

APPLIED MULTIVARIATE STATISTICAL CONCEPTS

Debbie L. Hahs-Vaughn

APPLIED MULTIVARIATE STATISTICAL CONCEPTS

More comprehensive than other texts, this new book covers the classic and cutting edge multivariate techniques used in today's research. Ideal for courses on multivariate statistics, advanced statistics, or quantitative techniques taught in psychology, education, and sociology, the book also appeals to researchers with no training in multivariate methods. Through clear writing, engaging pedagogy, and examples using real data, Hahs-Vaughn walks students through the most-used methods to learn why and how to apply each technique, and annotated screenshots from SPSS and other packages are integrated throughout. The conceptual approach with a high text-to-formula ratio helps readers master key concepts so they can implement and interpret results generated by today's sophisticated software. Yet each chapter includes a mathematical snapshot that highlights the technical components of each procedure, so only the most crucial equations are included. Designed for course flexibility, after the first 4 chapters, instructors can use chapters in any sequence or combination to fit the needs of their students.

Highlights include:

- Outlines, key concepts, and opening vignettes preview what's to come in each chapter
- Examples using real data from education, psychology, and other social sciences illustrate key concepts
- Extensive coverage of assumptions including tables, the effects of their violation, and how to test for each technique
- Conceptual, computational, and interpretative problems mirror the real-world problems students encounter in their studies and careers
- A focus on data screening and power analysis with attention on the special needs of each particular method
- Instructions for using SPSS via screenshots and annotated output along with HLM, Mplus, LISREL, and G*Power where appropriate, to demonstrate how to interpret results
- Templates for writing research questions and APA-style write-ups of results
- Propensity score analysis chapter demonstrates the use of the increasingly popular technique
- A review of matrix algebra for those who want an introduction (prerequisites include an introduction to factorial ANOVA, ANCOVA, and simple linear regression, but knowledge of matrix algebra is not assumed)
- www.routledge.com/9780415842365 provides the text's datasets pre-formatted for use in SPSS for students, and answers to all chapter problems, Power Points, and test items for instructors

Debbie L. Hahs-Vaughn is a Professor in the University of Central Florida's Methodology, Measurement, & Analysis Program in the College of Education and Human Performance. Her research focuses on methodological and substantive analyses of complex

survey data, program evaluation, and practitioner use of research to inform practice. Dr. Hahs-Vaughn has co-authored two textbooks with Dr. Richard Lomax: *An Introduction to Statistical Concepts*, 3rd edition, and *Statistical Concepts*, 4th edition; has authored more than 40 articles; and has served as the Executive Editor of the Measurement, Statistics, and Research Design section of the *Journal of Experimental Education*.

She is the recipient of numerous awards including, among others, the 2015 AERA Educational Statisticians SIG Annual Service Award, 2013 Researcher of the Year Award from the Florida Educational Research Association, 2014 and 2009 Excellence in Research Awards (UCF), 2013 and 2008 UCF Teaching Incentive Program Award, and 2007 Excellence in Graduate Teaching Award. Dr. Hahs-Vaughn received her Ph.D. in Educational Research from the University of Alabama.

“Hahs-Vaughn provides a strong foundation for learning advanced statistical techniques by first explaining ‘why’ each analysis is used and then supporting the ‘how’ each statistical application is conducted with a review of basic theoretical concepts and a summary of the mathematical background for each statistical analysis. Her approach provides exactly the right balance of theory to practice for understanding and applying multivariate statistical analyses.”

—Robyn Cooper, Drake University

“Ideal for students in a wide variety of social science disciplines, this book approaches multivariate statistics with an appealing mix of conceptual and technical content. Easy-to-follow and interesting demonstrations of applications to real-world problems make it an ideal teaching tool and will keep students engaged. The writing is clear, concise, and informative in a way that students at the advanced undergraduate and graduate levels will really appreciate. Unlike many other multivariate statistics textbooks, this book includes important concepts such as cluster analysis and propensity score analysis, which are very important areas that are too infrequently covered.”

—W. Holmes Finch, Ball State University

“The text provides comprehensive coverage of multivariate statistical techniques, with explanations that are both concise and clear. The step-by-step instructions and annotated outputs will continue to serve as excellent resources for students even after completing the course.”

—Sylvie Mrug, University of Alabama at Birmingham

“*Applied Multivariate Statistical Concepts* is a great addition to . . . textbooks in the social and behavioral sciences for graduate students and researchers. The author took extreme care in selecting key pedagogical methods and statistical procedures in current use. Moreover, she makes sure that students will have the necessary tools to see their projects completed from start to finish by providing innovative instructional and learning strategies specific to relevant fields of study. Students will be challenged by the topics but also guided throughout the research enterprise with step-by-step instructions including the appropriate use of various statistical software applications on real data sets.”

—Arturo Olivárez, Jr., University of Texas at El Paso

APPLIED MULTIVARIATE STATISTICAL CONCEPTS

Debbie L. Hahs-Vaughn

First published 2017
by Routledge
711 Third Avenue, New York, NY 10017

and by Routledge
2 Park Square, Milton Park, Abingdon, Oxon, OX14 4RN

Routledge is an imprint of the Taylor & Francis Group, an informa business

© 2017 Taylor & Francis

The right of Debbie L. Hahs-Vaughn to be identified as author of this work has been asserted by her in accordance with sections 77 and 78 of the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this book may be reprinted or reproduced or utilised in any form or by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying and recording, or in any information storage or retrieval system, without permission in writing from the publishers.

Trademark notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

Names: Hahs-Vaughn, Debbie L.

Title: Applied multivariate statistical concepts / by Debbie Hahs-Vaughn.

Description: New York, NY : Routledge, 2017.

Identifiers: LCCN 2016017149 | ISBN 9780415842365 (pbk. : alk. paper) | ISBN 9780415842358 (hardback : alk. paper) | ISBN 9781315816685 (ebook)

Subjects: LCSH: Multivariate analysis—Textbooks. | Mathematical statistics—Textbooks.

Classification: LCC QA278. H336 2017 | DDC 519.5/35—dc23

LC record available at <https://lccn.loc.gov/2016017149>

ISBN: 978-0-415-84235-8 (hbk)

ISBN: 978-0-415-84236-5 (pbk)

ISBN: 978-1-315-81668-5 (ebk)

Typeset in Times New Roman
by Apex CoVantage, LLC

Visit the eResources: <https://www.routledge.com/9780415842365>

BRIEF CONTENTS

<i>Preface</i>	<i>xi</i>
<i>Acknowledgments</i>	<i>xiii</i>
1 Multivariate Statistics	1
2 Univariate and Bivariate Statistics Review	9
3 Data Screening	35
4 Multiple Linear Regression	57
5 Logistic Regression	117
6 Multivariate Analysis of Variance: Single Factor, Factorial, and Repeated Measures Designs	169
7 Discriminant Analysis	273
8 Cluster Analysis	335
9 Exploratory Factor Analysis	362
10 Path Analysis, Confirmatory Factor Analysis, and Structural Equation Modeling	441
11 Multilevel Linear Modeling	505
12 Propensity Score Analysis	571
Appendix A: An Introduction to Matrix Algebra	599
Appendix B: Answers to Odd-Numbered Conceptual & Computational Questions	609
<i>Index</i>	635



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

DETAILED CONTENTS

<i>Preface</i>	<i>xi</i>	
<i>Acknowledgments</i>	<i>xiii</i>	
1	Multivariate Statistics	1
1.1	What Are Multivariate Statistics?	2
1.2	Decision Rules	2
1.3	Coverage of the Textbook	4
1.4	Layout of the Textbook	6
1.5	Overarching Goal of the Textbook	6
2	Univariate and Bivariate Statistics Review	9
2.1	Fundamental Concepts	10
2.2	Foundational Univariate Statistics	14
2.3	Foundational Bivariate Statistics	21
3	Data Screening	35
3.1	Independence	36
3.2	Variance	41
3.3	Normality	43
3.4	Linearity	51
3.5	Noncollinearity	52
4	Multiple Linear Regression	57
4.1	What Multiple Linear Regression Is and How It Works	58
4.2	Mathematical Introduction Snapshot	83
4.3	Computing Multiple Linear Regression Using SPSS	87
4.4	Data Screening	96
4.5	Power Using G*Power	104
4.6	Research Question Template and Example Write-Up	107

5	Logistic Regression	117
5.1	What Logistic Regression Is and How It Works	118
5.2	Mathematical Introduction Snapshot	138
5.3	Computing Logistic Regression Using SPSS	139
5.4	Data Screening	150
5.5	Power Using G*Power	160
5.6	Research Question Template and Example Write-Up	163
6	Multivariate Analysis of Variance: Single Factor, Factorial, and Repeated Measures Designs	169
6.1	What Multivariate Analysis of Variance Is and How It Works	170
6.2	Mathematical Introduction Snapshot	188
6.3	Computing MANOVA Using SPSS	191
6.4	Data Screening	227
6.5	Power Using G*Power	251
6.6	Research Question Template and Example Write-Up	260
7	Discriminant Analysis	273
7.1	What Discriminant Analysis Is and How It Works	275
7.2	Mathematical Introduction Snapshot	292
7.3	Computing Discriminant Analysis Using SPSS	293
7.4	Data Screening	315
7.5	Power Using G*Power	324
7.6	Research Question Template and Example Write-Up	328
8	Cluster Analysis	335
8.1	What Cluster Analysis Is and How It Works	336
8.2	Mathematical Introduction Snapshot	345
8.3	Computing Cluster Analysis Using SPSS	346
8.4	Data Screening	358
8.5	Research Question Template and Example Write-Up	358
9	Exploratory Factor Analysis	362
9.1	What Exploratory Factor Analysis Is and How It Works	363
9.2	Mathematical Introduction Snapshot	381
9.3	Computing EFA Using SPSS	383
9.4	Data Screening	427
9.5	Research Question Template and Example Write-Up	433
10	Path Analysis, Confirmatory Factor Analysis, and Structural Equation Modeling	441
10.1	What Path Analysis and Confirmatory Factor Analysis Are and How They Work	442
10.2	Mathematical Introduction Snapshot	465
10.3	Computing Path Analysis and Confirmatory Factor Analysis Using LISREL	466

10.4	Data Screening	494
10.5	Power	494
10.6	Research Question Template and Example Write-Up	496
11	Multilevel Linear Modeling	505
11.1	What Multilevel Linear Modeling Is and How It Works	507
11.2	Mathematical Introduction Snapshot	530
11.3	Computing Multilevel Modeling Using HLM	531
11.4	Data Screening	553
11.5	Power Using Optimal Design	559
11.6	Research Question Template and Example Write-Up	561
12	Propensity Score Analysis	571
12.1	What Propensity Score Analysis Is and How It Works	572
12.2	Mathematical Introduction Snapshot	582
12.3	Computing Propensity Score Analysis Using R	582
12.4	Example Write-Up	594
Appendix A: An Introduction to Matrix Algebra		599
A.1	Matrices	600
A.2	Calculations With Matrices	600
A.3	Types of Matrices	602
A.4	Matrices and Multivariate Statistics	605
Appendix B: Answers to Odd-Numbered Conceptual & Computational Questions		609
<i>Index</i>		635



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

PREFACE

This book is designed for a single course or two-course sequence in applied multivariate statistics for graduate students broadly in the social sciences, including education and the behavioral sciences. A basic understanding of applied multivariate methods is becoming more and more necessary for students pursuing advanced degrees. As a result, more and more graduate students are required, or strongly recommended, to take this course than ever before. In short, we live in a multivariate world consisting of many variables operating simultaneously and thus we need to be able to read, critique, and conduct multivariate research. After the first three introductory chapters, the text is designed to be quite flexible in terms of topical coverage. Thus, no specific sequence of chapters is necessary for the use of this text, nor are any of these chapters absolutely essential. The only exception may be coverage of multiple regression (chapter 4) prior to introducing latent variable modeling or multilevel linear modeling. This will allow researchers within and across disciplines to utilize those multivariate methods most appropriate for their discipline and for their course.

An important aspect of the text is clear writing, which is sacrificed in many texts for formal statistical presentation of formulae. This text has a very high text-to-formula ratio to increase understanding of the important concepts. The rationale is that because of the sophisticated statistics software now on the market, it is the concepts, and the implementation and interpretation of results from the software, that are most important, with hand computations and extensive formulae being generally unnecessary. Thus, only the most crucial equations are included. However, each chapter includes a ‘mathematical snapshot’ that provides an overview of more technical components of the procedure. An overview of matrix algebra is provided in the appendix for those who desire an introduction to and more in-depth understanding of this area.

This textbook instead takes a conceptual and applied approach, including features in every chapter that help students understand both why and how to apply each statistical method to their work. Features include:

- Each chapter opens with a **Chapter Outline, Lists of Key Concepts and Objectives**, and **Chapter-Opening Vignettes** related to the concepts covered;

- **Realistic Examples from the Social Sciences** are used to illustrate the concepts along with **Real Data Sets** so students can connect the models to their own work;
- Each chapter concludes with numerous **Conceptual, Computational, and Interpretive (Critical Thinking) Problems** (with odd-numbered answers to conceptual questions in the back of the book, and even-numbered answers in the instructor's manual), **Templates for Writing Research Questions**, and **Example Write-Ups of Results** that provide manuscript-ready models;
- **Step-by-Step Instructions on Using Statistical Software** are thoroughly integrated into the text (input and output) along with screenshots and annotated output, primarily using SPSS but also including other software where appropriate;
- Clear **Coverage of Assumptions**, including how to test them and the effects of their violation, is included in all relevant chapters;
- **Data Screening** is covered in both a discrete introductory chapter and in a section of each chapter to cover the special needs of each particular method;
- Plenty of **Tables and Figures**;
- An **Extensive Reference List**;
- A free **eResource** at www.routledge.com/9780415842365 includes supplemental content to assist both students and instructors, including an instructor's manual of answers to all end-of-chapter problems; all datasets that are utilized in the text in a format that can be used in SPSS and most other relevant statistical packages; PowerPoint slides; and test items.

ACKNOWLEDGMENTS

Thank you to Routledge/Taylor & Francis, particularly Debra Riegert, for shepherding our introductory textbook, encouraging the pursuit of this text, and for your incredible support and flexibility along the way. You are a joy to work with! Thanks also to Alison Daltroy for the arduous task of summarizing hundreds of pages of recommendations from reviewers, encouraging words to keep plugging, and a ton of assistance with formatting and other preproduction work for which I am so grateful. I am incredibly appreciative of the insightful suggestions provided by so many anonymous reviewers of this text as it went through development, as well as of Charles R. Ciorba (University of Oklahoma), Robyn Cooper (Drake University), Pamela Davis-Kean (University of Michigan), George M. Diekhoff (Midwestern State University), Brian Distelberg (Loma Linda University), John Hulland (University of Georgia), Arturo Olivarez (The University of Texas at El Paso), Rodney K. Schutz (Georgia State University), Christopher Sink (Seattle Pacific University), David Stockburger (Missouri State University), Robert Triscari (Florida Gulf Coast University), and William B. Ware (University of North Carolina, Chapel Hill). It is a much richer tool because of the valuable suggestions offered.

I would be remiss if I did not thank the multitude of students at both the University of Alabama (Roll Tide!) and the University of Central Florida (Go Knights!) that have helped *me* to increase my depth and breadth of statistics, both while as a student and while as a professor.

There are a few individuals, in particular, to whom I am forever indebted. Dr. Richard Lomax, my advisor at the University of Alabama, now at The Ohio State University, who has been a role model for and mentor to me since he was *my* introductory statistics professor. Thank you for beginning this multivariate journey with me and encouraging me during the process. Last but never least, I am most grateful to my *Ohana*: Mark, Malani, Molly (woof woof), and Camel (glub glub). I love you and it is because of your love and support that I had the fortitude, courage, endurance, resilience, and other similar juicy words to cope with and complete such a monumental project. Thank you one and all. ☺



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Chapter 1

MULTIVARIATE STATISTICS

CHAPTER OUTLINE

1.1	What Are Multivariate Statistics?	2
1.2	Decision Rules	2
1.3	Coverage of the Textbook	4
1.3.1	Multiple Regression	4
1.3.2	Logistic Regression	4
1.3.3	Multivariate Analysis of Variance	4
1.3.4	Discriminant Analysis	5
1.3.5	Cluster Analysis	5
1.3.6	Exploratory Factor Analysis	5
1.3.7	Path Analysis, Confirmatory Factor Analysis, and Structural Equation Modeling	5
1.3.8	Multilevel Linear Modeling	5
1.3.9	Propensity Score Analysis	5
1.4	Layout of the Textbook	6
1.5	Overarching Goal of the Textbook	6

1.1 WHAT ARE MULTIVARIATE STATISTICS?

Statistics can be categorized in various ways; one such way is based on the number of variables employed. Univariate statistics refers to statistics that employ just one variable. Bivariate statistics employ two variables. ‘Multivariate,’ by definition, means multiple variables. Even though, technically, ‘multiple variables’ may refer to only two variables, in the context of statistics, any procedure that simultaneously employs *more than two* variables is a multivariate procedure. As we’ll see, the multiple variables may come in the form of multiple independent variables, multiple dependent variables, or both. Many times, these multivariate procedures have a bivariate cousin, such as simple regression (with only one independent and one dependent variable) which segues into multiple regression (with two or more independent variables) or analysis of variance (with only one independent and one dependent variable) which segues into multivariate analysis of variance (with one or more independent variables and two or more dependent variables).

We can also think of multivariate statistics in terms of classic versus modern. Classic multivariate techniques include, for example, multivariate analysis of variance, multiple linear and logistic regression, discriminant analysis, cluster analysis, and exploratory factor analysis. Modern multivariate techniques, however, are (at least chronologically) statistical techniques that have recently become more common, such as multilevel linear modeling, structural equation modeling, and propensity score analysis.

By default of additional variables, multivariate statistics are more advanced, which may sometimes operationally be defined as more complex and/or more difficult to understand. However, the short-term pain of learning more advanced procedures equates to long-term gain of having a more powerful toolbox with which to work. An added bonus: *Multivariate statistical procedures are just really fun.* The types of questions that can be explored are nearly limitless, and the interpretations of the data are incredibly rich. Once you have one or two multivariate statistics under your belt, you’ll never want to run a *t* test or correlation again!

1.2 DECISION RULES

The decision rules for when to use which procedure are pretty straightforward in bivariate statistics. For example, if the goal is to determine mean differences, a *t* test or analysis of variance would be considered as a starting point. A relationship between two variables would be examined using a correlation or simple regression. The decision rules in multivariate statistics are not always as straightforward, an artifact of the complexity of these procedures and the additional elements that they entail. Regardless, some decision rules for consideration in planning your analysis are offered in Figure 1.1.

As seen in the figure, the primary purpose of the research question will guide the decision first. In the context of this textbook, questions of relationship/prediction, mean difference, and classification are covered. From there, the analytic decision is based

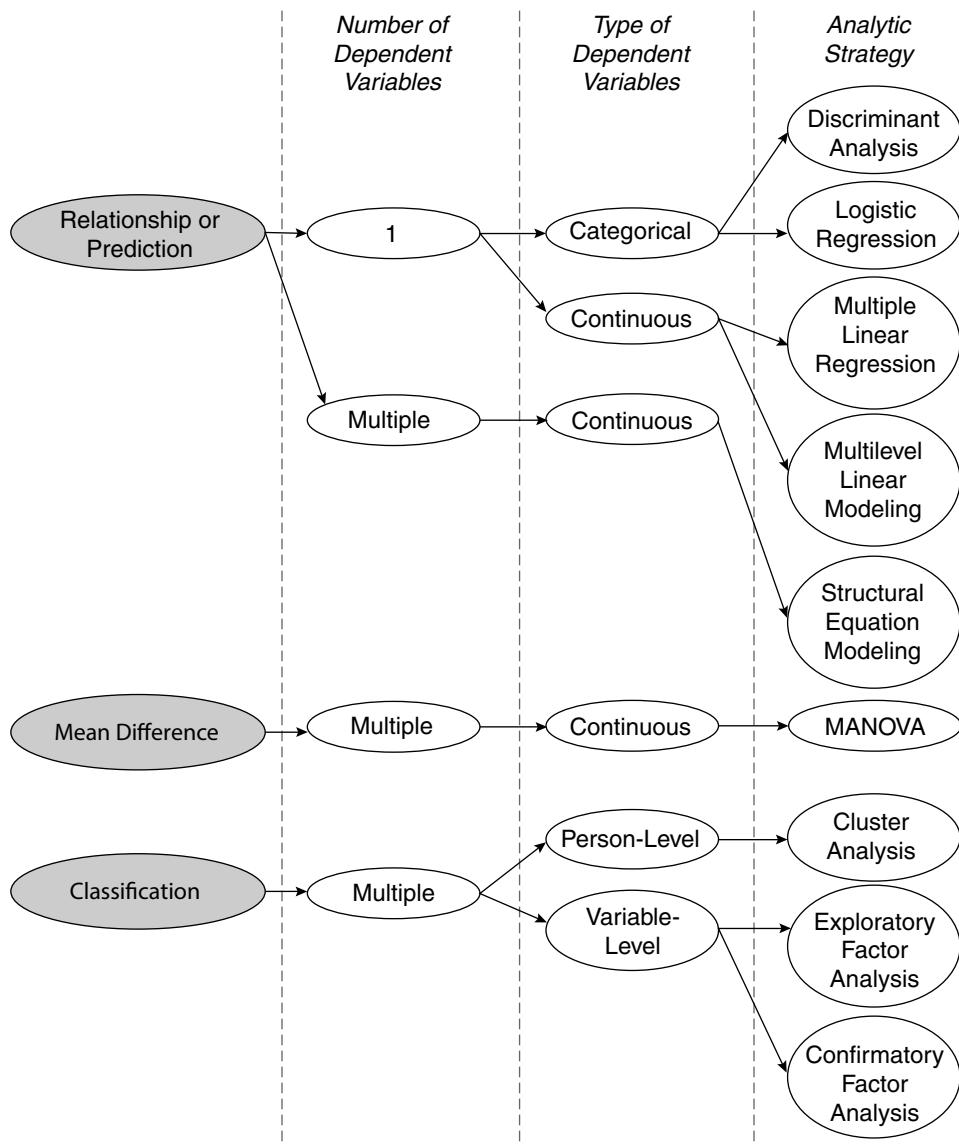


FIGURE 1.1
Decision Rules

on the number and types of variables in your model. For example, if the goal of my research is prediction and I have only one dependent categorical variable, there are two potential multivariate procedures that may be applicable: discriminant analysis and logistic regression. As we will later learn, the decision on which to use is based on both the type of independent variables and the extent to which the assumptions of the test are met. Thus, there is not necessarily one clear-cut decision on which multivariate procedure to select when the goal of my research is prediction and I have only one dependent categorical variable. Regardless, the intent of the diagram in Figure 1.1

is to point you in the right direction, understanding that the final procedure selected may entail more probing than a simple decision tree will allow for. What has not been illustrated on this decision tree is propensity score analysis (PSA). PSA is a preprocessor, so to speak, and allows for the matching of units prior to the primary analytic procedure being applied. Any of the procedures illustrated here can be applied to data that has been matched using propensity score.

1.3 COVERAGE OF THE TEXTBOOK

Chapters 2 and 3 cover data screening and offer a refresher on univariate and bivariate statistics. Data screening is a helpful precursor to many of the common assumptions that will be discussed in more detail in the text. Univariate and bivariate statistics will reacquaint you to several of the multivariate cousins with which you are already familiar. Chapters 3 through 11 are devoted to some of the most common multivariate analytical procedures: multiple regression; logistic regression; multivariate analysis of variance; discriminant analysis; cluster analysis; exploratory factor analysis; path analysis, confirmatory factor analysis, and structural equation modeling; and multilevel linear modeling. The final chapter is an introduction to propensity score analysis, a statistical procedure for matching. The book is written with a conceptual focus, not a technical one. However, there is an appendix that provides an introduction to matrix algebra, so that there is some scaffolding for the mathematics that underlie multivariate statistics.

1.3.1 Multiple Regression

Multiple linear regression is the multivariate extension of simple linear regression, which allows for prediction. Whereas simple linear regression includes one independent and one dependent variable, multiple linear regression includes *two or more* independent variables and one continuous dependent variable.

1.3.2 Logistic Regression

There are many outcomes that are worthy of investigation to determine predictions that are not continuous. In the case of categorical outcomes, such as pass/fail or win/lose, logistic regression is appropriate. When the outcome includes only two categories, binary logistic regression is appropriate. When the outcome includes more than two categories, multinomial logistic regression is appropriate. Logistic regression shares commonality with discriminant analysis but is more flexible in terms of the types of independent variables that can be applied to the logistic models. Both continuous and categorical (binary) variables can be used as independent variables in logistic regression.

1.3.3 Multivariate Analysis of Variance

Multivariate analysis of variance (MANOVA) is the multivariate cousin of analysis of variance. There are multiple MANOVA procedures, including single factor (i.e.,

one independent variable), factorial (two or more independent variables), and repeated measures designs that examine repeated outcomes.

1.3.4 Discriminant Analysis

When the goal is predicting group membership from continuous independent variables, discriminant analysis is appropriate. Discriminant analysis is often considered MANOVA in reverse.

1.3.5 Cluster Analysis

Cluster analysis is considered a ‘person-level’ approach to examining data, in that groups of units are classified based on the variables. Units (e.g., people or things) that share commonalities on the variables are clustered together. Cluster analysis may be helpful to determine profiles of users or consumers, for example.

1.3.6 Exploratory Factor Analysis

In comparison to cluster analysis, exploratory factor analysis (EFA) is a ‘variable-level’ approach to examining data, in that variables are grouped together based on the units. Variables that share commonalities on the units are grouped together. EFA may be helpful in determining evidence of construct validity, for example. EFA is a technique that is appropriate when you are in search of a structure to the data and there is minimal, if no, theoretical basis on which to derive the model.

1.3.7 Path Analysis, Confirmatory Factor Analysis, and Structural Equation Modeling

Path analysis, confirmatory factor analysis, and structural equation modeling are all procedures that fall under the general (and large) umbrella of latent variable modeling. These techniques are appropriate when you are in search of a structure to the data that is grounded in theory.

1.3.8 Multilevel Linear Modeling

Multilevel linear modeling is an extension of multiple regression that is appropriate when the data are nested in structure. Most data, if you think of it, are often nested. For example, children within classroom, employee within organization, athlete within team, resident within neighborhood. . . . The examples are nearly endless. Ignoring the context of the nesting may mask relationships, and multilevel linear modeling provides for a more complete contextualization of the situation.

1.3.9 Propensity Score Analysis

Propensity score analysis is a preparatory step to analyses that allows units to be matched on covariates and should be considered in situations where group analysis

is of interest but randomization was not feasible. Thus, propensity score matching is a very powerful tool that can be used to circumvent self-selection. After the units are matched, any procedure (descriptive or inferential) that is appropriate given the data can be applied.

1.4 LAYOUT OF THE TEXTBOOK

With the exception of the first few chapters, each chapter has been organized with the attempt to be consistent in how the material is presented. For example, each chapter begins with a chapter outline followed by key concepts. The first section of each chapter outlines what that particular procedure is and how it works. A mathematical introduction snapshot follows that provides some foundational information on the mathematics that underlie the procedure. Computing the procedure using statistical software is next. For many of the procedures, this entails using SPSS; however, some procedures apply more-specialized software (e.g., HLM for multilevel modeling). The chapter concludes with a section on data screening, a template for writing a research question for the requisite procedure, and an example write-up of results that mimics what would be appropriate for a manuscript.

