

> SPSS Categories® 13.0

Jacqueline J. Meulman
Willem J. Heiser
SPSS Inc



For more information about SPSS® software products, please visit our Web site at <http://www.spss.com> or contact

SPSS Inc.

233 South Wacker Drive, 11th Floor
Chicago, IL 60606-6412

Tel: (312) 651-3000

Fax: (312) 651-3668

SPSS is a registered trademark and the other product names are the trademarks of SPSS Inc. for its proprietary computer software. No material describing such software may be produced or distributed without the written permission of the owners of the trademark and license rights in the software and the copyrights in the published materials.

The SOFTWARE and documentation are provided with RESTRICTED RIGHTS. Use, duplication, or disclosure by the Government is subject to restrictions as set forth in subdivision (c) (1) (ii) of The Rights in Technical Data and Computer Software clause at 52.227-7013. Contractor/manufacturer is SPSS Inc., 233 South Wacker Drive, 11th Floor, Chicago, IL 60606-6412.

General notice: Other product names mentioned herein are used for identification purposes only and may be trademarks of their respective companies.

TableLook is a trademark of SPSS Inc.

Windows is a registered trademark of Microsoft Corporation.

DataDirect, DataDirect Connect, INTERSOLV, and SequeLink are registered trademarks of DataDirect Technologies.

Portions of this product were created using LEADTOOLS © 1991–2000, LEAD Technologies, Inc. ALL RIGHTS RESERVED.

LEAD, LEADTOOLS, and LEADVIEW are registered trademarks of LEAD Technologies, Inc.

Sax Basic is a trademark of Sax Software Corporation. Copyright © 1993–2004 by Polar Engineering and Consulting.

All rights reserved.

Portions of this product were based on the work of the FreeType Team (<http://www.freetype.org>).

A portion of the SPSS software contains zlib technology. Copyright © 1995–2002 by Jean-loup Gailly and Mark Adler. The zlib software is provided “as is,” without express or implied warranty.

A portion of the SPSS software contains Sun Java Runtime libraries. Copyright © 2003 by Sun Microsystems, Inc. All rights reserved. The Sun Java Runtime libraries include code licensed from RSA Security, Inc. Some portions of the libraries are licensed from IBM and are available at <http://oss.software.ibm.com/icu4j/>.

SPSS Categories® 13.0

Copyright © 2004 by SPSS Inc.

All rights reserved.

Printed in the United States of America.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher.

1 2 3 4 5 6 7 8 9 0 07 06 05 04

ISBN 1-56827-351-7

Preface

SPSS 13.0 is a comprehensive system for analyzing data. The Categories optional add-on module provides the additional analytic techniques described in this manual. The Categories add-on module must be used with the SPSS 13.0 Base system and is completely integrated into that system.

Installation

To install the Categories add-on module, run the License Authorization Wizard using the authorization code that you received from SPSS Inc. For more information, see the installation instructions supplied with the SPSS Base system.

Compatibility

SPSS is designed to run on many computer systems. See the installation instructions that came with your system for specific information on minimum and recommended requirements.

Serial Numbers

Your serial number is your identification number with SPSS Inc. You will need this serial number when you contact SPSS Inc. for information regarding support, payment, or an upgraded system. The serial number was provided with your Base system.

Customer Service

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

Technical Support

The services of SPSS Technical Support are available to registered customers. Customers may contact Technical Support for assistance in using SPSS or for installation help for one of the supported hardware environments. To reach Technical Support, see the SPSS Web site at <http://www.spss.com>, or contact your local office, listed on the SPSS Web site at <http://www.spss.com/worldwide>. Be prepared to identify yourself, your organization, and the serial number of your system.

Additional Publications

Additional copies of SPSS product manuals may be purchased directly from SPSS Inc. Visit the SPSS Web Store at <http://www.spss.com/estore>, or contact your local SPSS office, listed on the SPSS Web site at <http://www.spss.com/worldwide>. For telephone orders in the United States and Canada, call SPSS Inc. at 800-543-2185. For telephone orders outside of North America, contact your local office, listed on the SPSS Web site.

The *SPSS Statistical Procedures Companion*, by Marija Norušis, has been published by Prentice Hall. A new version of this book, updated for SPSS 13.0, is planned. The *SPSS Advanced Statistical Procedures Companion*, also based on SPSS 13.0, is forthcoming. The *SPSS Guide to Data Analysis* for SPSS 13.0 is also in development. Announcements of publications available exclusively through Prentice Hall will be available on the SPSS Web site at <http://www.spss.com/estore> (select your home country, and then click Books).

Tell Us Your Thoughts

Your comments are important. Please let us know about your experiences with SPSS products. We especially like to hear about new and interesting applications using the SPSS system. Please send e-mail to suggest@spss.com or write to SPSS Inc.,

Attn.: Director of Product Planning, 233 South Wacker Drive, 11th Floor, Chicago, IL 60606-6412.

About This Manual

This manual documents the graphical user interface for the procedures included in the Categories add-on module. Illustrations of dialog boxes are taken from SPSS for Windows. Dialog boxes in other operating systems are similar. Detailed information about the command syntax for features in this module is provided in the *SPSS Command Syntax Reference*, available from the Help menu.

Contacting SPSS

If you would like to be on our mailing list, contact one of our offices, listed on our Web site at <http://www.spss.com/worldwide>.

Acknowledgments

The optimal scaling procedures and their SPSS implementation were developed by the Data Theory Scaling System Group (DTSS), consisting of members of the departments of Education and Psychology of the Faculty of Social and Behavioral Sciences at Leiden University.

Willem Heiser, Jacqueline Meulman, Gerda van den Berg, and Patrick Groenen were involved with the original 1990 procedures. Jacqueline Meulman and Peter Neufeglise participated in the development of procedures for categorical regression, correspondence analysis, categorical principal components analysis, and multidimensional scaling. In addition, Anita van der Kooij contributed especially to CATREG, CORRESPONDENCE, and CATPCA, and Frank Busing and Willem Heiser, to the PROXSCAL procedure. The development of PROXSCAL has profited from technical comments and suggestions from Jacques Commandeur and Patrick Groenen.

Contents

Part I: User's Guide

1	<i>Introduction to SPSS Optimal Scaling Procedures for Categorical Data</i>	1
	What Is Optimal Scaling?	1
	Why Use Optimal Scaling?	2
	Optimal Scaling Level and Measurement Level	2
	Selecting the Optimal Scaling Level	3
	Transformation Plots	4
	Category Codes	5
	Which Procedure Is Best for Your Application?	8
	Categorical Regression	9
	Categorical Principal Components Analysis	10
	Nonlinear Canonical Correlation Analysis	11
	Correspondence Analysis	12
	Multiple Correspondence Analysis	14
	Multidimensional Scaling	15
	Aspect Ratio in Optimal Scaling Charts	16
	Recommended Readings	16
2	<i>Categorical Regression (CATREG)</i>	19
	Categorical Regression	19
	Define Scale in Categorical Regression	21
	Categorical Regression Discretization	23

Categorical Regression Missing Values	24
Categorical Regression Options	25
Categorical Regression Output	26
Categorical Regression Save	28
Categorical Regression Transformation Plots	29
CATREG Command Additional Features.	29

3 Categorical Principal Components Analysis (CATPCA) 31

Define Scale and Weight in CATPCA.	34
Categorical Principal Components Analysis Discretization	36
Categorical Principal Components Analysis Missing Values	38
Categorical Principal Components Analysis Options	39
Categorical Principal Components Analysis Output.	42
Categorical Principal Components Analysis Save	44
Categorical Principal Components Analysis Object Plots.	45
Categorical Principal Components Analysis Category Plots.	46
Categorical Principal Components Analysis Loading Plots	47
CATPCA Command Additional Features.	48

4 Nonlinear Canonical Correlation Analysis (OVERALS) 49

Nonlinear Canonical Correlation Analysis	49
Define Range and Scale	53
Define Range	54

Nonlinear Canonical Correlation Analysis Options	54
OVERALS Command Additional Features.	56

5 Correspondence Analysis 59

Define Row Range in Correspondence Analysis	61
Define Column Range in Correspondence Analysis	62
Correspondence Analysis Model	63
Correspondence Analysis Statistics	66
Correspondence Analysis Plots	67
CORRESPONDENCE Command Additional Features.	68

6 Multiple Correspondence Analysis 71

Define Variable Weight in Multiple Correspondence Analysis	73
Multiple Correspondence Analysis Discretization	74
Multiple Correspondence Analysis Missing Values	75
Multiple Correspondence Analysis Options.	77
Multiple Correspondence Analysis Output	80
Multiple Correspondence Analysis Save	82
Multiple Correspondence Analysis Object Plots	82
Multiple Correspondence Analysis Variable Plots	83
MULTIPLE CORRESPONDENCE Command Additional Features	85

7 Multidimensional Scaling (PROXSCAL) 87

Proximities in Matrices across Columns	90
--	----

Proximities in Columns	91
Proximities in One Column	92
Create Proximities from Data	93
Measures Dialog Box	94
Define a Multidimensional Scaling Model	95
Multidimensional Scaling Restrictions	97
Multidimensional Scaling Options	98
Multidimensional Scaling Plots, Version 1	99
Multidimensional Scaling Plots, Version 2	101
Multidimensional Scaling Output	101
PROXSCAL Command Additional Features	103

Part II: Examples

8 *Categorical Regression* 107

Example: Carpet Cleaner Data	107
A Standard Linear Regression Analysis	108
A Categorical Regression Analysis	117
Example: Ozone Data	130
Discretizing Variables	131
Selection of Transformation Type	132
Optimality of the Quantifications	145
Effects of Transformations	148
Recommended Readings	157

9 Categorical Principal Components Analysis 159

Example: Examining Interrelations of Social Systems	159
Running the Analysis	160
Number of Dimensions	165
Quantifications	167
Object Scores	169
Component Loadings	170
Additional Dimensions	172
Example: Symptomatology of Eating Disorders	175
Running the Analysis	177
Transformation Plots	190
Model Summary	194
Component Loadings	194
Object Scores	196
Examining the Structure of the Course of Illness	198
Recommended Readings	213

10 Nonlinear Canonical Correlation Analysis 217

Example: An Analysis of Survey Results	217
Examining the Data	218
Accounting for Similarity between Sets	226
Component Loadings	230
Transformation Plots	231
Single versus Multiple Category Coordinates	235
Centroids and Projected Centroids	236
An Alternative Analysis	241
General Suggestions	247
Recommended Readings	248

11 Correspondence Analysis **249**

Normalization	250
Example: Smoking Behavior by Job Category	251
Running the Analysis	252
Correspondence Table	255
Dimensionality	256
Biplot	257
Profiles and Distances	258
Row and Column Scores	260
Permutations of the Correspondence Table	262
Confidence Statistics	263
Supplementary Profiles	264
Example: Perceptions of Coffee Brands	270
Running the Analysis	271
Dimensionality	276
Contributions	277
Plots	279
Symmetrical Normalization	281
Example: Flying Mileage between Cities	283
Correspondence Table	292
Row and Column Scores	293
Recommended Readings	294

12 Multiple Correspondence Analysis **297**

Example: Characteristics of Hardware	297
Running the Analysis	298
Model Summary	301
Object Scores	302
Discrimination Measures	304

Category Quantifications	305
A More Detailed Look at Object Scores	307
Omission of Outliers	310
Recommended Readings	315

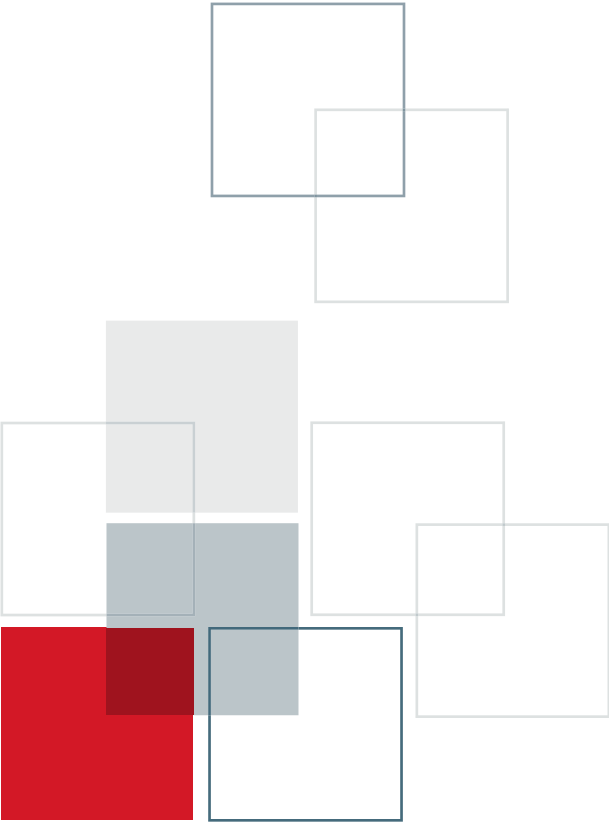
13 Multidimensional Scaling 317

Example: An Examination of Kinship Terms	317
Choosing the Number of Dimensions	318
A Three-Dimensional Solution	328
A Three-Dimensional Solution with Nondefault Transformations	335
Discussion	340
Recommended Readings	341

***Bibliography* 343**

***Index* 351**

Part 1: User's Guide



Introduction to SPSS Optimal Scaling Procedures for Categorical Data

SPSS Categories procedures use optimal scaling to analyze data that are difficult or impossible for standard statistical procedures to analyze. This chapter describes what each procedure does, the situations in which each procedure is most appropriate, the relationships between the procedures, and the relationships of these procedures to their standard statistical counterparts.

Note: These procedures and their SPSS implementation were developed by the Data Theory Scaling System Group (DTSS), consisting of members of the departments of Education and Psychology, Faculty of Social and Behavioral Sciences, Leiden University.

What Is Optimal Scaling?

The idea behind optimal scaling is to assign numerical quantifications to the categories of each variable, thus allowing standard procedures to be used to obtain a solution on the quantified variables.

The optimal scale values are assigned to categories of each variable based on the optimizing criterion of the procedure in use. Unlike the original labels of the nominal or ordinal variables in the analysis, these scale values have metric properties.

In most Categories procedures, the optimal quantification for each scaled variable is obtained through an iterative method called **alternating least squares** in which, after the current quantifications are used to find a solution, the quantifications are updated using that solution. The updated quantifications are then used to find a new

solution, which is used to update the quantifications, and so on, until some criterion is reached that signals the process to stop.

Why Use Optimal Scaling?

Categorical data are often found in marketing research, survey research, and research in the social and behavioral sciences. In fact, many researchers deal almost exclusively with categorical data.

While adaptations of most standard models exist specifically to analyze categorical data, they often do not perform well for data sets that feature:

- Too few observations
- Too many variables
- Too many values per variable

By quantifying categories, optimal scaling techniques avoid problems in these situations. Moreover, they are useful even when specialized techniques are appropriate.

Rather than interpreting parameter estimates, the interpretation of optimal scaling output is often based on graphical displays. Optimal scaling techniques offer excellent exploratory analyses, which complement other SPSS models well. By narrowing the focus of your investigation, visualizing your data through optimal scaling can form the basis of an analysis that centers on interpretation of model parameters.

Optimal Scaling Level and Measurement Level

This can be a very confusing concept when you first use Categories procedures. When specifying the level, you specify not the level at which variables are *measured* but the level at which they are *scaled*. The idea is that the variables to be quantified may have nonlinear relations regardless of how they are measured.

For Categories purposes, there are three basic levels of measurement:

- The **nominal** level implies that a variable's values represent unordered categories. Examples of variables that might be nominal are region, zip code area, religious affiliation, and multiple choice categories.

- The **ordinal** level implies that a variable's values represent ordered categories. Examples include attitude scales representing degree of satisfaction or confidence and preference rating scores.
- The **numerical** level implies that a variable's values represent ordered categories with a meaningful metric so that distance comparisons between categories are appropriate. Examples include age in years and income in thousands of dollars.

For example, suppose that the variables *region*, *job*, and *age* are coded as shown in the following table.

Table 1-1
Coding scheme for *region*, *job*, and *age*

Region		Job		Age	
1	North	1	intern	20	twenty years old
2	South	2	sales rep	22	twenty-two years old
3	East	3	manager	25	twenty-five years old
4	West			27	twenty-seven years old

The values shown represent the categories of each variable. *Region* would be a nominal variable. There are four categories of *region*, with no intrinsic ordering. Values 1 through 4 simply represent the four categories; the coding scheme is completely arbitrary. *Job*, on the other hand, could be assumed to be an ordinal variable. The original categories form a progression from intern to manager. Larger codes represent a job higher on the corporate ladder. However, only the order information is known—nothing can be said about the distance between adjacent categories. In contrast, *age* could be assumed to be a numerical variable. In the case of *age*, the distances between the values are intrinsically meaningful. The distance between 20 and 22 is the same as the distance between 25 and 27, while the distance between 22 and 25 is greater than either of these.

Selecting the Optimal Scaling Level

It is important to understand that there are no intrinsic properties of a variable that automatically predefine what optimal scaling level you should specify for it. You can explore your data in any way that makes sense and makes interpretation easier. By

analyzing a numerical-level variable at the ordinal level, for example, the use of a nonlinear transformation may allow a solution in fewer dimensions.

The following two examples illustrate how the “obvious” level of measurement might not be the best optimal scaling level. Suppose that a variable sorts objects into age groups. Although age can be scaled as a numerical variable, it may be true that for people younger than 25 safety has a positive relation with age, whereas for people older than 60 safety has a negative relation with age. In this case, it might be better to treat age as a nominal variable.

As another example, a variable that sorts persons by political preference appears to be essentially nominal. However, if you order the parties from political left to political right, you might want the quantification of parties to respect this order by using an ordinal level of analysis.

Even though there are no predefined properties of a variable that make it exclusively one level or another, there are some general guidelines to help the novice user. With single-nominal quantification, you don’t usually know the order of the categories but you want the analysis to impose one. If the order of the categories is known, you should try ordinal quantification. If the categories are unorderable, you might try multiple-nominal quantification.

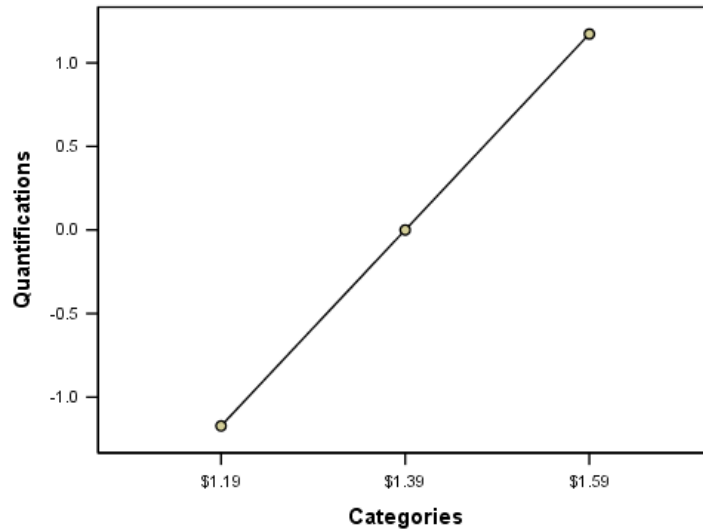
Transformation Plots

The different levels at which each variable can be scaled impose different restrictions on the quantifications. Transformation plots illustrate the relationship between the quantifications and the original categories resulting from the selected optimal scaling level. For example, a linear transformation plot results when a variable is treated as numerical. Variables treated as ordinal result in a nondecreasing transformation plot. Transformation plots for variables treated nominally that are U-shaped (or the inverse) display a quadratic relationship. Nominal variables could also yield transformation plots without apparent trends by changing the order of the categories completely. The following figure displays a sample transformation plot.

Transformation plots are particularly suited to determining how well the selected optimal scaling level performs. If several categories receive similar quantifications, collapsing these categories into one category may be warranted. Alternatively, if a variable treated as nominal receives quantifications that display an increasing trend, an ordinal transformation may result in a similar fit. If that trend is linear, numerical

treatment may be appropriate. However, if collapsing categories or changing scaling levels is warranted, the analysis will not change significantly.

Figure 1-1
Transformation plot of price (numerical)



Category Codes

Some care should be taken when coding categorical variables because some coding schemes may yield unwanted output or incomplete analyses. Possible coding schemes for *job* are displayed in the following table.

Table 1-2
Alternative coding schemes for job

Category	Scheme			
	A	B	C	D
intern	1	1	5	1
sales rep	2	2	6	5
manager	3	7	7	3

Some Categories procedures require that the range of every variable used be defined. Any value outside this range is treated as a missing value. The minimum category value is always 1. The maximum category value is supplied by the user. This value is not the *number* of categories for a variable—it is the *largest* category value. For example, in the table, scheme A has a maximum category value of 3 and scheme B has a maximum category value of 7, yet both schemes code the same three categories.

The variable range determines which categories will be omitted from the analysis. Any categories with codes outside the defined range are omitted from the analysis. This is a simple method for omitting categories but can result in unwanted analyses. An incorrectly defined maximum category can omit *valid* categories from the analysis. For example, for scheme B, defining the maximum category value to be 3 indicates that *job* has categories coded from 1 to 3; the *manager* category is treated as missing. Because no category has actually been coded 3, the third category in the analysis contains no cases. If you wanted to omit all manager categories, this analysis would be appropriate. However, if managers are to be included, the maximum category must be defined as 7, and missing values must be coded with values above 7 or below 1.

For variables treated as nominal or ordinal, the range of the categories does not affect the results. For nominal variables, only the label and not the value associated with that label is important. For ordinal variables, the order of the categories is preserved in the quantifications; the category values themselves are not important. All coding schemes resulting in the same category ordering will have identical results. For example, the first three schemes in the table are functionally equivalent if *job* is analyzed at an ordinal level. The order of the categories is identical in these schemes. Scheme D, on the other hand, inverts the second and third categories and will yield different results than the other schemes.

Although many coding schemes for a variable are functionally equivalent, schemes with small differences between codes are preferred because the codes have an impact on the amount of output produced by a procedure. All categories coded with values between 1 and the user-defined maximum are valid. If any of these categories are empty, the corresponding quantifications will be either system missing or zero, depending on the procedure. Although neither of these assignments affect the analyses, output is produced for these categories. Thus, for scheme B, *job* has four categories that receive system-missing values. For scheme C, there are also four categories receiving system-missing indicators. In contrast, for scheme A there are no system-missing quantifications. Using consecutive integers as codes for variables treated as nominal or ordinal results in much less output without affecting the results.

Coding schemes for variables treated as numerical are more restricted than the ordinal case. For these variables, the differences between consecutive categories are important. The following table displays three coding schemes for *age*.

Table 1-3
Alternative coding schemes for age

Category	Scheme		
	A	B	C
20	20	1	1
22	22	3	2
25	25	6	3
27	27	8	4

Any recoding of numerical variables must preserve the differences between the categories. Using the original values is one method for ensuring preservation of differences. However, this can result in many categories having system-missing indicators. For example, scheme A employs the original observed values. For all Categories procedures except for Correspondence Analysis, the maximum category value is 27 and the minimum category value is set to 1. The first 19 categories are empty and receive system-missing indicators. The output can quickly become rather cumbersome if the maximum category is much greater than 1 and there are many empty categories between 1 and the maximum.

To reduce the amount of output, recoding can be done. However, in the numerical case, the Automatic Recode facility should not be used. Coding to consecutive integers results in differences of 1 between all consecutive categories, and, as a result, all quantifications will be equally spaced. The metric characteristics deemed important when treating a variable as numerical are destroyed by recoding to consecutive integers. For example, scheme C in the table corresponds to automatically recoding *age*. The difference between categories 22 and 25 has changed from three to one, and the quantifications will reflect the latter difference.

An alternative recoding scheme that preserves the differences between categories is to subtract the smallest category value from every category and add 1 to each difference. Scheme B results from this transformation. The smallest category value, 20, has been subtracted from each category, and 1 was added to each result. The transformed codes have a minimum of 1, and all differences are identical to the original data. The maximum category value is now 8, and the zero quantifications before the first nonzero quantification are all eliminated. Yet, the nonzero

quantifications corresponding to each category resulting from scheme B are identical to the quantifications from scheme A.

Which Procedure Is Best for Your Application?

The techniques embodied in four of these procedures (Correspondence Analysis, Multiple Correspondence Analysis, Categorical Principal Components Analysis, and Nonlinear Canonical Correlation Analysis) fall into the general area of multivariate data analysis known as **dimension reduction**. That is, relationships between variables are represented in a few dimensions—say two or three—as often as possible. This enables you to describe structures or patterns in the relationships that would be too difficult to fathom in their original richness and complexity. In market research applications, these techniques can be a form of **perceptual mapping**. A major advantage of these procedures is that they accommodate data with different levels of optimal scaling.

Categorical Regression describes the relationship between a categorical response variable and a combination of categorical predictor variables. The influence of each predictor variable on the response variable is described by the corresponding regression weight. As in the other procedures, data can be analyzed with different levels of optimal scaling.

Multidimensional Scaling describes relationships between objects in as few dimensions as possible, starting either with a matrix of proximities between the objects or with the original data from which the proximities are computed.

Following are brief guidelines for each of the procedures:

- Use Categorical Regression to predict the values of a categorical dependent variable from a combination of categorical independent variables.
- Use Categorical Principal Components Analysis to account for patterns of variation in a single set of variables of mixed optimal scaling levels.
- Use Nonlinear Canonical Correlation Analysis to assess the extent to which two or more sets of variables of mixed optimal scaling levels are correlated.
- Use Correspondence Analysis to analyze two-way contingency tables or data that can be expressed as a two-way table, such as brand preference or sociometric choice data.

- Use Multiple Correspondence Analysis to analyze a categorical multivariate data matrix when you are willing to make no stronger assumption that all variables are analyzed at the nominal level.
- Use Multidimensional Scaling to analyze proximity data to find a least-squares representation of the objects in a low-dimensional space.

Categorical Regression

The use of Categorical Regression is most appropriate when the goal of your analysis is to predict a dependent (response) variable from a set of independent (predictor) variables. As with all optimal scaling procedures, scale values are assigned to each category of every variable such that these values are optimal with respect to the regression. The solution of a categorical regression maximizes the squared correlation between the transformed response and the weighted combination of transformed predictors.

Relation to other Categories procedures. Categorical regression with optimal scaling is comparable to optimal scaling canonical correlation analysis with two sets, one of which contains only the dependent variable. In the latter technique, similarity of sets is derived by comparing each set to an unknown variable that lies somewhere between all of the sets. In categorical regression, similarity of the transformed response and the linear combination of transformed predictors is assessed directly.

Relation to standard techniques. In standard linear regression, categorical variables can either be recoded as indicator variables or be treated in the same fashion as interval level variables. In the first approach, the model contains a separate intercept and slope for each combination of the levels of the categorical variables. This results in a large number of parameters to interpret. In the second approach, only one parameter is estimated for each variable. However, the arbitrary nature of the category codings makes generalizations impossible.

If some of the variables are not continuous, alternative analyses are available. If the response is continuous and the predictors are categorical, analysis of variance is often employed. If the response is categorical and the predictors are continuous, logistic regression or discriminant analysis may be appropriate. If the response and the predictors are both categorical, loglinear models are often used.

Regression with optimal scaling offers three scaling levels for each variable. Combinations of these levels can account for a wide range of nonlinear relationships for which any single “standard” method is ill-suited. Consequently, optimal scaling offers greater flexibility than the standard approaches with minimal added complexity.

In addition, nonlinear transformations of the predictors usually reduce the dependencies among the predictors. If you compare the eigenvalues of the correlation matrix for the predictors with the eigenvalues of the correlation matrix for the optimally scaled predictors, the latter set will usually be less variable than the former. In other words, in categorical regression, optimal scaling makes the larger eigenvalues of the predictor correlation matrix smaller and the smaller eigenvalues larger.

Categorical Principal Components Analysis

The use of Categorical Principal Components Analysis is most appropriate when you want to account for patterns of variation in a single set of variables of mixed optimal scaling levels. This technique attempts to reduce the dimensionality of a set of variables while accounting for as much of the variation as possible. Scale values are assigned to each category of every variable so that these values are optimal with respect to the principal components solution. Objects in the analysis receive component scores based on the quantified data. Plots of the component scores reveal patterns among the objects in the analysis and can reveal unusual objects in the data. The solution of a categorical principal components analysis maximizes the correlations of the object scores with each of the quantified variables for the number of components (dimensions) specified.

An important application of categorical principal components is to examine preference data, in which respondents rank or rate a number of items with respect to preference. In the usual SPSS data configuration, rows are individuals, columns are measurements for the items, and the scores across rows are preference scores (on a 0 to 10 scale, for example), making the data row-conditional. For preference data, you may want to treat the individuals as variables. Using the Transpose procedure, you can transpose the data. The raters become the variables, and all variables are declared ordinal. There is no objection to using more variables than objects in CATPCA.

Relation to other Categories procedures. If all variables are declared multiple nominal, categorical principal components analysis produces an analysis equivalent to a multiple correspondence analysis run on the same variables. Thus, categorical

principal components analysis can be seen as a type of multiple correspondence analysis in which some of the variables are declared ordinal or numerical.

Relation to standard techniques. If all variables are scaled on the numerical level, categorical principal components analysis is equivalent to standard principal components analysis.

More generally, categorical principal components analysis is an alternative to computing the correlations between non-numerical scales and analyzing them using a standard principal components or factor-analysis approach. Naive use of the usual Pearson correlation coefficient as a measure of association for ordinal data can lead to nontrivial bias in estimation of the correlations.

Nonlinear Canonical Correlation Analysis

Nonlinear Canonical Correlation Analysis is a very general procedure with many different applications. The goal of nonlinear canonical correlation analysis is to analyze the relationships between two or more sets of variables instead of between the variables themselves, as in principal components analysis. For example, you may have two sets of variables, where one set of variables might be demographic background items on a set of respondents and a second set might be responses to a set of attitude items. The scaling levels in the analysis can be any mix of nominal, ordinal, and numerical. Optimal scaling canonical correlation analysis determines the similarity among the sets by simultaneously comparing the canonical variables from each set to a compromise set of scores assigned to the objects.

Relation to other Categories procedures. If there are two or more sets of variables with only one variable per set, optimal scaling canonical correlation analysis is equivalent to optimal scaling principal components analysis. If all variables in a one-variable-per-set analysis are multiple nominal, optimal scaling canonical correlation analysis is equivalent to multiple correspondence analysis. If there are two sets of variables, one of which contains only one variable, optimal scaling canonical correlation analysis is equivalent to categorical regression with optimal scaling.

Relation to standard techniques. Standard canonical correlation analysis is a statistical technique that finds a linear combination of one set of variables and a linear combination of a second set of variables that are maximally correlated. Given this set of linear combinations, canonical correlation analysis can find subsequent

independent sets of linear combinations, referred to as canonical variables, up to a maximum number equal to the number of variables in the smaller set.

If there are two sets of variables in the analysis and all variables are defined to be numerical, optimal scaling canonical correlation analysis is equivalent to a standard canonical correlation analysis. Although SPSS does not have a canonical correlation analysis procedure, many of the relevant statistics can be obtained from multivariate analysis of variance.

Optimal scaling canonical correlation analysis has various other applications. If you have two sets of variables and one of the sets contains a nominal variable declared as single nominal, optimal scaling canonical correlation analysis results can be interpreted in a similar fashion to regression analysis. If you consider the variable to be multiple nominal, the optimal scaling analysis is an alternative to discriminant analysis. Grouping the variables in more than two sets provides a variety of ways to analyze your data.

Correspondence Analysis

The goal of correspondence analysis is to make biplots for correspondence tables. In a correspondence table, the row and column variables are assumed to represent unordered categories; therefore, the nominal optimal scaling level is always used. Both variables are inspected for their nominal information only. That is, the only consideration is the fact that some objects are in the same category while others are not. Nothing is assumed about the distance or order between categories of the same variable.

One specific use of correspondence analysis is the analysis of two-way contingency tables. If a table has r active rows and c active columns, the number of dimensions in the correspondence analysis solution is the minimum of r minus 1 or c minus 1, whichever is less. In other words, you could perfectly represent the row categories or the column categories of a contingency table in a space of dimensions. Practically speaking, however, you would like to represent the row and column categories of a two-way table in a low-dimensional space, say two dimensions, for the reason that two-dimensional plots are more easily comprehensible than multidimensional spatial representations.

When fewer than the maximum number of possible dimensions is used, the statistics produced in the analysis describe how well the row and column categories are represented in the low-dimensional representation. Provided that the quality of

representation of the two-dimensional solution is good, you can examine plots of the row points and the column points to learn which categories of the row variable are similar, which categories of the column variable are similar, and which row and column categories are similar to each other.

Relation to other Categories procedures. Simple correspondence analysis is limited to two-way tables. If there are more than two variables of interest, you can combine variables to create interaction variables. For example, for the variables *region*, *job*, and *age*, you can combine *region* and *job* to create a new variable *rejob* with the 12 categories shown in the following table. This new variable forms a two-way table with *age* (12 rows, 4 columns), which can be analyzed in correspondence analysis.

Table 1-4
Combinations of region and job

Category Code	Category Definition	Category Code	Category Definition
1	North, intern	7	East, intern
2	North, sales rep	8	East, sales rep
3	North, manager	9	East, manager
4	South, intern	10	West, intern
5	South, sales rep	11	West, sales rep
6	South, manager	12	West, manager

One shortcoming of this approach is that any pair of variables can be combined. We can combine *job* and *age*, yielding another 12-category variable. Or we can combine *region* and *age*, which results in a new 16-category variable. Each of these interaction variables forms a two-way table with the remaining variable. Correspondence analyses of these three tables will not yield identical results, yet each is a valid approach. Furthermore, if there are four or more variables, two-way tables comparing an interaction variable with another interaction variable can be constructed. The number of possible tables to analyze can get quite large, even for a few variables. You can select one of these tables to analyze, or you can analyze all of them. Alternatively, the Multiple Correspondence Analysis procedure can be used to examine all of the variables simultaneously without the need to construct interaction variables.

Relation to standard techniques. The SPSS Crosstabs procedure can also be used to analyze contingency tables, with independence as a common focus in the analyses. However, even in small tables, detecting the cause of departures from independence may be difficult. The utility of correspondence analysis lies in displaying such

patterns for two-way tables of any size. If there is an association between the row and column variables—that is, if the chi-square value is significant—correspondence analysis may help reveal the nature of the relationship.

Multiple Correspondence Analysis

Multiple Correspondence Analysis tries to produce a solution in which objects within the same category are plotted close together and objects in different categories are plotted far apart. Each object is as close as possible to the category points of categories that apply to the object. In this way, the categories divide the objects into homogeneous subgroups. Variables are considered homogeneous when they classify objects in the same categories into the same subgroups.

For a one-dimensional solution, multiple correspondence analysis assigns optimal scale values (category quantifications) to each category of each variable in such a way that overall, on average, the categories have maximum spread. For a two-dimensional solution, multiple correspondence analysis finds a second set of quantifications of the categories of each variable unrelated to the first set, attempting again to maximize spread, and so on. Because categories of a variable receive as many scorings as there are dimensions, the variables in the analysis are assumed to be multiple nominal in optimal scaling level.

Multiple correspondence analysis also assigns scores to the objects in the analysis in such a way that the category quantifications are the averages, or centroids, of the object scores of objects in that category.

Relation to other Categories procedures. Multiple correspondence analysis is also known as homogeneity analysis or dual scaling. It gives comparable, but not identical, results to correspondence analysis when there are only two variables. Correspondence analysis produces unique output summarizing the fit and quality of representation of the solution, including stability information. Thus, correspondence analysis is usually preferable to multiple correspondence analysis in the two-variable case. Another difference between the two procedures is that the input to multiple correspondence analysis is a data matrix, where the rows are objects and the columns are variables, while the input to correspondence analysis can be the same data matrix, a general proximity matrix, or a joint contingency table, which is an aggregated matrix in which both the rows and columns represent categories of variables. Multiple correspondence analysis can also be thought of as principal components analysis of data scaled at the multiple nominal level.

Relation to standard techniques. Multiple correspondence analysis can be thought of as the analysis of a multiway contingency table. Multiway contingency tables can also be analyzed with the SPSS Crosstabs procedure, but Crosstabs gives separate summary statistics for each category of each control variable. With multiple correspondence analysis, it is often possible to summarize the relationship between all of the variables with a single two-dimensional plot. An advanced use of multiple correspondence analysis is to replace the original category values with the optimal scale values from the first dimension and perform a secondary multivariate analysis. Since multiple correspondence analysis replaces category labels with numerical scale values, many different procedures that require numerical data can be applied after the multiple correspondence analysis. For example, the Factor Analysis procedure produces a first principal component that is equivalent to the first dimension of multiple correspondence analysis. The component scores in the first dimension are equal to the object scores, and the squared component loadings are equal to the discrimination measures. The second multiple correspondence analysis dimension, however, is not equal to the second dimension of factor analysis.

Multidimensional Scaling

The use of Multidimensional Scaling is most appropriate when the goal of your analysis is to find the structure in a set of distance measures between objects or cases. This is accomplished by assigning observations to specific locations in a conceptual low-dimensional space so that the distances between points in the space match the given (dis)similarities as closely as possible. The result is a least-squares representation of the objects in that low-dimensional space, which, in many cases, will help you further understand your data.

Relation to other Categories procedures. When you have multivariate data from which you create distances and then analyze with multidimensional scaling, the results are similar to analyzing the data using categorical principal components analysis with object principal normalization. This kind of PCA is also known as principal coordinates analysis.

Relation to standard techniques. The Categories Multidimensional Scaling procedure (PROXSCAL) offers several improvements upon the scaling procedure available in the Base system (ALSCAL). PROXSCAL offers an accelerated algorithm for certain models and allows you to put restrictions on the common space. Moreover,

PROXSCAL attempts to minimize normalized raw stress rather than S-stress (also referred to as **strain**). The normalized raw stress is generally preferred because it is a measure based on the distances, while the S-stress is based on the squared distances.

Aspect Ratio in Optimal Scaling Charts

Aspect ratio in optimal scaling plots is isotropic. In a two-dimensional plot, the distance representing one unit in dimension 1 is equal to the distance representing one unit in dimension 2. If you change the range of a dimension in a two-dimensional plot, the system changes the size of the other dimension to keep the physical distances equal. Isotropic aspect ratio cannot be overridden for the optimal scaling procedures.

Recommended Readings

See the following texts for general information on optimal scaling techniques:

Barlow, R. E., D. J. Bartholomew, D. J. Bremner, and H. D. Brunk. 1972. *Statistical inference under order restrictions*. New York: John Wiley and Sons.

Benzécri, J. P. 1969. Statistical analysis as a tool to make patterns emerge from data. In: *Methodologies of pattern recognition*, S. Watanabe, ed. New York: Academic Press, 35–74.

Bishop, Y. M., S. E. Feinberg, and P. W. Holland. 1975. *Discrete multivariate analysis*. Cambridge, Mass.: MIT Press.

De Leeuw, J. 1984. The Gifi system of nonlinear multivariate analysis. In: *Data analysis and informatics III*, E. Diday, and Coll., eds., 415–424.

De Leeuw, J. 1990. Multivariate analysis with optimal scaling. In: *Progress in Multivariate Analysis*, S. DasGupta, and J. Sethuraman, eds. Calcutta: Indian Statistical Institute.

De Leeuw, J., and J. Van Rijckevorsel. 1980. HOMALS and PRINCALS—Some generalizations of principal components analysis. In: *Data analysis and informatics*, E. Diday, and Coll., eds. Amsterdam: North-Holland, 231–242.

- De Leeuw, J., F. W. Young, and Y. Takane. 1976. Additive structure in qualitative data: An alternating least squares method with optimal scaling features. *Psychometrika*, 41, 471–503.
- Gifi, A. 1990. *Nonlinear multivariate analysis*. Chichester: John Wiley and Sons.
- Heiser, W. J., and J. J. Meulman. 1995. Nonlinear methods for the analysis of homogeneity and heterogeneity. In: *Recent advances in descriptive multivariate analysis*, W. J. Krzanowski, ed. Oxford: Oxford University Press, 51–89.
- Israels, A. 1987. *Eigenvalue techniques for qualitative data*. Leiden: DSWO Press.
- Krzanowski, W. J., and F. H. C. Marriott. 1994. *Multivariate analysis: Part I, distributions, ordination and inference*. London: Edward Arnold.
- Lebart, L., A. Morineau, and K. M. Warwick. 1984. *Multivariate descriptive statistical analysis*. New York: John Wiley and Sons.
- Max, J. 1960. Quantizing for minimum distortion. *Proceedings IEEE (Information Theory)*, 6, 7–12.
- Meulman, J. 1986. *A distance approach to nonlinear multivariate analysis*. Leiden: DSWO Press.
- Meulman, J. 1992. The integration of multidimensional scaling and multivariate analysis with optimal transformations of the variables. *Psychometrika*, 57, 539–565.
- Nishisato, S. 1980. *Analysis of categorical data: Dual scaling and its applications*. Toronto: University of Toronto Press.
- Nishisato, S. 1994. *Elements of dual scaling: An introduction to practical data analysis*. Hillsdale, N.J.: Lawrence Erlbaum Associates, Inc.
- Rao, C. R. 1973. *Linear statistical inference and its applications*. New York: John Wiley & Sons.
- Rao, C. R. 1980. Matrix approximations and reduction of dimensionality in multivariate statistical analysis. In: *Multivariate Analysis, Vol. 5*, P. R. Krishnaiah, ed. Amsterdam: North-Holland, 3–22.
- Roskam, E. E. 1968. *Metric analysis of ordinal data in psychology*. Voorschoten: VAM.
- Shepard, R. N. 1966. Metric structures in ordinal data. *Journal of Mathematical Psychology*, 3, 287–315.

Wolter, K. M. 1985. *Introduction to variance estimation*. Berlin: Springer-Verlag.

Young, F. W. 1981. Quantitative analysis of qualitative data. *Psychometrika*, 46, 357–387.

Categorical Regression (CATREG)

Categorical Regression

Categorical regression quantifies categorical data by assigning numerical values to the categories, resulting in an optimal linear regression equation for the transformed variables. Categorical regression is also known by the acronym CATREG, for *categorical regression*.

Standard linear regression analysis involves minimizing the sum of squared differences between a response (dependent) variable and a weighted combination of predictor (independent) variables. Variables are typically quantitative, with (nominal) categorical data recoded to binary or contrast variables. As a result, categorical variables serve to separate groups of cases, and the technique estimates separate sets of parameters for each group. The estimated coefficients reflect how changes in the predictors affect the response. Prediction of the response is possible for any combination of predictor values.

An alternative approach involves regressing the response on the categorical predictor values themselves. Consequently, one coefficient is estimated for each variable. However, for categorical variables, the category values are arbitrary. Coding the categories in different ways yield different coefficients, making comparisons across analyses of the same variables difficult.

CATREG extends the standard approach by simultaneously scaling nominal, ordinal, and numerical variables. The procedure quantifies categorical variables so that the quantifications reflect characteristics of the original categories. The procedure treats quantified categorical variables in the same way as numerical variables. Using nonlinear transformations allow variables to be analyzed at a variety of levels to find the best-fitting model.

Example. Categorical regression could be used to describe how job satisfaction depends on job category, geographic region, and amount of travel. You might find that high levels of satisfaction correspond to managers and low travel. The resulting

regression equation could be used to predict job satisfaction for any combination of the three independent variables.

Statistics and plots. Frequencies, regression coefficients, ANOVA table, iteration history, category quantifications, correlations between untransformed predictors, correlations between transformed predictors, residual plots, and transformation plots.

Data. CATREG operates on category indicator variables. The category indicators should be positive integers. You can use the Discretization dialog box to convert fractional-value variables and string variables into positive integers.

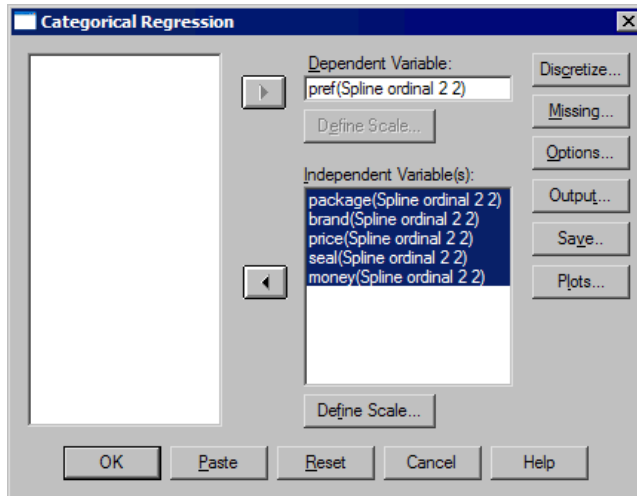
Assumptions. Only one response variable is allowed, but the maximum number of predictor variables is 200. The data must contain at least three valid cases, and the number of valid cases must exceed the number of predictor variables plus one.

Related procedures. CATREG is equivalent to categorical canonical correlation analysis with optimal scaling (OVERALS) with two sets, one of which contains only one variable. Scaling all variables at the numerical level corresponds to standard multiple regression analysis.

To Obtain a Categorical Regression

- ▶ From the menus choose:
 - Analyze
 - Regression
 - Optimal Scaling...

Figure 2-1
Categorical Regression dialog box

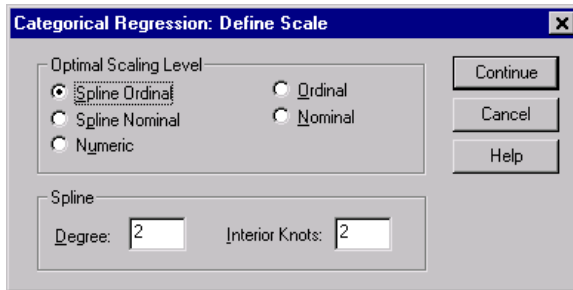


- ▶ Select the dependent variable and independent variable(s).
- ▶ Click OK.
Optionally, change the scaling level for each variable.

Define Scale in Categorical Regression

You can set the optimal scaling level for the dependent and independent variables. By default, they are scaled as second-degree monotonic splines (ordinal) with two interior knots. Additionally, you can set the weight for analysis variables.

Figure 2-2
Define Scale dialog box



Optimal Scaling Level. You can also select the scaling level for quantifying each variable.

- **Spline Ordinal.** The order of the categories of the observed variable is preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. The resulting transformation is a smooth monotonic piecewise polynomial of the chosen degree. The pieces are specified by the user-specified number and procedure-determined placement of the interior knots.
- **Spline Nominal.** The only information in the observed variable that is preserved in the optimally scaled variable is the grouping of objects in categories. The order of the categories of the observed variable is not preserved. Category points will be on a straight line (vector) through the origin. The resulting transformation is a smooth, possibly nonmonotonic, piecewise polynomial of the chosen degree. The pieces are specified by the user-specified number and procedure-determined placement of the interior knots.
- **Ordinal.** The order of the categories of the observed variable is preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. The resulting transformation fits better than the spline ordinal transformation but is less smooth.
- **Nominal.** The only information in the observed variable that is preserved in the optimally scaled variable is the grouping of objects in categories. The order of the categories of the observed variable is not preserved. Category points will be

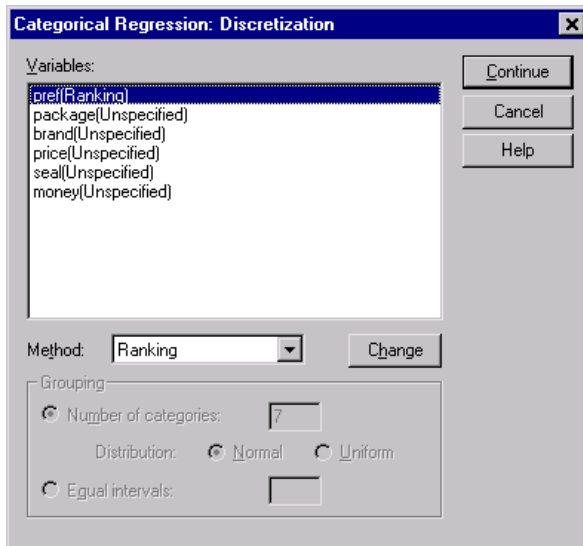
on a straight line (vector) through the origin. The resulting transformation fits better than the spline nominal transformation but is less smooth.

- **Numeric.** Categories are treated as ordered and equally spaced (interval level). The order of the categories and the equal distances between category numbers of the observed variable are preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. When all variables are at the numeric level, the analysis is analogous to standard principal components analysis.

Categorical Regression Discretization

The Discretization dialog box allows you to select a method of recoding your variables. Fractional-value variables are grouped into seven categories (or into the number of distinct values of the variable if this number is less than seven) with an approximately normal distribution unless otherwise specified. String variables are always converted into positive integers by assigning category indicators according to ascending alphanumeric order. Discretization for string variables applies to these integers. Other variables are left alone by default. The discretized variables are then used in the analysis.

Figure 2-3
Discretization dialog box



Method. Choose between grouping, ranking, and multiplying.

- **Grouping.** Recode into a specified number of categories or recode by interval.
- **Ranking.** The variable is discretized by ranking the cases.
- **Multiplying.** The current values of the variable are standardized, multiplied by 10, rounded, and have a constant added so that the lowest discretized value is 1.

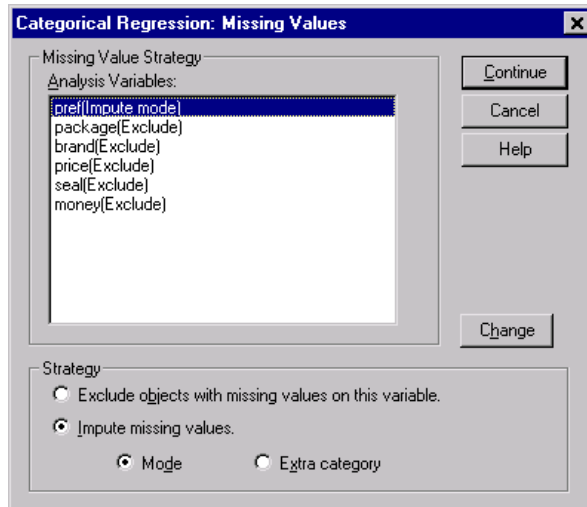
Grouping. The following options are available when discretizing variables by grouping:

- **Number of categories.** Specify a number of categories and whether the values of the variable should follow an approximately normal or uniform distribution across those categories.
- **Equal intervals.** Variables are recoded into categories defined by these equally sized intervals. You must specify the length of the intervals.

Categorical Regression Missing Values

The Missing Values dialog box allows you to choose the strategy for handling missing values in analysis variables and supplementary variables.

Figure 2-4
Missing Values dialog box



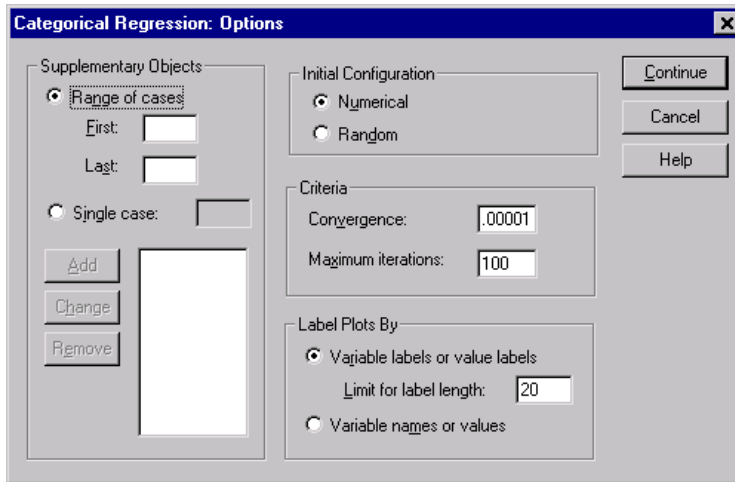
Strategy. Choose to exclude objects with missing values (listwise deletion) or impute missing values (active treatment).

- **Exclude objects with missing values on this variable.** Objects with missing values on the selected variable are excluded from the analysis. This strategy is not available for supplementary variables.
- **Impute missing values.** Objects with missing values on the selected variable have those values imputed. You can choose the method of imputation. Select Mode to replace missing values with the most frequent category. When there are multiple modes, the one with the smallest category indicator is used. Select Extra category to replace missing values with the same quantification of an extra category. This implies that objects with a missing value on this variable are considered to belong to the same (extra) category.

Categorical Regression Options

The Options dialog box allows you to select the initial configuration style, specify iteration and convergence criteria, select supplementary objects, and set the labeling of plots.

Figure 2-5
Options dialog box



Supplementary Objects. This allows you to specify the objects that you want to treat as supplementary. Simply type the number of a supplementary object and click Add. You can not weight supplementary objects (specified weights are ignored).

Initial Configuration. If no variables are treated as nominal, select the Numerical configuration. If at least one variable is treated as nominal, select the Random configuration.

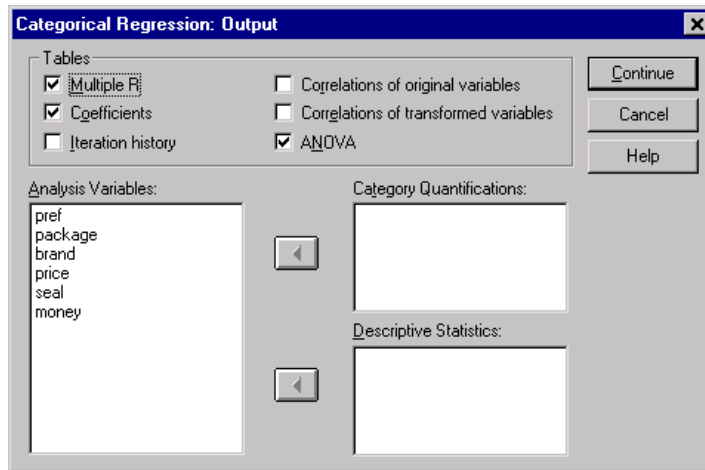
Criteria. You can specify the maximum number of iterations that the regression may go through in its computations. You can also select a convergence criterion value. The regression stops iterating if the difference in total fit between the last two iterations is less than the convergence value or if the maximum number of iterations is reached.

Label Plots By. Allows you to specify whether variables and value labels or variable names and values will be used in the plots. You can also specify a maximum length for labels.

Categorical Regression Output

The Output dialog box allows you to select the statistics to display in the output.

Figure 2-6
Output dialog box



Tables. Produces tables for:

- **Multiple R.** Includes R^2 , adjusted R^2 , and adjusted R^2 taking the optimal scaling into account.
- **Coefficients.** This option gives three tables: a Coefficients table that includes betas, standard error of the betas, t values, and significance; a Coefficients-Optimal Scaling table with the standard error of the betas taking the optimal scaling degrees of freedom into account; and a table with the zero-order, part, and partial correlation, Pratt's relative importance measure for the transformed predictors, and the tolerance before and after transformation.
- **Iteration history.** For each iteration, including the starting values for the algorithm, the multiple R and regression error are shown. The increase in multiple R is listed starting from the first iteration.
- **Correlations of the original variables.** A matrix showing the correlations between the untransformed variables is displayed.

- **Correlations of the transformed variables.** A matrix showing the correlations between the transformed variables is displayed.
- **ANOVA.** This option includes regression and residual sums of squares, mean squares, and F . Two ANOVA tables are displayed: one with degrees of freedom for the regression equal to the number of predictor variables and one with degrees of freedom for the regression taking the optimal scaling into account.

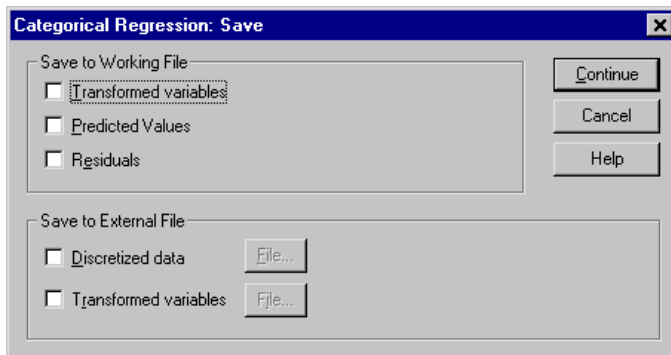
Category Quantifications. Tables showing the transformed values of the selected variables are displayed.

Descriptive Statistics. Tables showing the frequencies, missing values, and modes of the selected variables are displayed.

Categorical Regression Save

The Save dialog box allows you to save results to the working file or an external file.

Figure 2-7
Save dialog box



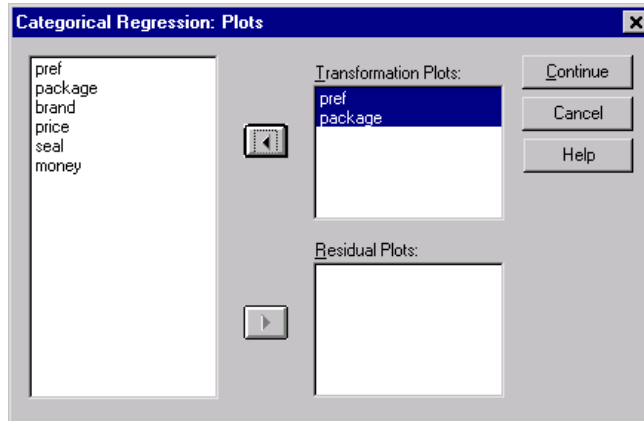
Save to Working File. You can save the transformed values of the variables, model-predicted values, and residuals to the working file.

Save to External File. You can save the discretized data and transformed variables to external files.

Categorical Regression Transformation Plots

The Plots dialog box allows you to specify the variables that will produce transformation and residual plots.

Figure 2-8
Plots dialog box



Transformation Plots. For each of these variables, the category quantifications are plotted against the original category values. Empty categories appear on the horizontal axis but do not affect the computations. These categories are identified by breaks in the line connecting the quantifications.

Residual Plots. For each of these variables, residuals (computed for the dependent variable predicted from all predictor variables except the predictor variable in question) are plotted against category indicators and the optimal category quantifications multiplied with beta against category indicators.

CATREG Command Additional Features

You can customize your categorical regression if you paste your selections into a syntax window and edit the resulting CATREG command syntax. SPSS command language also allows you to:

- Specify rootnames for the transformed variables when saving them to the working data file (with the SAVE subcommand).

Categorical Principal Components Analysis (CATPCA)

This procedure simultaneously quantifies categorical variables while reducing the dimensionality of the data. Categorical principal components analysis is also known by the acronym CATPCA, for *c*ategorical principal components analysis.

The goal of principal components analysis is to reduce an original set of variables into a smaller set of uncorrelated components that represent most of the information found in the original variables. The technique is most useful when a large number of variables prohibits effective interpretation of the relationships between objects (subjects and units). By reducing the dimensionality, you interpret a few components rather than a large number of variables.

Standard principal components analysis assumes linear relationships between numeric variables. On the other hand, the optimal-scaling approach allows variables to be scaled at different levels. Categorical variables are optimally quantified in the specified dimensionality. As a result, nonlinear relationships between variables can be modeled.

Example. Categorical principal components analysis could be used to graphically display the relationship between job category, job division, region, amount of travel (high, medium, and low), and job satisfaction. You might find that two dimensions account for a large amount of variance. The first dimension might separate job category from region, whereas the second dimension might separate job division from amount of travel. You also might find that high job satisfaction is related to a medium amount of travel.

Statistics and plots. Frequencies, missing values, optimal scaling level, mode, variance accounted for by centroid coordinates, vector coordinates, total per variable and per dimension, component loadings for vector-quantified variables, category quantifications and coordinates, iteration history, correlations of the transformed

variables and eigenvalues of the correlation matrix, correlations of the original variables and eigenvalues of the correlation matrix, object scores, category plots, joint category plots, transformation plots, residual plots, projected centroid plots, object plots, biplots, triplots, and component loadings plots.

Data. String variable values are always converted into positive integers by ascending alphanumeric order. User-defined missing values, system-missing values, and values less than 1 are considered missing; you can recode or add a constant to variables with values less than 1 to make them nonmissing.

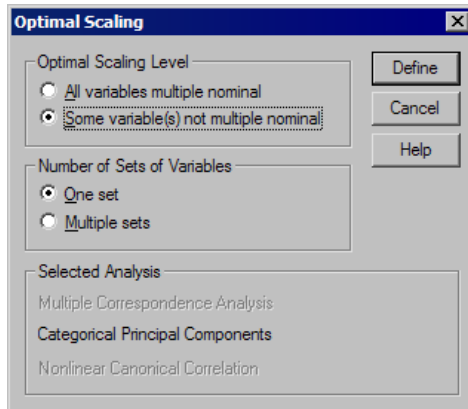
Assumptions. The data must contain at least three valid cases. The analysis is based on positive integer data. The discretization option will automatically categorize a fractional-valued variable by grouping its values into categories with a close to “normal” distribution and will automatically convert values of string variables into positive integers. You can specify other discretization schemes.

Related procedures. Scaling all variables at the numeric level corresponds to standard principal components analysis. Alternate plotting features are available by using the transformed variables in a standard linear principal components analysis. If all variables have multiple nominal scaling levels, categorical principal components analysis is identical to multiple correspondence analysis. If sets of variables are of interest, categorical (nonlinear) canonical correlation analysis should be used.

To Obtain a Categorical Principal Components Analysis

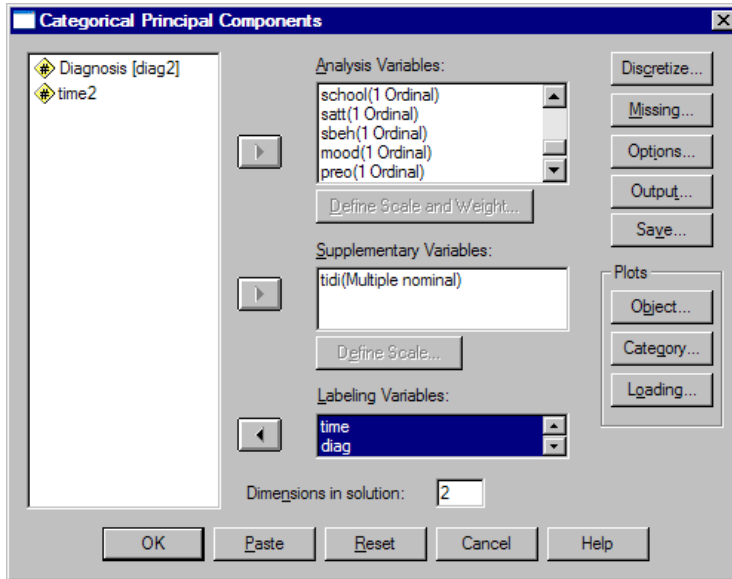
- ▶ From the menus choose:
 - Analyze
 - Data Reduction
 - Optimal Scaling...

Figure 3-1
Optimal Scaling dialog box



- ▶ Select Some variable(s) not multiple nominal.
- ▶ Select One set.
- ▶ Click Define.

Figure 3-2
Categorical Principal Components dialog box



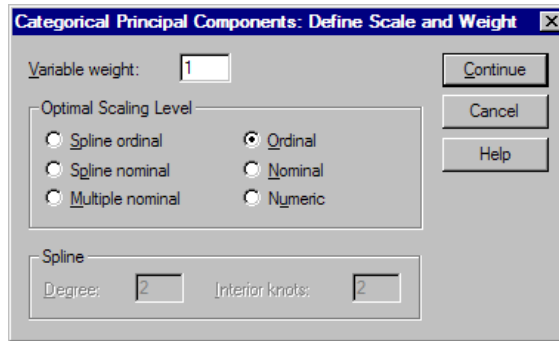
- ▶ Select at least two analysis variables and specify the number of dimensions in the solution.
- ▶ Click OK.

You may optionally specify supplementary variables, which are fitted into the solution found, or labeling variables for the plots.

Define Scale and Weight in CATPCA

You can set the optimal scaling level for analysis variables and supplementary variables. By default, they are scaled as second-degree monotonic splines (ordinal) with two interior knots. Additionally, you can set the weight for analysis variables.

Figure 3-3
Define Scale and Weight dialog box



Variable weight. You can choose to define a weight for each variable. The value specified must be a positive integer. The default value is 1.

Optimal Scaling Level. You can also select the scaling level to be used to quantify each variable.

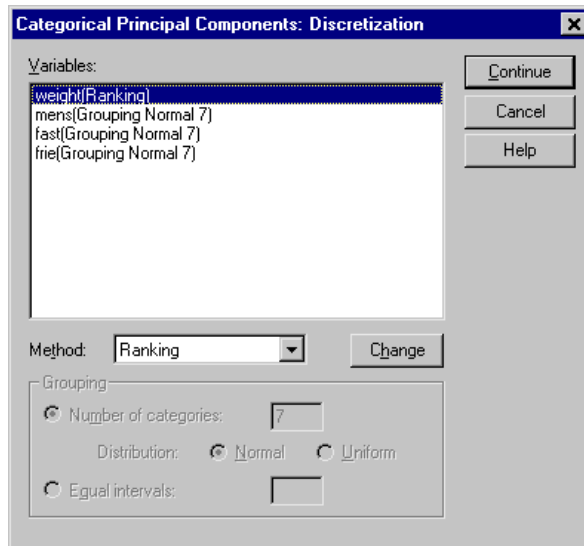
- **Spline ordinal.** The order of the categories of the observed variable is preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. The resulting transformation is a smooth monotonic piecewise polynomial of the chosen degree. The pieces are specified by the user-specified number and procedure-determined placement of the interior knots.
- **Spline nominal.** The only information in the observed variable that is preserved in the optimally scaled variable is the grouping of objects in categories. The order of the categories of the observed variable is not preserved. Category points will be on a straight line (vector) through the origin. The resulting transformation is a smooth, possibly nonmonotonic, piecewise polynomial of the chosen degree. The pieces are specified by the user-specified number and procedure-determined placement of the interior knots.
- **Multiple nominal.** The only information in the observed variable that is preserved in the optimally scaled variable is the grouping of objects in categories. The order of the categories of the observed variable is not preserved. Category points will be in the centroid of the objects in the particular categories. *Multiple* indicates that different sets of quantifications are obtained for each dimension.

- **Ordinal.** The order of the categories of the observed variable is preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. The resulting transformation fits better than the spline ordinal transformation but is less smooth.
- **Nominal.** The only information in the observed variable that is preserved in the optimally scaled variable is the grouping of objects in categories. The order of the categories of the observed variable is not preserved. Category points will be on a straight line (vector) through the origin. The resulting transformation fits better than the spline nominal transformation but is less smooth.
- **Numeric.** Categories are treated as ordered and equally spaced (interval level). The order of the categories and the equal distances between category numbers of the observed variable are preserved in the optimally scaled variable. Category points will be on a straight line (vector) through the origin. When all variables are at the numeric level, the analysis is analogous to standard principal components analysis.

Categorical Principal Components Analysis Discretization

The Discretization dialog box allows you to select a method of recoding your variables. Fractional-valued variables are grouped into seven categories (or into the number of distinct values of the variable if this number is less than seven) with an approximately normal distribution, unless specified otherwise. String variables are always converted into positive integers by assigning category indicators according to ascending alphanumeric order. Discretization for string variables applies to these integers. Other variables are left alone by default. The discretized variables are then used in the analysis.

Figure 3-4
Discretization dialog box



Method. Choose between grouping, ranking, and multiplying.

- **Grouping.** Recode into a specified number of categories or recode by interval.
- **Ranking.** The variable is discretized by ranking the cases.
- **Multiplying.** The current values of the variable are standardized, multiplied by 10, rounded, and have a constant added such that the lowest discretized value is 1.

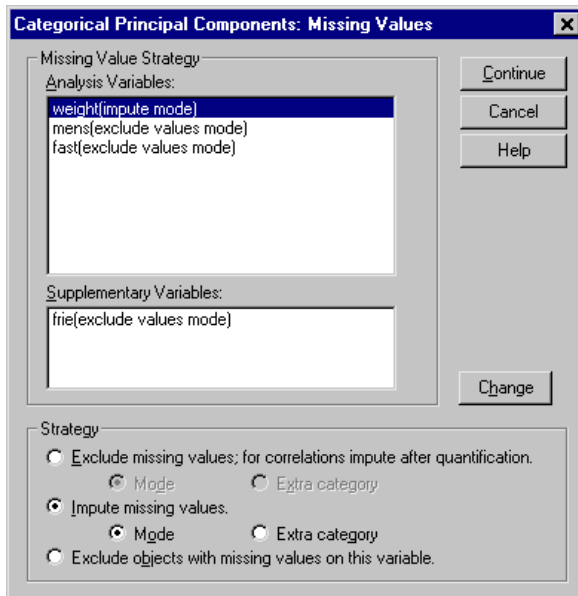
Grouping. The following options are available when you are discretizing variables by grouping:

- **Number of categories.** Specify a number of categories and whether the values of the variable should follow an approximately normal or uniform distribution across those categories.
- **Equal intervals.** Variables are recoded into categories defined by these equally sized intervals. You must specify the length of the intervals.

