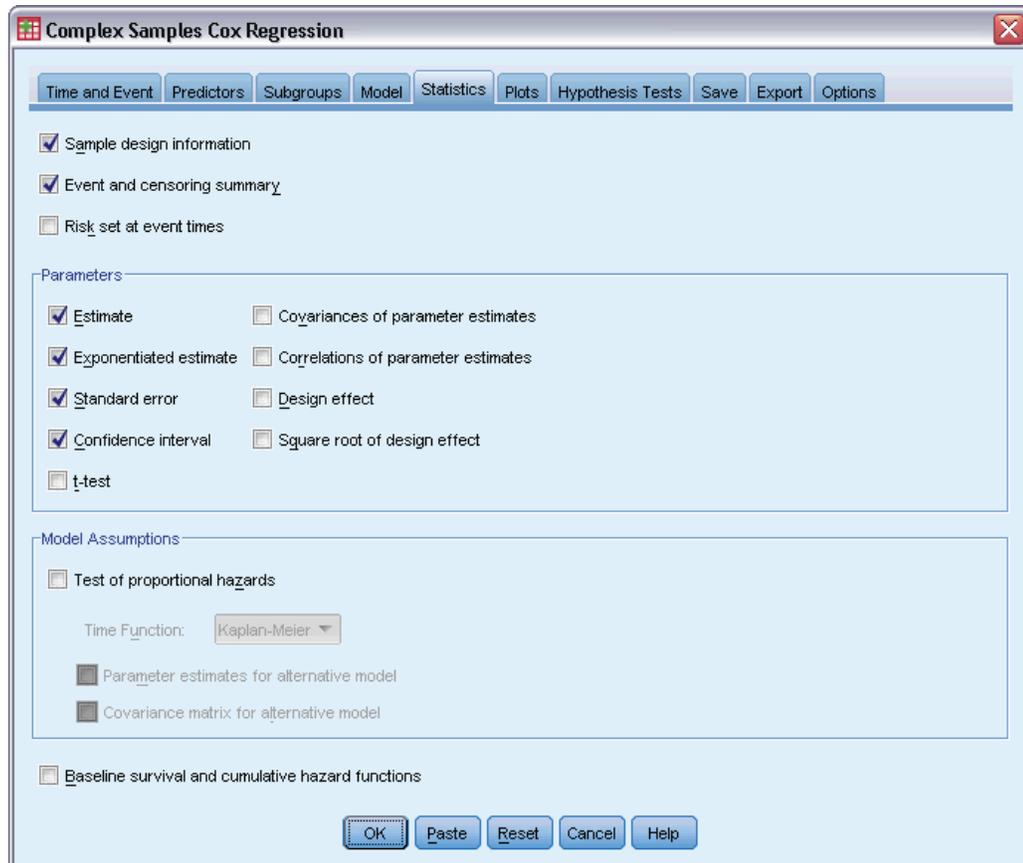
- ▶ Select *Patient ID [patid]* as the subject identifier.
- ▶ Click the Predictors tab.

Figure 22-45  
Cox Regression dialog box, Predictors tab



- ▶ Select *History of myocardial infarction [mi]* through *History of hemorrhagic stroke [hs]* as factors.
- ▶ Click the Statistics tab.

Figure 22-46  
Cox Regression dialog box, Statistics tab



- ▶ Select Estimate, Exponentiated estimate, Standard error, and Confidence interval in the Parameters group.
- ▶ Click the Plots tab.

Figure 22-47  
Cox Regression dialog box, Statistics tab

Complex Samples Cox Regression

Time and Event Predictors Subgroups Model Statistics **Plots** Hypothesis Tests Save Export Options

Plots

Survival function  Log-minus-log of survival function

Hazard function  One minus survival function

Display confidence intervals in selected plots

Plot Factors at:

Factor	Level	Separate Lines
History of myocardial infarction	(Highest level)	<input checked="" type="checkbox"/>
History of ischemic stroke	1.0	<input type="checkbox"/>
History of hemorrhagic stroke	0.0	<input type="checkbox"/>

Plot Covariates at:

Covariate	Value

By default, covariates in the model are evaluated at their means, and factors in the model are evaluated at their highest levels. You can change the value at which any model predictor is evaluated and plot separate lines for each level of one factor variable.

OK Paste Reset Cancel Help

- ▶ Select Log-minus-log of survival function.
- ▶ Check Separate Lines for *History of myocardial infarction*.
- ▶ Select 1.0 as the level for *History of ischemic stroke*.
- ▶ Select 0.0 as the level for *History of hemorrhagic stroke*.
- ▶ Click the Options tab.

Figure 22-48  
Cox Regression dialog box, Options tab

- ▶ Select Breslow as the tie-breaking method in the Estimation group.
- ▶ Click OK.

## Sample Design Information

Figure 22-49  
Sample design information

			N
Unweighted Counts	Valid	Subjects	2421
		Cases	3310
	Invalid Cases		0
	Total Cases		3310
Population Subject Size			2421.000
Stage 1	Strata	1	
	Units	2421	
Sampling Design Degrees of Freedom			2420

This table contains information on the sample design pertinent to the estimation of the model.

- There are multiple cases for some subjects, and all 3,310 cases are used in the analysis.
- The design has a single stratum and 2,421 units (one for each subject). The sampling design degrees of freedom are estimated by  $2421-1=2420$ .

## Tests of Model Effects

Figure 22-50  
Tests of model effects

Source	df1	df2	Wald F	Sig.
mi	3.000	2418.000	452.873	.000
is	2.000	2419.000	1064.936	.000
hs	2.000	2419.000	739.197	.000

Survival Time Variable: Length of stay for rehabilitation, Time to first event post-attack  
Event Status Variable: First event post-attack = 4  
Subject ID Variable: Patient ID  
Model: mi, is, hs

The significance value for each effect is near 0, suggesting that they all contribute to the model.

## Parameter Estimates

Figure 22-51  
Parameter estimates

Parameter	B	Std. Error	95% Confidence Interval		Exp(B)	95% Confidence Interval for Exp(B)	
			Lower	Upper		Lower	Upper
[mi=0]	-6.381	.283	-6.935	-5.827	.002	.001	.003
[mi=1]	-5.589	.284	-6.147	-5.032	.004	.002	.007
[mi=2]	-2.119	.344	-2.794	-1.445	.120	.061	.236
[mi=3]	.000 <sup>a</sup>	.	.	.	1.000	.	.
[is=1]	-6.421	.202	-6.817	-6.024	.002	.001	.002
[is=2]	-2.803	.222	-3.239	-2.366	.061	.039	.094
[is=3]	.000 <sup>a</sup>	.	.	.	1.000	.	.
[hs=0]	-6.148	.355	-6.844	-5.453	0	.001	.004
[hs=1]	-2.232	.373	-2.963	-1.502	.107	.052	.223
[hs=2]	.000 <sup>a</sup>	.	.	.	1.000	.	.