1.5 OVERARCHING GOAL OF THE TEXTBOOK

You may or may not be acquainted with my co-authored introductory statistics textbooks (Lomax & Hahs-Vaughn, 2012a, 2012b). If you are not, I encourage you to pick one up, particularly as a resource for the univariate and bivariate statistical concepts that underlie multivariate statistics. If you are familiar with one or both of my previous texts, you will likely recall that they are (as the title suggests) focused on *concepts*. In other words, if you desire to be deeply entrenched in the mathematical statistics of multivariate statistics, this textbook will not live up to your expectations. On the other hand, if you seek to understand the concepts and desire to apply them (either as a consumer or producer of research), this book will be a valuable source. Treat it as a compass, and enjoy the journey.

PROBLEMS

Conceptual Problems

1. Which one of the following procedures would be appropriate to consider if the interest is in a relationship and there are multiple dependent variables?
 - a. Logistic regression
 - b. MANOVA
 - c. Multiple linear regression
 - d. Structural equation modeling

2. Which one of the following procedures would be appropriate to consider if the interest is in a relationship and there is one dependent categorical variable?
 - a. Logistic regression
 - b. MANOVA
 - c. Multiple linear regression
 - d. Structural equation modeling
3. Which one of the following procedures would be appropriate to consider if the interest is in a relationship and there is one dependent continuous variable?
 - a. Logistic regression
 - b. MANOVA
 - c. Multiple linear regression
 - d. Propensity score analysis
4. Which one of the following procedures would be appropriate to consider if the interest is in a mean difference and there are multiple dependent variables?
 - a. Logistic regression
 - b. MANOVA
 - c. Multiple linear regression
 - d. Structural equation modeling
5. Which one of the following procedures would be appropriate to consider if the interest is in examining contextual relationships of individuals that are within groups?
 - a. Cluster analysis
 - b. Discriminant analysis
 - c. Multilevel linear modeling
 - d. Confirmatory factor analysis
6. Which one of the following procedures would be appropriate to consider if the interest is in developing profiles of people?
 - a. Cluster analysis
 - b. Discriminant analysis
 - c. Multilevel linear modeling
 - d. Confirmatory factor analysis
7. Which one of the following procedures would be appropriate to consider if the interest is in classifying variables into groups based on a strong theoretical basis?
 - a. Cluster analysis
 - b. Discriminant analysis
 - c. Multilevel linear modeling
 - d. Confirmatory factor analysis
8. Which one of the following procedures would be appropriate to consider if the interest is in predicting groups based on variables?
 - a. Cluster analysis
 - b. Discriminant analysis
 - c. Multilevel linear modeling
 - d. Confirmatory factor analysis

9. Which one of the following procedures would be appropriate to consider if the interest is matching units prior to conducting a test of inference?
 - a. Exploratory factor analysis
 - b. Logistic regression
 - c. MANOVA
 - d. Propensity score analysis
10. Which one of the following procedures would be appropriate to consider if the interest is in determining mean differences of multiple outcomes?
 - a. Exploratory factor analysis
 - b. Logistic regression
 - c. MANOVA
 - d. Propensity score analysis

REFERENCES

- Lomax, R. G., & Hahs-Vaughn, D. L. (2012a). *An introduction to statistical concepts* (3rd ed.). New York: Taylor & Francis.
- Lomax, R. G., & Hahs-Vaughn, D. L. (2012b). *Statistical concepts: A second course* (4th ed.). New York: Taylor & Francis.

Chapter 2

UNIVARIATE AND BIVARIATE STATISTICS REVIEW

CHAPTER OUTLINE

2.1 Fundamental Concepts	10
2.1.1 Hypothesis Testing	10
2.1.2 Types of Decision Errors	11
2.1.3 Statistical Versus Practical Significance	13
2.2 Foundational Univariate Statistics	14
2.2.1 Histogram	14
2.2.2 Box-and-Whisker Plot	14
2.2.3 Scatterplot	16
2.2.4 Measures of Central Tendency	17
2.2.5 Measures of Dispersion	19
2.3 Foundational Bivariate Statistics	21
2.3.1 Independent and Dependent Samples t Test	21
2.3.2 Analysis of Variance	23
2.3.3 Two-Factor ANOVA	25
2.3.4 Covariance	27
2.3.5 Pearson Product-Moment Correlation Coefficient	28
2.3.6 Simple Linear Regression	28

KEY CONCEPTS

1. Covariance
2. Correlation
3. Effect size
4. Distributional shape

5. Graphical representation
6. Hypothesis test
7. Power
8. Tests of means
9. Tests of relationships
10. Type I and Type II errors

By this point in your quantitative statistics career, you've been introduced to a large number of inferential statistics procedures, and you are likely quite comfortable with concepts such as hypothesis testing, significance, power, and similar, as well as procedures to test mean differences and relationships. Nonetheless, it does not hurt to have a brief review of some common univariate and bivariate statistics, as the concepts behind them and the statistics themselves are the foundation for multivariate statistics. We'll begin with a refresher on a number of basic fundamental concepts, and then touch on a few foundational univariate and bivariate statistics. The goal is that, by the end of this chapter, you will be refreshed on a number of basic concepts that are important as we move into multivariate statistics.

2.1 FUNDAMENTAL CONCEPTS

This section is not meant to serve as either a comprehensive review or an exhaustive summary of fundamental concepts in quantitative statistics. Rather, this section is meant to refresh and reacquaint you with some key concepts that have been and will continue to be important as we move into multivariate statistics.

2.1.1 Hypothesis Testing

Hypothesis testing is a decision-making process where two possible decisions are weighed in a statistical fashion. Sample data are used to help us select one of these decisions. The two types of hypotheses competing against one another are known as the **null or statistical hypothesis**, denoted by H_0 , and the **scientific, alternative, or research hypothesis**, denoted by H_1 . The null or statistical hypothesis is a statement about the value of an unknown population parameter. The null hypothesis is basically set up by the researcher in an attempt to reject the null hypothesis in favor of our own personal scientific, alternative, or research hypothesis. Thus, we are trying to reject the null hypothesis and find evidence in favor of our scientific hypothesis, H_1 .

Based on the sample data, hypothesis testing involves making a decision as to whether the null or the research hypothesis is supported. Because we are dealing with sample statistics in our decision-making process and trying to make an inference back to the population parameter(s), there is always some risk of making an incorrect decision. This uncertainty is due to sampling error, which can be described by a probability statement. That is, because the decision is made based on sample data, the sample

may not be very representative of the population and therefore leads us to an incorrect decision. If we had population data, we would always make the correct decision about a population parameter. Because we usually do not, we use inferential statistics to help make decisions from sample data and infer those results back to the population. The nature of such decision errors and the probabilities we can attribute to them are described in the next section.

2.1.2 Types of Decision Errors

For our statistical decision, there are two distinct possibilities. Either we fail to reject H_0 or we reject H_0 , based on our sample data. If we look inside of Table 2.1, we see four different outcomes based on a combination of our statistical decision and the state of nature. Consider the first row of the table where H_0 is in actuality true. First, if H_0 is true and we fail to reject H_0 , then we have made a correct decision; that is, we have correctly failed to reject a true H_0 . The probability of this first outcome is known as $1 - \alpha$ (where α represents alpha). Second, if H_0 is true and we reject H_0 , then we have made a decision error known as a **Type I error**. That is, we have incorrectly rejected a true H_0 . Our sample data has led us to a different conclusion than the population data would have. The probability of this second outcome is known as α . Therefore, if H_0 is actually true, then our sample data lead us to one of two conclusions, either we correctly fail to reject H_0 , or we incorrectly reject H_0 . The sum of the probabilities for these two outcomes when H_0 is true is equal to 1 [i.e., $(1 - \alpha) + \alpha = 1$].

Consider now the second row of the table where H_0 is in actuality false. First, if H_0 is really false and we fail to reject H_0 , then we have made a decision error known as a **Type II error**. That is, we have incorrectly failed to reject a false H_0 . Our sample data has led us to a different conclusion than the population data would have. The probability of this outcome is known as β (beta). Second, if H_0 is really false and we reject H_0 , then we have made a correct decision; that is, we have correctly rejected a false H_0 . The probability of this second outcome is known as $1 - \beta$ or **power**. Therefore, if H_0 is actually false, then our sample data lead us to one of two conclusions, either we incorrectly fail to reject H_0 , or we correctly reject H_0 . The sum of the probabilities for these two outcomes when H_0 is false is equal to 1 [i.e., $\beta + (1 - \beta) = 1$].

■ TABLE 2.1

Statistical Decision Table

State of Nature (Reality)	Decision	
	Fail to reject H_0	Reject H_0
H_0 is true	Correct decision $(1 - \alpha)$	Type I error (α)
H_0 is false	Type II error (β)	Correct decision $(1 - \beta) = \text{power}$

One can never totally eliminate the possibility of both a Type I and a Type II error. No matter what decision we make, there is always some possibility of making a Type I and/or Type II error. Therefore, as researchers our job is to make conscious decisions in designing and conducting our study and in analyzing the data so that the possibility of decision error is minimized.

2.1.2.1 Level of Significance (α)

The probability of a Type I error is α , which is also known as the level of significance or significance level. Alpha is used as the basis for helping to make statistical decisions. If there is a relatively high cost associated with a Type I error—for example, such that lives are lost, as in the medical field—then one would want to select a relatively small level of significance (e.g., .01 or smaller). A small alpha would translate to a very small probability of rejecting the null if it were really true (i.e., a small probability of making an incorrect decision). If there is a relatively low cost associated with a Type I error—for example, such that children have to eat the second-rate candy rather than the first—then selecting a larger level of significance may be appropriate (e.g., .05 or larger). Costs are not always known, however. A second consideration is the level of significance commonly used in your field of study. In many disciplines, the .05 level of significance has become the standard (although no one seems to have a really good rationale). This is true in many of the social and behavioral sciences. Thus, you would do well to consult the published literature in your field to see if some standard is commonly used and to consider it for your own research.

2.1.2.2 Type II Error (β) and Power ($1 - \beta$)

Type II error is the probability of failing to reject H_0 when H_0 is really false. In other words, in reality H_0 is false, yet we made a decision error and did not reject H_0 . The probability associated with a Type II error is denoted by β . **Power** is a related concept and is defined as the probability of rejecting H_0 when H_0 is really false. In other words, in reality H_0 is false, and we made the correct decision to reject H_0 . The probability associated with power is denoted by $1 - \beta$.

Power is determined by five different factors: (1) level of significance (i.e., alpha), (2) sample size, (3) population standard deviation, (4) difference between the true population mean μ and the hypothesized mean value μ_0 , and (5) directionality of the test (i.e., one- or two-tailed test).

First, power is determined by the level of significance, α . As alpha increases, power increases. Thus, if alpha increases from .05 to .10, then power will increase. Second, power is determined by sample size. As sample size n increases, power increases. Third, power is determined by the size of the population standard deviation, σ . Although not under the researcher's control, as the population standard deviation increases, power decreases. Thus, if the population standard deviation *increases*, meaning the variability in the population is larger, it will cause the standard error of the mean to increase,

as there is more sampling error with larger variability. If the population standard deviation *decreases*, meaning the variability in the population is smaller, it will cause the standard error of the mean to decrease, as there is less sampling error with smaller variability. When the standard error term decreases, the denominator is smaller and thus the test statistic value becomes larger (and thereby easier to reject the null hypothesis). Fourth, power is determined by the difference between the true population mean μ and the hypothesized mean value μ_0 . Although not always under the researcher's control, as the difference between the true population mean and the hypothesized mean value increases, power increases. Thus, if the difference between the true population mean and the hypothesized mean value is large, it will be easier to correctly reject H_0 . This would result in greater separation between the two sampling distributions. Finally, power is determined by directionality and type of statistical procedure—whether we conduct a one- or a two-tailed test, as well as the type of test of inference. There is greater power in a one-tailed test than in a two-tailed test and in parametric (as compared to nonparametric) procedures.

2.1.3 Statistical Versus Practical Significance

Are statistically significant results always practically significant? In other words, if a result is statistically significant, should we make a big deal out of this result in a practical sense? Sample size can drive the results of the hypothesis test, and it is possible to find statistical significance simply as an artifact of sample size. Holding all else constant, increasing the sample size will allow (at some point) nonstatistically significant results to become statistically significant. If we gather enough sample data, any small difference, no matter how small, can wind up being statistically significant. Thus, larger samples are more likely to yield statistically significant results. Practical significance is not entirely a statistical matter. It is also a matter for the substantive field under investigation. Therefore, the meaningfulness of a “small difference” is for the substantive area to determine. All that inferential statistics can really determine is *statistical significance*. However, we should always keep practical significance in mind when interpreting our findings. In recent years, a major debate has been ongoing in the statistical community about the role of significance testing. The debate centers on whether null hypothesis significance testing (NHST) best suits the needs of researchers. At one extreme, some argue that NHST is fine as is. At the other extreme, others argue that NHST should be totally abandoned. In the middle, yet others argue that NHST should be supplemented with measures of effect size. In this text, we have taken the middle road, believing that more information is a better choice.

This is where the notion of **effect size** comes in. While there are a number of different measures of effect size, the most commonly used measure is Cohen's δ (delta) or d (1988). For example, Cohen's d for the effect between groups is computed as follows:

$$d = \frac{\bar{Y}_1 - \bar{Y}_2}{s_p}$$

The numerator of the formula is the difference between the two sample means. The denominator is the pooled standard deviation. Interpreting, d indicates how many standard deviations the sample means differ. Thus if $d = 1.0$, the sample mean of group one is one standard deviation away from the sample mean of group two. Cohen has proposed the following subjective standards for the social and behavioral sciences as a convention for interpreting d : small effect size, $d = .2$; medium effect size, $d = .5$; large effect size, $d = .8$. Interpretation of effect size should always be made first based on a comparison to similar studies; what is considered a ‘small’ effect using Cohen’s recommendations may actually be quite large in comparison to other related studies that have been conducted. In lieu of a comparison to other studies, such as in those cases where there are no or minimal related studies, then Cohen’s subjective standards may be appropriate. Computing confidence intervals for effect sizes is also valuable. The benefit in creating confidence intervals for effect size values is similar to that of creating confidence intervals for parameter estimates—confidence intervals for the effect size provide an added measure of precision that is not obtained from knowledge of the effect size alone.

2.2 FOUNDATIONAL UNIVARIATE STATISTICS

You are likely familiar with a number of univariate statistics that will be useful in understanding and describing your data in the context of multivariate statistics. As with the fundamental concepts, this section is not meant to serve as either a comprehensive review or an exhaustive summary of foundational univariate statistics. Rather, this section is meant to refresh and reacquaint you with some key univariate statistics that will be useful as we move into multivariate statistics.

2.2.1 Histogram

A graphical tool appropriate for data that are at least ordinal (i.e., ordinal, interval, or ratio) is the **histogram** (see Figure 2.1). Because the data are at least theoretically continuous (even though they may be measured in whole numbers), the main difference in the histogram (as compared to the bar graph) is that the bars touch one another, much like intervals touching one another as real limits. Along the X axis we plot the values of the variable X and along the Y axis the frequencies for each interval. The height of the bar corresponds to the number of frequencies for a particular value (or group of values in the case of grouped frequencies) of X . The histogram is a helpful tool for understanding the distributional shape.

2.2.2 Box-and-Whisker Plot

A simplified form of the frequency distribution is the **box-and-whisker plot** (often referred to simply as a ‘box plot’) (see Figure 2.2). The **box** in the center of the figure displays the middle 50% of the distribution of scores. The bottom, or **hinge**, of the box represents the 25th percentile (or Q_1). The top, or hinge, of the box represents the 75th percentile (or Q_3). The thick middle vertical line in the box represents the 50th

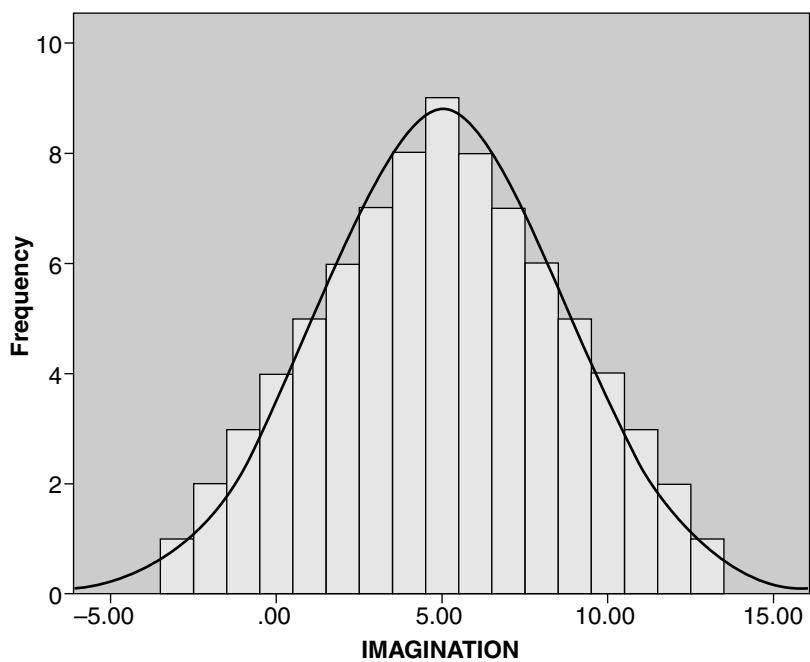


FIGURE 2.1

Histogram

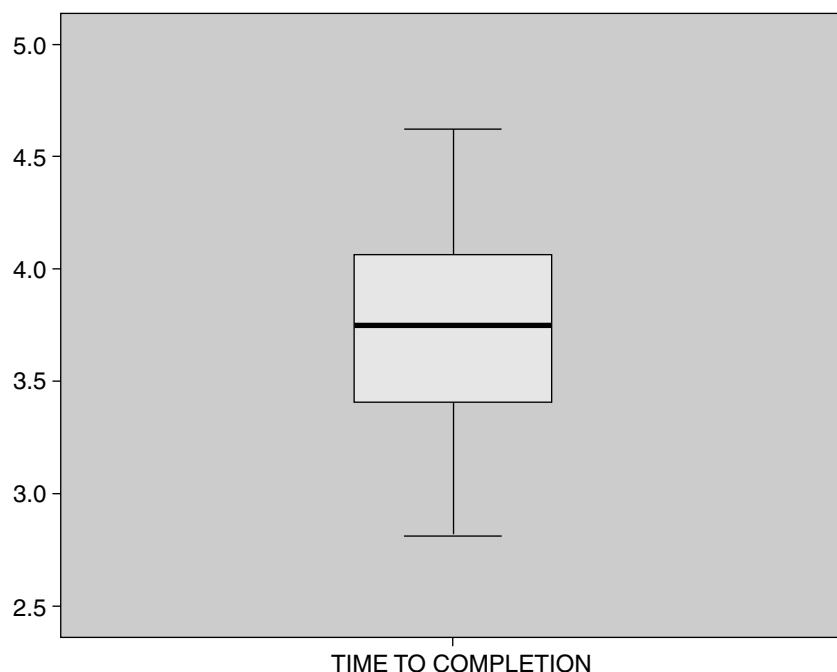


FIGURE 2.2

Box-and-Whisker Plot

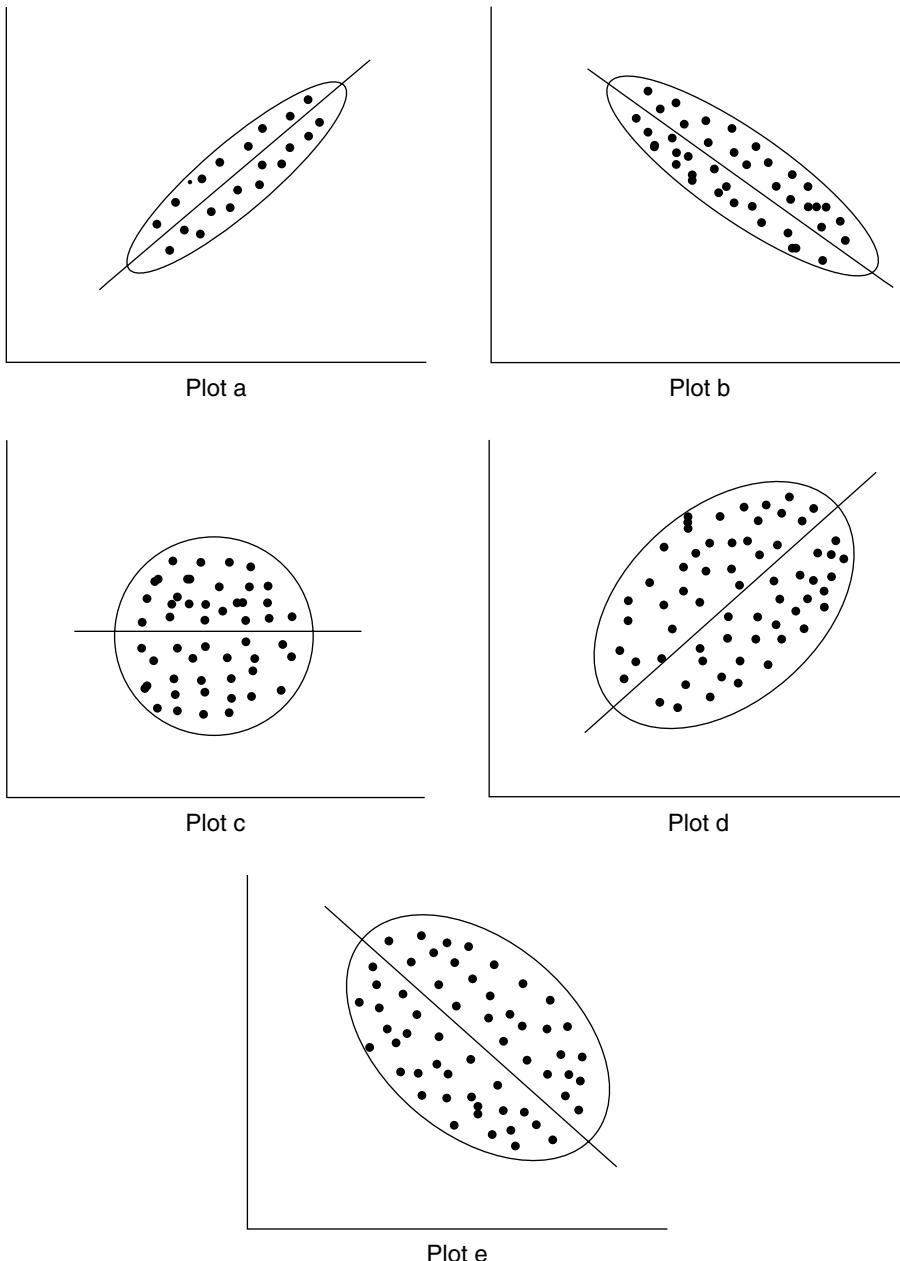
percentile (or Q_2). The lines extending from the box are known as the **whiskers**. The purpose of the whiskers is to display data outside of the middle 50%. The bottom whisker can extend down to the lowest score, or to the 5th or the 10th percentile to display more extreme low scores, and the top whisker correspondingly can extend up to the highest score, or to the 95th or 90th percentile to display more extreme high scores. The choice of where to extend the whiskers is the preference of the researcher and/or the software. Scores that fall beyond the end of the whiskers, known as **outliers** due to their extremeness relative to the bulk of the distribution, are often displayed by dots and/or asterisks. Box-and-whisker plots can be used to examine distributional attributes, such as skewness, kurtosis, and outliers.

2.2.3 Scatterplot

A graphical method for depicting the relationship between two variables (X, Y) is to plot the pair of scores on X and Y for each unit (e.g., individual) on a two-dimensional figure known as a **scatterplot** (or scattergram) (see Figure 2.3). Each individual has two scores in a two-dimensional coordinate system, denoted by (X, Y) . The X axis (the horizontal axis or abscissa) represents the values for variable X and the Y axis (the vertical axis or ordinate) represents the values for variable Y . Each point on the scatterplot represents a pair of scores (X, Y) for a particular individual. In essence, the scatterplot is actually a bivariate frequency distribution.

The scatterplot allows the researcher to evaluate both the direction and the strength of the relationship among X and Y . The **direction** of the relationship has to do with whether the relationship is positive or negative. A positive relationship occurs when scores on variable X increase (from left to right), and scores on variable Y also increase (from bottom to top). Examples of different scatterplots are shown in Figure 2.3. Figures 2.3 (a) and (d) both display positive relationships. A negative relationship, sometimes called an inverse relationship, occurs when scores on variable X increase (from left to right) as scores on variable Y decrease (from top to bottom). Figures 2.3 (b) and (e) are examples of negative relationships. There is no relationship between X and Y when for a large value of X , a large or a small value of Y can occur, and for a small value of X , a large or a small value of Y can also occur. In other words, X and Y are not related, as shown in Figure 2.3 (c).

The **strength** of the relationship among X and Y is determined by the scatter of the points (hence the name scatterplot). If the scatter is such that the points tend to fall close to the line, then this is indicative of a strong relationship between X and Y . Both Figures 2.3 (a) and (b) denote strong relationships. If the scatter is such that the points are widely scattered around the line, then this is indicative of a weak relationship between X and Y . Both Figures 2.3 (d) and (e) denote weak relationships. To summarize Figure 2.3, part (a) represents a strong positive relationship, part (b) a strong negative relationship, part (c) no relationship, part (d) a weak positive relationship, and part (e) a weak negative relationship. Thus, the scatterplot is useful for providing a quick visual indication of the nature of the relationship between variables X and Y .

**FIGURE 2.3**

Examples of Possible Scatterplots

2.2.4 Measures of Central Tendency

One method for summarizing a set of scores is to construct a single index or value that can somehow be used to represent the entire collection of scores. In this section, we consider the three most popular indices, known as **measures of central tendency**. Although other indices exist, the most popular ones are the mode, the median, and the mean.

2.2.4.1 Mode

The simplest method to use for measuring central tendency is the **mode**. The mode is defined as that value in a distribution of scores that occurs most frequently. In terms of the characteristics of the mode, the first characteristic of the mode is that it is simple to obtain. The mode is often used as a quick-and-dirty method for reporting central tendency. The second characteristic is that the mode does not always have a unique value. This is generally a disadvantage if the goal is to have a single index that could be used to represent the collection of scores. The mode cannot guarantee a single index. Third, the mode is not a function of all the scores in the distribution, and this is generally a disadvantage. The mode is strictly determined by which score or interval contains the most frequencies. Also, the location or value of the other scores is not taken into account. The fourth characteristic of the mode is that it is difficult to deal with mathematically. For example, the mode is not very stable from one sample to another, especially with small samples. We could have two nearly identical samples except for one score, which can alter the mode. A fifth and final characteristic is that the mode can be used with any type of measurement scale, from nominal to ratio, and is the only measure of central tendency appropriate for nominal data.

2.2.4.2 Median

The **median** is that score which divides a distribution of scores into two equal parts. In other words, one-half of the scores fall below the median and one-half of the scores fall above the median. The median is also known as the 50th percentile or Q_2 . The general characteristics of the median are as follows. First, the median is not influenced by extreme scores (scores far away from the middle of the distribution are known as **outliers**). Because the median is defined conceptually as the middle score, the actual size of an extreme score is not relevant. This characteristic is an advantage, particularly when extreme scores are observed. A second characteristic is the median is not a function of all of the scores. Because we already know that the median is not influenced by extreme scores, we know that the median does not take such scores into account. As you probably surmised, this characteristic is generally thought to be a disadvantage. If you really think about the first two characteristics, no measure could possibly possess both. That is, if a measure is a function of all of the scores, then extreme scores must also be taken into account. If a measure does not take extreme scores into account, like the median, then it cannot be a function of all of the scores. A third characteristic is that the median is difficult to deal with mathematically, a disadvantage as with the mode. The median is somewhat unstable from sample to sample, especially with small samples. As a fourth characteristic, the median always has a unique value, another advantage. This is unlike the mode, which does not always have a unique value. Finally, the fifth characteristic of the median is that it can be used with all types of measurement scales except the nominal. Nominal data cannot be ranked, and thus percentiles and the median are inappropriate.