Categorical Principal Components Analysis Missing Values

The Missing Values dialog box allows you to choose the strategy for handling missing values in analysis variables and supplementary variables.

Figure 3-5
Missing Values dialog box



Missing Value Strategy. Choose to exclude missing values (passive treatment), impute missing values (active treatment), or exclude objects with missing values (listwise deletion).

- **Exclude missing values; for correlations impute after quantification.** Objects with missing values on the selected variable do not contribute to the analysis for this variable. If all variables are given passive treatment, then objects with missing values on all variables are treated as supplementary. If correlations are specified in the Output dialog box, then (after analysis) missing values are imputed with the most frequent category, or mode, of the variable for the correlations of the original variables. For the correlations of the optimally scaled variables, you can choose the method of imputation. Select Mode to replace missing values with the mode of the optimally scaled variable. Select Extra category to replace missing values with

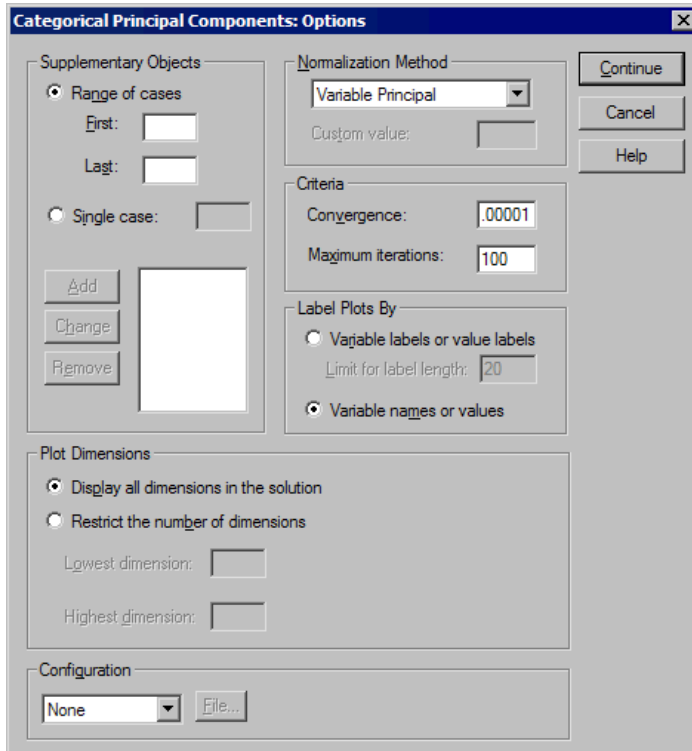
the quantification of an extra category. This implies that objects with a missing value on this variable are considered to belong to the same (extra) category.

- **Impute missing values.** Objects with missing values on the selected variable have those values imputed. You can choose the method of imputation. Select Mode to replace missing values with the most frequent category. When there are multiple modes, the one with the smallest category indicator is used. Select Extra category to replace missing values with the same quantification of an extra category. This implies that objects with a missing value on this variable are considered to belong to the same (extra) category.
- **Exclude objects with missing values on this variable.** Objects with missing values on the selected variable are excluded from the analysis. This strategy is not available for supplementary variables.

Categorical Principal Components Analysis Options

The Options dialog box allows you to select the initial configuration, specify iteration and convergence criteria, select a normalization method, choose the method for labeling plots, and specify supplementary objects.

Figure 3-6
Options dialog box



Supplementary Objects. Specify the case number of the object, or the first and last case numbers of a range of objects, that you want to make supplementary and then click Add. Continue until you have specified all of your supplementary objects. If an object is specified as supplementary, then case weights are ignored for that object.

Normalization Method. You can specify one of five options for normalizing the object scores and the variables. Only one normalization method can be used in a given analysis.

- **Variable Principal.** This option optimizes the association between variables. The coordinates of the variables in the object space are the component loadings (correlations with principal components, such as dimensions and object scores). This is useful when you are primarily interested in the correlation between the variables.

- **Object Principal.** This option optimizes distances between objects. This is useful when you are primarily interested in differences or similarities between the objects.
- **Symmetrical.** Use this normalization option if you are primarily interested in the relation between objects and variables.
- **Independent.** Use this normalization option if you want to examine distances between objects and correlations between variables separately.
- **Custom.** You can specify any real value in the closed interval $[-1, 1]$. A value of 1 is equal to the Object Principal method, a value of 0 is equal to the Symmetrical method, and a value of -1 is equal to the Variable Principal method. By specifying a value greater than -1 and less than 1, you can spread the eigenvalue over both objects and variables. This method is useful for making a tailor-made biplot or triplot.

Criteria. You can specify the maximum number of iterations the procedure can go through in its computations. You can also select a convergence criterion value. The algorithm stops iterating if the difference in total fit between the last two iterations is less than the convergence value or if the maximum number of iterations is reached.

Label Plots By. Allows you to specify whether variables and value labels or variable names and values will be used in the plots. You can also specify a maximum length for labels.

Plot Dimensions. Allows you to control the dimensions displayed in the output.

- **Display all dimensions in the solution.** All dimensions in the solution are displayed in a scatterplot matrix.
- **Restrict the number of dimensions.** The displayed dimensions are restricted to plotted pairs. If you restrict the dimensions, you must select the lowest and highest dimensions to be plotted. The lowest dimension can range from 1 to the number of dimensions in the solution minus 1 and is plotted against higher dimensions. The highest dimension value can range from 2 to the number of dimensions in the solution and indicates the highest dimension to be used in plotting the dimension pairs. This specification applies to all requested multidimensional plots.

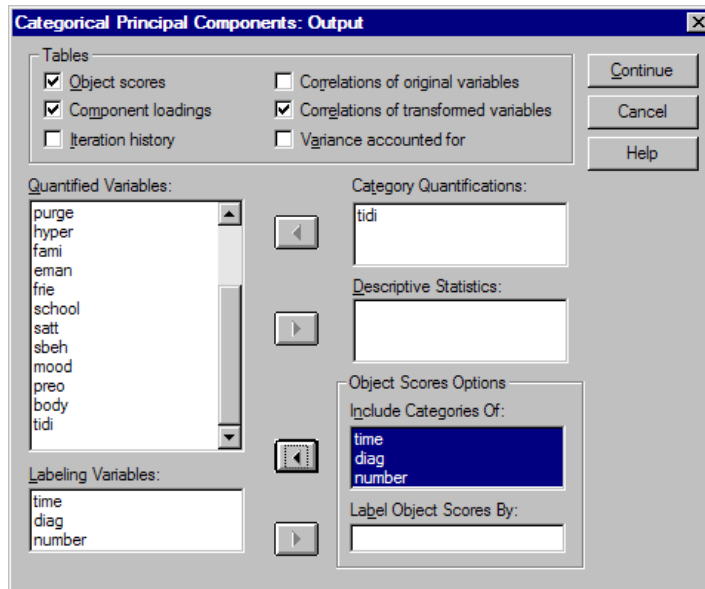
Configuration. You can read data from a file containing the coordinates of a configuration. The first variable in the file should contain the coordinates for the first dimension, the second variable should contain the coordinates for the second dimension, and so on.

- **Initial.** The configuration in the file specified will be used as the starting point of the analysis.
- **Fixed.** The configuration in the file specified will be used to fit in the variables. The variables that are fitted in must be selected as analysis variables, but because the configuration is fixed, they are treated as supplementary variables (so they do not need to be selected as supplementary variables).

Categorical Principal Components Analysis Output

The Output dialog box allows you to produce tables for object scores, component loadings, iteration history, correlations of original and transformed variables, the variance accounted for per variable and per dimension, category quantifications for selected variables, and descriptive statistics for selected variables.

Figure 3-7
Output dialog box



Object scores. Displays the object scores and has the following options:

- **Include Categories Of.** Displays the category indicators of the analysis variables selected.
- **Label Object Scores By.** From the list of variables specified as labeling variables, you can select one to label the objects.

Component loadings. Displays the component loadings for all variables that were not given multiple nominal scaling levels.

Iteration history. For each iteration, the variance accounted for, loss, and increase in variance accounted for are shown.

Correlations of original variables. Shows the correlation matrix of the original variables and the eigenvalues of that matrix.

Correlations of transformed variables. Shows the correlation matrix of the transformed (optimally scaled) variables and the eigenvalues of that matrix.

Variance accounted for. Displays the amount of variance accounted for by centroid coordinates, vector coordinates, and total (centroid and vector coordinates combined) per variable and per dimension.

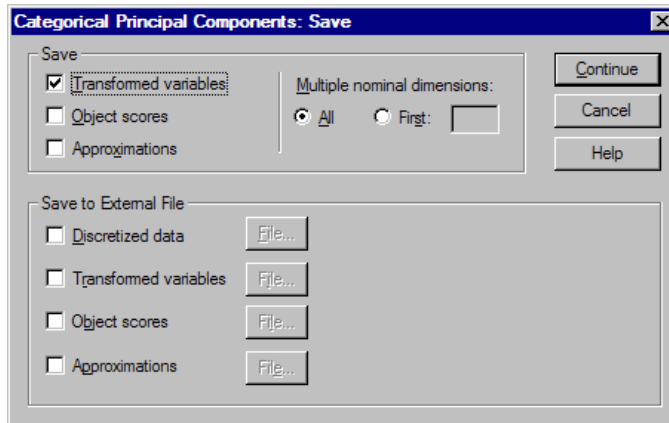
Category Quantifications. Gives the category quantifications and coordinates for each dimension of the variable(s) selected.

Descriptive Statistics. Displays frequencies, number of missing values, and mode of the variable(s) selected.

Categorical Principal Components Analysis Save

The Save dialog box allows you to add the transformed variables, object scores, and approximations to the working data file or as new variables in external files and save the discretized data as new variables in an external data file.

Figure 3-8
Save dialog box



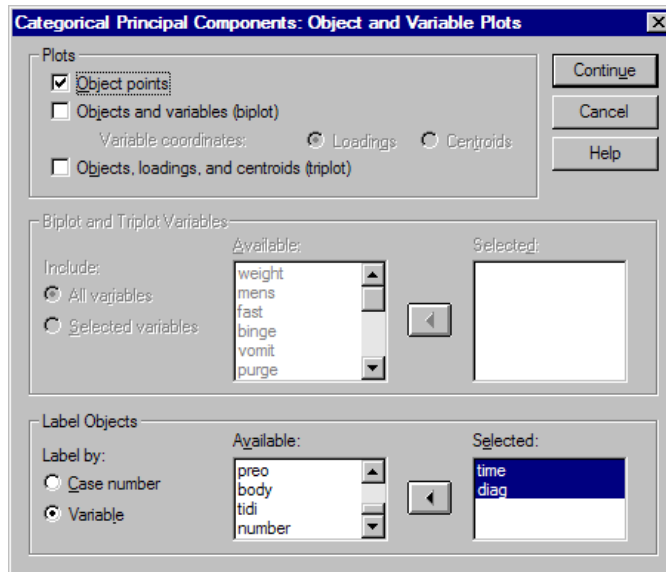
Save. Save selections to the working data file. If any variable has been given the multiple nominal scaling level, the number of dimensions to be saved must be specified.

Save to External File. Save selections to a new external file. Specify a filename for each selected option by clicking File. Each file specified must have a different name.

Categorical Principal Components Analysis Object Plots

The Object and Variable Plots dialog box allows you to specify the types of plots desired and the variables for which plots will be produced.

Figure 3-9
Object and Variable Plots dialog box



Object points. A plot of the object points is displayed.

Objects and variables (biplot). The object points are plotted with your choice of the variable coordinates—component loadings or variable centroids.

Objects, loadings, and centroids (triplot). The object points are plotted with the centroids of multiple nominal-scaling-level variables and the component loadings of other variables.

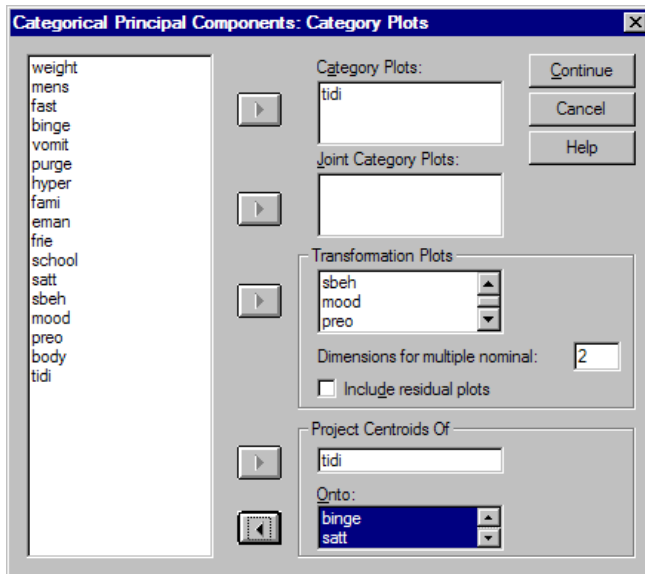
Biplot and Triplot Variables. You can choose to use all variables for the biplots and triplots or select a subset.

Label Objects. You can choose to have objects labeled with the categories of selected variables (you may choose category indicator values or value labels in the Options dialog box) or with their case numbers. One plot is produced per variable if Variable is selected.

Categorical Principal Components Analysis Category Plots

The Category Plots dialog box allows you to specify the types of plots desired and the variables for which plots will be produced.

Figure 3-10
Category Plots dialog box



Category Plots. For each variable selected, a plot of the centroid and vector coordinates is plotted. For variables with multiple nominal scaling levels, categories are in the centroids of the objects in the particular categories. For all other scaling levels, categories are on a vector through the origin.

Joint Category Plots. This is a single plot of the centroid and vector coordinates of each selected variable.

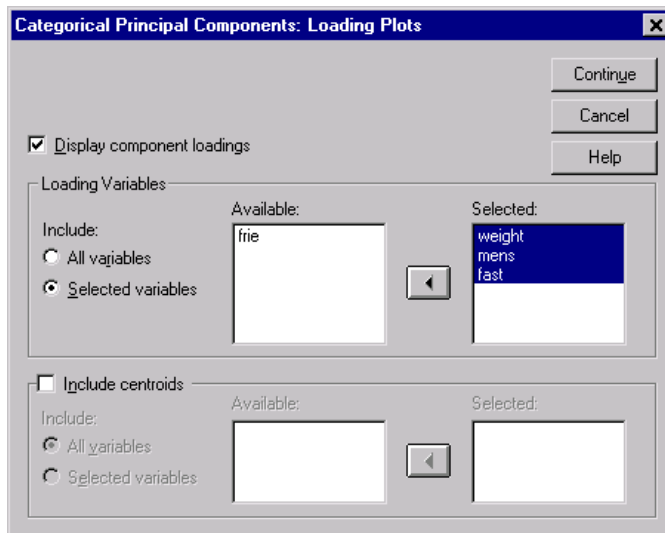
Transformation Plots. Displays a plot of the optimal category quantifications versus the category indicators. You can specify the number of dimensions desired for variables with multiple nominal scaling levels; one plot will be generated for each dimension. You can also choose to display residual plots for each variable selected.

Project Centroids Of. You may choose a variable and project its centroids onto selected variables. Variables with multiple nominal scaling levels cannot be selected to project on. When this plot is requested, a table with the coordinates of the projected centroids is also displayed.

Categorical Principal Components Analysis Loading Plots

The Loading Plots dialog box allows you to specify the variables which will be included in the plot, and whether or not to include centroids in the plot.

Figure 3-11
Loading Plots dialog box



Display component loadings. If selected, a plot of the component loadings is displayed.

Loading Variables. You can choose to use all variables for the component loadings plot or select a subset.

Include centroids. Variables with multiple nominal scaling levels do not have component loadings, but you may choose to include the centroids of those variables in the plot. You can choose to use all multiple nominal variables or select a subset.

CATPCA Command Additional Features

You can customize your categorical principal components analysis if you paste your selections into a syntax window and edit the resulting `CATPCA` command syntax. SPSS command language also allows you to:

- Specify rootnames for the transformed variables, object scores, and approximations when saving them to the working data file (with the `SAVE` subcommand).
- Specify a maximum length for labels for each plot separately (with the `PLOT` subcommand).
- Specify a separate variable list for residual plots (with the `PLOT` subcommand).

Nonlinear Canonical Correlation Analysis (OVERALS)

Nonlinear Canonical Correlation Analysis

Nonlinear canonical correlation analysis corresponds to categorical canonical correlation analysis with optimal scaling. The purpose of this procedure is to determine how similar sets of categorical variables are to one another. Nonlinear canonical correlation analysis is also known by the acronym OVERALS.

Standard canonical correlation analysis is an extension of multiple regression, where the second set does not contain a single response variable, but multiple ones. The goal is to explain as much as possible of the variance in the relationships among two sets of numerical variables in a low dimensional space. Initially, the variables in each set are linearly combined such that the linear combinations have a maximal correlation. Given these combinations, subsequent linear combinations are determined that are uncorrelated with the previous combinations and that have the largest correlation possible.

The optimal scaling approach expands the standard analysis in three crucial ways. First, OVERALS allows more than two sets of variables. Second, variables can be scaled as either nominal, ordinal, or numerical. As a result, nonlinear relationships between variables can be analyzed. Finally, instead of maximizing correlations between the variable sets, the sets are compared to an unknown compromise set defined by the object scores.

Example. Categorical canonical correlation analysis with optimal scaling could be used to graphically display the relationship between one set of variables containing job category and years of education and another set of variables containing region of residence and gender. You might find that years of education and region of residence

discriminate better than the remaining variables. You might also find that years of education discriminates best on the first dimension.

Statistics and plots. Frequencies, centroids, iteration history, object scores, category quantifications, weights, component loadings, single and multiple fit, object scores plots, category coordinates plots, component loadings plots, category centroids plots, transformation plots.

Data. Use integers to code categorical variables (nominal or ordinal scaling level). To minimize output, use consecutive integers beginning with 1 to code each variable. Variables scaled at the numerical level should not be recoded to consecutive integers. To minimize output, for each variable scaled at the numerical level, subtract the smallest observed value from every value and add 1. Fractional values are truncated after the decimal.

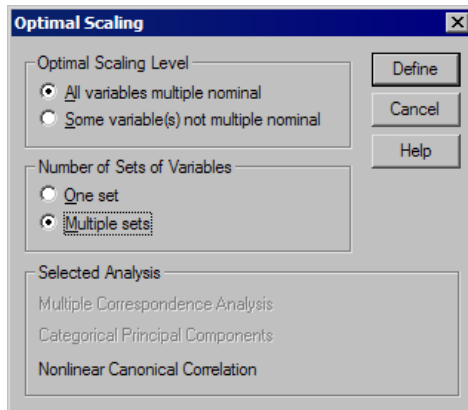
Assumptions. Variables can be classified into two or more sets. Variables in the analysis are scaled as multiple nominal, single nominal, ordinal, or numerical. The maximum number of dimensions used in the procedure depends on the optimal scaling level of the variables. If all variables are specified as ordinal, single nominal, or numerical, the maximum number of dimensions is the minimum of the number of observations minus 1 and the total number of variables. However, if only two sets of variables are defined, the maximum number of dimensions is the number of variables in the smaller set. If some variables are multiple nominal, the maximum number of dimensions is the total number of multiple nominal categories plus the number of nonmultiple nominal variables minus the number of multiple nominal variables. For example, if the analysis involves five variables, one of which is multiple nominal with four categories, the maximum number of dimensions is $(4 + 4 - 1)$, or 7. If you specify a number greater than the maximum, the maximum value is used.

Related procedures. If each set contains one variable, nonlinear canonical correlation analysis is equivalent to principal components analysis with optimal scaling. If each of these variables is multiple nominal, the analysis corresponds to multiple correspondence analysis. If two sets of variables are involved and one of the sets contains only one variable, the analysis is identical to categorical regression with optimal scaling.

To Obtain a Nonlinear Canonical Correlation Analysis

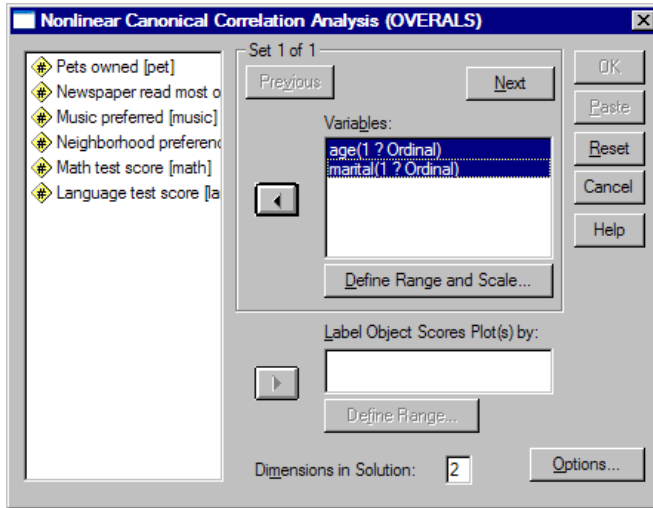
- ▶ From the menus choose:
 - Analyze
 - Data Reduction
 - Optimal Scaling...

Figure 4-1
Optimal Scaling dialog box



- ▶ Select either All variables multiple nominal or Some variable(s) not multiple nominal.
- ▶ Select Multiple sets.
- ▶ Click Define.

Figure 4-2
Nonlinear Canonical Correlation Analysis (OVERALS) dialog box



- ▶ Define at least two sets of variables. Select the variable(s) that you want to include in the first set. To move to the next set, click **Next**, and select the variables that you want to include in the second set. You can add additional sets as desired. Click **Previous** to return to the previously defined variable set.
- ▶ Define the value range and measurement scale (optimal scaling level) for each selected variable.
- ▶ Click **OK**.

Optionally, you can:

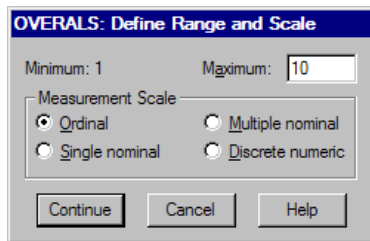
- Select one or more variables to provide point labels for object scores plots. Each variable produces a separate plot, with the points labeled by the values of that variable. You must define a range for each of these plot label variables. Using the dialog box, a single variable cannot be used both in the analysis and as a labeling variable. If labeling the object scores plot with a variable used in the analysis is desired, use the **Compute** facility on the **Transform** menu to create a copy of

that variable. Use the new variable to label the plot. Alternatively, command syntax can be used.

- Specify the number of dimensions you want in the solution. In general, choose as few dimensions as needed to explain most of the variation. If the analysis involves more than two dimensions, SPSS produces three-dimensional plots of the first three dimensions. Other dimensions can be displayed by editing the chart.

Define Range and Scale

Figure 4-3
Define Range and Scale dialog box



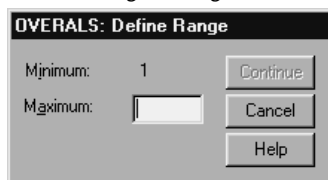
You must define a range for each variable. The maximum value specified must be an integer. Fractional data values are truncated in the analysis. A category value that is outside of the specified range is ignored in the analysis. To minimize output, use the Automatic Recode facility on the Transform menu to create consecutive categories beginning with 1 for variables treated as nominal or ordinal. Recoding to consecutive integers is not recommended for variables scaled at the numerical level. To minimize output for variables treated as numerical, for each variable, subtract the minimum value from every value and add 1.

You must also select the scaling to be used to quantify each variable.

- **Ordinal.** The order of the categories of the observed variable is preserved in the quantified variable.
- **Single nominal.** In the quantified variable, objects in the same category receive the same score.
- **Multiple nominal.** The quantifications can be different for each dimension.
- **Discrete numeric.** Categories are treated as ordered and equally spaced. The differences between category numbers and the order of the categories of the observed variable are preserved in the quantified variable.

Define Range

Figure 4-4
Define Range dialog box



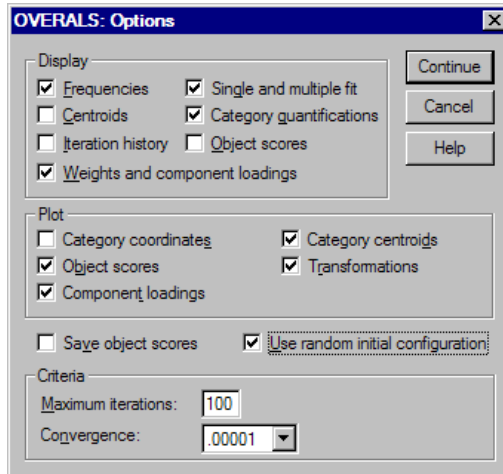
You must define a range for each variable. The maximum value specified must be an integer. Fractional data values are truncated in the analysis. A category value that is outside of the specified range is ignored in the analysis. To minimize output, use the Automatic Recode facility on the Transform menu to create consecutive categories beginning with 1.

You must also define a range for each variable used to label the object scores plots. However, labels for categories with data values outside of the defined range for the variable do appear on the plots.

Nonlinear Canonical Correlation Analysis Options

The Options dialog box allows you to select optional statistics and plots, save object scores as new variables in the working data file, specify iteration and convergence criteria, and specify an initial configuration for the analysis.

Figure 4-5
Options dialog box



Display. Available statistics include marginal frequencies (counts), centroids, iteration history, weights and component loadings, category quantifications, object scores, and single and multiple fit statistics.

- **Centroids.** Category quantifications, and the projected and the actual averages of the object scores for the objects (cases) included in each set for those belonging to the same category of the variable.
- **Weights and Component Loadings.** The regression coefficients in each dimension for every quantified variable in a set, where the object scores are regressed on the quantified variables, and the projection of the quantified variable in the object space. Provides an indication of the contribution each variable makes to the dimension within each set.
- **Single and Multiple Fit.** Measures of goodness of fit of the single- and multiple-category coordinates/category quantifications with respect to the objects.
- **Category Quantifications.** Optimal scale values assigned to the categories of a variable.
- **Object Scores.** Optimal score assigned to an object (case) in a particular dimension.

Plot. You can produce plots of category coordinates, object scores, component loadings, category centroids, and transformations.

Save object scores. You can save the object scores as new variables in the working data file. Object scores are saved for the number of dimensions specified in the main dialog box.

Use random initial configuration. A random initial configuration should be used if all or some of the variables are single nominal. If this option is not selected, a nested initial configuration is used.

Criteria. You can specify the maximum number of iterations the nonlinear canonical correlation analysis can go through in its computations. You can also select a convergence criterion value. The analysis stops iterating if the difference in total fit between the last two iterations is less than the convergence value or if the maximum number of iterations is reached.

OVERALS Command Additional Features

You can customize your nonlinear canonical correlation analysis if you paste your selections into a syntax window and edit the resulting `OVERALS` command syntax. SPSS command language also allows you to:

- Specify the dimension pairs to be plotted, rather than plotting all extracted dimensions (using the `NDIM` keyword on the `PLOT` subcommand).
- Specify the number of value label characters used to label points on the plots (with the `PLOT` subcommand).
- Designate more than five variables as labeling variables for object scores plots (with the `PLOT` subcommand).
- Select variables used in the analysis as labeling variables for the object scores plots (with the `PLOT` subcommand).
- Select variables to provide point labels for the quantification score plot (with the `PLOT` subcommand).
- Specify the number of cases to be included in the analysis, if you do not want to use all cases in the working data file (with the `NOOBSERVATIONS` subcommand).
- Specify rootnames for variables created by saving object scores (with the `SAVE` subcommand).
- Specify the number of dimensions to be saved, rather than saving all extracted dimensions (with the `SAVE` subcommand).

- Write category quantifications to a matrix file (using the `MATRIX` subcommand).
- Produce low-resolution plots that may be easier to read than the usual high-resolution plots (using the `SET` command).
- Produce centroid and transformation plots for specified variables only (with the `PLOT` subcommand).

Correspondence Analysis

One of the goals of correspondence analysis is to describe the relationships between two nominal variables in a correspondence table in a low-dimensional space, while simultaneously describing the relationships between the categories for each variable. For each variable, the distances between category points in a plot reflect the relationships between the categories with similar categories plotted close to each other. Projecting points for one variable on the vector from the origin to a category point for the other variable describe the relationship between the variables.

An analysis of contingency tables often includes examining row and column profiles and testing for independence via the chi-square statistic. However, the number of profiles can be quite large, and the chi-square test does not reveal the dependence structure. The Crosstabs procedure offers several measures of association and tests of association but cannot graphically represent any relationships between the variables.

Factor analysis is a standard technique for describing relationships between variables in a low-dimensional space. However, factor analysis requires interval data, and the number of observations should be five times the number of variables. Correspondence analysis, on the other hand, assumes nominal variables and can describe the relationships between categories of each variable, as well as the relationship between the variables. In addition, correspondence analysis can be used to analyze any table of positive correspondence measures.

Example. Correspondence analysis could be used to graphically display the relationship between staff category and smoking habits. You might find that with regard to smoking, junior managers differ from secretaries, but secretaries do not differ from senior managers. You might also find that heavy smoking is associated with junior managers, whereas light smoking is associated with secretaries.

Statistics and plots. Correspondence measures, row and column profiles, singular values, row and column scores, inertia, mass, row and column score confidence statistics, singular value confidence statistics, transformation plots, row point plots, column point plots, and biplots.

Data. Categorical variables to be analyzed are scaled nominally. For aggregated data or for a correspondence measure other than frequencies, use a weighting variable with positive similarity values. Alternatively, for table data, use syntax to read the table.

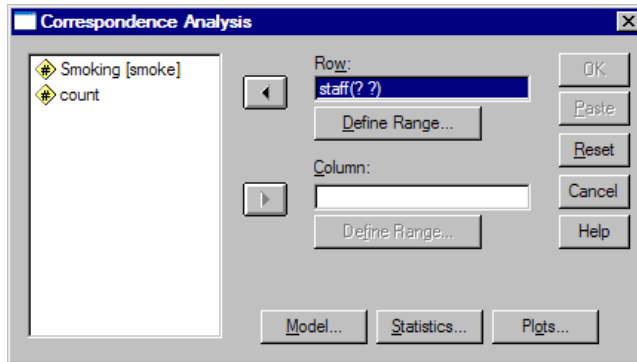
Assumptions. The maximum number of dimensions used in the procedure depends on the number of active rows and column categories and the number of equality constraints. If no equality constraints are used and all categories are active, the maximum dimensionality is one fewer than the number of categories for the variable with the fewest categories. For example, if one variable has five categories and the other has four, the maximum number of dimensions is three. Supplementary categories are not active. For example, if one variable has five categories, two of which are supplementary, and the other variable has four categories, the maximum number of dimensions is two. Treat all sets of categories that are constrained to be equal as one category. For example, if a variable has five categories, three of which are constrained to be equal, that variable should be treated as having three categories when determining the maximum dimensionality. Two of the categories are unconstrained, and the third category corresponds to the three constrained categories. If you specify a number of dimensions greater than the maximum, the maximum value is used.

Related procedures. If more than two variables are involved, use multiple correspondence analysis. If the variables should be scaled ordinally, use categorical principal components analysis.

To Obtain a Correspondence Analysis

- ▶ From the menus choose:
 - Analyze
 - Data Reduction...
 - Correspondence Analysis...

Figure 5-1
Correspondence Analysis dialog box

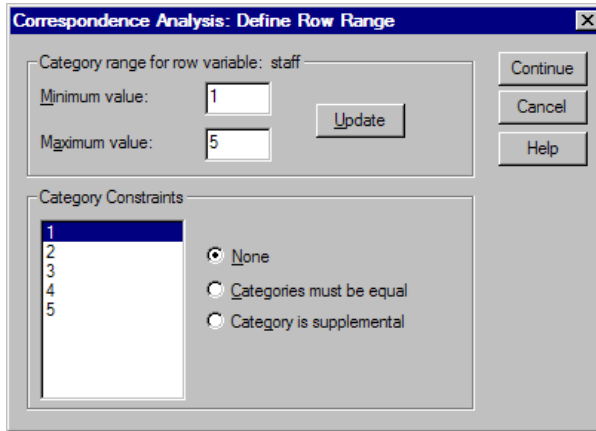


- ▶ Select a row variable.
- ▶ Select a column variable.
- ▶ Define the ranges for the variables.
- ▶ Click OK.

Define Row Range in Correspondence Analysis

You must define a range for the row variable. The minimum and maximum values specified must be integers. Fractional data values are truncated in the analysis. A category value that is outside of the specified range is ignored in the analysis.

Figure 5-2
Define Row Range dialog box



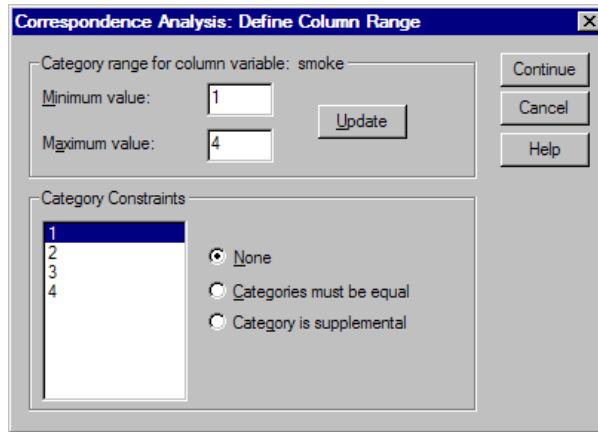
All categories are initially unconstrained and active. You can constrain row categories to equal other row categories, or you can define a row category as supplementary.

- **Categories must be equal.** Categories must have equal scores. Use equality constraints if the obtained order for the categories is undesirable or counterintuitive. The maximum number of row categories that can be constrained to be equal is the total number of active row categories minus 1. To impose different equality constraints on sets of categories, use syntax. For example, use syntax to constrain categories 1 and 2 to be equal and categories 3 and 4 to be equal.
- **Category is supplemental.** Supplementary categories do not influence the analysis but are represented in the space defined by the active categories. Supplementary categories play no role in defining the dimensions. The maximum number of supplementary row categories is the total number of row categories minus 2.

Define Column Range in Correspondence Analysis

You must define a range for the column variable. The minimum and maximum values specified must be integers. Fractional data values are truncated in the analysis. A category value that is outside of the specified range is ignored in the analysis.

Figure 5-3
Define Column Range dialog box



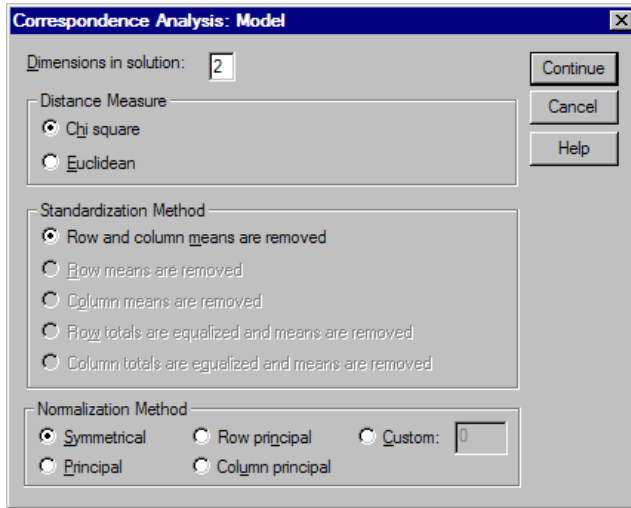
All categories are initially unconstrained and active. You can constrain column categories to equal other column categories, or you can define a column category as supplementary.

- **Categories must be equal.** Categories must have equal scores. Use equality constraints if the obtained order for the categories is undesirable or counterintuitive. The maximum number of column categories that can be constrained to be equal is the total number of active column categories minus 1. To impose different equality constraints on sets of categories, use syntax. For example, use syntax to constrain categories 1 and 2 to be equal and categories 3 and 4 to be equal.
- **Category is supplemental.** Supplementary categories do not influence the analysis but are represented in the space defined by the active categories. Supplementary categories play no role in defining the dimensions. The maximum number of supplementary column categories is the total number of column categories minus 2.

Correspondence Analysis Model

The Model dialog box allows you to specify the number of dimensions, the distance measure, the standardization method, and the normalization method.

Figure 5-4
Model dialog box



Dimensions in solution. Specify the number of dimensions. In general, choose as few dimensions as needed to explain most of the variation. The maximum number of dimensions depends on the number of active categories used in the analysis and on the equality constraints. The maximum number of dimensions is the smaller of:

- The number of active row categories minus the number of row categories constrained to be equal, plus the number of constrained row category sets.
- The number of active column categories minus the number of column categories constrained to be equal, plus the number of constrained column category sets.

Distance Measure. You can select the measure of distance among the rows and columns of the correspondence table. Choose one of the following alternatives:

- **Chi square.** Use a weighted profile distance, where the weight is the mass of the rows or columns. This measure is required for standard correspondence analysis.
- **Euclidean.** Use the square root of the sum of squared differences between pairs of rows and pairs of columns.

Standardization Method. Choose one of the following alternatives:

- **Row and column means are removed.** Both the rows and columns are centered. This method is required for standard correspondence analysis.

- **Row means are removed.** Only the rows are centered.
- **Column means are removed.** Only the columns are centered.
- **Row totals are equalized and means are removed.** Before centering the rows, the row margins are equalized.
- **Column totals are equalized and means are removed.** Before centering the columns, the column margins are equalized.

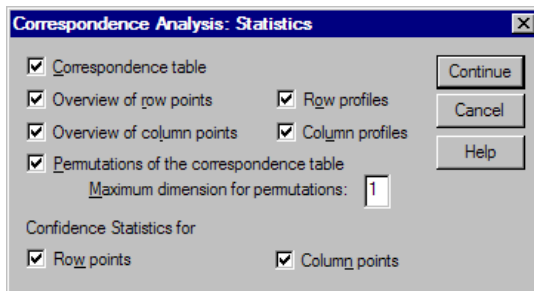
Normalization Method. Choose one of the following alternatives:

- **Symmetrical.** For each dimension, the row scores are the weighted average of the column scores divided by the matching singular value, and the column scores are the weighted average of row scores divided by the matching singular value. Use this method if you want to examine the differences or similarities between the categories of the two variables.
- **Principal.** The distances between row points and column points are approximations of the distances in the correspondence table according to the selected distance measure. Use this method if you want to examine differences between categories of either or both variables instead of differences between the two variables.
- **Row principal.** The distances between row points are approximations of the distances in the correspondence table according to the selected distance measure. The row scores are the weighted average of the column scores. Use this method if you want to examine differences or similarities between categories of the row variable.
- **Column principal.** The distances between column points are approximations of the distances in the correspondence table according to the selected distance measure. The column scores are the weighted average of the row scores. Use this method if you want to examine differences or similarities between categories of the column variable.
- **Custom.** You must specify a value between -1 and 1 . A value of -1 corresponds to column principal. A value of 1 corresponds to row principal. A value of 0 corresponds to symmetrical. All other values spread the inertia over both the row and column scores to varying degrees. This method is useful for making tailor-made biplots.

Correspondence Analysis Statistics

The Statistics dialog box allows you to specify the numerical output produced.

Figure 5-5
Statistics dialog box



Correspondence table. A crosstabulation of the input variables with row and column marginal totals.

Overview of row points. For each row category, the scores, mass, inertia, contribution to the inertia of the dimension, and the contribution of the dimension to the inertia of the point.

Overview of column points. For each column category, the scores, mass, inertia, contribution to the inertia of the dimension, and the contribution of the dimension to the inertia of the point.

Row profiles. For each row category, the distribution across the categories of the column variable.

Column profiles. For each column category, the distribution across the categories of the row variable.

Permutations of the correspondence table. The correspondence table reorganized such that the rows and columns are in increasing order according to the scores on the first dimension. Optionally, you can specify the maximum dimension number for which permuted tables will be produced. A permuted table for each dimension from 1 to the number specified is produced.

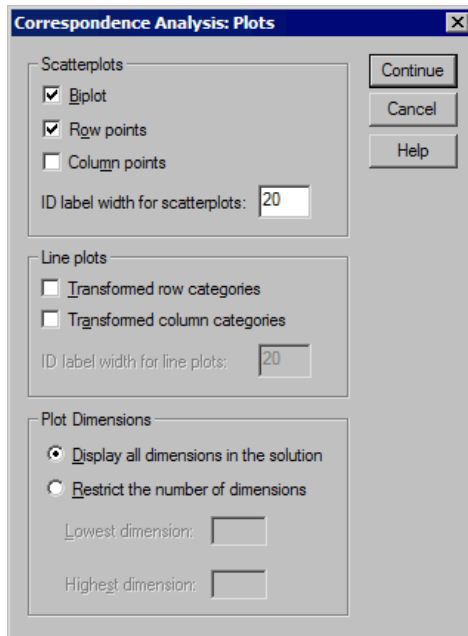
Confidence Statistics for Row points. Includes standard deviation and correlations for all nonsupplementary row points.

Confidence Statistics for Column points. Includes standard deviation and correlations for all nonsupplementary column points.

Correspondence Analysis Plots

The Plots dialog box allows you to specify which plots are produced.

Figure 5-6
Plots dialog box



Scatterplots. Produces a matrix of all pairwise plots of the dimensions. Available scatterplots include:

- **Biplot.** Produces a matrix of joint plots of the row and column points. If principal normalization is selected, the biplot is not available.
- **Row points.** Produces a matrix of plots of the row points.
- **Column points.** Produces a matrix of plots of the column points.

Optionally, you can specify how many value label characters to use when labeling the points. This value must be a non-negative integer less than or equal to 20.

Line plots. Produces a plot for every dimension of the selected variable. Available line plots include:

- **Transformed row categories.** Produces a plot of the original row category values against their corresponding row scores.
- **Transformed column categories.** Produces a plot of the original column category values against their corresponding column scores.

Optionally, you can specify how many value label characters to use when labeling the category axis. This value must be a non-negative integer less than or equal to 20.

Plot Dimensions. Allows you to control the dimensions displayed in the output.

- **Display all dimensions in the solution.** All dimensions in the solution are displayed in a scatterplot matrix.
- **Restrict the number of dimensions.** The displayed dimensions are restricted to plotted pairs. If you restrict the dimensions, you must select the lowest and highest dimensions to be plotted. The lowest dimension can range from 1 to the number of dimensions in the solution minus 1, and is plotted against higher dimensions. The highest dimension value can range from 2 to the number of dimensions in the solution, and indicates the highest dimension to be used in plotting the dimension pairs. This specification applies to all requested multidimensional plots.

CORRESPONDENCE Command Additional Features

You can customize your correspondence analysis if you paste your selections into a syntax window and edit the resulting `CORRESPONDENCE` command syntax. SPSS command language also allows you to:

- Specify table data as input instead of using casewise data (using the `TABLE = ALL` subcommand).
- Specify the number of value-label characters used to label points for each type of scatterplot matrix or biplot matrix (with the `PLOT` subcommand).
- Specify the number of value-label characters used to label points for each type of line plot (with the `PLOT` subcommand).
- Write a matrix of row and column scores to an SPSS matrix data file (with the `OUTFILE` subcommand).

- Write a matrix of confidence statistics (variances and covariances) for the singular values and the scores to an SPSS matrix data file (with the `OUTFILE` subcommand).
- Specify multiple sets of categories to be equal (with the `EQUAL` subcommand).

Multiple Correspondence Analysis

Multiple Correspondence Analysis quantifies nominal (categorical) data by assigning numerical values to the cases (objects) and categories so that objects within the same category are close together and objects in different categories are far apart. Each object is as close as possible to the category points of categories that apply to the object. In this way, the categories divide the objects into homogeneous subgroups. Variables are considered homogeneous when they classify objects in the same categories into the same subgroups.

Example. Multiple Correspondence Analysis could be used to graphically display the relationship between job category, minority classification, and gender. You might find that minority classification and gender discriminate between people but that job category does not. You might also find that the Latino and African-American categories are similar to each other.

Statistics and plots. Object scores, discrimination measures, iteration history, correlations of original and transformed variables, category quantifications, descriptive statistics, object points plots, biplots, category plots, joint category plots, transformation plots, and discrimination measures plots.

Data. String variable values are always converted into positive integers by ascending alphanumeric order. User-defined missing values, system-missing values, and values less than 1 are considered missing; you can recode or add a constant to variables with values less than 1 to make them nonmissing.

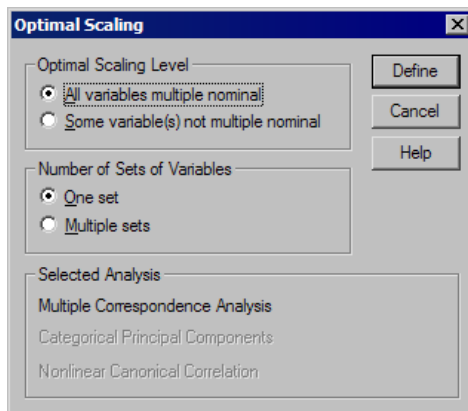
Assumptions. All variables have the multiple nominal scaling level. The data must contain at least three valid cases. The analysis is based on positive integer data. The discretization option will automatically categorize a fractional-valued variable by grouping its values into categories with a close to “normal” distribution and will automatically convert values of string variables into positive integers. You can specify other discretization schemes.

Related procedures. For two variables, Multiple Correspondence Analysis is analogous to Correspondence Analysis. If you believe that variables possess ordinal or numerical properties, Categorical Principal Components Analysis should be used. If sets of variables are of interest, Nonlinear Canonical Correlation Analysis should be used.

To Obtain a Multiple Correspondence Analysis

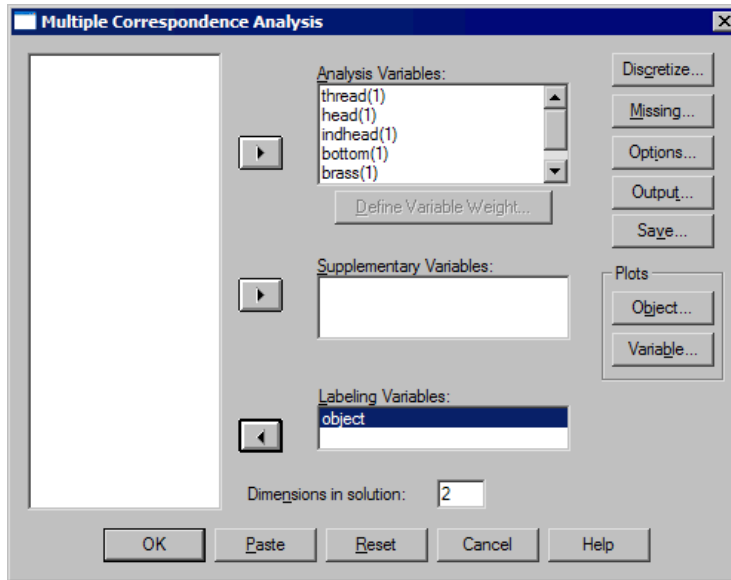
- ▶ From the menus choose:
Analyze
Data Reduction
Optimal Scaling...

Figure 6-1
Optimal Scaling dialog box



- ▶ Select All variables multiple nominal.
- ▶ Select One set.
- ▶ Click Define.

Figure 6-2
Multiple Correspondence Analysis dialog box



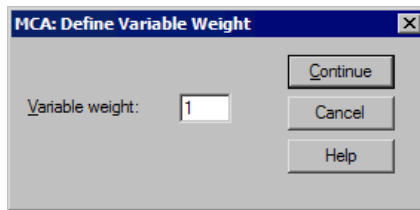
- ▶ Select at least two analysis variables and specify the number of dimensions in the solution.
- ▶ Click OK.

You may optionally specify supplementary variables, which are fitted into the solution found, or labeling variables for the plots.

Define Variable Weight in Multiple Correspondence Analysis

You can set the weight for analysis variables.

Figure 6-3
Define Variable Weight dialog box

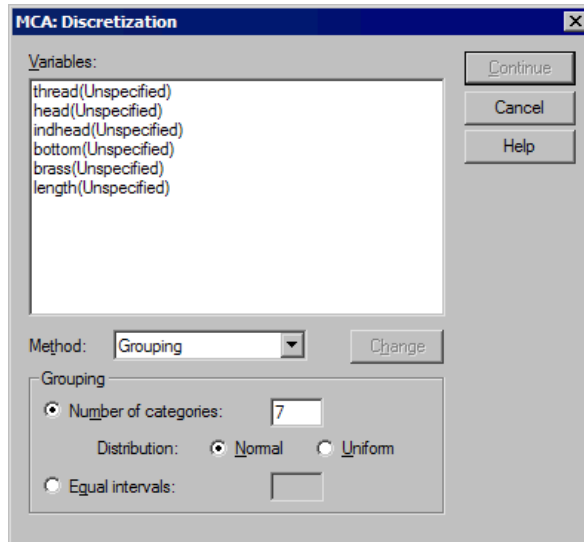


Variable weight. You can choose to define a weight for each variable. The value specified must be a positive integer. The default value is 1.

Multiple Correspondence Analysis Discretization

The Discretization dialog box allows you to select a method of recoding your variables. Fractional-valued variables are grouped into seven categories (or into the number of distinct values of the variable if this number is less than seven) with an approximately normal distribution unless otherwise specified. String variables are always converted into positive integers by assigning category indicators according to ascending alphanumeric order. Discretization for string variables applies to these integers. Other variables are left alone by default. The discretized variables are then used in the analysis.

Figure 6-4
Discretization dialog box



Method. Choose between grouping, ranking, and multiplying.

- **Grouping.** Recode into a specified number of categories or recode by interval.
- **Ranking.** The variable is discretized by ranking the cases.
- **Multiplying.** The current values of the variable are standardized, multiplied by 10, rounded, and have a constant added so that the lowest discretized value is 1.

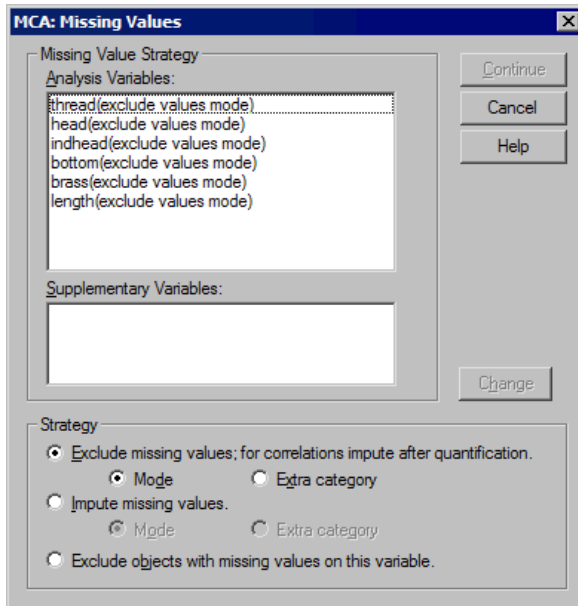
Grouping. The following options are available when discretizing variables by grouping:

- **Number of categories.** Specify a number of categories and whether the values of the variable should follow an approximately normal or uniform distribution across those categories.
- **Equal intervals.** Variables are recoded into categories defined by these equally sized intervals. You must specify the length of the intervals.

Multiple Correspondence Analysis Missing Values

The Missing Values dialog box allows you to choose the strategy for handling missing values in analysis variables and supplementary variables.

Figure 6-5
Missing Values dialog box



Missing Value Strategy. Choose to exclude missing values (passive treatment), impute missing values (active treatment), or exclude objects with missing values (listwise deletion).

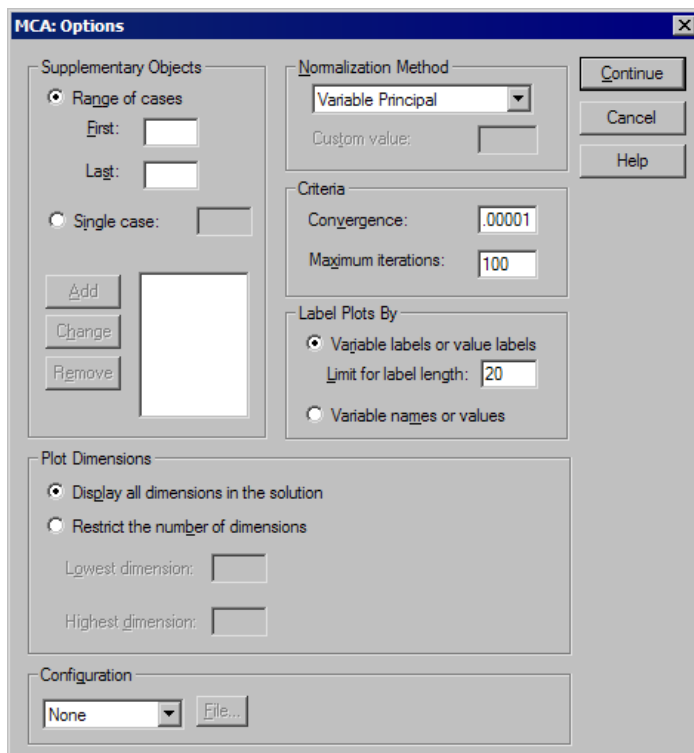
- **Exclude missing values; for correlations impute after quantification.** Objects with missing values on the selected variable do not contribute to the analysis for this variable. If all variables are given passive treatment, then objects with missing values on all variables are treated as supplementary. If correlations are specified in the Output dialog box, then (after analysis) missing values are imputed with the most frequent category, or mode, of the variable for the correlations of the original variables. For the correlations of the optimally scaled variables, you can choose the method of imputation. Select Mode to replace missing values with the mode of the optimally scaled variable. Select Extra category to replace missing values with the quantification of an extra category. This implies that objects with a missing value on this variable are considered to belong to the same (extra) category.

- **Impute missing values.** Objects with missing values on the selected variable have those values imputed. You can choose the method of imputation. Select Mode to replace missing values with the most frequent category. When there are multiple modes, the one with the smallest category indicator is used. Select Extra category to replace missing values with the same quantification of an extra category. This implies that objects with a missing value on this variable are considered to belong to the same (extra) category.
- **Exclude objects with missing values on this variable.** Objects with missing values on the selected variable are excluded from the analysis. This strategy is not available for supplementary variables.

Multiple Correspondence Analysis Options

The Options dialog box allows you to select the initial configuration, specify iteration and convergence criteria, select a normalization method, choose the method for labeling plots, and specify supplementary objects.

Figure 6-6
Options dialog box



Supplementary Objects. Specify the case number of the object, or the first and last case numbers of a range of objects, that you want to make supplementary, and then click Add. Continue until you have specified all of your supplementary objects. If an object is specified as supplementary, then case weights are ignored for that object.

Normalization Method. You can specify one of five options for normalizing the object scores and the variables. Only one normalization method can be used in a given analysis.

- **Variable Principal.** This option optimizes the association between variables. The coordinates of the variables in the object space are the component loadings (correlations with principal components, such as dimensions and object scores). This is useful when you are interested primarily in the correlation between the variables.

- **Object Principal.** This option optimizes distances between objects. This is useful when you are interested primarily in differences or similarities between the objects.
- **Symmetrical.** Use this normalization option if you are interested primarily in the relation between objects and variables.
- **Independent.** Use this normalization option if you want to examine distances between objects and correlations between variables separately.
- **Custom.** You can specify any real value in the closed interval $[-1, 1]$. A value of 1 is equal to the Object Principal method, a value of 0 is equal to the Symmetrical method, and a value of -1 is equal to the Variable Principal method. By specifying a value greater than -1 and less than 1, you can spread the eigenvalue over both objects and variables. This method is useful for making a tailor-made biplot or triplot.

Criteria. You can specify the maximum number of iterations the procedure can go through in its computations. You can also select a convergence criterion value. The algorithm stops iterating if the difference in total fit between the last two iterations is less than the convergence value or if the maximum number of iterations is reached.

Label Plots By. Allows you to specify whether variables and value labels or variable names and values will be used in the plots. You can also specify a maximum length for labels.

Plot Dimensions. Allows you to control the dimensions displayed in the output.

- **Display all dimensions in the solution.** All dimensions in the solution are displayed in a scatterplot matrix.
- **Restrict the number of dimensions.** The displayed dimensions are restricted to plotted pairs. If you restrict the dimensions, you must select the lowest and highest dimensions to be plotted. The lowest dimension can range from 1 to the number of dimensions in the solution minus 1 and is plotted against higher dimensions. The highest dimension value can range from 2 to the number of dimensions in the solution and indicates the highest dimension to be used in plotting the dimension pairs. This specification applies to all requested multidimensional plots.

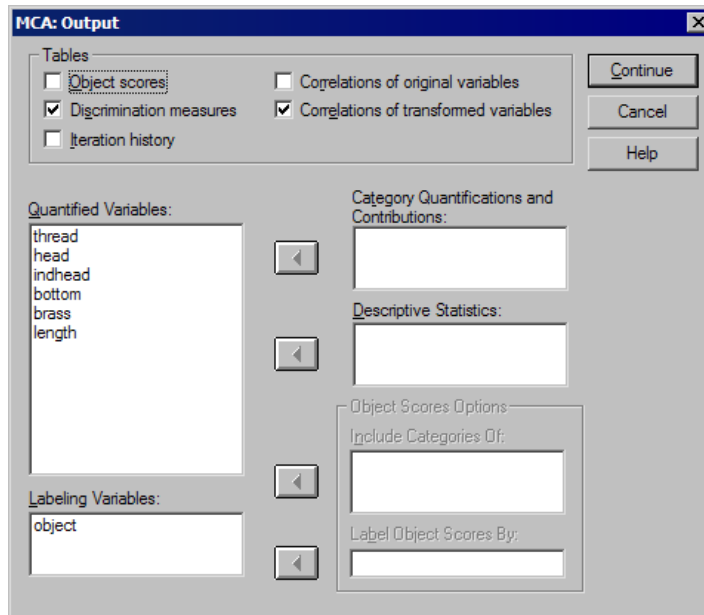
Configuration. You can read data from a file containing the coordinates of a configuration. The first variable in the file should contain the coordinates for the first dimension, the second variable should contain the coordinates for the second dimension, and so on.

- **Initial.** The configuration in the file specified will be used as the starting point of the analysis.
- **Fixed.** The configuration in the file specified will be used to fit in the variables. The variables that are fitted in must be selected as analysis variables, but, because the configuration is fixed, they are treated as supplementary variables (so they do not need to be selected as supplementary variables).

Multiple Correspondence Analysis Output

The Output dialog box allows you to produce tables for object scores, discrimination measures, iteration history, correlations of original and transformed variables, category quantifications for selected variables, and descriptive statistics for selected variables.

Figure 6-7
Output dialog box



Object scores. Displays the object scores, including mass, inertia, and contributions, and has the following options:

- **Include Categories Of.** Displays the category indicators of the analysis variables selected.
- **Label Object Scores By.** From the list of variables specified as labeling variables, you can select one to label the objects.

Discrimination measures. Displays the discrimination measures per variable and per dimension.

Iteration history. For each iteration, the variance accounted for, loss, and increase in variance accounted for are shown.

Correlations of original variables. Shows the correlation matrix of the original variables and the eigenvalues of that matrix.

Correlations of transformed variables. Shows the correlation matrix of the transformed (optimally scaled) variables and the eigenvalues of that matrix.

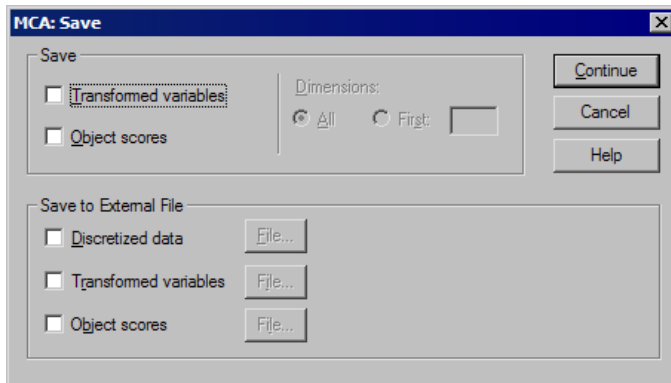
Category Quantifications and Contributions. Gives the category quantifications (coordinates), including mass, inertia, and contributions, for each dimension of the variable(s) selected.

Descriptive Statistics. Displays frequencies, number of missing values, and mode of the variable(s) selected.

Multiple Correspondence Analysis Save

The Save dialog box allows you to add the transformed variables and object scores to the working data file or as new variables in external files and to save the discretized data as new variables in an external data file.

Figure 6-8
Save dialog box



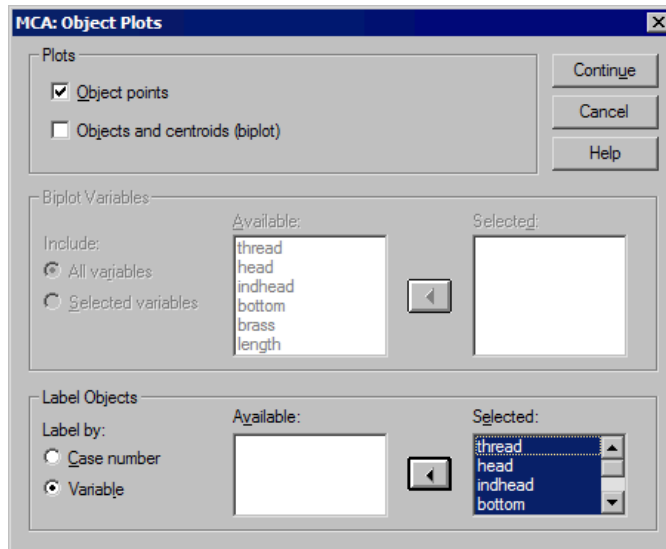
Save. Save selections to the working data file. The number of dimensions to be saved must be specified.

Save to External File. Save selections to a new external file. Specify a filename for each selected option by clicking File. Each file specified must have a different name.

Multiple Correspondence Analysis Object Plots

The Object Plots dialog box allows you to specify the types of plots desired and the variables for which plots will be produced.

Figure 6-9
Object Plots dialog box



Object points. A plot of the object points is displayed.

Objects and centroids (biplot). The object points are plotted with the variable centroids.

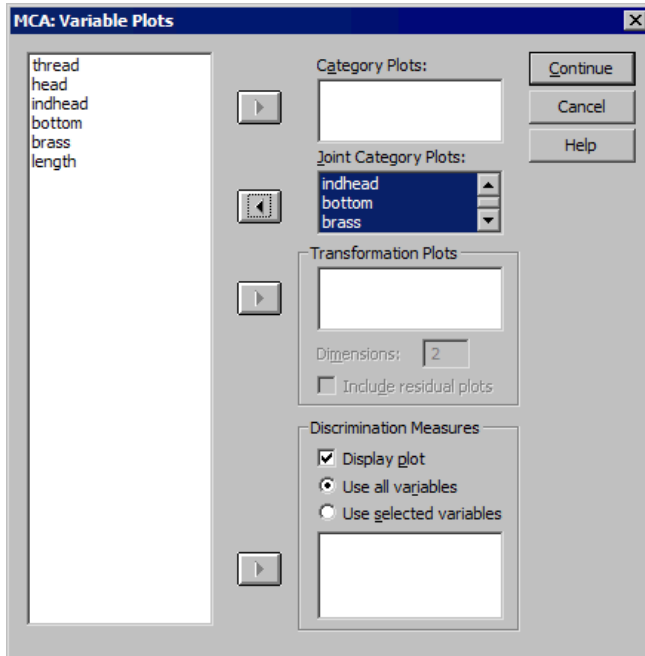
Biplot Variables. You can choose to use all variables for the biplots or select a subset.

Label Objects. You can choose to have objects labeled with the categories of selected variables (you may choose category indicator values or value labels in the Options dialog box) or with their case numbers. One plot is produced per variable if Variable is selected.

Multiple Correspondence Analysis Variable Plots

The Variable Plots dialog box allows you to specify the types of plots desired and the variables for which plots will be produced.

Figure 6-10
Variable Plots dialog box



Category Plots. For each variable selected, a plot of the centroid coordinates is plotted. Categories are in the centroids of the objects in the particular categories.

Joint Category Plots. This is a single plot of the centroid coordinates of each selected variable.

Transformation Plots. Displays a plot of the optimal category quantifications versus the category indicators. You can specify the number of dimensions desired; one plot will be generated for each dimension. You can also choose to display residual plots for each variable selected.

Discrimination Measures. Produces a single plot of the discrimination measures for the selected variables.

MULTIPLE CORRESPONDENCE Command Additional Features

You can customize your Multiple Correspondence Analysis if you paste your selections into a syntax window and edit the resulting `MULTIPLE CORRESPONDENCE` command syntax. The SPSS command language also allows you to:

- Specify rootnames for the transformed variables, object scores, and approximations when saving them to the working data file (with the `SAVE` subcommand).
- Specify a maximum length for labels for each plot separately (with the `PLOT` subcommand).
- Specify a separate variable list for residual plots (with the `PLOT` subcommand).

Multidimensional Scaling (PROXSCAL)

Multidimensional scaling attempts to find the structure in a set of proximity measures between objects. This is accomplished by assigning observations to specific locations in a conceptual low-dimensional space such that the distances between points in the space match the given (dis)similarities as closely as possible. The result is a least-squares representation of the objects in that low-dimensional space, which, in many cases, will help you to further understand your data.

Example. Multidimensional scaling can be very useful in determining perceptual relationships. For example, when considering your product image, you can conduct a survey in order to obtain a data set that describes the perceived similarity (or proximity) of your product to those of your competitors. Using these proximities and independent variables (such as price), you can try to determine which variables are important to how people view these products, and adjust your image accordingly.

Statistics and plots. Iteration history, stress measures, stress decomposition, coordinates of the common space, object distances within the final configuration, individual space weights, individual spaces, transformed proximities, transformed independent variables, stress plots, common space scatterplots, individual space weight scatterplots, individual spaces scatterplots, transformation plots, Shepard residual plots, and independent variables transformation plots.

Data. Data can be supplied in the form of proximity matrices or variables that are converted into proximity matrices. The matrices may be formatted in columns or across columns. The proximities may be treated on the ratio, interval, ordinal, or spline scaling levels.

Assumptions. At least three variables must be specified. The number of dimensions may not exceed the number of objects minus one. Dimensionality reduction is omitted if combined with multiple random starts. If only one source is specified, all models are equivalent to the identity model; therefore, the analysis defaults to the identity model.

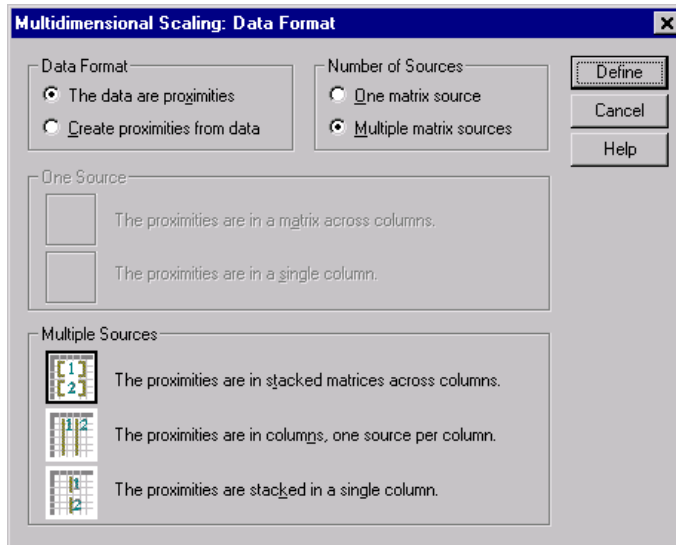
Related procedures. Scaling all variables at the numerical level corresponds to standard multidimensional scaling analysis.

To Obtain a Multidimensional Scaling

- ▶ From the menus choose:
 - Analyze
 - Scale
 - Multidimensional Scaling (PROXSCAL)...

This opens the Data Format dialog box.

Figure 7-1
Data Format dialog box



- ▶ Specify the format of your data:

Data Format. Specify whether your data consist of proximity measures or you want to create proximities from the data.

Number of Sources. If your data are proximities, specify whether you have a single source or multiple sources of proximity measures.

One Source. If there is one source of proximities, specify whether your data set is formatted with the proximities in a matrix across the columns or in a single column with two separate variables to identify the row and column of each proximity.

- **Proximities in a matrix across columns.** The proximity matrix is spread across a number of columns equal to the number of objects. This leads to the Proximities in Matrices across Columns dialog box.
- **Proximities in a single column.** The proximity matrix is collapsed into a single column, or variable. Two additional variables, identifying the row and column for each cell, are necessary. This leads to the Proximities in One Column dialog box.