Survival Time Variable: Length of stay for rehabilitation, Time to first event post-attack  
Event Status Variable: First event post-attack = 4  
Subject ID Variable: Patient ID  
Model: mi, is, hs

a. Set to zero because this parameter is redundant.

b. Tie breaking method: Breslow

The procedure uses the last category of each factor as the reference category; the effect of other categories is relative to the reference category. Note that while the estimate is useful for statistical testing, the exponentiated estimate,  $\text{Exp}(B)$ , is more easily interpreted as the predicted change in the hazard relative to the reference category.

- The value of  $\text{Exp}(B)$  for  $[mi=0]$  means that the hazard of death for a patient with no prior myocardial infarctions (mi) is 0.002 times that of a patient with three prior mi's.

- The confidence intervals for  $[mi=0]$  and  $[mi=1]$  do not overlap with the interval for  $[mi=2]$ , and none of them include 0. Therefore, it appears that the hazard for patients with one or no prior mi's is distinguishable from the hazard for patients with two prior mi's, which in turn is distinguishable from the hazard for patients with three prior mi's.

Similar relationships hold for the levels of *is* and *hs*, where increasing the number of prior incidents increases the hazard of death.

## Pattern Values

Figure 22-52  
Pattern values

		Survival Time Interval		History of myocardial infarction	History of ischemic stroke	History of hemorrhagic stroke
		Start	End			
Reference Pattern	1	.000	a	Three	Three	Two
Pattern 1.1	1	.000	a	None	One	None
Pattern 1.2	1	.000	a	One	One	None
Pattern 1.3	1	.000	a	Two	One	None
Pattern 1.4	1	.000	a	Three	One	None

Unspecified predictor is assigned the value of this predictor at the reference pattern.  
Each Survival Time Interval is defined as Start < Survival Time <= End.  
Model: mi, is, hs.

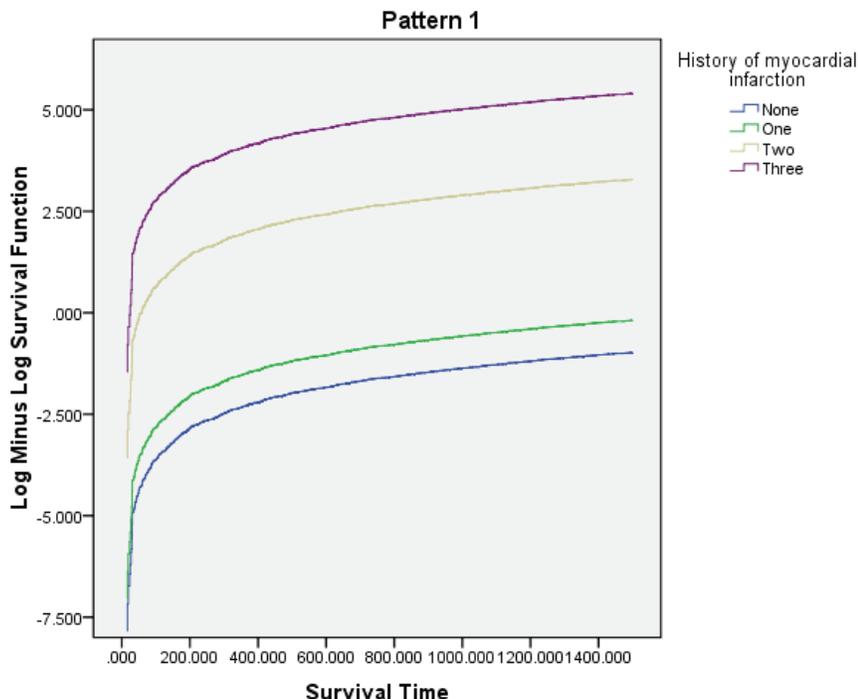
a. Unbounded

The pattern values table lists the values that define each predictor pattern. In addition to the predictors in the model, the start and end times for the survival interval are displayed. For analyses run from the dialogs, the start and end times will always be 0 and unbounded, respectively; through syntax you can specify piecewise constant predictor paths.

- The reference pattern is set at the reference category for each factor and the mean value of each covariate (there are no covariates in this model). For this dataset, the combination of factors shown for the reference model cannot occur, so we will ignore the log-minus-log plot for the reference pattern.
- Patterns 1.1 through 1.4 differ only on the value of *History of myocardial infarction*. A separate pattern (and separate line in the requested plot) is created for each value of *History of myocardial infarction* while the other variables are held constant.

## Log-Minus-Log Plot

Figure 22-53  
Log-minus-log plot



This plot displays the log-minus-log of the survival function,  $\ln(-\ln(\text{survival}))$ , versus the survival time. This particular plot displays a separate curve for each category of *History of myocardial infarction*, with *History of ischemic stroke* fixed at *One* and *History of hemorrhagic stroke* fixed at *None*, and is a useful visualization of the effect of *History of myocardial infarction* on the survival function. As seen in the parameter estimates table, it appears that the survival for patients with one or no prior mi's is distinguishable from the survival for patients with two prior mi's, which in turn is distinguishable from the survival for patients with three prior mi's.

## Summary

You have fit a Cox regression model for post-stroke survival that estimates the effects of the changing post-stroke patient history. This is just a beginning, as researchers would undoubtedly want to include other potential predictors in the model. Moreover, in further analysis of this dataset you might consider more significant changes to the model structure. For example, the current model assumes that the effect of a patient history-altering event can be quantified by a multiplier to the baseline hazard. Instead, it may be reasonable to assume that the shape of the baseline hazard is altered by the occurrence of a nondeath event. To accomplish this, you could stratify the analysis based on *Event index*.

# Sample Files

The sample files installed with the product can be found in the *Samples* subdirectory of the installation directory. There is a separate folder within the *Samples* subdirectory for each of the following languages: English, French, German, Italian, Japanese, Korean, Polish, Russian, Simplified Chinese, Spanish, and Traditional Chinese.

Not all sample files are available in all languages. If a sample file is not available in a language, that language folder contains an English version of the sample file.

## **Descriptions**

Following are brief descriptions of the sample files used in various examples throughout the documentation.

- **accidents.sav.** This is a hypothetical data file that concerns an insurance company that is studying age and gender risk factors for automobile accidents in a given region. Each case corresponds to a cross-classification of age category and gender.
- **adl.sav.** This is a hypothetical data file that concerns efforts to determine the benefits of a proposed type of therapy for stroke patients. Physicians randomly assigned female stroke patients to one of two groups. The first received the standard physical therapy, and the second received an additional emotional therapy. Three months following the treatments, each patient's abilities to perform common activities of daily life were scored as ordinal variables.
- **advert.sav.** This is a hypothetical data file that concerns a retailer's efforts to examine the relationship between money spent on advertising and the resulting sales. To this end, they have collected past sales figures and the associated advertising costs.
- **aflatoxin.sav.** This is a hypothetical data file that concerns the testing of corn crops for aflatoxin, a poison whose concentration varies widely between and within crop yields. A grain processor has received 16 samples from each of 8 crop yields and measured the aflatoxin levels in parts per billion (PPB).
- **anorectic.sav.** While working toward a standardized symptomatology of anorectic/bulimic behavior, researchers (Van der Ham, Meulman, Van Strien, and Van Engeland, 1997) made a study of 55 adolescents with known eating disorders. Each patient was seen four times over four years, for a total of 220 observations. At each observation, the patients were scored for each of 16 symptoms. Symptom scores are missing for patient 71 at time 2, patient 76 at time 2, and patient 47 at time 3, leaving 217 valid observations.
- **bankloan.sav.** This is a hypothetical data file that concerns a bank's efforts to reduce the rate of loan defaults. The file contains financial and demographic information on 850 past and prospective customers. The first 700 cases are customers who were previously given

loans. The last 150 cases are prospective customers that the bank needs to classify as good or bad credit risks.