2.2.4.3 Mean

The final measure of central tendency to be considered is the **mean**, sometimes known as the arithmetic mean or “average” (although the term average is used rather loosely by laypeople). Statistically, we define the mean as the sum of all of the scores divided by the number of scores.

The population mean is denoted by μ (Greek letter mu) and computed as follows:

$$\mu = \frac{\sum_{i=1}^N X_i}{N}$$

For sample data, the sample mean is denoted by \bar{X} (read “X bar”) and computed as follows:

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n}$$

Here are the general characteristics of the mean. First, the mean is a function of every score, a definite advantage in terms of a measure of central tendency representing all of the data. If you look at the numerator of the mean, you see that all of the scores are clearly taken into account in the sum. The second characteristic of the mean is that it is influenced by extreme scores. Because the numerator sum takes all of the scores into account, it also includes the extreme scores, which is a disadvantage. Third, the mean always has a unique value, another advantage. Fourth, the mean is easy to deal with mathematically. The mean is the most stable measure of central tendency from sample to sample, and because of that is the measure most often used in inferential statistics. Finally, the fifth characteristic of the mean is that it is only appropriate for interval and ratio measurement scales. This is because the mean implicitly assumes equal intervals, which of course the nominal and ordinal scales do not possess.

2.2.4.4 Summary of Measures of Central Tendency

To summarize the measures of central tendency then:

1. The mode is the only appropriate measure for nominal data.
2. The median and mode are both appropriate for ordinal data (and conceptually the median fits the ordinal scale as both deal with ranked scores).
3. All three measures are appropriate for interval and ratio data.

2.2.5 Measures of Dispersion

Another method for summarizing a set of scores is to construct an index or value that can be used to describe the amount of spread amongst the collection of scores. In other

words, we need measures that can be used to determine whether the scores fall fairly close to the central tendency measure, are fairly well spread out, or are somewhere in between. In this section we consider **measures of dispersion** (i.e., the extent to which the scores are dispersed or spread out) that are most applicable to multivariate statistics, including the variance and the standard deviation.

The **population variance**, which is denoted as σ^2 (i.e., sigma squared), is computed by the following formula:

$$\sigma^2 = \frac{\sum_{i=1}^N (X_i - \mu)^2}{N}$$

Conceptually, the variance is a measure of the area of a distribution. That is, the more spread out the scores, the more area or space the distribution takes up and the larger is the variance. The variance may also be thought of as an average distance from the mean. The variance has nice mathematical properties and is useful for deriving other statistics, such as inferential statistics.

The two quantities derived by the summation operations in the numerator are computed in much different ways and generally yield different values.

There are a few bothersome aspects about the variance. Say you are measuring the height of children in inches. The raw scores are measured in terms of inches, the mean is measured in terms of inches, but the variance is measured in terms of inches squared. Squaring the scale is bothersome to some as the scale is no longer in the original units of measure, but rather a squared unit of measure—making interpretation a bit difficult. To generate a deviational measure in the original scale of inches, we can take the square root of the variance. This is known as the **standard deviation**. The population standard deviation is defined as the positive square root of the population variance and is denoted by σ (i.e., $\sigma = +\sqrt{\sigma^2}$). The standard deviation, then, is measured in the original scale of inches.

What are the major characteristics of the population variance and standard deviation? First, the variance and standard deviation are a function of every score, an advantage. Second, therefore, the variance and standard deviation are affected by extreme scores, a disadvantage. Third, the variance and standard deviation are only appropriate for interval and ratio measurement scales. Like the mean, this is due to the implicit requirement of equal intervals. A fourth and final characteristic of the variance and standard deviation is they are quite useful for deriving other statistics, particularly in inferential statistics, another advantage.

Most of the time we are interested in computing the sample variance and standard deviation; we also often have large samples of data with multiple frequencies for many

of the scores. Here we consider these last aspects of the measures of dispersion. In order to obtain an unbiased sample estimate of the population variance, the following adjustments have to be made in the formula:

$$s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}$$

In terms of the notation, s^2 is the **sample variance**, n has been substituted for N , and \bar{X} has been substituted for μ . These changes are relatively minor and expected. The major change is in the denominator, where instead of N for the definitional formula we have $n - 1$, and instead of N^2 for the computational formula we have $n(n - 1)$. This turns out to be the correction that early statisticians discovered was necessary to obtain an unbiased estimate of the population variance.

2.3 FOUNDATIONAL BIVARIATE STATISTICS

As we begin to delve into multivariate statistics, you'll discover that bivariate statistics are scaffolding your learning as they lay the foundation for multivariate procedures. This section does not provide comprehensive coverage of all the bivariate statistics that pave the way for multivariate statistics but does provide a basic refresher to the most applicable bivariate statistics.

2.3.1 Independent and Dependent Samples *t* Test

Independent and dependent samples *t* tests are scaffolding to understanding, for example, analysis of variance and later multivariate analysis of variance—and their repeated measures counterparts. Two samples are **independent** when the method of sample selection is such that those individuals selected for sample 1 do not have any relationship to those individuals selected for sample 2. In other words, the selection of individuals to be included in the two samples are unrelated or uncorrelated such that they have absolutely nothing to do with one another. The independence condition leads us to consider, for example, the **independent samples *t* test**. (This should not, however, be confused with the assumption of independence.)

Two samples are **dependent** when the method of sample selection is such that those individuals selected for sample 1 *do* have a relationship to those individuals selected for sample 2. In other words, the selections of individuals to be included in the two samples *are* related or correlated. You might think of the samples as being selected simultaneously such that there are actually pairs of individuals. The dependence condition leads us to consider the **dependent samples *t* test**, alternatively known as the **correlated samples *t* test** or the **paired samples *t* test**.

The test statistic for the independent t test is known as t and is denoted by the following formula:

$$t = \frac{\bar{Y}_1 - \bar{Y}_2}{S_{\bar{Y}_1 - \bar{Y}_2}}$$

where \bar{Y}_1 and \bar{Y}_2 are the means for sample 1 and sample 2, respectively, and $S_{\bar{Y}_1 - \bar{Y}_2}$ is the **standard error of the difference between two means**. This standard error is the standard deviation of the sampling distribution of the difference between two means and is computed as follows:

$$S_{\bar{Y}_1 - \bar{Y}_2} = s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

where s_p is the pooled standard deviation computed as

$$s_p = \sqrt{\frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}}$$

and where s_1^2 and s_2^2 are the sample variances for groups 1 and 2, respectively, and n_1 and n_2 are the sample sizes for groups 1 and 2, respectively. Conceptually, the standard error $S_{\bar{Y}_1 - \bar{Y}_2}$ is a pooled standard deviation weighted by the two sample sizes; more specifically, the two sample variances are weighted by their respective sample sizes and then pooled. The assumptions for the independent t test include normality, independence, and homogeneity of variance.

For the dependent t test, although there are several methods for computing the test statistic t , the most direct method and the one most closely aligned conceptually with the one-sample t test is as follows:

$$t = \frac{\bar{d}}{S_{\bar{d}}}$$

where \bar{d} is the mean difference, and $S_{\bar{d}}$ is the standard error of the mean difference. The standard error of the mean difference is computed by

$$S_{\bar{d}} = \frac{s_d}{\sqrt{n}}$$

where s_d is the standard deviation of the difference scores (i.e., like any other standard deviation, only this one is computed from the difference scores rather than raw scores), and n is the total number of pairs. The dependent t test shares the same assumptions as the independent t test: normality, independence, and homogeneity of variance.

2.3.2 Analysis of Variance

Analysis of variance allows us to compare the means of more than two independent samples, and it does this by partitioning the variation. The partitioning of the sums of squares in ANOVA is also an important concept in regression analysis. In part, this is because ANOVA and regression are both forms of the same general linear model (GLM). Let us begin with the total sum of squares in Y , denoted as SS_{total} . The term SS_{total} represents the amount of total variation in Y . The next step is to partition the total variation into variation between the groups (i.e., the categories or levels of the independent variable), denoted by SS_{betw} , and variation within the groups (i.e., units or cases within each category or level of the independent variable), denoted by SS_{with} . In the one-factor analysis of variance, we therefore partition SS_{total} as follows:

$$SS_{\text{total}} = SS_{\text{betw}} + SS_{\text{with}}$$

or

$$\sum_{i=1}^n \sum_{j=1}^J (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^n \sum_{j=1}^J (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{i=1}^n \sum_{j=1}^J (Y_{ij} - \bar{Y}_{.j})^2$$

where SS_{total} is the total sum of squares due to variation among all of the observations without regard to group membership, SS_{betw} is the between-groups sum of squares due to the variation between the groups, and SS_{with} is the within-groups sum of squares due to the variation within the groups combined across groups. We refer to this particular formulation of the partitioned sums of squares as the **definitional (or conceptual) formula**, because each term literally defines a form of variation. The assumptions of ANOVA are concerned with independence, homogeneity of variance, and normality.

2.3.2.1 ANOVA Summary Table

An important result of the analysis is the **ANOVA summary table**. The purpose of the summary table is to literally summarize the analysis of variance. A general form of the summary table is shown in Table 2.2. The first column lists the sources of variation in the model. As we already know, in the one-factor model the total variation is partitioned into between-groups variation and within-groups variation. The second column notes the sums of squares terms computed for each source (i.e., SS_{betw} , SS_{with} , and SS_{total}).

TABLE 2.2

Analysis of Variance Summary Table

Source	SS	df	MS	F
Between groups	SS_{betw}	$J - 1$	MS_{betw}	$MS_{\text{betw}} / MS_{\text{with}}$
Within groups	SS_{with}	$N - J$	MS_{with}	
Total	SS_{total}	$N - 1$		

The third column gives the degrees of freedom for each source. Recall that, in general, the degrees of freedom have to do with the number of observations that are free to vary. For example, if a sample mean and all of the sample observations except for one are known, then the final observation is not free to vary. That is, the final observation is predetermined to be a particular value. For instance, say the mean is 10 and there are three observations: 7, 11, and an unknown observation. Based on that information, first, the sum of the three observations must be 30 for the mean to be 10. Second, the sum of the known observations is 18. Therefore the unknown observation must be 12. Otherwise the sample mean would not be exactly equal to 10.

For the between-groups source, the definitional formula is concerned with the deviation of each group mean from the overall mean. There are J group means (where J represents the number of groups or categories or levels of the independent variable), so the df_{betw} (also known as the degrees of freedom numerator) must be $J - 1$. Why? If we have J group means and we know the overall mean, then only $J - 1$ of the group means are free to vary. In other words, if we know the overall mean and all but one of the group means, then the final unknown group mean is predetermined. For the within-groups source, the definitional formula is concerned with the deviation of each observation from its respective group mean. There are n observations (i.e., cases or units) in each group; consequently, there are $n - 1$ degrees of freedom in each group and J groups. Why are there $n - 1$ degrees of freedom in each group? If there are n observations in each group, then only $n - 1$ of the observations are free to vary. In other words, if we know one group mean and all but one of the observations for that group, then the final unknown observation for that group is predetermined. There are J groups, so the df_{with} (also known as the degrees of freedom denominator) is $J(n - 1)$, or more simply as $N - J$. Thus, we lose one degree of freedom for each group. For the total source, the definitional formula is concerned with the deviation of each observation from the overall mean. There are N total observations; therefore the df_{total} must be $N - 1$. Why? If there are N total observations and we know the overall mean, then only $N - 1$ of the observations are free to vary. In other words, if we know the overall mean, and all but one of the N observations, then the final unknown observation is predetermined.

Why is the number of degrees of freedom important in the analysis of variance? Suppose two researchers have conducted similar studies, except Researcher A uses 20 observations per group and Researcher B uses 10 observations per group. Each researcher obtains an SS_{with} value of 15. Would it be fair to say that this particular result for the two studies is the same? Such a comparison would be unfair because SS_{with} is influenced by the number of observations per group. A fair comparison would be to weight the SS_{with} terms by their respective number of degrees of freedom. Similarly, it would not be fair to compare the SS_{betw} terms from two similar studies based on different numbers of groups. A fair comparison would be to weight the SS_{betw} terms by their respective number of degrees of freedom. The method of weighting a sum of squares term by the respective number of degrees of freedom on which it is based yields what is called a **mean squares** term. Thus, $MS_{\text{betw}} = SS_{\text{betw}} / df_{\text{betw}}$ and $MS_{\text{with}} = SS_{\text{with}} / df_{\text{with}}$, as shown in the fourth column of Table 2.2. They

are referred to as mean squares because they represent a summed quantity that is weighted by the number of observations used in the sum itself, like the mean. The mean squares terms are also variance estimates because they represent the sum of the squared deviations from a mean divided by their degrees of freedom, like the sample variance s^2 .

The last column in the ANOVA summary table, the F value, is the summary test statistic of the summary table. The F value is computed by taking the ratio of the two mean squares or variance terms. Thus, for the one-factor ANOVA fixed-effects model, the F value is computed as $F = MS_{\text{betw}} / MS_{\text{with}}$. When developed by Sir Ronald A. Fisher in the 1920s, this test statistic was originally known as the variance ratio because it represents the ratio of two variance estimates. Later, the variance ratio was renamed the F ratio by George W. Snedecor (who worked out the table of F values, discussed momentarily) in honor of Fisher (F for Fisher).

The F ratio tells us whether there is more variation *between* groups than there is *within* groups, which is required if we are to reject H_0 . Thus, if there is more variation *between* groups than there is *within* groups, then MS_{betw} will be larger than MS_{with} . As a result of this, the F ratio of $MS_{\text{betw}} / MS_{\text{with}}$ will be greater than 1. If, on the other hand, the amount of variation *between* groups is about the same as there is *within* groups, then MS_{betw} and MS_{with} will be about the same, and the F ratio will be approximately 1. Thus, we want to find large F values in order to reject the null hypothesis. The degrees of freedom are $df_{\text{betw}} = J - 1$ for the numerator of the F ratio and $df_{\text{with}} = N - J$ for the denominator of the F ratio. This is the omnibus F test, which, again, simply provides evidence of the extent to which there is at least one statistically significant mean difference between the groups.

In the case of a statistically significant omnibus F test, some multiple comparison procedure should be used to determine where the mean differences are in the groups. When there are only two groups, it is obvious where the mean difference falls, that is, between groups 1 and 2. A researcher can simply look at the descriptive statistics to determine which group had the higher mean relative to the other group. Consider the situation where there are more than two samples (e.g., three types of interventions), and H_0 has already been rejected in the omnibus test. This situation is one where a **multiple comparison procedure** (MCP) would be quite informative to determine which means or combination of means are different. Also consider the situation where the researcher is not even interested in the ANOVA omnibus test, but is only interested in comparisons involving particular means (e.g., certain medications are more effective than a placebo). This is a situation where an MCP is useful for evaluating those specific comparisons.

2.3.3 Two-Factor ANOVA

The first characteristic of the two-factor ANOVA model should be obvious by now; this model considers the effect of two factors or independent variables on a dependent

variable. Each factor consists of two or more levels (or categories). This yields what we call a **factorial design** because more than a single factor is included. We see then that the two-factor ANOVA is an extension of the one-factor ANOVA. Why would a researcher want to complicate things by considering a second factor? Three reasons come to mind. First, the researcher may have a genuine interest in studying the second factor. Rather than studying each factor separately in two analyses, the researcher includes both factors in the same analysis. This allows a test of not only the effect of each individual factor, known as **main effects**, but also the effect of both factors collectively. This latter effect is known as an **interaction effect** and provides information about whether the two factors are operating independent of one another (i.e., no interaction exists) or whether the two factors are operating together to produce some additional impact (i.e., an interaction exists). If two separate analyses were conducted, one for each independent variable, no information would be obtained about the interaction effect. As becomes evident, assuming a factorial ANOVA with two independent variables, the researcher will test three hypotheses: one for each factor or main effect individually and a third for the interaction between the factors. Factorial ANOVA models with more than two independent variables will, accordingly, test for additional main effects and interactions.

A second reason for including an additional factor is an attempt to reduce the error (or within groups) variation, which is variation that is unexplained by the first factor. The use of a second factor provides a more precise estimate of error variance. For this reason, a two-factor design is generally more powerful than two one-factor designs, as the second factor and the interaction serve to control for additional extraneous variability. A third reason for considering two factors simultaneously is to provide greater generalizability of the results and to provide a more efficient and economical use of observations and resources. Thus, the results can be generalized to more situations, and the study will be more cost efficient in terms of time and money.

In addition, for the two-factor ANOVA every level of the first factor (hereafter known as factor A) is paired with every level of the second factor (hereafter known as factor B). In other words, every combination of factors A and B is included in the design of the study, yielding what is referred to as a **fully crossed design**. If some combinations are not included, then the design is not fully crossed and may form some sort of a nested design. Individuals (or objects or subjects) are randomly assigned to one combination of the two factors. In other words, each individual responds to only one combination of the factors. If individuals respond to more than one combination of the factors, this would be some sort of repeated measures design. Considering models where all factors are fixed, the overall design is known as a fixed-effects model. If one or both factors are random, then the design is not a fixed-effects model. It is also a condition for factorial ANOVA that the dependent variable is measured at least at the interval level and the independent variables are categorical (either nominal or ordinal).

For simplicity sake, we impose the restriction that the number of observations is the same for each factor combination. This yields what is known as an orthogonal design, where the effects due to the factors (separately and collectively) are independent or

unrelated. In addition, there must be at least two observations per factor combination so as to have within-groups variation.

In summary, the characteristics of the two-factor analysis of variance fixed-effects model are as follows: (a) two independent variables (both of which are categorical) each with two or more levels, (b) the levels of both independent variables are fixed by the researcher, (c) subjects are randomly assigned to only one combination of these levels, (d) the two factors are fully crossed, and (e) the dependent variable is measured at least at the interval level. In the context of experimental design, the two-factor analysis of variance is often referred to as the **completely randomized factorial design**.

2.3.4 Covariance

Covariance is one type of statistical method for measuring the relationship among variables X and Y . The covariance conceptually is the shared variance (or covariance) among X and Y . The covariance and correlation share commonalities, as the correlation is simply the standardized covariance. The population covariance is denoted by σ_{XY} and the conceptual formula is given as follows:

$$\sigma_{XY} = \frac{\sum_{i=1}^N (X_i - \mu_X)(Y_i - \mu_Y)}{N}$$

where X_i and Y_i are the scores for unit (i.e., individual) i on variables X and Y , respectively, and μ_X and μ_Y are the population means for variables X and Y , respectively. This equation looks similar to the formula for the variance, where deviation scores from the mean are computed for each individual. The formula for the covariance is essentially an average of the paired deviation score products. If variables X and Y are positively related, then the deviation scores will tend to be of the same sign, their products will tend to be positive, and the covariance will be a positive value (i.e., $\sigma_{XY} > 0$). If variables X and Y are negatively related, then the deviation scores will tend to be of opposite signs, their products will tend to be negative, and the covariance will be a negative value (i.e., $\sigma_{XY} < 0$). Finally, if variables X and Y are not related, then the deviation scores will consist of both the same and opposite signs, their products will be both positive and negative and sum to zero, and the covariance will be a zero value (i.e., $\sigma_{XY} = 0$).

The sample covariance is denoted by s_{XY} , and the formula becomes as follows:

$$s_{XY} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

where \bar{X} and \bar{Y} are the sample means for variables X and Y , respectively, and n is sample size. Note that the denominator becomes $n - 1$ so as to yield an unbiased sample estimate of the population covariance (i.e., similar to what we did in the sample variance situation).

Like the variance, the value of the covariance depends on the scales of the variables involved. Thus, interpretation of the magnitude of a single covariance is difficult, as it can take on literally any value. We see shortly that the correlation coefficient takes care of this problem. The covariance is commonly utilized in multivariate techniques, such as structural equation modeling and multilevel modeling.

2.3.5 Pearson Product-Moment Correlation Coefficient

The **Pearson product-moment correlation coefficient**, developed by the famous statistician Karl Pearson, and simply referred to as the Pearson here, is another method for measuring the relationship among X and Y . The *sign* of the Pearson denotes the direction of the relationship (e.g., positive or negative), and the *value* of the Pearson denotes the strength of the relationship. The Pearson falls on a scale from -1.00 to $+1.00$, where -1.00 indicates a perfect negative relationship, 0 indicates no relationship, and $+1.00$ indicates a perfect positive relationship. Values near $.50$ or $-.50$ are considered as moderate relationships, values near 0 as weak relationships, and values near $+1.00$ or -1.00 as strong relationships (although these are subjective terms). Cohen (1988) also offers recommendations for interpreting the value of the correlation. As you may see as you read more statistics and research methods textbooks, there are other guidelines offered for interpreting the value of the correlation.

The Pearson can be considered in several different forms, where the population value is denoted by ρ_{XY} (rho) and the sample value by r_{XY} . One form of the Pearson is in terms of the covariance and the standard deviations and is given as:

$$\rho_{XY} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

This form is useful when the covariance and standard deviations are already known.

The sample correlation is denoted by r_{XY} . The formulas are essentially the same for the sample correlation, r_{XY} , and the population correlation, ρ_{XY} , except that n is substituted for N .

Unlike the sample variance and covariance, the sample correlation has no correction for bias.

2.3.6 Simple Linear Regression

It is helpful to understand simple linear regression, as it is the foundation for multiple procedures including multiple regression, logistic regression, and multilevel modeling. The **sample regression model** for predicting Y from X is computed as:

$$Y_i = b_{YX} X_i + a_{YX} + e_i$$

where Y and X are as before (i.e., the dependent and independent variables respectively), b_{YX} is the sample slope for Y predicted by X , a_{YX} is the sample intercept for Y predicted by X , e_i is sample residuals or errors of prediction (the part of Y_i not predictable from X_i), and i represents an index for a case (an individual or object). The index i can take on values from 1 to n , where n is the size of the sample, and is written as $i = 1, \dots, n$.

The **sample prediction model** is computed as follows:

$$Y'_i = b_{YX} X_i + a_{YX}$$

where Y'_i is the predicted value of Y for a specific value of X . We define the sample prediction error as the difference between the *actual score* obtained by individual i (i.e., Y_i) and the *predicted score* based on the X score for that individual (i.e., Y'_i). In other words, the residual is that part of Y that is *not* predicted by X . The goal of the prediction model is to include an independent variable X that minimizes the residual; this means that the independent variable does a nice job of predicting the outcome. Computationally, the residual (or error) is computed as:

$$e_i = Y_i - Y'_i$$

The difference between the regression and prediction models is the same as previously discussed, except now we are dealing with a sample rather than a population.

The sample slope (b_{YX}) and intercept (a_{YX}) can be determined by

$$b_{YX} = \left(r_{XY} \frac{s_Y}{s_X} \right)$$

and

$$a_{YX} = \bar{Y} - b_{YX} \bar{X}$$

where s_Y and s_X are the sample standard deviations for Y and X respectively, r_{XY} is the sample correlation between X and Y (again the Pearson correlation coefficient, rho), and \bar{Y} and \bar{X} are the sample means for Y and X , respectively. The sample slope (b_{YX}) is referred to alternately as (a) the expected or predicted change in Y for a one-unit change in X , and (b) the unstandardized or raw regression coefficient. The sample intercept (a_{YX}) is referred to alternately as (a) the point at which the regression line intersects (or crosses) the Y axis, and (b) the value of Y when X is zero.

Let us interpret example slope and intercept values. A slope of 0.5250 means that if your score on the independent variable is increased by one point, then your predicted score on dependent variable will be increased by 0.5250 points or about half a point. An intercept of 8.8625 means that if your score on the independent variable is zero, then your score on the dependent variable is 8.8625. The sample simple linear regression model, given these values, becomes

$$Y_i = b_{YX} X_i + a_{YX} + e_i = .5250 X_i + 8.8625 + e_i$$

If your score on the independent variable is 63, then your predicted score on the outcome is the following:

$$Y'_i = .5250(63) + 8.8625 = 41.9375$$

Thus, based on the prediction model developed, your predicted score on the dependent variable is approximately 42; however, as we know, predictions are generally not perfect.

2.3.6.1 Standardized Regression Model

Up until now, the computations in simple linear regression have involved the use of raw scores. For this reason, we call this the *unstandardized regression model*. The slope estimate is an unstandardized or raw regression slope because it is the predicted change in Y raw score units for a one raw score unit change in X . We can also express regression in standard z score units for both X and Y as

$$z(X_i) = \frac{X_i - \bar{X}}{s_X}$$

and

$$z(Y_i) = \frac{Y_i - \bar{Y}}{s_Y}$$

In both cases, the numerator is the difference between the observed score and the mean, and the denominator is the standard deviation (and dividing by the standard deviation standardizes the value). The means and variances of both standardized variables (i.e., z_X and z_Y) are 0 and 1, respectively.

The sample standardized linear prediction model becomes the following where $z(Y'_i)$ is the standardized predicted value of Y :

$$z(Y'_i) = b_{YX}^* z(X_i) = r_{XY} z(X_i)$$

Thus the standardized regression slope, b_{YX}^* , sometimes referred to as a **beta weight**, is equal to r_{XY} . No intercept term is necessary in the prediction model as the mean of the z scores for both X and Y is zero (i.e., $a_{YX}^* = \bar{z}_Y - b_{YX}^* \bar{z}_X = 0$). In summary, *the standardized slope is equal to the correlation coefficient and the standardized intercept is equal to zero*.

A slope of .9177 would be interpreted as the expected increase in the dependent variable in z score (i.e., standardized score) units for a one z score (i.e., standardized score) unit increase in the independent variable. A one z score unit increase is also the same as a one standard deviation increase because the standard deviation of z is equal to one.

When should you consider use of the standardized versus unstandardized regression analyses? According to Pedhazur (1997), the standardized regression slope b^* is not very stable from sample to sample. Thus, in simple regression most researchers prefer the use of b .

2.3.6.2 Prediction Errors

Previously we mentioned that perfect prediction of Y from X is extremely unlikely, only occurring with a perfect correlation between X and Y (i.e., $r_{XY} = \pm 1.0$). When developing the regression model, the values of the outcome, Y , are known. Once the slope and intercept have been estimated, we can then use the prediction model to predict the outcome (Y) from the independent variable (X) when the values of Y are unknown. We have already defined the predicted values of Y as Y' . In other words, a predicted value Y' can be computed by plugging the obtained value for X into the prediction model. It can be shown that $Y'_i = Y_i$ for all i only when there is perfect prediction. However, this is extremely unlikely in reality, particularly in simple linear regression using a single predictor.

We can determine a value of Y' for each of the i cases (individuals or objects) from the prediction model. In comparing the actual Y values to the predicted Y values, we obtain the *residuals* as the difference between the observed (Y_i) and predicted values (Y'_i), computed as follows:

$$e_i = Y_i - Y'_i$$

for all $i = 1, \dots, n$ individuals or objects in the sample. The residuals, e_i , are also known as **errors of estimate**, or **prediction errors**, and are that portion of Y_i that is not predictable from X_i . The residual terms are random values that are unique to each individual or object.

A regression example is shown graphically in the **scatterplot** of Figure 2.4, where the straight diagonal line represents the regression line. Individuals falling above the regression line have positive residuals (e.g., observation 1) (in other words, the difference between the observed score, represented as open circle 1 on the graph, is greater in value than the predicted value, which is represented by the regression line) and individuals falling below the regression line have negative residuals (e.g., observation 3) (in other words, the difference between the observed score, represented as open circle 3 on the graph, is less in value than the predicted value, which is represented by the regression line). The residual is, very simply, the vertical distance between the observed score (represented by the open circles or ‘dots’ in the scatterplot, Figure 2.4) and the regression line. In Figure 2.4, we see that half of the points fall above the regression line and half below the regression line. It can be shown that the mean of the residuals is always zero (i.e., $\bar{e} = 0$), as the sum of the residuals is always zero. This results from the fact that the mean of the observed criterion scores is equal to the mean of the predicted criterion scores (i.e., $\bar{Y} = \bar{Y}'$).