Multiple Sources. If there are multiple sources of proximities, specify whether the data set is formatted with the proximities in stacked matrices across columns, in multiple columns with one source per column, or in a single column.

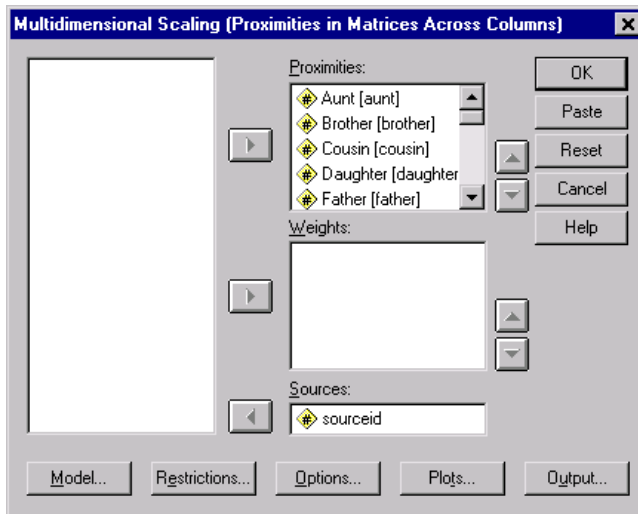
- **Proximities in stacked matrices across columns.** The proximity matrices are spread across a number of columns equal to the number of objects and are stacked above one another across a number of rows equal to the number of objects times the number of sources. This leads to the Proximities in Matrices across Columns dialog box.
- **Proximities in columns, one source per column.** The proximity matrices are collapsed into multiple columns, or variables. Two additional variables, identifying the row and column for each cell, are necessary. This leads to the Proximities in Columns dialog box.
- **Proximities in single column.** The proximity matrices are collapsed into a single column, or variable. Three additional variables, identifying the row, column, and source for each cell, are necessary. This leads to the Proximities in One Column dialog box.

- ▶ Click Define.

Proximities in Matrices across Columns

If you select the proximities in matrices data model for either one source or multiple sources in the Data Format dialog box, the main dialog box will appear as follows:

Figure 7-2
Proximities in Matrices across Columns dialog box



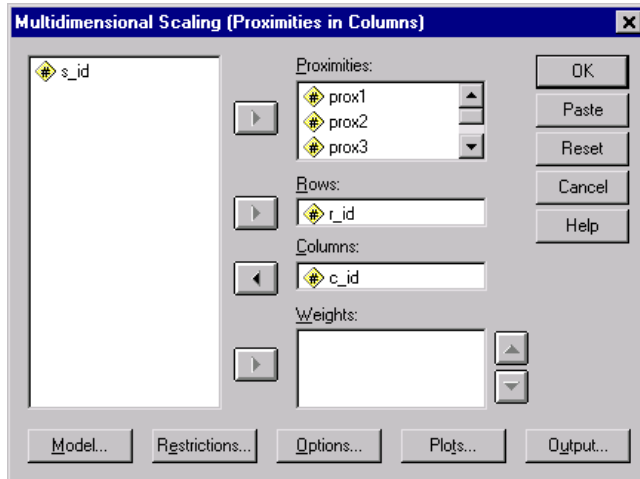
- ▶ Select three or more proximities variables. Please be sure that the order of the variables in the list matches the order of the columns of the proximities.
- ▶ Optionally, select a number of weights variables equal to the number of proximities variables. Again, be sure that the order of the weights matches the order of the proximities they weight.
- ▶ If there are multiple sources, optionally, select a sources variable. The number of cases in each proximities variable should equal the number of proximities variables times the number of sources.

Additionally, you can define a model for the multidimensional scaling, place restrictions on the common space, set convergence criteria, specify the initial configuration to be used, and choose plots and output.

Proximities in Columns

If you select the multiple columns model for multiple sources in the Data Format dialog box, the main dialog box will appear as follows:

Figure 7-3
Proximities in Columns dialog box



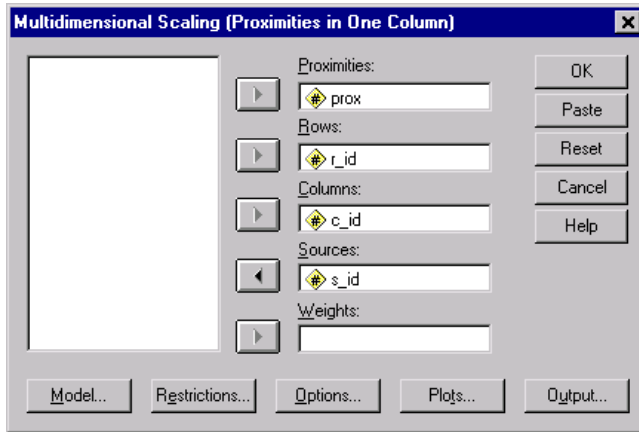
- ▶ Select two or more proximities variables. Each variable is assumed to be a matrix of proximities from a separate source.
- ▶ Select a rows variable. This defines the row locations for the proximities in each proximities variable.
- ▶ Select a columns variable. This defines the column locations for the proximities in each proximities variable. Cells of the proximity matrix that are not given a row/column designation are treated as missing.
- ▶ Optionally, select a number of weights variables equal to the number of proximities variables.

Additionally, you can define a model for the multidimensional scaling, place restrictions on the common space, set convergence criteria, specify the initial configuration to be used, and choose plots and output.

Proximities in One Column

If you select the one column model for either one source or multiple sources in the Data Format dialog box, the main dialog box will appear as follows:

Figure 7-4
Proximities in One Column dialog box



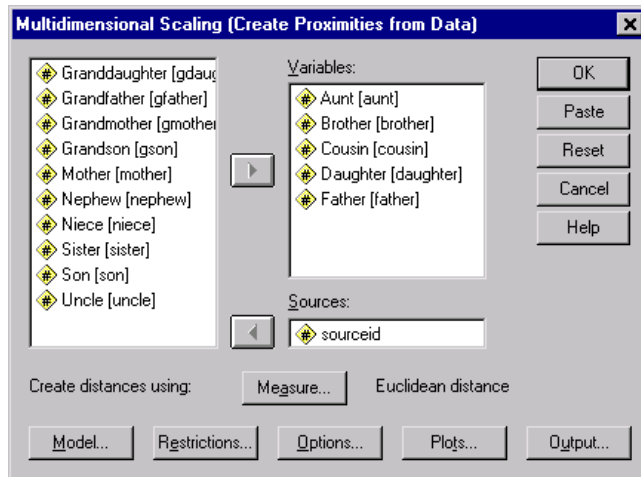
- ▶ Select a proximities variable. It is assumed to be one or more matrices of proximities.
- ▶ Select a rows variable. This defines the row locations for the proximities in the proximities variable.
- ▶ Select a columns variable. This defines the column locations for the proximities in the proximities variable.
- ▶ If there are multiple sources, select a sources variable. For each source, cells of the proximity matrix that are not given a row/column designation are treated as missing.
- ▶ Optionally, select a weights variable.

Additionally, you can define a model for the multidimensional scaling, place restrictions on the common space, set convergence criteria, specify the initial configuration to be used, and choose plots and output.

Create Proximities from Data

If you choose to create proximities from the data in the Data Format dialog box, the main dialog box will appear as follows:

Figure 7-5
Create Proximities from Data dialog box

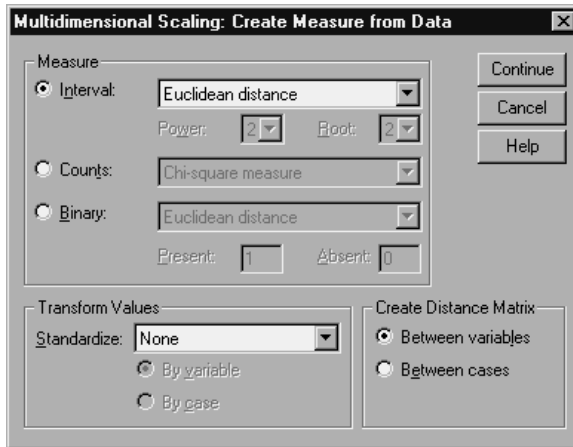


- ▶ If you create distances between variables (see the Measures dialog box), select at least three variables. These will be used to create the proximity matrix (or matrices, if there are multiple sources). If you create distances between cases, only one variable is needed.
- ▶ If there are multiple sources, select a sources variable.
- ▶ Optionally, choose a measure for creating proximities.

Additionally, you can define a model for the multidimensional scaling, place restrictions on the common space, set convergence criteria, specify the initial configuration to be used, and choose plots and output.

Measures Dialog Box

Figure 7-6
Create Measure from Data dialog box



Multidimensional scaling uses dissimilarity data to create a scaling solution. If your data are multivariate data (values of measured variables), you must create dissimilarity data in order to compute a multidimensional scaling solution. You can specify the details of creating dissimilarity measures from your data.

Measure. Allows you to specify the dissimilarity measure for your analysis. Select one alternative from the Measure group corresponding to your type of data, and then select one of the measures from the drop-down list corresponding to that type of measure. Available alternatives are:

- **Interval.** Euclidean distance, Squared Euclidean distance, Chebychev, Block, Minkowski, or Customized.
- **Counts.** Chi-square measure or Phi-square measure.
- **Binary.** Euclidean distance, Squared Euclidean distance, Size difference, Pattern difference, Variance, or Lance and Williams.

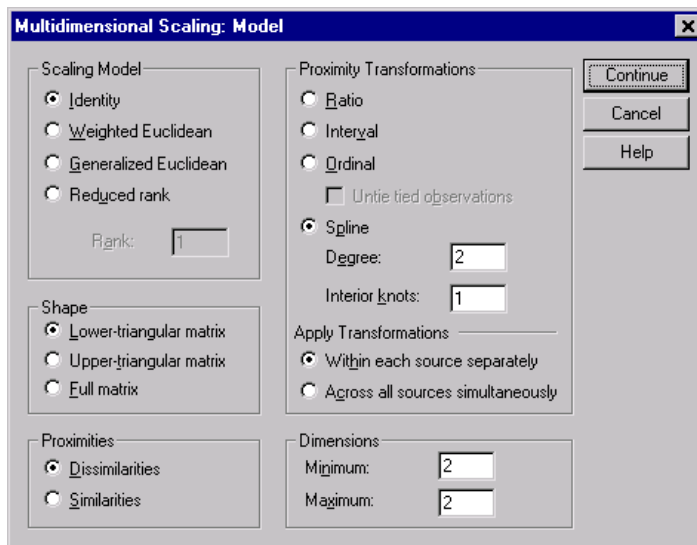
Create Distance Matrix. Allows you to choose the unit of analysis. Alternatives are Between variables or Between cases.

Transform Values. In certain cases, such as when variables are measured on very different scales, you may want to standardize values before computing proximities (not applicable to binary data). Select a standardization method from the Standardize drop-down list (if no standardization is required, select None).

Define a Multidimensional Scaling Model

The Model dialog box allows you to specify a scaling model, its minimum and maximum number of dimensions, the structure of the proximity matrix, the transformation to use on the proximities, and whether proximities are transformed within each source separately, or unconditionally on the source.

Figure 7-7
Model dialog box



Scaling Model. Choose from the following alternatives:

- **Identity.** All sources have the same configuration.
- **Weighted Euclidean.** This model is an individual differences model. Each source has an individual space in which every dimension of the common space is weighted differentially.

- **Generalized Euclidean.** This model is an individual differences model. Each source has an individual space that is equal to a rotation of the common space, followed by a differential weighting of the dimensions.
- **Reduced rank.** This is a Generalized Euclidean model for which you can specify the rank of the individual space. You must specify a rank that is greater than or equal to 1 and less than the maximum number of dimensions.

Shape. Specify whether the proximities should be taken from the lower-triangular part or the upper-triangular part of the proximity matrix. You may specify that the full matrix be used, in which case the weighted sum of the upper-triangular part and the lower-triangular part will be analyzed. In any case, the complete matrix should be specified, including the diagonal, though only the specified parts will be used.

Proximities. Specify whether your proximity matrix contains measures of similarity or dissimilarity.

Proximity Transformations. Choose from the following alternatives:

- **Ratio.** The transformed proximities are proportional to the original proximities. This is allowed only for positively valued proximities.
- **Interval.** The transformed proximities are proportional to the original proximities, plus an intercept term. The intercept assures all transformed proximities to be positive.
- **Ordinal.** The transformed proximities have the same order as the original proximities. You may specify whether tied proximities should be kept tied or allowed to become untied.
- **Spline.** The transformed proximities are a smooth nondecreasing piecewise polynomial transformation of the original proximities. You may specify the degree of the polynomial and the number of interior knots.

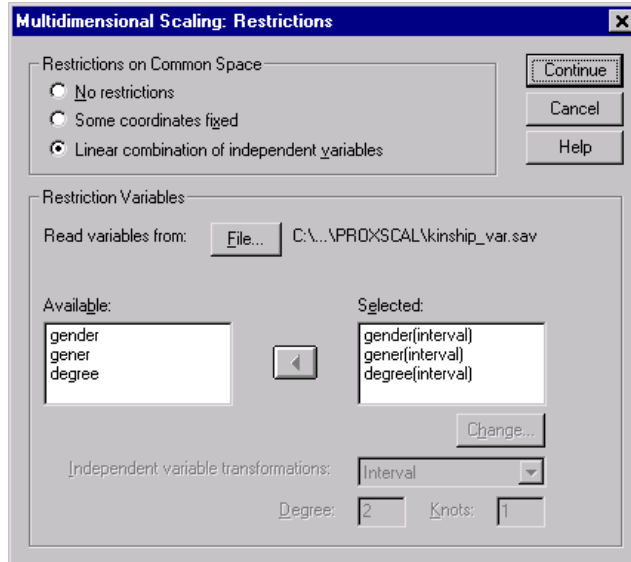
Apply Transformations. Specify whether only proximities within each source are compared with each other or whether the comparisons are unconditional on the source.

Dimensions. By default, a solution is computed in two dimensions (Minimum = 2, Maximum = 2). You may choose an integer minimum and maximum from 1 to the number of objects minus 1 as long as the minimum is less than or equal to the maximum. The procedure computes a solution in the maximum dimensions and then reduces the dimensionality in steps until the lowest is reached.

Multidimensional Scaling Restrictions

The Restrictions dialog box allows you to place restrictions on the common space.

Figure 7-8
Restrictions dialog box



Restrictions on Common Space. Specify the type of restriction desired.

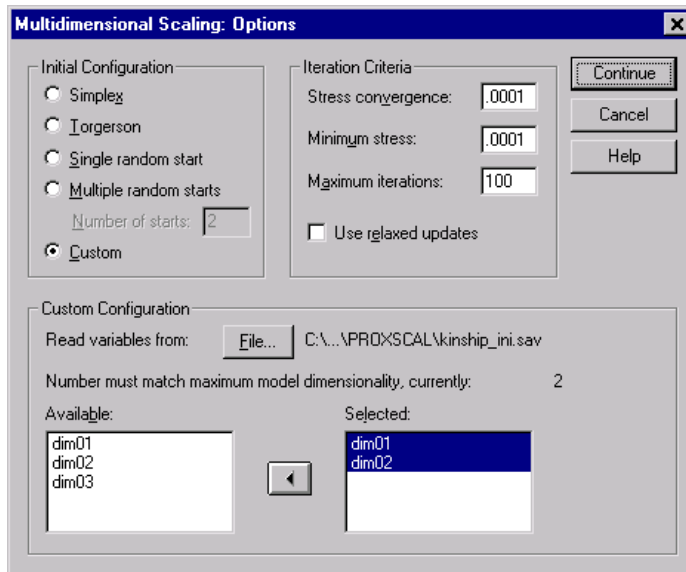
- **No restrictions.** No restrictions are placed on the common space.
- **Some coordinates fixed.** The first variable selected contains the coordinates of the objects on the first dimension; the second variable corresponds to coordinates on the second dimension; and so on. A missing value indicates that a coordinate on a dimension is free. The number of variables selected must equal the maximum number of dimensions requested.
- **Linear combination of independent variables.** The common space is restricted to be a linear combination of the variables selected.

Restriction Variables. Select the variables that define the restrictions on the common space. If you specified a linear combination, you may specify an interval, nominal, ordinal, or spline transformation for the restriction variables. In either case, the number of cases for each variable must equal the number of objects.

Multidimensional Scaling Options

The Options dialog box allows you to select the initial configuration style, specify iteration and convergence criteria, and select standard or relaxed updates.

Figure 7-9
Options dialog box



Initial Configuration. Choose one of the following alternatives:

- **Simplex.** Objects are placed at the same distance from each other in the maximum dimension. One iteration is taken to improve this high-dimensional configuration, followed by a dimension reduction operation to obtain an initial configuration that has the maximum number of dimensions that you specified in the Model dialog box.
- **Torgerson.** A classical scaling solution is used as the initial configuration.
- **Single random start.** A configuration is chosen at random.

- **Multiple random starts.** Several configurations are chosen at random, and the one with the lowest normalized raw stress is used as the initial configuration.
- **Custom.** You may select variables that contain the coordinates of your own initial configuration. The number of variables selected should equal the maximum number of dimensions specified, with the first variable corresponding to coordinates on dimension 1, the second variable corresponding to coordinates on dimension 2, and so on. The number of cases in each variable should equal the number of objects.

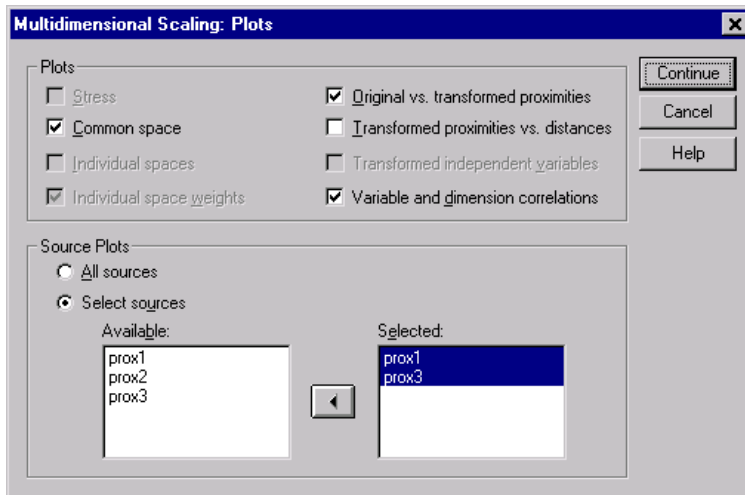
Iteration Criteria. Specify the iteration criteria values.

- **Stress convergence.** The algorithm will stop iterating when the difference in consecutive normalized raw stress values is less than the number specified here, which must lie between 0.0 and 1.0.
- **Minimum stress.** The algorithm will stop when the normalized raw stress falls below the number specified here, which must lie between 0.0 and 1.0.
- **Maximum iterations.** The algorithm will perform the number of iterations specified here, unless one of the above criteria is satisfied first.
- **Use relaxed updates.** Relaxed updates will speed up the algorithm; these cannot be used with models other than the identity model or with restrictions.

Multidimensional Scaling Plots, Version 1

The Plots dialog box allows you to specify which plots will be produced. If you have the Proximities in Columns data format, the following Plots dialog box is displayed. For Individual space weights, Original vs. transformed proximities, and Transformed proximities vs. distances plots, you may specify the sources for which the plots should be produced. The list of available sources is the list of proximities variables in the main dialog box.

Figure 7-10
Plots dialog box, version 1



Stress. A plot is produced of normalized raw stress versus dimensions. This plot is produced only if the maximum number of dimensions is larger than the minimum number of dimensions.

Common space. A scatterplot matrix of coordinates of the common space is displayed.

Individual spaces. For each source, the coordinates of the individual spaces are displayed in scatterplot matrices. This is possible only if one of the individual differences models is specified in the Model dialog box.

Individual space weights. A scatterplot is produced of the individual space weights. This is possible only if one of the individual differences models is specified in the Model dialog box. For the weighted Euclidean model, the weights are printed in plots, with one dimension on each axis. For the generalized Euclidean model, one plot is produced per dimension, indicating both rotation and weighting of that dimension. The reduced rank model produces the same plot as the generalized Euclidean model but reduces the number of dimensions for the individual spaces.

Original vs. transformed proximities. Plots are produced of the original proximities versus the transformed proximities.

Transformed proximities vs. distances. The transformed proximities versus the distances are plotted.

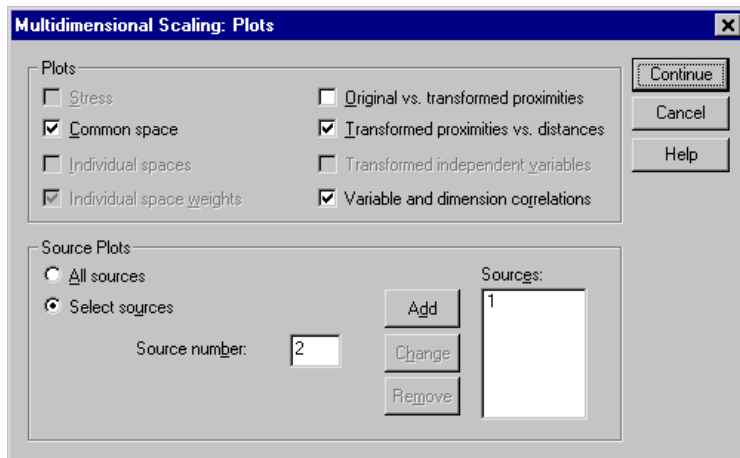
Transformed independent variables. Transformation plots are produced for the independent variables.

Variable and dimension correlations. A plot of correlations between the independent variables and the dimensions of the common space is displayed.

Multidimensional Scaling Plots, Version 2

The Plots dialog box allows you to specify which plots will be produced. If your data format is anything other than Proximities in Columns, the following Plots dialog box is displayed. For Individual space weights, Original vs. transformed proximities, and Transformed proximities vs. distances plots, you may specify the sources for which the plots should be produced. The source numbers entered must be values of the sources variable specified in the main dialog box and range from 1 to the number of sources.

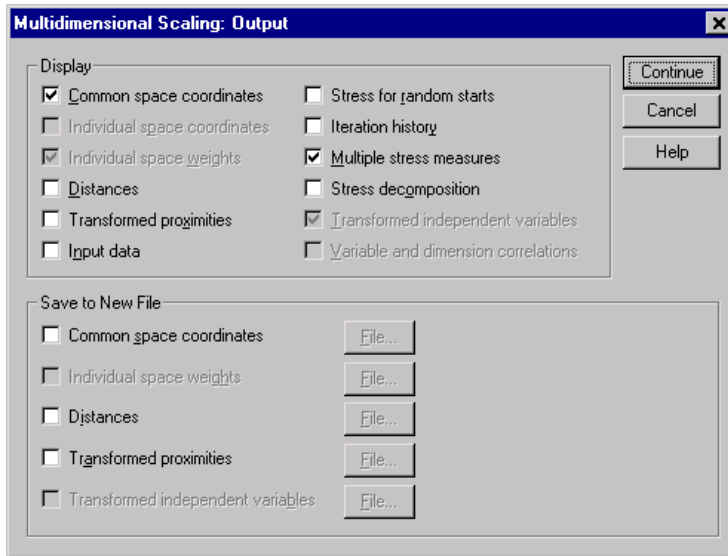
Figure 7-11
Plots dialog box, version 2



Multidimensional Scaling Output

The Output dialog box allows you to control the amount of displayed output and save some of it to separate files.

Figure 7-12
Output dialog box



Display. Select one or more of the following for display:

- **Common space coordinates.** Displays the coordinates of the common space.
- **Individual space coordinates.** The coordinates of the individual spaces are displayed only if the model is not the identity model.
- **Individual space weights.** Displays the individual space weights only if one of the individual differences models is specified. Depending on the model, the space weights are decomposed in rotation weights and dimension weights, which are also displayed.
- **Distances.** Displays the distances between the objects in the configuration.
- **Transformed proximities.** Displays the transformed proximities between the objects in the configuration.
- **Input data.** Includes the original proximities and, if present, the data weights, the initial configuration, and the fixed coordinates or the independent variables.
- **Stress for random starts.** Displays the random number seed and normalized raw stress value of each random start.
- **Iteration history.** Displays the history of iterations of the main algorithm.

- **Multiple stress measures.** Displays different stress values. The table contains values for normalized raw stress, Stress-I, Stress-II, S-Stress, Dispersion Accounted For (DAF), and Tucker's Coefficient of Congruence.
- **Stress decomposition.** Displays an objects and sources decomposition of final normalized raw stress, including the average per object and the average per source.
- **Transformed independent variables.** If a linear combination restriction was selected, the transformed independent variables and the corresponding regression weights are displayed.
- **Variable and dimension correlations.** If a linear combination restriction was selected, the correlations between the independent variables and the dimensions of the common space are displayed.

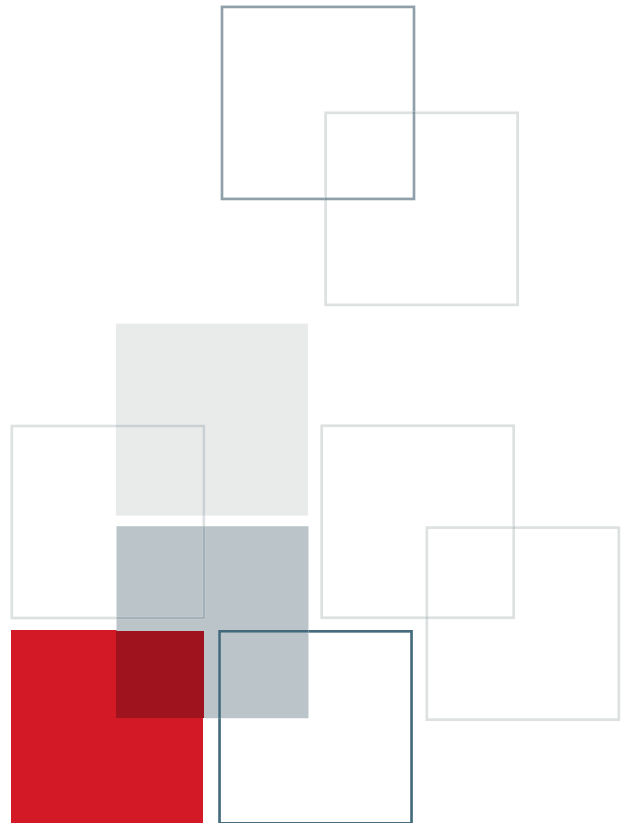
Save to New File. You can save the common space coordinates, individual space weights, distances, transformed proximities, and transformed independent variables to separate SPSS data files.

PROXSCAL Command Additional Features

You can customize your multidimensional scaling of proximities analysis if you paste your selections into a syntax window and edit the resulting PROXSCAL command syntax. SPSS command language also allows you to:

- Specify separate variable lists for transformations and residuals plots (with the PLOT subcommand).
- Specify separate source lists for individual space weights, transformations, and residuals plots (with the PLOT subcommand).
- Specify a subset of the independent variables transformation plots to be displayed (with the PLOT subcommand).

Part 2: Examples



Categorical Regression

The goal of categorical regression with optimal scaling is to describe the relationship between a response variable and a set of predictors. By quantifying this relationship, values of the response can be predicted for any combination of predictors.

In this chapter, two examples serve to illustrate the analyses involved in optimal scaling regression. The first example uses a small data set to illustrate the basic concepts. The second example uses a much larger set of variables and observations in a practical example.

Example: Carpet Cleaner Data

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. The following table displays the variables used in the carpet-cleaner study, with their variable labels and values.

Table 8-1
Explanatory variables in the carpet-cleaner study

Variable name	Variable label	Value label
<i>package</i>	Package design	A*, B*, C*
<i>brand</i>	Brand name	K2R, Glory, Bissell
<i>price</i>	Price	\$1.19, \$1.39, \$1.59

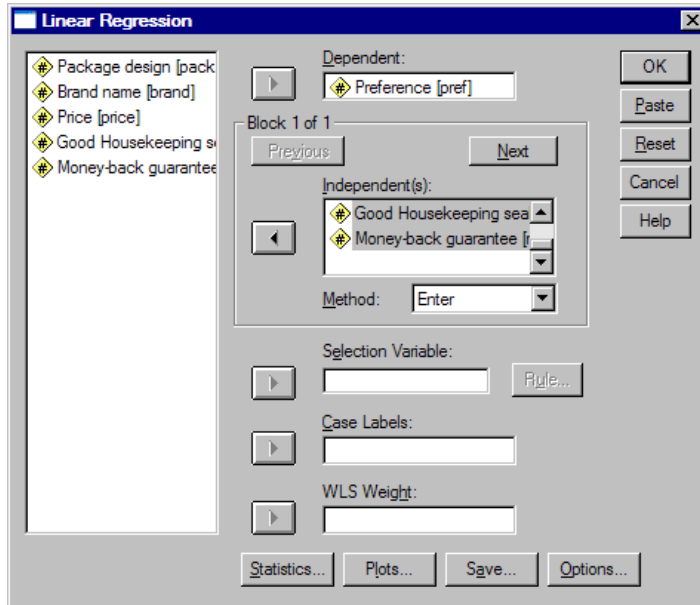
Variable name	Variable label	Value label
<i>seal</i>	<i>Good Housekeeping seal</i>	No, yes
<i>money</i>	Money-back guarantee	No, yes

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. Using categorical regression, you will explore how the five factors are related to preference. This data set can be found in *carpet.sav*, found in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

A Standard Linear Regression Analysis

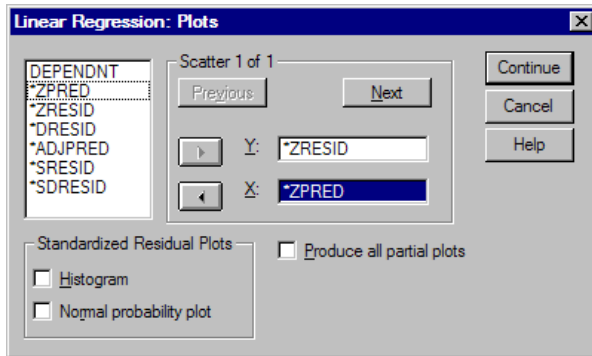
- ▶ To produce standard linear regression output, from the menus choose:
 - Analyze
 - Regression
 - Linear...

Figure 8-1
Linear Regression dialog box



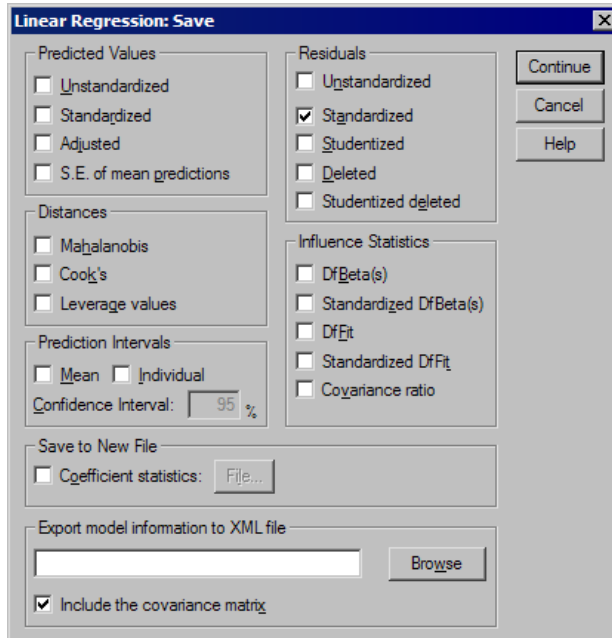
- ▶ Select *Preference* as the dependent variable.
- ▶ Select *Package design* through *Money-back guarantee* as independent variables.
- ▶ Click *Plots*.

Figure 8-2
Plots dialog box



- ▶ Select **ZRESID* as the y-axis variable.
- ▶ Select **ZPRED* as the x-axis variable.
- ▶ Click Continue.
- ▶ Click Save in the Linear Regression dialog box.

Figure 8-3
Save dialog box



- ▶ Select Standardized in the Residuals group.
- ▶ Click Continue.
- ▶ Click OK in the Linear Regression dialog box.

Model Summary

Figure 8-4
Model summary for standard linear regression

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.841 ^a	.707	.615	3.99810

a. Predictors: (Constant), Money-back guarantee, Price, Good Housekeeping seal, Brand name, Package design

The standard approach for describing the relationships in this problem is linear regression. The most common measure of how well a regression model fits the data is R^2 . This statistic represents how much of the variance in the response is explained by the weighted combination of predictors. The closer R^2 is to 1, the better the model fits. Regressing *Preference* on the five predictors results in an R^2 of 0.707, indicating that approximately 71% of the variance in the preference rankings is explained by the predictor variables in the linear regression.

Coefficients

The standardized coefficients are shown in the table. The sign of the coefficient indicates whether the predicted response increases or decreases when the predictor increases, all other predictors being constant. For categorical data, the category coding determines the meaning of an increase in a predictor. For instance, an increase in *Money-back guarantee*, *Package design*, or *Good Housekeeping seal* will result in a decrease in predicted preference ranking. *Money-back guarantee* is coded 1 for *no money-back guarantee* and 2 for *money-back guarantee*. An increase in *Money-back guarantee* corresponds to the addition of a money-back guarantee. Thus, adding a money-back guarantee reduces the predicted preference ranking, which corresponds to an increased predicted preference.

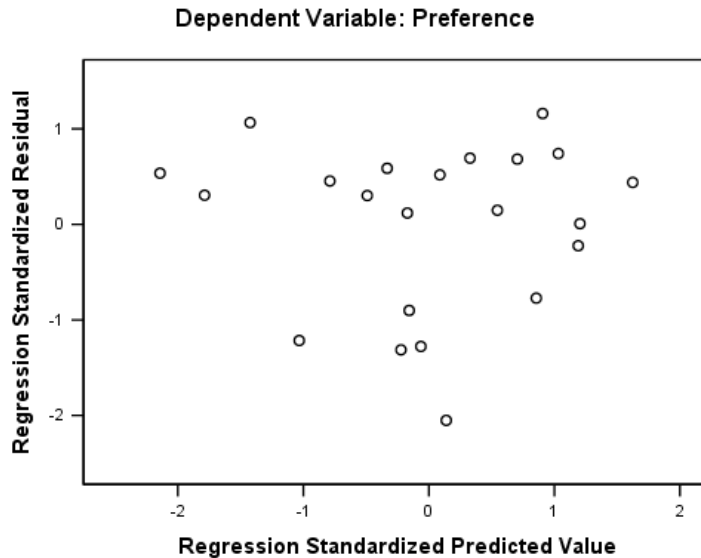
Figure 8-5
Regression coefficients

Model		Unstandardized Coefficients		Standardized Coefficients	t	Sig.
		B	Std. Error	Beta		
1	(Constant)	22.529	5.177		4.352	.000
	Package design	-4.159	1.036	-.560	-4.015	.001
	Brand name	.429	1.054	.056	.407	.689
	Price	2.703	1.009	.366	2.681	.016
	Good Housekeeping seal	-4.314	1.780	-.330	-2.423	.028
	Money-back guarantee	-2.779	1.921	-.197	-1.447	.167

The value of the coefficient reflects the amount of change in the predicted preference ranking. Using standardized coefficients, interpretations are based on the standard deviations of the variables. Each coefficient indicates the number of standard deviations that the predicted response changes for a one standard deviation change in a predictor, all other predictors remaining constant. For example, a one standard deviation change in *Brand name* yields an increase in predicted preference of 0.056 standard deviations. The standard deviation of *Preference* is 6.44, so *Preference* increases by $0.056 \times 6.44 = 0.361$. Changes in *Package design* yield the greatest changes in predicted preference.

Residual Scatterplots

Figure 8-6
Residuals versus predicted values

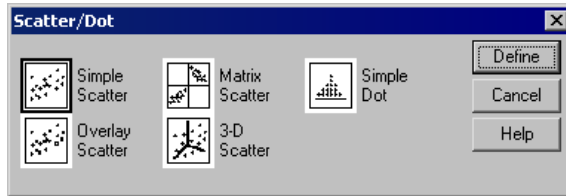


The standardized residuals are plotted against the standardized predicted values. No patterns should be present if the model fits well. Here you see a U-shape in which both low and high standardized predicted values have positive residuals. Standardized predicted values near 0 tend to have negative residuals.

- ▶ To produce a scatterplot of the residuals by the predictor *Package design*, from the menus choose:

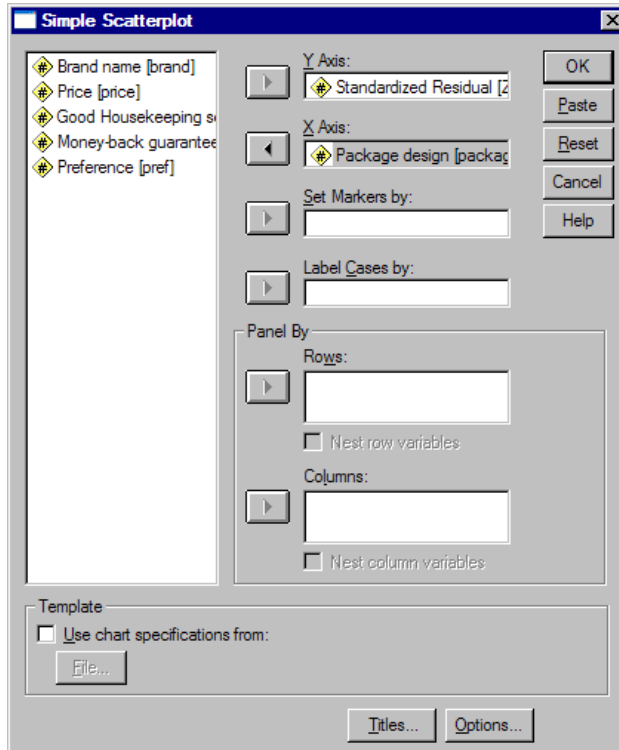
Graphs
Scatter/Dot...

Figure 8-7
Scatter/Dot dialog box



- Click Define.

Figure 8-8
Simple Scatterplot dialog box

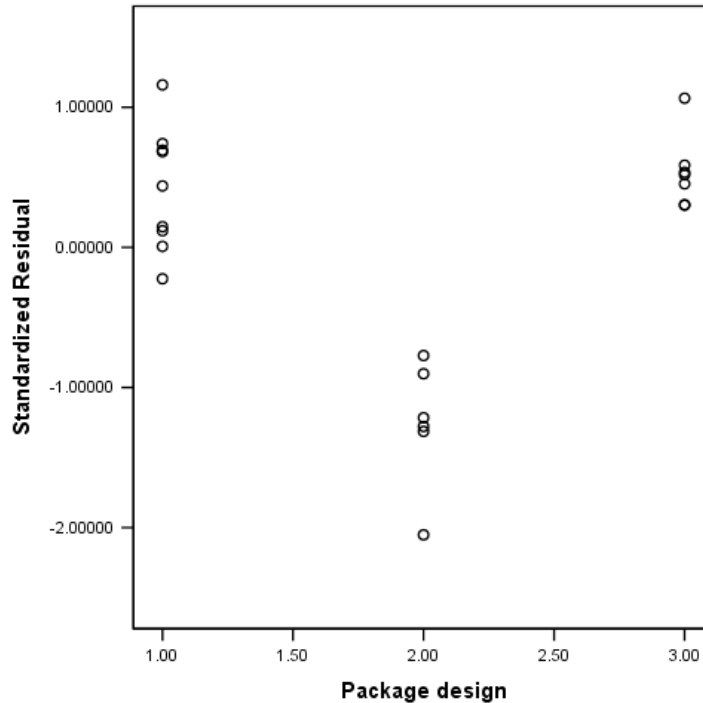


- Select *Standardized Residual* as the y-axis variable and *Package design* as the x-axis variable.

- Click OK.

Figure 8-9

Residuals versus package design



The U-shape is more pronounced in the plot of the standardized residuals against package. Every residual for Design B* is negative, whereas all but one of the residuals is positive for the other two designs. Because the linear regression model fits one parameter for each variable, the relationship cannot be captured by the standard approach.

A Categorical Regression Analysis

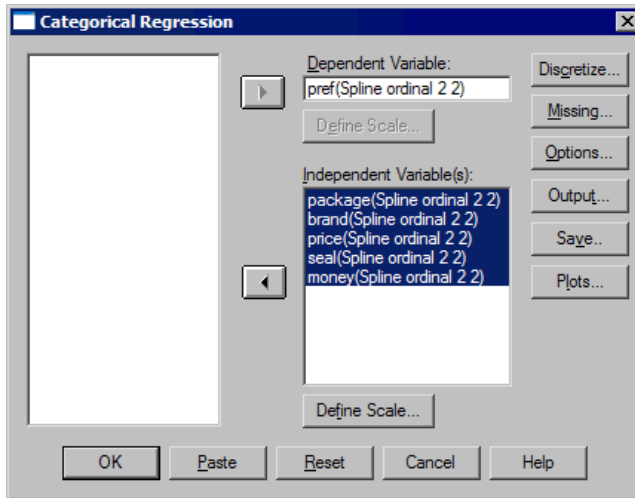
The categorical nature of the variables and the nonlinear relationship between *Preference* and *Package design* suggest that regression on optimal scores may perform better than standard regression. The U-shape of the residual plots indicates that a nominal treatment of *Package design* should be used. All other predictors will be treated at the numerical scaling level.

The response variable warrants special consideration. You want to predict the values of *Preference*. Thus, recovering as many properties of its categories as possible in the quantifications is desirable. Using an ordinal or nominal scaling level ignores the differences between the response categories. However, linearly transforming the response categories preserves category differences. Consequently, scaling the response numerically is generally preferred and will be employed here.

Running the Analysis

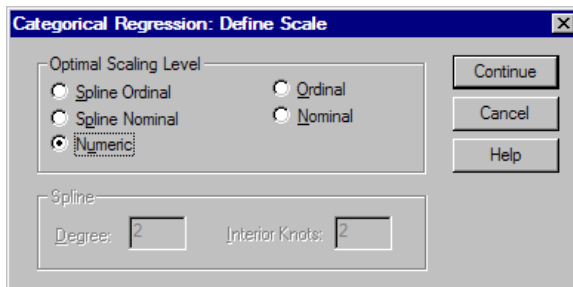
- ▶ To run a Categorical Regression analysis, from the menus choose:
 - Analyze
 - Regression
 - Optimal Scaling...

Figure 8-10
Categorical Regression dialog box



- ▶ Select *Preference* as the dependent variable.
- ▶ Select *Package design* through *Money-back guarantee* as independent variables.
- ▶ Select *Preference* and click *Define Scale*.

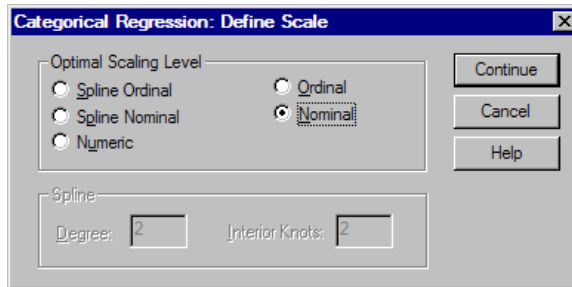
Figure 8-11
Define Scale dialog box



- ▶ Select *Numeric* as the optimal scaling level.
- ▶ Click *Continue*.

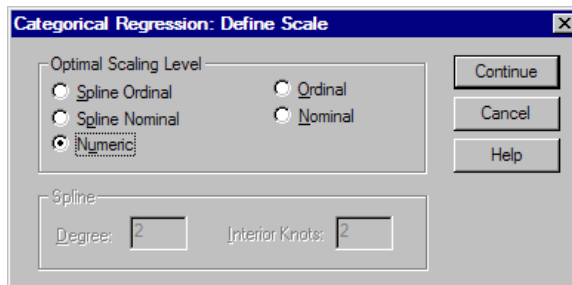
- ▶ Select *Package design* and click Define Scale in the Categorical Regression dialog box.

Figure 8-12
Define Scale dialog box



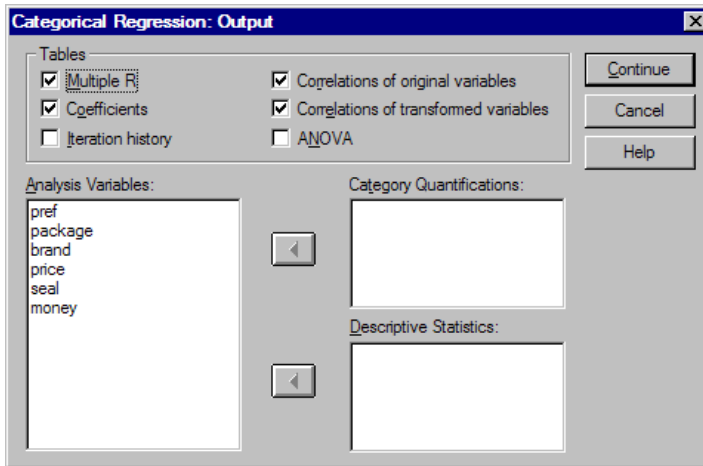
- ▶ Select Nominal as the optimal scaling level.
- ▶ Click Continue.
- ▶ Select *Brand name* through *Money-back guarantee* and click Define Scale in the Categorical Regression dialog box.

Figure 8-13
Define Scale dialog box



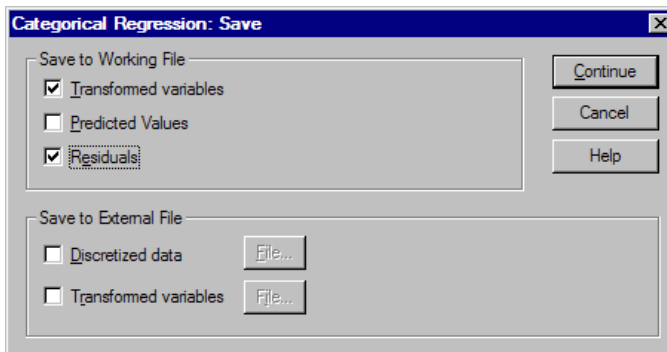
- ▶ Select Numeric as the optimal scaling level.
- ▶ Click Continue.
- ▶ Click Output in the Categorical Regression dialog box.

Figure 8-14
Output dialog box



- ▶ Select Correlations of original variables and Correlations of transformed variables.
- ▶ Deselect ANOVA.
- ▶ Click Continue.
- ▶ Click Save in the Categorical Regression dialog box.

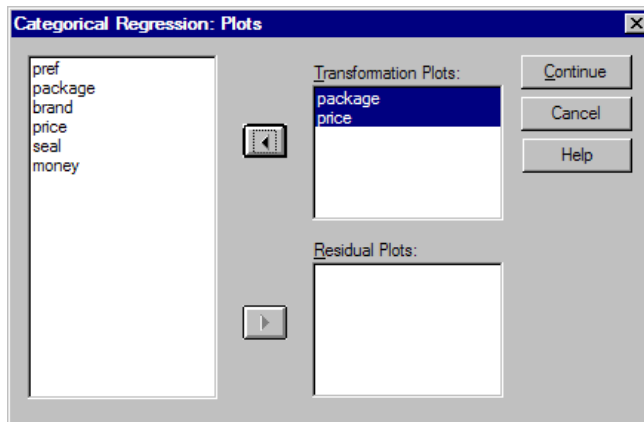
Figure 8-15
Save dialog box



- ▶ Select Transformed variables and Residuals.

- ▶ Click Continue.
- ▶ Click Plots in the Categorical Regression dialog box.

Figure 8-16
Plots dialog box



- ▶ Choose to create transformation plots for *package* and *price*.
- ▶ Click Continue.
- ▶ Click OK in the Categorical Regression dialog box.

Intercorrelations

The intercorrelations among the predictors are useful for identifying multicollinearity in the regression. Variables that are highly correlated will lead to unstable regression estimates. However, due to their high correlation, omitting one of them from the model only minimally affects prediction. The variance in the response that can be explained by the omitted variable is still explained by the remaining correlated variable. However, zero-order correlations are sensitive to outliers and also cannot identify multicollinearity due to a high correlation between a predictor and a combination of other predictors.

Figure 8-17
Original predictor correlations

	Package design	Brand name	Price	Good Housekeeping seal	Money-back guarantee
Package design	1.000	-.189	-.126	.081	.066
Brand name	-.189	1.000	.065	-.042	-.034
Price	-.126	.065	1.000	.000	.000
Good Housekeeping seal	.081	-.042	.000	1.000	-.039
Money-back guarantee	.066	-.034	.000	-.039	1.000
Dimension	1	2	3	4	5
Eigenvalue	1.291	1.038	.980	.905	.785

Figure 8-18
Transformed predictor correlations

	Package design	Brand name	Price	Good Housekeeping seal	Money-back guarantee
Package design	1.000	-.156	-.089	.032	.102
Brand name	-.156	1.000	.065	-.042	-.034
Price	-.089	.065	1.000	.000	.000
Good Housekeeping seal	.032	-.042	.000	1.000	-.039
Money-back guarantee	.102	-.034	.000	-.039	1.000
Dimension	1	2	3	4	5
Eigenvalue	1.248	1.043	.983	.905	.821

The intercorrelations of the predictors for both the untransformed and transformed predictors are displayed. All values are near 0, indicating that multicollinearity between individual variables is not a concern.

Notice that the only correlations that change involve *Package design*. Because all other predictors are treated numerically, the differences between the categories and the order of the categories are preserved for these variables. Consequently, the correlations cannot change.

Model Fit and Coefficients

The Categorical Regression procedure yields an R^2 of 0.948, indicating that almost 95% of the variance in the transformed preference rankings is explained by the regression on the optimally transformed predictors. Transforming the predictors improves the fit over the standard approach.

Figure 8-19
Model summary for categorical regression

Multiple R	R Square	Adjusted R Square
.974	.948	.927

Dependent Variable: Preference
 Predictors: Package design Brand name Price
 Good Housekeeping seal Money-back guarantee

The following table shows the standardized regression coefficients. Categorical regression standardizes the variables, so only standardized coefficients are reported. These values are divided by their corresponding standard errors, yielding an *F* test for each variable. However, the test for each variable is contingent upon the other predictors being in the model. In other words, the test determines if omission of a predictor variable from the model with all other predictors present significantly worsens the predictive capabilities of the model. These values should not be used to omit several variables at one time for a subsequent model. Moreover, alternating least squares optimizes the quantifications, implying that these tests must be interpreted conservatively.

Figure 8-20
Standardized coefficients for transformed predictors

	Standardized Coefficients		df	F	Sig.
	Beta	Std. Error			
Package design	-.748	.060	2	155.289	.000
Brand name	.045	.060	1	.578	.459
Price	.371	.059	1	39.312	.000
Good Housekeeping seal	-.350	.059	1	35.299	.000
Money-back guarantee	-.159	.059	1	7.175	.017

Dependent Variable: Preference

The largest coefficient occurs for *Package design*. A one standard deviation increase in *Package design* yields a 0.748 standard deviation decrease in predicted preference ranking. However, *Package design* is treated nominally, so an increase in the quantifications need not correspond to an increase in the original category codes.

Standardized coefficients are often interpreted as reflecting the importance of each predictor. However, regression coefficients cannot fully describe the impact of a predictor or the relationships between the predictors. Alternative statistics must

be used in conjunction with the standardized coefficients to fully explore predictor effects.

Correlations and Importance

To interpret the contributions of the predictors to the regression, it is not sufficient to only inspect the regression coefficients. In addition, the correlations, partial correlations, and part correlations should be inspected. The following table contains these correlational measures for each variable.

The zero-order correlation is the correlation between the transformed predictor and the transformed response. For this data, the largest correlation occurs for *Package design*. However, if you can explain some of the variation in either the predictor or the response, you will get a better representation of how well the predictor is doing.

Figure 8-21
Zero-order, part, and partial correlations (transformed variables)

	Correlations			Importance	Tolerance	
	Zero-Order	Partial	Part		After Transformation	Before Transformation
Package design	-.816	-.955	-.733	.644	.959	.942
Brand name	.206	.193	.045	.010	.971	.961
Price	.440	.851	.369	.172	.989	.982
Good Housekeeping seal	-.370	-.838	-.349	.137	.996	.991
Money-back guarantee	-.223	-.569	-.158	.037	.987	.993

Dependent Variable: Preference

Other variables in the model can confound the performance of a given predictor in predicting the response. The partial correlation coefficient removes the linear effects of other predictors from both the predictor and the response. This measure equals the correlation between the residuals from regressing the predictor on the other predictors and the residuals from regressing the response on the other predictors. The squared partial correlation corresponds to the proportion of the variance explained relative to the residual variance of the response remaining after removing the effects of the other variables. For example, *Package design* has a partial correlation of -0.955 . Removing the effects of the other variables, *Package design* explains $(-0.955)^2 = 0.91 = 91\%$ of the variation in the preference rankings. Both *Price* and *Good Housekeeping seal* also explain a large portion of variance if the effects of the other variables are removed.

As an alternative to removing the effects of variables from both the response and a predictor, you can remove the effects from just the predictor. The correlation between the response and the residuals from regressing a predictor on the other predictors is the part correlation. Squaring this value yields a measure of the proportion of variance explained relative to the total variance of response. If you remove the effects of *Brand name*, *Good Housekeeping seal*, *Money back guarantee*, and *Price* from *Package design*, the remaining part of *Package design* explains $(-0.733)^2 = 0.54 = 54\%$ of the variation in preference rankings.

Importance

In addition to the regression coefficients and the correlations, Pratt's measure of relative importance (Pratt, 1987) aids in interpreting predictor contributions to the regression. Large individual importances relative to the other importances correspond to predictors that are crucial to the regression. Also, the presence of suppressor variables is signaled by a low importance for a variable that has a coefficient of similar size to the important predictors.

In contrast to the regression coefficients, this measure defines the importance of the predictors additively—that is, the importance of a set of predictors is the sum of the individual importances of the predictors. Pratt's measure equals the product of the regression coefficient and the zero-order correlation for a predictor. These products add to R^2 , so they are divided by R^2 , yielding a sum of 1. The set of predictors *Package design* and *Brand name*, for example, have an importance of 0.654. The largest importance corresponds to *Package design*, with *Package design*, *Price*, and *Good Housekeeping seal* accounting for 95% of the importance for this combination of predictors.

Multicollinearity

Large correlations between predictors will dramatically reduce a regression model's stability. Correlated predictors result in unstable parameter estimates. Tolerance reflects how much the independent variables are linearly related to one another. This measure is the proportion of a variable's variance not accounted for by other independent variables in the equation. If the other predictors can explain a large amount of a predictor's variance, that predictor is not needed in the model. A tolerance value near 1 indicates that the variable cannot be predicted very well from the other predictors. In contrast, a variable with a very low tolerance contributes little

information to a model, and can cause computational problems. Moreover, large negative values of Pratt's importance measure indicate multicollinearity.

All of the tolerance measures are very high. None of the predictors are predicted very well by the other predictors and multicollinearity is not present.

Transformation Plots

Plotting the original category values against their corresponding quantifications can reveal trends that might not be noticed in a list of the quantifications. Such plots are commonly referred to as transformation plots. Attention should be given to categories that receive similar quantifications. These categories affect the predicted response in the same manner. However, the transformation type dictates the basic appearance of the plot.

Variables treated as numerical result in a linear relationship between the quantifications and the original categories, corresponding to a straight line in the transformation plot. The order and the difference between the original categories is preserved in the quantifications.

The order of the quantifications for variables treated as ordinal correspond to the order of the original categories. However, the differences between the categories are not preserved. As a result, the transformation plot is nondecreasing but need not be a straight line. If consecutive categories correspond to similar quantifications, the category distinction may be unnecessary and the categories could be combined. Such categories result in a plateau on the transformation plot. However, this pattern can also result from imposing an ordinal structure on a variable that should be treated as nominal. If a subsequent nominal treatment of the variable reveals the same pattern, combining categories is warranted. Moreover, if the quantifications for a variable treated as ordinal fall along a straight line, a numerical transformation may be more appropriate.

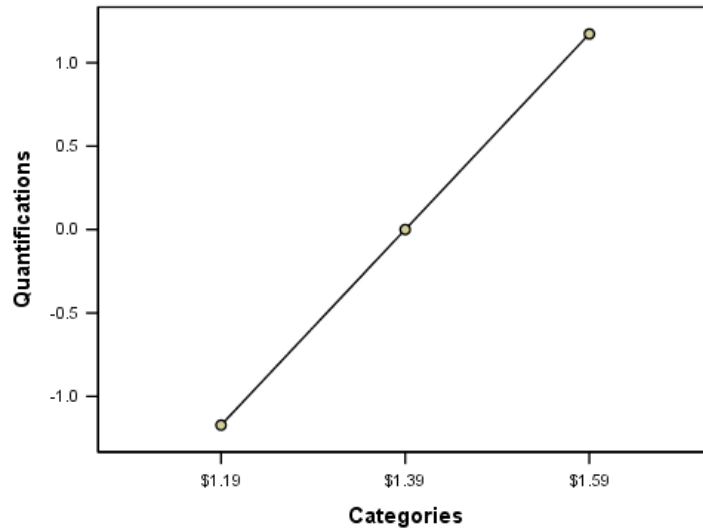
For variables treated as nominal, the order of the categories along the horizontal axis corresponds to the order of the codes used to represent the categories. Interpretations of category order or of the distance between the categories is unfounded. The plot can assume any nonlinear or linear form. If an increasing trend is present, an ordinal treatment should be attempted. If the nominal transformation plot displays a linear trend, a numerical transformation may be more appropriate.

The following figure displays the transformation plot for *Price*, which was treated as numerical. Notice that the order of the categories along the straight line correspond to the order of the original categories. Also, the difference between the quantifications

for \$1.19 and \$1.39 (-1.173 and 0) is the same as the difference between the quantifications for \$1.39 and \$1.59 (0 and 1.173). The fact that categories 1 and 3 are the same distance from category 2 is preserved in the quantifications.

Figure 8-22

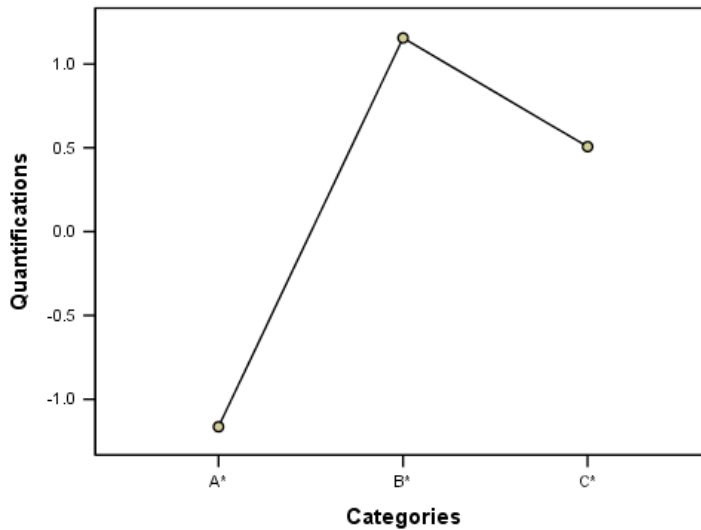
Transformation plot of Price (numerical)



The nominal transformation of *Package design* yields the following transformation plot. Notice the distinct nonlinear shape in which the second category has the largest quantification. In terms of the regression, the second category decreases predicted preference ranking, whereas the first and third categories have the opposite effect.

Figure 8-23

Transformation plot of Package design (nominal)

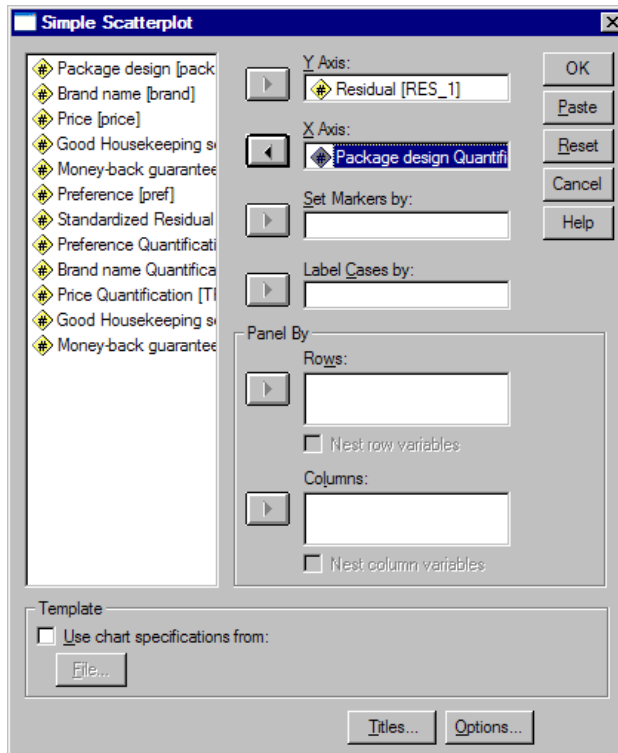


Residual Analysis

Using the transformed data and residuals that you saved to the working file allows you to create a scatterplot of the predicted values by the transformed values of *Package design*.

To obtain such a scatterplot, recall the Simple Scatterplot dialog box and click Reset to clear your previous selections and restore the default options.

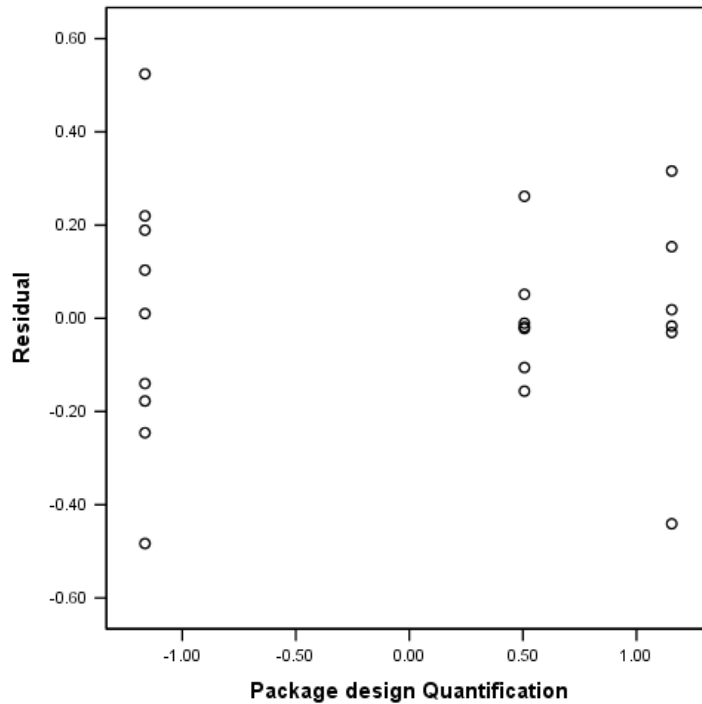
Figure 8-24
Simple Scatterplot dialog box



- ▶ Select *Residual* as the y-axis variable.
- ▶ Select *Package design Quantification* as the x-axis variable.
- ▶ Click OK.

The scatterplot shows the standardized residuals plotted against the optimal scores for *Package design*. All of the residuals are within two standard deviations of 0. A random scatter of points replaces the U-shape present in the scatterplot from the standard linear regression. Predictive abilities are improved by optimally quantifying the categories.

Figure 8-25
Residuals for Categorical Regression



Example: Ozone Data

In this example, you will use a larger set of data to illustrate the selection and effects of optimal scaling transformations. The data include 330 observations on six meteorological variables previously analyzed by Breiman and Friedman (Breiman and Friedman, 1985), and Hastie and Tibshirani (Hastie and Tibshirani, 1990), among others. The following table describes the original variables. Your categorical regression attempts to predict the ozone concentration from the remaining variables. Previous researchers found nonlinearities among these variables, which hinder standard regression approaches.

Table 8-2
Original variables

Variable	Description
<i>ozon</i>	daily ozone level; categorized into one of 38 categories
<i>ibh</i>	inversion base height
<i>dpg</i>	pressure gradient (mm Hg)
<i>vis</i>	visibility (miles)
<i>temp</i>	temperature (degrees F)
<i>doy</i>	day of the year

This data set can be found in *ozone.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

Discretizing Variables

If a variable has more categories than is practically interpretable, you should modify the categories using the Discretization dialog box to reduce the category range to a more manageable number.

The variable *Day of the year* has a minimum value of 3 and a maximum value of 365. Using this variable in a categorical regression corresponds to using a variable with 365 categories. Similarly, *Visibility (miles)* ranges from 0 to 350. To simplify interpretation of analyses, discretize these variables into equal intervals of length 10.

The variable *Inversion base height* ranges from 111 to 5000. A variable with this many categories results in very complex relationships. However, discretizing this variable into equal intervals of length 100 yields roughly 50 categories. Using a 50-category variable rather than a 5000-category variable simplifies interpretations significantly.

Pressure gradient (mm Hg) ranges from -69 to 107. The procedure omits any categories coded with negative numbers from the analysis, but discretizing this variable into equal intervals of length 10 yields roughly 19 categories.

Temperature (degrees F) ranges from 25 to 93 on the Fahrenheit scale. In order to analyze the data as if it were on the Celsius scale, discretize this variable into equal intervals of length 1.8.

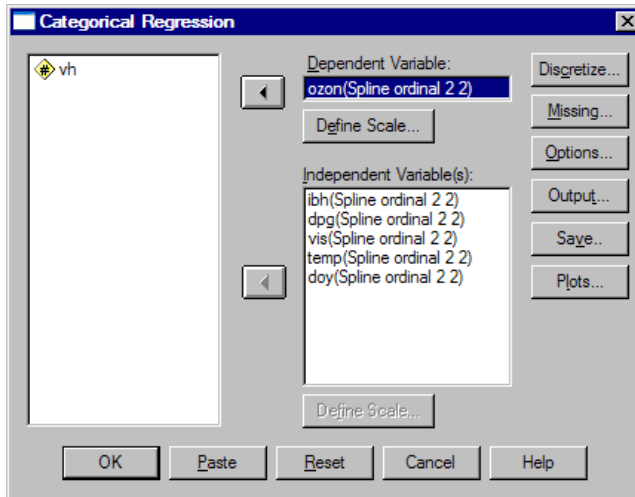
Different discretizations for variables may be desired. The choices used here are purely subjective. If you desire fewer categories, choose larger intervals. For example, *Day of the year* could have been divided into months of the year or seasons.

Selection of Transformation Type

Each variable can be analyzed at one of several different levels. However, because prediction of the response is the goal, you should scale the response “as is” by employing the numerical optimal scaling level. Consequently, the order and the differences between categories will be preserved in the transformed variable.

- ▶ To run a Categorical Regression analysis, from the menus choose:
Analyze
Regression
Optimal Scaling...

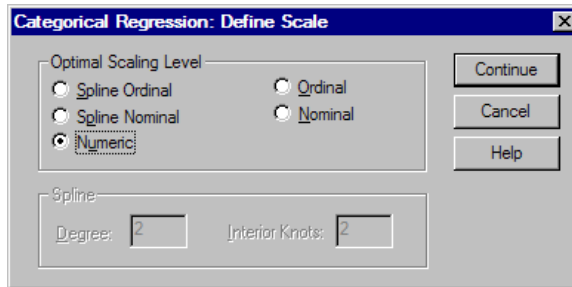
Figure 8-26
Categorical Regression dialog box



- ▶ Select *Daily ozone level* as the dependent variable.
- ▶ Select *Inversion base height* through *Day of the year* as independent variables.

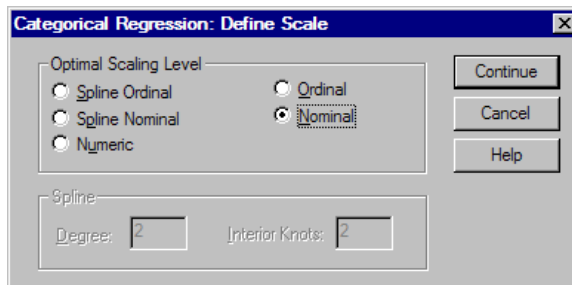
- ▶ Select *Daily ozone level* and click Define Scale.

Figure 8-27
Define Scale dialog box



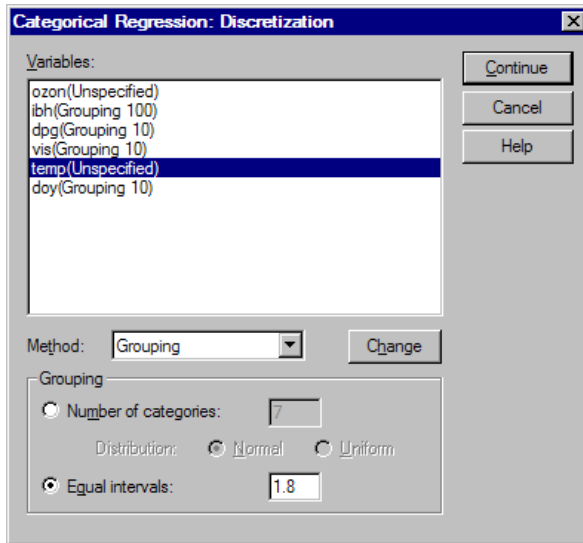
- ▶ Select Numeric as the optimal scaling level.
- ▶ Click Continue.
- ▶ Select *Inversion base height* through *Day of the year*, and click Define Scale in the Categorical Regression dialog box.