- **bankloan\_binning.sav.** This is a hypothetical data file containing financial and demographic information on 5,000 past customers.
- **behavior.sav.** In a classic example (Price and Bouffard, 1974), 52 students were asked to rate the combinations of 15 situations and 15 behaviors on a 10-point scale ranging from 0=“extremely appropriate” to 9=“extremely inappropriate.” Averaged over individuals, the values are taken as dissimilarities.
- **behavior\_ini.sav.** This data file contains an initial configuration for a two-dimensional solution for *behavior.sav*.
- **brakes.sav.** This is a hypothetical data file that concerns quality control at a factory that produces disc brakes for high-performance automobiles. The data file contains diameter measurements of 16 discs from each of 8 production machines. The target diameter for the brakes is 322 millimeters.
- **breakfast.sav.** In a classic study (Green and Rao, 1972), 21 Wharton School MBA students and their spouses were asked to rank 15 breakfast items in order of preference with 1=“most preferred” to 15=“least preferred.” Their preferences were recorded under six different scenarios, from “Overall preference” to “Snack, with beverage only.”
- **breakfast-overall.sav.** This data file contains the breakfast item preferences for the first scenario, “Overall preference,” only.
- **broadband\_1.sav.** This is a hypothetical data file containing the number of subscribers, by region, to a national broadband service. The data file contains monthly subscriber numbers for 85 regions over a four-year period.
- **broadband\_2.sav.** This data file is identical to *broadband\_1.sav* but contains data for three additional months.
- **car\_insurance\_claims.sav.** A dataset presented and analyzed elsewhere (McCullagh and Nelder, 1989) concerns damage claims for cars. The average claim amount can be modeled as having a gamma distribution, using an inverse link function to relate the mean of the dependent variable to a linear combination of the policyholder age, vehicle type, and vehicle age. The number of claims filed can be used as a scaling weight.
- **car\_sales.sav.** This data file contains hypothetical sales estimates, list prices, and physical specifications for various makes and models of vehicles. The list prices and physical specifications were obtained alternately from *edmunds.com* and manufacturer sites.
- **car\_sales\_uprepared.sav.** This is a modified version of *car\_sales.sav* that does not include any transformed versions of the fields.
- **carpet.sav.** In a popular example (Green and Wind, 1973), a company interested in marketing a new carpet cleaner wants to examine the influence of five factors on consumer preference—package design, brand name, price, a *Good Housekeeping* seal, and a money-back guarantee. There are three factor levels for package design, each one differing in the location of the applicator brush; three brand names (*K2R*, *Glory*, and *Bissell*); three price levels; and two levels (either no or yes) for each of the last two factors. Ten consumers rank 22 profiles defined by these factors. The variable *Preference* contains the rank of the average rankings for each profile. Low rankings correspond to high preference. This variable reflects an overall measure of preference for each profile.

- **carpet\_prefs.sav.** This data file is based on the same example as described for *carpet.sav*, but it contains the actual rankings collected from each of the 10 consumers. The consumers were asked to rank the 22 product profiles from the most to the least preferred. The variables *PREF1* through *PREF22* contain the identifiers of the associated profiles, as defined in *carpet\_plan.sav*.
- **catalog.sav.** This data file contains hypothetical monthly sales figures for three products sold by a catalog company. Data for five possible predictor variables are also included.
- **catalog\_seasfac.sav.** This data file is the same as *catalog.sav* except for the addition of a set of seasonal factors calculated from the Seasonal Decomposition procedure along with the accompanying date variables.
- **cellular.sav.** This is a hypothetical data file that concerns a cellular phone company's efforts to reduce churn. Churn propensity scores are applied to accounts, ranging from 0 to 100. Accounts scoring 50 or above may be looking to change providers.
- **ceramics.sav.** This is a hypothetical data file that concerns a manufacturer's efforts to determine whether a new premium alloy has a greater heat resistance than a standard alloy. Each case represents a separate test of one of the alloys; the heat at which the bearing failed is recorded.
- **cereal.sav.** This is a hypothetical data file that concerns a poll of 880 people about their breakfast preferences, also noting their age, gender, marital status, and whether or not they have an active lifestyle (based on whether they exercise at least twice a week). Each case represents a separate respondent.
- **clothing\_defects.sav.** This is a hypothetical data file that concerns the quality control process at a clothing factory. From each lot produced at the factory, the inspectors take a sample of clothes and count the number of clothes that are unacceptable.
- **coffee.sav.** This data file pertains to perceived images of six iced-coffee brands (Kennedy, Riquier, and Sharp, 1996). For each of 23 iced-coffee image attributes, people selected all brands that were described by the attribute. The six brands are denoted AA, BB, CC, DD, EE, and FF to preserve confidentiality.
- **contacts.sav.** This is a hypothetical data file that concerns the contact lists for a group of corporate computer sales representatives. Each contact is categorized by the department of the company in which they work and their company ranks. Also recorded are the amount of the last sale made, the time since the last sale, and the size of the contact's company.
- **creditpromo.sav.** This is a hypothetical data file that concerns a department store's efforts to evaluate the effectiveness of a recent credit card promotion. To this end, 500 cardholders were randomly selected. Half received an ad promoting a reduced interest rate on purchases made over the next three months. Half received a standard seasonal ad.
- **customer\_dbase.sav.** This is a hypothetical data file that concerns a company's efforts to use the information in its data warehouse to make special offers to customers who are most likely to reply. A subset of the customer base was selected at random and given the special offers, and their responses were recorded.
- **customer\_information.sav.** A hypothetical data file containing customer mailing information, such as name and address.
- **customer\_subset.sav.** A subset of 80 cases from *customer\_dbase.sav*.