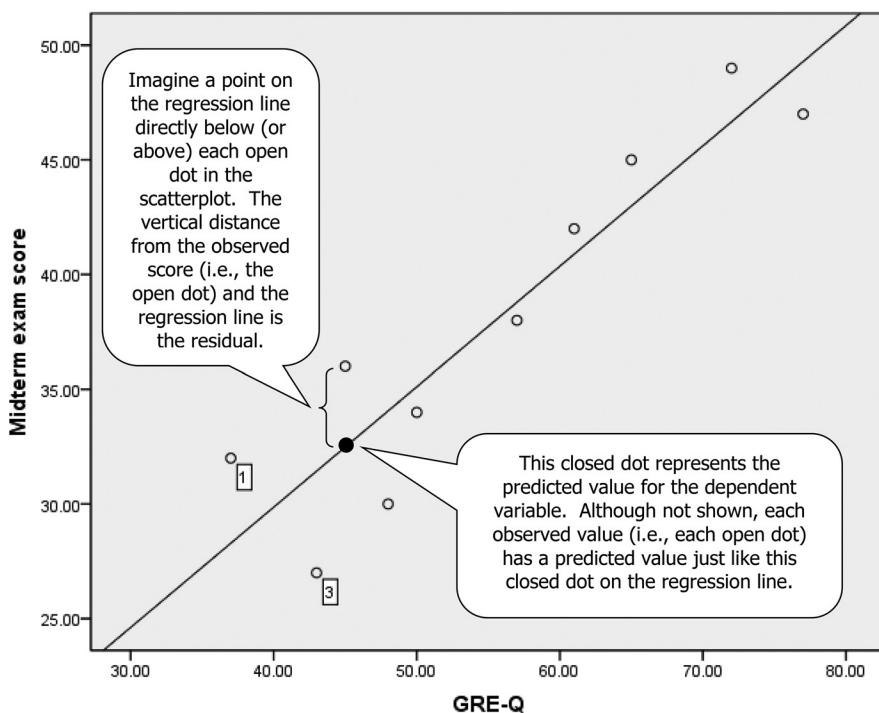


FIGURE 2.4
Regression Line Example

2.3.6.3 Least Squares Criterion

How was one particular method selected for determining the slope and intercept? Obviously, some standard procedure has to be used. Thus, there are statistical criteria that help us decide which method to use in determining the slope and intercept. The criterion usually used in linear regression analysis (and in all general linear models, for that matter) is the **least squares criterion**. According to the least squares criterion, the sum of the squared prediction errors or residuals is smallest. That is, we want to find that regression line, defined by a particular slope and intercept, which results in the smallest sum of the squared residuals (recall that the residual is the difference between the observed and predicted values for the outcome). Since the residual is the vertical difference between the observed and predicted value, the regression line is simply the line that minimizes that vertical distance. Given the value that we place on the accuracy of prediction, this is the most logical choice of a method for estimating the slope and intercept.

In summary then, the least squares criterion gives us a particular slope and intercept, and thus a particular regression line, such that the sum of the squared residuals is smallest. We often refer to this particular method for determining the slope and intercept as **least squares estimation**, because b and a represent sample estimates of the population parameters β and α obtained using the least squares criterion.

PROBLEMS

Conceptual Problems

1. In hypothesis testing, the probability of failing to reject H_0 when H_0 is false is denoted by
 - a. α
 - b. $1 - \alpha$
 - c. β
 - d. $1 - \beta$
2. The probability of making a Type II error when rejecting H_0 at the .05 level of significance is which one of the following?
 - a. 0
 - b. .05
 - c. between .05 and .95
 - d. .95
3. Which one of the following is a correct interpretation of d ?
 - a. alpha level
 - b. confidence interval
 - c. effect size
 - d. observed probability
 - e. power
4. Which one of the following is best used to visually examine the relationship between two variables?
 - a. Box-and-whisker plot
 - b. Histogram
 - c. Scatterplot
 - d. Simple regression
5. Which of the following is NOT an assumption of the independent t test?
 - a. normality
 - b. independence
 - c. equal sample sizes
 - d. homogeneity of variance
6. The mathematic ability of 10 preschool children was measured when they entered their first year of preschool and then again in the spring of their kindergarten year. To test for pre to post mean differences, which of the following tests would be used?
 - a. independent t test
 - b. dependent t test
 - c. simple regression
 - d. z test

7. The regression line for predicting salary of principals from cumulative GPA in graduate school is found to be $Y' = 35000 X + 37000$. What does the value of 37000 represent?
 - a. Average cumulative GPA
 - b. The criterion value
 - c. The mean salary of principals when cumulative GPA is zero
 - d. The standardized regression coefficient given an intercept of zero
8. If the relationship between two variables is linear, then which one of the following is correct?
 - a. all of the points must fall on a curved line
 - b. the relationship is best represented by a curved line
 - c. all of the points must fall on a straight line
 - d. the relationship is best represented by a straight line
9. In a one-factor ANOVA, H_0 asserts that
 - a. all of the population means are equal.
 - b. the between-groups variance estimate and the within-groups variance estimate are both estimates of the same population residual variance.
 - c. the within-groups sum of squares is equal to the between-groups sum of squares.
 - d. both a and b
10. For $J = 2$ and $\alpha = .05$, if the result of the independent t test is significant, then the result of the one-factor fixed-effects ANOVA is uncertain. True or false?

REFERENCES

- Cohen, J. (1988). Statistical power analysis for the behavioral sciences (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.
- Pedhazur, E. J. (1997). Multiple regression in behavioral research (3rd ed.). Fort Worth, TX: Harcourt Brace.

Chapter 3

DATA SCREENING

CHAPTER OUTLINE

3.1 Independence	36
3.1.1 Screening for Independence	37
3.2 Variance	41
3.2.1 Screening for Homoscedasticity	42
3.2.2 Screening for Homogeneity of Variance-Covariance	43
3.3 Normality	43
3.3.1 Screening for Univariate Normality	45
3.3.2 Screening for Multivariate Normality	47
3.4 Linearity	51
3.4.1 Screening for Linearity	51
3.5 Noncollinearity	52
3.5.1 Screening for Noncollinearity	53

KEY CONCEPTS

1. Homogeneity
2. Homoscedasticity
3. Kurtosis
4. Multicollinearity
5. Residual
6. Skewness

Before you ever run a test of inference, you should screen your data. The initial screen should be purely descriptive in nature and undertaken to make sure that all data points are legitimate. A frequency distribution, for example, will alert you to values that are outside the range of plausible (e.g., a value of 6 on a scale that ranges from 1 to 5 is a red flag). A histogram or boxplot will give a visual indication of the shape of the distribution. A scatterplot will help you understand how variables relate (e.g., linear or nonlinear). Beyond that, data screening is critical to examining the extent to which assumptions of inferential tests are met. As we all know, the confidence in interpreting the results of tests of inference can become low to nonexistent, depending on the extent to which assumptions are violated. Testing assumptions is a task that you are likely well acquainted with by this point in your statistical career. This chapter considers several assumptions that are common to many of the multivariate procedures covered in this text, such as independence, noncollinearity, and normality, among others. Please note that this chapter does not replace the coverage of these concepts within their respective chapter. Indeed, talking about testing for assumptions outside of the context of the respective statistical procedure is a bit like being told that you need to prepare a venue for a party but you know neither what type of venue it is nor what type of party you are hosting. Context is critical. Thus, although each chapter has coverage of the assumptions that are relevant for that particular procedure, this particular chapter pulls out some common data screening advice and provides a bit of scaffolding. Our objectives are that, by the end of this chapter, you will be able to understand and evaluate several assumptions that are common to multivariate procedures.

3.1 INDEPENDENCE

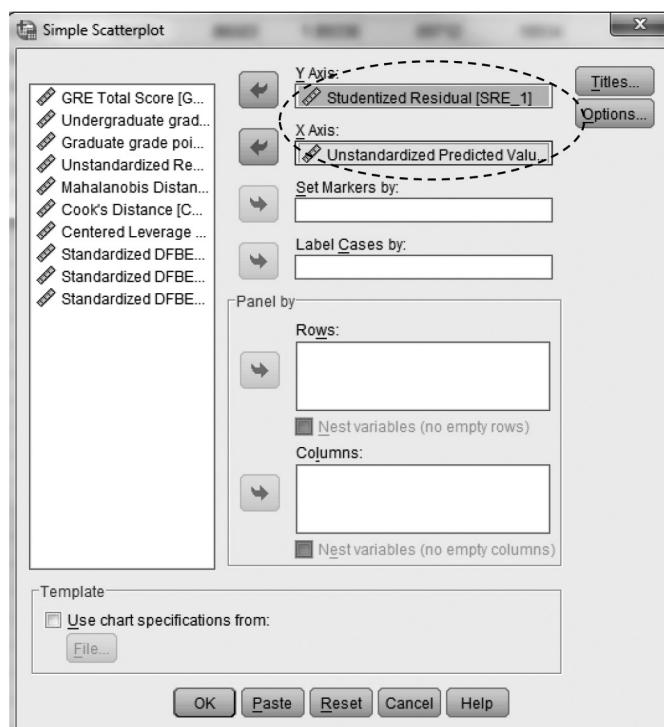
The first assumption is concerned with **independence** of the observations. We discuss this first as it is applicable broadly to inferential procedures, although testing for independence may be approached in different ways given different analytic procedures. As you may recall, there is a connection between sampling method and independence. Simple random sampling is defined as the process of selecting sample observations from a population so that each observation has an equal and independent probability of being selected. If the sampling process is truly random, then (a) each observation in the population has an equal chance of being included in the sample, and (b) each observation selected into the sample is independent of (or not affected by) every other selection. Independence implies that each observation is selected without regard to any other observation sampled. We also would fail to have equal and independent probability of selection if the sampling procedure employed was something other than a simple random sample—because it is only with a simple random sample that we have met the conditions (a) and (b) presented earlier. (Although there are statistical means to correct for nonsimple random samples, such as weighting and variance correction methods, they are beyond the scope of this textbook.) This concept of *independence* is an important assumption. If we have independence, then generalizations from the sample back to the population can be made (you may remember this as *external validity* which was likely introduced in your research methods course). Violations of this assumption can detrimentally impact standard error values and thus any resulting hypothesis tests. In particular, even small violations of this assumption can result in a quite dramatically increased actual alpha level as compared to nominal alpha level (Barcikowski, 1981;

Scariano & Davenport, 1987). At minimum, when there is dependence, tests should be conducted at a decreased alpha level (e.g., .01), given that the actual alpha will likely be higher. Lack of independence affects the estimated standard errors of the model. For serious violations, transformations or other estimation procedures can be considered.

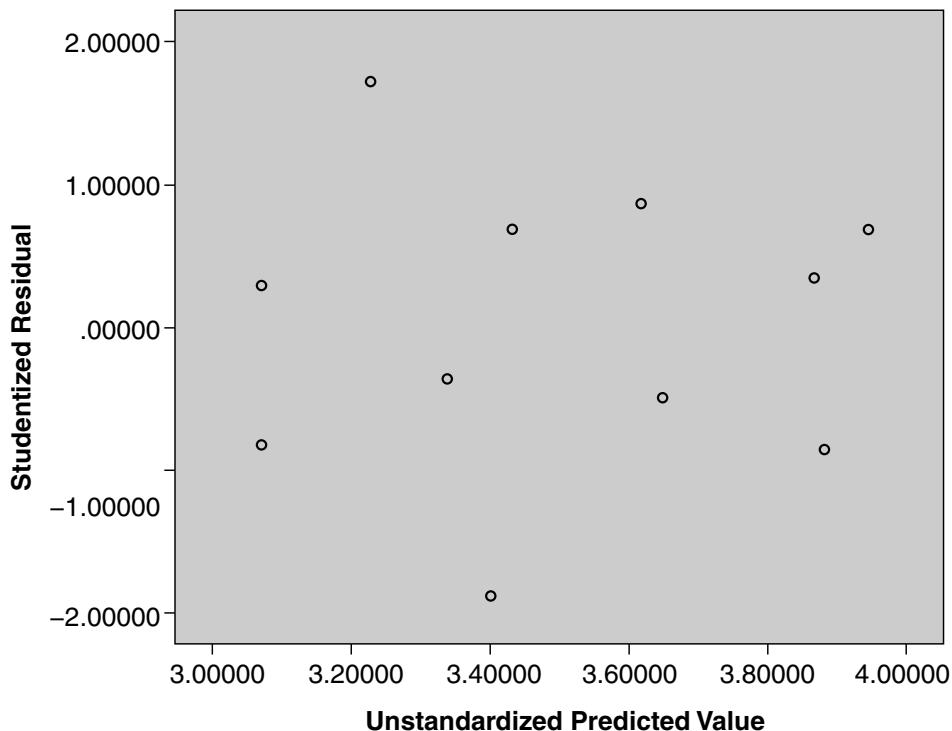
3.1.1 Screening for Independence

In regression based procedures, such as multiple regression, the simplest procedure for assessing independence is to examine residual plots versus the predicted values of the dependent variable and of residuals versus each independent variable (alternatively, one can look at plots of observed values of the dependent variable versus predicted values of the dependent variable and of observed values of the dependent variable versus each independent variable). If the independence assumption is satisfied, the residuals should fall into a random display of points. If the assumption is violated, the residuals will fall into some sort of pattern.

In this example, we will create scatterplots of the (1) residuals against unstandardized predicted values and (2) residuals against each independent variable, and will do so working with the studentized residuals. We will illustrate this using data from the illustration in the multiple regression chapter. From the “Simple Scatterplot” dialog screen, click the studentized residual variable and move it into the “Y Axis” box by clicking on the arrow. Click the unstandardized predicted values and move them into the “X Axis” box by clicking on the arrow. Then click “OK”. Repeat these steps to plot the studentized residual to each independent variable.

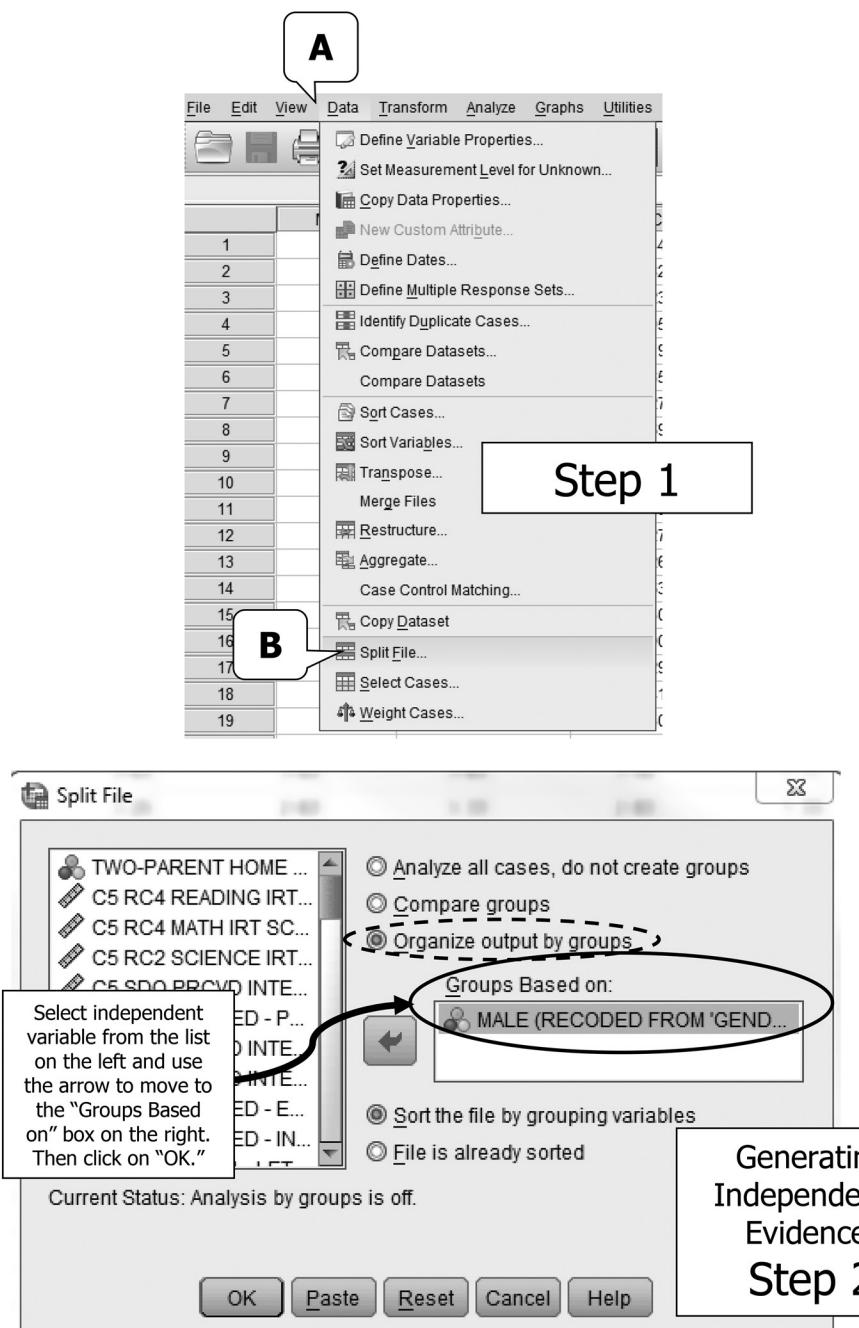


If the assumption of independence is met, the points should fall randomly within a band of -2.0 to $+2.0$. In this scatterplot, we have evidence of independence as all points for all graphs are within an absolute value of 2.0 and fall relatively randomly.



For procedures that examine mean group differences, like MANOVA, if subjects have been randomly assigned to conditions (or to the different combinations of the levels of the independent variables in the MANOVA), the assumption of independence has been met. In many cases, we use independent variables that do not allow random assignment, such as preexisting or self-selected characteristics. In those cases, we can plot residuals against levels of our independent variables in a scatterplot to get an idea of whether or not there are patterns in the data. Patterns, or lack thereof, provide an indication of whether we have met this assumption. Let's examine the data from our MANOVA example. Given that we have multiple independent variables in the factorial MANOVA, we will split the scatterplot by levels of one independent variable ('MALE') and then generate a bivariate scatterplot for 'TWOPARENT' by residual.

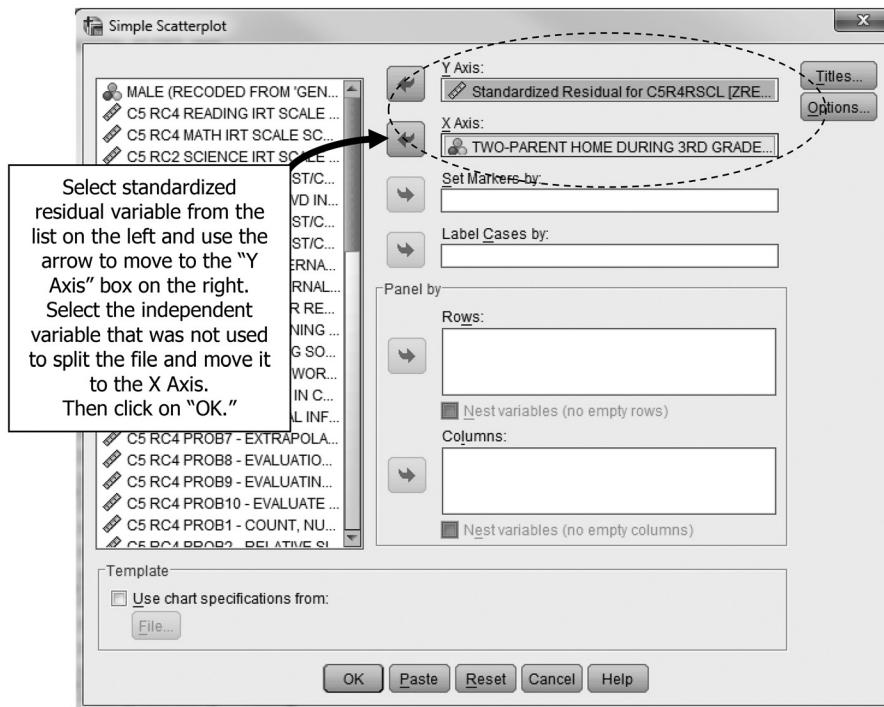
The first step is to split our file by the levels of one of our independent variables (e.g., 'MALE'). To do that, go to 'Data' in the top pulldown menu and then select 'Split File' (see screenshot Step 1). Next, select the radio button for 'Organize output by groups' and move the independent variable of interest ('MALE' in this example) into the dialog box (see screenshot Step 2). Then click on 'OK.'



3.1.1.1 Generating the Scatterplot

By this point in your statistical career, you've likely generated many scatterplots, thus the general steps for accessing a simple scatterplot in SPSS will not be reiterated here. From the "Simple Scatterplot" dialog screen, click one of the residual variables (let's

start with reading since it's the first residual in our dataset) and move it into the "Y Axis" box by clicking on the arrow. Click the independent variable that was not used to split the file (e.g., 'TWOPARENT') and move it into the "X Axis" box by clicking on the arrow. Then click "OK."

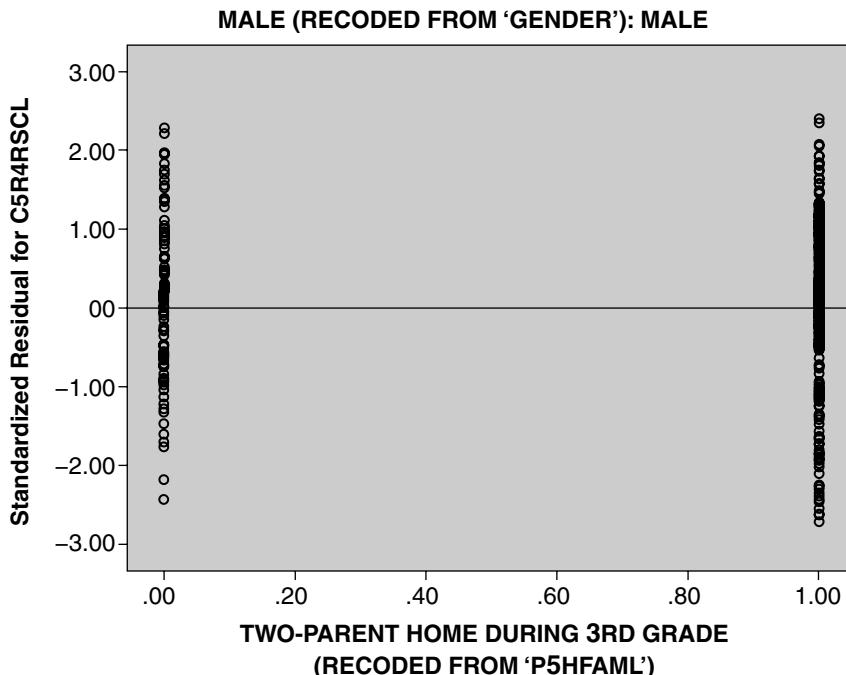


3.1.1.2 Interpreting Independence Evidence

In examining the scatterplots for evidence of independence, the points should fall relatively randomly above and below a horizontal line at zero and within the band of +2.0 to -2.0. (You can add a reference line to the graph using Chart Editor. To add a reference line, double click on the graph in the output to activate the chart editor. Select 'Options' in the top pull-down menu, then 'Y axis reference line.' This will bring up the 'Properties' dialog box. Change the value of the position to be '0.' Then click on 'Apply' and 'Close' to generate the graph with a horizontal line at zero.)

In this example, our scatterplot for each level of two-parent home generally suggests evidence of independence with a relatively random display of residuals above and below the horizontal line at zero for both males and females (i.e., each category of the independent variable that was used to split the file) and generally within a band of an absolute value of +2.0. We repeat this process for each standardized residual. Additional graphs are not presented here, but all are similar to the one provided here, suggesting a relatively random display of residuals above and below zero. Thus,

since we have not met the assumption of independence through random assignment of cases to groups, this gives us some assurance that independence is a reasonable assumption.



3.2 VARIANCE

There are a number of statistical assumptions related to variance or dispersion of the data. **Homogeneity of variance** is one that should already be familiar as it is an applicable assumption for the independent t test, ANOVA, and ANCOVA. Homogeneity of variance indicates that the variance of the dependent variable is approximately equal at all levels or categories of the independent variable. In other words, it is evident when the conditional distributions have the same constant variance for all values of X . In other words, homogeneity of variance is when the dependent variable has the same variance for all values of the independent variable(s) (i.e., uniform variance). Residual plots can be used to screen for homogeneity of variance, the consistency of the variance of the conditional distributions. If the homogeneity of variance assumption is violated, estimates of the standard errors are larger, and the conditional distributions may also be nonnormal. Solutions include variance stabilizing transformations (such as the square root or log of Y) or the application of other estimation methods.

Homoscedasticity, on the other hand, is applicable in the context of multiple linear regression and is evident when the conditional distributions have the same constant variance for all values of X . In other words, homoscedasticity is evident when there is

approximately equal variance in scores for one continuous variable as compared to the variance in scores for another continuous variable.

Homogeneity of variance-covariance matrices is the multivariate cousin of homogeneity of variance and is an assumption that is applicable to multivariate procedures such as MANOVA. In plain language, when this assumption is met in MANOVA, the variance-covariance matrix is approximately equal at all levels of the independent variable. This also means that variation or dispersion between groups on the collective dependent variables is equal (or at least not statistically different). Box's M , which simultaneously examines the KJ group variances and covariances, can be used to test this assumption in MANOVA, and nonstatistically significant Box's M indicates the assumption has been met. MANOVA is not robust to violations of this assumption, and this worsens under the following conditions: as the number of dependent variables increases and as the imbalance in the sample sizes per cell increases (particularly when the largest group size is two or more times the size of the smallest group) (Huberty & Olejnik, 2006). If and/or when this assumption is violated, use Levene's test for univariate analysis to determine the dependent variable that has the heterogeneous variance and apply a variance stabilizing transformation such as natural log or square root. Because violations of this assumption and that of multivariate normality often go hand-in-hand, a transformation to correct for nonnormality may also correct for unequal variance-covariance matrices. To replace the omnibus test, the Yao test is an option for examining specific contrasts (Huberty & Olejnik, 2006).

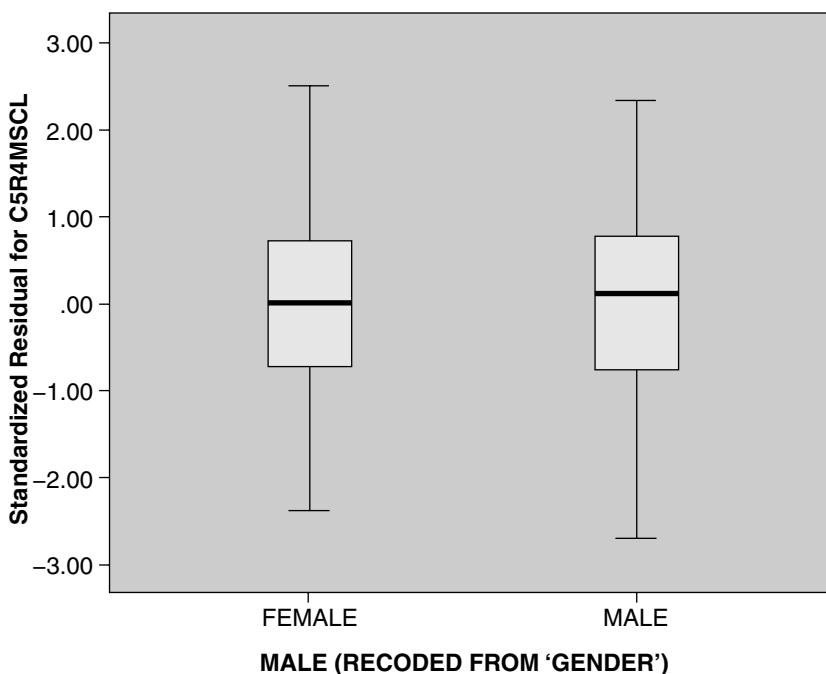
Violations of this assumption are most problematic when the observed probability values are within close range of the alpha level. This is why: The alpha level is too conservative in cases where the cells that have the larger sample sizes also have the larger variance (thus, results that are not statistically significant *may* have been had there been equal variance-covariance matrices—so if you *do* find statistical significance in this situation, there is little reason for concern for violation of the assumption). The alpha level is too liberal in cases where the cells that have the smaller sample sizes have the larger variances (thus, results that are statistically significant *may not* have been had there been equal variance-covariance matrices—this situation is extremely problematic; so if you *do not* find statistical significance in this situation, there is little reason for concern for violation of the assumption).