Figure 8-28
Define Scale dialog box



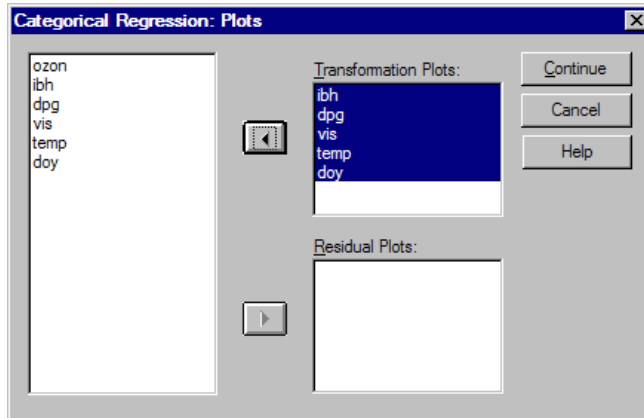
- ▶ Select Nominal as the optimal scaling level.
- ▶ Click Continue.
- ▶ Click Discretize in the Categorical Regression dialog box.

Figure 8-29
Discretization dialog box



- ▶ Select *ibh*.
- ▶ Select Equal intervals and type 100 as the interval length.
- ▶ Click Change.
- ▶ Select *dpq*, *vis*, and *doy*.
- ▶ Type 10 as the interval length.
- ▶ Click Change.
- ▶ Select *temp*.
- ▶ Type 1.8 as the interval length.
- ▶ Click Change.
- ▶ Click Continue.
- ▶ Click Plots in the Categorical Regression dialog box.

Figure 8-30
Plots dialog box



- ▶ Select transformation plots for *ibh* through *doy*.
- ▶ Click Continue.
- ▶ Click OK in the Categorical Regression dialog box.

Figure 8-31
Model summary

Multiple R	R Square	Adjusted R Square
.941	.886	.793

Dependent Variable: Daily ozone level

Predictors: Inversion base height Pressure gradient (mm Hg)

Visibility (miles) Temperature (degrees F) Day of the year

Treating all predictors as nominal yields an R^2 of 0.886. This large amount of variance accounted for is not surprising because nominal treatment imposes no restrictions on the quantifications. However, interpreting the results can be quite difficult.

Figure 8-32
Regression coefficients (all predictors nominal)

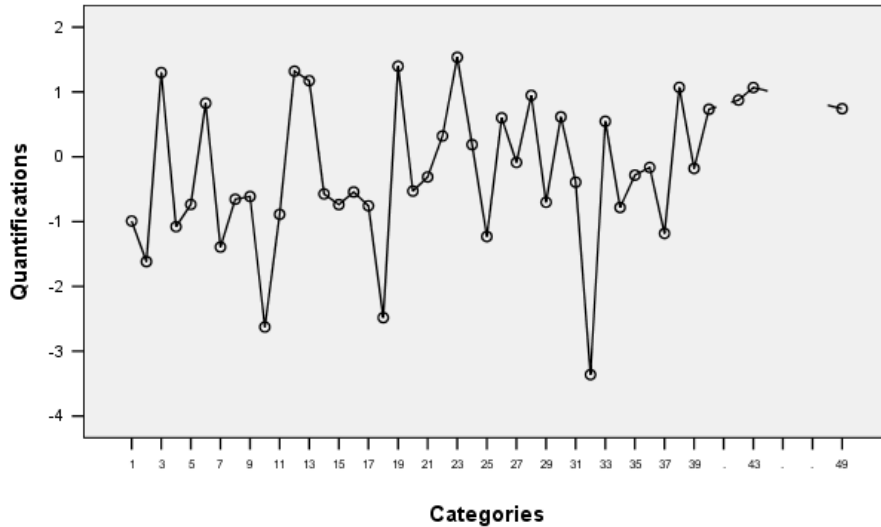
	Standardized Coefficients		df	F	Sig.
	Beta	Std. Error			
Inversion base height	-.309	.026	42	144.187	.000
Pressure gradient (mm Hg)	.307	.028	16	123.227	.000
Visibility (miles)	-.216	.026	17	69.528	.000
Temperature (degrees F)	.588	.027	36	468.542	.000
Day of the year	-.408	.029	36	203.250	.000

Dependent Variable: Daily ozone level

This table shows the standardized regression coefficients of the predictors. A common mistake made when interpreting these values involves focusing on the coefficients while neglecting the quantifications. You cannot assert that the large positive value of the *Temperature* coefficient implies that as temperature increases, predicted *Ozone* increases.

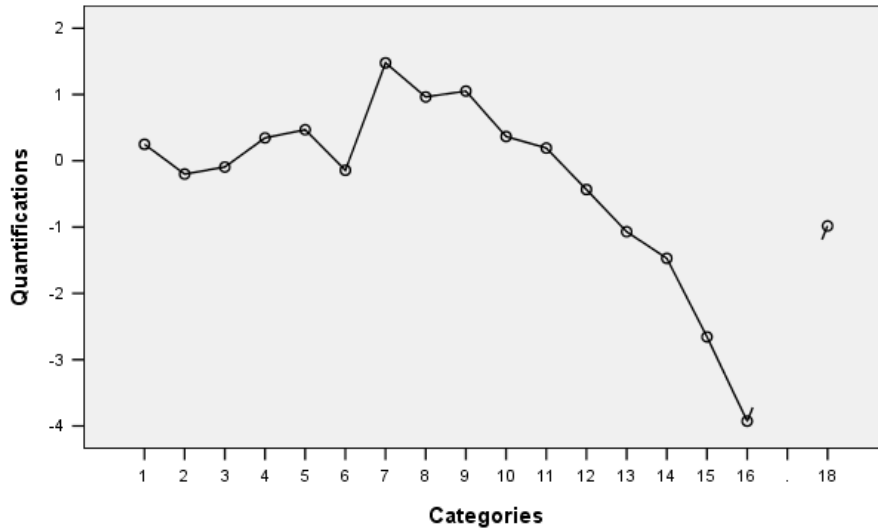
Similarly, the negative coefficient for *Inversion base height* does not suggest that as *Inversion base height* increases, predicted *Ozone* decreases. All interpretations must be relative to the transformed variables. As the quantifications for *Temperature* increase, or as the quantifications for *Inversion base height* decrease, predicted *Ozone* increases. To examine the effects of the original variables, you must relate the categories to the quantifications.

Figure 8-33
Transformation plot of Inversion base height (nominal)



The transformation plot of *Inversion base height* shows no apparent pattern. As evidenced by the jagged nature of the plot, moving from low categories to high categories yields fluctuations in the quantifications in both directions. Thus, describing the effects of this variable requires focusing on the individual categories. Imposing ordinal or linear restrictions on the quantifications for this variable might significantly reduce the fit.

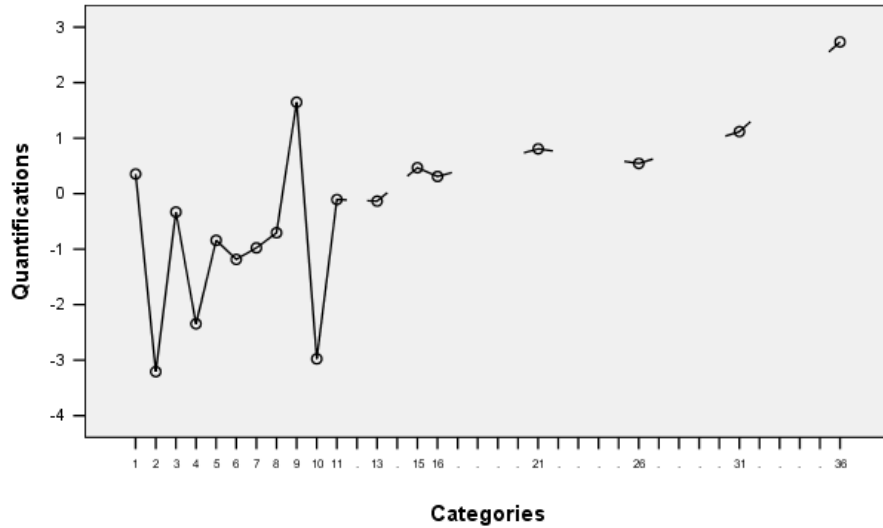
Figure 8-34
Transformation plot of *Pressure gradient* (nominal)



This figure displays the transformation plot of *Pressure gradient*. The initial discretized categories (1 through 6) receive small quantifications and thus have minimal contributions to the predicted response. The next three categories receive somewhat higher, positive values, resulting in a moderate increase in predicted ozone.

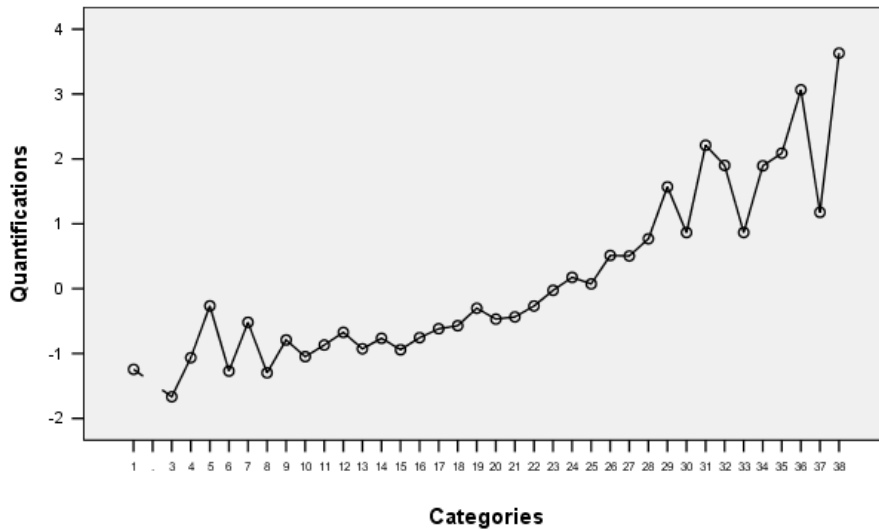
The quantifications decrease up to category 16, where *Pressure gradient* has its greatest decreasing effect on predicted ozone. Although the line increases after this category, using an ordinal scaling level for *Pressure gradient* may not significantly reduce the fit, while simplifying the interpretations of the effects. However, the importance measure of 0.04 and the regression coefficient for *Pressure gradient* indicates that this variable is not very useful in the regression.

Figure 8-35
Transformation plot of Visibility (nominal)



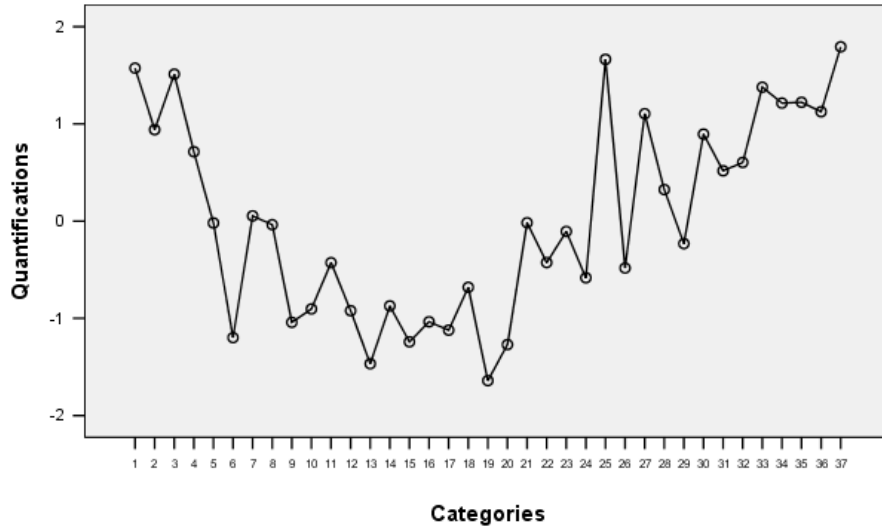
The transformation plot of *Visibility*, like that for *Inversion base height*, shows no apparent pattern. Imposing ordinal or linear restrictions on the quantifications for this variable might significantly reduce the fit.

Figure 8-36
Transformation plot of Temperature (nominal)



The transformation plot of *Temperature* displays an alternative pattern. As the categories increase, the quantifications tend to increase. As a result, as *Temperature* increases, predicted ozone tends to increase. This pattern suggests scaling *Temperature* at the ordinal level.

Figure 8-37
Transformation plot of *Day of the year* (nominal)

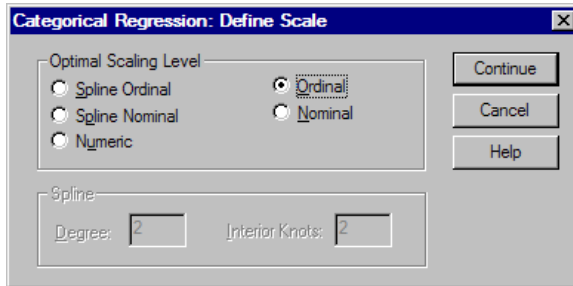


This figure shows the transformation plot of *Day of the year*. The quantifications tend to decrease up to category 19, at which point they tend to increase, yielding a U-shape. Considering the sign of the regression coefficient for *Day of the year*, the initial categories (1 through 5) receive quantifications that have a decreasing effect on predicted ozone. From category 6 onward, the effect of the quantifications on predicted ozone gets more increasing, reaching a maximum around category 19.

Beyond category 19, the quantifications tend to decrease the predicted ozone. Although the line is quite jagged, the general shape is still identifiable. Thus, the transformation plots suggest scaling *Temperature* at the ordinal level while keeping all other predictors nominally scaled.

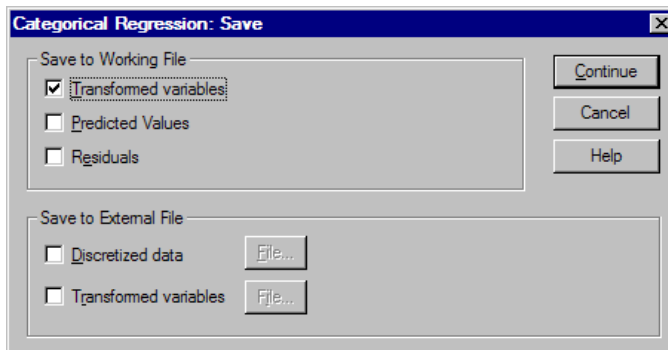
To recompute the regression, scaling *Temperature* at the ordinal level, recall the Categorical Regression dialog box.

Figure 8-38
Define Scale dialog box



- ▶ Select *Temperature* and click Define Scale.
- ▶ Select Ordinal as the optimal scaling level.
- ▶ Click Continue.
- ▶ Click Save in the Categorical Regression dialog box.

Figure 8-39
Save dialog box



- ▶ Select Transformed variables in the Save to Working File group.
- ▶ Click Continue.
- ▶ Click OK in the Categorical Regression dialog box.

Figure 8-40

Model summary for regression with *Temperature* (ordinal)

Multiple R	R Square	Adjusted R Square
.935	.875	.791

Dependent Variable: Daily ozone level

Predictors: Inversion base height Pressure gradient (mm Hg)

Visibility (miles) Temperature (degrees F) Day of the year

This model results in an R^2 of 0.875, so the variance accounted for decreases negligibly when the quantifications for *Temperature* are restricted to be ordered.

Figure 8-41

Regression coefficients with *Temperature* (ordinal)

	Standardized Coefficients		df	F	Sig.
	Beta	Std. Error			
Inversion base height	-.300	.026	42	132.628	.000
Pressure gradient (mm Hg)	.294	.028	16	111.757	.000
Visibility (miles)	-.221	.026	17	72.360	.000
Temperature (degrees F)	.619	.027	20	517.262	.000
Day of the year	-.372	.028	36	173.158	.000

Dependent Variable: Daily ozone level

This table displays the coefficients for the model in which *Temperature* is scaled as ordinal. Comparing the coefficients to those for the model in which *Temperature* is scaled as nominal, no large changes occur.

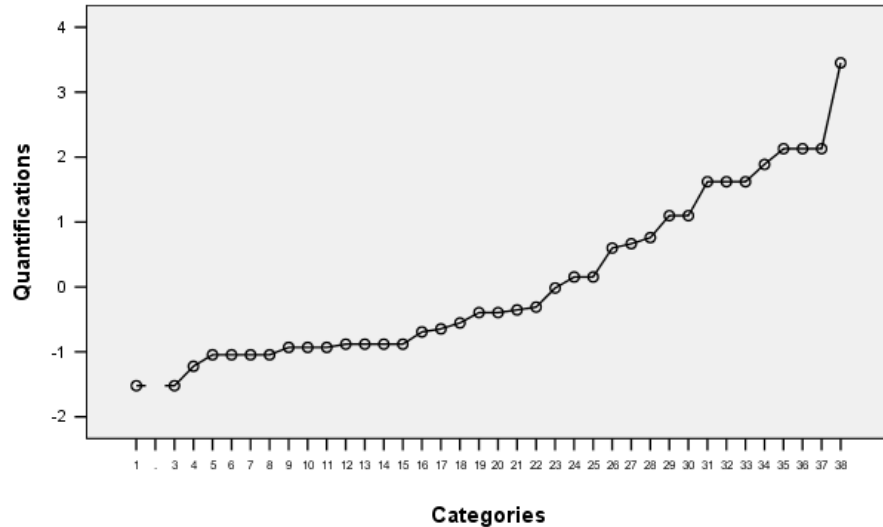
Figure 8-42
Correlations, importance, and tolerance

	Correlations			Importance	Tolerance	
	Zero-Order	Partial	Part		After Transformation	Before Transformation
Inversion base height	-.435	-.633	-.290	.150	.931	.596
Pressure gradient (mm Hg)	.135	.601	.266	.045	.819	.859
Visibility (miles)	-.352	-.517	-.214	.089	.939	.752
Temperature (degrees F)	.806	.850	.573	.571	.855	.580
Day of the year	-.340	-.683	-.331	.145	.791	.801

Dependent Variable: Daily ozone level

Moreover, the importance measures suggest that *Temperature* is still much more important to the regression than the other variables. Now, however, as a result of the ordinal scaling level of *Temperature* and the positive regression coefficient, you can assert that as *Temperature* increases, predicted ozone increases.

Figure 8-43
Transformation plot of Temperature (ordinal)



The transformation plot illustrates the ordinal restriction on the quantifications for *Temperature*. The jagged line from the nominal transformation is replaced here by a smooth ascending line. Moreover, no long plateaus are present, indicating that collapsing categories is not needed.

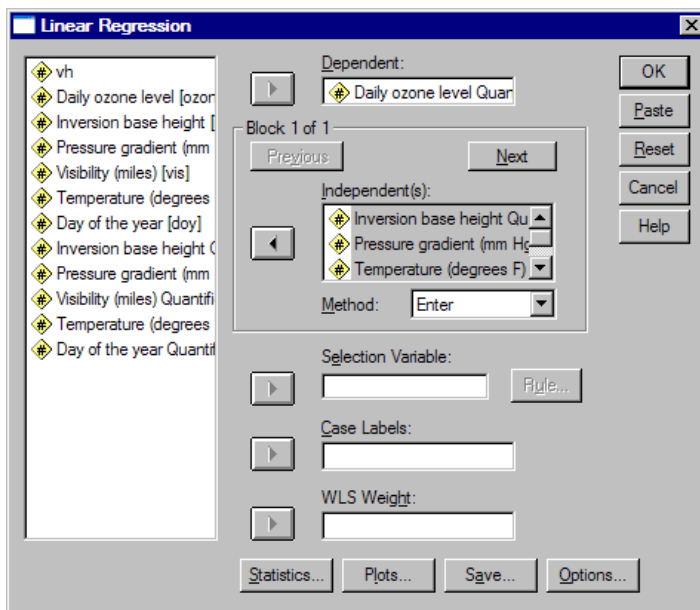
Optimality of the Quantifications

The transformed variables from a categorical regression can be used in a standard linear regression, yielding identical results. However, the quantifications are optimal only for the model that produced them. Using a subset of the predictors in linear regression does not correspond to an optimal scaling regression on the same subset.

For example, the categorical regression that you have computed has an R^2 of 0.875. You have saved the transformed variables, so in order to fit a linear regression using only *Temperature*, *Pressure gradient*, and *Inversion base height* as predictors, from the menus choose:

Analyze
Regression
Linear...

Figure 8-44
Linear Regression dialog box



- ▶ Select *Daily ozone level Quantification* as the dependent variable.
- ▶ Select *Inversion base height Quantification*, *Pressure gradient (mm Hg) Quantification*, and *Temperature (degrees F) Quantification* as independent variables.
- ▶ Click OK.

Figure 8-45

Model summary for regression with subset of optimally scaled predictors

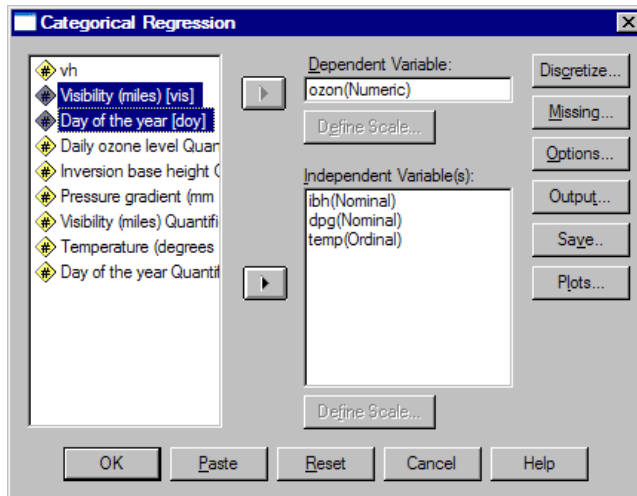
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.856 ^a	.733	.730	.51993

a. Predictors: (Constant), Temperature (degrees F)
Quantification, Pressure gradient (mm Hg)
Quantification, Inversion base height Quantification

Using the quantifications for the response, *Temperature*, *Pressure gradient*, and *Inversion base height* in a standard linear regression results in a fit of 0.733. To compare this to the fit of a categorical regression using just those three predictors, recall the Categorical Regression dialog box.

Figure 8-46

Categorical Regression dialog box



- ▶ Deselect *Visibility (miles)* and *Day of the year* as independent variables.
- ▶ Click OK.

Figure 8-47

Model summary for categorical regression on three predictors

Multiple R	R Square	Adjusted R Square
.893	.798	.740

Dependent Variable: Daily ozone level
 Predictors: Inversion base height Pressure gradient (mm Hg) Temperature (degrees F)

The categorical regression analysis has a fit of 0.798, which is better than the fit of 0.733. This demonstrates the property of the scalings that the quantifications obtained in the original regression are only optimal when all five variables are included in the model.

Effects of Transformations

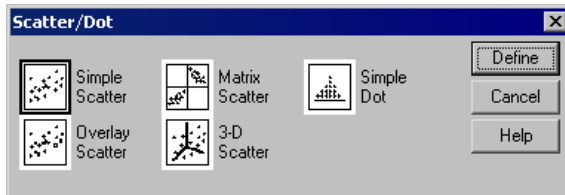
Transforming the variables makes a nonlinear relationship between the original response and the original set of predictors linear for the transformed variables. However, when there are multiple predictors, pairwise relationships are confounded by the other variables in the model.

To focus your analysis on the relationship between *Daily ozone level* and *Day of the year*, begin by looking at a scatterplot. From the menus choose:

Graphs
 Scatter/Dot...

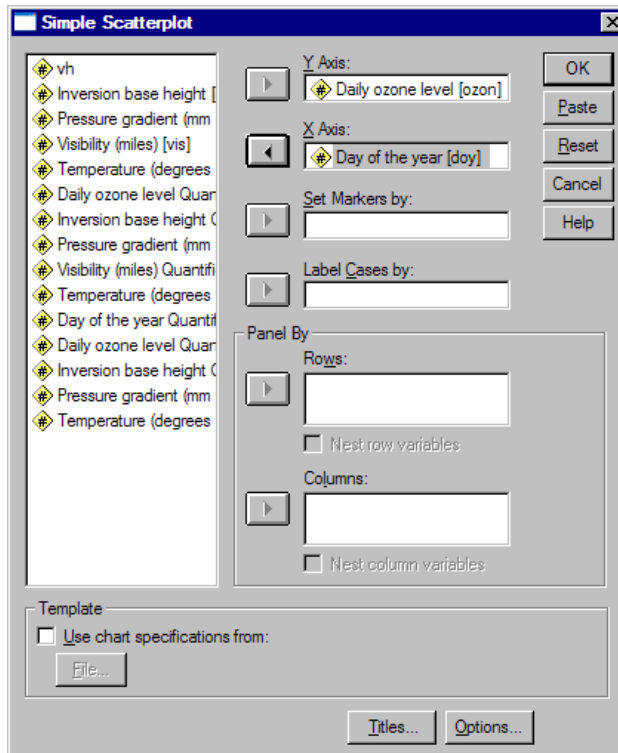
Figure 8-48

Scatter/Dot dialog box



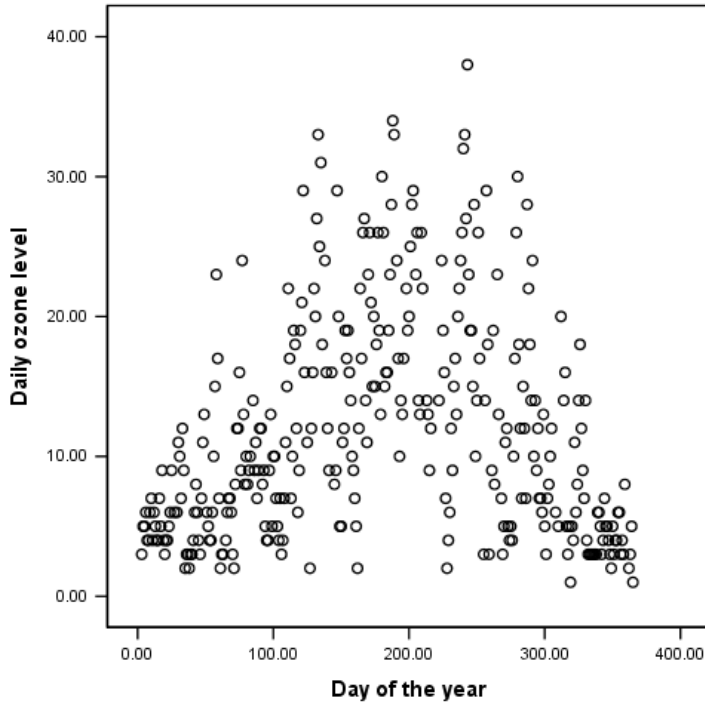
- Click Define.

Figure 8-49
Simple Scatterplot dialog box



- ▶ Select *Daily ozone level* as the y-axis variable and *Day of the year* as the x-axis variable.
- ▶ Click OK.

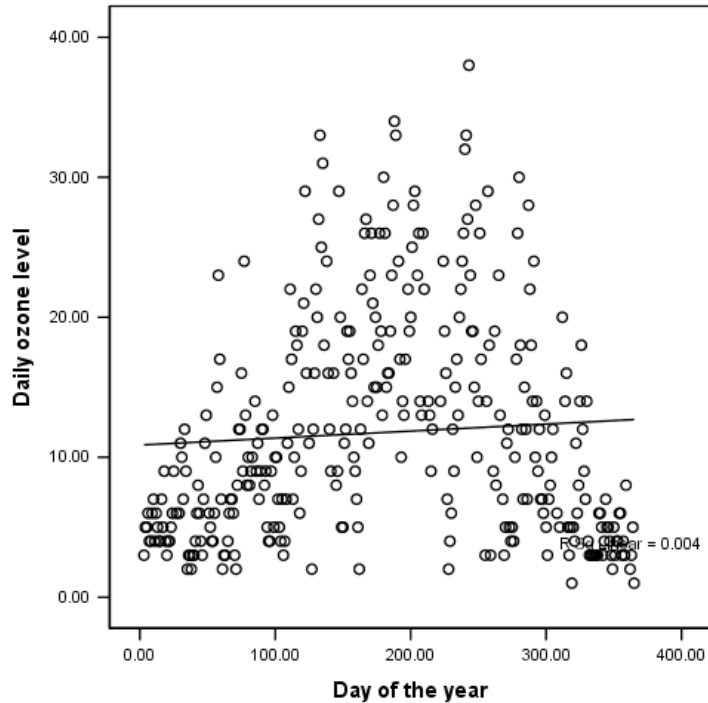
Figure 8-50
Scatterplot of Daily ozone level and Day of the year



This figure illustrates the relationship between *Daily ozone level* and *Day of the year*. As *Day of the year* increases to approximately 200, *Daily ozone level* increases. However, for *Day of the year* values greater than 200, *Daily ozone level* decreases. This inverted U pattern suggests a quadratic relationship between the two variables. A linear regression cannot capture this relationship.

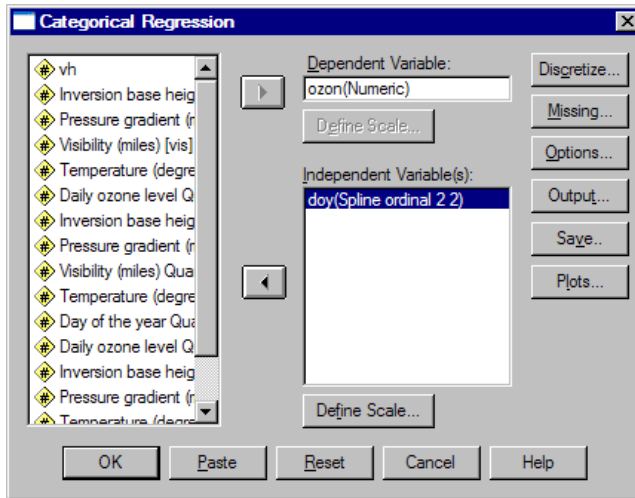
- ▶ To see a best-fit line overlaid on the points in the scatterplot, activate the graph by double-clicking on it.
- ▶ Select a point in the Chart Editor.
- ▶ Click the Add Fit Line at Total tool, and close the Chart Editor.

Figure 8-51
Scatterplot showing best-fit line



A linear regression of *Daily ozone level* on *Day of the year* yields an R^2 of 0.004. This fit suggests that *Day of the year* has no predictive value for *Daily ozone level*. This is not surprising, given the pattern in the figure. By using optimal scaling, however, you can linearize the quadratic relationship and use the transformed *Day of the year* to predict the response.

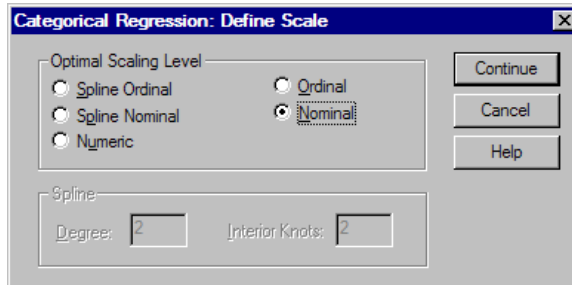
Figure 8-52
Categorical Regression dialog box



To obtain a categorical regression of *Daily ozone level* on *Day of the year*, recall the Categorical Regression dialog box.

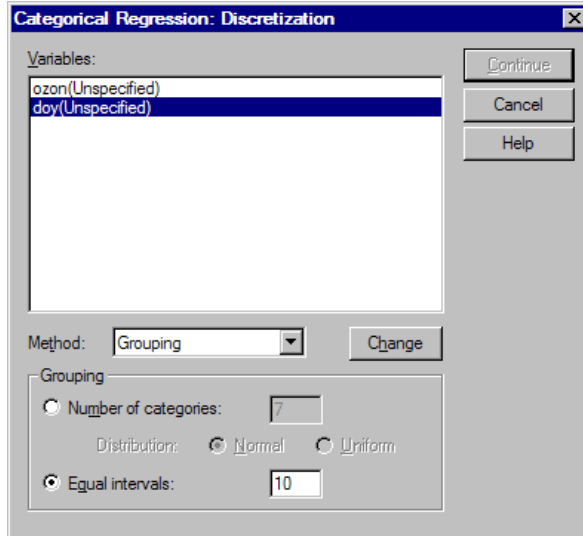
- ▶ Deselect *Inversion base height* through *Temperature (degrees F)* as independent variables.
- ▶ Select *Day of the year* as an independent variable.
- ▶ Click Define Scale.

Figure 8-53
Define Scale dialog box



- ▶ Select Nominal as the optimal scaling level.
- ▶ Click Continue.
- ▶ Click Discretize in the Categorical Regression dialog box.

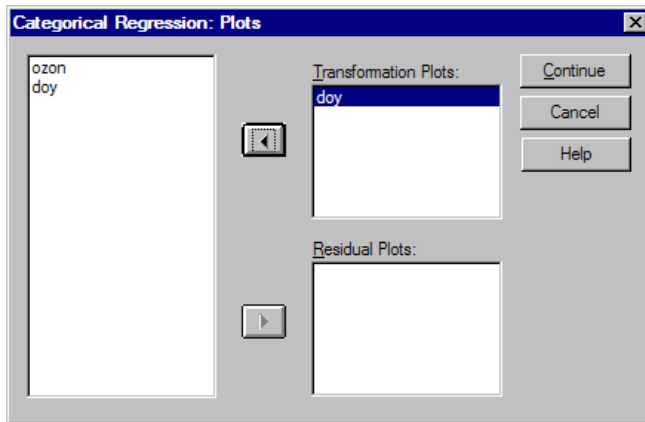
Figure 8-54
Discretization dialog box



- ▶ Select *doy*.
- ▶ Select Equal intervals.

- ▶ Type 10 as the interval length.
- ▶ Click Change.
- ▶ Click Continue.
- ▶ Click Plots in the Categorical Regression dialog box.

Figure 8-55
Plots dialog box



- ▶ Select *doy* for transformation plots.
- ▶ Click Continue.
- ▶ Click OK in the Categorical Regression dialog box.

Figure 8-56
Model summary for categorical regression of Daily ozone level on Day of the year

Multiple R	R Square	Adjusted R Square
.741	.549	.494

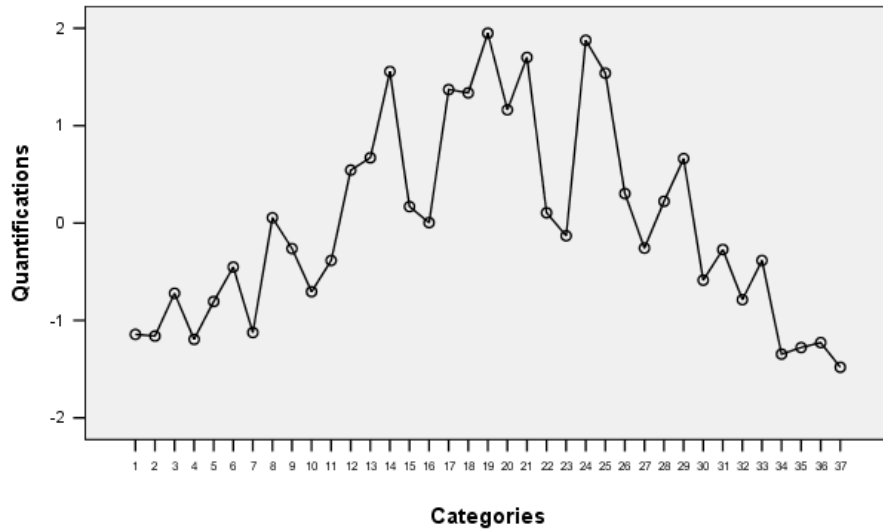
Dependent Variable: Daily ozone level
Predictors: Day of the year

The optimal scaling regression treats *Daily ozone level* as numerical and *Day of the year* as nominal. This results in an R^2 of 0.549. Although only 55% of the variation in *Daily ozone level* is accounted for by the categorical regression, this is a substantial

improvement over the original regression. Transforming *Day of the year* allows for the prediction of *Daily ozone level*.

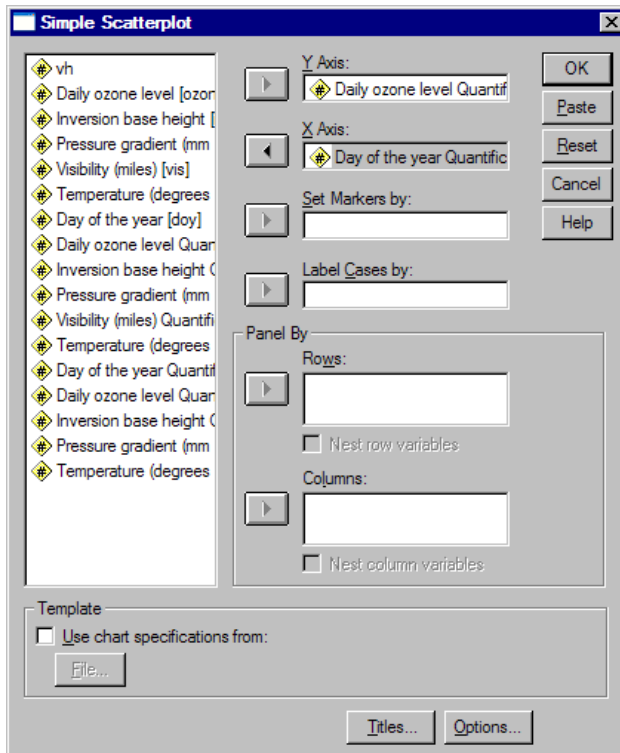
Figure 8-57

Transformation plot of Day of the year (nominal)



This figure displays the transformation plot of *Day of the year*. The extremes of *Day of the year* both receive negative quantifications, whereas the central values have positive quantifications. By applying this transformation, the low and high *Day of the year* values have similar effects on predicted *Daily ozone level*.

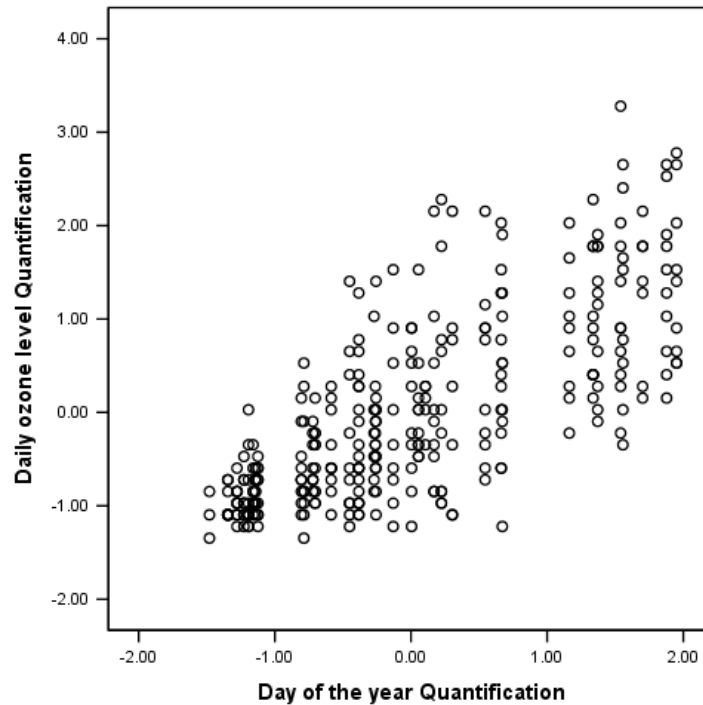
Figure 8-58
Simple Scatterplot dialog box



To see a scatterplot of the transformed variables, recall the Simple Scatterplot dialog box, and click Reset to clear your previous selections.

- ▶ Select *Daily ozone level Quantification [TRA1_3]* as the y-axis variable and *Day of the year Quantification [TRA2_3]* as the x-axis variable.
- ▶ Click OK.

Figure 8-59
Scatterplot of the transformed variables



This figure depicts the relationship between the transformed variables. An increasing trend replaces the inverted U. The regression line has a positive slope, indicating that as transformed *Day of the year* increases, predicted *Daily ozone level* increases. Using optimal scaling linearizes the relationship and allows interpretations that would otherwise go unnoticed.

Recommended Readings

See the following texts for more information on categorical regression:

Hastie, T., R. Tibshirani, and A. Buja. 1994. Flexible discriminant analysis. *Journal of the American Statistical Association*, 89, 1255–1270.

Hayashi, C. 1952. On the prediction of phenomena from qualitative data and the quantification of qualitative data from the mathematico-statistical point of view. *Annals of the Institute of Statistical Mathematics*, 2, 93–96.

Kruskal, J. B. 1965. Analysis of factorial experiments by estimating monotone transformations of the data. *Journal of the Royal Statistical Society Series B*, 27, 251–263.

Meulman, J. J. 2003. Prediction and classification in nonlinear data analysis: Something old, something new, something borrowed, something blue. *Psychometrika*, 4, 493–517.

Ramsay, J. O. 1989. Monotone regression splines in action. *Statistical Science*, 4, 425–441.

Vander Kooij, A. J., and J. J. Meulman. 1997. MURALS: Multiple regression and optimal scaling using alternating least squares. In: *Softstat '97*, F. Faulbaum, and W. Bandilla, eds. Stuttgart: Gustav Fisher, 99–106.

Winsberg, S., and J. O. Ramsay. 1980. Monotonic transformations to additivity using splines. *Biometrika*, 67, 669–674.

Winsberg, S., and J. O. Ramsay. 1983. Monotone spline transformations for dimension reduction. *Psychometrika*, 48, 575–595.

Young, F. W., J. De Leeuw, and Y. Takane. 1976. Regression with qualitative and quantitative variables: An alternating least squares method with optimal scaling features. *Psychometrika*, 41, 505–528.

Categorical Principal Components Analysis

Categorical principal components analysis can be thought of as a method of dimension reduction. A set of variables is analyzed to reveal major dimensions of variation. The original data set can then be replaced by a new, smaller data set with minimal loss of information. The method reveals relationships among variables, among cases, and among variables and cases.

The criterion used by categorical principal components analysis for quantifying the observed data is that the object scores (component scores) should have large correlations with each of the quantified variables. A solution is good to the extent that this criterion is satisfied.

Two examples of categorical principal components analysis will be presented. The first employs a rather small data set useful for illustrating the basic concepts and interpretations associated with the procedure. The second example examines a practical application.

Example: Examining Interrelations of Social Systems

This example examines Guttman's (Guttman, 1968) adaptation of a table by Bell (Bell, 1961). The data are also discussed by Lingoes (Lingoes, 1968).

Bell presented a table to illustrate possible social groups. Guttman 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).

The following table shows the variables in the data set resulting from the classification into seven social groups used in the Guttman-Bell data, with their variable labels and the value labels (categories) associated with the levels of each variable. This data set can be found in *guttman.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS. In addition to selecting variables to be included in the computation of the categorical principal components analysis, you can select variables that are used to label objects in plots. In this example, the first five variables in the data are included in the analysis, while *cluster* is used exclusively as a labeling variable. When you specify a categorical principal components analysis, you must specify the optimal scaling level for each analysis variable. In this example, an ordinal level is specified for all analysis variables.

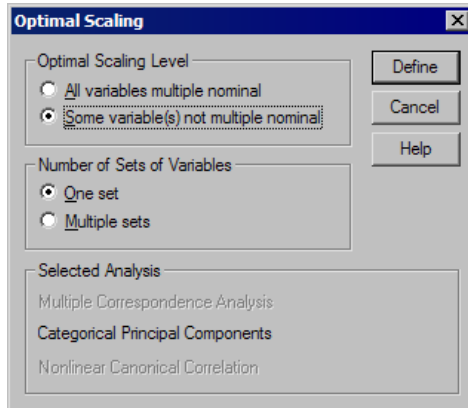
Table 9-1
Variables in the Guttman-Bell data set

Variable name	Variable label	Value label
<i>intnsity</i>	Intensity of interaction	Slight, low, moderate, high
<i>frquency</i>	Frequency of interaction	Slight, nonrecurring, infrequent, frequent
<i>blonging</i>	Feeling of belonging	None, slight, variable, high
<i>proxmity</i>	Physical proximity	Distant, close
<i>formlity</i>	Formality of relationship	No relationship, formal, informal
<i>cluster</i>		Crowds, audiences, public, mobs, primary groups, secondary groups, modern community

Running the Analysis

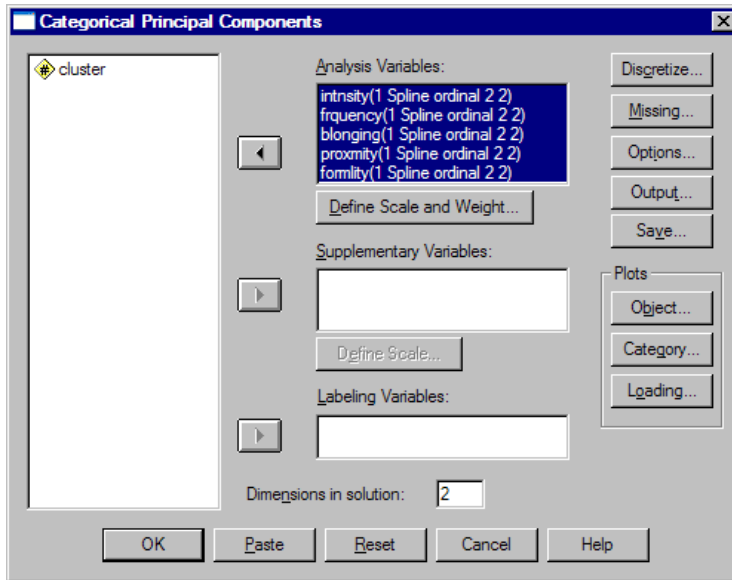
- ▶ To produce categorical principal components output for this data set, from the menus choose:
 - Analyze
 - Data Reduction
 - Optimal Scaling...

Figure 9-1
Optimal Scaling dialog box



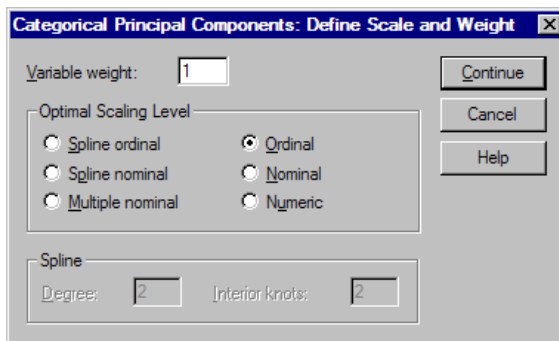
- ▶ Select Some variable(s) not multiple nominal in the Optimal Scaling Level group.
- ▶ Click Define.

Figure 9-2
Categorical Principal Components dialog box



- ▶ Select *Intensity of interaction* through *Formality of relationship* as analysis variables.
- ▶ Click Define Scale and Weight.

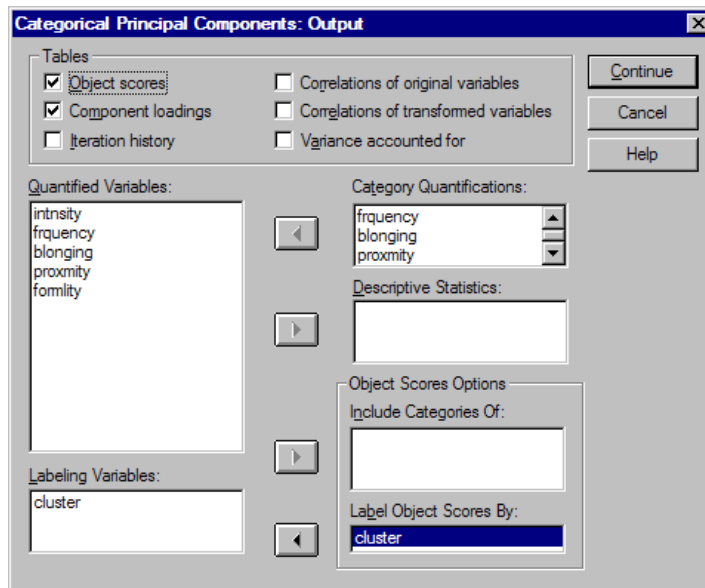
Figure 9-3
Define Scale and Weight dialog box



- ▶ Select Ordinal in the Optimal Scaling Level group.

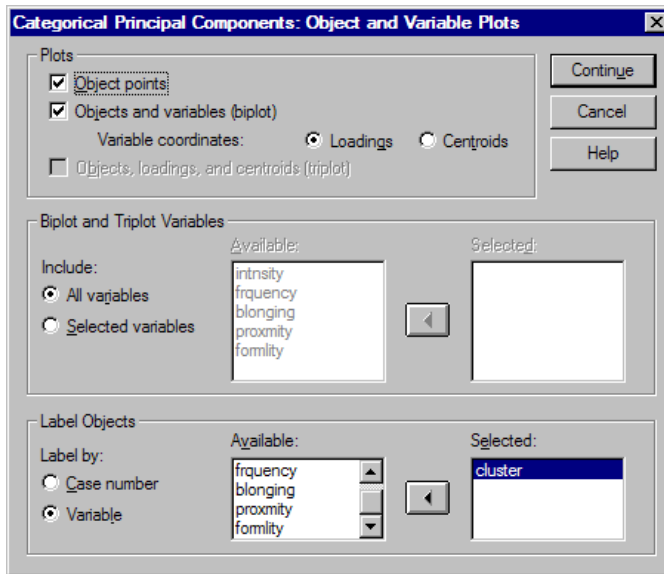
- ▶ Click Continue.
- ▶ Select *cluster* as a labeling variable in the Categorical Principal Components Analysis dialog box.
- ▶ Click Output.

Figure 9-4
Output dialog box



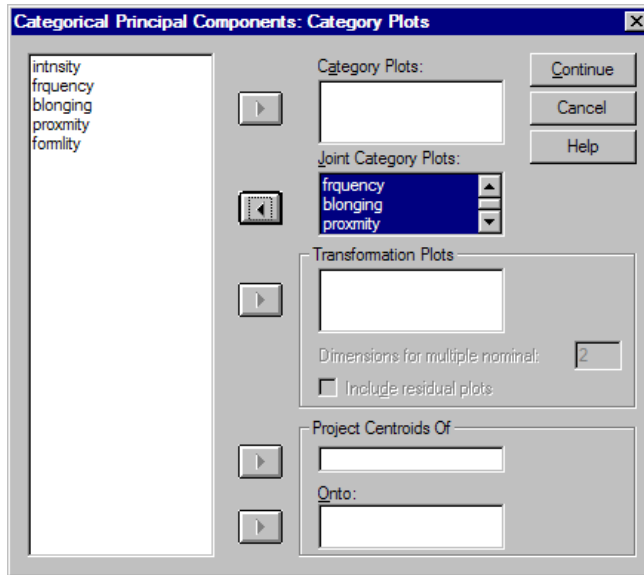
- ▶ Select Object scores and deselect Correlations of transformed variables in the Tables group.
- ▶ Choose to produce category quantifications for *intnsity* (*Intensity of interaction*) through *formlity* (*Formality of relationship*).
- ▶ Choose to label object scores by *cluster*.
- ▶ Click Continue.
- ▶ Click Object in the Plots group of the Categorical Principal Components dialog box.

Figure 9-5
Object and Variable Plots dialog box



- ▶ Select Objects and variables (biplot) in the Plots group.
- ▶ Choose to label objects by Variable in the Label Objects group, and then select *cluster* as the variable to label objects by.
- ▶ Click Continue.
- ▶ Click Category in the Plots group of the Categorical Principal Components dialog box.

Figure 9-6
Category Plots dialog box



- ▶ Choose to produce joint category plots for *intnsity* (*Intensity of interaction*) through *formlity* (*Formality of relationship*).
- ▶ Click Continue.
- ▶ Click OK in the Categorical Principal Components dialog box.

Number of Dimensions

These figures show some of the initial output for the categorical principal components analysis. After the iteration history of the algorithm, the model summary, including the eigenvalues of each dimension, is displayed. These eigenvalues are equivalent to those of classical principal components analysis. They are measures of how much variance is accounted for by each dimension.

Figure 9-7
Iteration history

Iteration Number	Variance Accounted For		Loss		
	Total	Increase	Total	Centroid Coordinates	Restriction of Centroid to Vector Coordinates
0	4.515315	.000000	5.484685	4.075583	1.409101
31 ^a	4.726009	.000008	5.273991	4.273795	1.000196

a. The iteration process stopped because the convergence test value was reached.

Figure 9-8
Model summary

Dimension	Cronbach's Alpha	Variance Accounted For	
		Total (Eigenvalue)	% of Variance
1	.881	3.389	67.774
2	.315	1.337	26.746
Total	.986 ^a	4.726	94.520

a. Total Cronbach's Alpha is based on the total Eigenvalue.

The eigenvalues can be used as an indication of how many dimensions are needed. In this example, the default number of dimensions, 2, was used. Is this the right number? As a general rule, when all variables are either single nominal, ordinal, or numerical, the eigenvalue for a dimension should be larger than 1. Since the two-dimensional solution accounts for 94.52% of the variance, a third dimension probably would not add much more information.

For multiple nominal variables, there is no easy rule of thumb to determine the appropriate number of dimensions. If the number of variables is replaced by the total number of categories minus the number of variables, the above rule still holds. But this rule alone would probably allow more dimensions than are needed. When choosing the number of dimensions, the most useful guideline is to keep the number small enough so that meaningful interpretations are possible. The model summary table also shows Cronbach's alpha (a measure of reliability), which is maximized by the procedure.

Quantifications

For each variable, the quantifications, the vector coordinates, and the centroid coordinates for each dimension are presented. The quantifications are the values assigned to each category. The centroid coordinates are the average of the object scores of objects in the same category. The vector coordinates are the coordinates of the categories when they are required to be on a line, representing the variable in the object space. This is required for variables with the ordinal and numerical scaling level.

Figure 9-9
Quantifications for Intensity of interaction

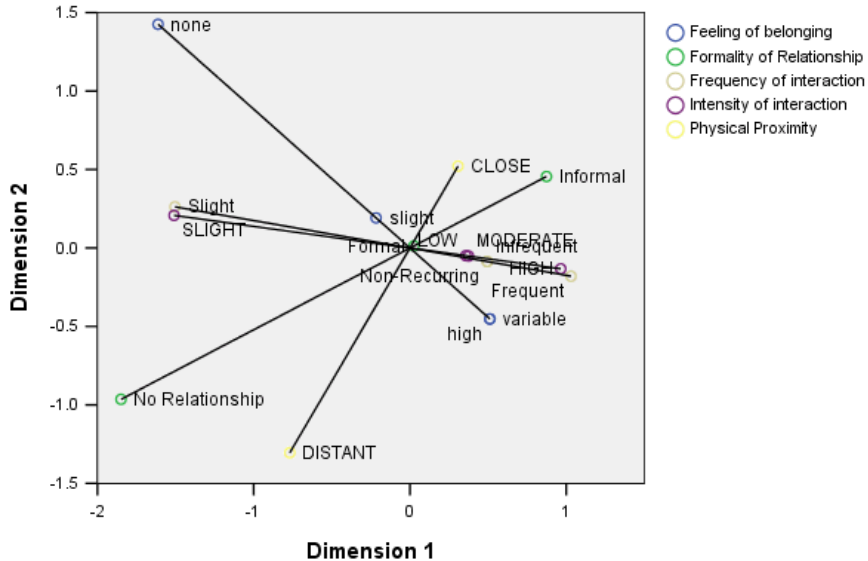
Category	Frequency	Quantification	Centroid Coordinates		Vector Coordinates	
			Dimension		Dimension	
			1	2	1	2
SLIGHT	2	-1.530	-1.496	.308	-1.510	.208
LOW	2	.362	.392	.202	.358	-.049
MODERATE	1	.379	.188	-1.408	.374	-.051
HIGH	2	.978	1.010	.194	.965	-.133

Variable Principal Normalization.

Glancing at the quantifications in the joint plot of the category points, you can see that some of the categories of some variables were not clearly separated by the categorical principal components analysis as cleanly as would have been expected if the level had been truly ordinal. Variables *Intensity of interaction* and *Frequency of interaction*, for example, have equal or almost equal quantifications for their two middle categories. This kind of result might suggest trying alternative categorical principal components analyses, perhaps with some categories collapsed, or perhaps with a different level of analysis, such as (multiple) nominal.

Figure 9-10

Joint plot category points



The joint plot of category points resembles the plot for the component loadings, but it also shows where the endpoints are located that correspond to the lowest quantifications (for example, *slight* for *Intensity of interaction* and *none* for *Feeling of belonging*). The two variables measuring interaction, *Intensity of interaction* and *Frequency of interaction*, appear very close together and account for much of the variance in dimension 1. *Formality of Relationship* also appears close to *Physical Proximity*.

By focusing on the category points, you can see the relationships even more clearly. Not only are *Intensity of interaction* and *Frequency of interaction* close, but the directions of their scales are similar; that is, slight intensity is close to slight frequency, and frequent interaction is near high intensity of interaction. You also see that close physical proximity seems to go hand-in-hand with an informal type of relationship, and physical distance is related to no relationship.

Object Scores

You can also request a listing and plot of object scores. The plot of the object scores can be useful for detecting outliers, detecting typical groups of objects, or revealing some special patterns.

The object scores table shows the listing of object scores labeled by social group for the Guttman-Bell data. By examining the values for the object points, you can identify specific objects in the plot.

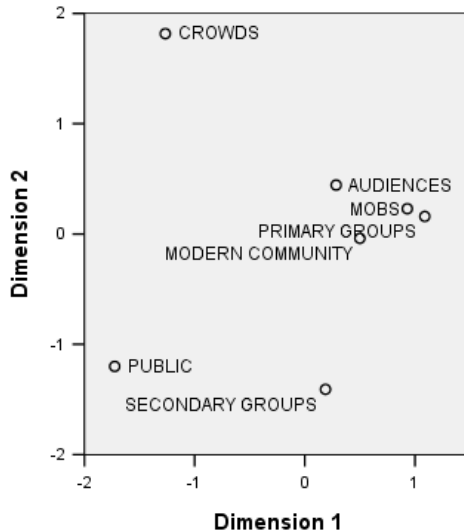
Figure 9-11
Object scores

cluster	Dimension	
	1	2
CROWDS	-1.266	1.816
AUDIENCES	.284	.444
PUBLIC	-1.726	-1.201
MOBS	.931	.229
PRIMARY GROUPS	1.089	.159
SECONDARY GROUPS	.188	-1.408
MODERN COMMUNITY	.500	-0.39

Variable Principal Normalization.

The first dimension appears to separate *CROWDS* and *PUBLIC*, which have relatively large negative scores, from *MOBS* and *PRIMARY GROUPS*, which have relatively large positive scores. The second dimension has three clumps: *PUBLIC* and *SECONDARY GROUPS* with large negative values, *CROWDS* with large positive values, and the other social groups in between. This is easier to see by inspecting the plot of the object scores.

Figure 9-12
Object scores plot



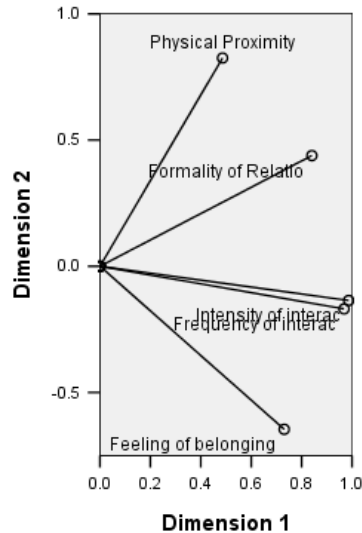
In the plot, you see *PUBLIC* and *SECONDARY GROUPS* at the bottom, *CROWDS* at the top, and the other social groups in the middle. Examining patterns among individual objects depends on the additional information available for the units of analysis. In this case, you know the classification of the objects. In other cases, you can use supplementary variables to label the objects. You can also see that the categorical principal components analysis does not separate *MOBS* from *PRIMARY GROUPS*. Although most people usually don't think of their families as mobs, on the variables used, these two groups received the same score on four of the five variables! Obviously, you might want to explore possible shortcomings of the variables and categories used. For example, high intensity of interaction and informal relationships probably mean different things to these two groups. Alternatively, you might consider a higher dimensional solution.

Component Loadings

This figure shows the plot of component loadings. The vectors (lines) are relatively long, indicating again that the first two dimensions account for most of the variance of all of the quantified variables. On the first dimension, all variables have high (positive)

component loadings. The second dimension is correlated mainly with the quantified variables *Feeling of belonging* and *Physical Proximity*, in opposite directions. This means that objects with a large negative score in dimension 2 will have a high score in feeling of belonging and a low score in physical proximity. The second dimension, therefore, reveals a contrast between these two variables while having little relation with the quantified variables *Intensity of interaction* and *Frequency of interaction*.

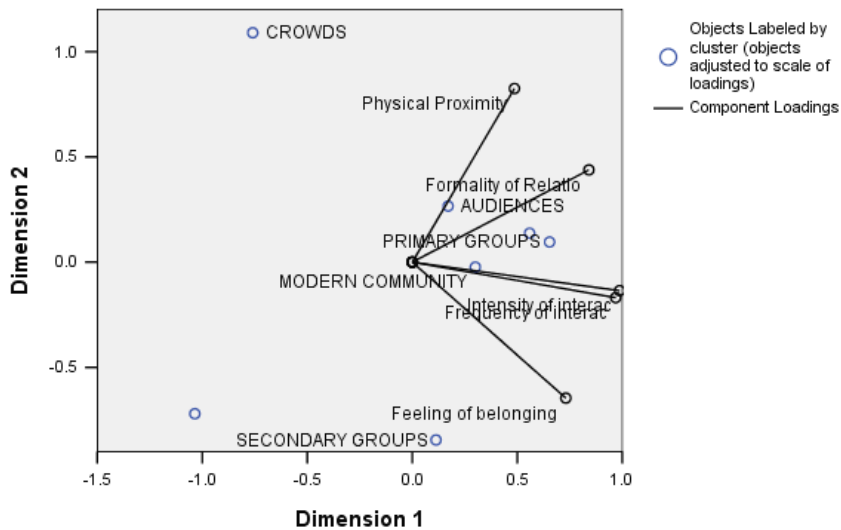
Figure 9-13
Component loadings



To examine the relation between the objects and the variables, look at the biplot of objects and component loadings. The vector of a variable points into the direction of the highest category of the variable. For example, for *Physical Proximity* and *Feeling of belonging* the highest categories are *close* and *high*, respectively. Therefore, *CROWDS* are characterized by close physical proximity and no feeling of belonging,

and *SECONDARY GROUPS*, by distant physical proximity and a high feeling of belonging.

Figure 9-14
Biplot



Additional Dimensions

Increasing the number of dimensions will increase the amount of variation accounted for and may reveal differences concealed in lower dimensional solutions. As noted previously, in two dimensions *MOBS* and *PRIMARY GROUPS* cannot be separated. However, increasing the dimensionality may allow the two groups to be differentiated.

Running the Analysis

- ▶ To obtain a three-dimensional solution, recall the Categorical Principal Components dialog box.
- ▶ Type 3 as the number of dimensions in the solution.

- Click OK in the Categorical Principal Components dialog box.

Model Summary

Figure 9-15
Model summary

Dimension	Cronbach's Alpha	Variance Accounted For	
		Total (Eigenvalue)	% of Variance
1	.885	3.424	68.480
2	-.232	.844	16.871
3	-.459	.732	14.649
Total	1.000 ^a	5.000	99.999

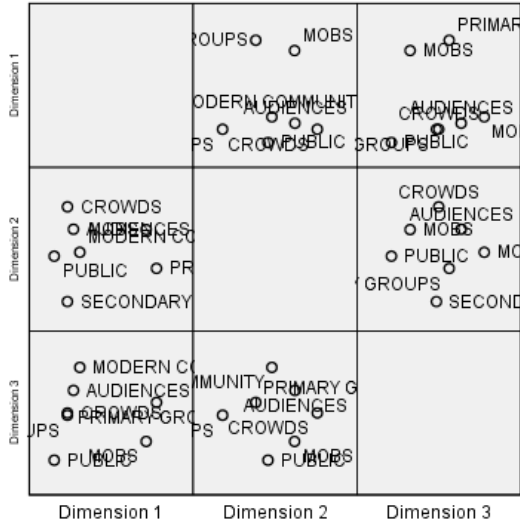
a. Total Cronbach's Alpha is based on the total Eigenvalue.

A three-dimensional solution has eigenvalues of 3.424, 0.844, and 0.732, accounting for nearly all of the variance.

Object Scores

The object scores for the three-dimensional solution are plotted in a scatterplot matrix. In a scatterplot matrix, every dimension is plotted against every other dimension in a series of two-dimensional scatterplots. Note that the first two eigenvalues in three dimensions are not equal to the eigenvalues in the two-dimensional solution; in other words, the solutions are not nested. Because the eigenvalues in dimensions 2 and 3 are now smaller than 1 (giving a Cronbach's alpha that is negative), you should prefer the two-dimensional solution. The three-dimensional solution is included for purposes of illustration.

Figure 9-16
 Three-dimensional object scores scatterplot matrix



The top row of plots reveals that the first dimension separates *PRIMARY GROUPS* and *MOBS* from the other groups. Notice that the order of the objects along the vertical axis does not change in any of the plots in the top row; each of these plots employs dimension 1 as the y axis.

The middle row of plots allows for interpretation of dimension 2. The second dimension has changed slightly from the two-dimensional solution. Previously, the second dimension had three distinct clumps, but now the objects are more spread out along the axis.

The third dimension helps to separate *MOBS* from *PRIMARY GROUPS*, which did not occur in the two-dimensional solution.

Look more closely at the dimension 2 versus dimension 3 and dimension 1 versus dimension 2 plots. On the plane defined by dimensions 2 and 3, the objects form a rough rectangle, with *CROWDS*, *MODERN COMMUNITY*, *SECONDARY GROUPS*, and *PUBLIC* at the vertices. On this plane, *MOBS* and *PRIMARY GROUPS* appear to be convex combinations of *PUBLIC-CROWDS* and *SECONDARY GROUPS-MODERN COMMUNITY*, respectively. However, as previously mentioned, they are separated from the other groups along dimension 1. *AUDIENCES* is not separated from the other groups along dimension 1 and appears to be a combination of *CROWDS* and *MODERN COMMUNITY*.

Component Loadings

Figure 9-17
Three-dimensional component loadings

	Dimension		
	1	2	3
Intensity of interaction	.980	-.005	-.201
Frequency of interaction	.521	-.643	.561
Feeling of belonging	.980	-.002	-.197
Physical Proximity	.519	.656	.549
Formality of Relationship	.981	.004	-.193

Knowing how the objects are separated does not reveal which variables correspond to which dimensions. This is accomplished using the component loadings. The first dimension corresponds primarily to *Feeling of belonging*, *Intensity of interaction*, and *Formality of Relationship*; the second dimension separates *Frequency of interaction* and *Physical Proximity*; and the third dimension separates these from the others.

Example: Symptomatology of Eating Disorders

Eating disorders are debilitating illnesses associated with disturbances in eating behavior, severe body image distortion, and an obsession with weight that affects the mind and body simultaneously. Millions of people are affected each year, with adolescents particularly at risk. Treatments are available and most are helpful when the condition is identified early.

A health professional can attempt to diagnose an eating disorder through a psychological and medical evaluation. However, it can be difficult to assign a patient to one of several different classes of eating disorders because there is no standardized symptomatology of anorectic/bulimic behavior. Are there symptoms that clearly differentiate patients into the four groups? Which symptoms do they have in common?

In order to try to answer these questions, researchers (Van der Ham et al., 1997) made a study of 55 adolescents with known eating disorders, as shown in the following table.

Table 9-2
Patient diagnoses

Diagnosis	Number of Patients
Anorexia nervosa	25
Anorexia with bulimia nervosa	9
Bulimia nervosa after anorexia	14
Atypical eating disorder	7
Total	55

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 the 16 symptoms outlined in the following table. 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. The data can be found in *anorectic.sav*, located in the *\tutorial\sample_files* subdirectory of the directory in which you installed SPSS.