- **debate.sav.** This is a hypothetical data file that concerns paired responses to a survey from attendees of a political debate before and after the debate. Each case corresponds to a separate respondent.
- **debate\_aggregate.sav.** This is a hypothetical data file that aggregates the responses in *debate.sav*. Each case corresponds to a cross-classification of preference before and after the debate.
- **demo.sav.** This is a hypothetical data file that concerns a purchased customer database, for the purpose of mailing monthly offers. Whether or not the customer responded to the offer is recorded, along with various demographic information.
- **demo\_cs\_1.sav.** This is a hypothetical data file that concerns the first step of a company's efforts to compile a database of survey information. Each case corresponds to a different city, and the region, province, district, and city identification are recorded.
- **demo\_cs\_2.sav.** This is a hypothetical data file that concerns the second step of a company's efforts to compile a database of survey information. Each case corresponds to a different household unit from cities selected in the first step, and the region, province, district, city, subdivision, and unit identification are recorded. The sampling information from the first two stages of the design is also included.
- **demo\_cs.sav.** This is a hypothetical data file that contains survey information collected using a complex sampling design. Each case corresponds to a different household unit, and various demographic and sampling information is recorded.
- **dmdata.sav.** This is a hypothetical data file that contains demographic and purchasing information for a direct marketing company. *dmdata2.sav* contains information for a subset of contacts that received a test mailing, and *dmdata3.sav* contains information on the remaining contacts who did not receive the test mailing.
- **dietstudy.sav.** This hypothetical data file contains the results of a study of the "Stillman diet" (Rickman, Mitchell, Dingman, and Dalen, 1974). Each case corresponds to a separate subject and records his or her pre- and post-diet weights in pounds and triglyceride levels in mg/100 ml.
- **dvdplayer.sav.** This is a hypothetical data file that concerns the development of a new DVD player. Using a prototype, the marketing team has collected focus group data. Each case corresponds to a separate surveyed user and records some demographic information about them and their responses to questions about the prototype.
- **german\_credit.sav.** This data file is taken from the "German credit" dataset in the Repository of Machine Learning Databases (Blake and Merz, 1998) at the University of California, Irvine.
- **grocery\_1month.sav.** This hypothetical data file is the *grocery\_coupons.sav* data file with the weekly purchases "rolled-up" so that each case corresponds to a separate customer. Some of the variables that changed weekly disappear as a result, and the amount spent recorded is now the sum of the amounts spent during the four weeks of the study.
- **grocery\_coupons.sav.** This is a hypothetical data file that contains survey data collected by a grocery store chain interested in the purchasing habits of their customers. Each customer is followed for four weeks, and each case corresponds to a separate customer-week and records information about where and how the customer shops, including how much was spent on groceries during that week.

- **guttman.sav.** Bell (Bell, 1961) presented a table to illustrate possible social groups. Guttman (Guttman, 1968) used a portion of this table, in which five variables describing such things as social interaction, feelings of belonging to a group, physical proximity of members, and formality of the relationship were crossed with seven theoretical social groups, including crowds (for example, people at a football game), audiences (for example, people at a theater or classroom lecture), public (for example, newspaper or television audiences), mobs (like a crowd but with much more intense interaction), primary groups (intimate), secondary groups (voluntary), and the modern community (loose confederation resulting from close physical proximity and a need for specialized services).
- **health\_funding.sav.** This is a hypothetical data file that contains data on health care funding (amount per 100 population), disease rates (rate per 10,000 population), and visits to health care providers (rate per 10,000 population). Each case represents a different city.
- **hivassay.sav.** This is a hypothetical data file that concerns the efforts of a pharmaceutical lab to develop a rapid assay for detecting HIV infection. The results of the assay are eight deepening shades of red, with deeper shades indicating greater likelihood of infection. A laboratory trial was conducted on 2,000 blood samples, half of which were infected with HIV and half of which were clean.
- **hourlywagedata.sav.** This is a hypothetical data file that concerns the hourly wages of nurses from office and hospital positions and with varying levels of experience.
- **insurance\_claims.sav.** This is a hypothetical data file that concerns an insurance company that wants to build a model for flagging suspicious, potentially fraudulent claims. Each case represents a separate claim.
- **insure.sav.** This is a hypothetical data file that concerns an insurance company that is studying the risk factors that indicate whether a client will have to make a claim on a 10-year term life insurance contract. Each case in the data file represents a pair of contracts, one of which recorded a claim and the other didn't, matched on age and gender.
- **judges.sav.** This is a hypothetical data file that concerns the scores given by trained judges (plus one enthusiast) to 300 gymnastics performances. Each row represents a separate performance; the judges viewed the same performances.
- **kinship\_dat.sav.** Rosenberg and Kim (Rosenberg and Kim, 1975) set out to analyze 15 kinship terms (aunt, brother, cousin, daughter, father, granddaughter, grandfather, grandmother, grandson, mother, nephew, niece, sister, son, uncle). They asked four groups of college students (two female, two male) to sort these terms on the basis of similarities. Two groups (one female, one male) were asked to sort twice, with the second sorting based on a different criterion from the first sort. Thus, a total of six "sources" were obtained. Each source corresponds to a  $15 \times 15$  proximity matrix, whose cells are equal to the number of people in a source minus the number of times the objects were partitioned together in that source.
- **kinship\_ini.sav.** This data file contains an initial configuration for a three-dimensional solution for *kinship\_dat.sav*.
- **kinship\_var.sav.** This data file contains independent variables *gender*, *gener(ation)*, and *degree* (of separation) that can be used to interpret the dimensions of a solution for *kinship\_dat.sav*. Specifically, they can be used to restrict the space of the solution to a linear combination of these variables.
- **marketvalues.sav.** This data file concerns home sales in a new housing development in Algonquin, Ill., during the years from 1999–2000. These sales are a matter of public record.

- **nhis2000\_subset.sav.** The National Health Interview Survey (NHIS) is a large, population-based survey of the U.S. civilian population. Interviews are carried out face-to-face in a nationally representative sample of households. Demographic information and observations about health behaviors and status are obtained for members of each household. This data file contains a subset of information from the 2000 survey. National Center for Health Statistics. National Health Interview Survey, 2000. Public-use data file and documentation. [ftp://ftp.cdc.gov/pub/Health\\_Statistics/NCHS/Datasets/NHIS/2000/](ftp://ftp.cdc.gov/pub/Health_Statistics/NCHS/Datasets/NHIS/2000/). Accessed 2003.
- **ozone.sav.** The data include 330 observations on six meteorological variables for predicting ozone concentration from the remaining variables. Previous researchers (Breiman and Friedman, 1985), (Hastie and Tibshirani, 1990), among others found nonlinearities among these variables, which hinder standard regression approaches.
- **pain\_medication.sav.** This hypothetical data file contains the results of a clinical trial for anti-inflammatory medication for treating chronic arthritic pain. Of particular interest is the time it takes for the drug to take effect and how it compares to an existing medication.
- **patient\_los.sav.** This hypothetical data file contains the treatment records of patients who were admitted to the hospital for suspected myocardial infarction (MI, or “heart attack”). Each case corresponds to a separate patient and records many variables related to their hospital stay.
- **patlos\_sample.sav.** This hypothetical data file contains the treatment records of a sample of patients who received thrombolytics during treatment for myocardial infarction (MI, or “heart attack”). Each case corresponds to a separate patient and records many variables related to their hospital stay.
- **poll\_cs.sav.** This is a hypothetical data file that concerns pollsters’ efforts to determine the level of public support for a bill before the legislature. The cases correspond to registered voters. Each case records the county, township, and neighborhood in which the voter lives.
- **poll\_cs\_sample.sav.** This hypothetical data file contains a sample of the voters listed in *poll\_cs.sav*. The sample was taken according to the design specified in the *poll\_csplan* plan file, and this data file records the inclusion probabilities and sample weights. Note, however, that because the sampling plan makes use of a probability-proportional-to-size (PPS) method, there is also a file containing the joint selection probabilities (*poll\_jointprob.sav*). The additional variables corresponding to voter demographics and their opinion on the proposed bill were collected and added to the data file after the sample was taken.
- **property\_assess.sav.** This is a hypothetical data file that concerns a county assessor’s efforts to keep property value assessments up to date on limited resources. The cases correspond to properties sold in the county in the past year. Each case in the data file records the township in which the property lies, the assessor who last visited the property, the time since that assessment, the valuation made at that time, and the sale value of the property.
- **property\_assess\_cs.sav.** This is a hypothetical data file that concerns a state assessor’s efforts to keep property value assessments up to date on limited resources. The cases correspond to properties in the state. Each case in the data file records the county, township, and neighborhood in which the property lies, the time since the last assessment, and the valuation made at that time.
- **property\_assess\_cs\_sample.sav.** This hypothetical data file contains a sample of the properties listed in *property\_assess\_cs.sav*. The sample was taken according to the design specified in the *property\_assess\_csplan* plan file, and this data file records the inclusion probabilities