3.2.1 Screening for Homoscedasticity

Screening for homoscedasticity can be done using the same plots that were used to examine independence. To examine the extent to which homogeneity was met, we plot (1) residuals (we use the studentized residuals in this example) against unstandardized predicted values and (2) studentized residuals against each independent variable. Evidence of meeting the assumption of homoscedasticity is a plot where the spread of residuals appears fairly constant over the range of unstandardized predicted values (i.e., a random display of points) and observed values of the independent variables. If the display of residuals increases or decreases across the plot, then there may be an indication that the assumption of homoscedasticity has been violated.

3.2.2 Screening for Homogeneity of Variance-Covariance

In MANOVA, homogeneity of variance-covariance can be tested using Box's M test. Statistically significant results indicate a violation of the assumption. A more subjective and visual examination of the homogeneity of variance-covariance assumption can be accomplished via spread-versus-level plots and boxplots of the standardized residuals to the factors. Spread-versus-level plots were requested and reviewed in the output (see chapter 6, table 6.3) and suggest reasonable homogeneity. Boxplots of standardized residuals to the factors suggest evidence of homogeneity, as there are no substantial differences in the box lengths or whisker lengths for the predictors by group. In the graph presented here, an illustration from the coverage of this assumption in the MANOVA chapter, the standardized residual suggests very similar box lengths and whisker lengths by group, providing another form of evidence for homogeneity of variance-covariance.



3.3 NORMALITY

The assumption of normality, that the conditional distributions of the scores on the dependent variable, or the prediction errors, are **normal** in shape, is an assumption that is applicable to many multivariate procedures. In the multivariate world, normality is an assumption that is usually considered both in the context of *univariate* as well as in the context of *multivariate*. We will first touch on **univariate normality**. Violation of the normality assumption may be the result of outliers. Various recommendations are used to crudely detect outliers from a residual plot or scatterplot. The simplest outlier detection procedure and a commonly used rule is to define an outlier as an observation more than two or three standard errors from the mean (i.e., a large distance from the mean).

The outlier observation may be a result of (a) a simple recording or data entry error, (b) an error in observation, (c) an improperly functioning instrument, (d) inappropriate use of administration instructions, or (e) a true outlier. If the outlier is the result of an error, correct the error if possible and redo the regression analysis. If the error cannot be corrected, then the observation could be deleted. If the outlier represents an accurate observation, then this observation may contain important theoretical information, and one would be more hesitant to delete it (or perhaps seek out similar observations).

The following can be used to examine the extent to which **univariate normality** is present: frequency distributions, normal probability (Q-Q) plots, formal tests of normality, and skewness and kurtosis statistics. Graphical displays, such as a histogram, frequency distribution, and boxplot, can assist in checking for symmetry as can review of skewness and kurtosis statistics. **Nonzero kurtosis** (i.e., a distribution that is either flat, platykurtic, or has a sharp peak, leptokurtic) is less problematic; however **nonzero skewness** (i.e., a distribution that is not symmetric with either a positive or negative skew) usually has much more impact on parameter estimates. Thus, finding asymmetrical distributions is a must. One suggestion is to be concerned if the skewness value is larger than an absolute value of 2.0 in magnitude. Kurtosis statistics beyond an absolute value of 7.0 are considered nonnormal.

Another useful graphical technique is the normal probability plot (or Q-Q plot). With normally distributed data, the points on the normal probability plot will fall along a straight diagonal line, whereas nonnormal data will not. Because many of the tools for examining univariate normality are subjective (e.g., graphs), the application of multiple data screening tools is recommended.

Formal tests of normality can also be used to gauge normality. The Kolmogorov-Smirnov (K-S) (Chakravart, Laha, & Roy, 1967) with Lilliefors significance (Lilliefors, 1967), and the Shapiro-Wilk (S-W) (Shapiro & Wilk, 1965) are tests that provide evidence of the extent to which our sample distribution is statistically different from a normal distribution. The K-S test tends to be conservative, whereas the S-W test is usually considered the more powerful of the two for testing normality and is recommended for use with small sample sizes ($n < 50$). Nonstatistically significant results are desirable and indicate that the sample distribution is not statistically significantly different from what would be expected from a normal distribution.

In the event that univariate nonnormality is detected, transformations can be used to normalize the data. The challenge with transformed variables comes in interpretation, given that they are measured along some scale other than that of the original variables.

Multivariate normality is met when the linear combination of variables is normally distributed. *Univariate normality is a necessary condition for multivariate normality*, rejection of the assumption of univariate normality means that multivariate normality is also rejected. Univariate normality is *not* a sufficient condition for multivariate normality, however, and *meeting univariate normality should not imply tenable multivariate normality*. In terms of multivariate normality, we can employ the use of a macro in SPSS (DeCarlo, 1997) to examine a number of multivariate normality indices

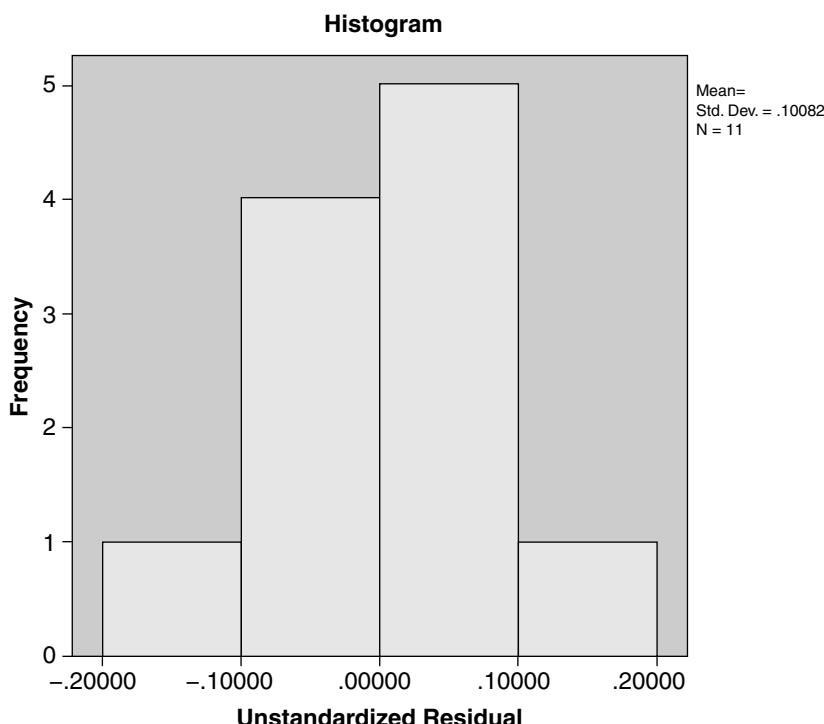
including: (a) multivariate kurtosis (Mardia, 1970); (b) multivariate skewness and kurtosis based on Small's (1980) multivariate extension of univariate skewness and kurtosis (Looney, 1995); (c) multivariate normality omnibus test (Looney, 1995); (d) largest squared and plot of squared Mahalanobis distance; and (e) critical values for hypothesis test for a single multivariate outlier using Mahalanobis distance (Penny, 1996).

3.3.1 Screening for Univariate Normality

A basic place to start in screening for univariate normality is by examining the unstandardized residuals. The data for this example comes from the illustration in the multiple regression chapter.

Descriptives		
	Statistic	Std. Error
Unstandardized Residual		
Skewness	-.336	.661
Kurtosis	.484	1.279

In this example, the skewness statistic of the residuals is $-.336$ and kurtosis is $.484$ —both being within the range of an absolute value of 2.0 and 7.0, respectively, suggesting some evidence of normality. Given the very small sample size, the histogram below reflects as normal a distribution as might be expected.



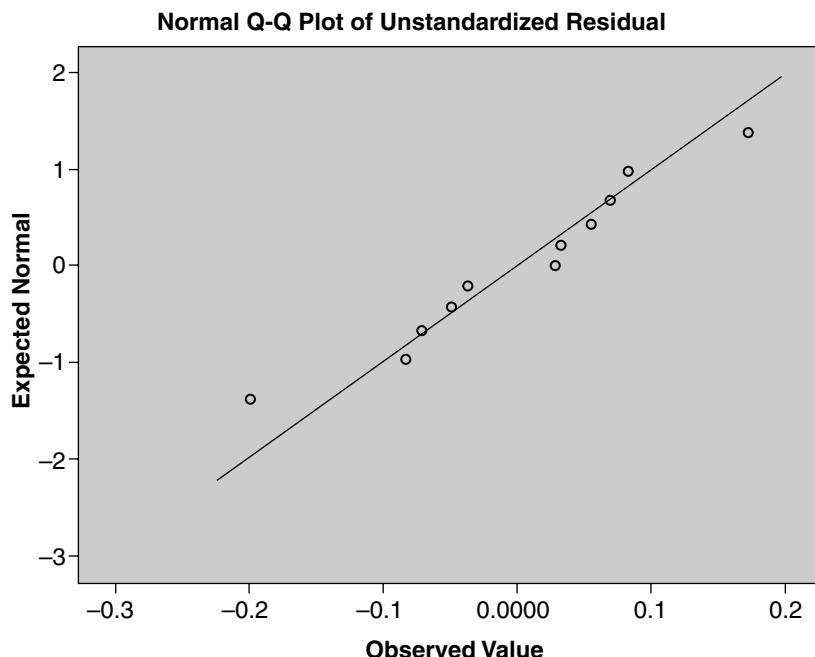
There are a few other statistics that can be used to gauge normality. The formal test of normality, the Shapiro-Wilk test (*SW*) (Shapiro & Wilk, 1965), provides evidence of the extent to which our sample distribution is statistically different from a normal distribution. The output for the Shapiro-Wilk test is presented below and suggests that our sample distribution for the residual is *not* statistically significantly different than what would be expected from a normal distribution, as the *p* value is greater than α ($p = .918$).

	Tests of Normality			Shapiro-Wilk			
	Kolmogorov-Smirnov ^a	Statistic	df	Sig.	Statistic	df	Sig.
Unstandardized Residual		.155	11	.200*	.973	11	.918

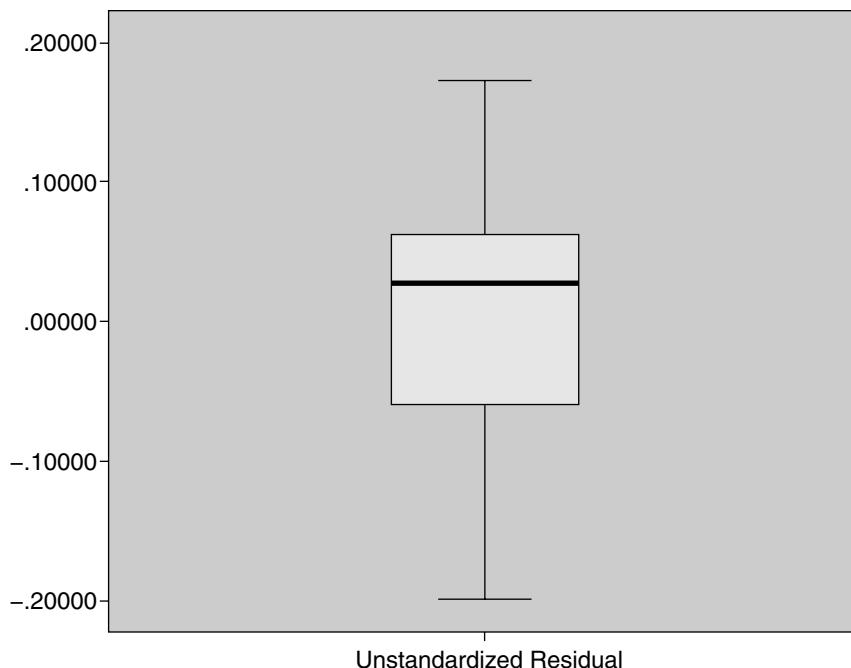
a. Lilliefors Significance Correction

*. This is a lower bound of the true significance.

Quantile-quantile (Q-Q) plots are also often examined to determine evidence of normality. Q-Q plots graph quantiles of the theoretical normal distribution against quantiles of the sample distribution. Points that fall on or close to the diagonal line suggest evidence of normality. The Q-Q plot of residuals suggests relative normality.



Examination of the boxplot below also suggests a relatively normal distribution of residuals with no outliers.

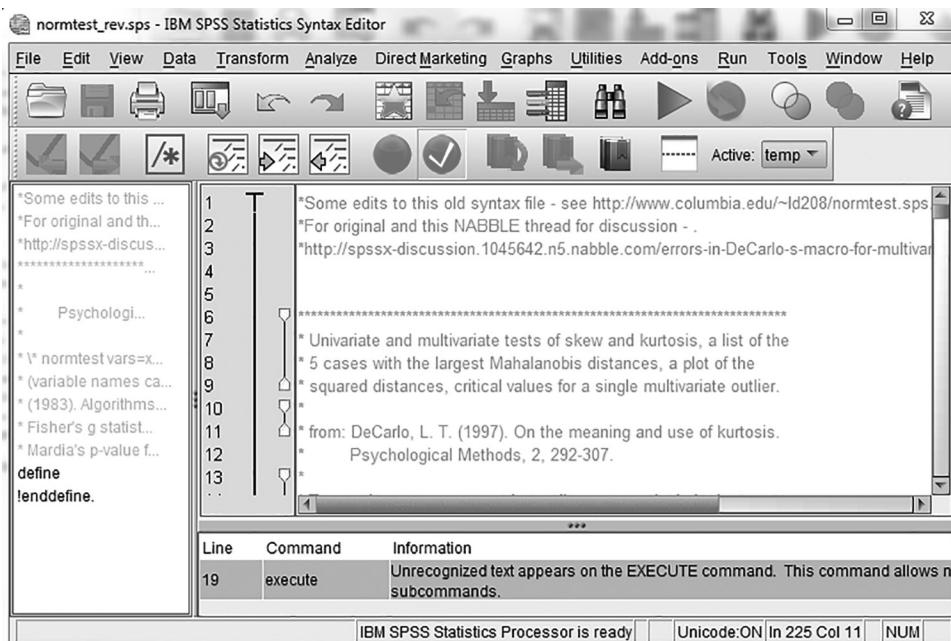


Considering the forms of evidence we have examined, skewness and kurtosis statistics, the Shapiro-Wilk test, histogram, the Q-Q plot, and the boxplot, all suggest univariate normality is a reasonable assumption.

3.3.2 Screening for Multivariate Normality

Although univariate normality is a necessary condition of multivariate normality, it is not a sufficient condition. SPSS does not have built-in capability to examine *multivariate normality*. However, using the data example from MANOVA, we will enlist the use of a macro that will do just that. The macro is based on the work of DeCarlo (1997) and is extremely simple to use.

First, copy the macro language (the original language is accessible from <http://www.columbia.edu/~ld208/normtest.sps> and printed in the 1997 article as well; there is a bug in that macro that is fixed in the language accessible from this site: https://dl.dropbox-usercontent.com/u/3385251/Mult_Macro_fixAndy.sps) and paste it into a SPSS syntax file. To open a new syntax file, in SPSS go to File in the top pull-down menu, then New then Syntax. Your syntax file with the macro language should look something like this:



Save this syntax file by going to **File** then **Save as** in the top pull-down menu in SPSS.

Second, open the SPSS data file that includes the variables you want to examine. In this case, it's the ECLS-K data from the MANOVA chapter. The dataset you will be working with needs to be open when running the syntax.

Third, open another new syntax file. To open a new syntax file, in SPSS go to **File** in the top pull-down menu, then **New**, then **Syntax**. Type the following two-line command:

```
INCLUDE file='C:\normtest.sps'.
normtest vars=x1,x2,x3,x4 /.
```

The first line of this syntax (i.e., `INCLUDE file='C:\normtest.sps'`) tells SPSS where to find the macro file (the file that was saved to syntax in step 1). Thus, in the first line of this new syntax file, change 'normtest' to the name of *your* syntax file that holds the macro. The second line (i.e., '`normtest vars=x1,x2,x3,x4 /.`') specifies the variables to examine. Thus, change the 'x1' to the name of the first variable in *your dataset* to examine, 'x2' to the name of the second variable in your dataset, and so forth. In this illustration, we are reviewing the standardized residuals (ZRE_1, ZRE_2, ZRE_3) so the second line of our syntax specifies: `normtest vars= ZRE_1, ZRE_2, ZRE_3 /.`

Fourth, save this two-line syntax file by going to **File** then **Save as** in the top pull-down menu in SPSS. Then click on **Run** in the top horizontal menu in SPSS to generate the output (which follows). In summary, you will have two syntax (i.e., sps) files: one that contains the macro and a second that contains only two lines of syntax that is used to

run the macro on the dataset that is opened in SPSS. [Note: If you run into problems in running the syntax, such as getting errors in generating the graphs, run SPSS as an administrator (if you are able) and that should solve the problem.]

Figure 3.1 shows the **output** from the multivariate normality macro. Generally, there is evidence of a violation of multivariate normality based on the multivariate skewness and kurtosis tests, as well as the omnibus multivariate normality test.

Output from DeCarlo (1997) SPSS Macro for Multivariate Normality

Run MATRIX procedure:			
Number of observations:	1393		
Number of variables:	3		
Measures and tests of skew:			
ZRE_1 ZRE_2 ZRE_3	g1 -.2301 -.2092 .0370	sqrt(b1) -.2298 -.2090 .0370	z (b1) -3.4764 -3.1679 .5659
			p-value .0005 .0015 .5715
Measures and tests of kurtosis:			
ZRE_1 ZRE_2 ZRE_3	g2 -.4173 -.5561 -.4846	b2-3 -.4201 -.5585 -.4871	z (b2) -4.0089 -5.9430 -4.8948
			p-value .0001 .0000 .0000
Omnibus tests of normality (both chisq, 2 df):			
D'Agostino & Pearson K sq ZRE_1 ZRE_2 ZRE_3	K sq 28.1562 45.3556 24.2788	p-value .0000 .0000 .0000	Jarque & Bera LM test LM 22.5059 28.2421 14.0897
			p-value .0000 .0000 .0009
***** Multivariate Statistics *****			
< Tests of multivariate skew: >			
Small's test (chisq)			
chi (b1p) 22.1106	df 3.0000	p-value .0001	
Srivastava's test			
chi (b1p) 10.3201	df 3.0000	p-value .0160	
< Tests of multivariate kurtosis: >			
A variant of Small's test (chisq)			
VQ2 53.4114	df 3.0000	p-value .0000	Statistically significant multivariate skew and kurtosis indicate violation of multivariate normality. In this case, there is a violation of multivariate skew for both Small's and Srivastava's tests and a violation of multivariate kurtosis for the variant of Small's and Mardia's test. The only evidence of multivariate normality is based on Srivastava's multivariate kurtosis.
Srivastava's test			
b2p 2.8906	N(b2p) -1.4436	p-value .1488	

FIGURE 3.1

Output from DeCarlo (1997) SPSS Macro for Multivariate Normality

```
Mardia's test
      b2p      N(b2p)    p-value
    13.8359     -3.9662     .0001
```

Omnibus test of multivariate normality:

```
(based on Small's test, chisq)
      VQ3      df    p-value
    75.5220     6.0000     .0000
```

The omnibus test of *multivariate* normality also indicates a violation.

----- END MATRIX -----

Critical values (Bonferroni) for a single multivar. outlier:

```
critical F(.05/n) =23.08  df = 3,1389
critical F(.01/n) =26.36  df = 3,1389
```

5 observations with largest Mahalanobis distances:

rank = 1	case# = 1063	Mahal D sq =	18.05
rank = 2	case# = 1153	Mahal D sq =	17.49
rank = 3	case# = 231	Mahal D sq =	14.99
rank = 4	case# = 221	Mahal D sq =	13.07
rank = 5	case# = 1165	Mahal D sq =	12.36

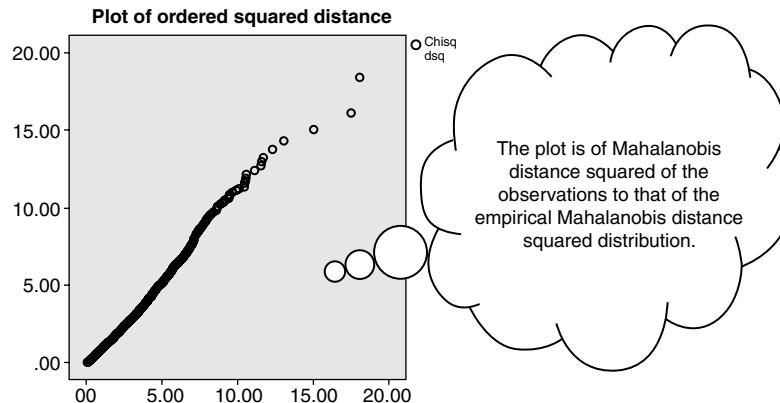


FIGURE 3.1

Continued

In addition to this output, a new SPSS data file is created that includes the requisite variables generated from the multivariate normality macro in the event additional exploration is desired.

	Name	Type	Width	Decimals	Label	Values	Missing
1	case	Numeric	8	2		None	None
2	rnk	Numeric	8	2		None	None
3	top	Numeric	8	2		None	None
4	dsq	Numeric	8	2		None	None
5	pvar	Numeric	8	2		None	None
6	ddf	Numeric	8	2		None	None
7	ncase	Numeric	8	2		None	None
8	a01	Numeric	8	2		None	None
9	a05	Numeric	8	2		None	None
10	f01	Numeric	8	2		None	None
11	f05	Numeric	8	2		None	None
12	fc01	Numeric	8	2		None	None
13	fc05	Numeric	8	2		None	None
14	chisq	Numeric	8	2		None	None

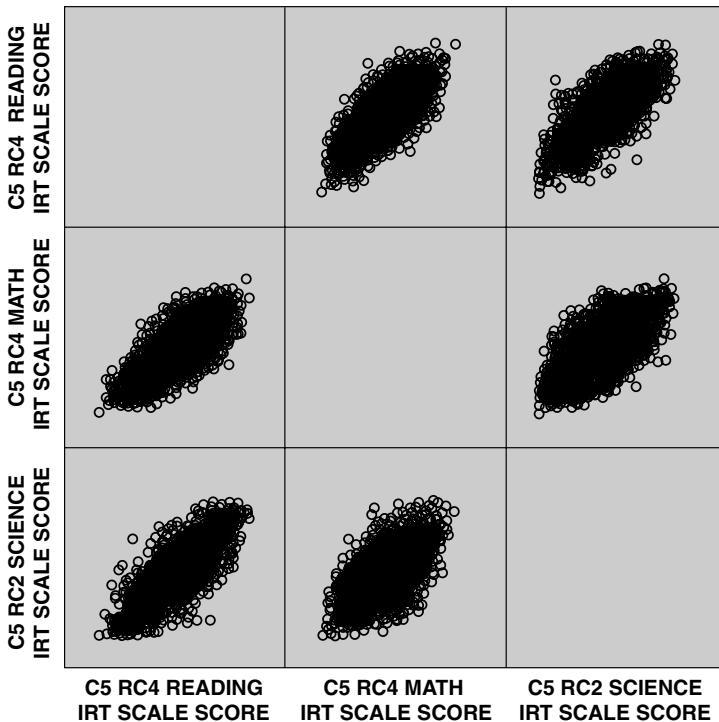
3.4 LINEARITY

Linearity is an assumption that is common to most of the multivariate procedures that we discuss, including multiple and logistic regression, MANOVA, repeated measures MANOVA, discriminant analysis, and exploratory and confirmatory factor analysis. Linearity, very simply, means that there is a linear relationship. In a very general context, linearity in a model refers to there being linearity in the parameters of the model. In the context of regression-based procedures, the linear relationship is between the observed scores on the dependent variable and the values of the independent variables. In the context of MANOVA, when all pairs of dependent variables are bivariate normally distributed (i.e., linear), the assumption of linearity is met.

3.4.1 Screening for Linearity

Screening for linearity differs depending on the analytic approach. For regression, violation of the linearity assumption can be detected through residual plots—the same

plots used to examine independence and homogeneity. The residuals should be located within a band of $+2 s_{\text{res}}$ (or standard errors), indicating no systematic pattern of points. In MANOVA, for example, linearity of the dependent variables can be examined by scatterplots of all pairs of dependent variables. The linearity assumption means that a straight line provides a reasonable fit to the data (as seen in this figure, illustrated later in the MANOVA chapter). Procedures to deal with nonlinearity include transformations and other analytic approaches that do not require linearity (e.g., logistic regression).



3.5 NONCOLLINEARITY

Noncollinearity is an assumption unique to multivariate procedures. Collinearity (also known as multicollinearity) occurs when there is a very strong linear relationship between two or more of the predictors, in the case of regression-based procedures or two or more variables included in factor analytic models (e.g., exploratory and confirmatory factor analysis). The presence of severe collinearity is problematic in several respects, and the issues it can create vary depending on the procedure. As an example, in multiple regression, collinearity can lead to instability of the estimated regression coefficients across samples. In cluster analysis, as another example, collinearity reduces the ability to identify distinct profiles. Singularity is a special case of multicollinearity; it is perfect multicollinearity and occurs when two or more items/variables perfectly

predict and are therefore perfectly redundant. This can occur in factor analysis as well as multiple regression. Trust it to say that collinearity is important to avoid.

In multiple regression, the simplest way to detect collinearity is to conduct a series of special regression analyses, one for each independent variable, where that predictor is predicted by all of the remaining independent variables (i.e., the dependent variable is not involved). If any of the resultant R_k^2 values are close to one (greater than .9 is a good rule), then there may be a collinearity problem. Note, however, that large R^2 values may also be due to small sample size. Collinearity can occur if the number of predictors is greater than or equal to n . Another statistical method for detecting collinearity in regression is to compute a variance inflation factor (VIF) for each predictor, which is equal to $1 / (1 - R_k^2)$. The VIF is defined as the inflation that occurs for each regression coefficient above the ideal situation of uncorrelated predictors. Many suggest that the largest VIF should be less than 10 in order to satisfy this assumption (Myers, 1990; Stevens, 2009; Wetherill, 1986).

There are several possible methods for dealing with a collinearity problem. First, and the easiest, is that one can remove one or more of the correlated predictors. Second, ridge regression techniques can be used (e.g., Hoerl & Kennard, 1970a, 1970b; Marquardt & Snee, 1975; Myers, 1986; Wetherill, 1986). Third, principal component scores resulting from principal component analysis can be utilized rather than raw scores on each variable (e.g., Kleinbaum, Kupper, Muller, & Nizam, 1998; Myers, 1986; Weisberg, 1985; Wetherill, 1986). Fourth, transformations of the variables can be used to remove or reduce the extent of the problem. The final solution, and probably our last choice, is to use a univariate procedure, as collinearity cannot exist in models that include only a single predictor.