Table 9-3
Modified Morgan-Russell subscales measuring well-being

Variable name	Variable label	Lower end (score1)	Upper end (score 3 or 4)
<i>weight</i>	Body weight	Outside normal range	Normal
<i>mens</i>	Menstruation	Amenorrhea	Regular periods
<i>fast</i>	Restriction of food intake (fasting)	Less than 1200 calories	Normal/regular meals
<i>binge</i>	Binge eating	Greater than once a week	No bingeing
<i>vomit</i>	Vomiting	Greater than once a week	No vomiting
<i>purge</i>	Purging	Greater than once a week	No purging
<i>hyper</i>	Hyperactivity	Not able to be at rest	No hyperactivity
<i>fami</i>	Family relations	Poor	Good
<i>eman</i>	Emancipation from family	Very dependent	Adequate
<i>frie</i>	Friends	No good friends	Two or more good friends
<i>school</i>	School/employment record	Stopped school/work	Moderate to good record
<i>satt</i>	Sexual attitude	Inadequate	Adequate

Variable name	Variable label	Lower end (score1)	Upper end (score 3 or 4)
<i>sbeh</i>	Sexual behavior	Inadequate	Can enjoy sex
<i>mood</i>	Mental state (mood)	Very depressed	Normal
<i>preo</i>	Preoccupation with food and weight	Complete	No preoccupation
<i>body</i>	Body perception	Disturbed	Normal

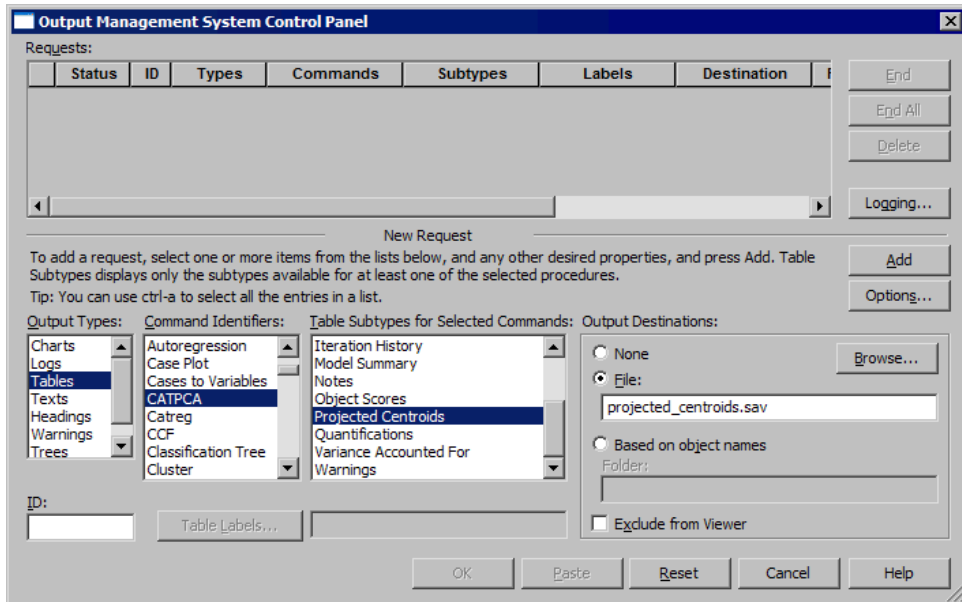
Principal components analysis is ideal for this situation, since the purpose of the study is to ascertain the relationships between symptoms and the different classes of eating disorders. Moreover, categorical principal components analysis is likely to be more useful than classical principal components analysis because the symptoms are scored on an ordinal scale.

Running the Analysis

In order to properly examine the structure of the course of illness for each diagnosis, you will want to make the results of the projected centroids table available as data for scatterplots. You can accomplish this using the Output Management System.

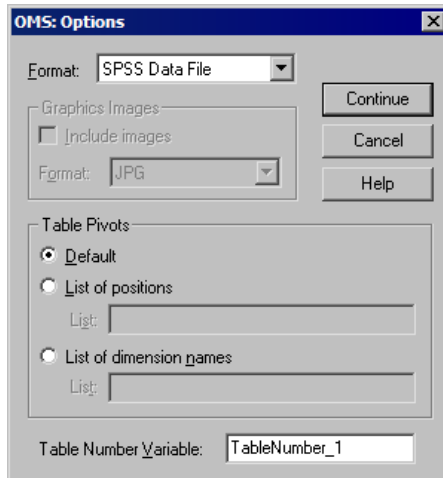
- ▶ To begin an OMS request, from the menus choose:
 - Utilities
 - OMS Control Panel...

Figure 9-18
Output Management System Control Panel



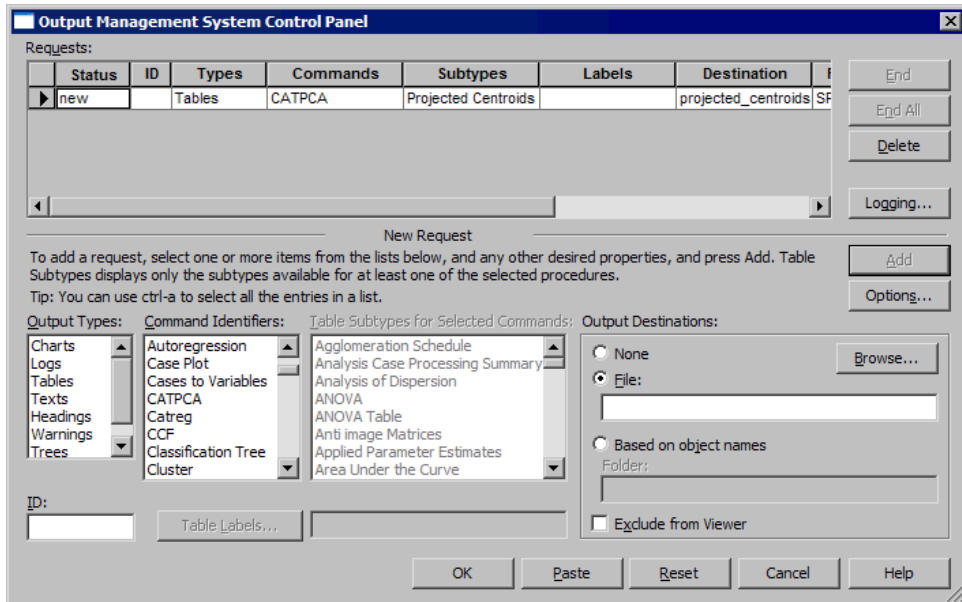
- ▶ Select Tables as the output type.
- ▶ Select CATPCA as the command.
- ▶ Select Projected Centroids as the table type.
- ▶ Select File in the Output Destinations group and type projected_centroids.sav as the file name.
- ▶ Click Options.

Figure 9-19
Options dialog box



- ▶ Select SPSS Data File as the output format.
- ▶ Type TableNumber_1 as the table number variable.
- ▶ Click Continue.

Figure 9-20
Output Management System Control Panel



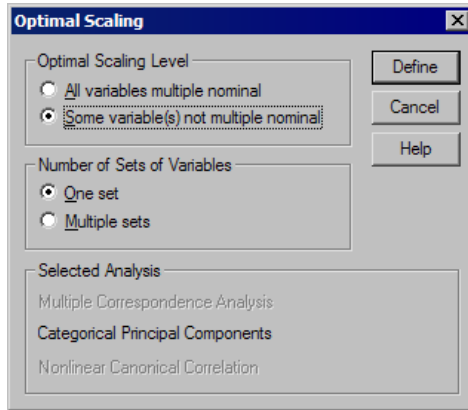
- ▶ Click Add.
- ▶ Click OK, and then click OK to confirm the OMS session.

The Output Management System is now set to write the results of the projected centroids table to the file *projected_centroids.sav*.

- ▶ To produce categorical principal components output for this data set, from the menus choose:

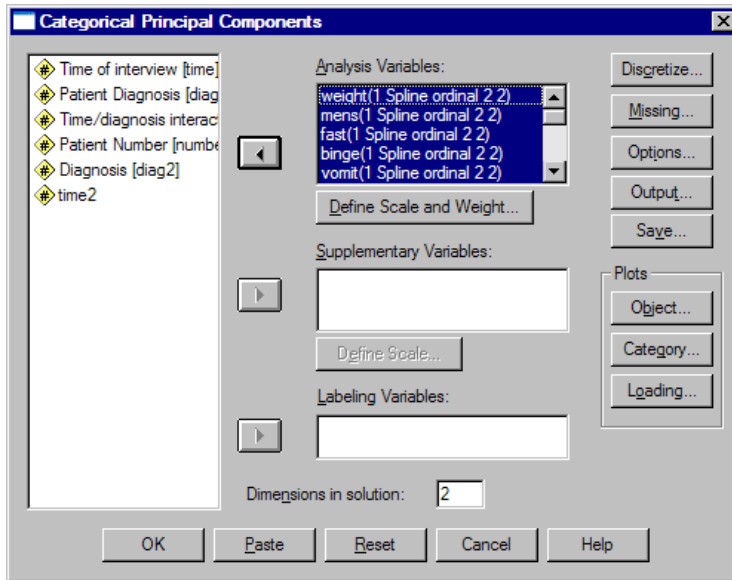
Analyze
Data Reduction
Optimal Scaling...

Figure 9-21
Optimal Scaling dialog box



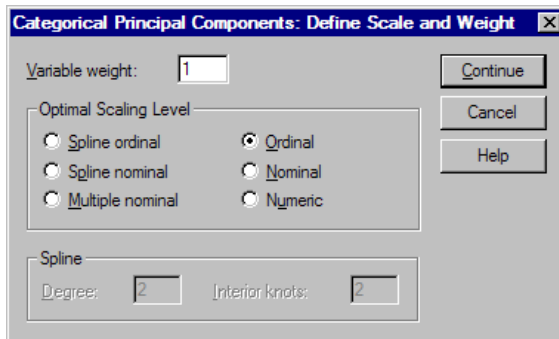
- ▶ Select Some variable(s) not multiple nominal in the Optimal Scaling Level group.
- ▶ Click Define.

Figure 9-22
Categorical Principal Components dialog box



- ▶ Select *Body weight* through *Body perception* as analysis variables.
- ▶ Click Define Scale and Weight.

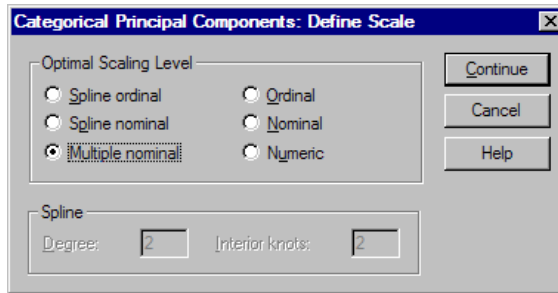
Figure 9-23
Define Scale and Weight dialog box



- ▶ Select Ordinal as the optimal scaling level.

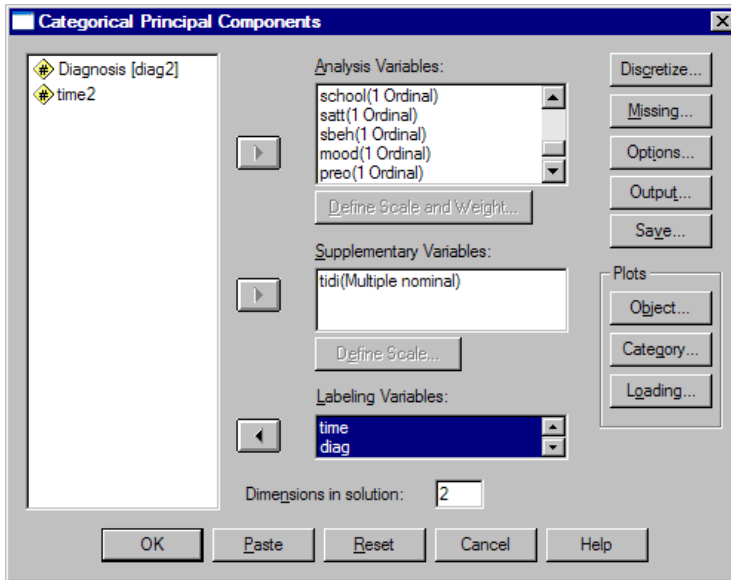
- ▶ Click Continue.
- ▶ Select *Time/diagnosis interaction* as a supplementary variable and click Define Scale in the Categorical Principal Components dialog box.

Figure 9-24
Define Scale dialog box



- ▶ Select Multiple nominal as the optimal scaling level.
- ▶ Click Continue.

Figure 9-25
Categorical Principal Components dialog box



- ▶ Select *Time of interview* through *Patient number* as labeling variables.
- ▶ Click Options.

Figure 9-26
Options dialog box

Categorical Principal Components: Options

Supplementary Objects

Range of cases

First:

Last:

Single case:

Add

Change

Remove

Normalization Method

Variable Principal

Custom value:

Criteria

Convergence:

Maximum iterations:

Label Plots By

Variable labels or value labels

Limit for label length:

Variable names or values

Plot Dimensions

Display all dimensions in the solution

Restrict the number of dimensions

Lowest dimension:

Highest dimension:

Configuration

None

File...

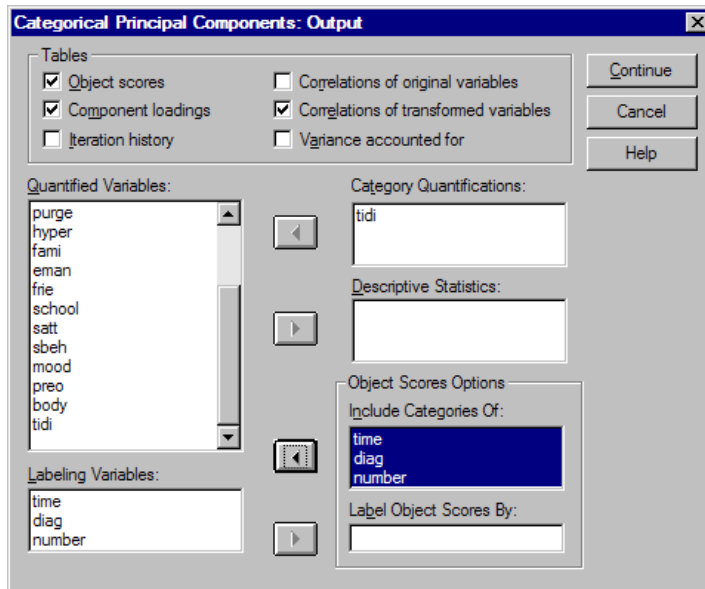
Continue

Cancel

Help

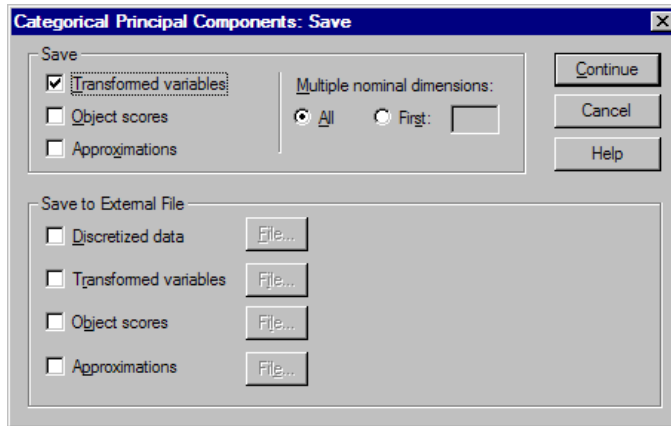
- ▶ Choose to label plots by Variable names or values.
- ▶ Click Continue.
- ▶ Click Output in the Categorical Principal Components dialog box.

Figure 9-27
Output dialog box



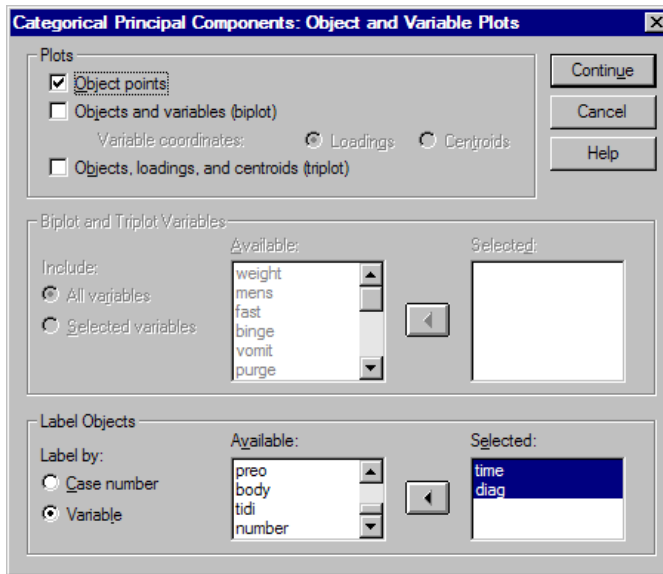
- ▶ Select Object scores in the Tables group.
- ▶ Request category quantifications for *tidi*.
- ▶ Choose to include categories of *time*, *diag*, and *number*.
- ▶ Click Continue.
- ▶ Click Save in the Categorical Principal Components dialog box.

Figure 9-28
Save dialog box



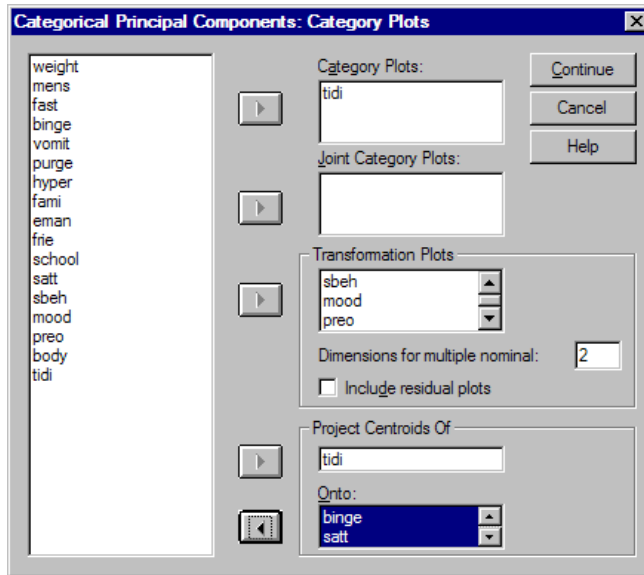
- ▶ Choose to save Transformed variables to the working file.
- ▶ Click Continue.
- ▶ Click Object in the Categorical Principal Components dialog box.

Figure 9-29
Object and Variable Plots dialog box



- ▶ Choose to label objects by Variable.
- ▶ Select *time* and *diag* as the variables to label objects by.
- ▶ Click Continue.
- ▶ Click Category in the Categorical Principal Components dialog box.

Figure 9-30
Category Plots dialog box

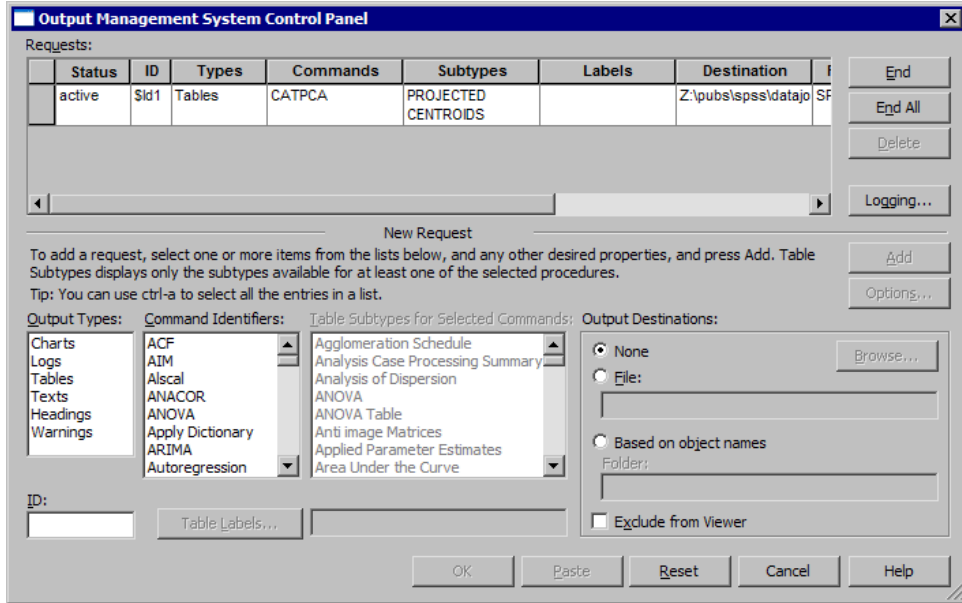


- ▶ Request category plots for *tidi*.
- ▶ Request transformation plots for *weight* through *body*.
- ▶ Choose to project centroids of *tidi* onto *binge*, *satt*, and *preo*.
- ▶ Click Continue.
- ▶ Click OK in the Categorical Principal Components dialog box.

The procedure results in scores for the subjects (with mean 0 and unit variance) and quantifications of the categories that maximize the mean squared correlation of the subject scores and the transformed variables. In the present analysis, the category quantifications were constrained to reflect the ordinal information.

Finally, to write the projected centroids table information to *projected_centroids.sav*, you need to end the OMS request. Recall the OMS Control Panel.

Figure 9-31
Output Management System Control Panel

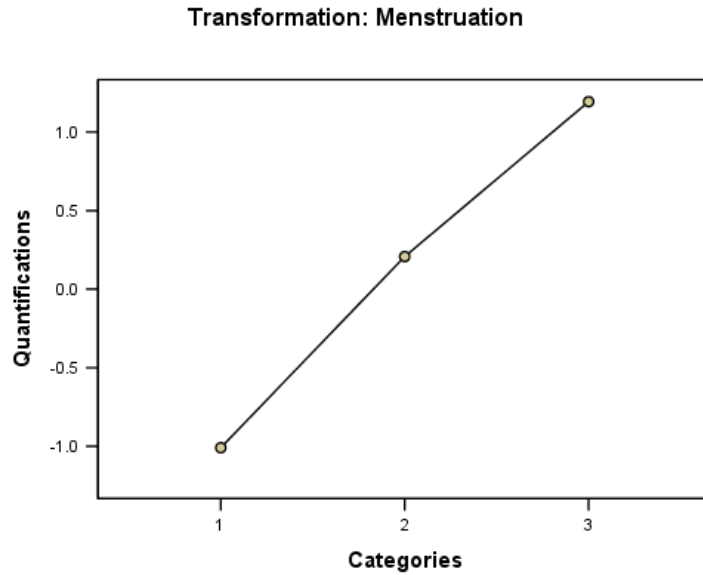


- ▶ Click End.
- ▶ Click OK, and then click OK to confirm.

Transformation Plots

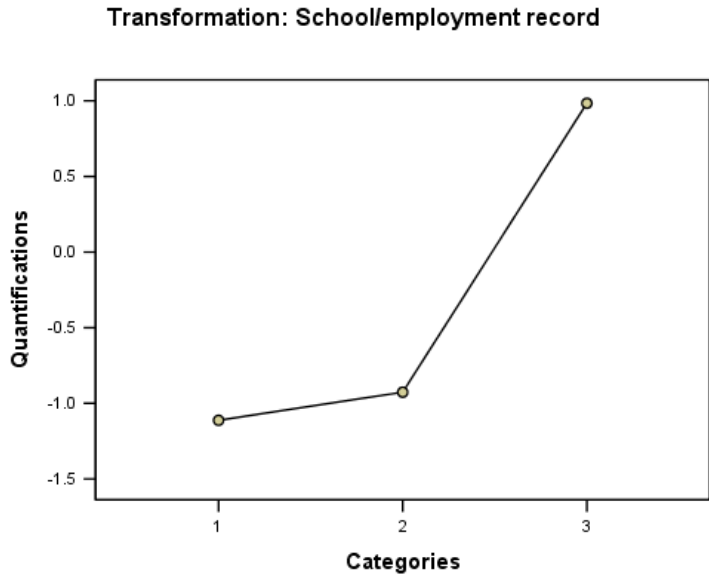
The transformation plots display the original category number on the horizontal axes; the vertical axes give the optimal quantifications.

Figure 9-32
Transformation plot for menstruation



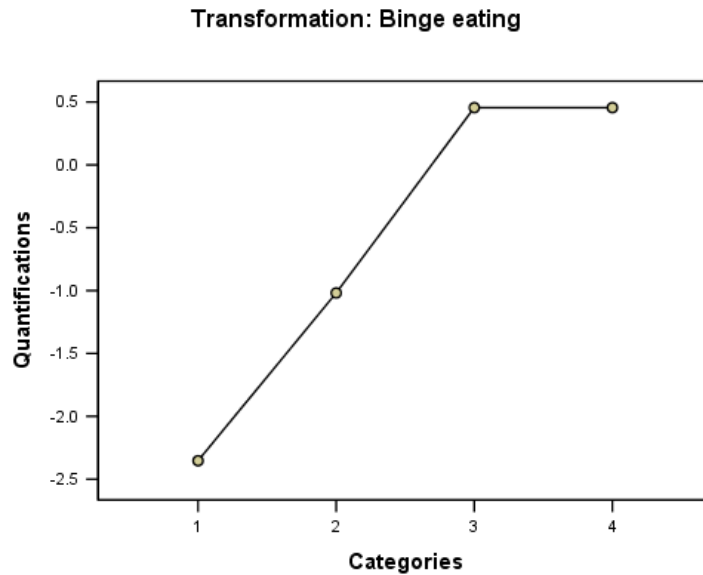
Some variables, like *Menstruation*, obtained nearly linear transformations, so in this analysis you may interpret them as numerical.

Figure 9-33
Transformation plot for *School/employment record*



The quantifications for other variables like *School/employment record* did not obtain linear transformations and should be interpreted at the ordinal scaling level. The difference between the second and third categories is much more important than that between the first and second categories.

Figure 9-34
Transformation plot for *Binge eating*



An interesting case arises in the quantifications for *Binge eating*. The transformation obtained is linear for categories 1 through 3, but the quantified values for categories 3 and 4 are equal. This result shows that scores of 3 and 4 do not differentiate between patients and suggests that you could use the numerical scaling level in a two-component solution by recoding 4's as 3's.

Model Summary

Figure 9-35
Model summary

Dimension	Cronbach's Alpha	Variance Accounted For	
		Total (Eigenvalue)	% of Variance
1	.874	5.550	34.690
2	.522	1.957	12.234
Total	.925 ^a	7.508	46.924

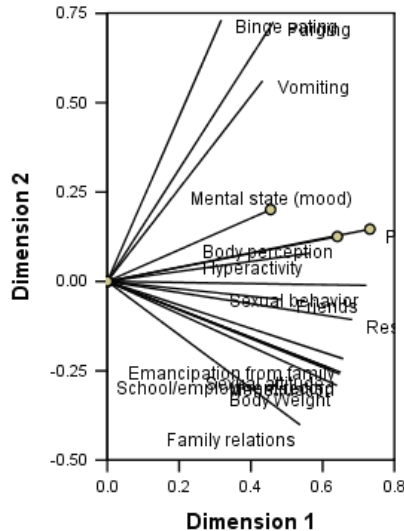
a. Total Cronbach's Alpha is based on the total Eigenvalue.

To see how well your model fits the data, look at the model summary. About 47% of the total variance is explained by the two-component model, 35% by the first dimension and 12% by the second. So, almost half of the variability on the individual objects level is explained by the two-component model.

Component Loadings

To begin to interpret the two dimensions of your solution, look at the component loadings. All variables have a positive component loading in the first dimension, which means that there is a common factor that correlates positively with all of the variables.

Figure 9-36
Component loadings plot



The second dimension separates the variables. The variables *Binge eating*, *Vomiting*, and *Purging* form a bundle having large positive loadings in the second dimension. These symptoms are typically considered to be representative of bulimic behavior.

The variables *Emancipation from family*, *School/employment record*, *Sexual attitude*, *Body weight*, and *Menstruation* form another bundle, and you can include *Restriction of food intake (fasting)* and *Family relations* in this bundle, because their vectors are close to the main cluster, and these variables are considered to be anorectic symptoms (fasting, weight, menstruation) or are psychosocial in nature (emancipation, school/work record, sexual attitude, family relations). The vectors of this bundle are orthogonal (perpendicular) to the vectors of binge, vomit, and purge, which means that this set of variables is uncorrelated with the set of bulimic variables.

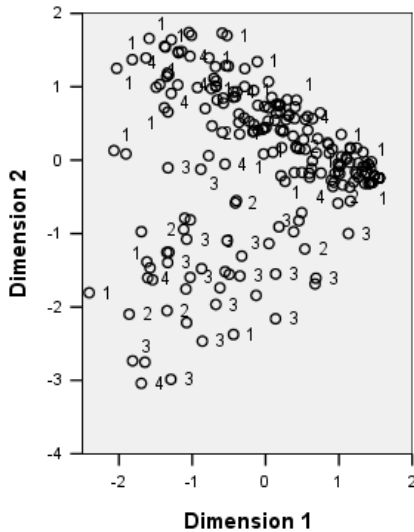
The variables *Friends*, *Mental state (mood)*, and *Hyperactivity* do not appear to fit very well into the solution. You can see this in the plot by observing the lengths of each vector. The length of a given variable's vector corresponds to its fit, and these variables have the shortest vectors. Based on a two-component solution, you would probably drop these variables from a proposed symptomatology for eating disorders. They may, however, fit better in a higher dimensional solution.

The variables *Sexual behavior*, *Preoccupation with food and weight*, and *Body perception* form another theoretic group of symptoms, pertaining to how the patient experiences his or her body. While correlated with the two orthogonal bundles of variables, these variables have fairly long vectors and are strongly associated with the first dimension and therefore may provide some useful information about the “common” factor.

Object Scores

The following figure shows a plot of the object scores, in which the subjects are labeled with their diagnosis category.

Figure 9-37
Object scores plot labeled by diagnosis

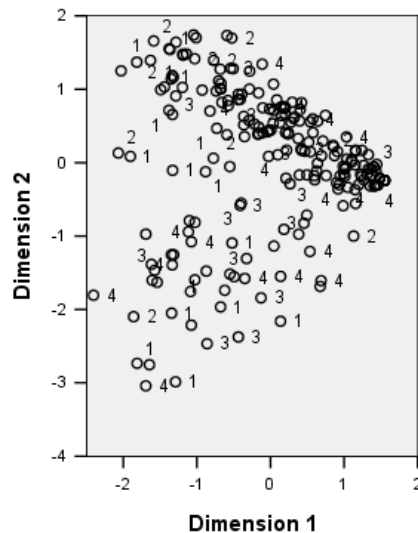


This plot does not help to interpret the first dimension, because patients are not separated by diagnosis along it. However, there is some information about the second dimension. Anorexia subjects (1) and patients with atypical eating disorder (4) form a group, located above subjects with some form of bulimia (2 and 3). Thus, the second dimension separates bulimic patients from others, as you have also seen in the

previous section (the variables in the bulimic bundle have large positive component loadings in the second dimension). This makes sense, given that the component loadings of the symptoms that are traditionally associated with bulimia have large values in the second dimension.

This figure shows a plot of the object scores, in which the subjects are labeled with their time of diagnosis.

Figure 9-38
Object scores labeled by time of interview



Labeling the object scores by time reveals that the first dimension has a relation to time because there seems to be a progression of times of diagnosis from the 1's mostly to the left and others to the right. Note that you can connect the time points in this plot by saving the object scores and creating a scatterplot using the dimension 1 scores on the x axis, the dimension 2 scores on the y axis, and setting the markers using the patient numbers.

Comparing the object scores plot labeled by time with the one labeled by diagnosis can give you some insight into unusual objects. For example, in the plot labeled by time, there is a patient whose diagnosis at time 4 lies to the left of all other points in the plot. This is unusual because the general trend of the points is for the later times to lie further to the right. Interestingly, this point that seems out of place in time also has an unusual diagnosis, in that the patient is an anorectic whose scores

place the patient in the bulimic cluster. By looking in the table of object scores, you find that this is patient 43, diagnosed with anorexia nervosa, whose object scores are shown in the following table.

Table 9-4
Object scores for patient 43

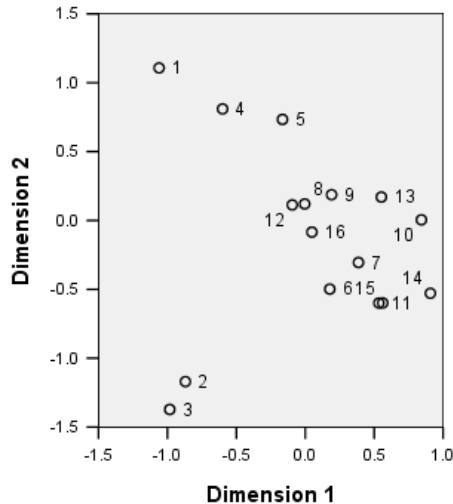
Time	Dimension 1	Dimension 2
1	-2.031	1.250
2	-2.067	0.131
3	-1.575	-1.467
4	-2.405	-1.807

The patient's scores at time 1 are prototypical for anorectics, with the large negative score in dimension 1 corresponding to poor body image and the positive score in dimension 2 corresponding to anorectic symptoms or poor psychosocial behavior. However, unlike the majority of patients, there is little or no progress in dimension 1. In dimension 2, there is apparently some progress toward "normal" (around 0, between anorectic and bulimic behavior), but then the patient shifts to exhibit bulimic symptoms.

Examining the Structure of the Course of Illness

To find out more about how the two dimensions were related to the four diagnosis categories and the four time points, a supplementary variable *Time/diagnosis interaction* was created by a cross-classification of the four categories of *Patient diagnosis* and the four categories of *Time of interview*. Thus, *Time/diagnosis interaction* has 16 categories, where the first category indicates the anorexia nervosa patients at their first visit. The fifth category indicates the anorexia nervosa patients at time point 2, and so on, with the sixteenth category indicating the atypical eating disorder patients at time point 4. The use of the supplementary variable *Time/diagnosis interaction* allows for the study of the courses of illness for the different groups over time. The variable was given a multiple nominal scaling level, and the category points are displayed in the following figure.

Figure 9-39
Category points for time/diagnosis interaction

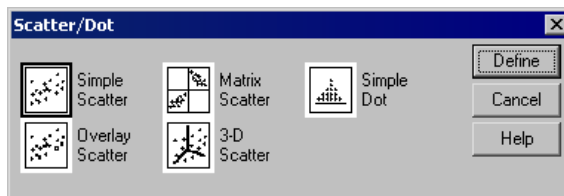


Some of the structure is apparent from this plot: the diagnosis categories at time point 1 clearly separate anorexia nervosa and atypical eating disorder from anorexia nervosa with bulimia nervosa and bulimia nervosa after anorexia nervosa in the second dimension. After that, it's a little more difficult to see the patterns.

However, you can make the patterns more easily visible by creating a scatterplot based on the quantifications. To do this, from the menu choose:

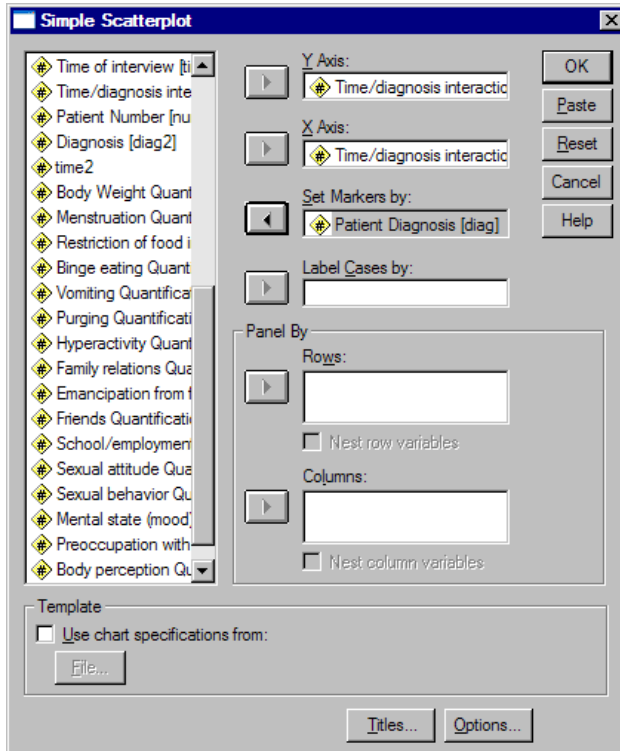
Graphs
Scatter/Dot...

Figure 9-40
Scatter/Dot dialog box



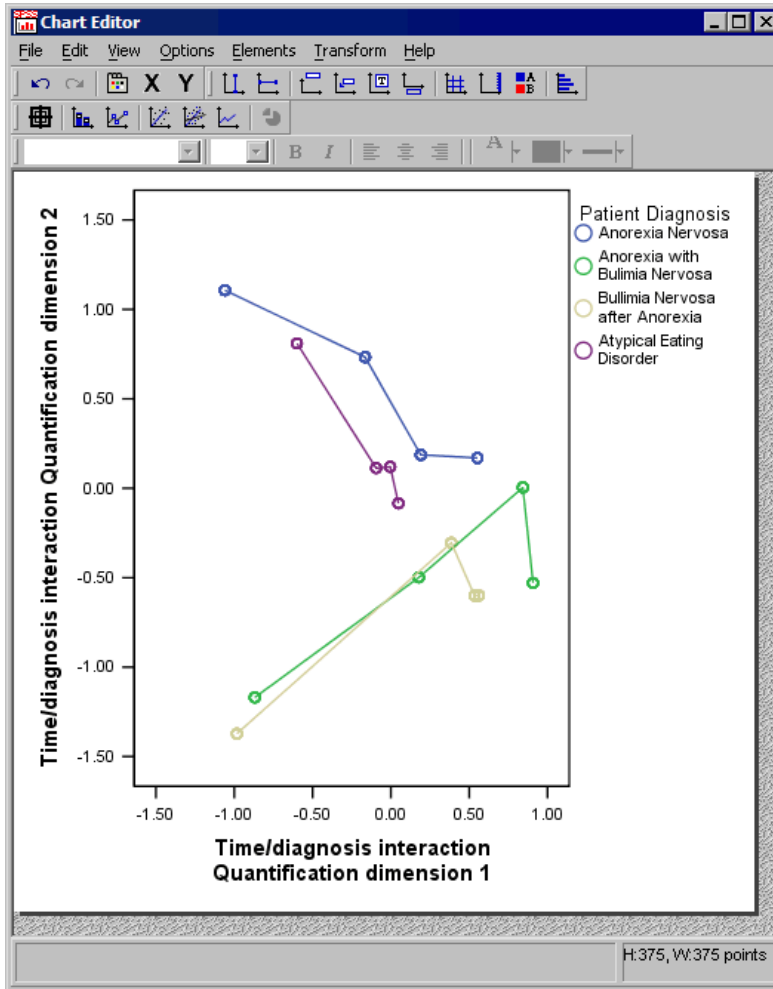
► Click Define.

Figure 9-41
Simple Scatterplot dialog box



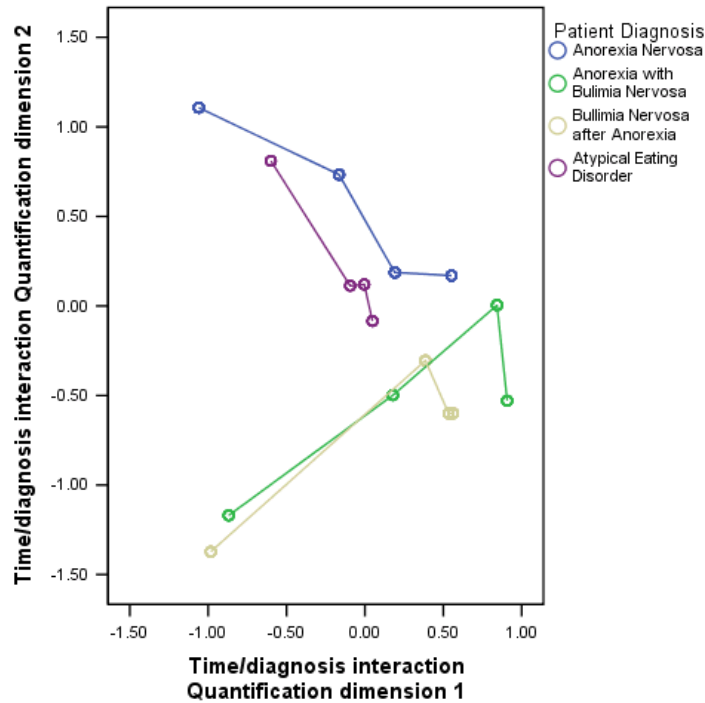
- ▶ Select *Time/diagnosis interaction Quantification dimension 2* as the y-axis variable and *Time/diagnosis interaction Quantification dimension 1* as the x-axis variable.
- ▶ Choose to set markers by *Patient Diagnosis*.
- ▶ Click OK.

Figure 9-42
Structures of the courses of illness



- ▶ Then, to connect the points, double-click on the graph, and click the Add interpolation line tool in the Chart Editor.
- ▶ Close the Chart Editor.

Figure 9-43
Structures of the courses of illness



By connecting the category points for each diagnostic category across time, the patterns immediately suggest that the first dimension is related to time and the second, to diagnosis, as you previously determined from the object scores plots.

However, this plot further shows that, over time, the illnesses tend to become more alike. Moreover, for all groups, the progress is greatest between time points 1 and 2; the anorectic patients show some more progress from 2 to 3, but the other groups show little progress.

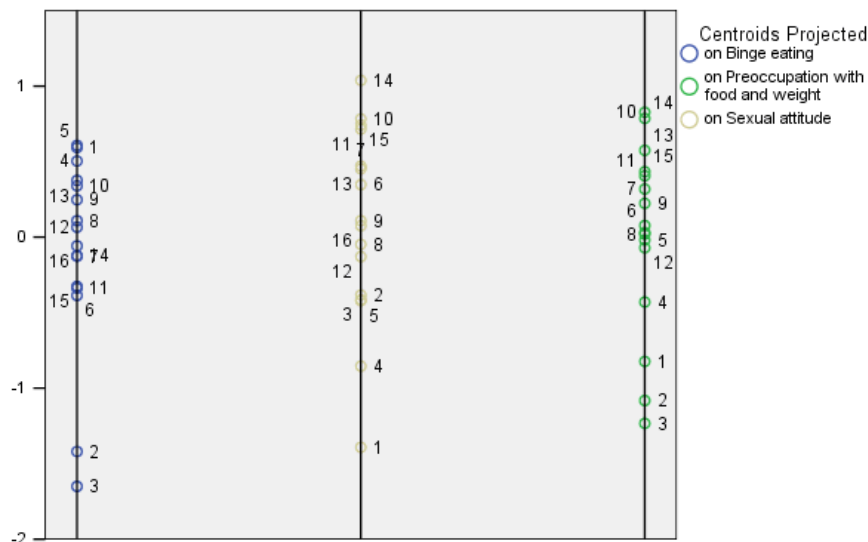
Differential Development for Selected Variables

One variable from each bundle of symptoms identified by the component loadings was selected as “representative” of the bundle. Binge eating was selected from the bulimic bundle; sexual attitude, from the anorectic/psychosocial bundle; and body preoccupation, from the third bundle.

In order to examine the possible differential courses of illness, the projections of *Time/diagnosis interaction on Binge eating, Sexual attitude, and Preoccupation with food and weight* were computed and plotted in the following figure.

Figure 9-44

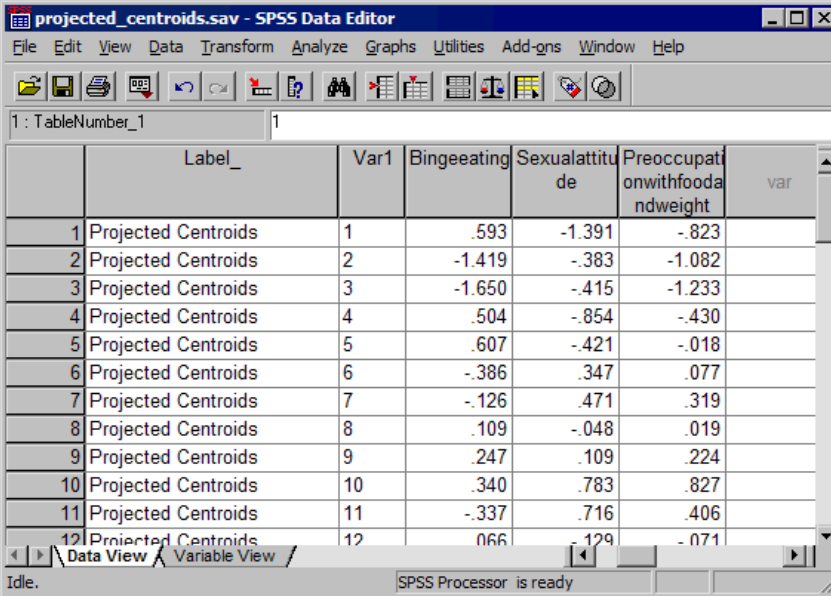
Projected centroids of Time/diagnosis interaction on Binge eating, Sexual attitude, and Preoccupation with food and weight



This plot shows that at the first time point, the symptom binge eating separates bulimic patients (2 and 3) from others (1 and 4); sexual attitude separates anorectic and atypical patients (1 and 4) from others (2 and 3); and body preoccupation does not really separate the patients. In many applications, this plot would be sufficient to describe the relationship between the symptoms and diagnosis, but because of the complication of multiple time points, the picture becomes muddled.

In order to view these projections over time, you need to be able to plot the contents of the projected centroids table. This is made possible by the OMS request that saved this information to *projected_centroids.sav*.

Figure 9-45
Projected_centroids.sav

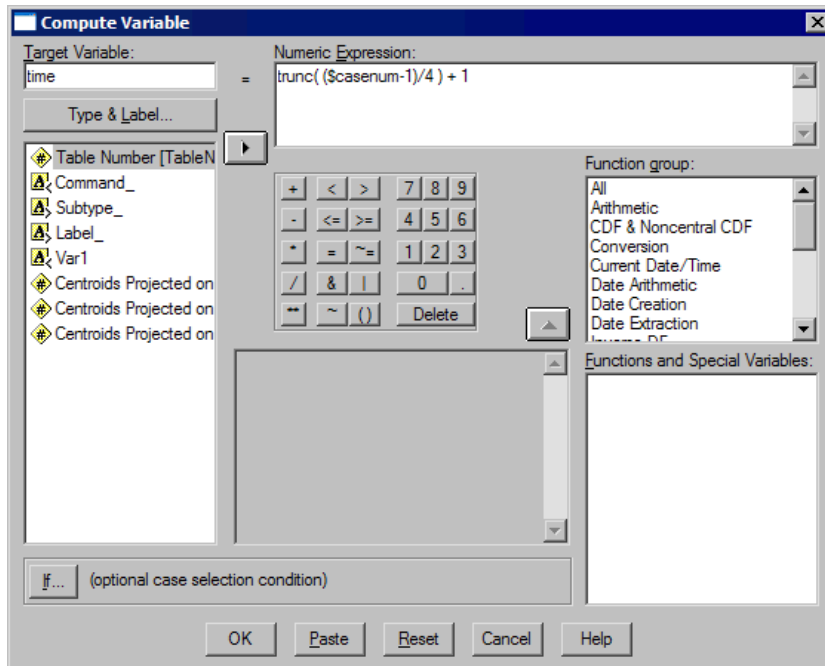


	Label_	Var1	Bingeeating	Sexualattitude	Preoccupationwithfoodandweight	var
1	Projected Centroids	1	.593	-1.391	- .823	
2	Projected Centroids	2	-1.419	-.383	-1.082	
3	Projected Centroids	3	-1.650	-.415	-1.233	
4	Projected Centroids	4	.504	-.854	-.430	
5	Projected Centroids	5	.607	-.421	-.018	
6	Projected Centroids	6	-.386	.347	.077	
7	Projected Centroids	7	-.126	.471	.319	
8	Projected Centroids	8	.109	-.048	.019	
9	Projected Centroids	9	.247	.109	.224	
10	Projected Centroids	10	.340	.783	.827	
11	Projected Centroids	11	-.337	.716	.406	
12	Projected Centroids	12	.066	-.129	-.071	

The variables *Bingeeating*, *Sexualattitude*, and *Preoccupationwithfoodandweight* contain the values of the centroids projected on each of the symptoms of interest. The case number (1 through 16) corresponds to the time/diagnosis interaction. You will need to compute new variables that separate out the Time and Diagnosis values.

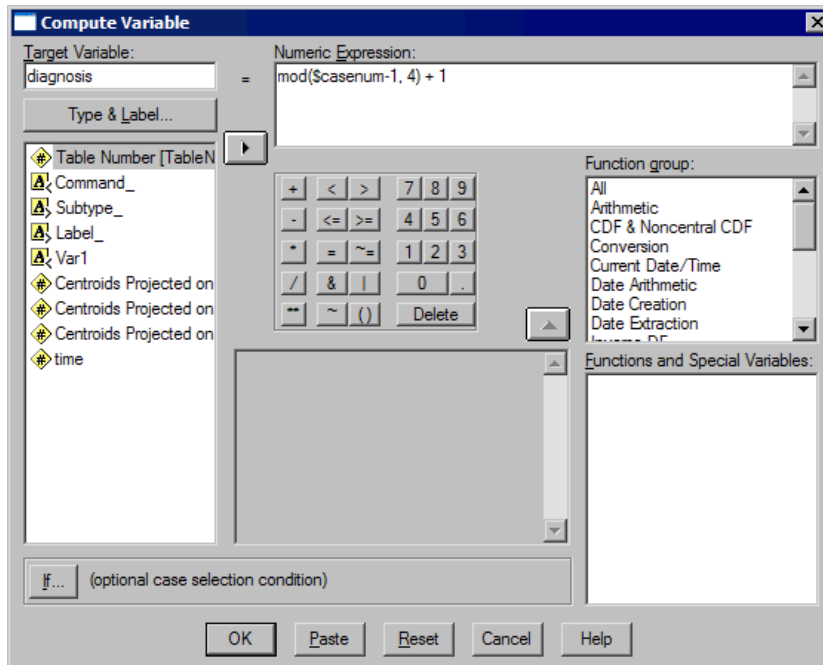
- ▶ From the menus choose:
 - Transform
 - Compute...

Figure 9-46
Compute Variable dialog box



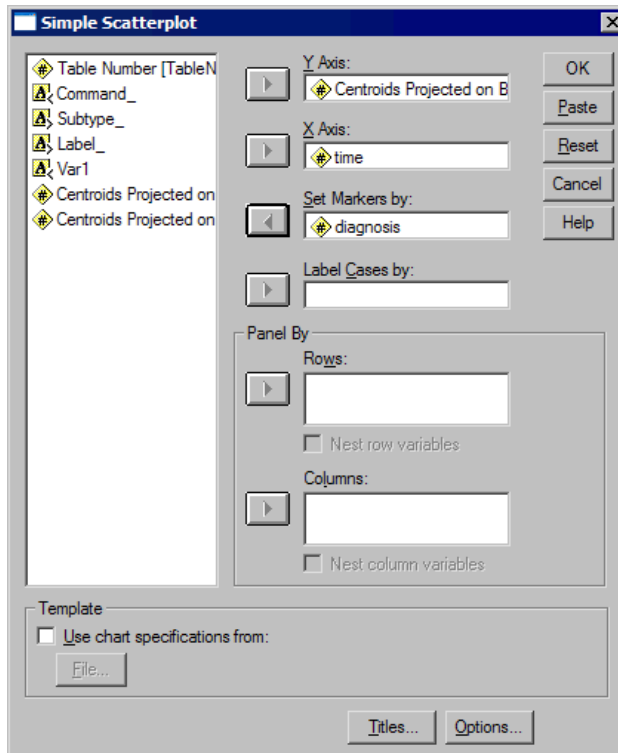
- ▶ Type *time* as the target variable.
- ▶ Type $\text{trunc}((\$casenum-1)/4) + 1$ as the numeric expression.
- ▶ Click OK.

Figure 9-47
Compute Variable dialog box



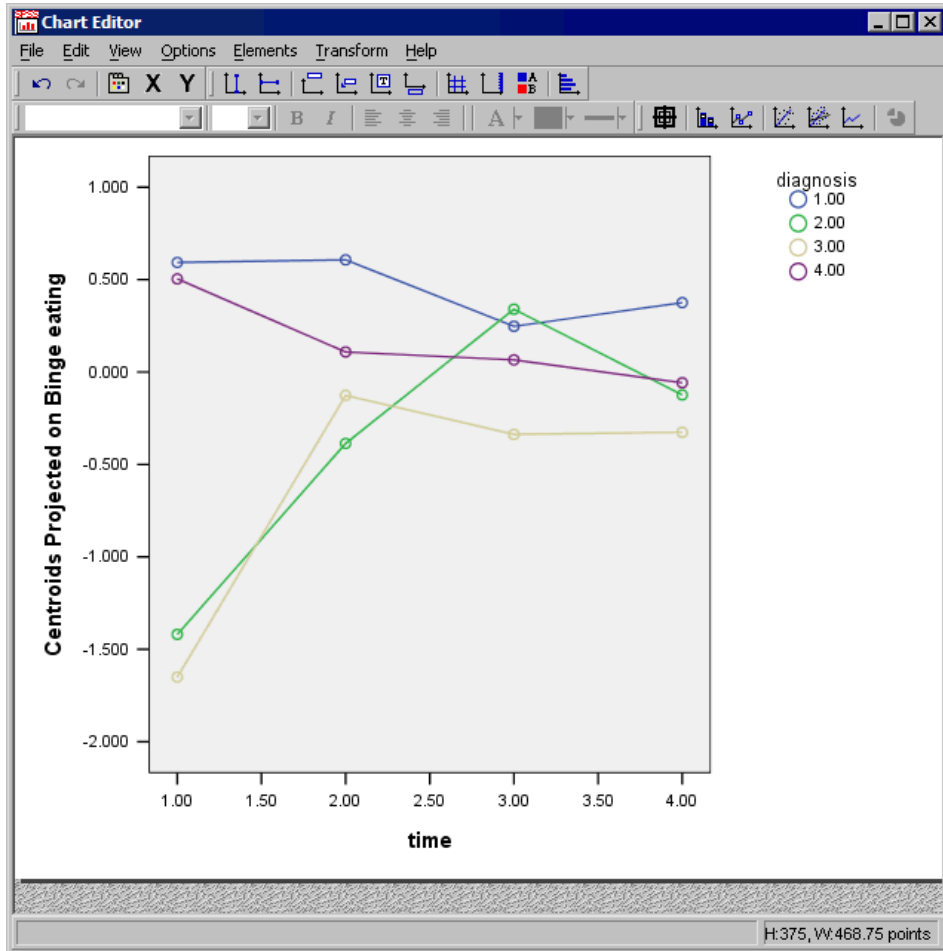
- ▶ Recall the Compute Variable dialog box.
- ▶ Type *diagnosis* as the target variable.
- ▶ Type $\text{mod}(\$casenum-1, 4) + 1$ as the numeric expression.
- ▶ Click OK.

Figure 9-48
Simple Scatterplot dialog box



- ▶ Finally, to view the projected centroids of time of diagnosis on binging over time, recall the Simple Scatterplot dialog box and click Reset to clear your previous selections.
- ▶ Select *Centroids Projected on Binge eating* as the y-axis variable and *time* as the x-axis variable.
- ▶ Choose to set markers by *diagnosis*.
- ▶ Click OK.

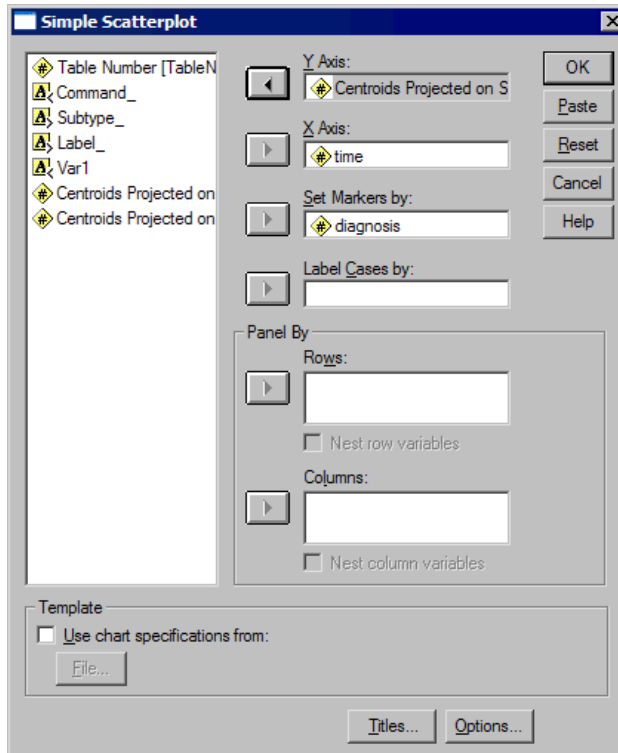
Figure 9-49
Projected centroids of Time of diagnosis on Binge eating over time



- ▶ Then, to connect the points, double-click on the graph, and click the Add interpolation line tool in the Chart Editor.
- ▶ Close the Chart Editor.

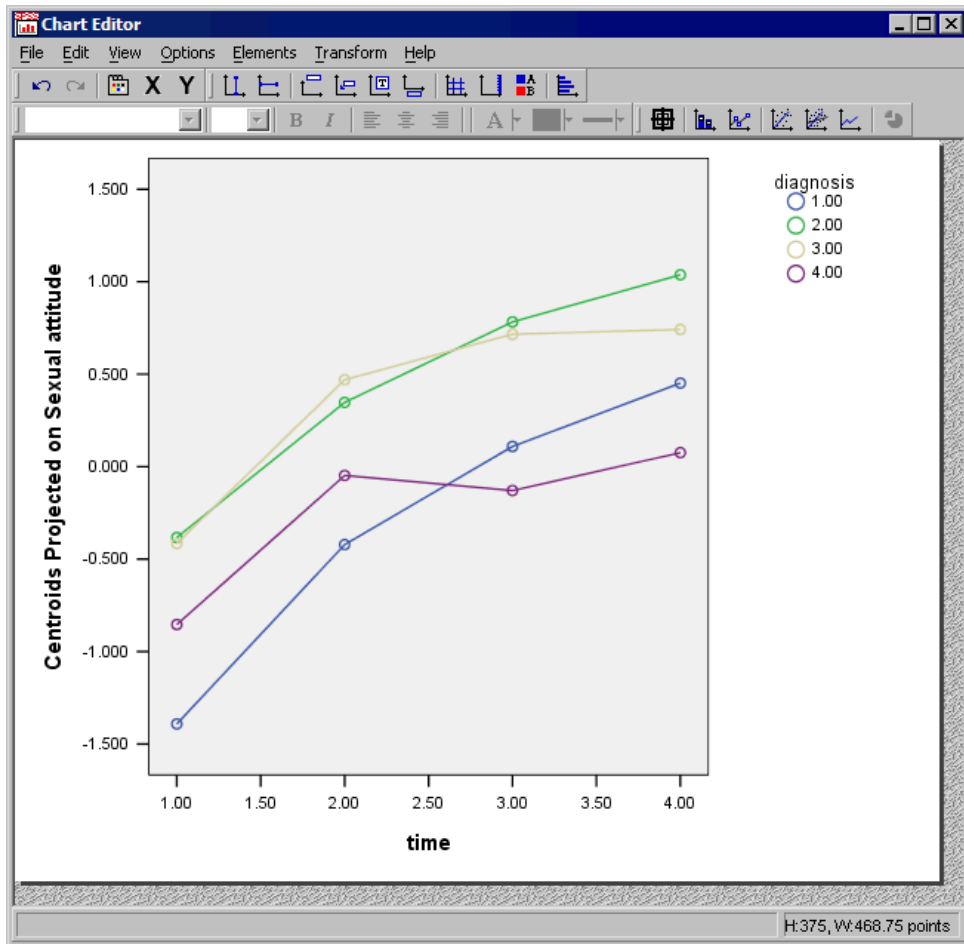
With respect to binge eating, it is clear that the anorectic groups have different starting values from the bulimic groups. This difference shrinks over time, as the anorectic groups hardly change, while the bulimic groups show progress.

Figure 9-50
Simple Scatterplot dialog box



- ▶ Recall the Simple Scatterplot dialog box.
- ▶ Deselect *Centroids Projected on Binge eating* as the y-axis variable and select *Centroids Projected on Sexual attitude* as the y-axis variable.
- ▶ Click OK.

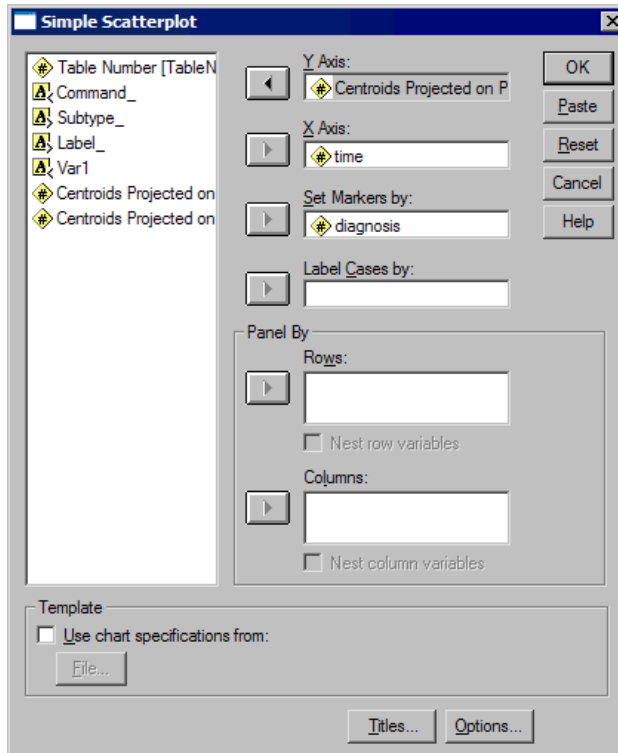
Figure 9-51
Projected centroids of Time of diagnosis on Sexual attitude over time



- ▶ Then, to connect the points, double-click on the graph, and click the Add interpolation line tool in the Chart Editor.
- ▶ Close the Chart Editor.

With respect to sexual attitude, the four trajectories are more or less parallel over time, and all groups show progress. The bulimic groups, however, have higher (better) scores than the anorectic group.

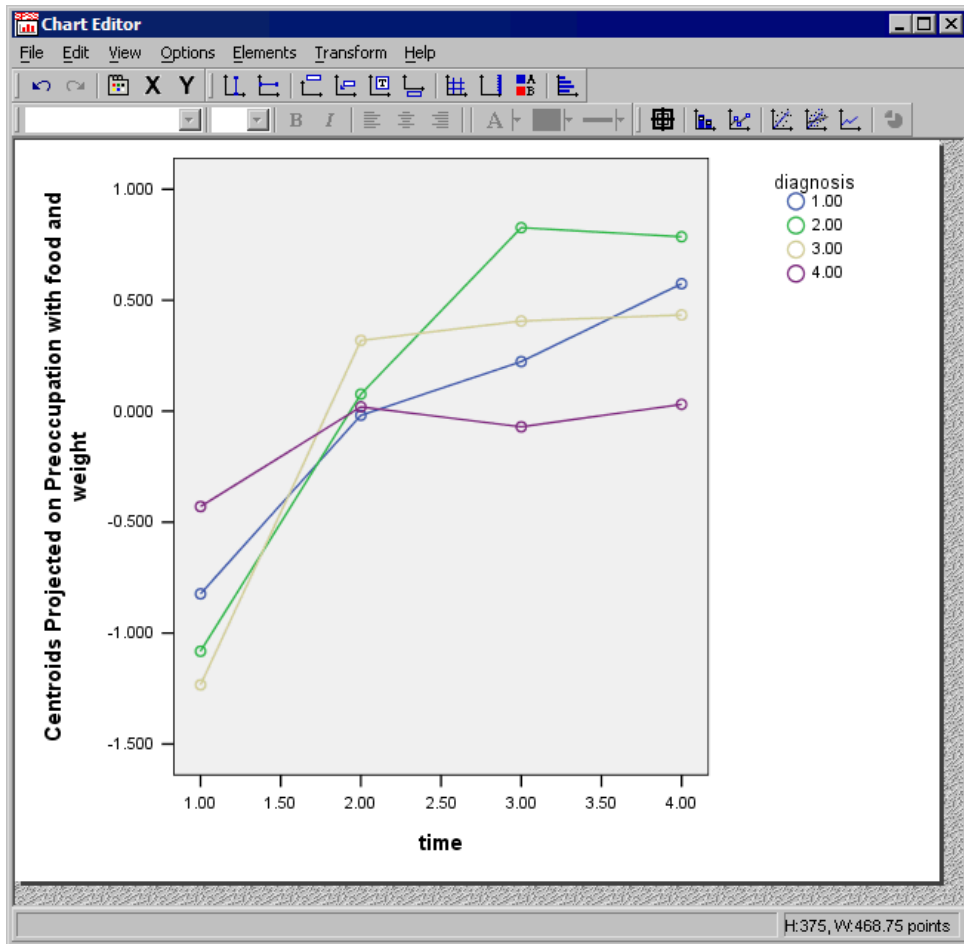
Figure 9-52
Simple Scatterplot dialog box



- ▶ Recall the Simple Scatterplot dialog box.
- ▶ Deselect *Centroids Projected on Sexual attitude* as the y-axis variable and select *Centroids Projected on Preoccupation with food and weight* as the y-axis variable.
- ▶ Click OK.

Figure 9-53

Projected centroids of Time of diagnosis on Body preoccupation over time



- ▶ Then, to connect the points, double-click on the graph, and click the Add interpolation line tool in the Chart Editor.
- ▶ Close the Chart Editor.

Body preoccupation is a variable that represents the core symptoms, which are shared by the four different groups. Apart from the atypical eating disorder patients, the anorectic group and the two bulimic groups have very similar levels both at the beginning and at the end.

Recommended Readings

See the following texts for more information on categorical principal components analysis:

Buja, A. 1990. Remarks on functional canonical variates, alternating least squares methods and ACE. *Annals of Statistics*, 18, 1032–1069.

De Haas, M., J. A. Algera, H. F. J. M. Van Tuijl, and J. J. Meulman. 2000. Macro and micro goal setting: In search of coherence. *Applied Psychology*, 49, 579–595.

De Leeuw, J. 1982. Nonlinear principal components analysis. In: *COMPSTAT Proceedings in Computational Statistics*, Vienna: PhysicaVerlag, 77–89.

Eckart, C., and G. Young. 1936. The approximation of one matrix by another one of lower rank. *Psychometrika*, 1, 211–218.

Gabriel, K. R. 1971. The biplot graphic display of matrices with application to principal components analysis. *Biometrika*, 58, 453–467.

Gifi, A. 1985. *PRINCALS. Research Report UG-85-02*. Leiden: Department of Data Theory, University of Leiden.

Gower, J. C., and J. J. Meulman. 1993. The treatment of categorical information in physical anthropology. *International Journal of Anthropology*, 8, 43–51.

Heiser, W. J., and J. J. Meulman. 1994. Homogeneity analysis: Exploring the distribution of variables and their nonlinear relationships. In: *Correspondence analysis in the social sciences: Recent developments and applications*, M. Greenacre, and J. Blasius, eds. New York: AcademicPress, 179–209.

Kruskal, J. B. 1978. Factor analysis and principal components analysis: Bilinear methods. In: *International encyclopedia of statistics*, W. H. Kruskal, and J. M. Tanur, eds. New York: The Free Press, 307–330.

- Kruskal, J. B., and R. N. Shepard. 1974. A nonmetric variety of linear factor analysis. *Psychometrika*, 39, 123–157.
- Meulman, J. 1993. Principal coordinates analysis with optimal transformations of the variables: Minimizing the sum of squares of the smallest eigenvalues. *British Journal of Mathematical and Statistical Psychology*, 46, 287–300.
- Meulman, J. J., and P. Verboon. 1993. Points of view analysis revisited: Fitting multidimensional structures to optimal distance components with cluster restrictions on the variables. *Psychometrika*, 58, 7–35.
- Meulman, J. J., A. J. van der Kooij, and A. Babinec. 2000. New features of categorical principal components analysis for complicated data sets, including data mining. In: *Classification, Automation and New Media*, W. Gaul, and G. Ritter, eds. Berlin: Springer-Verlag, 207–217.
- Meulman, J. J., A. J. van der Kooij, and W. J. Heiser. 2004. Principal components analysis with nonlinear optimal scaling transformations for ordinal and nominal data. In: *Handbook of Quantitative Methodology for the Social Sciences*, D. Kaplan, ed. Thousand Oaks, Calif.: SagePublications, Inc, 49–70.
- Theunissen, N. C. M., J. J. Meulman, A. L. Den Ouden, H. M. Koopman, G. H. Verrips, S. P. Verloove-Vanhorick, and J. M. Wit. 2003. Changes can be studied when the measurement instrument is different at different time points. *Health Services and Outcomes Research Methodology*, 4, 109–126.
- Tucker, L. R. 1960. Intra-individual and inter-individual multidimensionality. In: *Psychological Scaling: Theory & Applications*, H. Gulliksen, and S. Messick, eds. New York: John Wiley & Sons, 155–167.
- Vlek, C., and P. J. Stallen. 1981. Judging risks and benefits in the small and in the large. *Organizational Behavior and Human Performance*, 28, 235–271.
- Wagenaar, W. A. 1988. *Paradoxes of gambling behaviour*. London: Lawrence Erlbaum Associates, Inc.
- Young, F. W., Y. Takane, and J. De Leeuw. 1978. The principal components of mixed measurement level multivariate data: An alternating least squares method with optimal scaling features. *Psychometrika*, 43, 279–281.