and sample weights. The additional variable *Current value* was collected and added to the data file after the sample was taken.

- **recidivism.sav.** This is a hypothetical data file that concerns a government law enforcement agency's efforts to understand recidivism rates in their area of jurisdiction. Each case corresponds to a previous offender and records their demographic information, some details of their first crime, and then the time until their second arrest, if it occurred within two years of the first arrest.
- **recidivism\_cs\_sample.sav.** This is a hypothetical data file that concerns a government law enforcement agency's efforts to understand recidivism rates in their area of jurisdiction. Each case corresponds to a previous offender, released from their first arrest during the month of June, 2003, and records their demographic information, some details of their first crime, and the data of their second arrest, if it occurred by the end of June, 2006. Offenders were selected from sampled departments according to the sampling plan specified in *recidivism\_cs.csplan*; because it makes use of a probability-proportional-to-size (PPS) method, there is also a file containing the joint selection probabilities (*recidivism\_cs\_jointprob.sav*).
- **rfm\_transactions.sav.** A hypothetical data file containing purchase transaction data, including date of purchase, item(s) purchased, and monetary amount of each transaction.
- **salesperformance.sav.** This is a hypothetical data file that concerns the evaluation of two new sales training courses. Sixty employees, divided into three groups, all receive standard training. In addition, group 2 gets technical training; group 3, a hands-on tutorial. Each employee was tested at the end of the training course and their score recorded. Each case in the data file represents a separate trainee and records the group to which they were assigned and the score they received on the exam.
- **satisf.sav.** This is a hypothetical data file that concerns a satisfaction survey conducted by a retail company at 4 store locations. 582 customers were surveyed in all, and each case represents the responses from a single customer.
- **screws.sav.** This data file contains information on the characteristics of screws, bolts, nuts, and tacks (Hartigan, 1975).
- **shampoo\_ph.sav.** This is a hypothetical data file that concerns the quality control at a factory for hair products. At regular time intervals, six separate output batches are measured and their pH recorded. The target range is 4.5–5.5.
- **ships.sav.** A dataset presented and analyzed elsewhere (McCullagh et al., 1989) that concerns damage to cargo ships caused by waves. The incident counts can be modeled as occurring at a Poisson rate given the ship type, construction period, and service period. The aggregate months of service for each cell of the table formed by the cross-classification of factors provides values for the exposure to risk.
- **site.sav.** This is a hypothetical data file that concerns a company's efforts to choose new sites for their expanding business. They have hired two consultants to separately evaluate the sites, who, in addition to an extended report, summarized each site as a "good," "fair," or "poor" prospect.
- **smokers.sav.** This data file is abstracted from the 1998 National Household Survey of Drug Abuse and is a probability sample of American households. (<http://dx.doi.org/10.3886/ICPSR02934>) Thus, the first step in an analysis of this data file should be to weight the data to reflect population trends.

- **stroke\_clean.sav.** This hypothetical data file contains the state of a medical database after it has been cleaned using procedures in the Data Preparation option.
- **stroke\_invalid.sav.** This hypothetical data file contains the initial state of a medical database and contains several data entry errors.
- **stroke\_survival.** This hypothetical data file concerns survival times for patients exiting a rehabilitation program post-ischemic stroke face a number of challenges. Post-stroke, the occurrence of myocardial infarction, ischemic stroke, or hemorrhagic stroke is noted and the time of the event recorded. The sample is left-truncated because it only includes patients who survived through the end of the rehabilitation program administered post-stroke.
- **stroke\_valid.sav.** This hypothetical data file contains the state of a medical database after the values have been checked using the Validate Data procedure. It still contains potentially anomalous cases.
- **survey\_sample.sav.** This data file contains survey data, including demographic data and various attitude measures. It is based on a subset of variables from the 1998 NORC General Social Survey, although some data values have been modified and additional fictitious variables have been added for demonstration purposes.
- **telco.sav.** This is a hypothetical data file that concerns a telecommunications company's efforts to reduce churn in their customer base. Each case corresponds to a separate customer and records various demographic and service usage information.
- **telco\_extra.sav.** This data file is similar to the *telco.sav* data file, but the "tenure" and log-transformed customer spending variables have been removed and replaced by standardized log-transformed customer spending variables.
- **telco\_missing.sav.** This data file is a subset of the *telco.sav* data file, but some of the demographic data values have been replaced with missing values.
- **testmarket.sav.** This hypothetical data file concerns a fast food chain's plans to add a new item to its menu. There are three possible campaigns for promoting the new product, so the new item is introduced at locations in several randomly selected markets. A different promotion is used at each location, and the weekly sales of the new item are recorded for the first four weeks. Each case corresponds to a separate location-week.
- **testmarket\_1month.sav.** This hypothetical data file is the *testmarket.sav* data file with the weekly sales "rolled-up" so that each case corresponds to a separate location. Some of the variables that changed weekly disappear as a result, and the sales recorded is now the sum of the sales during the four weeks of the study.
- **tree\_car.sav.** This is a hypothetical data file containing demographic and vehicle purchase price data.
- **tree\_credit.sav.** This is a hypothetical data file containing demographic and bank loan history data.
- **tree\_missing\_data.sav** This is a hypothetical data file containing demographic and bank loan history data with a large number of missing values.
- **tree\_score\_car.sav.** This is a hypothetical data file containing demographic and vehicle purchase price data.
- **tree\_textdata.sav.** A simple data file with only two variables intended primarily to show the default state of variables prior to assignment of measurement level and value labels.