3.5.1 Screening for Noncollinearity

Detecting multicollinearity can be done by reviewing the VIF and tolerance statistics in multiple regression. Tolerance is calculated as $(1 - R^2)$ and values close to zero (a recommendation is .10 or less) suggest potential multicollinearity problems. A tolerance of .10 suggests that 90% (or more) of the variance in one of the independent variables can be explained by another independent variable. VIF is the ‘variance inflation factor’ and is the reciprocal of tolerance, where $VIF = \frac{1}{tolerance}$. VIF values greater than 10 (which correspond to a tolerance of .10) suggest potential multicollinearity.

Collinearity diagnostics (see SPSS output below) can also be reviewed. Multiple eigenvalues close to zero indicate independent variables that have strong intercorrelations. The condition index is calculated as the square root of the ratio of the largest eigenvalue to each preceding eigenvalue. A general recommendation for interpreting condition indices is that values in the range of 10 to 30 should be of concern, greater than 30 indicates trouble, and greater than 100 indicates disaster (Belsley, 1991). In

this case, drawn from our multiple regression chapter, both the eigenvalues and condition indices suggest possible problems with multicollinearity.

Collinearity Diagnostics^a

Model	Dimension	Eigenvalue	Condition Index	Variance Proportions		
				(Constant)	GRE Total Score	Undergraduate grade point average
1	1	2.981	1.000	.00	.00	.00
	2	.012	15.727	.03	.86	.40
	3	.007	20.537	.97	.13	.60

a. Dependent Variable: Graduate grade point average

Multicollinearity can also be examined by computing regression models where each independent variable is considered the outcome and is predicted by the remaining independent variables (the dependent variable is not included in these models). If any of the resultant R_k^2 values are close to one (greater than .9 is a good criterion), then there may be a collinearity problem. For the example data in multiple regression, $R_{12}^2 = .091$ and therefore collinearity is not a concern.

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate
1	.301	.091	-.010	16.41926

a. Dependent Variable: Graduate grade point average

PROBLEMS

Conceptual Problems

- More so than other assumptions, the assumption of normality is inherently tied to the sampling design of your study. True or false?
- Which one of the following would suggest the assumption of independence has been met?
 - Nonstatistically significant Box's M
 - Q-Q plot shows points adhering closely to the diagonal line
 - Scatterplot of residuals to predicted values shows random display of points
 - Kurtosis statistic that is within an absolute value of 7
- Which one of the following would suggest the assumption of normality has been met?
 - Nonstatistically significant Levene's test
 - Nonstatistically significant Box's M
 - Scatterplot of residuals to predicted values shows random display of points
 - Kurtosis statistic that is within an absolute value of 7

4. Which one of the following would suggest the assumption of homogeneity of variance-covariance has been met?
 - a. Nonstatistically significant Levene's test
 - b. Nonstatistically significant Box's M
 - c. Scatterplot of residuals to predicted values shows random display of points
 - d. Kurtosis statistic that is within an absolute value of 7
5. Which one of the following would suggest the assumption of homogeneity of variance has been met?
 - a. Nonstatistically significant Levene's test
 - b. Nonstatistically significant Box's M
 - c. Scatterplot of residuals to predicted values shows random display of points
 - d. Kurtosis statistic that is within an absolute value of 7
6. Which one of the following would be considered a relatively normal distribution?
 - a. Skewness = .10; kurtosis = 50
 - b. Skewness = 1.0; kurtosis = 5.0
 - c. Skewness = 5.0; kurtosis = 1.0
 - d. Skewness = 10.0; kurtosis = 0.5
7. The multivariate version of equal variances that is applicable to multiple linear regression is which one of the following?
 - a. Homogeneity of variance
 - b. Homogeneity of variance-covariance
 - c. Homoscedasticity
 - d. Sphericity
8. The assumption of linearity is met in MANOVA under which one of the following conditions?
 - a. When all pairs of dependent variables are bivariate normally distributed
 - b. When points fall on or closely to the diagonal line in a Q-Q plot
 - c. When Shapiro-Wilk test is not statistically significant
 - d. When the dependent variable is linearly related to the independent variable
9. Multicollinearity is suggested by which one of the following?
 - a. Statistically significant Box's M
 - b. Statistically significant multivariate skew
 - c. Tolerance of .50
 - d. VIF value of 30
10. Generally, in designs that employ simple random sampling, which one of the following assumptions have been met?
 - a. Homogeneity
 - b. Independence
 - c. Linearity
 - d. Normality

REFERENCES

- Barcikowski, R. S. (1981). Statistical power with group mean as the unit of analysis. *Journal of Educational Statistics*, 6(3), 267–285.
- Belsley, D. A. (1991). A guide to using the collinearity diagnostics. *Computer Science in Economics and Management*, 4, 33–50.
- Chakravart, I. M., Laha, R. G., & Roy, J. (1967). *Handbook of methods of applied statistics* (Vol. 1). New York: Wiley.
- DeCarlo, L. T. (1997). On the meaning and use of kurtosis. *Psychological Methods*, 2, 292–307.
- Hoerl, A. E., & Kennard, R. W. (1970a). Ridge regression: Application to non-orthogonal models. *Technometrics*, 12, 591–612.
- Hoerl, A. E., & Kennard, R. W. (1970b). Ridge regression: Biased estimation for non-orthogonal models. *Technometrics*, 12, 55–67.
- Huberty, C. J., & Olejnik, S. (2006). *Applied MANOVA and discriminant analysis* (2nd ed.). Hoboken, NJ: Wiley.
- Kleinbaum, D. G., Kupper, L. L., Muller, K. E., & Nizam, A. (1998). *Applied regression analysis and other multivariable models* (3rd ed.). Pacific Grove, CA: Duxbury.
- Lilliefors, H. (1967). On the Kolmogorov-Smirnov test for normality with mean and variance unknown. *Journal of the American Statistical Association*, 62, 399–402.
- Looney, S. W. (1995). How to use tests for univariate normality to assess multivariate normality. *American Statistician*, 49, 64–70.
- Mardia, K. V. (1970). Measures of multivariate skewness and kurtosis with applications. *Biometrika*, 57, 519–530.
- Marquardt, D. W., & Snee, R. D. (1975). Ridge regression in practice. *The American Statistician*, 29, 3–19.
- Myers, R. H. (1986). *Classical and modern regression with applications*. Boston, MA: Duxbury.
- Myers, R. H. (1990). *Classical and modern regression with applications* (2nd ed.). Boston, MA: Duxbury.
- Penny, K. I. (1996). Appropriate critical values when testing for a single multivariate outlier by using the Mahalanobis distance. *Applied Statistics*, 45, 73–81.
- Scariano, S. M., & Davenport, J. M. (1987). The effects of violation of independence assumptions in the one-way ANOVA. *The American Statistician*, 41(2), 123–129.
- Shapiro, S. S., & Wilk, M. B. (1965). An analysis of variance test for normality (complete samples). *Biometrika*, 52(3 and 4), 591–611.
- Small, N. J. H. (1980). Marginal skewness and kurtosis in testing multivariate normality. *Applied Statistics*, 29, 85–87.
- Stevens, J. P. (2009). *Applied multivariate statistics for the social sciences* (5th ed.). New York: Psychology Press.
- Weisberg, S. (1985). *Applied linear regression* (2nd ed.). New York: Wiley.
- Wetherill, G. B. (1986). *Regression analysis with applications*. London: Chapman & Hall.

Chapter 4

MULTIPLE LINEAR REGRESSION

CHAPTER OUTLINE

4.1 What Multiple Linear Regression Is and How It Works	58
4.1.1 Characteristics	59
4.1.2 Sample Size	76
4.1.3 Power	76
4.1.4 Effect Size	76
4.1.5 Assumptions	77
4.2 Mathematical Introduction Snapshot	83
4.3 Computing Multiple Linear Regression Using SPSS	87
4.4 Data Screening	96
4.4.1 Independence	97
4.4.2 Homoscedasticity	97
4.4.3 Linearity	98
4.4.4 Normality	98
4.4.5 Noncollinearity	102
4.5 Power Using G*Power	104
4.5.1 Post Hoc Power for Multiple Linear Regression Using G*Power	104
4.5.2 A Priori Power for Multiple Linear Regression Using G*Power	107
4.6 Research Question Template and Example Write-Up	107

KEY CONCEPTS

1. Partial and semipartial (part) correlations
2. Standardized and unstandardized regression coefficients
3. Coefficient of multiple determination and multiple correlation

Modeling prediction is one of the most common methods of quantitative analysis. This leads us to multiple regression analysis, where we are able to model two or more predictors to predict or explain the criterion variable. Here we adopt the usual notation where the X 's are defined as the independent or predictor variables, and Y as the dependent or criterion variable.

For example, an admissions officer might use Graduate Record Exam (GRE) scores to predict graduate-level grade point averages (GPA) to make admissions decisions for a sample of applicants to your favorite local university or college. The admissions office may decide that including only one variable omits a number of other factors that relate to GPA. Other potentially useful predictors might be undergraduate grade point average, ratings of recommendation letters, scored writing samples, and/or an evaluation from a personal interview. The research question of interest would now be, How well do the GRE, undergraduate GPA, recommendation ratings, writing sample scores, and/or interview scores (the independent or predictor variables) predict performance in graduate school (the dependent or criterion variable)? This is an example of a situation where multiple regression analysis using multiple predictor variables might be the method of choice.

This chapter considers the concepts of partial, semipartial, and multiple correlations, standardized and unstandardized regression coefficients, and the coefficient of multiple determination, as well as introduces a number of other types of regression models. Our objectives are that, by the end of this chapter, you will be able to (a) determine and interpret the results of partial and semipartial correlations, (b) understand the concepts underlying multiple linear regression, (c) determine and interpret the results of multiple linear regression, (d) understand and evaluate the assumptions of multiple linear regression, and (e) have a basic understanding of other types of regression models.

4.1 WHAT MULTIPLE LINEAR REGRESSION IS AND HOW IT WORKS

Throughout the book, we will be following a group of superbly talented, creative, and energetic graduate research assistants (Challie Lenge, Ott Lier, Addie Venture, and Oso Wyse) working in their institution's statistics and research lab. The group is supervised and mentored by a research methodology faculty member who empowers the group to lead their projects to infinity and beyond, so to speak. With each chapter, we will find the group, or a subset of members thereof, delving into a fantastical multivariate journey. As we first become acquainted with the group, we find Addie Venture,

who has completed a number of quantitative statistics courses and has developed into quite a statistics guru, being sought after from across her university campus.

Dr. Golly, the assistant dean in the Graduate Student Services office, seeks advice from Addie on a special project. Dr. Golly is interested in estimating the extent to which graduate grade point average can be predicted by scores on the overall Graduate Record Exam (GRE-total) and undergraduate grade point average. From her recent statistical trek in regression, Addie knows that questions delving into relationships and prediction with continuous outcomes and multiple predictors can be examined using multiple regression. Addie suggests the following research question to Dr. Golly: *Can graduate grade point average be predicted by scores on the overall Graduate Record Exam (GRE-total) and undergraduate grade point average?* Addie determines that a multiple linear regression is the appropriate statistical procedure to use to answer Dr. Golly's question. Excited for the first project of the semester, Addie then proceeds to assist Dr. Golly in analyzing the data and interpreting the results.

4.1.1 Characteristics

Prior to a discussion of regression analysis, we need to consider two related concepts in correlational analysis, partial and semipartial correlations. Multiple regression analysis involves the use of two or more predictor variables and one criterion variable; thus, at a minimum three variables are involved in the analysis. If we think about these variables in the context of the Pearson correlation, we have a problem because this correlation can only be used to relate two variables at a time. How do we incorporate additional variables into a correlational analysis? The answer is through partial and semipartial correlations, and later in this chapter, multiple correlations.

4.1.1.1 Partial Correlation

First we discuss the concept of **partial correlation**. The simplest situation consists of three variables, which we label X_1 , X_2 , and X_3 . Here, an example of a partial correlation would be the correlation between X_1 and X_2 where X_3 is held constant (i.e., controlled or partialed out). That is, the influence of X_3 is removed from both X_1 and X_2 (both have been adjusted for X_3). Thus the partial correlation here represents the linear relationship between X_1 and X_2 independent of the linear influence of X_3 . This particular partial correlation is denoted by $r_{12.3}$, where the X 's are not shown for simplicity and the dot indicates that the variables preceding it are to be correlated and the variable(s) following it are to be partialed out. We compute $r_{12.3}$ as follows:

$$r_{12.3} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}}$$

Let us take an example of a situation where a partial correlation might be computed. Say a researcher is interested in the relationship between height (X_1) and weight (X_2).

The sample consists of individuals ranging in age (X_3) from 6 months to 65 years. The sample correlations are for height (X_1) and weight (X_2), $r_{12} = .7$; height (X_1) and age (X_3), $r_{13} = .1$; and weight (X_2) and age (X_3), $r_{23} = .6$. We compute the correlation between height and weight, controlling for age, $r_{12,3}$, as follows:

$$r_{12,3} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} = \frac{.7 - (.1)(.6)}{\sqrt{(1 - .01)(1 - .36)}} = .8040$$

We see here that the bivariate correlation between height and weight, ignoring age ($r_{12} = .7$), is smaller than the partial correlation between height and weight controlling for age ($r_{12,3} = .8040$). That is, the relationship between height and weight is stronger when age is held constant (i.e., for a particular age) than it is across all ages. Although we often talk about holding a particular variable constant, in reality variables such as age cannot be held constant artificially.

Holding age constant would be an *experimental control*—controlling for the effects of age by collecting height and weight data from everyone who has the same age. It is important to note that this is not the same as achieving *statistical control*—controlling for the effects of age by correlating the residuals of a regression to predict height from age with the residuals from a regression to predict weight from age.

Some rather interesting partial correlation results can occur in particular situations. At one extreme, if both the correlation between height (X_1) and age (X_3), r_{13} , and weight (X_2) and age (X_3), r_{23} , equal zero, then the correlation between height (X_1) and weight (X_2) will equal the partial correlation between height and weight controlling for age, $r_{12} = r_{12,3}$. That is, if the variable being partialled out is uncorrelated with each of the other two variables, then the partialing process will logically not have any effect. At the other extreme, if either r_{13} or r_{23} equals 1, then $r_{12,3}$ cannot be calculated, as the denominator is equal to zero (in other words, at least one of the terms in the denominator is equal to zero, which results in the product of the two terms in the denominator equaling zero and thus a denominator of zero—and you cannot divide by zero). Thus, in this situation (where either r_{13} or r_{23} is perfectly correlated at 1.0), the partial correlation (i.e., $r_{12,3}$, partial correlation between height and weight controlling for age) is not defined. Later in this chapter, we refer to this as perfect collinearity, which is a serious problem. In between these extremes, it is possible for the partial correlation to be greater than or less than its corresponding bivariate correlation (including a change in sign), and even for the partial correlation to be equal to zero when its bivariate correlation is not. For significance tests of partial and semipartial correlations, we refer you to your favorite statistical software.

4.1.1.2 Semipartial (Part) Correlation

Next, the concept of **semipartial correlation** (also called a **part correlation**) is discussed. The simplest situation consists again of three variables, which we label X_1 , X_2 , and X_3 . Here, an example of a semipartial correlation would be the correlation between X_1 and X_2 where X_3 is removed from X_2 only. That is, the influence of X_3 is removed

from X_2 only. Thus, the semipartial correlation here represents the linear relationship between X_1 and X_2 after that portion of X_2 that can be linearly predicted from X_3 has been removed from X_2 . This particular semipartial correlation is denoted by $r_{1(2,3)}$, where the X 's are not shown for simplicity and within the parentheses the dot indicates that the variable(s) following it are to be removed from the variable preceding it. Another use of the semipartial correlation is when we want to examine the predictive power in the prediction of Y from X_1 after removing X_2 from the prediction. A method for computing $r_{1(2,3)}$ is as follows:

$$r_{1(2,3)} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{23}^2)}}$$

Let us take an example of a situation where a semipartial correlation might be computed. Say a researcher is interested in the relationship between GPA (X_1) and GRE scores (X_2). The researcher would like to remove the influence of intelligence (IQ: X_3) from GRE scores, but not from GPA. The simple bivariate correlation between GPA and GRE is $r_{12} = .5$; between GPA and IQ is $r_{13} = .3$; and between GRE and IQ is $r_{23} = .7$. We compute the semipartial correlation that removes the influence of intelligence (IQ: X_3) from GRE scores (X_2), but not from GPA (X_1) (i.e., $r_{1(2,3)}$) as follows:

$$r_{1(2,3)} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{23}^2)}} = \frac{.5 - (.3)(.7)}{\sqrt{1 - .49}} = .4061$$

Thus the bivariate correlation between GPA (X_1) and GRE scores (X_2) ignoring IQ (X_3) ($r_{12} = .50$) is larger than the semipartial correlation between GPA and GRE controlling for IQ in GRE ($r_{1(2,3)} = .4061$). As was the case with partial correlations, various values of a semipartial correlation can be obtained depending on the combination of the bivariate correlations. For more information on partial and semipartial correlations, see Hays (1988), Glass and Hopkins (1996), or Pedhazur (1997).

Now that we have considered the correlational relationships among two or more variables (i.e., partial and semipartial correlations), let us move on to an examination of the multiple regression model where there are two or more predictor variables.

Let us take the concepts we have learned in this chapter and place them into the context of multiple linear regression. For purposes of brevity, we do not consider the population situation because the sample situation is invoked 99.44% of the time. In this section, we discuss the unstandardized and standardized multiple regression models, the coefficient of multiple determination, multiple correlation, tests of significance, and statistical assumptions.

4.1.1.3 Unstandardized Regression Model

The sample multiple linear regression model for predicting Y from m predictors $X_{1,2,\dots,m}$ is

$$Y_i = b_1 X_{1i} + b_2 X_{2i} + \dots + b_m X_{mi} + a + e_i$$

where Y is the criterion variable (also known as the dependent variable); the X_k 's are the predictor (or independent) variables where $k = 1, \dots, m$; b_k is the sample partial slope of the regression line for Y as predicted by X_k ; a is the sample intercept of the regression line for Y as predicted by the set of X_k 's; e_i are the residuals or errors of prediction (the part of Y not predictable from the X_k 's); and i represents an index for an individual or object. The index i can take on values from 1 to n where n is the size of the sample (i.e., $i = 1, \dots, n$). The term **partial slope** is used because it represents the slope of Y for a particular X_k in which we have partialled out the influence of the other X_k 's, much as we did with the partial correlation.

The sample prediction model is

$$Y_i = b_1 X_{1i} + b_2 X_{2i} + \dots + b_m X_{mi} + a$$

where Y_i is the predicted value of Y for specific values of the X_k 's, and the other terms are as before. There is only one difference between the regression and prediction models. The regression model explicitly includes prediction error as e_i whereas the prediction model includes prediction error implicitly as part of the predicted score Y_i (i.e., there is some error in the predicted values). The goal of the prediction model is to include an independent variable X that minimizes the residual; this means that the independent variable does a nice job of predicting the outcome. We can compute residuals, the e_i , for each of the i individuals or objects by comparing the actual Y values with the predicted Y values as

$$e_i = Y_i - Y'_i$$

for all $i = 1, \dots, n$ individuals or objects in the sample.

Determining the sample partial slopes and the intercept in the multiple predictor case is rather complicated. To keep it simple, we use a two-predictor model for illustrative purposes. Generally, we rely on statistical software for implementing multiple regression analysis. For the two-predictor case, the sample partial slopes (b_1 and b_2) and the intercept (a) can be determined as follows:

$$b_1 = \frac{(r_{Y1} - r_{Y2} r_{12}) s_Y}{(1 - r_{12}^2) s_1}$$

$$b_2 = \frac{(r_{Y2} - r_{Y1} r_{12}) s_Y}{(1 - r_{12}^2) s_2}$$

$$a = \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2$$

The sample partial slope b_1 is referred to alternately as (a) the expected or predicted change in Y for a one unit change in X_1 with X_2 held constant (or for individuals with

the same score on X_2), and (b) the unstandardized or raw regression coefficient for X_1 . Similar statements may be made for b_2 . Note the similarity of the partial slope equation to the semipartial correlation. The sample intercept is referred to as the value of the dependent variable Y when the values of the independent variables X_1 and X_2 are both zero.

An alternative method for computing the sample partial slopes that involves the use of a partial correlation is as follows:

$$b_1 = r_{Y1.2} \frac{s_Y \sqrt{1 - r_{Y2}^2}}{s_1 \sqrt{1 - r_{12}^2}}$$

$$b_2 = r_{Y2.1} \frac{s_Y \sqrt{1 - r_{Y1}^2}}{s_2 \sqrt{1 - r_{12}^2}}$$

What statistical criterion is used to arrive at the particular values for the partial slopes and intercept? The criterion usually used in multiple linear regression analysis (and in all general linear models [GLM] for that matter) is the least squares criterion. The least squares criterion arrives at those values for the partial slopes and intercept such that the sum of the squared prediction errors or residuals is smallest. That is, we want to find that regression model, defined by a particular set of partial slopes and an intercept, which has the smallest sum of the squared residuals. We often refer to this particular method for calculating the slope and intercept as least squares estimation, because a and the b_k 's represent sample estimates of the population parameters α and the β_k 's, which are obtained using the least squares criterion. Recall from simple linear regression that the residual is simply the vertical distance from the observed value of Y to the predicted value of Y , and the line of best fit minimizes this distance. This concept still applies to multiple linear regression with the exception that we are now in a three-dimensional (or more) plane, given there are multiple independent variables.

4.1.1.4 Standardized Regression Model

Up until this point in the chapter, everything in multiple linear regression analysis has involved the use of raw scores. For this reason, we referred to the model as the unstandardized regression model. Often we may want to express the regression in terms of standard z score units rather than in raw score units. The means and variances of the standardized variables (e.g., z_1 , z_2 , z_Y) are 0 and 1, respectively. The sample standardized linear prediction model becomes the following:

$$z(Y'_i) = b_1^* z_{1i} + b_2^* z_{2i} + \dots + b_m^* z_{mi}$$

where b_k^* represents a sample standardized partial slope (sometimes called beta weights) and the other terms are as before. As was the case in simple linear regression, no intercept term is necessary in the standardized prediction model, as the mean of the z scores for all variables is 0. (Recall that the intercept is the value of

the dependent variable when the scores on the independent variables are all zero. Thus, in a standardized prediction model, the dependent variable will equal zero when the values of the independent variables are equal to their means—i.e., zero). The sample standardized partial slopes are computed, in general, by the following equation:

$$b_k^* = b_k \frac{s_k}{s_Y}$$

For the two-predictor case, the standardized partial slopes can be calculated by

$$b_1^* = b_1 \frac{s_1}{s_Y}$$

$$\text{or } b_1^* = \frac{r_{Y1} - r_{Y2} r_{12}}{(1 - r_{12}^2)}$$

and

$$b_2^* = b_2 \frac{s_2}{s_Y}$$

$$\text{or } b_2^* = \frac{r_{Y2} - r_{Y1} r_{12}}{(1 - r_{12}^2)}$$

If the two predictors are uncorrelated (i.e., $r_{12} = 0$), then the standardized partial slopes are equal to the simple bivariate correlations between the dependent variable and the independent variables (i.e., $b_1^* = r_{Y1}$ and $b_2^* = r_{Y2}$) because the rest of the equation goes away, as we see here. In the ‘mathematical introduction snapshot’ in Section 4.2 we provide an illustration of this using the example data in the chapter.

$$b_1^* = \frac{r_{Y1} - r_{Y2} r_{12}}{(1 - r_{12}^2)} = \frac{r_{Y1} - r_{Y1}(0)}{(1 - 0)} = r_{Y1}$$

When would you want to use the standardized versus unstandardized regression analyses? According to Pedhazur (1997), b_k^* is sample specific and is not very stable across different samples due to the variance of X_k changing (as the variance of X_k increases, the value of b_k^* also increases, all else being equal). For example, at Ivy-Covered University, b_k^* would vary across different graduating classes (or samples), while b_k would be much more consistent across classes. Thus, most researchers prefer the use of b_k to compare the influence of a particular predictor variable across different samples and/or populations. Pedhazur also states that the b_k^* are of “limited value” (p. 321), but could be reported along with the b_k . As Pedhazur and others have reported, the b_k^* can be deceptive in determining the relative importance of the predictors as they are affected by the variances and covariances of both the included predictors and the predictors not included in the model. Thus, we recommend the b_k for general purpose use.

4.1.1.5 Coefficient of Multiple Determination and Multiple Correlation

An obvious question now is, How well is the criterion variable predicted or explained by the set of predictor variables? For our example, we are interested in how well the graduate grade point averages (the dependent variable) are predicted by the GRE total scores and the undergraduate grade point averages. In other words, what is the utility of the set of predictor variables?

The simplest method involves the partitioning of the familiar total sum of squares in Y , which we denote as SS_{total} . In multiple linear regression analysis, we can write SS_{total} as follows:

$$SS_{total} = [n \sum Y_i^2 - (\sum Y_i)^2] / n$$

$$\text{or } SS_{total} = (n - 1) s_Y^2$$

where we sum over Y from $i = 1, \dots, n$. Next we can conceptually partition SS_{total} as

$$SS_{total} = SS_{reg} + SS_{res}$$

$$\Sigma(Y_i - \bar{Y})^2 = \Sigma(Y'_i - \bar{Y}')^2 + \Sigma(Y_i - Y'_i)^2$$

where SS_{reg} is the regression sum of squares due to the prediction of Y from the X_k 's (often written as $SS_{Y'}$), and SS_{res} is the sum of squares due to the residuals.

Before we consider computation of SS_{reg} and SS_{res} , let us look at the coefficient of multiple determination. Recall the coefficient of determination that is applicable to simple linear regression, r_{XY}^2 . We now consider the multiple predictor version of r_{XY}^2 , here denoted as $R_{Y,1,\dots,m}^2$. The subscript tells us that Y is the criterion (or dependent) variable and that X_1, \dots, m are the predictor (or independent) variables (with m representing the total number of independent variables). The simplest procedure for computing R^2 is as follows:

$$R_{Y,1,\dots,m}^2 = b_1^* r_{Y1} + b_2^* r_{Y2} + \dots + b_m^* r_{Ym}$$

The coefficient of multiple determination tells us the proportion of total variation in the dependent variable Y that is predicted from the set of predictor variables (i.e., X_1, \dots, m 's). Often we see the coefficient in terms of SS as

$$R_{Y,1,\dots,m}^2 = SS_{reg} / SS_{total}$$

Thus, one method for computing the sums of squares regression and residual, SS_{reg} and SS_{res} , is from the coefficient of multiple determination, R^2 (an index that can also be used a measure of effect size) as follows:

$$SS_{reg} = R^2 SS_{total}$$

$$SS_{res} = (1 - R^2) SS_{total} = SS_{total} - SS_{reg}$$

Note also that $R^2_{Y_1, \dots, Y_m}$ is referred to as the *multiple correlation coefficient* so as not to confuse it with a simple bivariate correlation coefficient. In the ‘mathematical introduction snapshot’ in Section 4.2 we provide an illustration using the example data in the chapter.

It should be noted that R^2 is sensitive to sample size and to the number of predictor variables. As sample size and/or the number of predictor variables increase, R^2 will increase as well. R^2 is a biased estimate of the population multiple correlation due to sampling error in the bivariate correlations and in the standard deviations of X and Y . Because R^2 systematically overestimates the population multiple correlation, an adjusted coefficient of multiple determination has been devised. The adjusted R^2 (R^2_{adj}) is calculated as follows:

$$R^2_{adj} = 1 - (1 - R^2) \left(\frac{n - 1}{n - m - 1} \right)$$

Thus, R^2_{adj} adjusts for sample size and for the number of predictors in the model; this allows us to compare models fitted to the same set of data with different numbers of predictors or with different samples of data. The difference between R^2 and R^2_{adj} is called **shrinkage**.