Zeijl, E., Y. te Poel, M. du Bois-Reymond, J. Ravestloot, and J. J. Meulman. 2000. The role of parents and peers in the leisure activities of young adolescents. *Journal of Leisure Research*, 32, 281–302.

Nonlinear Canonical Correlation Analysis

The purpose of nonlinear canonical correlation analysis is to determine how similar two or more sets of variables are to one another. As in linear canonical correlation analysis, the aim is to account for as much of the variance in the relationships among the sets as possible in a low-dimensional space. Unlike linear canonical analysis, however, nonlinear canonical correlation analysis does not assume an interval level of measurement or that the relationships are linear. Another important difference is that nonlinear canonical correlation analysis establishes the similarity between the sets by simultaneously comparing linear combinations of the variables in each set to an unknown set—the object scores.

Example: An Analysis of Survey Results

The example in this chapter is from a survey (Verdegaal, 1985). The responses of 15 subjects to eight variables were recorded. The variables, variable labels, and value labels (categories) in the data set are shown in the following table.

Table 10-1
Survey data

Variable name	Variable label	Value label
<i>age</i>	Age in years	20–25, 26–30, 31–35, 36–40, 41–45, 46–50, 51–55, 56–60, 61–65, 66–70
<i>marital</i>	Marital status	Single, Married, Other
<i>pet</i>	Pets owned	No, Cat(s), Dog(s), Other than cat or dog, Various domestic animals

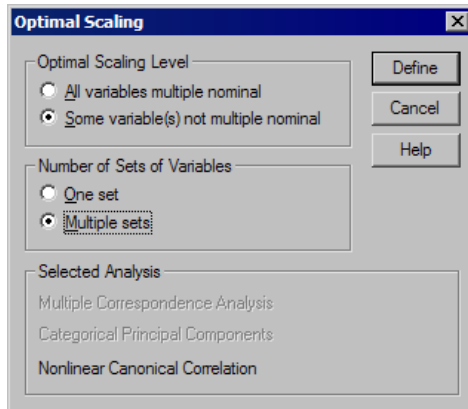
Variable name	Variable label	Value label
<i>news</i>	Newspaper read most often	None, Telegraaf, Volkskrant, NRC, Other
<i>music</i>	Music preferred	Classical, New wave, Popular, Variety, Don't like music
<i>live</i>	Neighborhood preference	Town, Village, Countryside
<i>math</i>	Math test score	0–5, 6–10, 11–15
<i>language</i>	Language test score	0–5, 6–10, 11–15, 16–20

This data set can be found in *verd1985.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS. The variables of interest are the first six, and they 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. This analysis requests a random initial configuration. By default, the initial configuration is numerical. However, when some of the variables are treated as single nominal with no possibility of ordering, it is best to choose a random initial configuration. This is the case with most of the variables in this study.

Examining the Data

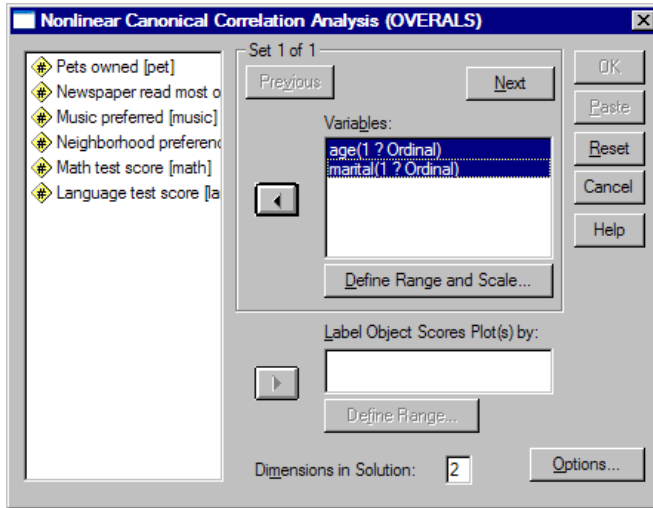
- ▶ To obtain a nonlinear canonical correlation analysis for this data set, from the menus choose:
Analyze
Data Reduction
Optimal Scaling...

Figure 10-1
Optimal Scaling dialog box



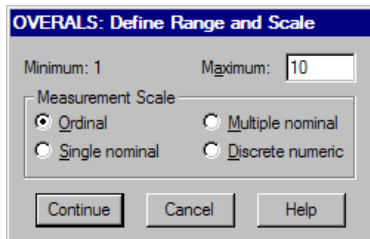
- ▶ Select Some variable(s) not multiple nominal in the Optimal Scaling Level group.
- ▶ Select Multiple sets in the Number of Sets of Variables group.
- ▶ Click Define.

Figure 10-2
Nonlinear Canonical Correlation Analysis dialog box



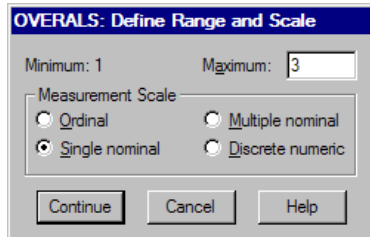
- ▶ Select *Age in years* and *Marital status* as variables for the first set.
- ▶ Select *age* and click Define Range and Scale.

Figure 10-3
Define Range and Scale dialog box



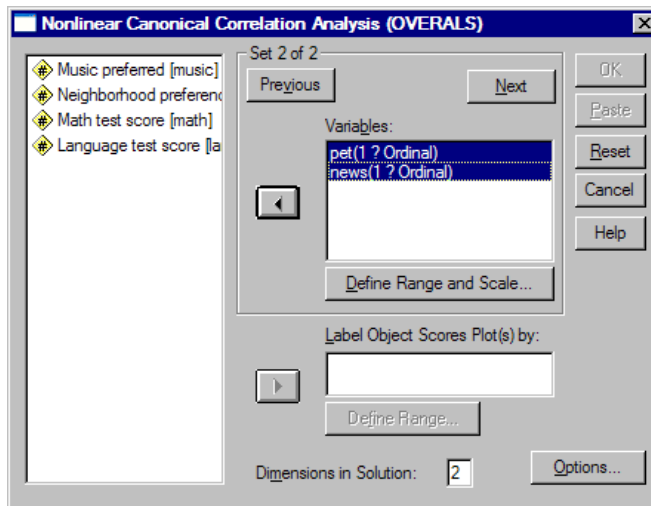
- ▶ Type 10 as the maximum value for this variable.
- ▶ Click Continue.
- ▶ Select *marital* and click Define Range and Scale in the Nonlinear Canonical Correlation Analysis dialog box.

Figure 10-4
Define Range and Scale dialog box



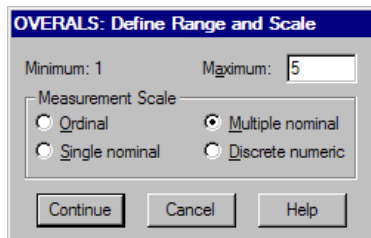
- ▶ Type 3 as the maximum value for this variable.
- ▶ Select Single nominal as the measurement scale.
- ▶ Click Continue.
- ▶ Click Next in the Nonlinear Canonical Correlation Analysis dialog box to define the next variable set.

Figure 10-5
Nonlinear Canonical Correlation Analysis dialog box



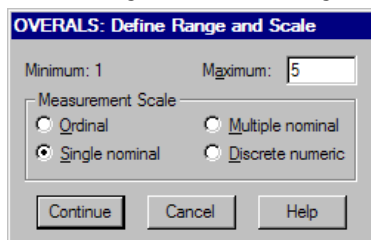
- ▶ Select *Pets owned* and *Newspaper read most often* as variables for the second set.
- ▶ Select *pet* and click Define Range and Scale.

Figure 10-6
Define Range and Scale dialog box



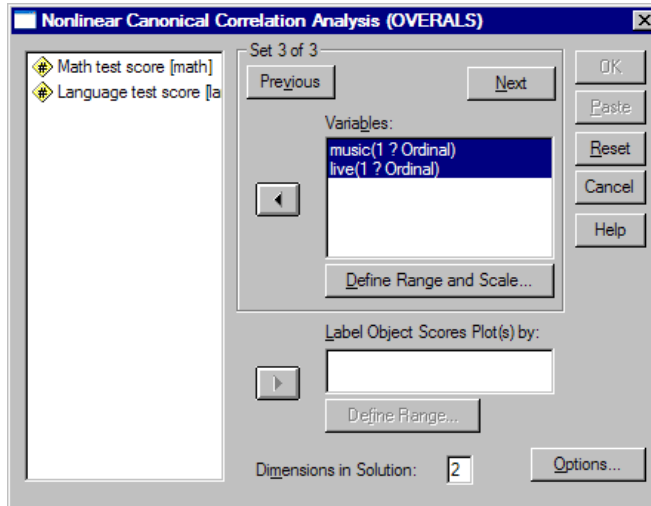
- ▶ Type 5 as the maximum value for this variable.
- ▶ Select Multiple nominal as the measurement scale.
- ▶ Click Continue.
- ▶ Select *news* and click Define Range and Scale in the Nonlinear Canonical Correlation Analysis dialog box.

Figure 10-7
Define Range and Scale dialog box



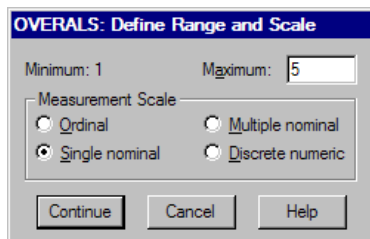
- ▶ Type 5 as the maximum value for this variable.
- ▶ Select Single nominal as the measurement scale.
- ▶ Click Continue.
- ▶ Click Next in the Nonlinear Canonical Correlation Analysis dialog box to define the last variable set.

Figure 10-8
Nonlinear Canonical Correlation Analysis dialog box



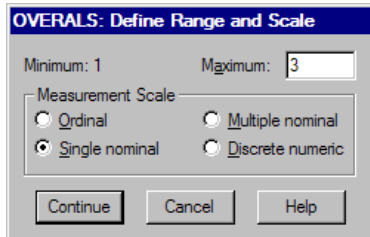
- ▶ Select *Music preferred* and *Neighborhood preference* as variables for the third set.
- ▶ Select *music* and click Define Range and Scale.

Figure 10-9
Define Range and Scale dialog box



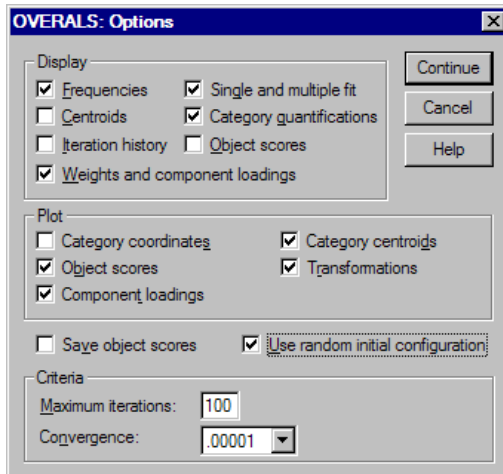
- ▶ Type 5 as the maximum value for this variable.
- ▶ Select Single nominal as the measurement scale.
- ▶ Click Continue.
- ▶ Select *live* and click Define Range and Scale in the Nonlinear Canonical Correlation Analysis dialog box.

Figure 10-10
Define Range and Scale dialog box



- ▶ Type 3 as the maximum value for this variable.
- ▶ Select Single nominal as the measurement scale.
- ▶ Click Continue.
- ▶ Click Options in the Nonlinear Canonical Correlation Analysis dialog box.

Figure 10-11
Options dialog box



- ▶ Deselect Centroids and select Weights and component loadings in the Display group.
- ▶ Select Category centroids and Transformations in the Plot group.
- ▶ Select Use random initial configuration.

- ▶ Click Continue.
- ▶ Click OK in the Nonlinear Canonical Correlation Analysis dialog box.

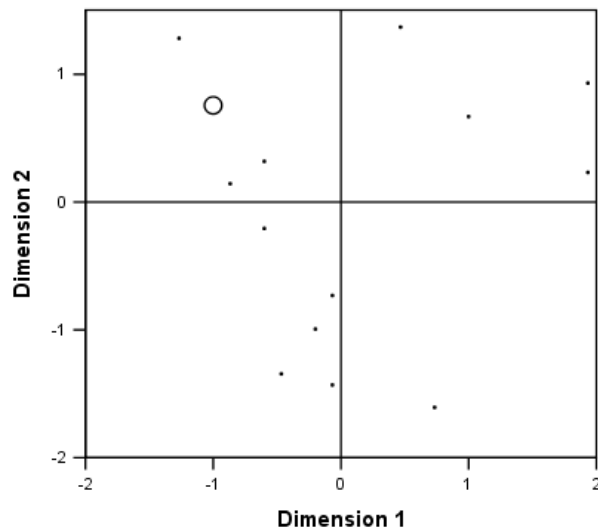
After a list of the variables with their levels of optimal scaling, categorical canonical correlation analysis with optimal scaling produces a table showing the frequencies of objects in categories. This table is especially important if there are missing data, since almost-empty categories are more likely to dominate the solution. In this example, there are no missing data.

A second preliminary check is to examine the plot of object scores for outliers. Outliers have such different quantifications from the other objects that they will be at the boundaries of the plot, thus dominating one or more dimensions.

If you find outliers, you can handle them in one of two ways. First, you can simply eliminate them from the data and run the nonlinear canonical correlation analysis again. Second, you can try recoding the extreme responses of the outlying objects by collapsing (merging) some categories.

As shown in the plot of object scores, there were no outliers for the survey data.

Figure 10-12
Object scores



Cases weighted by number of objects.

Accounting for Similarity between Sets

There are several ways to measure the association between sets in a nonlinear canonical correlation analysis, each of which is detailed in a separate table or set of tables.

Summary of Analysis

The fit and loss values tell you how well the nonlinear canonical correlation analysis solution fits the optimally quantified data with respect to the association between the sets. The summary of analysis table shows the fit value, loss values, and eigenvalues for the survey example.

Figure 10-13
Summary of analysis

		Dimension		Sum
		1	2	
Loss	Set 1	.240	.183	.423
	Set 2	.184	.408	.593
	Set 3	.171	.205	.376
	Mean	.199	.265	.464
Eigenvalue		.801	.735	
Fit				1.536

Loss is partitioned across dimensions and sets. For each dimension and set, loss represents the proportion of variation in the object scores that cannot be accounted for by the weighted combination of variables in the set. The average loss is labeled Mean. In this example, the average loss over sets is 0.464. Notice that more loss occurs for the second dimension than for the first.

The eigenvalue for each dimension equals 1 minus the average loss for the dimension and indicates how much of the relationship is shown by each dimension. The eigenvalues add up to the total fit. For Verdegaal's data, $0.801 / 1.536 = 52\%$ of the actual fit is accounted for by the first dimension.

The maximum fit value equals the number of dimensions and, if obtained, indicates that the relationship is perfect. The average loss value over sets and dimensions tells you the difference between the maximum fit and the actual fit. Fit plus the average

loss equals the number of dimensions. Perfect similarity rarely happens and usually capitalizes on trivial aspects in the data.

Another popular statistic with two sets of variables is the canonical correlation. Since the canonical correlation is related to the eigenvalue and thus provides no additional information, it is not included in the nonlinear canonical correlation analysis output. For two sets of variables, the canonical correlation per dimension is obtained by the formula:

$$\rho_d = 2 \times E_d - 1$$

where d is the dimension number and E is the eigenvalue.

You can generalize the canonical correlation for more than two sets with the formula:

$$\rho_d = ((K \times E_d) - 1)/(K - 1)$$

where d is the dimension number, K is the number of sets, and E is the eigenvalue. For our example,

$$\rho_1 = ((3 \times 0.801) - 1)/2 = 0.702$$

and

$$\rho_2 = ((3 \times 0.735) - 1)/2 = 0.603$$

Weights and Component Loadings

Another measure of association is the multiple correlation between linear combinations from each set and the object scores. If no variables in a set are multiple nominal, you can compute this by multiplying the weight and component loading of each variable within the set, adding these products, and taking the square root of the sum.

Figure 10-14
Weights

Set	Dimension	
	1	2
1 Age in years	.680	.789
Marital status	.296	-1.016
2 Newspaper read most often	-.845	-.361
3 Music preferred	.631	-.749
Neighborhood preference	-.484	-.780

Figure 10-15
Component loadings

Set	Dimension	
	1	2
1 Age in years ^{a,b}	.834	.259
Marital status ^{a,b}	.651	-.604
2 Pets owned ^{d,e} Dimension 1	.397	-.431
Dimension 2	-.277	.680
Newspaper read most often ^{a,b}	-.667	-.391
3 Music preferred ^{a,b}	.786	-.500
Neighborhood preference ^{a,b}	-.687	-.540

- a. Optimal Scaling Level: Ordinal
- b. Projections of the Single Quantified Variables in the Object Space
- c. Optimal Scaling Level: Single Nominal
- d. Optimal Scaling Level: Multiple Nominal
- e. Projections of the Multiple Quantified Variables in the Object Space

These figures give the weights and component loadings for the variables in this example. The multiple correlation (R) for the first weighted sum of optimally scaled variables (*Age in years* and *Marital status*) with the first dimension of object scores is:

$$\begin{aligned}
 R &= \sqrt{(0.701 \times 0.841 + (-0.273 \times -0.631))} \\
 &= \sqrt{(0.5895 + 0.1723)} \\
 &= 0.873
 \end{aligned}$$

For each dimension, $1 - \text{loss} = R^2$. For example, from the Summary of analysis table, $1 - 0.238 = 0.762$, which is 0.873 squared (plus some rounding error). Consequently, small loss values indicate large multiple correlations between weighted sums of

optimally scaled variables and dimensions. Weights are not unique for multiple nominal variables. For multiple nominal variables, use 1 – loss per set.

Partitioning Fit and Loss

The loss of each set is partitioned by the nonlinear canonical correlation analysis in several ways. The fit table presents the multiple fit, single fit, and single loss tables produced by the nonlinear canonical correlation analysis for the survey example. Note that multiple fit minus single fit equals single loss.

Figure 10-16
Partitioning fit and loss

Set		Multiple Fit			Single Fit			Single Loss		
		Dimension		Sum	Dimension		Sum	Dimension		Sum
		1	2		1	2		1	2	
1	Age in years ^a	.494	.676	1.170	.462	.622	1.085	.032	.054	.085
	Marital status ^b	.089	1.033	1.122	.088	1.03	1.120	.001	.000	.001
2	Pets owned ^c	.402	.439	.841						
	Newspaper read most often ^b	.724	.187	.911	.714	.130	.844	.010	.057	.067
3	Music preferred ^b	.421	.577	.998	.398	.561	.960	.022	.016	.039
	Neighborhood preference ^b	.234	.609	.843	.234	.608	.843	.000	.000	.000

a. Optimal Scaling Level: Ordinal

b. Optimal Scaling Level: Single Nominal

c. Optimal Scaling Level: Multiple Nominal

Single loss indicates the loss resulting from restricting variables to one set of quantifications (that is, single nominal, ordinal, or nominal). If single loss is large, it is better to treat the variables as multiple nominal. In this example, however, single fit and multiple fit are almost equal, which means that the multiple coordinates are almost on a straight line in the direction given by the weights.

Multiple fit equals the variance of the multiple category coordinates for each variable. These measures are analogous to the discrimination measures found in homogeneity analysis. You can examine the multiple fit table to see which variables discriminate best. For example, look at the multiple fit table for *Marital status* and *Newspaper read most often*. The fit values, summed across the two dimensions, are 1.122 for *Marital status* and 0.911 for *Newspaper read most often*. This tells us that a

person's marital status provides greater discriminatory power than the newspaper they subscribe to.

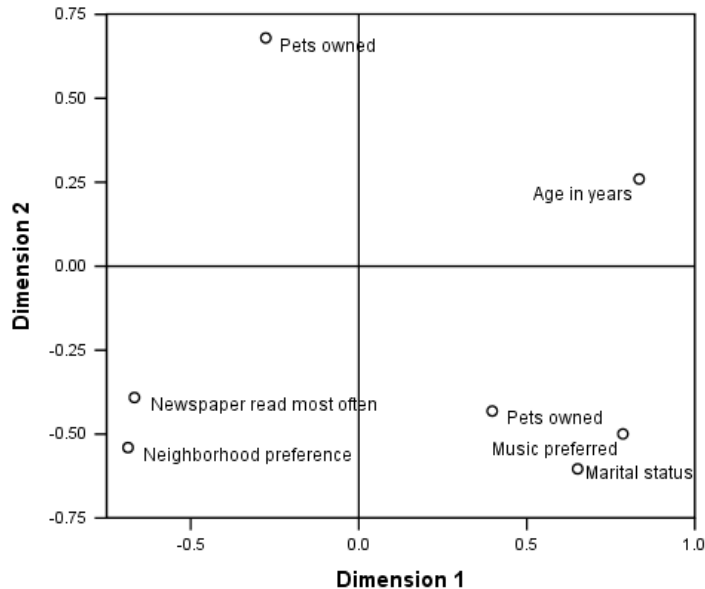
Single fit corresponds to the squared weight for each variable and equals the variance of the single category coordinates. As a result, the weights equal the standard deviations of the single category coordinates. Examining how the single fit is broken down across dimensions, we see that the variable *Newspaper read most often* discriminates mainly on the first dimension and that the variable *Marital status* discriminates almost totally on the second. In other words, the categories of *Newspaper read most often* are further apart in the first dimension than in the second, whereas the pattern is reversed for *Marital status*. In contrast, *Age in years* discriminates in both the first and second dimensions; thus, the spread of the categories is equal along both dimensions.

Component Loadings

The following figure shows the plot of component loadings for the survey data. When there are no missing data, the component loadings are equivalent to the Pearson correlations between the quantified variables and the object scores.

The distance from the origin to each variable point approximates the importance of that variable. The canonical variables are not plotted but can be represented by horizontal and vertical lines drawn through the origin.

Figure 10-17
Component loadings



The relationships between variables are apparent. There are two directions that do not coincide with the horizontal and vertical axes. One direction is determined by *Age in years*, *Newspaper read most often*, and *Neighborhood preference*. The other direction is defined by the variables *Marital status*, *Music preferred*, and *Pets owned*. The *Pets owned* variable is a multiple nominal variable, so there are two points plotted for it. Each quantification is interpreted as a single variable.

Transformation Plots

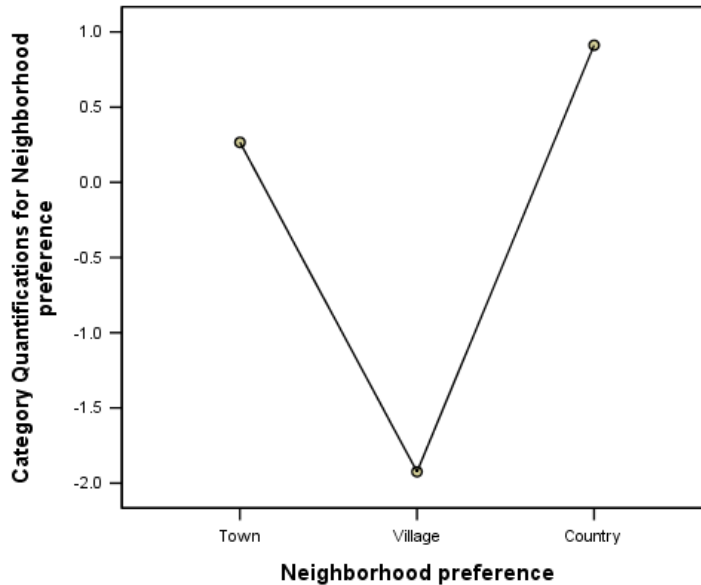
The different levels at which each variable can be scaled impose restrictions on the quantifications. Transformation plots illustrate the relationship between the quantifications and the original categories resulting from the selected optimal scaling level.

The transformation plot for *Neighborhood preference*, which was treated as nominal, displays a U-shaped pattern, in which the middle category receives the lowest quantification and the extreme categories receive values similar to each other.

This pattern indicates a quadratic relationship between the original variable and the transformed variable. Using an alternative optimal scaling level is not suggested for *Neighborhood preference*.

Figure 10-18

Transformation plot for Neighborhood preference (nominal)



The quantifications for *Newspaper read most often*, in contrast, correspond to an increasing trend across the three categories that have observed cases. The first category receives the lowest quantification, the second category receives a higher

value, and the third category receives the highest value. Although the variable is scaled as nominal, the category order is retrieved in the quantifications.

Figure 10-19

Transformation plot for Newspaper read most often (nominal)

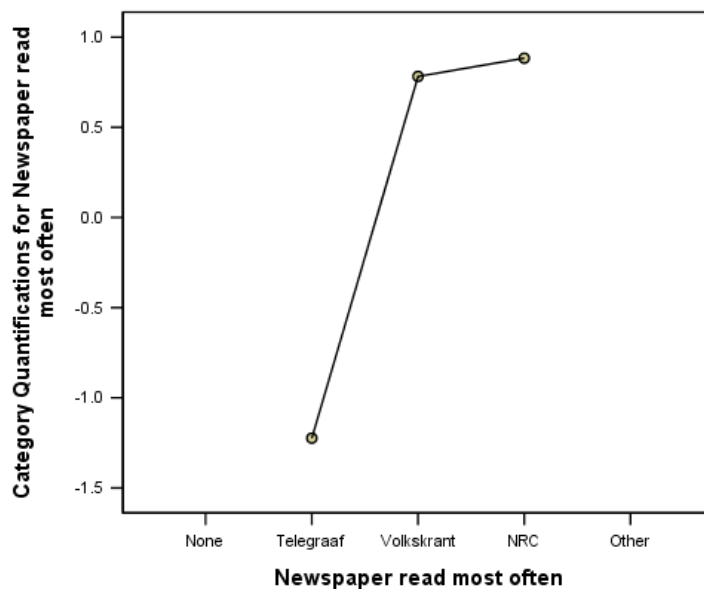
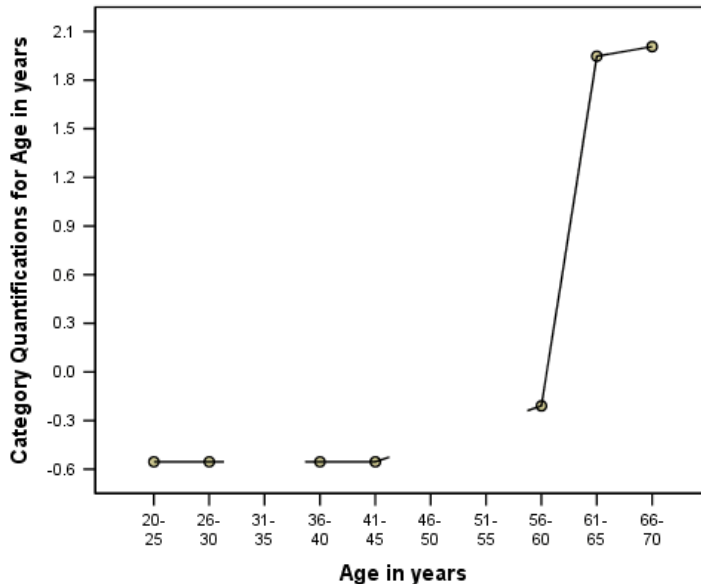


Figure 10-20
Transformation plot for *Age in years* (ordinal)



The transformation plot for *Age in years* displays an S-shaped curve. The four youngest observed categories all receive the same negative quantification, whereas the two oldest categories receive similar positive values. Consequently, collapsing all of the younger ages into one common category (that is, below 50) and collapsing the two oldest categories into one may be attempted. However, the exact equality of the quantifications for the younger groups indicates that restricting the order of the quantifications to the order of the original categories may not be desirable. Because the quantifications for the 26–30, 36–40, and 41–45 groups cannot be lower than the quantification for the 20–25 group, these values are set equal to the boundary value. Allowing these values to be smaller than the quantification for the youngest group (that is, treating age as nominal) may improve the fit. So although age may be considered an ordinal variable, treating it as such does not appear appropriate in this case. Moreover, treating age as numerical, and thus maintaining the distances between the categories, would substantially reduce the fit.

Single versus Multiple Category Coordinates

For every variable treated as single nominal, ordinal, or numerical, quantifications, single category coordinates, and multiple category coordinates are determined. These statistics for *Age in years* are presented.

Figure 10-21
Coordinates for *Age in years*

	Marginal Frequency	Quantification	Single Category Coordinates		Multiple Category Coordinates	
			Dimension		Dimension	
			1	2	1	2
20-25	3	-.554	-.377	-.437	-.192	-.139
26-30	5	-.554	-.377	-.437	-.404	-.623
31-35	0	.000				
36-40	1	-.554	-.377	-.437	-.318	-.733
41-45	1	-.554	-.377	-.437	-.356	-.534
46-50	0	.000				
51-55	0	.000				
56-60	2	-.209	-.142	-.165	-.435	.087
61-65	1	1.947	1.324	1.536	1.710	1.204
66-70	2	2.006	1.364	1.583	1.215	1.711
Missing	0					

Every category for which no cases were recorded receives a quantification of 0. For *Age in years*, this includes the 31–35, 46–50, and 51–55 categories. These categories are not restricted to be ordered with the other categories and do not affect any computations.

For multiple nominal variables, each category receives a different quantification on each dimension. For all other transformation types, a category has only one quantification, regardless of the dimensionality of the solution. The single category coordinates represent the locations of the categories on a line in the object space and equal the quantifications multiplied by the weights. For example, in the table for *Age in years*, the single category coordinates for category 8 (–0.192, –0.207) are the quantification multiplied by the weights.

The multiple category coordinates for variables treated as single nominal, ordinal, or numerical represent the coordinates of the categories in the object space before ordinal or linear constraints are applied. These values are unconstrained minimizers of the loss. For multiple nominal variables, these coordinates represent the quantifications of the categories.

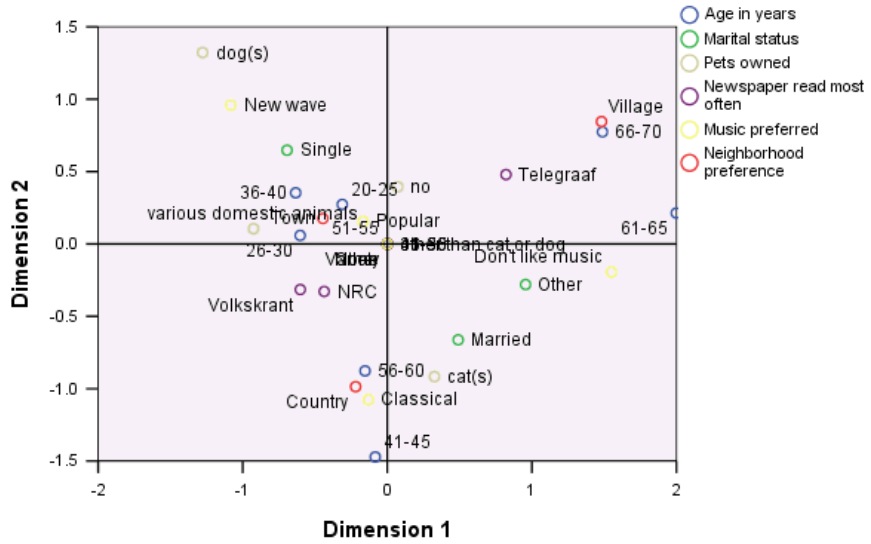
The effects of imposing constraints on the relationship between the categories and their quantifications are revealed by comparing the single with the multiple category coordinates. On the first dimension, the multiple category coordinates for *Age in years* decrease to category 2 and remain relatively at the same level until category 9, at which point a dramatic increase occurs. A similar pattern is evidenced for the second dimension. These relationships are removed in the single category coordinates, in which the ordinal constraint is applied. On both dimensions, the coordinates are now nondecreasing. The differing structure of the two sets of coordinates suggests that a nominal treatment may be more appropriate.

Centroids and Projected Centroids

The plot of centroids labeled by variables should be interpreted in the same way as the category quantifications plot in homogeneity analysis or the multiple category coordinates in nonlinear principal components analysis. By itself, such a plot shows how well variables separate groups of objects (the centroids are in the center of gravity of the objects).

Notice that the categories for *Age in years* are not separated very clearly. The younger age categories are grouped together at the left of the plot. As suggested previously, ordinal may be too strict a scaling level to impose on *Age in years*.

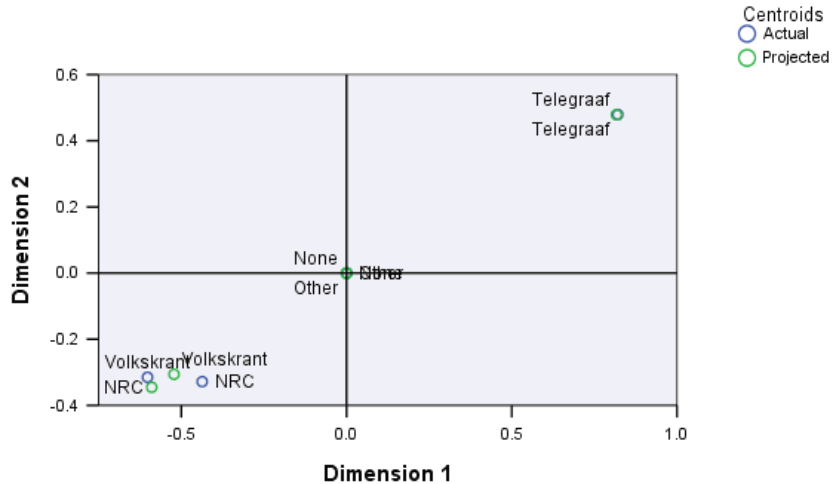
Figure 10-22
Centroids labeled by variables



When you request centroid plots, individual centroid and projected centroid plots for each variable labeled by value labels are also produced. The projected centroids are on a line in the object space.

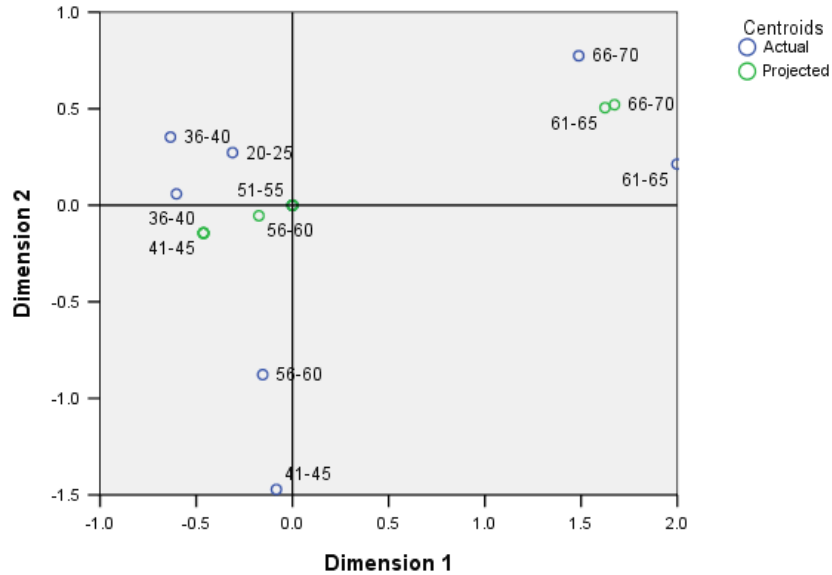
Figure 10-23

Centroids and projected centroids for Newspaper read most often



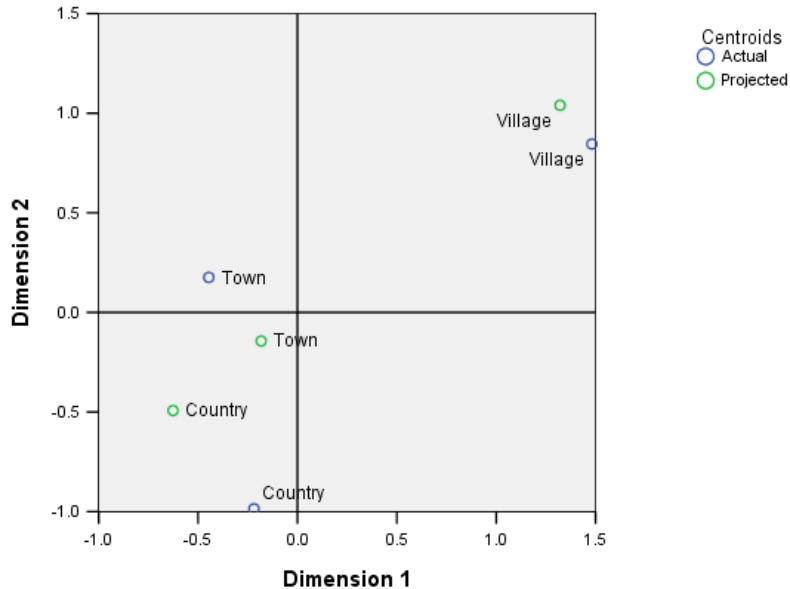
The actual centroids are projected onto the vectors defined by the component loadings. These vectors have been added to the centroid plots to aid in distinguishing the projected from the actual centroids. The projected centroids fall into one of four quadrants formed by extending two perpendicular reference lines through the origin. The interpretation of the direction of single nominal, ordinal, or numerical variables is obtained from the position of the projected centroids. For example, the variable *Newspaper read most often* is specified as single nominal. The projected centroids show that *Volkskrant* and *NRC* are contrasted with *Telegraaf*.

Figure 10-24
Centroids and projected centroids for Age in years



The problem with *Age in years* is evident from the projected centroids. Treating *Age in years* as ordinal implies that the order of the age groups must be preserved. To satisfy this restriction, all age groups below age 45 are projected into the same point. Along the direction defined by *Age in years*, *Newspaper read most often*, and *Neighborhood preference*, there is no separation of the younger age groups. Such a finding suggests treating the variable as nominal.

Figure 10-25
Centroids and projected centroids for Neighborhood preference



To understand the relationships among variables, find out what the specific categories (values) are for clusters of categories in the centroid plots. The relationships among *Age in years*, *Newspaper read most often*, and *Neighborhood preference* can be described by looking at the upper right and lower left of the plots. In the upper right, the age groups are the older respondents; they read the newspaper *Telegraaf* and prefer living in a village. Looking at the lower-left corner of each plot, you see that the younger to middle-aged respondents read the *Volkskrant* or *NRC* and want to live in the country or in a town. However, separating the younger groups is very difficult.

The same types of interpretations can be made about the other direction (*Music preferred*, *Marital status*, and *Pets owned*) by focusing on the upper left and the lower right of the centroid plots. In the upper left corner, we find that single people tend to have dogs and like new wave music. The married and other categories for marital have cats; the former group prefers classical music and the latter group does not like music.

An Alternative Analysis

The results of the analysis suggest that treating *Age in years* as ordinal does not appear appropriate. Although *Age in years* is measured at an ordinal level, its relationships with other variables are not monotonic. To investigate the effects of changing the optimal scaling level to single nominal, you may rerun the analysis.

Running the Analysis

- ▶ Recall the Nonlinear Canonical Correlation Analysis dialog box.
- ▶ Select *age* and click Define Range and Scale.
- ▶ Select Single nominal as the scaling range.
- ▶ Click Continue.
- ▶ Click OK in the Nonlinear Canonical Correlation Analysis dialog box.

The eigenvalues for a two-dimensional solution are 0.806 and 0.757, respectively, with a total fit of 1.564.

Figure 10-26
Eigenvalues for the two-dimensional solution

		Dimension		Sum
		1	2	
Loss	Set 1	.249	.115	.363
	Set 2	.176	.408	.584
	Set 3	.157	.205	.363
	Mean	.194	.243	.436
Eigenvalue		.806	.757	
Fit				1.564

The multiple- and single-fit tables show that *Age in years* is still a highly discriminating variable, as evidenced by the sum of the multiple-fit values. In contrast to the earlier results, however, examination of the single-fit values reveals the discrimination to be almost entirely along the second dimension.

Figure 10-27
Partitioning fit and loss

Set	Multiple Fit			Single Fit			Single Loss			
	Dimension		Sum	Dimension		Sum	Dimension		Sum	
	1	2		1	2		1	2		
1	Age in years ^a	.246	1.197	1.443	.195	1.188	1.384	.051	.008	.059
	Marital status ^a	.273	1.136	1.409	.272	1.135	1.407	.001	.000	.002
2	Pets owned ^b	.530	.392	.921						
	Newspaper read most often ^a	.639	.185	.824	.631	.149	.780	.008	.036	.044
3	Music preferred ^a	.604	.438	1.041	.603	.437	1.040	.000	.001	.001
	Neighborhood preference ^a	.075	.822	.897	.075	.822	.897	.000	.000	.000

a. Optimal Scaling Level: Single Nominal

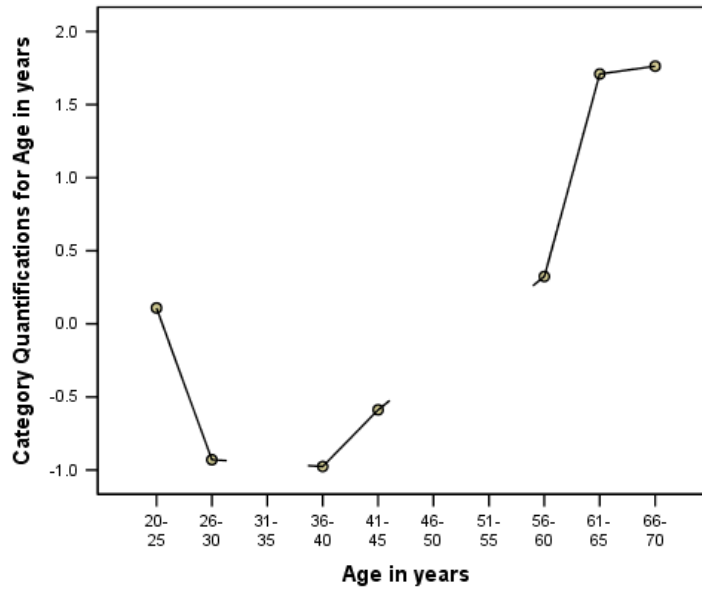
b. Optimal Scaling Level: Multiple Nominal

Turn to the transformation plot for *Age in years*. The quantifications for a nominal variable are unrestricted, so the nondecreasing trend displayed when *Age in years* was treated ordinally is no longer present. There is a decreasing trend until the age of 40 and an increasing trend thereafter, corresponding to a U-shaped (quadratic)

relationship. The two older categories still receive similar scores, and subsequent analyses may involve combining these categories.

Figure 10-28

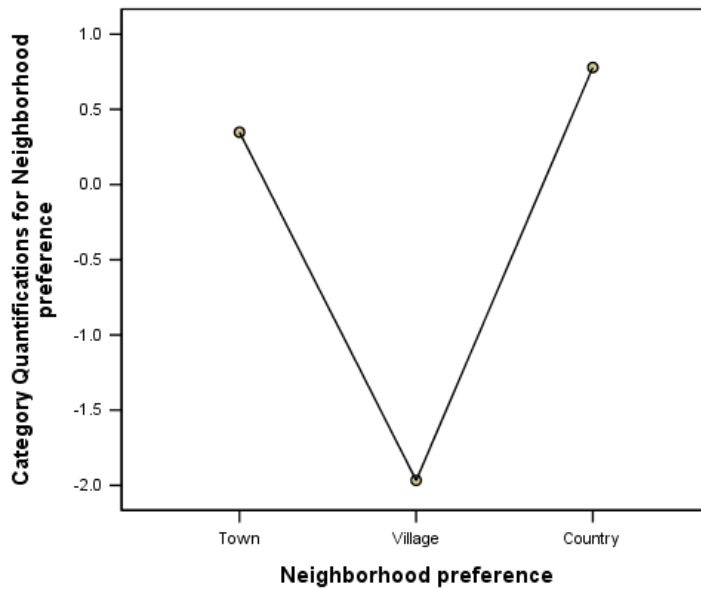
Transformation plot for Age in years (nominal)



The transformation plot for *Neighborhood preference* is shown here. Treating *Age in years* as nominal does not affect the quantifications for *Neighborhood preference* to any significant degree. The middle category receives the smallest quantification, with the extremes receiving large positive values.

Figure 10-29

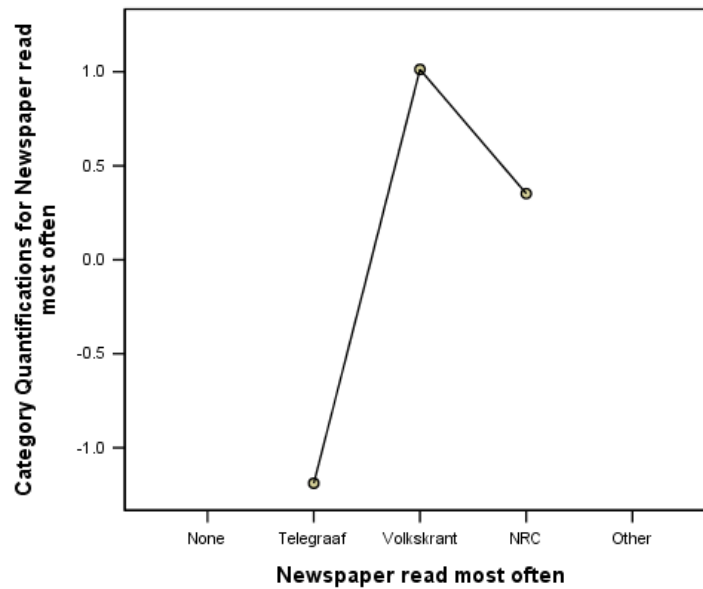
Transformation plot for Neighborhood preference (age nominal)



A change is found in the transformation plot for *Newspaper read most often*. Previously, an increasing trend was present in the quantifications, possibly suggesting an ordinal treatment for this variable. However, treating *Age in years* as nominal removes this trend from the news quantifications.

Figure 10-30

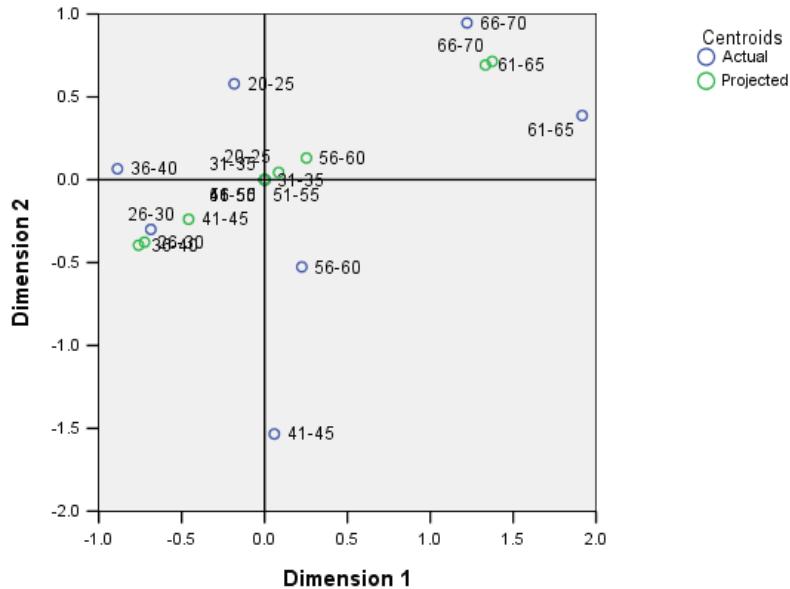
Transformation plot for Newspaper read most often (age nominal)



This is the centroid plot for *Age in years*. Notice that the categories do not fall in chronological order along the line joining the projected centroids. The 20–25 group is situated in the middle rather than at the end. The spread of the categories is much improved over the ordinal counterpart presented previously.

Figure 10-31

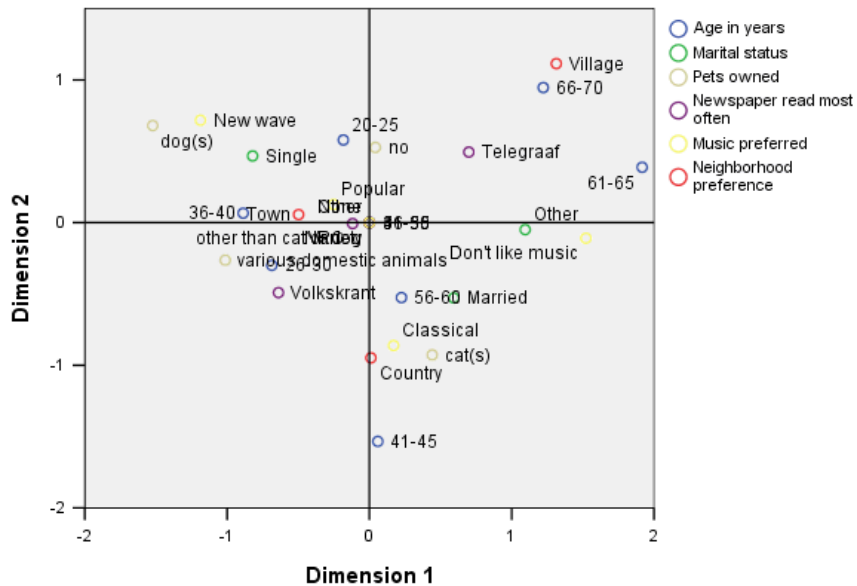
Centroids and projected centroids for Age in years (nominal)



Interpretation of the younger age groups is now possible from the centroid plot. The *Volkskrant* and the *NRC* categories are also further apart than in the previous analysis, allowing for separate interpretations of each. The groups between the ages of 26 and 45 read the *Volkskrant* and prefer country living. The 20–25 and 56–60 age groups read the *NRC*; the former group prefers to live in a town, and the latter group prefers country living. The oldest groups read the *Telegraaf* and prefer village living.

Interpretation of the other direction (*Music preferred, Marital status, and Pets owned*) is basically unchanged from the previous analysis. The only obvious difference is that people with a marital status of *Other* have either cats or no pets.

Figure 10-32
Centroids labeled by variables (age nominal)



General Suggestions

Once you have examined the initial results, you will probably want to refine your analysis by changing some of the specifications for the nonlinear canonical correlation analysis. Here are some tips for structuring your analysis:

- Create as many sets as possible. Put an important variable that you want to predict in a separate set by itself.
- Put variables that you consider predictors together in a single set. If there are many predictors, try to partition them into several sets.

- Put each multiple nominal variable in a separate set by itself.
- If variables are highly correlated to each other and you don't want this relationship to dominate the solution, put those variables together in the same set.

Recommended Readings

See the following texts for more information on nonlinear canonical correlation analysis:

Carroll, J. D. 1968. Generalization of canonical correlation analysis to three or more sets of variables. In: *Proceedings of the 76th annual convention of the American Psychological Association*, 3, , 227–228.

De Leeuw, J. 1984. *Canonical analysis of categorical data*. Leiden: DSWO Press.

Horst, P. 1961. Generalized canonical correlations and their applications to experimental data. *Journal of Clinical Psychology*, 17, 331–347.

Horst, P. 1961. Relations among m sets of measures. *Psychometrika*, 26, 129–149.

Kettenring, J. R. 1971. Canonical analysis of several sets of variables. *Biometrika*, 58, 433–460.

Van der Burg, E. 1988. *Nonlinear canonical correlation and some related techniques*. Leiden: DSWO Press.

Van der Burg, E., and J. De Leeuw. 1983. Nonlinear canonical correlation. *British Journal of Mathematical and Statistical Psychology*, 36, 54–80.

Van der Burg, E., J. De Leeuw, and R. Verdegaal. 1988. Homogeneity analysis with k sets of variables: An alternating least squares method with optimal scaling features. *Psychometrika*, 53, 177–197.

Verboon, P., and I. A. Van der Lans. 1994. Robust canonical discriminant analysis. *Psychometrika*, 59, 485–507.

Correspondence Analysis

A **correspondence table** is any two-way table whose cells contain some measurement of correspondence between the rows and the columns. The measure of correspondence can be any indication of the similarity, affinity, confusion, association, or interaction between the row and column variables. A very common type of correspondence table is a crosstabulation, where the cells contain frequency counts.

Such tables can be obtained easily with the Crosstabs procedure. However, a crosstabulation does not always provide a clear picture of the nature of the relationship between the two variables. This is particularly true if the variables of interest are nominal (with no inherent order or rank) and contain numerous categories. Crosstabulation may tell you that the observed cell frequencies differ significantly from the expected values in a 10×9 crosstabulation of *occupation* and *breakfast cereal*, but it may be difficult to discern which occupational groups have similar tastes or what those tastes are.

Correspondence Analysis allows you to examine the relationship between two nominal variables graphically in a multidimensional space. It computes row and column scores and produces plots based on the scores. Categories that are similar to each other appear close to each other in the plots. In this way, it is easy to see which categories of a variable are similar to each other or which categories of the two variables are related. The Correspondence Analysis procedure also allows you to fit supplementary points into the space defined by the active points.

If the ordering of the categories according to their scores is undesirable or counterintuitive, order restrictions can be imposed by constraining the scores for some categories to be equal. For example, suppose you expect the variable *smoking behavior*, with categories *none*, *light*, *medium* and *heavy*, to have scores that correspond to this ordering. However, if the analysis orders the categories *none*, *light*, *heavy*, and *medium*, constraining the scores for *heavy* and *medium* to be equal preserves the ordering of the categories in their scores.

The interpretation of correspondence analysis in terms of distances depends on the normalization method used. The Correspondence Analysis procedure can be used to analyze either the differences between categories of a variable or the differences between variables. With the default normalization, it analyzes the differences between the row and column variables.

The correspondence analysis algorithm is capable of many kinds of analyses. Centering the rows and columns and using chi-square distances corresponds to standard correspondence analysis. However, using alternative centering options combined with Euclidean distances allows for an alternative representation of a matrix in a low-dimensional space.

Three examples will be presented. The first employs a relatively small correspondence table and illustrates the concepts inherent in correspondence analysis. The second example demonstrates a practical marketing application. The final example uses a table of distances in a multidimensional scaling approach.

Normalization

Normalization is used to distribute the inertia over the row scores and column scores. Some aspects of the correspondence analysis solution, such as the singular values, the inertia per dimension, and the contributions, do not change under the various normalizations. The row and column scores and their variances are affected. Correspondence analysis has several ways to spread the inertia. The three most common include spreading the inertia over the row scores only, spreading the inertia over the column scores only, or spreading the inertia symmetrically over both the row scores and the column scores.

Row principal. In row principal normalization, the Euclidean distances between the row points approximate chi-square distances between the rows of the correspondence table. The row scores are the weighted average of the column scores. The column scores are standardized to have a weighted sum of squared distances to the centroid of 1. Since this method maximizes the distances between row categories, you should use row principal normalization if you are primarily interested in seeing how categories of the row variable differ from each other.

Column principal. On the other hand, you might want to approximate the chi-square distances between the columns of the correspondence table. In that case, the column scores should be the weighted average of the row scores. The row scores are standardized to have a weighted sum of squared distances to the centroid of 1. This method maximizes the distances between column categories and should be used if you are primarily concerned with how categories of the column variable differ from each other.

Symmetrical. You can also treat the rows and columns symmetrically. This normalization spreads inertia equally over the row and column scores. Note that neither the distances between the row points nor the distances between the column points are approximations of chi-square distances in this case. Use this method if you are primarily interested in the differences or similarities between the two variables. Usually, this is the preferred method to make biplots.

Principal. A fourth option is called principal normalization, in which the inertia is spread twice in the solution—once over the row scores and once over the column scores. You should use this method if you are interested in the distances between the row points and the distances between the column points separately but not in how the row and column points are related to each other. Biplots are not appropriate for this normalization option and are therefore not available if you have specified the principal normalization method.

Example: Smoking Behavior by Job Category

The aim of correspondence analysis is to show the relationships between the rows and columns of a correspondence table. You will use a hypothetical table introduced by Greenacre (Greenacre, 1984) to illustrate the basic concepts. This information is collected in *smoking.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

The table of interest is formed by the crosstabulation of smoking behavior by job category. The variable *Staff Group* contains the job categories *Sr Managers*, *Jr Managers*, *Sr Employees*, *Jr Employees*, and *Secretaries*, which will be used to create the solution, plus the category *National Average*, which can be used as supplementary to the analysis. The variable *Smoking* contains the behaviors *None*, *Light*, *Medium*,

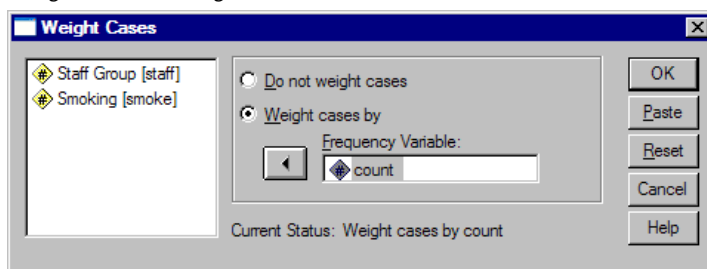
and *Heavy*, which will be used to create the solution, plus the categories *No Alcohol* and *Alcohol*, which can be used as supplementary to the analysis.

Running the Analysis

- ▶ Before you can run the Correspondence Analysis procedure, the setup of the data requires that the cases be weighted by the variable *count*. To do this, from the menus choose:

Data
Weight Cases...

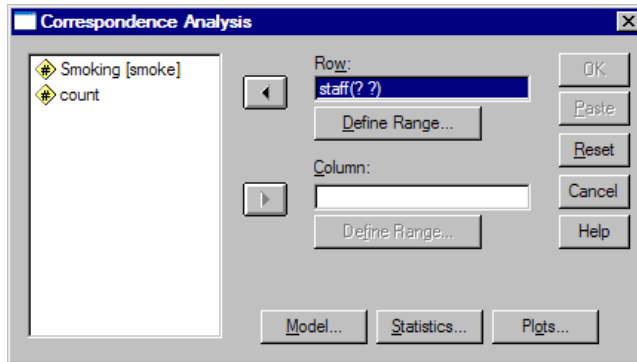
Figure 11-1
Weight Cases dialog box



- ▶ Weight cases by *count*.
- ▶ Click OK.
- ▶ Then, to obtain a correspondence analysis in two dimensions using row principal normalization, from the menus choose:

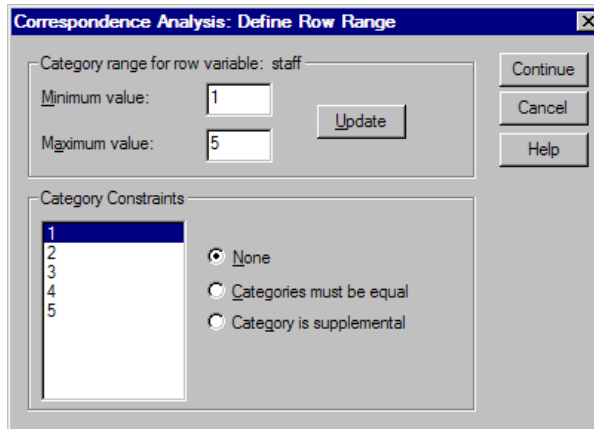
Analyze
Data Reduction
Correspondence Analysis...

Figure 11-2
Correspondence Analysis dialog box



- ▶ Select *Staff Group* as the row variable.
- ▶ Click Define Range.

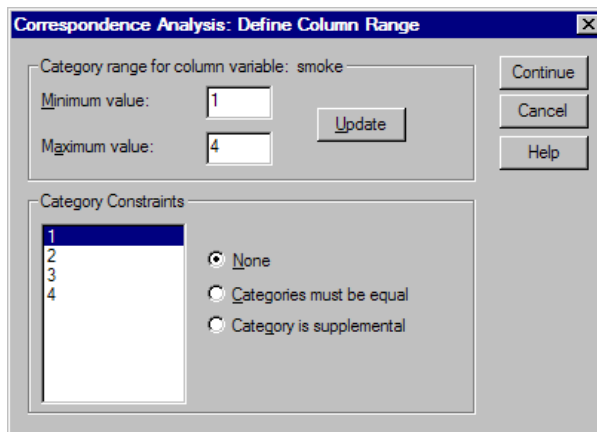
Figure 11-3
Define Row Range dialog box



- ▶ Type 1 as the minimum value.
- ▶ Type 5 as the maximum value.
- ▶ Click Update.

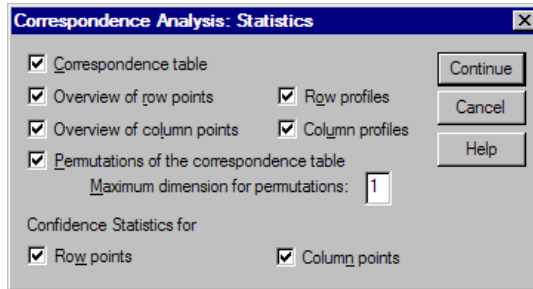
- ▶ Click Continue.
- ▶ Select *Smoking* as the column variable.
- ▶ Click Define Range in the Correspondence Analysis dialog box.

Figure 11-4
Define Column Range dialog box



- ▶ Type 1 as the minimum value.
- ▶ Type 4 as the maximum value.
- ▶ Click Update.
- ▶ Click Continue.
- ▶ Click Statistics in the Correspondence Analysis dialog box.

Figure 11-5
Statistics dialog box



- ▶ Select Row profiles and Column profiles.
- ▶ Select Permutations of the correspondence table.
- ▶ Select confidence statistics for Row points and Column points.
- ▶ Click Continue.
- ▶ Click OK in the Correspondence Analysis dialog box.

Correspondence Table

The correspondence table shows the distribution of smoking behavior for five levels of job category. The rows of the correspondence table represent the job categories. The columns represent the smoking behavior.

Figure 11-6
Correspondence table

Staff Group	Smoking				
	None	Light	Medium	Heavy	Active Margin
Sr Managers	4	2	3	2	11
Jr Managers	4	3	7	4	18
Sr Employees	25	10	12	4	51
Jr Employees	18	24	33	13	88
Secretaries	10	6	7	2	25
Active Margin	61	45	62	25	193

The marginal row totals show that the company has far more employees, both junior and senior, than managers and secretaries. However, the distribution of senior and junior positions for the managers is approximately the same as the distribution of senior and junior positions for the employees. Looking at the column totals, you see that there are similar numbers of nonsmokers and medium smokers. Furthermore, heavy smokers are outnumbered by each of the other three categories. But what, if anything, do any of these job categories have in common regarding smoking behavior? And what is the relationship between job category and smoking?

Dimensionality

Ideally, you want a correspondence analysis solution that represents the relationship between the row and column variables in as few dimensions as possible. But it is frequently useful to look at the maximum number of dimensions to see the relative contribution of each dimension. The maximum number of dimensions for a correspondence analysis solution equals the number of active rows minus 1 or the number of active columns minus 1, whichever is less. An active row or column is one for which a distinct set of scores is found. Supplementary rows or columns are not active. In the present example, the maximum number of dimensions is $\min(5,4) - 1 = 3$.

The first dimension displays as much of the inertia (a measure of the variation in the data) as possible, the second is orthogonal to the first and displays as much of the remaining inertia as possible, and so on. It is possible to split the total inertia into components attributable to each dimension. You can then evaluate the inertia shown by a particular dimension by comparing it to the total inertia. For example, the first dimension displays 87.8% ($0.075/0.085$) of the total inertia, whereas the second dimension displays only 11.8% ($0.010/0.085$).

Figure 11-7
Inertia per dimension

Dimension	Singular Value	Inertia	Chi Square	Sig.	Proportion of Inertia		Confidence Singular Value		
					Accounted for	Cumulative	Standard Deviation	Correlation	
								2	3
1	.273	.075			.878	.878	.070	.020	-.013
2	.100	.010			.118	.995	.076		-.059
3	.020	.000			.005	1.000	.072		
Total		.085	16.442	.172 ^a	1.000	1.000			

a. 12 degrees of freedom

If you decide that the first p dimensions of a q dimensional solution show enough of the total inertia, then you do not have to look at higher dimensions. In this example, the two-dimensional solution is sufficient, since the third dimension represents less than 1.0% of the total inertia.

The singular values can be interpreted as the correlation between the row and column scores. They are analogous to the Pearson correlation coefficient (r) in correlation analysis. For each dimension, the singular value squared (eigenvalue) equals the inertia and thus is another measure of the importance of that dimension.

Biplot

Correspondence analysis generates a variety of plots that graphically illustrate the underlying relationships between categories and between variables. This is the scatterplot of the row and column scores for the two-dimensional solution.

Figure 11-8
 Plot of row and column scores (symmetrical normalization)



The interpretation of the plot is fairly simple—row/column points that are close together are more alike than points that are far apart. The second dimension separates managers from other employees, while the first separates senior from junior, with secretaries in between.

The symmetrical normalization makes it easy to examine the relationship between job category and smoking. For example, managers are near the *Heavy* smoking category, while senior employees are closest to *None*. Junior employees seem to be associated with *Medium* or *Light* smoking, and secretaries are not strongly associated with any particular smoking behavior (but are far from *Heavy*).

Profiles and Distances

To determine the distance between categories, correspondence analysis considers the marginal distributions as well as the individual cell frequencies. It computes row and column profiles, which give the row and column proportions for each cell, based on the marginal totals.

Figure 11-9
Row profiles (symmetrical normalization)

Staff Group	Smoking				
	None	Light	Medium	Heavy	Active Margin
Sr Managers	.364	.182	.273	.182	1.000
Jr Managers	.222	.167	.389	.222	1.000
Sr Employees	.490	.196	.235	.078	1.000
Jr Employees	.205	.273	.375	.148	1.000
Secretaries	.400	.240	.280	.080	1.000
Mass	.316	.233	.321	.130	

The row profiles indicate the proportion of the row category in each column category. For example, among the senior employees, most are nonsmokers and very few are heavy smokers. In contrast, among the junior managers, most are medium smokers and very few are light smokers.

The column profiles indicate the proportion of the column in each row category. For example, most of the light smokers are junior employees. Similarly, most of the medium and heavy smokers are junior employees. Recall that the sample contains predominantly junior employees. It is not surprising that this staff category dominates the smoking categories.

Figure 11-10
Column profiles

Staff Group	Smoking				
	None	Light	Medium	Heavy	Mass
Sr Managers	.066	.044	.048	.080	.057
Jr Managers	.066	.067	.113	.160	.093
Sr Employees	.410	.222	.194	.160	.264
Jr Employees	.295	.533	.532	.520	.456
Secretaries	.164	.133	.113	.080	.130
Active Margin	1.000	1.000	1.000	1.000	

Mass is a measure that indicates the influence of an object based on its marginal frequency. Mass affects the **centroid**, which is the weighted mean row or column profile. The row centroid is the mean row profile. Points with a large mass, like junior employees, pull the centroid strongly to their location. A point with a small mass, like senior managers, pulls the row centroid only slightly to its location.

If you prefer to think of difference in terms of distance, then the greater the difference between row profiles, the greater the distance between points in a plot. For example, when using row principal normalization, the final configuration is one in which Euclidean distances between row points in the full dimensional space equal the chi-square distances between rows of the correspondence table. In a reduced space, the Euclidean distances approximate the chi-square distances. In turn, the chi-square distances are weighted profile distances. These weighted distances are based on mass.

Likewise, under column principal normalization, the Euclidean distances between column points in the full dimensional space equal the chi-square distances between columns of the correspondence table. Note, however, that under symmetric normalization these quantities are not equal.

The total inertia is defined as the weighted sum of all squared distances to the origin divided by the total over all cells, where the weights are the masses. Rows with a small mass influence the inertia only when they are far from the centroid. Rows with a large mass influence the total inertia, even when they are located close to the centroid. The same applies to columns.

Row and Column Scores

The row and column scores are the coordinates of the row and column points in the biplot.