- **tv-survey.sav.** This is a hypothetical data file that concerns a survey conducted by a TV studio that is considering whether to extend the run of a successful program. 906 respondents were asked whether they would watch the program under various conditions. Each row represents a separate respondent; each column is a separate condition.
- **ulcer\_recurrence.sav.** This file contains partial information from a study designed to compare the efficacy of two therapies for preventing the recurrence of ulcers. It provides a good example of interval-censored data and has been presented and analyzed elsewhere (Collett, 2003).
- **ulcer\_recurrence\_recoded.sav.** This file reorganizes the information in *ulcer\_recurrence.sav* to allow you model the event probability for each interval of the study rather than simply the end-of-study event probability. It has been presented and analyzed elsewhere (Collett et al., 2003).
- **verd1985.sav.** This data file concerns a survey (Verdegaal, 1985). The responses of 15 subjects to 8 variables were recorded. The variables of interest are divided into three sets. Set 1 includes *age* and *marital*, set 2 includes *pet* and *news*, and set 3 includes *music* and *live*. *Pet* is scaled as multiple nominal and *age* is scaled as ordinal; all of the other variables are scaled as single nominal.
- **virus.sav.** This is a hypothetical data file that concerns the efforts of an Internet service provider (ISP) to determine the effects of a virus on its networks. They have tracked the (approximate) percentage of infected e-mail traffic on its networks over time, from the moment of discovery until the threat was contained.
- **wheeze\_steubenville.sav.** This is a subset from a longitudinal study of the health effects of air pollution on children (Ware, Dockery, Spiro III, Speizer, and Ferris Jr., 1984). The data contain repeated binary measures of the wheezing status for children from Steubenville, Ohio, at ages 7, 8, 9 and 10 years, along with a fixed recording of whether or not the mother was a smoker during the first year of the study.
- **workprog.sav.** This is a hypothetical data file that concerns a government works program that tries to place disadvantaged people into better jobs. A sample of potential program participants were followed, some of whom were randomly selected for enrollment in the program, while others were not. Each case represents a separate program participant.

# Notices

Licensed Materials – Property of SPSS Inc., an IBM Company. © Copyright SPSS Inc. 1989, 2010.

Patent No. 7,023,453

**The following paragraph does not apply to the United Kingdom or any other country where such provisions are inconsistent with local law:** SPSS INC., AN IBM COMPANY, PROVIDES THIS PUBLICATION “AS IS” WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. Some states do not allow disclaimer of express or implied warranties in certain transactions, therefore, this statement may not apply to you.

This information could include technical inaccuracies or typographical errors. Changes are periodically made to the information herein; these changes will be incorporated in new editions of the publication. SPSS Inc. may make improvements and/or changes in the product(s) and/or the program(s) described in this publication at any time without notice.

Any references in this information to non-SPSS and non-IBM Web sites are provided for convenience only and do not in any manner serve as an endorsement of those Web sites. The materials at those Web sites are not part of the materials for this SPSS Inc. product and use of those Web sites is at your own risk.

When you send information to IBM or SPSS, you grant IBM and SPSS a nonexclusive right to use or distribute the information in any way it believes appropriate without incurring any obligation to you.

Information concerning non-SPSS products was obtained from the suppliers of those products, their published announcements or other publicly available sources. SPSS has not tested those products and cannot confirm the accuracy of performance, compatibility or any other claims related to non-SPSS products. Questions on the capabilities of non-SPSS products should be addressed to the suppliers of those products.

This information contains examples of data and reports used in daily business operations. To illustrate them as completely as possible, the examples include the names of individuals, companies, brands, and products. All of these names are fictitious and any similarity to the names and addresses used by an actual business enterprise is entirely coincidental.

#### COPYRIGHT LICENSE:

This information contains sample application programs in source language, which illustrate programming techniques on various operating platforms. You may copy, modify, and distribute these sample programs in any form without payment to SPSS Inc., for the purposes of developing,

using, marketing or distributing application programs conforming to the application programming interface for the operating platform for which the sample programs are written. These examples have not been thoroughly tested under all conditions. SPSS Inc., therefore, cannot guarantee or imply reliability, serviceability, or function of these programs. The sample programs are provided “AS IS”, without warranty of any kind. SPSS Inc. shall not be liable for any damages arising out of your use of the sample programs.

### **Trademarks**

IBM, the IBM logo, and [ibm.com](http://www.ibm.com) are trademarks of IBM Corporation, registered in many jurisdictions worldwide. A current list of IBM trademarks is available on the Web at <http://www.ibm.com/legal/copytrade.shtml>.

SPSS is a trademark of SPSS Inc., an IBM Company, registered in many jurisdictions worldwide.

Adobe, the Adobe logo, PostScript, and the PostScript logo are either registered trademarks or trademarks of Adobe Systems Incorporated in the United States, and/or other countries.

Intel, Intel logo, Intel Inside, Intel Inside logo, Intel Centrino, Intel Centrino logo, Celeron, Intel Xeon, Intel SpeedStep, Itanium, and Pentium are trademarks or registered trademarks of Intel Corporation or its subsidiaries in the United States and other countries.

Linux is a registered trademark of Linus Torvalds in the United States, other countries, or both.

Microsoft, Windows, Windows NT, and the Windows logo are trademarks of Microsoft Corporation in the United States, other countries, or both.

UNIX is a registered trademark of The Open Group in the United States and other countries.

Java and all Java-based trademarks and logos are trademarks of Sun Microsystems, Inc. in the United States, other countries, or both.

This product uses WinWrap Basic, Copyright 1993-2007, Polar Engineering and Consulting, <http://www.winwrap.com>.

Other product and service names might be trademarks of IBM, SPSS, or other companies.

Adobe product screenshot(s) reprinted with permission from Adobe Systems Incorporated.

Microsoft product screenshot(s) reprinted with permission from Microsoft Corporation.