When n is small relative to m , the amount of bias can be large, as R^2 can be expected to be large by chance alone. In this case, the adjustment will be quite large, as it should be. In addition, with small samples, the regression coefficients (i.e., the b_k 's) may not be very good estimates of the population values. When n is large relative to m , bias will be minimized and generalizations are likely to be better about the population values.

For the example data, we determine the adjusted multiple coefficient of determination R^2_{adj} to be as follows:

$$R^2_{adj} = 1 - (1 - R^2) \left(\frac{n - 1}{n - m - 1} \right) = 1 - (1 - .9089) \left(\frac{11 - 1}{11 - 2 - 1} \right) = .8861$$

In this case, the adjusted multiple coefficient of determination indicates a very small adjustment in comparison to R^2 .

4.1.1.6 Significance Tests

Here we describe two procedures used in multiple linear regression analysis. These involve testing the significance of the overall regression model and of each individual partial slope (or regression coefficient).

Test of Significance of the Overall Regression Model

The first test is the test of significance of the overall regression model, or alternatively the test of significance of the coefficient of multiple determination. This is a test of all of the b_k 's simultaneously, an examination of overall model fit of the independent variables in aggregate. The null and alternative hypotheses, respectively, are as follows:

$$H_0: \beta_1 = \beta_2 = \dots = \beta_k = 0$$

$$H_1: \text{not all the } \beta_k = 0$$

If H_0 is rejected, then one or more of the individual regression coefficients (i.e., the b_k) is statistically significantly different from zero (if the assumptions are satisfied, as discussed later). If H_0 is not rejected, then none of the individual regression coefficients will be significantly different from zero.

The test is based on the following test statistic:

$$F = \frac{R^2 / m}{(1 - R^2) / (n - m - 1)}$$

where F indicates that this is an F statistic, m is the number of predictors or independent variables, and n is the sample size. The F test statistic is compared to the F critical value, always a one-tailed test (by default, this value can never be negative given the terms in the equation, so this will always be a nondirectional test) and at the designated level of significance, with degrees of freedom being m and $(n - m - 1)$. That is, the tabled critical value is $\alpha F_{m,(n - m - 1)}$. The test statistic can also be written in equivalent form as

$$F = \frac{SS_{reg} / df_{reg}}{SS_{res} / df_{res}} = \frac{MS_{reg}}{MS_{res}}$$

Where the degrees of freedom regression equals the number of independent variables, $df_{reg} = m$, and degrees of freedom residual equals the difference between the sample size, number of independent variables, and one, $df_{res} = (n - m - 1)$.

Test of Significance of b_k

The second test is the test of the statistical significance of each individual partial slope or regression coefficient, b_k . That is, are the individual unstandardized regression coefficients statistically significantly different from zero? This is actually the same as the test of b_k^* , so we need not develop a separate test for b_k^* . The null and alternative hypotheses, respectively, are as follows:

$$H_0: \beta_k = 0$$

$$H_1: \beta_k \neq 0$$

where β_k is the population partial slope for X_k .

In multiple regression, it is necessary to compute a standard error for each regression coefficient b_k . The variance error of estimate is similarly defined for multiple linear regression and computed as:

$$S_{res}^2 = \frac{SS_{res}}{df_{res}} = MS_{res}$$

where $df_{res} = (n - m - 1)$. Degrees of freedom are lost as we have to estimate the population partial slopes and intercept, the β_k 's and α , respectively, from the sample data. The variance error of estimate indicates the amount of variation among the residuals. The standard error of estimate is simply the positive square root of the variance error of estimate and is the standard deviation of the residuals or errors of estimate. We call it the **standard error of estimate**, denoted as s_{res} .

Finally, we need to compute a standard error for each b_k . Denote the standard error of b_k as $s(b_k)$ and define it as

$$s(b_k) = \frac{s_{res}}{\sqrt{(n-1)s_k^2(1-R_k^2)}}$$

where s_k^2 is the sample variance for predictor X_k , and R_k^2 is the squared multiple correlation between X_k and the remaining X_i 's. R_k^2 represents the overlap between that predictor (X_k) and the remaining predictors. In the case of two predictors, the squared multiple correlation R_k^2 is equal to the simple bivariate correlation between the two independent variables r_{12}^2 .

The test statistic for testing the significance of the regression coefficients, b_k , is as follows:

$$t = \frac{b_k}{s(b_k)}$$

The test statistic t is compared to the critical values of t , a two-tailed test for a non-directional H_1 , at the designated level of significance, and with degrees of freedom $(n - m - 1)$. Thus, the tabled critical values are $+/-_{(a/2)} t_{(n-m-1)}$ for a two-tailed test.

We can also form a confidence interval around b_k as follows:

$$CI(b_k) = b_k + {}_{(a/2)} t_{(n-m-1)} s(b_k)$$

Recall that the null hypothesis tested is $H_0: \beta_k = 0$. Therefore, if the confidence interval contains zero, then the regression coefficient b_k is not statistically significantly different from zero at the specified α level. This is interpreted to mean that in $(1 - \alpha)\%$ of

the sample confidence intervals that would be formed from multiple samples, β_k will be included. In the ‘mathematical introduction snapshot’ in Section 4.2, we provide an illustration using the example data in the chapter.

Other Tests

One can also form confidence intervals for the predicted mean of Y and the prediction intervals for individual values of Y .

4.1.1.7 Methods of Entering Predictors

The multiple predictor model that we have considered thus far can be viewed as **simultaneous regression**. That is, all of the predictors to be used are entered (or selected) simultaneously, such that all of the regression parameters are estimated simultaneously; here the set of predictors has been selected *a priori*. In computing these regression models, we have used the default setting in SPSS of the method of entry as ‘Enter,’ which enters the set of independent variables in aggregate. There are other methods of entering the independent variables where the predictor variables are entered (or selected) systematically; here the set of predictors has not been selected *a priori*. This class of models is referred to as **sequential regression** (also known as **variable selection procedures**). This section introduces a brief description of the following sequential regression procedures: backward elimination, forward selection, stepwise selection, all possible subsets regression, and hierarchical regression.

Backward Elimination

First, consider the backward elimination procedure. Here, variables are eliminated from the model based on their minimal contribution to the prediction of the criterion variable. In the first stage of the analysis, all potential predictors are included in the model. In the second stage, that predictor is deleted from the model that makes the smallest contribution to the prediction of the dependent variable. This can be done by eliminating that variable having the smallest t or F statistic, such that it is making the smallest contribution to R^2_{adj} . In subsequent stages, the predictor is deleted that makes the next smallest contribution to the prediction of the outcome Y . The analysis continues until each of the remaining predictors in the model is a significant predictor of Y . This could be determined by comparing the t or F statistics for each predictor to the critical value, at a preselected level of significance. Some computer programs use as a stopping rule the maximum F -to-remove criterion, where the procedure is stopped when all of the selected predictors’ F values are greater than the specified F criterion. Another stopping rule is where the researcher stops at a predetermined number of predictors (see Hocking, 1976; Thompson, 1978). In SPSS, this is the **backward** method of entering predictors.

Forward Selection

In the forward selection procedure, variables are added or selected into the model based on their maximal contribution to the prediction of the criterion variable. Initially, none of the potential predictors is included in the model. In the first stage, the predictor is added to the model that makes the largest contribution to the prediction of the dependent variable. This can be done by selecting that variable having the largest t or F statistic, such that it is making the largest contribution to R^2_{adj} . In subsequent stages, the predictor is selected that makes the next largest contribution to the prediction of Y . The analysis continues until each of the selected predictors in the model is a significant predictor of the outcome Y , whereas none of the unselected predictors is a significant predictor. This could be determined by comparing the t or F statistics for each predictor to the critical value, at a preselected level of significance. Some computer programs use as a stopping rule the minimum F -to-enter criterion, where the procedure is stopped when all of the unselected predictors' F values are less than the specified F criterion. For the same set of data and at the same level of significance, the backward elimination and forward selection procedures may not necessarily result in the exact same final model due to the differences in how variables are selected. In SPSS, this is the **forward** method of entering predictors.

Stepwise Selection

The stepwise selection procedure is a modification of the forward selection procedure with one important difference. Predictors that have been selected into the model can at a later step be deleted from the model; thus, the modification conceptually involves a backward elimination mechanism. This situation can occur for a predictor when a significant contribution at an earlier step later becomes a nonsignificant contribution, given the set of other predictors in the model. Thus, a predictor loses its significance due to new predictors being added to the model.

The stepwise selection procedure is as follows. Initially, none of the potential predictors is included in the model. In the first step, the predictor is added to the model that makes the largest contribution to the explanation of the dependent variable. This can be done by selecting that variable having the largest t or F statistic, such that it is making the largest contribution to R^2_{adj} . In subsequent stages, the predictor is selected that makes the next largest contribution to the prediction of Y . Those predictors that have been entered at earlier stages are also checked to see if their contribution remains significant. If not, then that predictor is eliminated from the model. The analysis continues until each of the predictors remaining in the model is a significant predictor of Y , while none of the other predictors is a significant predictor. This could be determined by comparing the t or F statistics for each predictor to the critical value, at a specified level of significance. Some computer programs use as stopping rules the minimum F -to-enter and maximum F -to-remove criteria, where the F -to-enter value selected is usually equal to or slightly greater than the F -to-remove value selected (to prevent a predictor from continuously being entered and removed). For the same set of data

and at the same level of significance, the backward elimination, forward selection, and stepwise selection procedures may not necessarily result in the exact same final model, due to differences in how variables are selected. In SPSS, this is the **stepwise** method of entering predictors.

All Possible Subsets Regression

Another sequential regression procedure is known as all possible subsets regression. Let us say, for example, that there are five potential predictors. In this procedure, all possible one-, two-, three-, and four-variable models are analyzed (with five predictors there is only a single five-predictor model). Thus, there will be 5 one-predictor models, 10 two-predictor models, 10 three-predictor models, and 5 four-predictor models. The best k predictor model can be selected as the model that yields the largest R^2_{adj} . For example, the best three-predictor model would be that model of the 10 estimated that yields the largest R^2_{adj} . With today's powerful computers, this procedure is easier and more cost efficient than in the past. However, the researcher is not advised to consider this procedure, or for that matter any of the other sequential regression procedures, when the number of potential predictors is large. Here, the researcher is allowing number crunching to take precedence over thoughtful analysis. Also, the number of models will be equal to 2^n , so that for 10 predictors there are 1,024 possible subsets. Obviously, examining that number of models is not a thoughtful analysis.

Hierarchical Regression

In hierarchical regression, the researcher specifies a priori a sequence for the individual predictor variables (not to be confused with hierarchical linear models, which is a regression approach for analyzing nested data collected at multiple levels, such as child, classroom, and school). The analysis proceeds in a forward selection, backward elimination, or stepwise selection mode according to a researcher specified, theoretically based sequence, rather than an unspecified statistically based sequence. This variable selection method is different from those previously discussed, in that the researcher determines the order of entry from a careful consideration of the available theory research, instead of the software dictating the sequence.

A type of hierarchical regression is known as **setwise regression** (also called **blockwise**, **chunkwise**, or **forced stepwise regression**). Here the researcher specifies a priori a sequence for sets of predictor variables. This procedure is similar to hierarchical regression in that the researcher determines the order of entry of the predictors. The difference is that the setwise method uses sets of predictor variables at each stage rather than one individual predictor variable at a time. The sets of variables are determined by the researcher so that variables within a set share some common theoretical ground (e.g., home background variables in one set and aptitude variables in another set). Variables within a set are selected according to one of the sequential regression procedures. The variables selected for a particular set are then entered in the specified

theoretically based sequence. In SPSS, this is conducted by entering predictors in **blocks** and selecting their desired method of entering variables in each block (e.g., simultaneously, forward, backward, stepwise).

Commentary on Sequential Regression Procedures

Let us make some comments and recommendations about the sequential regression procedures. First, numerous statisticians have noted problems with stepwise methods (i.e., backward elimination, forward selection, and stepwise selection) (e.g., Derksen & Keselman, 1992; Huberty, 1989; Mickey, Dunn, & Clark, 2004; Miller, 1984, 1990; Wilcox, 2003). These problems include the following: (a) selecting noise rather than important predictors; (b) highly inflated R^2 and R^2_{adj} values; (c) confidence intervals for partial slopes that are too narrow; (d) p values that are not trustworthy; (e) important predictors being barely edged out of the model, making it possible to miss the true model; and (f) potentially heavy capitalization on chance given the number of models analyzed. Second, theoretically based regression models have become the norm in many disciplines (and the stepwise methods of entry are driven by mathematics of the models rather than theory). Thus, hierarchical regression either has or will dominate the landscape of the sequential regression procedures. Therefore, we strongly encourage you to consider more extended discussions of hierarchical regression (e.g., Bernstein, 1988; Cohen & Cohen, 1983; Pedhazur, 1997; Schafer, 1991; Tabachnick & Fidell, 2007).

If you are working in an area of inquiry where research evidence is scarce or nonexistent, then you are conducting exploratory research. Thus, you are probably trying to simply identify the key variables. Here, hierarchical regression is not appropriate, as a theoretically driven sequence cannot be developed since there is no theory to guide its development. We therefore recommend the use of all possible subsets regression (e.g., Kleinbaum, Kupper, Muller, & Nizam, 1998). For additional information on the sequential regression procedures, see Cohen and Cohen (1983), Kleinbaum et al. (1998), Miller (1990), Pedhazur (1997), and Weisberg (1985).

4.1.1.8 Nonlinear Relationships

Here we discuss how to deal with nonlinearity. We formally introduce several multiple regression models for when the criterion variable does not have a linear relationship with the predictor variables.

First, consider polynomial regression models. In polynomial models, powers of the predictor variables (e.g., squared, cubed) are used. In general, a sample polynomial regression model that includes one quadratic term is as follows:

$$Y = b_1 X + b_2 X^2 + \dots + b_m X^m + a + e$$

where the independent variable X is taken from the first power through the m^{th} power, and the i subscript for observations has been deleted to simplify matters. If the model consists only of X taken to the first power, then this is a **simple linear regression model** (or **first-degree polynomial**; this is a straight line and what we have studied to this point). A **second-degree polynomial** includes X taken to the second power (or **quadratic model**; this is a curve with one bend in it rather than a straight line). A **third-degree polynomial** includes X taken to the third power (or **cubic model**; this is a curve with two bends in it).

A polynomial model with multiple predictors can also be utilized. An example of a second-degree polynomial model with two predictors is illustrated in the following equation:

$$Y = b_1 X_1 + b_2 X_1^2 + b_3 X_2 + b_4 X_2^2 + a + e$$

It is important to note that whenever a higher-order polynomial is included in a model (e.g., quadratic, cubic, and more), the first-order polynomial must also be included in the model. In other words, it is not appropriate to include a quadratic term X^2 without also including the first-order polynomial X . For more information on polynomial regression models, see Bates and Watts (1988), Kleinbaum et al. (1998), Pedhazur (1997), Seber and Wild (1989), and Weisberg (1985). Alternatively, one might transform the criterion variable and/or the predictor variables to obtain a more linear form, as previously discussed.

4.1.1.9 Interactions

Another type of model involves the use of an interaction term, a term with which you may be familiar from factorial ANOVA. These can be implemented in any type of regression model. We can write a simple two-predictor interaction-type model as

$$Y = b_1 X_1 + b_2 X_2 + b_3 X_1 X_2 + a + e$$

where $X_1 X_2$ represents the interaction of predictor variables 1 and 2. An interaction can be defined as occurring when the relationship between Y and X_1 depends on the level of X_2 . In other words, X_2 is a **moderator variable**. For example, suppose one were to use years of education and age to predict political attitude. The relationship between education and attitude might be moderated by age. In other words, the relationship between education and attitude may be different for older versus younger individuals. If age were a moderator, we would expect there to be an interaction between age and education in a regression model. Note that if the predictors are very highly correlated, collinearity is likely. For more information on interaction models, see Cohen and Cohen (1983), Berry and Feldman (1985), Kleinbaum et al. (1998), Meyers, Gamst, and Guarino (2006), and Weinberg and Abramowitz (2002).

4.1.1.10 Categorical Predictors

So far we have only considered continuous predictors—*independent variables that are interval or ratio in scale*. There may be times, however, that you wish to use a categorical predictor—an independent variable that is nominal or ordinal in scale. For example, gender, grade level (e.g., freshman, sophomore, junior, senior), highest education earned (less than high school, high school graduate, etc.) are all categorical variables that may be very interesting and theoretically appropriate to include in either a simple or multiple regression model. Given their scale (i.e., nominal or ordinal), however, we must recode the values prior to analysis so that they are on a scale of zero and one. This is called ‘dummy coding,’ as this type of recoding makes the model work. For example, males might be coded as zero and females coded as one. When there are more than two categories to the categorical predictor, multiple dummy-coded variables must be created—*specifically one minus the number of levels or categories of the categorical variable*. Thus, in the case of grade level where there are four categories (freshman, sophomore, junior, senior), three of the four categories would be dummy coded and included in the regression model as predictors. The category that is ‘left out’ is the reference category, or that category to which all other levels are compared. The easiest way to understand this is perhaps to examine the data. In the screenshot that follows, the first column represents grade level where 1 = freshman, 2 = sophomore, 3 = junior, and 4 = senior. Dummy coding three of the four grade levels (check SPSS on if this is needed), with ‘senior’ as the reference category, will result in three additional columns (columns 2, 3, and 4 in the screenshot).

	Grade	Freshman	Sophomore	Junior
1	1.00	1.00	.0	.0
2	1.00	1.00	.0	.0
3	1.00	1.00	.0	.0
4	2.00	.0	1.00	.0
5	2.00	.0	1.00	.0
6	2.00	.0	1.00	.0
7	3.00	.0	.0	1.00
8	3.00	.0	.0	1.00
9	3.00	.0	.0	1.00
10	4.00	.0	.0	.0
11	4.00	.0	.0	.0
12	4.00	.0	.0	.0

In terms of generating the analysis and the point and click using SPSS to compute the regression model, nothing changes. The steps are the same regardless of whether the predictors are continuous or categorical. Now let us discuss *why* dummy coding works in this situation. The point biserial correlation is appropriate when one variable is dichotomous and the other variable is interval or ratio. The point biserial correlation is a variant of the Pearson product-moment correlation, and we can use the Pearson as a variant of the point biserial. Thus, while we will not have a linear relationship between a continuous outcome and a binary variable, the mathematics that underlie the model will hold.

Consider an example output for predicting grade point average (GPA) based on grade level, where ‘senior’ is the reference category. We see that the intercept (i.e., ‘constant’) is statistically significant, as is ‘freshman.’ The interpretation of the intercept remains the same regardless of the scale of the predictors. The intercept represents grade point average (the dependent variable) when all the predictors are zero. In this case, this means that grade point average is 3.267 for *seniors* (the reference category). The only statistically significant predictor is ‘freshman.’ This is interpreted to say that mean GPA decreases by .800 points for freshmen *as compared to seniors*. The nonstatistically significant regression coefficients for ‘sophomore’ and ‘junior’ indicate that mean GPA is similar for these grade levels as compared to seniors. The interpretation for dummy variable predictors is always in reference to the category that was ‘left out.’ In this case, that was ‘seniors.’

Coefficients^a

Model	Unstandardized Coefficients		Beta	t	Sig.
	B	Std. Error			
1 (Constant)	3.267	.183		17.892	.000
Freshman	-.800	.258	-.704	-3.098	.015
Sophomore	.233	.258	.205	.904	.393
Junior	.200	.258	.176	.775	.461

a. Dependent Variable: GPA

It is important to note that even though ‘sophomore’ and ‘junior’ were not statistically significant, they should be retained in the model as they represent (along with ‘freshman’) a group. Dropping one or more dummy-coded indicator variables that represent a group will change the reference category. For example, if ‘sophomore’ and ‘junior’ were dropped from the model, the interpretation would then become the mean GPA for freshmen *as compared to all other grade levels*. Thus, careful thought needs to be put into dropping one or more indicators that are part of a set.

4.1.2 Sample Size

There is a fair body of research that has examined minimum sample size in the context of multiple linear regression and some consensus that sample-size considerations differ depending on the goal of your research—either testing a hypothesis test or estimating a parameter (Algina & Olejnik, 2000; Maxwell, 2000)—with larger sample sizes needed for estimation (e.g., Pedhazur, 1997). Using simulation research, recent research suggests that the squared multiple correlation coefficient does have a relationship with overall sample size and the ratio of the sample size to predictors. Larger sample sizes are needed as the squared multiple correlation coefficient diminishes (Knofszynski, 2008). More specifically, as the squared multiple correlation coefficient nears zero, there is a quicker increase in sample size and this pattern is constant across varying numbers of predictors; however, the sample size does not dramatically increase as the number of predictors increases. For example, with a squared multiple correlation coefficient of .10 and three predictors, a sample size of 1,800 is needed to achieve “excellent prediction level” (Knofszynski, 2008, p. 438). In comparison, with a square multiple correlation coefficient of .50, again with three predictors, a sample size of 220 is needed to achieve “excellent” (Knofszynski, 2008). As we know, sample size and power are inextricably intertwined, and attempting to separate the two is futile. The best recommendation is to estimate power using power software and to consult current advances based on simulation research such as Knofszynski (2008).

4.1.3 Power

With a large number of predictors, power is reduced, and there is an increased likelihood of a Type I error across the total number of significance tests (i.e., one for each predictor and overall, as we show in the next section). In multiple regression, power is a function of sample size, the number of predictors, the level of significance, and the size of the population effect (i.e., for a given predictor, or overall). To determine how large a sample you need relative to the number of predictors, we suggest that you consult power tables (e.g., Cohen, 1988) or power software (e.g., Murphy & Myors, 2004; Power and Precision; G*Power).

4.1.4 Effect Size

Effect size in multiple linear regression can be gauged by the coefficient of multiple determination or multiple correlation coefficient, introduced previously. The coefficient of multiple determination indicates the proportion of total variation in the dependent variable Y that is predicted from the set of predictor variables. There is no objective gold standard as to how large the coefficient of determination needs to be in order to say that a meaningful proportion of variation has been predicted. The coefficient is determined not just by the quality of the predictor variables included in the

model, but also by the quality of relevant predictor variables not included in the model, as well as by the amount of total variation in the dependent variable Y . According to the subjective standard of Cohen (1988), a small effect size is defined as $R_m^2 = .02$, a medium effect size as $R_m^2 = .13$, and a large effect size as $R_m^2 = .26$. The squared multiple correlation coefficient can also be used to compute a globalized f^2 , which is

$$f^2 = \frac{R_m^2}{1 - R_m^2}.$$

A computation of f^2 that allows for a localized effect is $f^2 = \frac{R_{AB}^2 - R_A^2}{1 - R_{AB}^2}$,

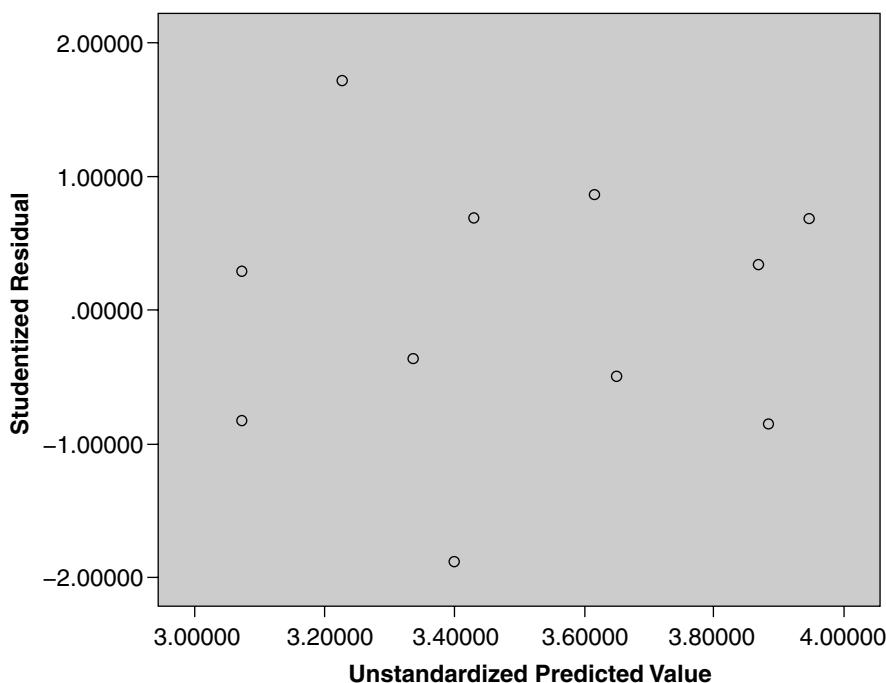
where R_{AB}^2 equals the proportion of variance accounted for by the model, R_A^2 equals the proportion of variance accounted for by the set of predictors excluding B , and R_B^2 equals the proportion of variance accounted for by predictor B (i.e., that variable that is of interest for the localized effect). The numerator, therefore, reflects the unique proportion of variance for which predictor B accounts. For additional information on effect size measures in regression, we suggest you consider Steiger and Fouladi (1992), Mendoza and Stafford (2001), and Smithson (2001; which also includes some discussion of power).

4.1.5 Assumptions

For the most part, the assumptions of multiple linear regression analysis are the same as that with simple linear regression. The assumptions are concerned with (a) independence, (b) homoscedasticity, (c) normality, (d) linearity, (e) fixed X , and (f) noncollinearity. This section also mentions those techniques appropriate for evaluating each assumption.

4.1.5.1 Independence

The first assumption is concerned with **independence** of the observations. The simplest procedure for assessing independence is to examine residual plots of e versus the predicted values of the dependent variable Y' and of e versus each independent variable X_k (alternatively, one can look at plots of observed values of the dependent variable Y versus predicted values of the dependent variable Y' and of observed values of the dependent variable Y versus each independent variable X_k). If the independence assumption is satisfied, the residuals should fall into a random display of points. If the assumption is violated, the residuals will fall into some sort of pattern. Lack of independence affects the estimated standard errors of the model. For serious violations, one could consider generalized or weighted least squares as the method of estimation (e.g., Myers, 1986; Weisberg, 1985), or some type of transformation. The residual plots shown in Figure 4.1 do not suggest any independence problems for the GGPA example, where Figure 4.1(a) represents the residual e versus the predicted value of the dependent variable Y' , Figure 4.1(b) represents e versus GRETOT, and Figure 4.1(c) represents e versus UGPA.

**FIGURE 4.1**

Residual Plots for GRE-GPA Example

4.1.5.2 Homoscedasticity

The second assumption is **homoscedasticity**, where the conditional distributions have the same constant variance for all values of X . In the residual plots, the consistency of the variance of the conditional distributions may be examined. If the homoscedasticity assumption is violated, estimates of the standard errors are larger, and the conditional distributions may also be nonnormal. Solutions include variance stabilizing transformations (such as the square root or log of Y), generalized or weighted least squares (e.g., Myers, 1986; Weisberg, 1985), or robust regression (Kleinbaum et al., 1998; Myers, 1986; Wilcox, 1996, 2003; Wu, 1985). Due to the small sample size, homoscedasticity cannot really be assessed for the example data.