Figure 11-11
Row scores (symmetrical normalization)

Staff Group	Mass	Score in Dimension		Inertia	Contribution				
		1	2		Of Point to Inertia of Dimension		Of Dimension to Inertia of Point		
					1	2	1	2	Total
Sr Managers	.057	-.126	.612	.003	.003	.214	.092	.800	.893
Jr Managers	.093	.495	.769	.012	.084	.551	.526	.465	.991
Sr Employees	.284	-.728	.034	.038	.512	.003	.999	.001	1.000
Jr Employees	.456	.446	-.183	.026	.331	.152	.942	.058	1.000
Secretaries	.130	-.385	-.249	.006	.070	.081	.865	.133	.999
Active Total	1.000			.085	1.000	1.000			

Figure 11-12
Column scores (symmetrical normalization)

Smoking	Mass	Score in Dimension		Inertia	Contribution				
		1	2		Of Point to Inertia of Dimension		Of Dimension to Inertia of Point		Total
					1	2	1	2	
None	.316	-.752	.096	.049	.654	.029	.994	.006	1.000
Light	.233	.190	-.446	.007	.031	.463	.327	.657	.984
Medium	.321	.375	-.023	.013	.166	.002	.982	.001	.983
Heavy	.130	.562	.625	.016	.150	.506	.684	.310	.995
Active Total	1.000			.085	1.000	1.000			

The column scores are related to the row scores via the profiles and singular value (from the inertia per dimension table). Specifically, the row scores are the matrix product of the row profiles and column scores, scaled by the singular value for each dimension. For example, the score of -0.126 for senior managers on the first dimension equals:

$$\frac{(0.364 \times -0.752) + (0.182 \times 0.190) + (0.273 \times 0.375) + (0.182 \times 0.562)}{0.273}$$

For row principal normalization, the singular value does not figure into this equation. The row points are in the weighted centroid of the active column points, where the weights correspond to the entries in the row profiles table. When the row points are the weighted average of the column points and the maximum dimensionality is used, the Euclidean distance between a row point and the origin equals the chi-square distance between the row and the average row, which in turn is equal to the inertia of a row. Because the chi-square statistic is equivalent to the total inertia times the sum of all cells of the correspondence table, you can think of the orientation of the row points as a pictorial representation of the chi-square statistic. A corresponding interpretation exists for column principal normalization but not for symmetrical.

Contributions

It is possible to compute the inertia displayed by a particular dimension. The scores on each dimension correspond to an orthogonal projection of the point onto that dimension. Thus, the inertia for a dimension equals the weighted sum of the squared distances from the scores on the dimension to the origin. However, whether this

applies to row or column scores (or both) depends on the normalization method used. Each row and column point contributes to the inertia. Row and column points that contribute substantially to the inertia of a dimension are important to that dimension. The contribution of a point to the inertia of a dimension is the weighted squared distance from the projected point to the origin divided by the inertia for the dimension.

The diagnostics that measure the contributions of points are an important aid in the interpretation of a correspondence analysis solution. Dominant points in the solution can easily be detected. For example, senior employees and junior employees are dominant points in the first dimension, contributing 84% of the inertia. Among the column points, none contributes 65% of the inertia for the first dimension alone.

The contribution of a point to the inertia of the dimensions depends on both the mass and the distance from the origin. Points that are far from the origin and have a large mass contribute most to the inertia of the dimension. Because supplementary points do not play any part in defining the solution, they do not contribute to the inertia of the dimensions.

In addition to examining the contribution of the points to the inertia per dimension, you can examine the contribution of the dimensions to the inertia per point. You can examine how the inertia of a point is spread over the dimensions by computing the percentage of the point inertia contributed by each dimension. Notice that the contributions of the dimensions to the point inertias do not all sum to one. In a reduced space, the inertia that is contributed by the higher dimensions is not represented. Using the maximum dimensionality would reveal the unaccounted inertia amounts.

The first two dimensions contribute all of the inertia for senior employees and junior employees and virtually all of the inertia for junior managers and secretaries. For senior managers, 11% of the inertia is not contributed by the first two dimensions. Two dimensions contribute a very large proportion of the inertia of the row points.

Similar results occur for the column points. For every active column point, two dimensions contribute at least 98% of the inertia. The third dimension contributes very little to these points.

Permutations of the Correspondence Table

Sometimes it is useful to order the categories of the rows and the columns. For example, you might have reason to believe that the categories of a variable correspond to a certain order, but you don't know the precise order. This ordination problem is

found in various disciplines—the seriation problem in archaeology, the ordination problem in phytosociology, and Guttman’s scalogram problem in the social sciences. Ordering can be achieved by taking the row and column scores as ordering variables. If you have row and column scores in p dimensions, p permuted tables can be made. When the first singular value is large, the first table will show a particular structure, with larger-than-expected relative frequencies close to the “diagonal.”

The following table shows the permutation of the correspondence table along the first dimension. Looking at the row scores for dimension 1, you can see that the ranking from lowest to highest is senior employees (−0.728), secretaries (−0.385), senior managers (−0.126), junior employees (0.446), and junior managers (0.495). Looking at the column scores for dimension 1, you see that the ranking is none, light, medium, and then heavy. These rankings are reflected in the ordering of the rows and columns of the table.

Figure 11-13
Permutation of the correspondence table

Staff Group	Smoking				
	None	Light	Medium	Heavy	Active Margin
Sr Employees	25	10	12	4	51
Secretaries	10	6	7	2	25
Sr Managers	4	2	3	2	11
Jr Employees	18	24	33	13	88
Jr Managers	4	3	7	4	18
Active Margin	61	45	62	25	193

Confidence Statistics

Assuming that the table to be analyzed is a frequency table and that the data are a random sample from an unknown population, the cell frequencies follow a multinomial distribution. From this, it is possible to compute the standard deviations and correlations of the singular values, row scores, and column scores.

In a one-dimensional correspondence analysis solution, you can compute a confidence interval for each score in the population. If the standard deviation is large, correspondence analysis is very uncertain of the location of the point in the population. On the other hand, if the standard deviation is small, then the correspondence analysis is fairly certain that this point is located very close to the point given by the solution.

In a multidimensional solution, if the correlation between dimensions is large, it may not be possible to locate a point in the correct dimension with much certainty. In such cases, multivariate confidence intervals must be calculated using the variance/covariance matrix that can be written to a file.

The confidence statistics for the row and column scores are shown. The standard deviations for the two manager categories are larger than the others, likely due to their relatively small numbers. The standard deviation for heavy smokers is also larger for the same reason. If you look at the correlations between the dimensions for the scores, you see that the correlations are generally small for the row and column scores with the exception of junior employees, with a correlation of 0.611.

Figure 11-14
Confidence statistics for row scores

Staff Group	Standard Deviation in Dimension		Correlation
	1	2	1-2
Sr Managers	.614	.917	.101
Jr Managers	.461	.511	.007
Sr Employees	.110	.157	.107
Jr Employees	.118	.124	.611
Secretaries	.158	.153	-.360

Figure 11-15
Confidence statistics for column scores

Smoking	Standard Deviation in Dimension		Correlation
	1	2	1-2
None	.118	.145	.402
Light	.281	.292	.054
Medium	.179	.332	.020
Heavy	.361	.441	-.155

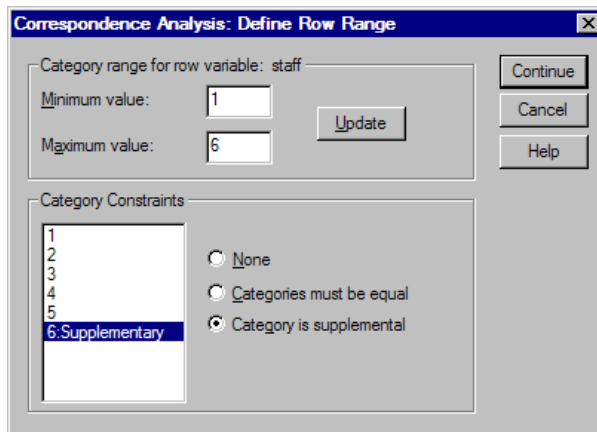
Supplementary Profiles

In correspondence analysis, additional categories can be represented in the space describing the relationships between the active categories. A supplementary profile defines a profile across categories of either the row or column variable and does not influence the analysis in any way. The data file contains one supplementary row and two supplementary columns.

The national average of people in each smoking category defines a supplementary row profile. The two supplementary columns define two column profiles across the categories of staff. The supplementary profiles define a point in either the row space or the column space. Because you will focus on both the rows and the columns separately, you will use principal normalization.

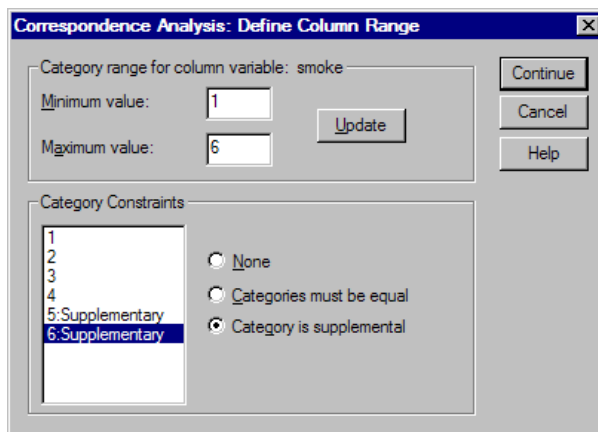
Running the Analysis

Figure 11-16
Define Row Range dialog box



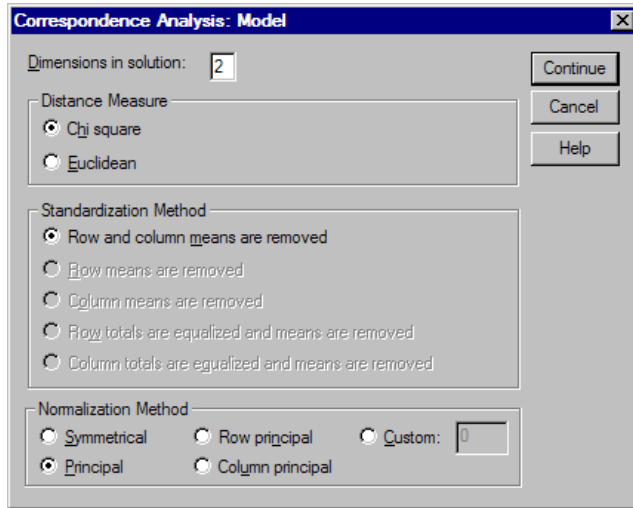
- ▶ To add the supplementary categories and obtain a principal normalization solution, recall the Correspondence Analysis dialog box.
- ▶ Select *staff* and click Define Range.
- ▶ Type 6 as the maximum value and click Update.
- ▶ Select 6 in the Category Constraints list and select Category is supplemental.
- ▶ Click Continue.
- ▶ Select *smoke* and click Define Range in the Correspondence Analysis dialog box.

Figure 11-17
Define Column Range dialog box



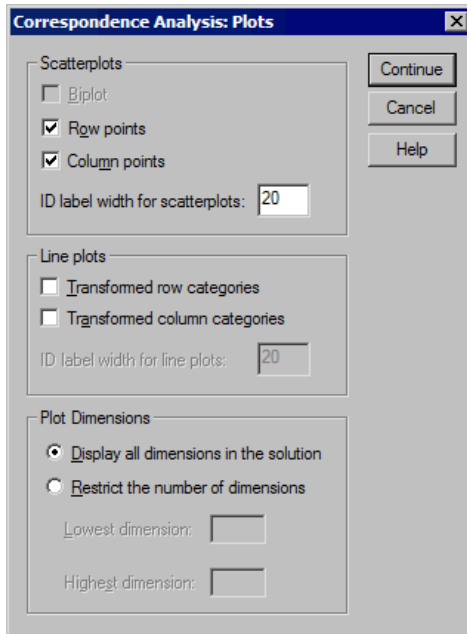
- ▶ Type 6 as the maximum value and click Update.
- ▶ Select 5 in the Category Constraints list and select Category is supplemental.
- ▶ Select 6 in the Category Constraints list and select Category is supplemental.
- ▶ Click Continue.
- ▶ Click Model in the Correspondence Analysis dialog box.

Figure 11-18
Model dialog box



- ▶ Select Principal as the normalization method.
- ▶ Click Continue.
- ▶ Click Plots in the Correspondence Analysis dialog box.

Figure 11-19
Plots dialog box

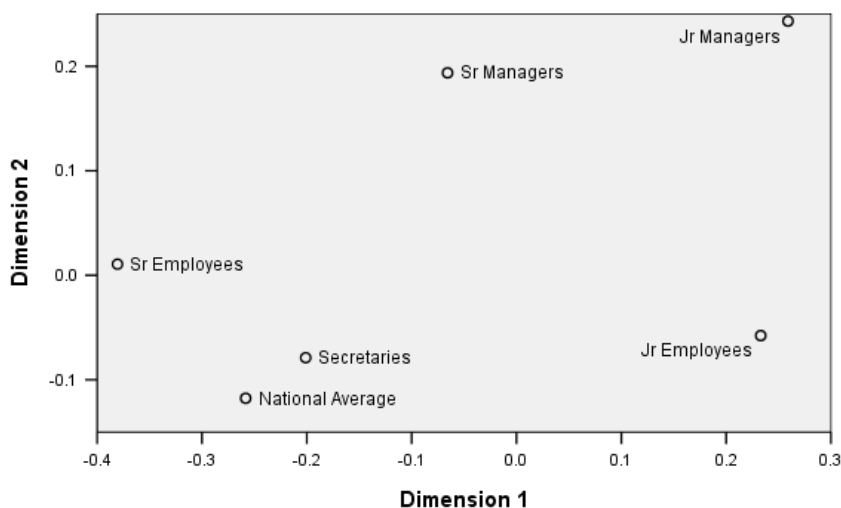


- ▶ Select Row points and Column points in the Scatterplots group.
- ▶ Click Continue.
- ▶ Click OK in the Correspondence Analysis dialog box.

The row points plot shows the first two dimensions for the row points with the supplementary point for *National Average*. *National Average* lies far from the origin, indicating that the sample is not representative of the nation in terms of smoking behavior. Secretaries and senior employees are close to the national average, whereas

junior managers are not. Thus, secretaries and senior employees have smoking behaviors similar to the national average, but junior managers do not.

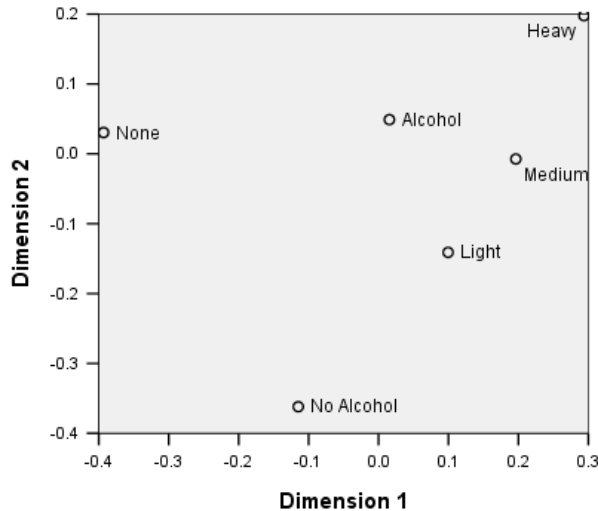
Figure 11-20
Row points (principal normalization)



The column points plot displays the column space with the two supplementary points for alcohol consumption. *Alcohol* lies near the origin, indicating a close correspondence between the alcohol profile and the average column profile. However, *No Alcohol* differs from the average column profile, illustrated by the large distance from the origin. The closest point to *No Alcohol* is *Light*. The light smokers profile is most similar to the nondrinkers. Among the smokers, *Medium* is next closest and *Heavy* is farthest. Thus, there is a progression in similarity to nondrinking from light

to heavy smoking. However, the relatively high proportion of secretaries in the *No Alcohol* group prevents any close correspondence to any of the smoking categories.

Figure 11-21
Column points (principal normalization)



Example: Perceptions of Coffee Brands

The previous example involved a small table of hypothetical data. Actual applications often involve much larger tables. In this example, you will use data pertaining to perceived images of six iced-coffee brands (Kennedy et al., 1996). This data set can be found in *coffee.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

For each of 23 iced-coffee image attributes, people selected all brands that were described by the attribute. The six brands are denoted as *AA*, *BB*, *CC*, *DD*, *EE*, and *FF* to preserve confidentiality.

Table 11-1
Iced-coffee attributes

Image attribute	Label	Image attribute	Label
good hangover cure	cure	fattening brand	fattening
low fat/calorie brand	low fat	appeals to men	men
brand for children	children	South Australian brand	South Australian
working class brand	working	traditional/old fashioned brand	traditional
rich/sweet brand	sweet	premium quality brand	premium
unpopular brand	unpopular	healthy brand	healthy
brand for fat/ugly people	ugly	high caffeine brand	caffeine
very fresh	fresh	new brand	new
brand for yuppies	yuppies	brand for attractive people	attractive
nutritious brand	nutritious	tough brand	tough
brand for women	women	popular brand	popular
minor brand	minor		

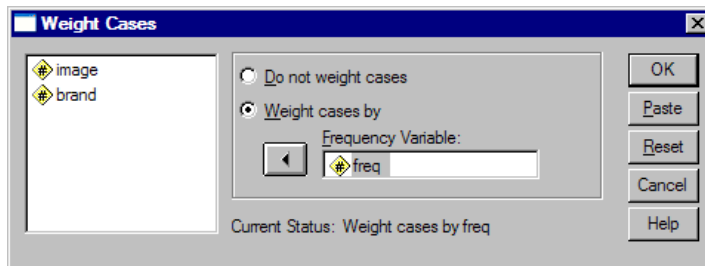
Initially, you will focus on how the attributes are related to each other and how the brands are related to each other. Using principal normalization spreads the total inertia once over the rows and once over the columns. Although this prevents biplot interpretation, the distances between the categories for each variable can be examined.

Running the Analysis

- The setup of the data requires that the cases be weighted by the variable *freq*. To do this, from the menus choose:

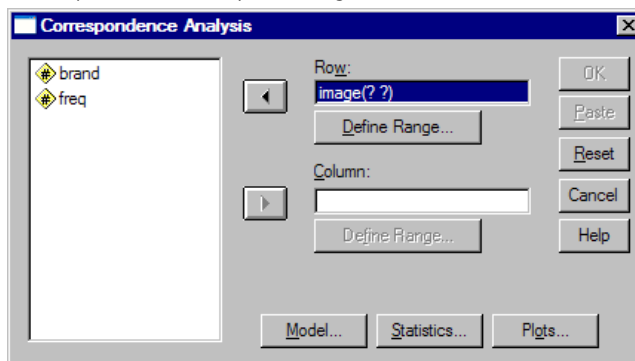
Data
Weight Cases...

Figure 11-22
Weight Cases dialog box



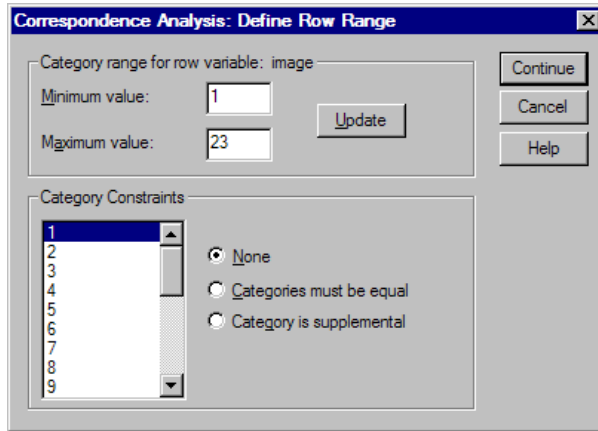
- ▶ Weight cases by *freq*.
- ▶ Click OK.
- ▶ To obtain an initial solution in five dimensions with principal normalization, from the menus choose:
Analyze
Data Reduction
Correspondence Analysis...

Figure 11-23
Correspondence Analysis dialog box



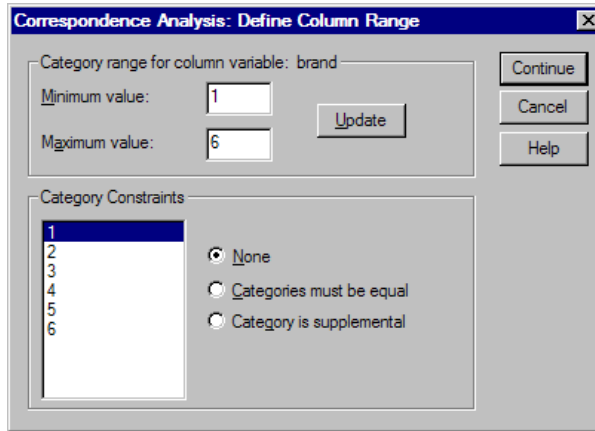
- ▶ Select *image* as the row variable.
- ▶ Click Define Range.

Figure 11-24
Define Row Range dialog box



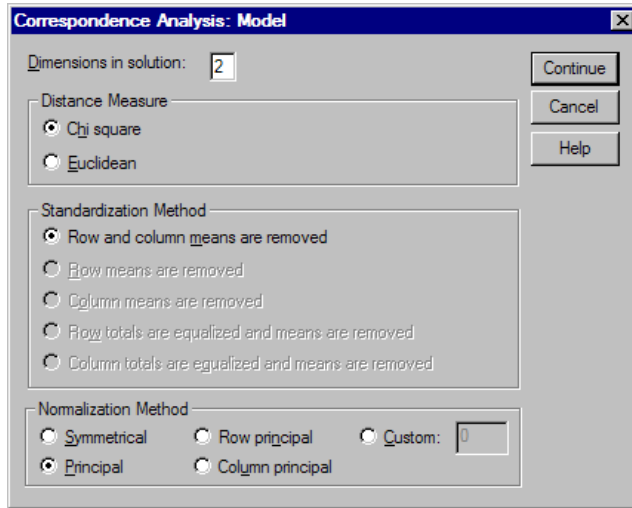
- ▶ Type 1 as the minimum value.
- ▶ Type 23 as the maximum value.
- ▶ Click Update.
- ▶ Click Continue.
- ▶ Select *brand* as the column variable.
- ▶ Click Define Range in the Correspondence Analysis dialog box.

Figure 11-25
Define Column Range dialog box



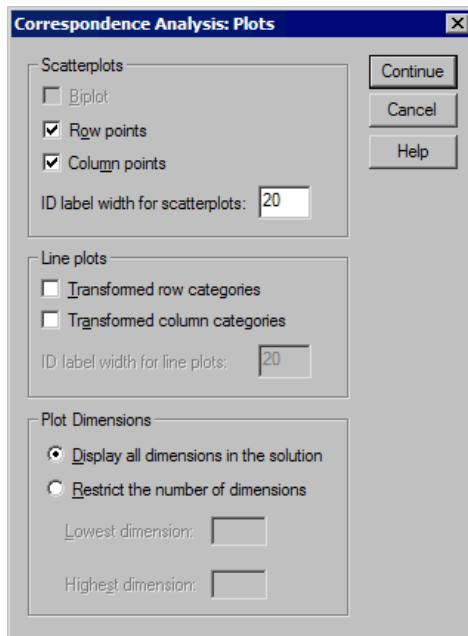
- ▶ Type 1 as the minimum value.
- ▶ Type 6 as the maximum value.
- ▶ Click Update.
- ▶ Click Continue.
- ▶ Click Model in the Correspondence Analysis dialog box.

Figure 11-26
Model dialog box



- ▶ Select Principal as the normalization method.
- ▶ Click Continue.
- ▶ Click Plots in the Correspondence Analysis dialog box.

Figure 11-27
Plots dialog box



- ▶ Select Row points and Column points in the Scatterplots group.
- ▶ Click Continue.
- ▶ Click OK in the Correspondence Analysis dialog box.

Dimensionality

The inertia per dimension shows the decomposition of the total inertia along each dimension. Two dimensions account for 83% of the total inertia. Adding a third dimension adds only 8.6% to the accounted-for inertia. Thus, you elect to use a two-dimensional representation.

Figure 11-28
Inertia per dimension

Dimension	Singular Value	Inertia	Chi Square	Sig.	Proportion of Inertia		Confidence Singular Value	
					Accounted for	Cumulative	Standard Deviation	Correlation 2
1	.711	.508			.829	.829	.009	.132
2	.399	.159			.198	.827	.014	
3	.263	.069			.086	.913		
4	.234	.055			.068	.982		
5	.121	.015			.018	1.000		
Total		.804	3746.97	.000 ^a	1.000	1.000		

a. 110 degrees of freedom

Contributions

The row points overview shows the contributions of the row points to the inertia of the dimensions and the contributions of the dimensions to the inertia of the row points. If all points contributed equally to the inertia, the contributions would be 0.043. *Healthy* and *low fat* both contribute a substantial portion to the inertia of the first dimension. *Men* and *tough* contribute the largest amounts to the inertia of the second dimension. Both *ugly* and *fresh* contribute very little to either dimension.

Figure 11-29
Attribute contributions

image	Mass	Score in Dimension		Inertia	Contribution				
		1	2		Of Point to Inertia of Dimension		Of Dimension to Inertia of Point		Total
					1	2	1	2	
fattening	.080	-.514	-.265	.033	.042	.035	.652	.173	.825
men	.051	-.852	.825	.072	.073	.219	.512	.480	.992
South Australian	.057	-.303	-.350	.046	.010	.044	.114	.152	.266
traditional	.040	-.703	-.532	.043	.039	.071	.454	.260	.715
premium	.042	-.444	-.582	.028	.016	.090	.296	.509	.805
healthy	.053	1.200	.174	.081	.152	.010	.953	.020	.973
caffeine	.047	-.452	.124	.014	.019	.005	.702	.053	.755
new	.047	.960	.147	.048	.086	.006	.893	.021	.914
attractive	.041	.657	-.056	.019	.035	.001	.911	.007	.918
tough	.039	-.850	1.002	.070	.056	.246	.404	.560	.964
popular	.060	-.697	-.042	.038	.058	.001	.771	.003	.774
cure	.026	-.389	.266	.009	.008	.011	.446	.209	.655
low fat	.052	1.305	.196	.094	.175	.013	.941	.021	.962
children	.024	-.352	-.513	.017	.006	.041	.179	.380	.559
working	.045	-.785	.477	.040	.055	.064	.693	.255	.948
sweet	.038	-.519	-.683	.048	.020	.112	.212	.368	.580
unpopular	.024	.489	.186	.010	.011	.005	.585	.085	.670
ugly	.030	.006	-.109	.003	.000	.002	.000	.131	.131
fresh	.036	-.096	-.100	.002	.001	.002	.196	.214	.410
yuppies	.034	.380	-.301	.012	.010	.019	.392	.246	.637
nutritious	.040	.722	.055	.022	.041	.001	.946	.006	.951
women	.054	.758	-.063	.032	.062	.001	.965	.007	.972
minor	.040	.579	.063	.023	.027	.001	.593	.007	.600
Active Total	1.000			.804	1.000	1.000			

Two dimensions contribute a large amount to the inertia for most row points. The large contributions of the first dimension to *healthy*, *new*, *attractive*, *low fat*, *nutritious*, and *women* indicate that these points are very well represented in one dimension. Consequently, the higher dimensions contribute little to the inertia of these points, which will lie very near the horizontal axis. The second dimension contributes most to *men*, *premium*, and *tough*. Both dimensions contribute very little to the inertia for *South Australian* and *ugly*, so these points are poorly represented.

The column points overview displays the contributions involving the column points. Brands *CC* and *DD* contribute the most to the first dimension, whereas *EE* and *FF* explain a large amount of the inertia for the second dimension. *AA* and *BB* contribute very little to either dimension.

Figure 11-30
Brand contributions

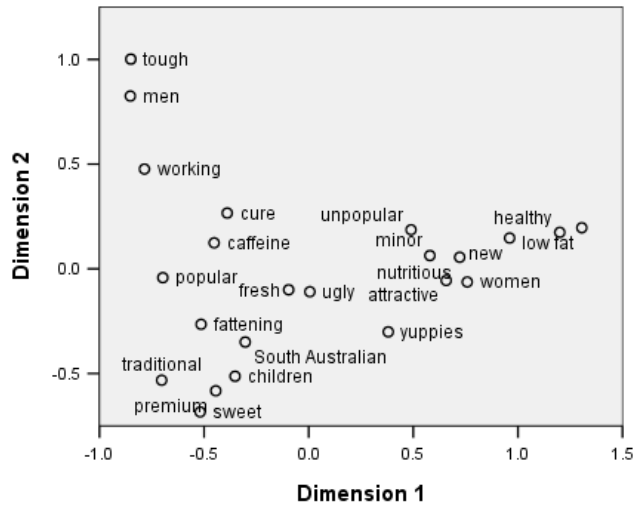
brand	Mass	Score in Dimension		Inertia	Contribution				
		1	2		Of Point to Inertia of Dimension		Of Dimension to Inertia of Point		Total
					1	2	1	2	
AA	.217	-.659	.048	.127	.187	.003	.744	.004	.748
BB	.131	-.284	-.404	.078	.021	.134	.135	.272	.407
CC	.185	.998	.078	.193	.362	.007	.951	.006	.957
DD	.162	.915	.101	.146	.267	.010	.928	.011	.939
EE	.152	-.651	.708	.153	.127	.477	.420	.494	.914
FF	.153	-.343	-.618	.107	.036	.369	.169	.550	.718
Active Total	1.000			.804	1.000	1.000			

In two dimensions, all brands but *BB* are well represented. *CC* and *DD* are represented well in one dimension. The second dimension contributes the largest amounts for *EE* and *FF*. Notice that *AA* is represented well in the first dimension but does not have a very high contribution to that dimension.

Plots

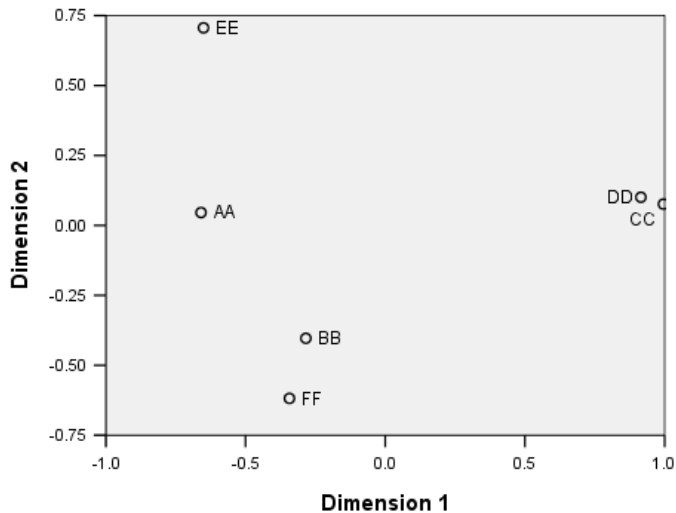
The row points plot shows that *fresh* and *ugly* are both very close to the origin, indicating that they differ little from the average row profile. Three general classifications emerge. Located in the upper left of the plot, *tough*, *men*, and *working* are all similar to each other. The lower left contains *sweet*, *fattening*, *children*, and *premium*. In contrast, *healthy*, *low fat*, *nutritious*, and *new* cluster on the right side of the plot.

Figure 11-31
Plot of image attributes (principal normalization)



Notice in the column points plot that all brands are far from the origin, so no brand is similar to the overall centroid. Brands *CC* and *DD* group together at the right, whereas brands *BB* and *FF* cluster in the lower half of the plot. Brands *AA* and *EE* are not similar to any other brand.

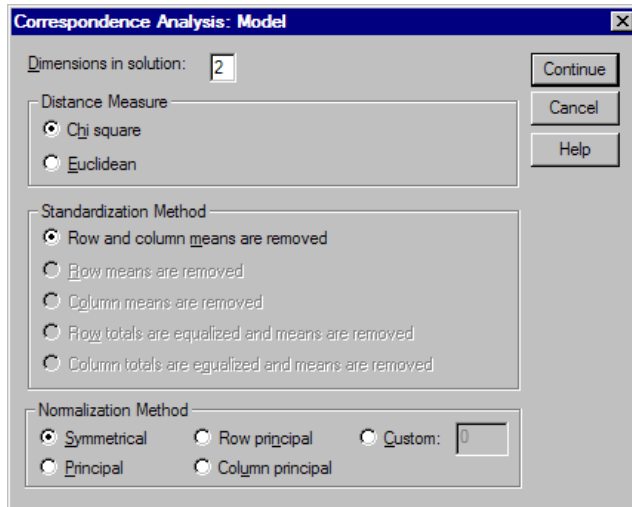
Figure 11-32
Plot of brands (principal normalization)



Symmetrical Normalization

How are the brands related to the image attributes? Principal normalization cannot address these relationships. To focus on how the variables are related to each other, use symmetrical normalization. Rather than spread the inertia twice (as in principal normalization), symmetrical normalization divides the inertia equally over both the rows and columns. Distances between categories for a single variable cannot be interpreted, but distances between the categories for different variables are meaningful.

Figure 11-33
Model dialog box

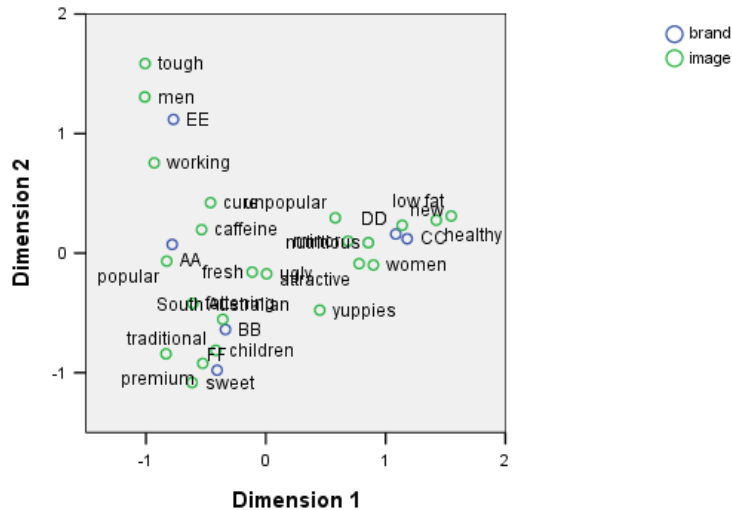


- ▶ To produce the following solution with symmetrical normalization, recall the Correspondence Analysis dialog box and click Model.
- ▶ Select Symmetrical as the normalization method.
- ▶ Click Continue.
- ▶ Click OK in the Correspondence Analysis dialog box.

In the upper left of the resulting biplot, brand *EE* is the only tough, working brand and appeals to men. Brand *AA* is the most popular and also viewed as the most highly caffeinated. The sweet, fattening brands include *BB* and *FF*. Brands *CC* and *DD*, while perceived as new and healthy, are also the most unpopular.

Figure 11-34

Biplot of the brands and the attributes (symmetrical normalization)



For further interpretation, you can draw a line through the origin and the two image attributes *men* and *yuppies*, and project the brands onto this line. The two attributes are opposed to each other, indicating that the association pattern of brands for *men* is reversed compared to the pattern for *yuppies*. That is, men are most frequently associated with brand *EE* and least frequently with brand *CC*, whereas yuppies are most frequently associated with brand *CC* and least frequently with brand *EE*.

Example: Flying Mileage between Cities

Correspondence analysis is not restricted to frequency tables. The entries can be any positive measure of correspondence. In this example, you use the flying mileages between ten American cities. This data set can be found in *flying.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

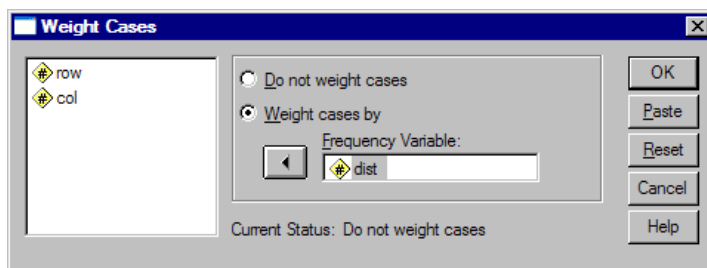
Table 11-2
City labels

City	Label	City	Label
Atlanta	Atl	Miami	Mia
Chicago	Chi	New York	NY
Denver	Den	San Francisco	SF
Houston	Hou	Seattle	Sea
Los Angeles	LA	Washington, DC	DC

- ▶ To view the flying mileages, first weight the cases by the variable *dist*. From the menus choose:

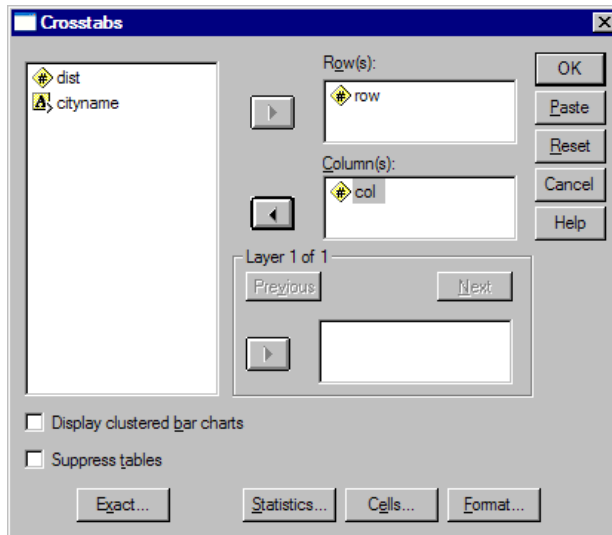
Data
Weight Cases...

Figure 11-35
Weight Cases dialog box



- ▶ Weight cases by *dist*.
- ▶ Click OK.
- ▶ Now, to view the mileages as a crosstabulation, from the menus choose:
Analyze
Descriptive Statistics
Crosstabs...

Figure 11-36
Crosstabs dialog box



- ▶ Select *row* as the row variable.
- ▶ Select *col* as the column variable.
- ▶ Click OK.

The following table contains the flying mileages between the cities. Notice that there is only one variable for both rows and columns and that the table is symmetric; the distance from Los Angeles to Miami is the same as the distance from Miami to Los

Angeles. Moreover, the distance between any city and itself is 0. The active margin reflects the total flying mileage from each city to all other cities.

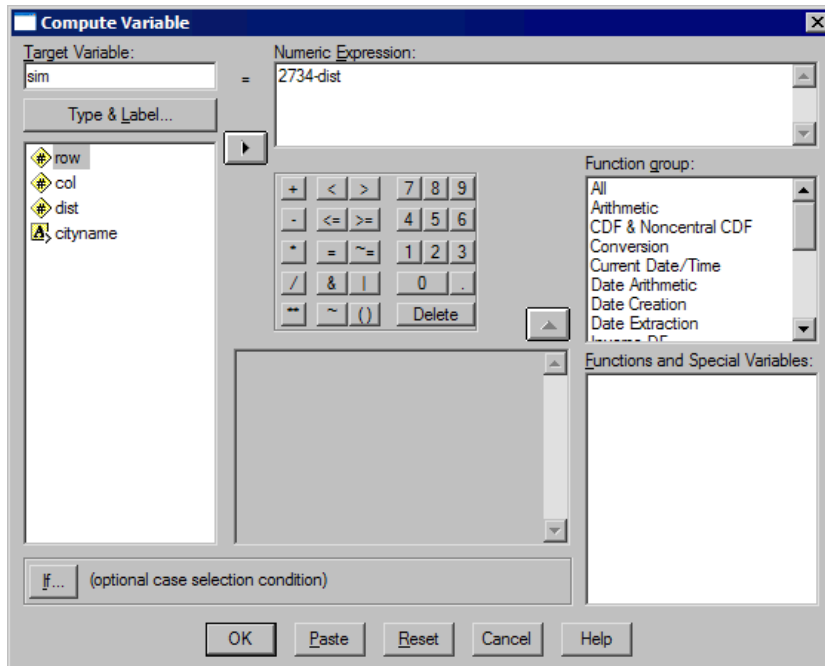
Figure 11-37
Flying mileages between 10 American cities

Count		col										Total
row	Atl	Chi	Den	Hou	LA	Mia	NY	SF	Sea	DC		
Atl	0	587	1212	701	1936	604	748	2139	2182	543	10652	
Chi	587	0	920	940	1745	1188	713	1858	1737	597	10285	
Den	1212	920	0	879	831	1726	1631	949	1021	1494	10663	
Hou	701	940	879	0	1374	968	1420	1645	1891	1220	11038	
LA	1936	1745	831	1374	0	2339	2451	347	959	2300	14282	
Mia	604	1188	1726	968	2339	0	1092	2594	2734	923	14168	
NY	748	713	1631	1420	2451	1092	0	2571	2408	205	13239	
SF	2139	1858	949	1645	347	2594	2571	0	678	2442	15223	
Sea	2182	1737	1021	1891	959	2734	2408	678	0	2329	15939	
DC	543	597	1494	1220	2300	923	205	2442	2329	0	12053	
Total	10652	10285	10663	11038	14282	14168	13239	15223	15939	12053	127542	

In general, distances are dissimilarities; large values indicate a large difference between the categories. However, correspondence analysis requires an association measure; thus, you need to convert dissimilarities into similarities. In other words, a large table entry must correspond to a small difference between the categories. Subtracting every table entry from the largest table entry converts the dissimilarities into similarities.

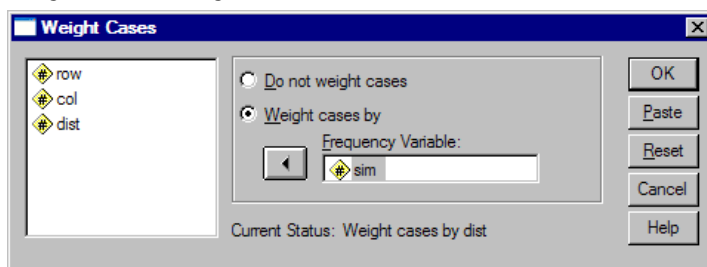
- ▶ To create the similarities and store them in a new variable, *sim*, from the menus choose:
Transform
Compute...

Figure 11-38
Compute Variable dialog box



- ▶ Type sim as the target variable.
- ▶ Type 2734-dist as the numeric expression.
- ▶ Click OK.

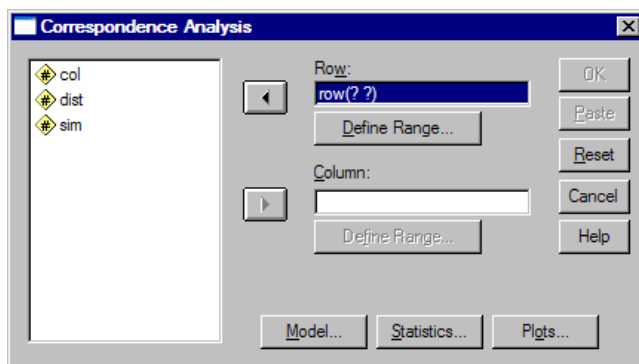
Figure 11-39
Weight Cases dialog box



Now reweight the cases by the similarity measure by recalling the Weight Cases dialog box:

- ▶ Weight cases by *sim*.
- ▶ Click OK.
- ▶ Finally, to obtain a correspondence analysis for the similarities, from the menus choose:
Analyze
Data Reduction
Correspondence Analysis...

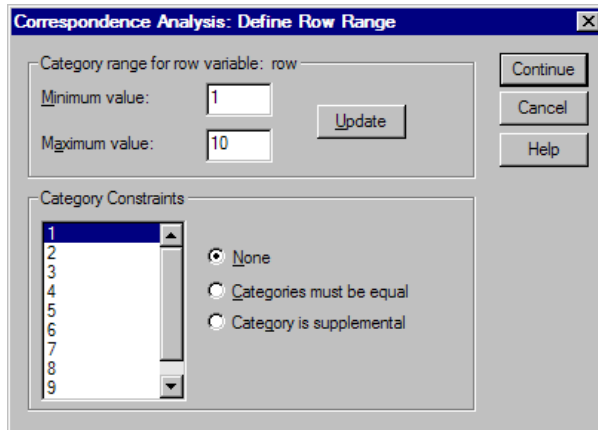
Figure 11-40
Correspondence Analysis dialog box



- ▶ Select *row* as the row variable.

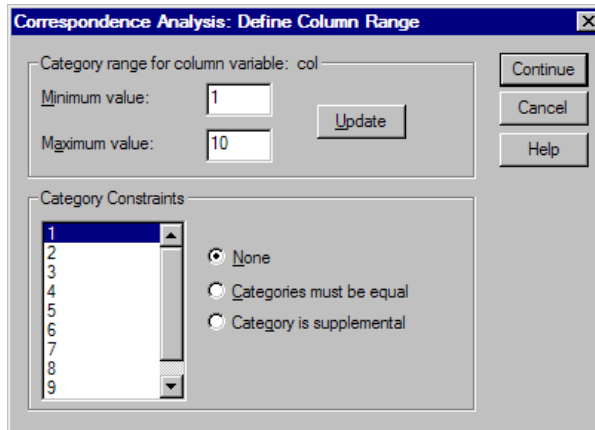
- ▶ Click Define Range.

Figure 11-41
Define Row Range dialog box



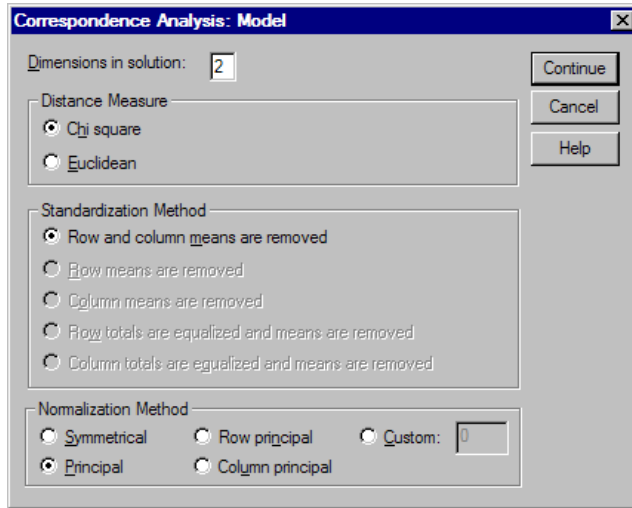
- ▶ Type 1 as the minimum value.
- ▶ Type 10 as the maximum value.
- ▶ Click Update.
- ▶ Click Continue.
- ▶ Select *col* as the column variable.
- ▶ Click Define Range in the Correspondence Analysis dialog box.

Figure 11-42
Define Column Range dialog box



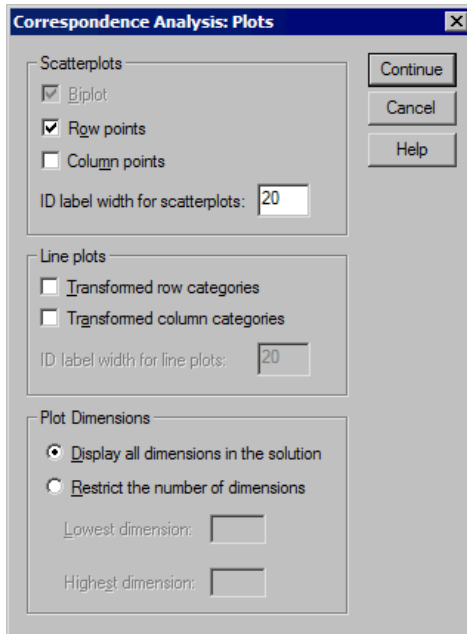
- ▶ Type 1 as the minimum value.
- ▶ Type 10 as the maximum value.
- ▶ Click Update.
- ▶ Click Continue.
- ▶ Click Model in the Correspondence Analysis dialog box.

Figure 11-43
Model dialog box



- ▶ Select Principal as the normalization method.
- ▶ Click Continue.
- ▶ Click Plots in the Correspondence Analysis dialog box.

Figure 11-44
Plots dialog box



- ▶ Select Row points in the Scatterplots group.
- ▶ Click Continue.
- ▶ Click OK in the Correspondence Analysis dialog box.

Correspondence Table

The new distance of 0 between Seattle and Miami indicates that they are most distant (least similar), whereas the distance of 2529 between New York and Washington, D.C., indicates that they are the least distant (most similar) pair of cities.

Figure 11-45
Correspondence table for similarities

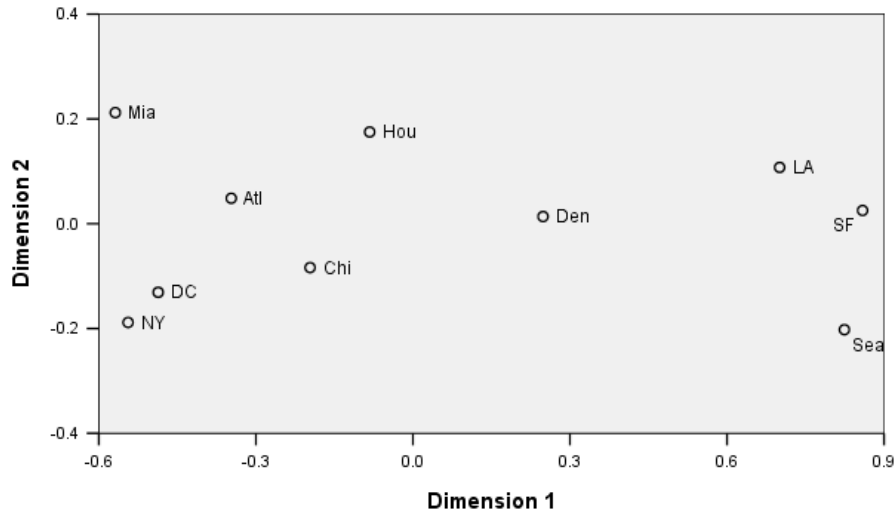
row	col										Active Margin
	Atl	Chi	Den	Hou	LA	Mia	NY	SF	Sea	DC	
Atl	2734	2147	1522	2033	798	2130	1986	595	552	2191	16688
Chi	2147	2734	1814	1794	989	1546	2021	876	997	2137	17055
Den	1522	1814	2734	1855	1903	1008	1103	1785	1713	1240	16677
Hou	2033	1794	1855	2734	1360	1766	1314	1089	843	1514	16302
LA	798	989	1903	1360	2734	395	283	2387	1775	434	13058
Mia	2130	1546	1008	1766	395	2734	1642	140	0	1811	13172
NY	1986	2021	1103	1314	283	1642	2734	163	326	2529	14101
SF	595	876	1785	1089	2387	140	163	2734	2056	292	12117
Sea	552	997	1713	843	1775	0	326	2056	2734	405	11401
DC	2191	2137	1240	1514	434	1811	2529	292	405	2734	15287
Active Margin	16688	17055	16677	16302	13058	13172	14101	12117	11401	15287	145858

Row and Column Scores

By using flying mileages instead of driving mileages, the terrain of the United States does not impact the distances. Consequently, all similarities should be representable in two dimensions. You center both the rows and columns and use principal normalization. Because of the symmetry of the correspondence table and

the principal normalization, the row and column scores are equal and the total inertia is in both, so it does not matter whether you inspect the row or column scores.

Figure 11-46
Points for 10 cities



The locations of the cities are very similar to their actual geographical locations, rotated about the origin. Cities that are further south have larger values along the second dimension, whereas cities that are further west have larger values along the first dimension.

Recommended Readings

See the following texts for more information on correspondence analysis:

Fisher, R. A. 1938. *Statistical methods for research workers*. Edinburgh: Oliver and Boyd.

Fisher, R. A. 1940. The precision of discriminant functions. *Annals of Eugenics*, 10, 422–429.

Gilula, Z., and S. J. Haberman. 1988. The analysis of multivariate contingency tables by restricted canonical and restricted association models. *Journal of the American Statistical Association*, 83, 760–771.

Multiple Correspondence Analysis

The purpose of multiple correspondence analysis, also known as homogeneity analysis, is to find quantifications that are optimal in the sense that the categories are separated from each other as much as possible. This implies that objects in the same category are plotted close to each other and objects in different categories are plotted as far apart as possible. The term **homogeneity** also refers to the fact that the analysis will be most successful when the variables are homogeneous; that is, when they partition the objects into clusters with the same or similar categories.

Example: Characteristics of Hardware

To explore how multiple correspondence analysis works, you will use data from Hartigan (Hartigan, 1975), which can be found in *screws.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS. This data set contains information on the characteristics of screws, bolts, nuts, and tacks. The following table shows the variables, along with their variable labels, and the value labels assigned to the categories of each variable in the Hartigan hardware data set.

Table 12-1
Hartigan hardware data set

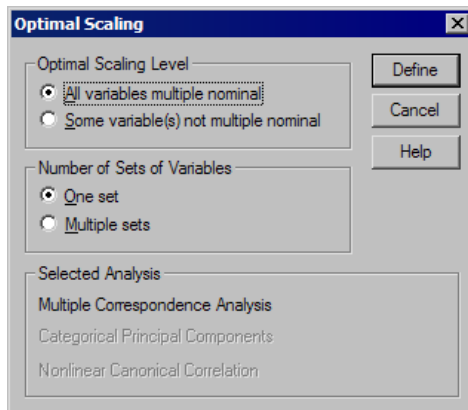
Variable name	Variable label	Value label
<i>thread</i>	Thread	Yes_Thread, No_Thread
<i>head</i>	Head form	Flat, Cup, Cone, Round, Cylinder
<i>indhead</i>	Indentation of head	None, Star, Slit
<i>bottom</i>	Bottom shape	sharp, flat
<i>length</i>	Length in half inches	1/2_in, 1_in, 1_1/2_in, 2_in, 2_1/2_in

Variable name	Variable label	Value label
<i>brass</i>	Brass	Yes_Br, Not_Br
<i>object</i>	Object	tack, nail1, nail2, nail3, nail4, nail5, nail6, nail7, nail8, screw1, screw2, screw3, screw4, screw5, bolt1, bolt2, bolt3, bolt4, bolt5, bolt6, tack1, tack2, nailb, screwb

Running the Analysis

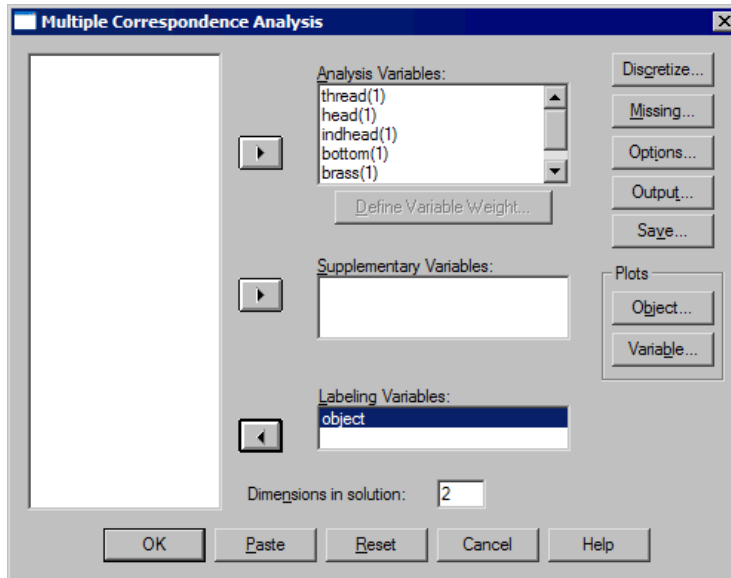
- ▶ To obtain a Multiple Correspondence Analysis, from the menus choose:
 Analyze
 Data Reduction
 Optimal Scaling...

Figure 12-1
Optimal Scaling dialog box



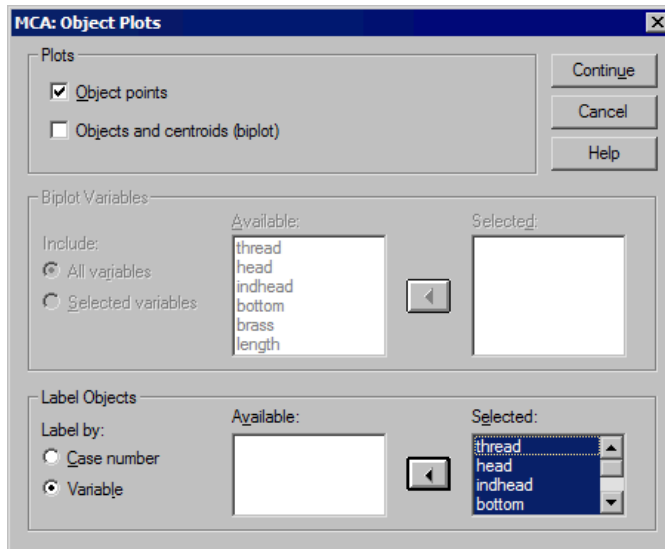
- ▶ Make sure All variables multiple nominal and One set are selected, and click Define.

Figure 12-2
Multiple Correspondence Analysis dialog box



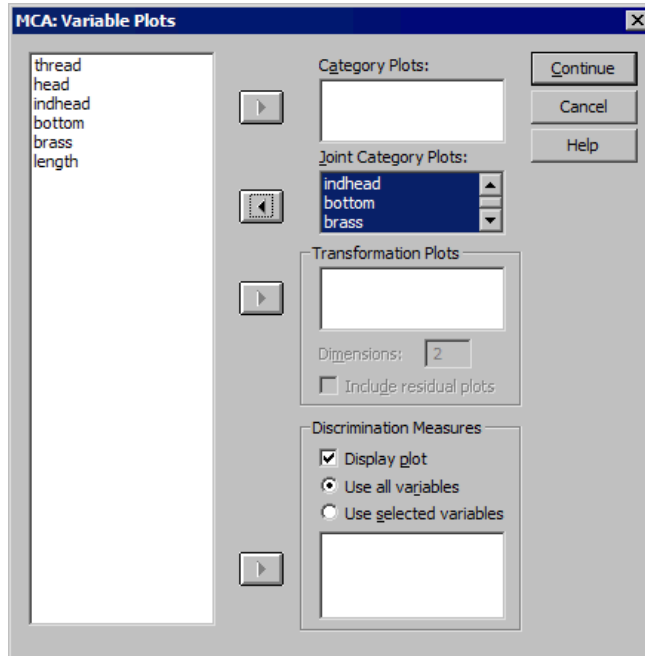
- ▶ Select *Thread* through *Length in half-inches* as analysis variables.
- ▶ Select *object* as a labeling variable.
- ▶ Click *Object* in the *Plots* group.

Figure 12-3
Object Plots dialog box



- ▶ Choose to label objects by Variable.
- ▶ Select *thread* through *object* as labeling variables.
- ▶ Click Continue, and then click Variable in the Plots group of the Multiple Correspondence Analysis dialog box.

Figure 12-4
Variable Plots dialog box



- ▶ Choose to produce a joint category plot for *thread* through *length*.
- ▶ Click Continue.
- ▶ Click OK in the Multiple Correspondence Analysis dialog box.

Model Summary

Homogeneity analysis can compute a solution for several dimensions. The maximum number of dimensions equals either the number of categories minus the number of variables with no missing data or the number of observations minus one, whichever is smaller. However, you should rarely use the maximum number of dimensions. A smaller number of dimensions is easier to interpret, and after a certain number of

dimensions, the amount of additional association accounted for becomes negligible. A one-, two-, or three-dimensional solution in homogeneity analysis is very common.

Figure 12-5
Model summary

Dimension	Cronbach's Alpha	Variance Accounted For		
		Total (Eigenvalue)	Inertia	% of Variance
1	.878	3.727	.621	62.123
2	.657	2.209	.368	36.809
Total		5.936	.989	
Mean	.796 ^a	2.968	.495	49.466

a. Mean Cronbach's Alpha is based on the mean Eigenvalue.

Nearly all of the variance in the data is accounted for by the solution, 62.1% by the first dimension and 36.8% by the second.

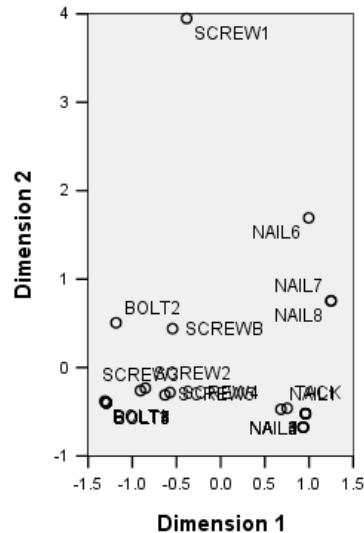
The two dimensions together provide an interpretation in terms of distances. If a variable discriminates well, the objects will be close to the categories to which they belong. Ideally, objects in the same category will be close to each other (that is, they should have similar scores), and categories of different variables will be close if they belong to the same objects (that is, two objects that have similar scores for one variable should also score close to each other for the other variables in the solution).

Object Scores

After examining the model summary, you should look at the object scores. You can specify one or more variables to label the object scores plot. Each labeling variable produces a separate plot labeled with the values of that variable. We'll take a look at the plot of object scores labeled by the variable object. This is just a case-identification variable and was not used in any computations.

The distance from an object to the origin reflects variation from the "average" response pattern. This average response pattern corresponds to the most frequent category for each variable. Objects with many characteristics corresponding to the most frequent categories lie near the origin. In contrast, objects with unique characteristics are located far from the origin.

Figure 12-6
Object scores plot labeled by object



Examining the plot, you see that the first dimension (the horizontal axis) discriminates the screws and bolts (which have threads) from the nails and tacks (which don't have threads). This is easily seen on the plot since screws and bolts are on one end of the horizontal axis and tacks and nails are on the other. To a lesser extent, the first dimension also separates the bolts (which have flat bottoms) from all the others (which have sharp bottoms).

The second dimension (the vertical axis) seems to separate *SCREW1* and *NAIL6* from all other objects. What *SCREW1* and *NAIL6* have in common are their values on variable length—they are the longest objects in the data. Moreover, *SCREW1* lies much farther from the origin than the other objects, suggesting that, taken as a whole, many of the characteristics of this object are not shared by the other objects.

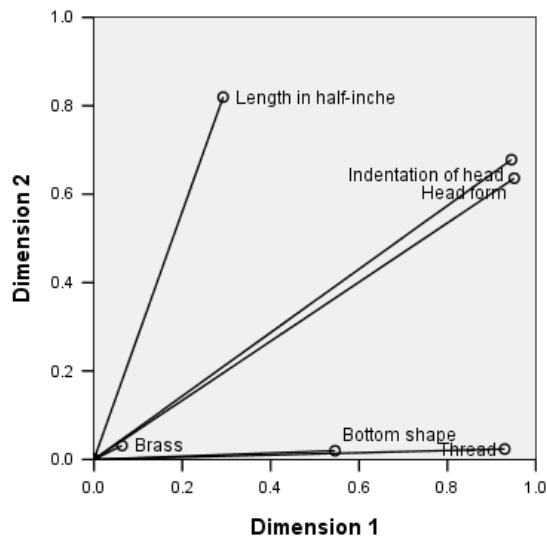
The object scores plot is particularly useful for spotting outliers. *SCREW1* might be considered an outlier. Later, we'll consider what happens if you drop this object.

Discrimination Measures

Before examining the rest of the object scores plots, let's see if the discrimination measures agree with what we've said so far. For each variable, a discrimination measure, which can be regarded as a squared component loading, is computed for each dimension. This measure is also the variance of the quantified variable in that dimension. It has a maximum value of 1, which is achieved if the object scores fall into mutually exclusive groups and all object scores within a category are identical. (*Note:* This measure may have a value greater than 1 if there are missing data.) Large discrimination measures correspond to a large spread among the categories of the variable and, consequently, indicate a high degree of discrimination between the categories of a variable along that dimension.

The average of the discrimination measures for any dimension equals the percentage of variance accounted for that dimension. Consequently, the dimensions are ordered according to average discrimination. The first dimension has the largest average discrimination, the second dimension has the second largest average discrimination, and so on, for all dimensions in the solution.

Figure 12-7
Plot of discrimination measures



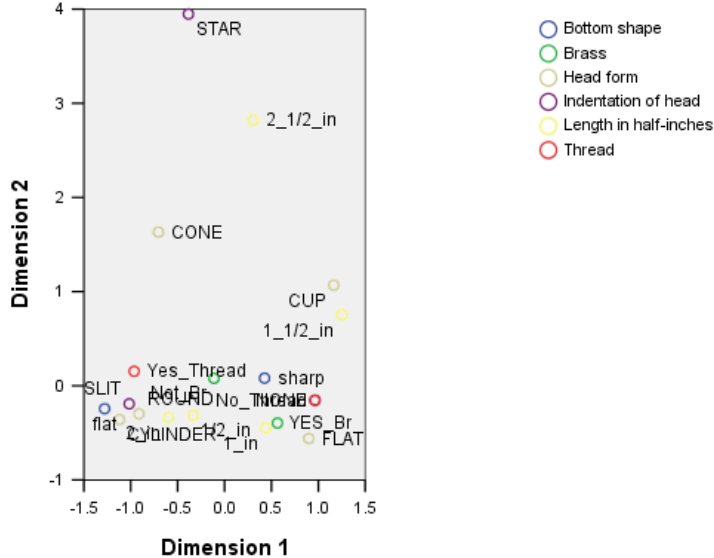
As noted on the object scores plot, the discrimination measures plot shows that the first dimension is related to variables *Thread* and *Bottom shape*. These variables have large discrimination measures on the first dimension and small discrimination measures on the second dimension. Thus, for both of these variables, the categories are spread far apart along the first dimension only. *Length in half-inches* has a large value on the second dimension but a small value on the first dimension. As a result, *length* is closest to the second dimension, agreeing with the observation from the object scores plot that the second dimension seems to separate the longest objects from the rest. *Indentation of head* and *Head form* have relatively large values on both dimensions, indicating discrimination in both the first and second dimensions. The variable *Brass*, located very close to the origin, does not discriminate at all in the first two dimensions. This makes sense, since all of the objects can be made of brass or not made of brass.

Category Quantifications

Recall that a discrimination measure is the variance of the quantified variable along a particular dimension. The discrimination measures plot contains these variances, indicating which variables discriminate along which dimension. However, the same variance could correspond to all of the categories being spread moderately far apart or to most of the categories being close together, with a few categories differing from this group. The discrimination plot cannot differentiate between these two conditions.

Category quantification plots provide an alternative method of displaying discrimination of variables that can identify category relationships. In this plot, the coordinates of each category on each dimension are displayed. Thus, you can determine which categories are similar for each variable.

Figure 12-8
Category quantifications



Length in half-inches has five categories, three of which group together near the top of the plot. The remaining two categories are in the lower half of the plot, with the *2_1/2_in* category very far from the group. The large discrimination for length along dimension 2 is a result of this one category being very different from the other categories of length. Similarly, for *Head form*, the category *STAR* is very far from the other categories and yields a large discrimination measure along the second dimension. These patterns cannot be illustrated in a plot of discrimination measures.

The spread of the category quantifications for a variable reflects the variance and thus indicates how well that variable is discriminated in each dimension. Focusing on dimension 1, the categories for *Thread* are far apart. However, along dimension 2, the categories for this variable are very close. Thus, *Thread* discriminates better in dimension 1 than in dimension 2. In contrast, the categories for *Head form* are spread far apart along both dimensions, suggesting that this variable discriminates well in both dimensions.

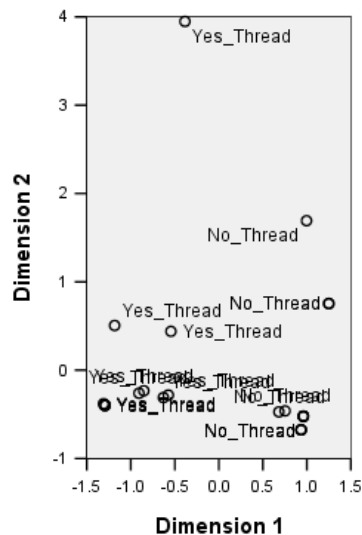
In addition to determining the dimensions along which a variable discriminates and how that variable discriminates, the category quantification plot also compares variable discrimination. A variable with categories that are far apart discriminates better than a variable with categories that are close together. For example, along

dimension 1, the two categories of *Brass* are much closer to each other than the two categories of *Thread*, indicating that *Thread* discriminates better than *Brass* along this dimension. However, along dimension 2, the distances are very similar, suggesting that these variables discriminate to the same degree along this dimension. The discrimination measures plot discussed previously identifies these same relationships by using variances to reflect the spread of the categories.

A More Detailed Look at Object Scores

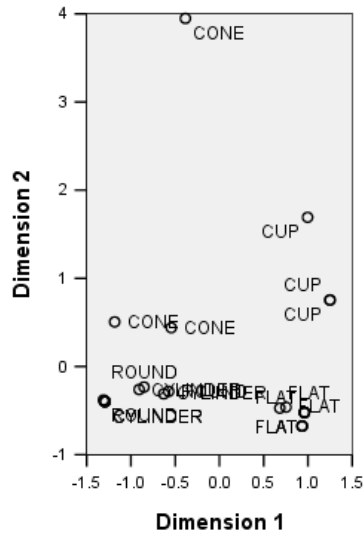
A greater insight into the data can be gained by examining the object scores plots labeled by each variable. Ideally, similar objects should form exclusive groups, and these groups should be far from each other.