---

# Bibliography

- Bell, E. H. 1961. *Social foundations of human behavior: Introduction to the study of sociology*. New York: Harper & Row.
- Blake, C. L., and C. J. Merz. 1998. "UCI Repository of machine learning databases." Available at <http://www.ics.uci.edu/~mlearn/MLRepository.html>.
- Breiman, L., and J. H. Friedman. 1985. Estimating optimal transformations for multiple regression and correlation. *Journal of the American Statistical Association*, 80, 580–598.
- Cochran, W. G. 1977. *Sampling Techniques*, 3rd ed. New York: John Wiley and Sons.
- Collett, D. 2003. *Modelling survival data in medical research*, 2 ed. Boca Raton: Chapman & Hall/CRC.
- Cox, D. R., and E. J. Snell. 1989. *The Analysis of Binary Data*, 2nd ed. London: Chapman and Hall.
- Green, P. E., and V. Rao. 1972. *Applied multidimensional scaling*. Hinsdale, Ill.: Dryden Press.
- Green, P. E., and Y. Wind. 1973. *Multiattribute decisions in marketing: A measurement approach*. Hinsdale, Ill.: Dryden Press.
- Guttman, L. 1968. A general nonmetric technique for finding the smallest coordinate space for configurations of points. *Psychometrika*, 33, 469–506.
- Hartigan, J. A. 1975. *Clustering algorithms*. New York: John Wiley and Sons.
- Hastie, T., and R. Tibshirani. 1990. *Generalized additive models*. London: Chapman and Hall.
- Kennedy, R., C. Riquier, and B. Sharp. 1996. Practical applications of correspondence analysis to categorical data in market research. *Journal of Targeting, Measurement, and Analysis for Marketing*, 5, 56–70.
- Kish, L. 1965. *Survey Sampling*. New York: John Wiley and Sons.
- Kish, L. 1987. *Statistical Design for Research*. New York: John Wiley and Sons.
- McCullagh, P., and J. A. Nelder. 1989. *Generalized Linear Models*, 2nd ed. London: Chapman & Hall.
- McFadden, D. 1974. Conditional logit analysis of qualitative choice behavior. In: *Frontiers in Economics*, P. Zarembka, ed. New York: Academic Press.
- Murthy, M. N. 1967. *Sampling Theory and Methods*. Calcutta, India: Statistical Publishing Society.
- Nagelkerke, N. J. D. 1991. A note on the general definition of the coefficient of determination. *Biometrika*, 78:3, 691–692.
- Price, R. H., and D. L. Bouffard. 1974. Behavioral appropriateness and situational constraints as dimensions of social behavior. *Journal of Personality and Social Psychology*, 30, 579–586.
- Rickman, R., N. Mitchell, J. Dingman, and J. E. Dalen. 1974. Changes in serum cholesterol during the Stillman Diet. *Journal of the American Medical Association*, 228, 54–58.

Rosenberg, S., and M. P. Kim. 1975. The method of sorting as a data-gathering procedure in multivariate research. *Multivariate Behavioral Research*, 10, 489–502.

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

Van der Ham, T., J. J. Meulman, D. C. Van Strien, and H. Van Engeland. 1997. Empirically based subgrouping of eating disorders in adolescents: A longitudinal perspective. *British Journal of Psychiatry*, 170, 363–368.

Verdegaal, R. 1985. *Meer sets analyse voor kwalitatieve gegevens (in Dutch)*. Leiden: Department of Data Theory, University of Leiden.

Ware, J. H., D. W. Dockery, A. Spiro III, F. E. Speizer, and B. G. Ferris Jr.. 1984. Passive smoking, gas cooking, and respiratory health of children living in six cities. *American Review of Respiratory Diseases*, 129, 366–374.

- adjusted chi-square
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- adjusted *F* statistic
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- adjusted residuals
  - in Complex Samples Crosstabs, 39
- aggregated residuals
  - in Complex Samples Cox Regression, 86
- analysis plan, 19
  
- baseline strata
  - in Complex Samples Cox Regression, 80
- Bonferroni
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- Breslow estimation method
  - in Complex Samples Cox Regression, 90
- Brewer's sampling method
  - in Sampling Wizard, 8
  
- chi-square
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- classification tables
  - in Complex Samples Logistic Regression, 57, 191
  - in Complex Samples Ordinal Regression, 68, 202
- clusters
  - in Analysis Preparation Wizard, 20
  - in Sampling Wizard, 6
- coefficient of variation (COV)
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34
  - in Complex Samples Frequencies, 30
  - in Complex Samples Ratios, 43
- column percentages
  - in Complex Samples Crosstabs, 39
- Complex Samples
  - hypothesis tests, 49, 59, 69
  - missing values, 31, 40
  - options, 32, 36, 41, 44
- Complex Samples Analysis Preparation Wizard, 140
  - public data, 140
  - related procedures, 154
  - sampling weights not available, 143
  - summary, 143, 154
- Complex Samples Cox Regression, 210
  - date and time variables, 74
  - define event, 77
  - hypothesis tests, 85
  - Kaplan-Meier analysis, 74
  - log-minus-log plot, 257
  - model, 81
  - model export, 88
  - options, 90
  - parameter estimates, 226, 255
  - pattern values, 256
  - piecewise-constant, time-dependent predictors, 226
  - plots, 84
  - predictors, 78
  - sample design information, 221, 254
  - save variables, 86
  - statistics, 82
  - subgroups, 80
  - test of proportional hazards, 222
  - tests of model effects, 222, 225, 255
  - time-dependent predictor, 79, 210
- Complex Samples Crosstabs, 37, 165
  - crosstabulation table, 168
  - related procedures, 170
  - relative risk, 165, 169–170
  - statistics, 39
- Complex Samples Descriptives, 33, 160
  - missing values, 35
  - public data, 160
  - related procedures, 164
  - statistics, 34, 163
  - statistics by subpopulation, 163
- Complex Samples Frequencies, 29, 155
  - frequency table, 158
  - frequency table by subpopulation, 158
  - related procedures, 159
  - statistics, 30
- Complex Samples General Linear Model, 45, 176
  - command additional features, 53
  - estimated means, 50
  - marginal means, 183
  - model, 47
  - model summary, 181
  - options, 52
  - parameter estimates, 182
  - related procedures, 185
  - save variables, 51
  - statistics, 48
  - tests of model effects, 181
- Complex Samples Logistic Regression, 54, 186
  - classification tables, 191
  - command additional features, 63
  - model, 56
  - odds ratios, 60, 193
  - options, 62
  - parameter estimates, 192

- pseudo  $R^2$  statistics, 190
  - reference category, 55
  - related procedures, 194
  - save variables, 61
  - statistics, 57
  - tests of model effects, 191
- Complex Samples Ordinal Regression, 64, 195
  - classification tables, 202
  - generalized cumulative model, 204
  - model, 66
  - odds ratios, 70, 203
  - options, 72
  - parameter estimates, 201
  - pseudo  $R^2$  statistics, 200, 208
  - related procedures, 209
  - response probabilities, 66
  - save variables, 71
  - statistics, 68
  - tests of model effects, 200
  - warnings, 207
- Complex Samples Ratios, 42, 171
  - missing values, 44
  - ratios, 174
  - related procedures, 175
  - statistics, 43
- Complex Samples Sampling Wizard, 93
  - PPS sampling, 123
  - related procedures, 139
  - sampling frame, full, 93
  - sampling frame, partial, 105
  - summary, 103, 135
- complex sampling
  - analysis plan, 19
  - sample plan, 4
- confidence intervals
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34, 163
  - in Complex Samples Frequencies, 30, 158
  - in Complex Samples General Linear Model, 48, 52
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
  - in Complex Samples Ratios, 43
- confidence level
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- contrasts
  - in Complex Samples General Linear Model, 50
- correlations of parameter estimates
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
- covariances of parameter estimates
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
- Cox-Snell residuals
  - in Complex Samples Cox Regression, 86
- crosstabulation table
  - in Complex Samples Crosstabs, 168
- cumulative probabilities
  - in Complex Samples Ordinal Regression, 71
- cumulative values
  - in Complex Samples Frequencies, 30
- degrees of freedom
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- design effect
  - in Complex Samples Cox Regression, 82
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34
  - in Complex Samples Frequencies, 30
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
  - in Complex Samples Ratios, 43
- deviance residuals
  - in Complex Samples Cox Regression, 86
- deviation contrasts
  - in Complex Samples General Linear Model, 50
- difference contrasts
  - in Complex Samples General Linear Model, 50
- Efron estimation method
  - in Complex Samples Cox Regression, 90
- estimated marginal means
  - in Complex Samples General Linear Model, 50
- expected values
  - in Complex Samples Crosstabs, 39
- $F$  statistic
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- Fisher scoring
  - in Complex Samples Ordinal Regression, 72
- generalized cumulative model
  - in Complex Samples Ordinal Regression, 204
- Helmert contrasts
  - in Complex Samples General Linear Model, 50
- inclusion probabilities
  - in Sampling Wizard, 12
- input sample weights
  - in Sampling Wizard, 6
- iteration history
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72