Figure 12-9
Object scores labeled with Thread



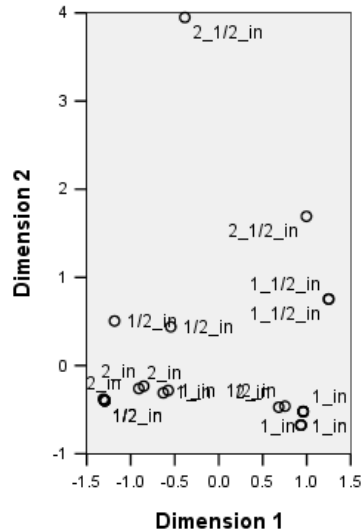
The plot labeled with *Thread* shows that the first dimension separates *Yes_Thread* and *No_Thread* perfectly. All of the objects with threads have negative object scores, whereas all of the nonthreaded objects have positive scores. Although the two categories do not form compact groups, the perfect differentiation between the categories is generally considered a good result.

Figure 12-10
Object scores labeled with Head form



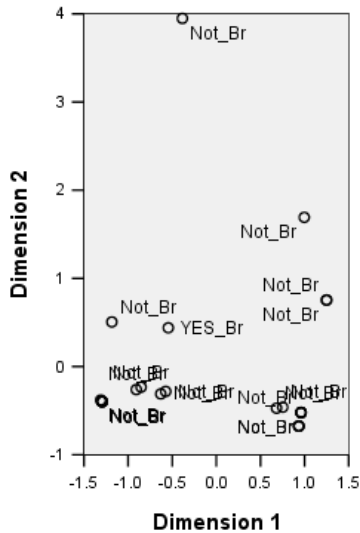
The plot labeled with *Head form* shows that this variable discriminates in both dimensions. The *FLAT* objects group together in the lower right corner of the plot, whereas the *CUP* objects group together in the upper right. *CONE* objects all lie in the upper left. However, these objects are more spread out than the other groups and, thus, are not as homogeneous. Finally, *CYLINDER* objects cannot be separated from *ROUND* objects, both of which lie in the lower left corner of the plot.

Figure 12-11
Object scores labeled with Length in half-inches



The plot labeled with *Length in half-inches* shows that this variable does not discriminate in the first dimension. Its categories display no grouping when projected onto a horizontal line. However, *Length in half-inches* does discriminate in the second dimension. The shorter objects correspond to positive scores, and the longer objects correspond to large negative scores.

Figure 12-12
Object scores labeled with Brass



The plot labeled with *Brass* shows that this variable has categories that can't be separated very well in the first or second dimensions. The object scores are widely spread throughout the space. The brass objects cannot be differentiated from the nonbrass objects.

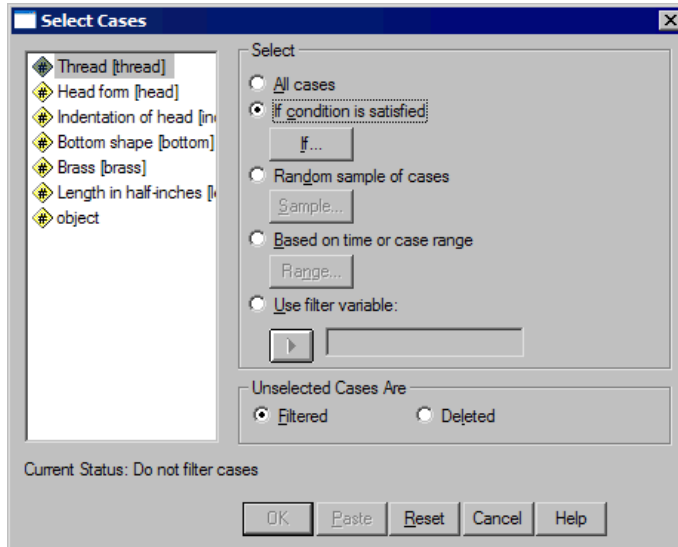
Omission of Outliers

In homogeneity analysis, outliers are objects that have too many unique features. As noted earlier, *SCREW1* might be considered an outlier.

To delete this object and run the analysis again, from the menus choose:

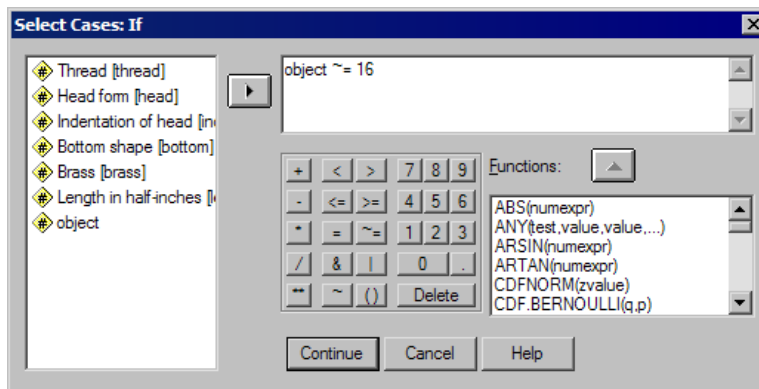
Data
Select Cases...

Figure 12-13
Select Cases dialog box



- ▶ Select If condition is satisfied.
- ▶ Click If.

Figure 12-14
If dialog box



- ▶ Type object ~= 16 as the condition.

- ▶ Click Continue.
- ▶ Click OK in the Select Cases dialog box.
- ▶ Finally, recall the Multiple Correspondence Analysis dialog box, and click OK.

Figure 12-15

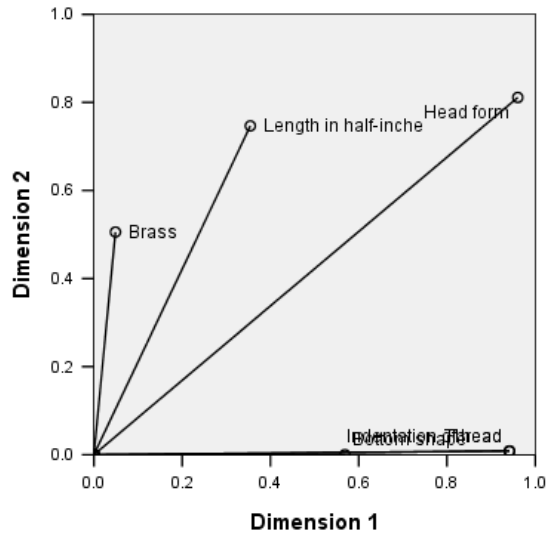
Model summary (outlier removed)

Dimension	Cronbach's Alpha	Variance Accounted For		
		Total (Eigenvalue)	Inertia	% of Variance
1	.885	3.815	.636	63.591
2	.623	2.081	.347	34.676
Total		5.896	.983	
Mean	.793 ^a	2.948	.491	49.133

^a. Mean Cronbach's Alpha is based on the mean Eigenvalue.

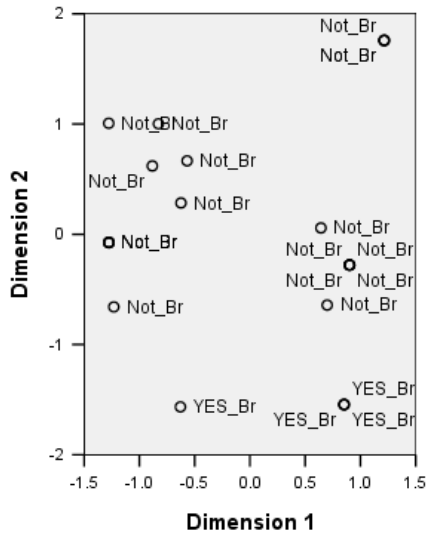
The eigenvalues shift slightly. The first dimension now accounts for a little more of the variance.

Figure 12-16
Discrimination measures



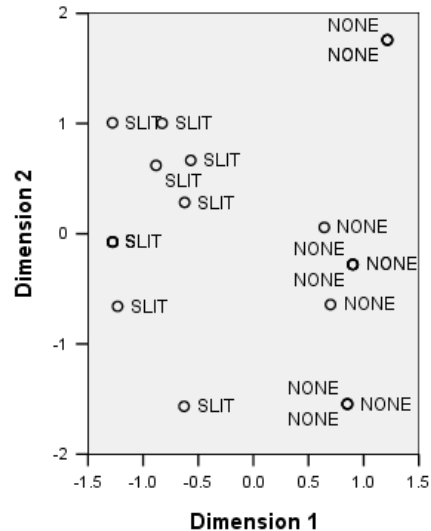
As shown in the discrimination plot, *Indentation of head* no longer discriminates in the second dimension, whereas *Brass* changes from no discrimination in either dimension to discrimination in the second dimension. Discrimination for the other variables is largely unchanged.

Figure 12-17
Object scores labeled with *Brass* (outlier removed)



The object scores plot labeled by *Brass* shows that the four brass objects all appear near the bottom of the plot (three objects occupy identical locations), indicating high discrimination along the second dimension. As was the case for *Thread* in the previous analysis, the objects do not form compact groups, but the differentiation of objects by categories is perfect.

Figure 12-18
Object scores labeled with *Indentation of head* (outlier removed)



The object scores plot labeled by *Indentation of head* shows that the first dimension discriminates perfectly between the non-indented objects and the indented objects, as in the previous analysis. In contrast to the previous analysis, however, the second dimension cannot now distinguish the two categories.

Thus, the omission of *SCREW1*, which is the only object with a star-shaped head, dramatically affects the interpretation of the second dimension. This dimension now differentiates objects based on *Brass*, *Head form*, and *Length in half-inches*.

Recommended Readings

See the following texts for more information on multiple correspondence analysis:

Benzécri, J. P. 1992. *Correspondence analysis handbook*. New York: Marcel Dekker.

Guttman, L. 1941. The quantification of a class of attributes: A theory and method of scale construction. In: *The prediction of personal adjustment*, P. Horst, ed. New York: Social Science Research Council, 319–348.

Meulman, J. 1982. *Homogeneity analysis of incomplete data*. Leiden: DSWO Press.

Meulman, J. 1996. Fitting a distance model to homogeneous subsets of variables: Points of view analysis of categorical data. *Journal of Classification*, 13, 249–266.

Meulman, J., and W. J. Heiser. 1997. Graphical display of interaction in multiway contingency tables by use of homogeneity analysis. In: *Visual display of categorical data*, M. Greenacre, and J. Blasius, eds. New York: Academic Press, 277–296.

Nishisato, S. 1984. Forced classification: A simple application of a quantification method. *Psychometrika*, 49, 25–36.

Tenenhaus, M., and F. W. Young. 1985. An analysis and synthesis of multiple correspondence analysis, optimal scaling, dual scaling, homogeneity analysis, and other methods for quantifying categorical multivariate data. *Psychometrika*, 50, 91–119.

Van Rijkevorsel, J. 1987. *The application of fuzzy coding and horseshoes in multiple correspondence analysis*. Leiden: DSWO Press.

Barlow, R. E., D. J. Bartholomew, D. J. Bremner, and H. D. Brunk. 1972. *Statistical inference under order restrictions*. New York: John Wiley and Sons.

Multidimensional Scaling

Given a set of objects, the goal of multidimensional scaling is to find a representation of the objects in a low-dimensional space. This solution is found using the **proximities** between the objects. The procedure minimizes the squared deviations between the original, possibly transformed, object proximities and their Euclidean distances in the low-dimensional space.

The purpose of the low-dimensional space is to uncover relationships between the objects. By restricting the solution to be a linear combination of independent variables, you may be able to interpret the dimensions of the solution in terms of these variables. In the following example, you will see how 15 different kinship terms can be represented in three dimensions and how that space can be interpreted with respect to the gender, generation, and degree of separation of each of the terms.

Example: An Examination of Kinship Terms

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 criteria from the first sort. Thus, a total of six “sources” were obtained, as outlined in the following table.

Table 13-1
Source structure of kinship data

Source	Gender	Condition	Sample size
1	Female	Single sort	85
2	Male	Single sort	85

Source	Gender	Condition	Sample size
3	Female	First sort	80
4	Female	Second sort	80
5	Male	First sort	80
6	Male	Second sort	80

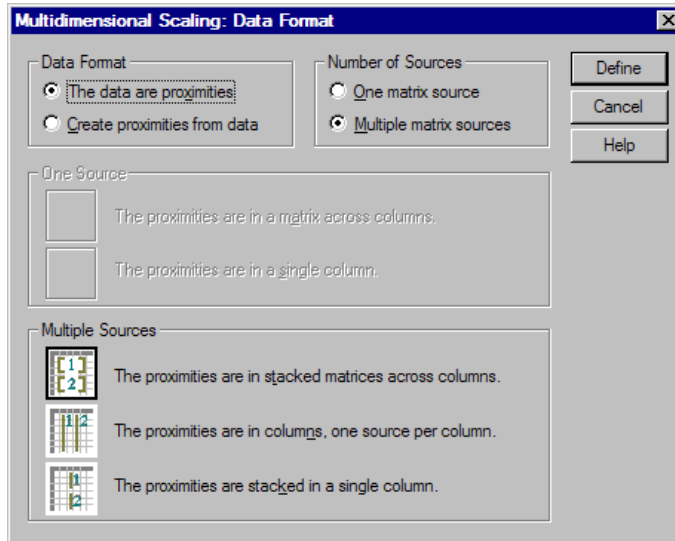
Each source corresponds to a 15×15 proximity matrix, whose cells are equal to the number of people in a source minus the number of times the objects were partitioned together in that source. This data set can be found in *kinship_dat.sav*, located in the `\tutorial\sample_files\` subdirectory of the directory in which you installed SPSS.

Choosing the Number of Dimensions

It is up to you to decide how many dimensions the solution should have. A good tool to help you make this decision is the scree plot. To create a scree plot, from the menus choose:

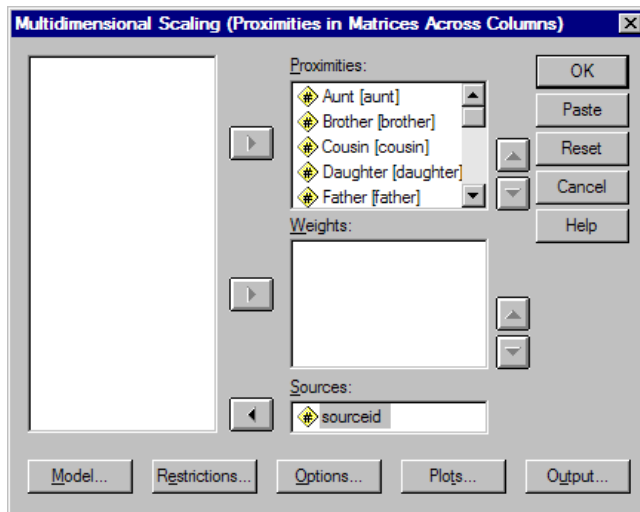
```
Analyze  
  Scale  
    Multidimensional Scaling (PROXSCAL)...
```

Figure 13-1
Data Format dialog box



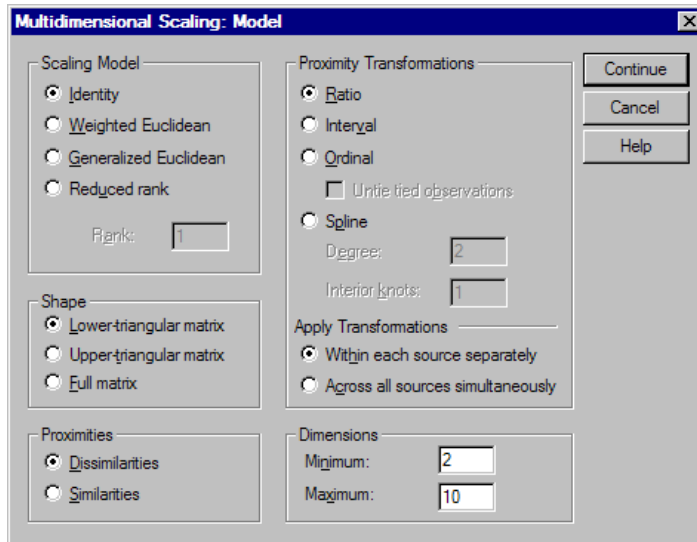
- ▶ Select Multiple matrix sources in the Number of Sources group.
- ▶ Click Define.

Figure 13-2
Multidimensional Scaling dialog box



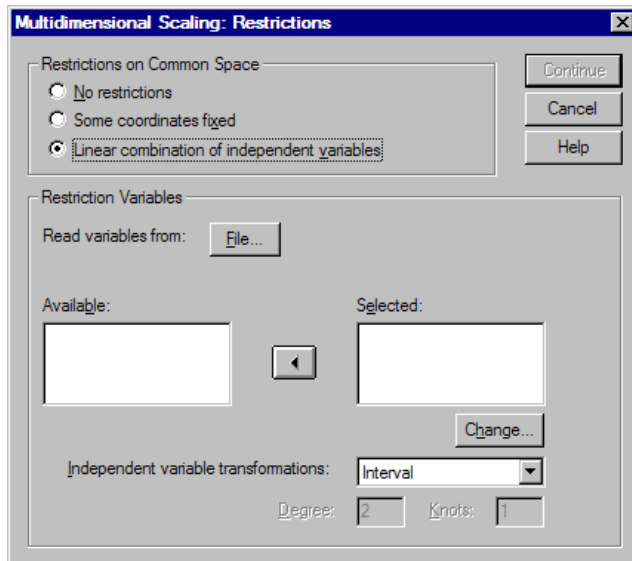
- ▶ Select *Aunt* through *Uncle* as proximities variables.
- ▶ Select *sourceid* as the variable identifying the source.
- ▶ Click Model.

Figure 13-3
Model dialog box



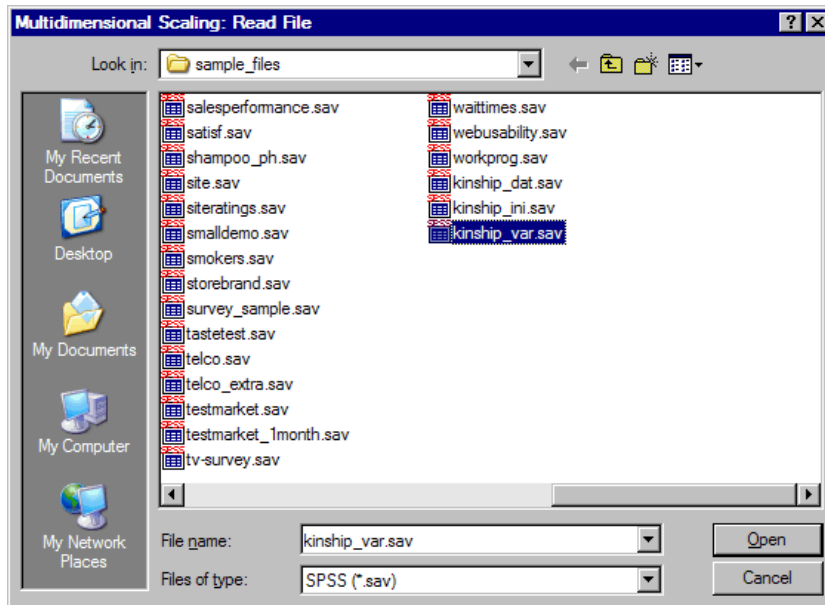
- ▶ Type 10 as the maximum number of dimensions.
- ▶ Click Continue.
- ▶ Click Restrictions in the Multidimensional Scaling dialog box.

Figure 13-4
Restrictions dialog box



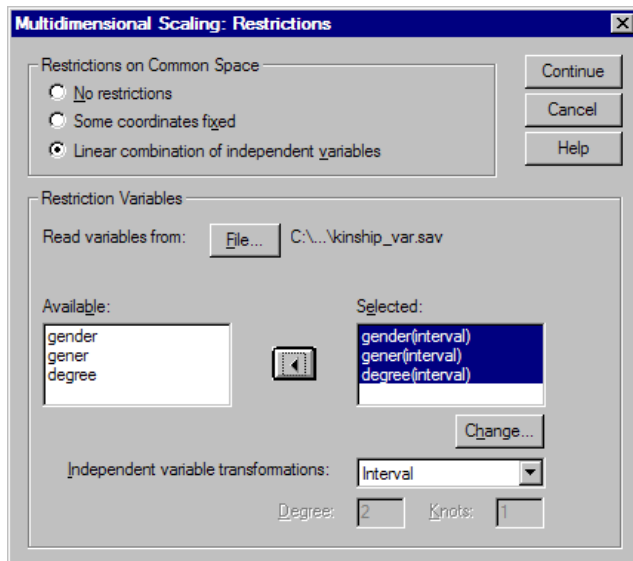
- ▶ Select Linear combination of independent variables.
- ▶ Click File to select the source of the independent variables.

Figure 13-5
Read File dialog box



- ▶ Select *kinship_var.sav*.
- ▶ Click Open.

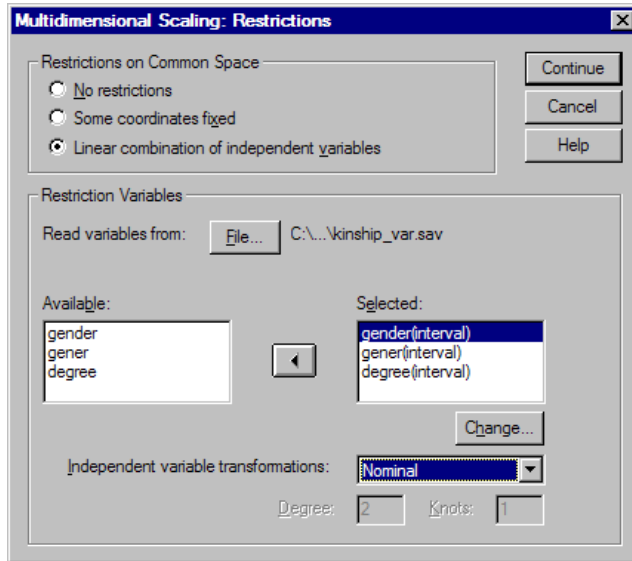
Figure 13-6
Restrictions dialog box



- Select *gender*, *gener*, and *degree* as restriction variables.

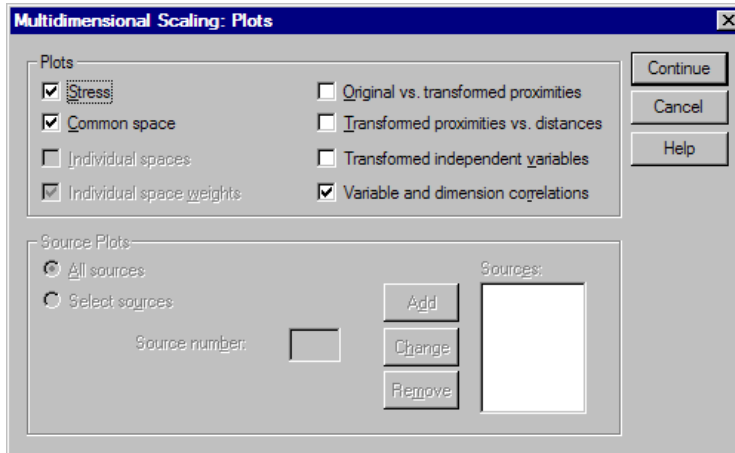
Note that the variable *gender* has a user-missing value—9 = missing (for cousin). The procedure treats this as a valid category. Thus, the default linear transformation is unlikely to be appropriate. Use a nominal transformation instead.

Figure 13-7
Restrictions dialog box



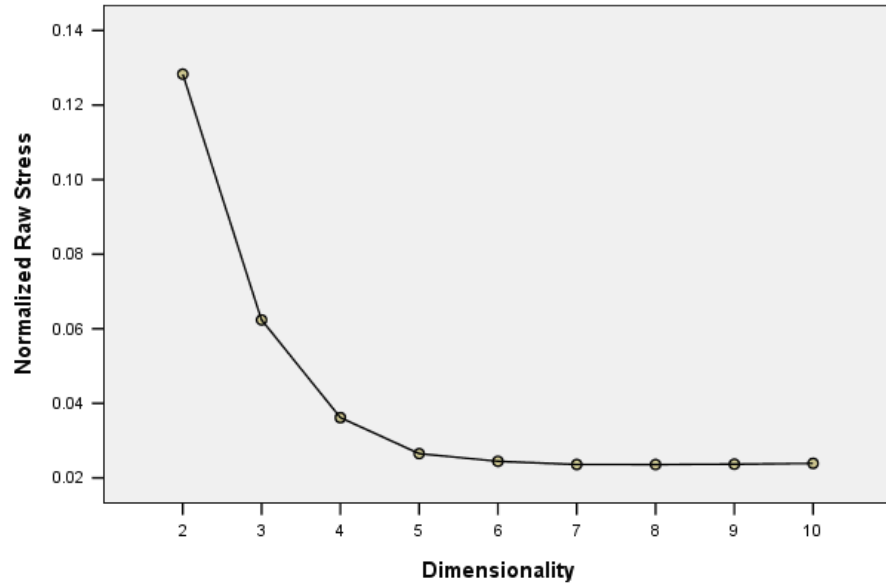
- ▶ Select *gender*.
- ▶ Select Nominal from the Independent variable transformations drop-down list.
- ▶ Click Change.
- ▶ Click Continue.
- ▶ Click Plots in the Multidimensional Scaling dialog box.

Figure 13-8
Plots dialog box



- ▶ Select Stress in the Plots group.
- ▶ Click Continue.
- ▶ Click OK in the Multidimensional Scaling dialog box.

Figure 13-9
Scree Plot



The procedure begins with a 10-dimensional solution and works down to a 2-dimensional solution. The scree plot shows the normalized raw stress of the solution at each dimension. You can see from the plot that increasing the dimensionality from 2 to 3 and from 3 to 4 offers large improvements in the stress. After 4, the improvements are rather small. You will choose to analyze the data using a 3-dimensional solution, since the results are easier to interpret.

A Three-Dimensional Solution

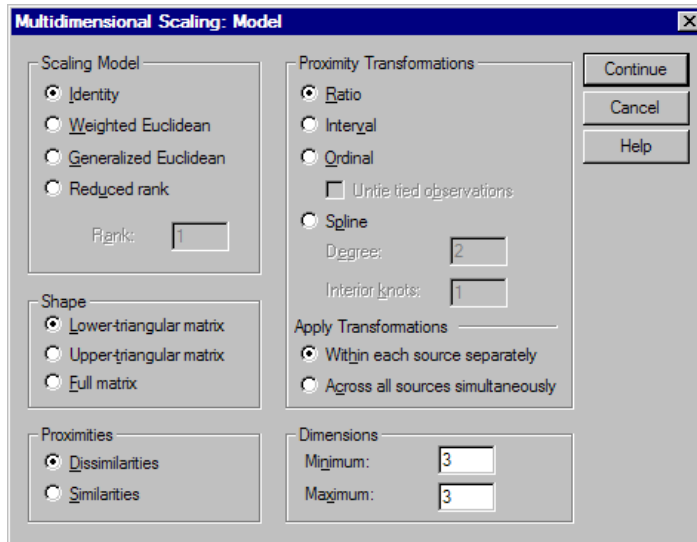
The independent variables *gender*, *gener* (generation), and *degree* (of separation) were constructed with the intention of using them to interpret the dimensions of the solution. The independent variables were constructed as follows:

<i>gender</i>	1 = male, 2 = female, 9 = missing (for cousin)
<i>gener</i>	The number of generations from you if the term refers to your kin, with lower numbers corresponding to older generations. Thus, grandparents are -2, grandchildren are 2, and siblings are 0.
<i>degree</i>	The number of degrees of separation along your family tree. Thus, your parents are up 1 node, while your children are down 1 node. Your siblings are up 1 node to your parents and then down 1 node to them, for 2 degrees of separation. Your cousin is 4 degrees away—2 up to your grandparents and then 2 down through your aunt/uncle to them.

The external variables can be found in *kinship_var.sav*. Additionally, an initial configuration from an earlier analysis is supplied in *kinship_ini.sav*.

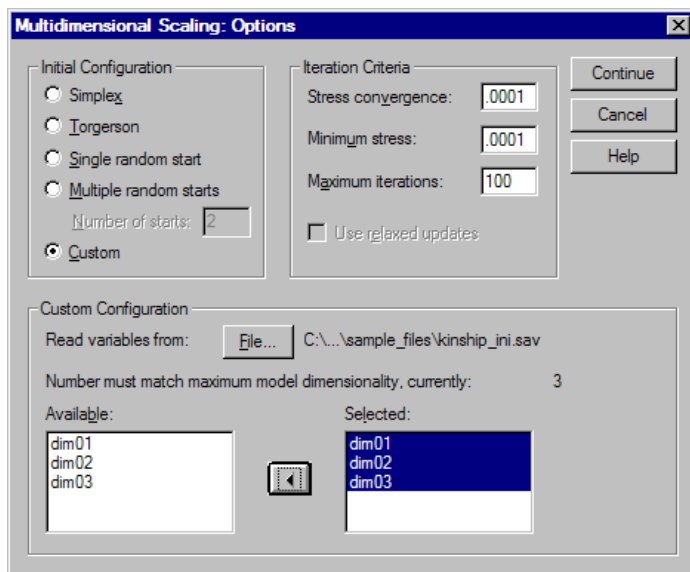
Running the Analysis

Figure 13-10
Model dialog box



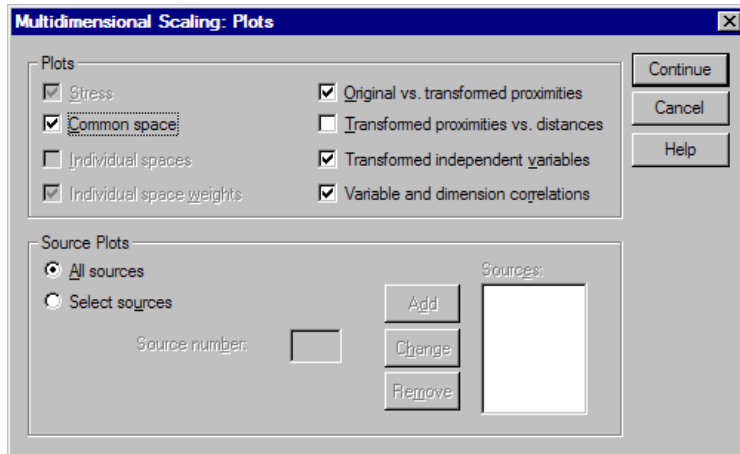
- ▶ To obtain a three-dimensional solution, recall the Multidimensional Scaling dialog box and click Model.
- ▶ Type 3 as the minimum and maximum number of dimensions.
- ▶ Click Continue.
- ▶ Click Options in the Multidimensional Scaling dialog box.

Figure 13-11
Options dialog box



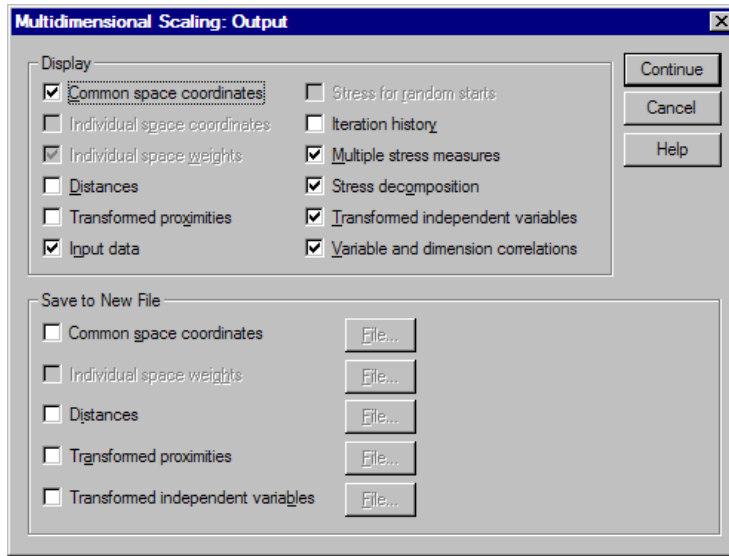
- ▶ Select Custom as the initial configuration.
- ▶ Select *kinship_ini.sav* as the file to read variables from.
- ▶ Select *dim01*, *dim02*, and *dim03* as variables.
- ▶ Click Continue.
- ▶ Click Plots in the Multidimensional Scaling dialog box.

Figure 13-12
Plots dialog box



- ▶ Select Original vs. transformed proximities and Transformed independent variables.
- ▶ Click Continue.
- ▶ Click Output in the Multidimensional Scaling dialog box.

Figure 13-13
Output dialog box



- ▶ Select Input data, Stress decomposition, and Variable and dimension correlations.
- ▶ Click Continue.
- ▶ Click OK in the Multidimensional Scaling dialog box.

Stress Measures

The stress and fit measures give an indication of how well the distances in the solution approximate the original distances.

Figure 13-14
Stress and fit measures

Normalized Raw Stress	.06234
Stress-I	.24968 ^a
Stress-II	.87849 ^a
S-Stress	.14716 ^b
Dispersion Accounted For (D.A.F.)	.93766
Tucker's Coefficient of Congruence	.96833

PROXSCAL minimizes Normalized Raw Stress.

a. Optimal scaling factor = 1.066.

b. Optimal scaling factor = .984.

Each of the four stress statistics measures the misfit of the data, while the dispersion accounted for and Tucker's coefficient of congruence measure the fit. Lower stress measures (to a minimum of 0) and higher fit measures (to a maximum of 1) indicate better solutions.

Figure 13-15
Decomposition of normalized raw stress

		Source						Mean
		SRC_1	SRC_2	SRC_3	SRC_4	SRC_5	SRC_6	
Object	Aunt	.0991	.0754	.0829	.0488	.0391	.0489	.0620
	Brother	.1351	.0974	.0496	.0813	.0613	.0597	.0807
	Cousin	.0325	.0336	.0480	.0290	.0327	.0463	.0370
	Daughter	.0700	.0370	.0516	.0229	.0326	.0207	.0391
	Father	.0751	.0482	.0521	.0225	.0272	.0298	.0425
	Granddaughter	.1410	.0736	.0801	.0707	.0790	.0366	.0802
	Grandfather	.1549	.1057	.0858	.0821	.0851	.0576	.0952
	Grandmother	.1550	.0979	.0858	.0844	.0816	.0627	.0946
	Grandson	.1374	.0772	.0793	.0719	.0791	.0382	.0805
	Mother	.0813	.0482	.0526	.0229	.0260	.0227	.0423
	Nephew	.0843	.0619	.0580	.0375	.0317	.0273	.0501
	Niece	.0850	.0577	.0503	.0353	.0337	.0260	.0480
	Sister	.1361	.0946	.0496	.0816	.0629	.0588	.0806
	Son	.0689	.0373	.0456	.0242	.0337	.0253	.0392
Uncle	.0977	.0761	.0678	.0489	.0383	.0498	.0631	
Mean		.1035	.0681	.0613	.0508	.0496	.0407	.0623

The decomposition of stress helps you to identify which sources and objects contribute the most to the overall stress of the solution. In this case, most of the stress among the sources is attributable to sources 1 and 2, while among the objects, most

of the stress is attributable to *Brother, Granddaughter, Grandfather, Grandmother, Grandson, and Sister*.

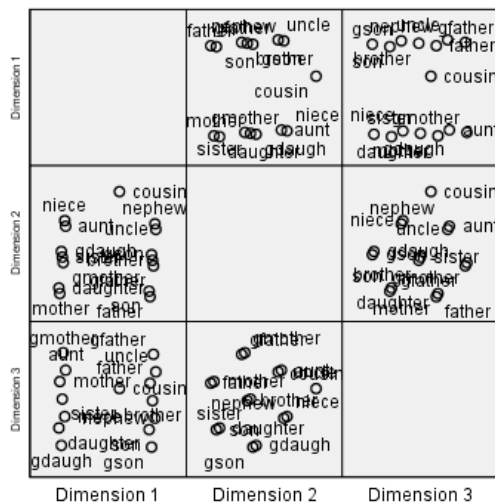
The two sources accountable for most of the stress are the two groups that sorted the terms only once. This suggests that the students considered multiple factors when sorting the terms, and those who were allowed to sort twice focused on a portion of those factors for the first sort and then considered the remaining factors during the second sort.

The objects that account for most of the stress are those with a *degree* of 2. These are relations who are not part of the “nuclear” family (*Mother, Father, Daughter, Son*), but are nonetheless closer than other relations. This middle position could easily cause some differential sorting of these terms.

Final Coordinates of the Common Space

The common space plot gives a visual representation of the relationships between the objects.

Figure 13-16
Common space coordinates



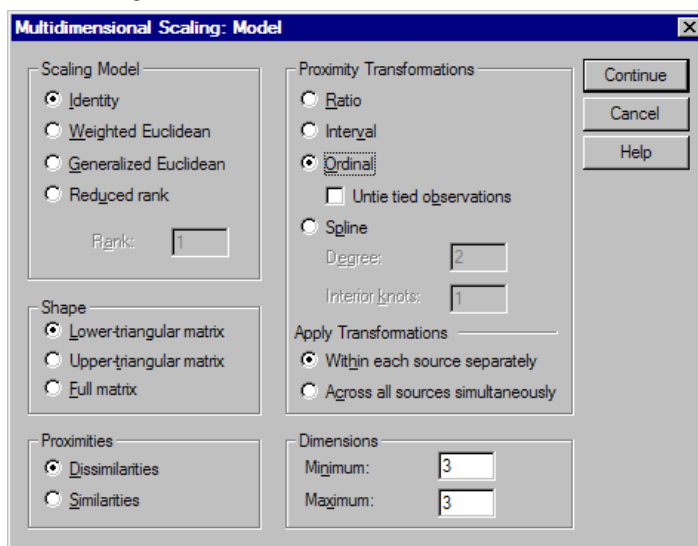
Look at the final coordinates for the objects in dimensions 1 and 3; this is the plot in the lower left corner of the scatterplot matrix. This plot shows that dimension 1 (on the x axis) is correlated with the variable *gender* and dimension 3 (on the y axis) is correlated with *gener*. From left to right, you see that dimension 1 separates the female and male terms, with the genderless term *Cousin* in the middle. From the bottom of the plot to the top, increasing values along the axis correspond to terms that are older.

Now look at the final coordinates for the objects in dimensions 2 and 3; this is the plot on the middle right side of the scatterplot matrix. From this plot, you can see that the second dimension (along the y axis) corresponds to the variable *degree*, with larger values along the axis corresponding to terms that are further from the “nuclear” family.

A Three-Dimensional Solution with Nondefault Transformations

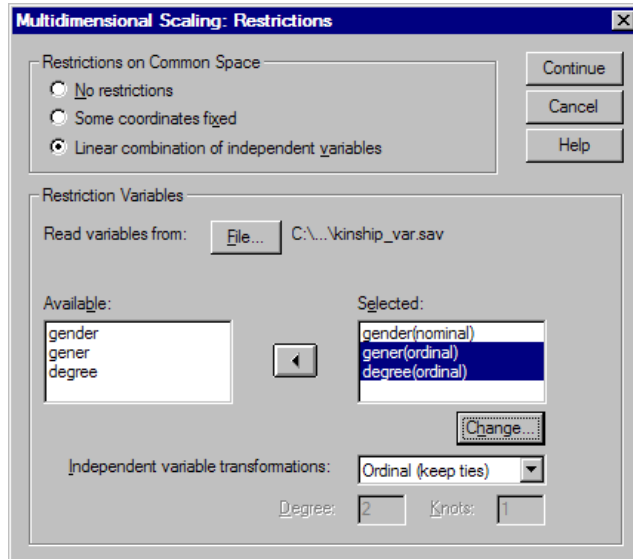
The previous solution was computed using the default ratio transformation for proximities and interval transformations for the independent variables *gener* and *degree*. The results are pretty good, but you may be able to do better by using other transformations. For example, the proximities, *gener*, and *degree* all have natural orderings, but they may be better modeled by an ordinal transformation than a linear transformation.

Figure 13-17
Model dialog box



- ▶ To rerun the analysis, scaling the proximities, *gener*, and *degree* at the ordinal level (keeping ties), recall the Multidimensional Scaling dialog box and click Model.
- ▶ Select Ordinal as the proximity transformation.
- ▶ Click Continue.
- ▶ Click Restrictions in the Multidimensional Scaling dialog box.

Figure 13-18
Restrictions dialog box



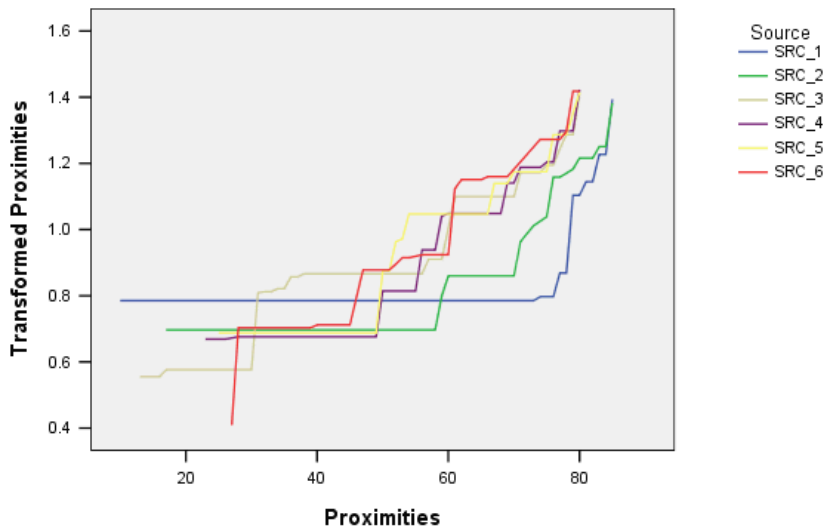
- ▶ Select *gener* and *degree*.
- ▶ Select Ordinal (keep ties) from the Independent variable transformations drop-down list.
- ▶ Click Change.
- ▶ Click Continue.
- ▶ Click OK in the Multidimensional Scaling dialog box.

Transformation Plots

The transformation plots are a good first check to see whether the original transformations were appropriate. If the plots are approximately linear, then the linear assumption is appropriate. If not, then you need to check the stress measures to see if there is an improvement in fit and the common space plot to see if the interpretation is more useful.

The independent variables each obtain approximately linear transformations, so it may be appropriate to interpret them as numerical. However, the proximities do not obtain a linear transformation, so it is possible that the ordinal transformation is more appropriate for the proximities.

Figure 13-19
Transformed proximities



Stress Measures

The stress for the current solution supports the argument for scaling the proximities at the ordinal level.

Figure 13-20
Stress and fit measures

Normalized Raw Stress	.03137
Stress-I	.17712 ^a
Stress-II	.61987 ^a
S-Stress	.07953 ^b
Dispersion Accounted For (D.A.F.)	.96863
Tucker's Coefficient of Congruence	.98419

PROXSAL minimizes Normalized Raw Stress.

a. Optimal scaling factor = 1.032.

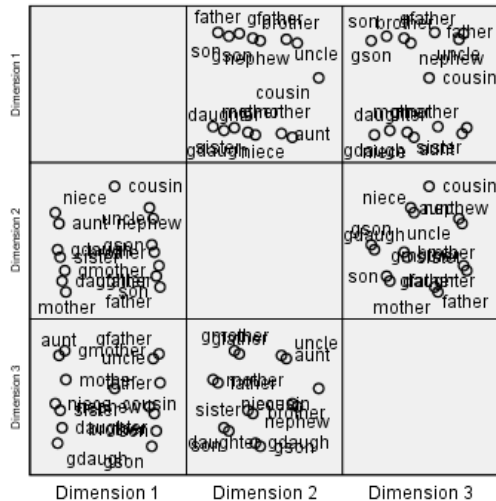
b. Optimal scaling factor = .980.

The normalized raw stress for the previous solution is 0.06234. Scaling the variables using nondefault transformations halves the stress to 0.03137.

Final Coordinates of the Common Space

The common space plots offer essentially the same interpretation of the dimensions as the previous solution.

Figure 13-21
Common space coordinates



Discussion

It is best to treat the proximities as ordinal variables, since there is a great improvement in the stress measures. As a next step, you may want to “untie” the ordinal variables—that is, allow equivalent values of the original variables to obtain different transformed values. For example, in the first source, the proximities between *Aunt* and *Son*, and *Aunt* and *Grandson*, are 85. The “tied” approach to ordinal variables forces the transformed values of these proximities to be equivalent, but there is no particular reason for you to assume that they should be. In this case, allowing the proximities to become untied frees you from an unnecessary restriction.

Recommended Readings

See the following texts for more information on multidimensional scaling:

Commandeur, J. J. F., and W. J. Heiser. 1993. *Mathematical derivations in the proximity scaling (PROXSCAL) of symmetric data matrices*. Leiden: Department of Data Theory, University of Leiden.

De Leeuw, J., and W. J. Heiser. 1980. Multidimensional scaling with restrictions on the configuration. In: *Multivariate Analysis, Vol. V*, P. R. Krishnaiah, ed. Amsterdam: North-Holland, 501–522.

Heiser, W. J. 1981. *Unfolding analysis of proximity data*. Leiden: Department of Data Theory, University of Leiden.

Heiser, W. J., and F. M. T. A. Busing. 2004. Multidimensional scaling and unfolding of symmetric and asymmetric proximity relations. In: *Handbook of Quantitative Methodology for the Social Sciences*, D. Kaplan, ed. Thousand Oaks, Calif.: SagePublications, Inc, 25–48.

Kruskal, J. B. 1964. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29, 1–28.

Kruskal, J. B. 1964. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 29, 115–129.

Shepard, R. N. 1962. The analysis of proximities: Multidimensional scaling with an unknown distance function I. *Psychometrika*, 27, 125–140.

Shepard, R. N. 1962. The analysis of proximities: Multidimensional scaling with an unknown distance function II. *Psychometrika*, 27, 219–246.

Bibliography

- Barlow, R. E., D. J. Bartholomew, D. J. Bremner, and H. D. Brunk. 1972. *Statistical inference under order restrictions*. New York: John Wiley and Sons.
- Bell, E. H. 1961. *Social foundations of human behavior: Introduction to the study of sociology*. New York: Harper & Row.
- Benzécri, J. P. 1992. *Correspondence analysis handbook*. New York: Marcel Dekker.
- Benzécri, J. P. 1969. Statistical analysis as a tool to make patterns emerge from data. In: *Methodologies of pattern recognition*, S. Watanabe, ed. New York: Academic Press, 35–74.
- Bishop, Y. M., S. E. Feinberg, and P. W. Holland. 1975. *Discrete multivariate analysis*. Cambridge, Mass.: MIT Press.
- Breiman, L., and J. H. Friedman. 1985. Estimating optimal transformations for multiple regression and correlation. *Journal of the American Statistical Association*, 80, 580–598.
- Buja, A. 1990. Remarks on functional canonical variates, alternating least squares methods and ACE. *Annals of Statistics*, 18, 1032–1069.
- Carroll, J. D. 1968. Generalization of canonical correlation analysis to three or more sets of variables. In: *Proceedings of the 76th annual convention of the American Psychological Association*, 3, , 227–228.
- Commandeur, J. J. F., and W. J. Heiser. 1993. *Mathematical derivations in the proximity scaling (PROXSCAL) of symmetric data matrices*. Leiden: Department of Data Theory, University of Leiden.
- De Haas, M., J. A. Algera, H. F. J. M. Van Tuijl, and J. J. Meulman. 2000. Macro and micro goal setting: In search of coherence. *Applied Psychology*, 49, 579–595.
- De Leeuw, J., and W. J. Heiser. 1980. Multidimensional scaling with restrictions on the configuration. In: *Multivariate Analysis, Vol. V*, P. R. Krishnaiah, ed. Amsterdam: North-Holland, 501–522.
- De Leeuw, J. 1982. Nonlinear principal components analysis. In: *COMPSTAT Proceedings in Computational Statistics*, Vienna: PhysicaVerlag, 77–89.
- De Leeuw, J. 1984. *Canonical analysis of categorical data*. Leiden: DSWO Press.

- De Leeuw, J. 1984. The Gifi system of nonlinear multivariate analysis. In: *Data analysis and informatics III*, E. Diday, and Coll., eds., 415–424.
- De Leeuw, J., and J. Van Rijckevorsel. 1980. HOMALS and PRINCALS—Some generalizations of principal components analysis. In: *Data analysis and informatics*, E. Diday, and Coll., eds. Amsterdam: North-Holland, 231–242.
- De Leeuw, J., F. W. Young, and Y. Takane. 1976. Additive structure in qualitative data: An alternating least squares method with optimal scaling features. *Psychometrika*, 41, 471–503.
- De Leeuw, J. 1990. Multivariate analysis with optimal scaling. In: *Progress in Multivariate Analysis*, S. DasGupta, and J. Sethuraman, eds. Calcutta: Indian Statistical Institute.
- Eckart, C., and G. Young. 1936. The approximation of one matrix by another one of lower rank. *Psychometrika*, 1, 211–218.
- Fisher, R. A. 1938. *Statistical methods for research workers*. Edinburgh: Oliver and Boyd.
- Fisher, R. A. 1940. The precision of discriminant functions. *Annals of Eugenics*, 10, 422–429.
- Gabriel, K. R. 1971. The biplot graphic display of matrices with application to principal components analysis. *Biometrika*, 58, 453–467.
- Gifi, A. 1985. *PRINCALS. Research Report UG-85-02*. Leiden: Department of Data Theory, University of Leiden.
- Gifi, A. 1990. *Nonlinear multivariate analysis*. Chichester: John Wiley and Sons.
- Gilula, Z., and S. J. Haberman. 1988. The analysis of multivariate contingency tables by restricted canonical and restricted association models. *Journal of the American Statistical Association*, 83, 760–771.
- Gower, J. C., and J. J. Meulman. 1993. The treatment of categorical information in physical anthropology. *International Journal of Anthropology*, 8, 43–51.
- Green, P. E., and Y. Wind. 1973. *Multiattribute decisions in marketing: A measurement approach*. Hinsdale, Ill.: Dryden Press.
- Greenacre, M. J. 1984. *Theory and applications of correspondence analysis*. London: Academic Press.
- Guttman, L. 1941. The quantification of a class of attributes: A theory and method of scale construction. In: *The prediction of personal adjustment*, P. Horst, ed. New York: Social Science Research Council, 319–348.

- 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., R. Tibshirani, and A. Buja. 1994. Flexible discriminant analysis. *Journal of the American Statistical Association*, 89, 1255–1270.
- Hastie, T., and R. Tibshirani. 1990. *Generalized additive models*. London: Chapman and Hall.
- Hayashi, C. 1952. On the prediction of phenomena from qualitative data and the quantification of qualitative data from the mathematico-statistical point of view. *Annals of the Institute of Statistical Mathematics*, 2, 93–96.
- Heiser, W. J., and F. M. T. A. Busing. 2004. Multidimensional scaling and unfolding of symmetric and asymmetric proximity relations. In: *Handbook of Quantitative Methodology for the Social Sciences*, D. Kaplan, ed. Thousand Oaks, Calif.: SagePublications, Inc, 25–48.
- Heiser, W. J. 1981. *Unfolding analysis of proximity data*. Leiden: Department of Data Theory, University of Leiden.
- Heiser, W. J., and J. J. Meulman. 1994. Homogeneity analysis: Exploring the distribution of variables and their nonlinear relationships. In: *Correspondence analysis in the social sciences: Recent developments and applications*, M. Greenacre, and J. Blasius, eds. New York: AcademicPress, 179–209.
- Heiser, W. J., and J. J. Meulman. 1995. Nonlinear methods for the analysis of homogeneity and heterogeneity. In: *Recent advances in descriptive multivariate analysis*, W. J. Krzanowski, ed. Oxford: Oxford University Press, 51–89.
- Horst, P. 1961. Generalized canonical correlations and their applications to experimental data. *Journal of Clinical Psychology*, 17, 331–347.
- Horst, P. 1961. Relations among m sets of measures. *Psychometrika*, 26, 129–149.
- Israels, A. 1987. *Eigenvalue techniques for qualitative data*. Leiden: DSWO Press.
- 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.
- Kettenring, J. R. 1971. Canonical analysis of several sets of variables. *Biometrika*, 58, 433–460.
- Kruskal, J. B. 1964. Multidimensional scaling by optimizing goodness of fit to a nonmetric hypothesis. *Psychometrika*, 29, 1–28.

- Kruskal, J. B. 1964. Nonmetric multidimensional scaling: A numerical method. *Psychometrika*, 29, 115–129.
- Kruskal, J. B. 1965. Analysis of factorial experiments by estimating monotone transformations of the data. *Journal of the Royal Statistical Society Series B*, 27, 251–263.
- Kruskal, J. B. 1978. Factor analysis and principal components analysis: Bilinear methods. In: *International encyclopedia of statistics*, W. H. Kruskal, and J. M. Tanur, eds. New York: The Free Press, 307–330.
- Kruskal, J. B., and R. N. Shepard. 1974. A nonmetric variety of linear factor analysis. *Psychometrika*, 39, 123–157.
- Krzanowski, W. J., and F. H. C. Marriott. 1994. *Multivariate analysis: Part I, distributions, ordination and inference*. London: Edward Arnold.
- Lebart, L., A. Morineau, and K. M. Warwick. 1984. *Multivariate descriptive statistical analysis*. New York: John Wiley and Sons.
- Lingoes, J. C. 1968. The multivariate analysis of qualitative data. *Multivariate Behavioral Research*, 3, 61–94.
- Max, J. 1960. Quantizing for minimum distortion. *Proceedings IEEE (Information Theory)*, 6, 7–12.
- Meulman, J. J. 2003. Prediction and classification in nonlinear data analysis: Something old, something new, something borrowed, something blue. *Psychometrika*, 4, 493–517.
- Meulman, J. 1982. *Homogeneity analysis of incomplete data*. Leiden: DSWO Press.
- Meulman, J. 1986. *A distance approach to nonlinear multivariate analysis*. Leiden: DSWO Press.
- Meulman, J. 1992. The integration of multidimensional scaling and multivariate analysis with optimal transformations of the variables. *Psychometrika*, 57, 539–565.
- Meulman, J. 1993. Principal coordinates analysis with optimal transformations of the variables: Minimizing the sum of squares of the smallest eigenvalues. *British Journal of Mathematical and Statistical Psychology*, 46, 287–300.
- Meulman, J. 1996. Fitting a distance model to homogeneous subsets of variables: Points of view analysis of categorical data. *Journal of Classification*, 13, 249–266.
- Meulman, J., and W. J. Heiser. 1997. Graphical display of interaction in multiway contingency tables by use of homogeneity analysis. In: *Visual display of categorical data*, M. Greenacre, and J. Blasius, eds. New York: Academic Press, 277–296.

- Meulman, J. J., and P. Verboon. 1993. Points of view analysis revisited: Fitting multidimensional structures to optimal distance components with cluster restrictions on the variables. *Psychometrika*, 58, 7–35.
- Meulman, J. J., A. J. van der Kooij, and A. Babinec. 2000. New features of categorical principal components analysis for complicated data sets, including data mining. In: *Classification, Automation and New Media*, W. Gaul, and G. Ritter, eds. Berlin: Springer-Verlag, 207–217.
- Meulman, J. J., A. J. van der Kooij, and W. J. Heiser. 2004. Principal components analysis with nonlinear optimal scaling transformations for ordinal and nominal data. In: *Handbook of Quantitative Methodology for the Social Sciences*, D. Kaplan, ed. Thousand Oaks, Calif.: SagePublications, Inc, 49–70.
- Nishisato, S. 1984. Forced classification: A simple application of a quantification method. *Psychometrika*, 49, 25–36.
- Nishisato, S. 1980. *Analysis of categorical data: Dual scaling and its applications*. Toronto: University of Toronto Press.
- Nishisato, S. 1994. *Elements of dual scaling: An introduction to practical data analysis*. Hillsdale, N.J.: Lawrence Erlbaum Associates, Inc.
- Pratt, J. W. 1987. Dividing the indivisible: Using simple symmetry to partition variance explained. In: *Proceedings of the Second International Conference in Statistics*, T. Pukkila, and S. Puntanen, eds. Tampere, Finland: University of Tampere, 245–260.
- Ramsay, J. O. 1989. Monotone regression splines in action. *Statistical Science*, 4, 425–441.
- Rao, C. R. 1980. Matrix approximations and reduction of dimensionality in multivariate statistical analysis. In: *Multivariate Analysis, Vol. 5*, P. R. Krishnaiah, ed. Amsterdam: North-Holland, 3–22.
- Rao, C. R. 1973. *Linear statistical inference and its applications*. New York: John Wiley & Sons.
- 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.
- Roskam, E. E. 1968. *Metric analysis of ordinal data in psychology*. Voorschoten: VAM.
- Shepard, R. N. 1962. The analysis of proximities: Multidimensional scaling with an unknown distance function I. *Psychometrika*, 27, 125–140.

- Shepard, R. N. 1962. The analysis of proximities: Multidimensional scaling with an unknown distance function II. *Psychometrika*, 27, 219–246.
- Shepard, R. N. 1966. Metric structures in ordinal data. *Journal of Mathematical Psychology*, 3, 287–315.
- Tenenhaus, M., and F. W. Young. 1985. An analysis and synthesis of multiple correspondence analysis, optimal scaling, dual scaling, homogeneity analysis, and other methods for quantifying categorical multivariate data. *Psychometrika*, 50, 91–119.
- Theunissen, N. C. M., J. J. Meulman, A. L. Den Ouden, H. M. Koopman, G. H. Verrips, S. P. Verloove-Vanhorick, and J. M. Wit. 2003. Changes can be studied when the measurement instrument is different at different time points. *Health Services and Outcomes Research Methodology*, 4, 109–126.
- Tucker, L. R. 1960. Intra-individual and inter-individual multidimensionality. In: *Psychological Scaling: Theory & Applications*, H. Gulliksen, and S. Messick, eds. New York: John Wiley & Sons, 155–167.
- Van der Burg, E. 1988. *Nonlinear canonical correlation and some related techniques*. Leiden: DSWO Press.
- Van der Burg, E., and J. De Leeuw. 1983. Nonlinear canonical correlation. *British Journal of Mathematical and Statistical Psychology*, 36, 54–80.
- Van der Burg, E., J. De Leeuw, and R. Verdegaal. 1988. Homogeneity analysis with k sets of variables: An alternating least squares method with optimal scaling features. *Psychometrika*, 53, 177–197.
- 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.
- Vander Kooij, A. J., and J. J. Meulman. 1997. MURALS: Multiple regression and optimal scaling using alternating least squares. In: *Softstat '97*, F. Faulbaum, and W. Bandilla, eds. Stuttgart: Gustav Fisher, 99–106.
- Van Rijckevorsel, J. 1987. *The application of fuzzy coding and horseshoes in multiple correspondence analysis*. Leiden: DSWO Press.
- Verboon, P., and I. A. Van der Lans. 1994. Robust canonical discriminant analysis. *Psychometrika*, 59, 485–507.
- Verdegaal, R. 1985. *Meer sets analyse voor kwalitatieve gegevens (in Dutch)*. Leiden: Department of Data Theory, University of Leiden.

- Vlek, C., and P. J. Stallen. 1981. Judging risks and benefits in the small and in the large. *Organizational Behavior and Human Performance*, 28, 235–271.
- Wagenaar, W. A. 1988. *Paradoxes of gambling behaviour*. London: Lawrence Erlbaum Associates, Inc.
- Winsberg, S., and J. O. Ramsay. 1980. Monotonic transformations to additivity using splines. *Biometrika*, 67, 669–674.
- Winsberg, S., and J. O. Ramsay. 1983. Monotone spline transformations for dimension reduction. *Psychometrika*, 48, 575–595.
- Wolter, K. M. 1985. *Introduction to variance estimation*. Berlin: Springer-Verlag.
- Young, F. W. 1981. Quantitative analysis of qualitative data. *Psychometrika*, 46, 357–387.
- Young, F. W., J. De Leeuw, and Y. Takane. 1976. Regression with qualitative and quantitative variables: An alternating least squares method with optimal scaling features. *Psychometrika*, 41, 505–528.
- Young, F. W., Y. Takane, and J. De Leeuw. 1978. The principal components of mixed measurement level multivariate data: An alternating least squares method with optimal scaling features. *Psychometrika*, 43, 279–281.
- Zeijl, E., Y. te Poel, M. du Bois-Reymond, J. Ravesloot, and J. J. Meulman. 2000. The role of parents and peers in the leisure activities of young adolescents. *Journal of Leisure Research*, 32, 281–302.

- ANOVA
 - in Categorical Regression, 26
- biplots
 - in Categorical Principal Components Analysis, 45
 - in Correspondence Analysis, 67, 257
 - in Multiple Correspondence Analysis, 82
- Categorical Principal Components Analysis, 31, 39, 159, 175
 - category points, 198
 - command additional features, 48
 - component loadings, 170, 175, 194
 - iteration history, 165
 - model summary, 165, 173, 194
 - object scores, 169, 173, 196
 - optimal scaling level, 34
 - quantifications, 167, 190
 - save variables, 44
- Categorical Regression, 19, 107
 - command additional features, 29
 - correlations, 122, 124
 - importance, 124
 - intercorrelations, 121
 - model fit, 122
 - optimal scaling level, 21
 - plots, 19
 - residuals, 128
 - save, 28
 - statistics, 19
 - transformation plots, 126
- category coordinates
 - in Nonlinear Canonical Correlation Analysis, 235
- category plots
 - in Categorical Principal Components Analysis, 46
 - in Multiple Correspondence Analysis, 83
- category points
 - in Categorical Principal Components Analysis, 198
- category quantifications
 - in Categorical Principal Components Analysis, 42
 - in Categorical Regression, 26
 - in Multiple Correspondence Analysis, 80, 305
 - in Nonlinear Canonical Correlation Analysis, 54
- centroids
 - in Nonlinear Canonical Correlation Analysis, 54, 236
- coefficients
 - in Categorical Regression, 122
- column principal normalization
 - in Correspondence Analysis, 250
- column scores
 - in Correspondence Analysis, 260
- column scores plots
 - in Correspondence Analysis, 279
- common space
 - in Multidimensional Scaling, 334, 339
- common space coordinates
 - in Multidimensional Scaling, 101
- common space plots
 - in Multidimensional Scaling, 99
- component loadings
 - in Categorical Principal Components Analysis, 42, 170, 175, 194
 - in Nonlinear Canonical Correlation Analysis, 54, 230
- component loadings plots
 - in Categorical Principal Components Analysis, 47

- confidence statistics
 - in Correspondence Analysis, 66, 263
- contributions
 - in Correspondence Analysis, 260, 277
- correlation matrix
 - in Categorical Principal Components Analysis, 42
 - in Multiple Correspondence Analysis, 80
- correlations
 - in Multidimensional Scaling, 101
- correlations plots
 - in Multidimensional Scaling, 99
- Correspondence Analysis, 59, 61, 62, 63, 66, 67, 249, 251, 270
 - biplots, 257
 - column scores plots, 279
 - command additional features, 68
 - confidence statistics, 263
 - contributions, 260, 277
 - correspondence tables, 255, 292
 - dimensions, 276
 - inertia per dimension, 256
 - normalization, 250
 - permutations, 262
 - plots, 59
 - profiles, 258
 - row and column scores, 260
 - row scores plots, 279, 293
 - statistics, 59
- correspondence tables
 - in Correspondence Analysis, 255, 292
- Cronbach's alpha
 - in Categorical Principal Components Analysis, 165
- descriptive statistics
 - in Categorical Regression, 26
- dimensions
 - in Correspondence Analysis, 63, 276
- discretization
 - in Categorical Principal Components Analysis, 36
 - in Categorical Regression, 23
 - in Multiple Correspondence Analysis, 74
- discrimination measures
 - in Multiple Correspondence Analysis, 80, 304
- discrimination measures plots
 - in Multiple Correspondence Analysis, 83
- distance measures
 - in Correspondence Analysis, 63
- distances
 - in Multidimensional Scaling, 101
- eigenvalues
 - in Categorical Principal Components Analysis, 165, 173, 194
 - in Nonlinear Canonical Correlation Analysis, 226
- fit
 - in Nonlinear Canonical Correlation Analysis, 54
- fit values
 - in Nonlinear Canonical Correlation Analysis, 226
- importance
 - in Categorical Regression, 124
- individual spaces plots
 - in Multidimensional Scaling, 99
- individual space weights
 - in Multidimensional Scaling, 101
- individual space weights plots
 - in Multidimensional Scaling, 99

- inertia
 - in Correspondence Analysis, 66, 256, 260
- initial configuration
 - in Categorical Regression, 25
 - in Multidimensional Scaling, 98
 - in Nonlinear Canonical Correlation Analysis, 54
- intercorrelations
 - in Categorical Regression, 121
- iteration criteria
 - in Multidimensional Scaling, 98
- iteration history
 - in Categorical Principal Components Analysis, 42, 165
 - in Multidimensional Scaling, 101
 - in Multiple Correspondence Analysis, 80

- joint category plots
 - in Categorical Principal Components Analysis, 46
 - in Multiple Correspondence Analysis, 83

- loss values
 - in Nonlinear Canonical Correlation Analysis, 226

- missing values
 - in Categorical Principal Components Analysis, 38
 - in Categorical Regression, 24
 - in Multiple Correspondence Analysis, 75
- model summary
 - in Multiple Correspondence Analysis, 301
- Multidimensional Scaling, 87, 90, 91, 92, 93, 94, 317, 317
 - command additional features, 103
 - common space, 334, 339
 - model, 95
 - options, 98
 - output, 101
 - plots, 87, 99, 101
 - restrictions, 97
 - statistics, 87
 - stress measures, 332, 338
 - transformation plots, 337
- Multiple Correspondence Analysis, 71, 77, 297
 - category quantifications, 305
 - command additional features, 85
 - discrimination measures, 304
 - model summary, 301
 - object scores, 302, 307
 - optimal scaling level, 73
 - outliers, 310
 - save variables, 82
- multiple *R*
 - in Categorical Regression, 26

- Nonlinear Canonical Correlation Analysis, 49, 53, 54, 217, 217
 - category coordinates, 235
 - centroids, 236
 - command additional features, 56
 - component loadings, 227, 230
 - plots, 49
 - quantifications, 231
 - statistics, 49
 - summary of analysis, 226
 - weights, 227
- normalization
 - in Correspondence Analysis, 63, 250

- object points plots
 - in Categorical Principal Components Analysis, 45
 - in Multiple Correspondence Analysis, 82
- object scores
 - in Categorical Principal Components Analysis, 42, 169, 173, 196
 - in Multiple Correspondence Analysis, 80, 302, 307
 - in Nonlinear Canonical Correlation Analysis, 54

- optimal scaling level
 - in Categorical Principal Components Analysis, 34
 - in Multiple Correspondence Analysis, 73
- outliers
 - in Multiple Correspondence Analysis, 310

- part correlations
 - in Categorical Regression, 124
- partial correlations
 - in Categorical Regression, 124
- permutations
 - in Correspondence Analysis, 262
- plots
 - in Categorical Regression, 29
 - in Correspondence Analysis, 67
 - in Multidimensional Scaling, 99, 101
 - in Nonlinear Canonical Correlation Analysis, 54
- principal normalization
 - in Correspondence Analysis, 250
- profiles
 - in Correspondence Analysis, 258
- projected centroids
 - in Nonlinear Canonical Correlation Analysis, 236
- projected centroids plots
 - in Categorical Principal Components Analysis, 46

- quantifications
 - in Categorical Principal Components Analysis, 167, 190
 - in Nonlinear Canonical Correlation Analysis, 231

- R^2
 - in Categorical Regression, 122
- regression coefficients
 - in Categorical Regression, 26

- relaxed updates
 - in Multidimensional Scaling, 98
- residuals
 - in Categorical Regression, 128
- restrictions
 - in Multidimensional Scaling, 97
- row principal normalization
 - in Correspondence Analysis, 250
- row scores
 - in Correspondence Analysis, 260
- row scores plots
 - in Correspondence Analysis, 279, 293

- standardization
 - in Correspondence Analysis, 63
- stress measures
 - in Multidimensional Scaling, 101, 332, 338
- stress plots
 - in Multidimensional Scaling, 99
- supplementary objects
 - in Categorical Regression, 25
- supplementary points
 - in Correspondence Analysis, 264
- symmetrical normalization
 - in Correspondence Analysis, 250

- transformation plots
 - in Categorical Principal Components Analysis, 46
 - in Categorical Regression, 126
 - in Multidimensional Scaling, 99, 337
 - in Multiple Correspondence Analysis, 83
- transformed independent variables
 - in Multidimensional Scaling, 101
- transformed proximities
 - in Multidimensional Scaling, 101
- triplots
 - in Categorical Principal Components Analysis, 45

-
- variable weight
in Categorical Principal Components Analysis, 34
in Multiple Correspondence Analysis, 73
- variance accounted for
in Categorical Principal Components Analysis, 42, 165, 194
- weights
in Nonlinear Canonical Correlation Analysis, 54, 227
- zero-order correlations
in Categorical Regression, 124