- iterations
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- least significant difference
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- legal notices, 267
- likelihood convergence
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- log-minus-log plot
  - in Complex Samples Cox Regression, 257
- marginal means
  - in GLM Univariate, 183
- martingale residuals
  - in Complex Samples Cox Regression, 86
- mean
  - in Complex Samples Descriptives, 34, 163
- measure of size
  - in Sampling Wizard, 8
- missing values
  - in Complex Samples, 31, 40
  - in Complex Samples Descriptives, 35
  - in Complex Samples General Linear Model, 52
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
  - in Complex Samples Ratios, 44
- Murthy's sampling method
  - in Sampling Wizard, 8
- Newton-Raphson method
  - in Complex Samples Ordinal Regression, 72
- odds ratios
  - in Complex Samples Crosstabs, 39, 165
  - in Complex Samples Logistic Regression, 60, 193
  - in Complex Samples Ordinal Regression, 70, 203
- parameter convergence
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- parameter estimates
  - in Complex Samples Cox Regression, 82
  - in Complex Samples General Linear Model, 48, 182
  - in Complex Samples Logistic Regression, 57, 192
  - in Complex Samples Ordinal Regression, 68, 201
- piecewise-constant, time-dependent predictors
  - in Complex Samples Cox Regression, 226
- plan file, 2
- polynomial contrasts
  - in Complex Samples General Linear Model, 50
- population size
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34
  - in Complex Samples Frequencies, 30, 158
  - in Complex Samples Ratios, 43
  - in Sampling Wizard, 12
- PPS sampling
  - in Sampling Wizard, 8
- predicted categories
  - in Complex Samples Logistic Regression, 61
  - in Complex Samples Ordinal Regression, 71
- predicted probability
  - in Complex Samples Logistic Regression, 61
  - in Complex Samples Ordinal Regression, 71
- predicted values
  - in Complex Samples General Linear Model, 51
- predictor patterns
  - in Complex Samples Cox Regression, 256
- proportional hazards test
  - in Complex Samples Cox Regression, 222
- pseudo  $R^2$  statistics
  - in Complex Samples Logistic Regression, 57, 190
  - in Complex Samples Ordinal Regression, 68, 200, 208
- public data
  - in Analysis Preparation Wizard, 140
  - in Complex Samples Descriptives, 160
- $R^2$  statistic
  - in Complex Samples General Linear Model, 48, 181
- ratios
  - in Complex Samples Ratios, 174
- reference category
  - in Complex Samples General Linear Model, 50
  - in Complex Samples Logistic Regression, 55
- relative risk
  - in Complex Samples Crosstabs, 39, 165, 169–170
- repeated contrasts
  - in Complex Samples General Linear Model, 50
- residuals
  - in Complex Samples Crosstabs, 39
  - in Complex Samples General Linear Model, 51
- response probabilities
  - in Complex Samples Ordinal Regression, 66
- risk difference
  - in Complex Samples Crosstabs, 39
- row percentages
  - in Complex Samples Crosstabs, 39
- Sampford's sampling method
  - in Sampling Wizard, 8
- sample design information
  - in Complex Samples Cox Regression, 82, 221, 254
- sample files
  - location, 258
- sample plan, 4
- sample proportion
  - in Sampling Wizard, 12
- sample size
  - in Sampling Wizard, 10, 12

- sample weights
  - in Analysis Preparation Wizard, 20
  - in Sampling Wizard, 12
- sampling
  - complex design, 4
- sampling estimation
  - in Analysis Preparation Wizard, 22
- sampling frame, full
  - in Sampling Wizard, 93
- sampling frame, partial
  - in Sampling Wizard, 105
- sampling method
  - in Sampling Wizard, 8
- Schoenfeld's partial residuals
  - in Complex Samples Cox Regression, 86
- score residuals
  - in Complex Samples Cox Regression, 86
- separation
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- sequential Bonferroni correction
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- sequential sampling
  - in Sampling Wizard, 8
- sequential Sidak correction
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- Sidak correction
  - in Complex Samples, 49, 59, 69
  - in Complex Samples Cox Regression, 85
- simple contrasts
  - in Complex Samples General Linear Model, 50
- simple random sampling
  - in Sampling Wizard, 8
- square root of design effect
  - in Complex Samples Cox Regression, 82
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34
  - in Complex Samples Frequencies, 30
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
  - in Complex Samples Ratios, 43
- standard error
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34, 163
  - in Complex Samples Frequencies, 30, 158
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
  - in Complex Samples Ratios, 43
- step-halving
  - in Complex Samples Logistic Regression, 62
  - in Complex Samples Ordinal Regression, 72
- stratification
  - in Analysis Preparation Wizard, 20
  - in Sampling Wizard, 6
- subpopulation
  - in Complex Samples Cox Regression, 80
- sum
  - in Complex Samples Descriptives, 34
- summary
  - in Analysis Preparation Wizard, 143, 154
  - in Sampling Wizard, 103, 135
- systematic sampling
  - in Sampling Wizard, 8
- t* test
  - in Complex Samples General Linear Model, 48
  - in Complex Samples Logistic Regression, 57
  - in Complex Samples Ordinal Regression, 68
- table percentages
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Frequencies, 30, 158
- test of parallel lines
  - in Complex Samples Ordinal Regression, 68, 204
- test of proportional hazards
  - in Complex Samples Cox Regression, 82
- tests of model effects
  - in Complex Samples Cox Regression, 255
  - in Complex Samples General Linear Model, 181
  - in Complex Samples Logistic Regression, 191
  - in Complex Samples Ordinal Regression, 200
- time-dependent predictor
  - in Complex Samples Cox Regression, 79, 210
- trademarks, 268
- unweighted count
  - in Complex Samples Crosstabs, 39
  - in Complex Samples Descriptives, 34
  - in Complex Samples Frequencies, 30
  - in Complex Samples Ratios, 43
- warnings
  - in Complex Samples Ordinal Regression, 207