4.1.5.3 Normality

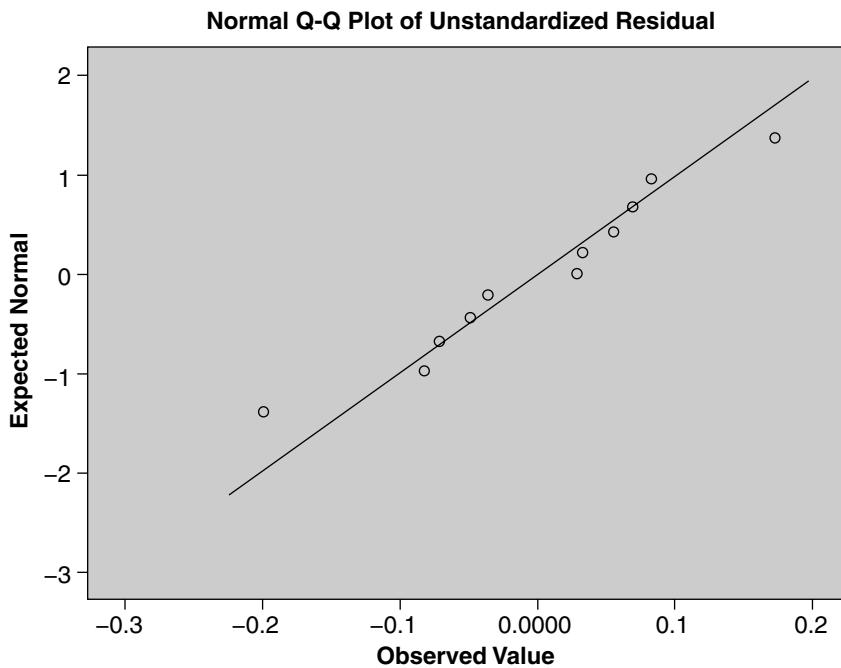
The third assumption is that the conditional distributions of the scores on Y , or the prediction errors, are **normal** in shape. Violation of the normality assumption may be the result of outliers. Various recommendations are used to crudely detect outliers from a residual plot or scatterplot. The simplest outlier detection procedure and a commonly used rule is to define an outlier as an observation more than two or three standard errors from the mean (i.e., a large distance from the mean). The outlier

observation may be a result of (a) a simple recording or data entry error, (b) an error in observation, (c) an improperly functioning instrument, (d) inappropriate use of administration instructions, or (e) a true outlier. If the outlier is the result of an error, correct the error if possible and redo the regression analysis. If the error cannot be corrected, then the observation could be deleted. If the outlier represents an accurate observation, then this observation may contain important theoretical information, and one would be more hesitant to delete it (or perhaps seek out similar observations).

Several methods for dealing with outliers are available. A simple procedure to use for single case outliers (i.e., just one outlier) is to perform two regression analyses, both with and without the outlier being included. A comparison of the regression results will provide some indication of the effects of the outlier. Other methods include robust regression (Kleinbaum et al., 1998; Myers, 1986; Wilcox, 1996, 2003; Wu, 1985), and nonparametric regression (Miller, 1997; Rousseeuw & Leroy, 1987; Wu, 1985).

The following can be used to detect normality violations: frequency distributions, normal probability (Q-Q) plots, and skewness statistics. The simplest procedure involves checking for symmetry in a histogram, frequency distribution, boxplot, or skewness and kurtosis statistics. Although **nonzero kurtosis** (i.e., a distribution that is either flat, platykurtic, or has a sharp peak, leptokurtic) will have minimal effect on the regression estimates, **nonzero skewness** (i.e., a distribution that is not symmetric with either a positive or negative skew) will have much more impact on these estimates. Thus, finding asymmetrical distributions is necessary. One suggestion is to be concerned if the skewness value is larger than 1.5 or 2.0 in magnitude.

Another useful graphical technique is the normal probability plot (or Q-Q plot). With normally distributed residuals, the points on the normal probability plot will fall along a straight diagonal line, whereas nonnormal data will not. There is a difficulty with this plot because there is no criterion with which to judge deviation from linearity. It is recommended that skewness and/or the normal probability plot be considered at a minimum when determining normality evidence. For the example data, the normal probability plot is shown in Figure 4.2, and even with a small sample looks good. Violation can lead to imprecision in the partial slopes and in the coefficient of determination. There are also several statistical procedures available for the detection of nonnormality (e.g., Andrews, 1971; Belsley, Kuh, & Welsch, 1980; D'Agostino, 1971; Ruppert & Carroll, 1980; Shapiro & Wilk, 1965; Wu, 1985); transformations can also be used to normalize the data. The most commonly used transformations to correct for nonnormality in regression analysis are to transform the dependent variable using the log (to correct for positive skew) or the square root (to correct for positive or negative skew). However, again there is the problem of dealing with transformed variables measured along some other scale than that of the original variables.

**FIGURE 4.2**

Normal Probability Plot for GRE-GPA Example

4.1.5.4 Linearity

The fourth assumption is **linearity**, that there is a linear relationship between the observed scores on the dependent variable Y and the values of the independent variables, X_k 's. If satisfied, then the sample partial slopes and intercept are unbiased estimators of the population partial slopes and intercept, respectively. The linearity assumption is important because regardless of the value of X_k , we always expect Y to increase by b_k units for a one-unit increase in X_k , controlling for the other X_k 's. If a nonlinear relationship exists, this means that the expected increase in Y depends on the value of X_k ; that is, the expected increase is not a constant value. Strictly speaking, linearity in a model refers to there being linearity in the parameters of the model (i.e., α and the β 's).

Violation of the linearity assumption can be detected through residual plots. The residuals should be located within a band of $\pm 2 s_{res}$ (or standard errors), indicating no systematic pattern of points. Residual plots for the GGPA example are shown in Figure 4.1. Even with a very small sample, we see a fairly random pattern of residuals, and therefore feel fairly confident that the linearity assumption has been satisfied. Note also that there are other types of residual plots developed especially for multiple regression analysis, such as the added variable and partial residual plots (Larsen & McCleary, 1972; Mansfield & Conerly, 1987; Weisberg, 1985). Procedures to deal with nonlinearity include transformations (of one or more of the X_k 's and/or of Y) and other regression models (discussed later in this chapter).

4.1.5.5 Fixed X

The fifth assumption is that the values of X_k are **fixed**, where the independent variables X_k are fixed variables rather than random variables. This results in the regression model being valid only for those particular values of X_k that were actually observed and used in the analysis. Thus, the same values of X_k would be used in replications or repeated samples.

Strictly speaking, the regression model and its parameter estimates are only valid for those values of X_k actually sampled. The use of a prediction model developed to predict the dependent variable Y , based on one sample of individuals, may be suspect for another sample of individuals. Depending on the circumstances, the new sample of individuals may actually call for a different set of parameter estimates. Generally, we may not want to make predictions about individuals having combinations of X_k scores outside of the range of values used in developing the prediction model; this is defined as *extrapolating* beyond the sample predictor data. On the other hand, we may not be quite as concerned in making predictions about individuals having combinations of X_k scores within the range of values used in developing the prediction model; this is defined as *interpolating* within the range of the sample predictor data.

It has been shown that when other assumptions are met, regression analysis performs just as well when X is a random variable (e.g., Glass & Hopkins, 1996; Myers & Well, 1995; Pedhazur, 1997; Wonnacott & Wonnacott, 1981). There is no such assumption about Y .

4.1.5.6 Noncollinearity

Considering simple and multiple linear regression, the final assumption is unique to multiple linear regression analysis (as compared to simple linear regression), but will be quite common throughout multivariate procedures that we cover. A violation of this assumption is known as collinearity, where there is a very strong linear relationship between two or more of the predictors. The presence of severe collinearity is problematic in several respects. First, it will lead to instability of the regression coefficients across samples, where the estimates will bounce around quite a bit in terms of magnitude and even occasionally result in changes in sign (perhaps opposite of expectation). This occurs because the standard errors of the regression coefficients become larger, making it more difficult to achieve statistical significance. Another result that may occur involves an overall regression that is significant, but none of the individual predictors is significant. Collinearity will also restrict the utility and generalizability of the estimated regression model.

Recall from earlier in the chapter the notion of partial regression coefficients, where the other predictors were held constant. In the presence of severe collinearity, the other predictors cannot really be held constant because they are so highly

intercorrelated. Collinearity may be indicated when there are large changes in estimated coefficients due to (a) a variable being added or deleted, and/or (b) an observation being added or deleted (Chatterjee & Price, 1977). Singularity is a special case of multicollinearity; it is perfect multicollinearity and occurs when two or more items/variables perfectly predict and are therefore perfectly redundant. This can occur in factor analysis (just as it did in multiple regression), for example, when a composite variable as well as its component variables are used as predictors in the same factor analytic model (e.g., including GRETOT, GRE-Quantitative, and GRE-Verbal as predictors).

How do we detect violations of this assumption? The simplest procedure is to conduct a series of special regression analyses, one for each X , where that predictor is predicted by all of the remaining X 's (i.e., the criterion variable is not involved). If any of the resultant R_k^2 values are close to one (greater than .9 is a good cut point), then there may be a collinearity problem. However, the large R^2 value may also be due to small sample size; thus, more data would be useful. For the example data, $R_{12}^2 = .091$ and therefore collinearity is not a concern.

Also, if the number of predictors is greater than or equal to n , then perfect collinearity is a possibility. Another statistical method for detecting collinearity is to compute a variance inflation factor (VIF) for each predictor, which is equal to $1 / (1 - R_k^2)$. The VIF is defined as the inflation that occurs for each regression coefficient above the ideal situation of uncorrelated predictors. Many suggest that the largest VIF should be less than 10 in order to satisfy this assumption (Myers, 1990; Stevens, 2009; Wetherill, 1986).

There are several possible methods for dealing with a collinearity problem. First, one can remove one or more of the correlated predictors. Second, ridge regression techniques can be used (e.g., Hoerl & Kennard, 1970a, 1970b; Marquardt & Snee, 1975; Myers, 1986; Wetherill, 1986). Third, principal component scores resulting from principal component analysis can be utilized rather than raw scores on each variable (e.g., Kleinbaum et al., 1998; Myers, 1986; Weisberg, 1985; Wetherill, 1986). Fourth, transformations of the variables can be used to remove or reduce the extent of the problem. The final solution, and probably our last choice, is to use simple linear regression, as collinearity cannot exist with a single predictor.

4.1.5.7 Summary of Assumptions

For the GGPA example, although sample size is quite small in terms of looking at conditional distributions, it would appear that all of our assumptions have been satisfied. All of the residuals are within two standard errors of zero, and there does not seem to be any systematic pattern in the residuals. The distribution of the residuals is nearly symmetric and the normal probability plot looks good. A summary of the assumptions and the effects of their violation for multiple linear regression analysis is presented in Table 4.1.

TABLE 4.1

Assumptions and Violation of Assumptions: Multiple Linear Regression Analysis

Assumption	Effect of Assumption Violation
Independence	<ul style="list-style-type: none"> Influences standard errors of the model
Homogeneity	<ul style="list-style-type: none"> Bias in s_{res}^2 May inflate standard errors and thus increase likelihood of a Type II error May result in nonnormal conditional distributions
Normality	<ul style="list-style-type: none"> Less precise slopes, intercept, and R^2
Linearity	<ul style="list-style-type: none"> Bias in slope and intercept Expected change in Y is not a constant and depends on value of X
Fixed X values	<ul style="list-style-type: none"> Extrapolating beyond the range of X combinations: prediction errors larger, may also bias slopes and intercept Interpolating within the range of X combinations: smaller effects than above; if other assumptions met, negligible effect
Noncollinearity of X 's	<ul style="list-style-type: none"> Regression coefficients can be quite unstable across samples (as standard errors are larger) R^2 may be significant, yet none of the predictors are significant Restricted generalizability of the model

4.2 MATHEMATICAL INTRODUCTION SNAPSHOT

Throughout the chapter, we have woven some of the mathematics of multiple linear regression. Now, let's consider the analysis illustrated using the data that Addie, our graduate researcher, is working. We use the GRE Quantitative + Verbal Total (GRE-TOT) and undergraduate grade point average (UGPA) to predict graduate grade point average (GGPA). GRETOT has a possible range of 40 to 160 points (if we remove the unnecessary last digit of zero), and GPA is defined as having a possible range of 0.00 to 4.00 points. Given the sample of 11 statistics students as shown in Table 4.2, let us work through a multiple linear regression analysis.

As sample statistics, we compute for GRETOT (X_1 or subscript 1) that the mean is $\bar{X}_1 = 112.7273$ and the variance is $s_1^2 = 266.8182$, for UGPA (X_2 or subscript 2) that the mean is $\bar{X}_2 = 3.1091$ and the variance is $s_2^2 = 0.1609$, and for GGPA (Y), a mean of $\bar{Y} = 3.5000$ and variance of $s_Y^2 = 0.1100$. In addition, we compute the bivariate correlation between the dependent variable (graduate GPA) and GRE total, $r_{Y1} = .7845$; between the dependent variable (graduate GPA) and undergraduate GPA, $r_{Y2} = .7516$; and between GRE total and undergraduate GPA, $r_{12} = .3011$. The sample partial slopes (b_1 and b_2) and intercept (a) are determined as follows:

$$b_1 = \frac{(r_{Y1} - r_{Y2} r_{12}) s_Y}{(1 - r_{12}^2) s_1} = \frac{[.7845 - (.7516)(.3011)].3317}{(1 - .3011^2)16.3346} = .0125$$

$$b_2 = \frac{(r_{Y2} - r_{Y1} r_{12}) s_Y}{(1 - r_{12}^2) s_2} = \frac{[.7516 - (.7845)(.3011)].3317}{(1 - .3011^2).4011} = .4687$$

TABLE 4.2

GRE-GPA Example Data

Student	GRE-Total (X_1)	Undergraduate GPA (X_2)	Graduate GPA (Y)
1	145	3.2	4.0
2	120	3.7	3.9
3	125	3.6	3.8
4	130	2.9	3.7
5	110	3.5	3.6
6	100	3.3	3.5
7	95	3.0	3.4
8	115	2.7	3.3
9	105	3.1	3.2
10	90	2.8	3.1
11	105	2.4	3.0

$$a = \bar{Y} - b_1 \bar{X}_1 - b_2 \bar{X}_2 = 3.5000 - (.0125)(112.7273) - (.4687)(3.1091) = .6337$$

Let us interpret the partial slope and intercept values. A partial slope of .0125 for GRE-TOT would mean that if your score on the GRETOT was increased by 1 point, then your graduate grade point average would be increased by .0125 points, controlling for undergraduate grade point average. Likewise, a partial slope of .4687 for UGPA would mean that if your undergraduate grade point average was increased by 1 point, then your graduate grade point average would be increased by .4687 points, controlling for GRETOT. An intercept of .6337 would mean that if your scores on the GRETOT and UGPA were both 0, then your graduate grade point average would be .6337. However, it is impossible to obtain a GRETOT score of 0 because you receive 40 points for putting your name on the answer sheet. In a similar way, an undergraduate student could not obtain a UGPA of 0 and be admitted to graduate school. This is not to say that the regression equation is incorrect, but just to point out how the interpretation of “GRETOT and UGPA were both 0” is a bit meaningless in context.

To put all of this together then, the sample multiple linear regression model is

$$Y_i = b_1 X_{1i} + b_2 X_{2i} + a + e_i = .0125 X_{1i} + .4687 X_{2i} + .6337 + e_i$$

In other words, if your score on the GRETOT was 130 and your UGPA was 3.5, then your predicted score on the GGPA would be computed as:

$$Y'_i = .0125 (130) + .4687 (3.5000) + .6337 = 3.8992$$

Based on the prediction equation, we predict your GGPA to be around 3.9; however, predictions are usually somewhat less than perfect, even with two predictors.

For the GGPA example, we compute the overall F test statistic as the following:

$$F = \frac{R^2 / m}{(1 - R^2) / (n - m - 1)} = \frac{.9089 / 2}{(.1 - .9089) / (11 - 2 - 1)} = 39.9078$$

or as

$$F = \frac{SS_{reg} / df_{reg}}{SS_{res} / df_{res}} = \frac{0.9998 / 2}{.1002 / 8} = 39.9122$$

The critical value, at the .05 level of significance, is $F_{.05, 2, 8} = 4.46$. The test statistic exceeds the critical value, so we reject H_0 and conclude that all of the partial slopes are not equal to zero at the .05 level of significance (the two F test statistics differ slightly due to rounding error).

For our graduate grade point average example, the standardized partial slopes are equal to

$$b_1^* = b_1 \frac{s_1}{s_Y} = .0125(16.3346 / .3317) = .6156$$

and

$$b_2^* = b_2 \frac{s_2}{s_Y} = .4687(.4011 / .3317) = .5668$$

The prediction model is then as follows:

$$z(Y_i) = .6156 z_{1i} + .5668 z_{2i}$$

The standardized partial slope of .6156 for GRETOT would be interpreted as the expected increase in GGPA in z score units for a one z score unit increase in the GRETOT, controlling for UGPA. A similar statement may be made for the standardized partial slope of UGPA. The b_k^* can also be interpreted as the expected standard deviation change in the dependent variable Y associated with a one standard deviation change in the independent variable X_k when the other X_k 's are held constant.

With the example of predicting GGPA from GRETOT and UGPA, let us examine the partitioning of the total sum of squares SS_{total} as follows:

$$SS_{total} = (n - 1) s_Y^2 = (10) .1100 = 1.1000$$

Next, we can determine the multiple correlation coefficient R^2 as

$$R^2_{Y, 1, \dots, m} = b_1^* r_{y1} + b_2^* r_{y2} + \dots + b_m^* r_{ym} = .6156 (.7845) + .5668 (.7516) = .9089$$

We can also partition SS_{total} into SS_{reg} and SS_{res} , where

$$SS_{reg} = R^2 SS_{total} = .9089 (1.1000) = 0.9998$$

$$SS_{res} = (1 - R^2) SS_{total} = (1 - .9089) 1.1000 = .1002$$

Finally, let us summarize these results for the example data. We found that the coefficient of multiple determination (R^2) was equal to .9089. Thus, the GRE total score and the undergraduate grade point average predicts around 91% of the variation in the graduate grade point average. This would be quite satisfactory for the college admissions officer in that there is little variation left to be explained, although this result is quite unlikely in actual research in education and the behavioral sciences. Obviously, there is a large effect size here.

Let us compute the second test statistic for the GGPA example. We specify the null hypothesis to be $\beta_k = 0$ (i.e., the slope is zero) and conduct two-tailed tests. First the variance error of estimate is

$$s_{res}^2 = \frac{SS_{res}}{df_{res}} = \frac{.1002}{8} = .0125$$

The standard error of estimate, s_{res} , is .1118. Next, the standard errors of the b_k are found to be

$$s(b_1) = \frac{s_{res}}{\sqrt{(n-1)s_1^2(1-r_{12}^2)}} = \frac{.1118}{\sqrt{(10)266.8182(1-.3011^2)}} = .0023$$

$$s(b_2) = \frac{s_{res}}{\sqrt{(n-1)s_2^2(1-r_{12}^2)}} = \frac{.1118}{\sqrt{(10)0.1609(1-.3011^2)}} = .0924$$

Finally, we find the t test statistics to be computed as follows:

$$t_1 = b_1 / s(b_1) = .0125 / .0023 = 5.4348$$

$$t_2 = b_2 / s(b_2) = .4687 / .0924 = 5.0725$$

To evaluate the null hypotheses, we compare these test statistics to the critical values of $\pm .025 t_8 = \pm 2.306$. Both test statistics exceed the critical value; consequently, H_0 is rejected in favor of H_1 for both predictors. We conclude that both partial slopes are indeed statistically significantly different from zero at the .05 level of significance.

Finally, let us compute the confidence intervals for the b_k 's as follows:

$$\begin{aligned} CI(b_1) &= b_1 \pm {}_{(\alpha/2)} t_{(n-m-1)} s(b_1) = b_1 \pm .025 t_8 s(b_1) \\ &= .0125 \pm 2.306 (.0023) = (.0072, .0178) \end{aligned}$$

$$\begin{aligned} CI(b_2) &= b_2 \pm {}_{(\alpha/2)} t_{(n-m-1)} s(b_2) = b_2 \pm .025 t_8 s(b_2) \\ &= .4687 \pm 2.306 (.0924) = (.2556, .6818) \end{aligned}$$

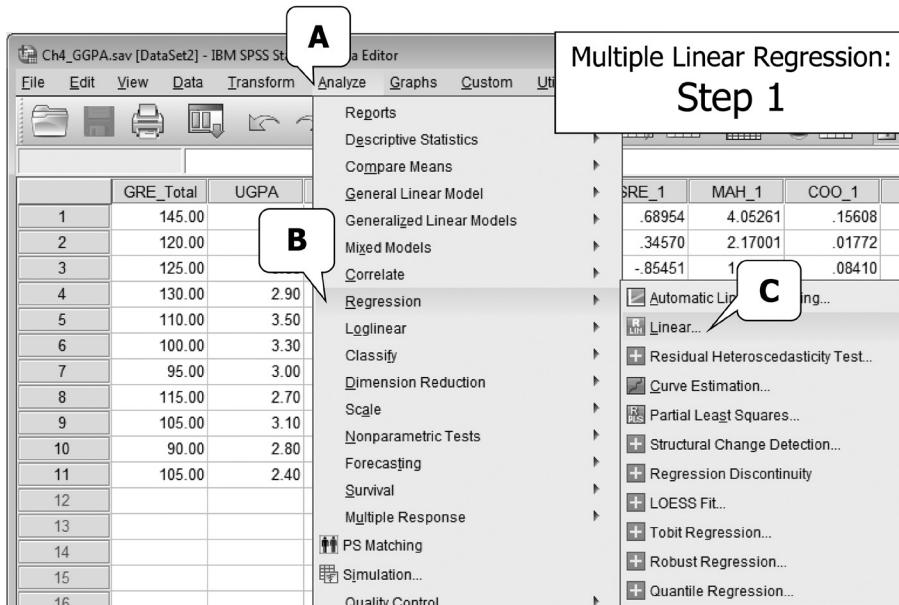
The intervals do not contain zero, the value specified in H_0 ; thus, we again conclude that both b_k 's are significantly different from zero at the .05 level of significance.

4.3 COMPUTING MULTIPLE LINEAR REGRESSION USING SPSS

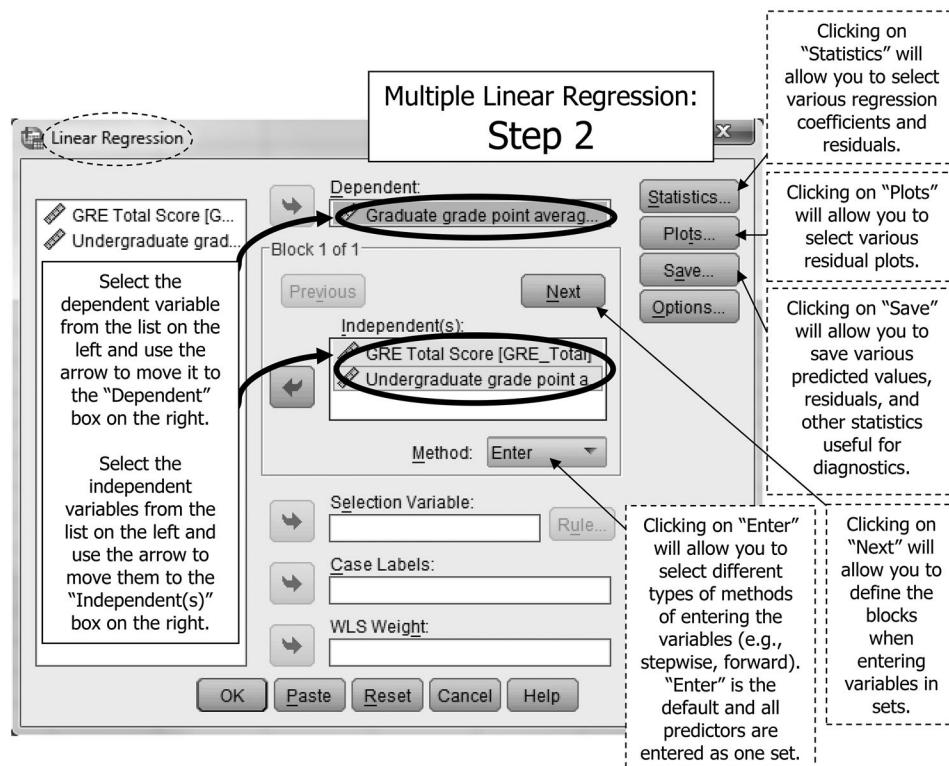
Next, we consider SPSS for the multiple linear regression model. Before we conduct the analysis, let us review the data. With one dependent variable and two independent variables, the dataset must consist of three variables or columns, one for each independent variable and one for the dependent variable. Each row still represents one individual, indicating the value of the independent variables for that particular case and their score on the dependent variable. As seen in the screenshot below, for a multiple linear regression analysis therefore, the SPSS data is in the form of three columns that represent the two independent variables (GRE total score and undergraduate GPA) and one dependent variable (graduate grade point average).

	GRE_Total	UGPA	GGPA
1	145.00	3.20	4.00
2	120.00	3.70	3.90
3	125.00	3.60	3.80
4	130.00	2.90	3.70
5	110.00	3.50	3.60
6	100.00	3.30	3.50
7	95.00	3.00	3.40
8	115.00	2.70	3.30
9	105.00	3.10	3.20
10	90.00	2.80	3.10
11	105.00	2.40	3.00

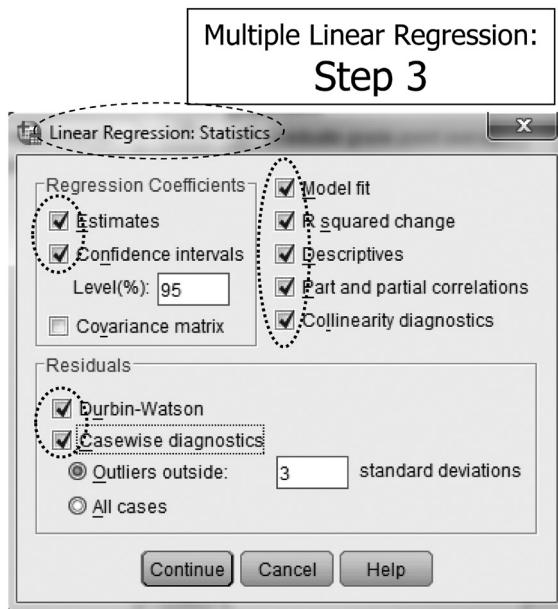
Step 1. To conduct a simple linear regression, go to “Analyze” in the top pull-down menu, then select “Regression,” and then select “Linear.” Following the screenshot below (Step 1) produces the “Linear Regression” dialog box.



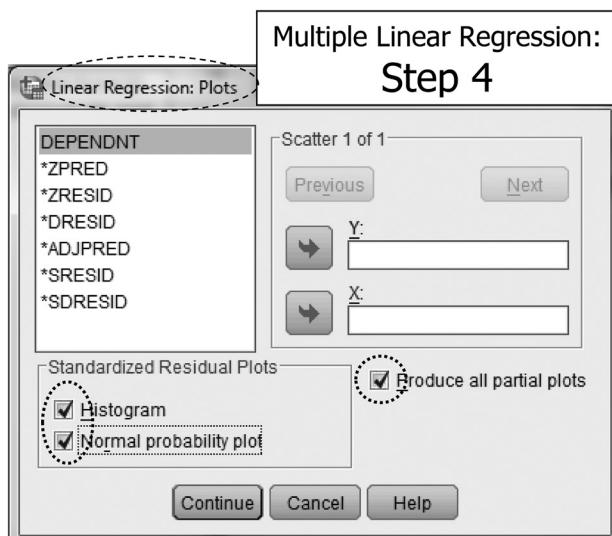
Step 2. Click the dependent variable (e.g., ‘GGPA’) and move it into the “Dependent” box by clicking the arrow button. Click the independent variables and move them into the “Independent(s)” box by clicking the arrow button (see screenshot Step 2).



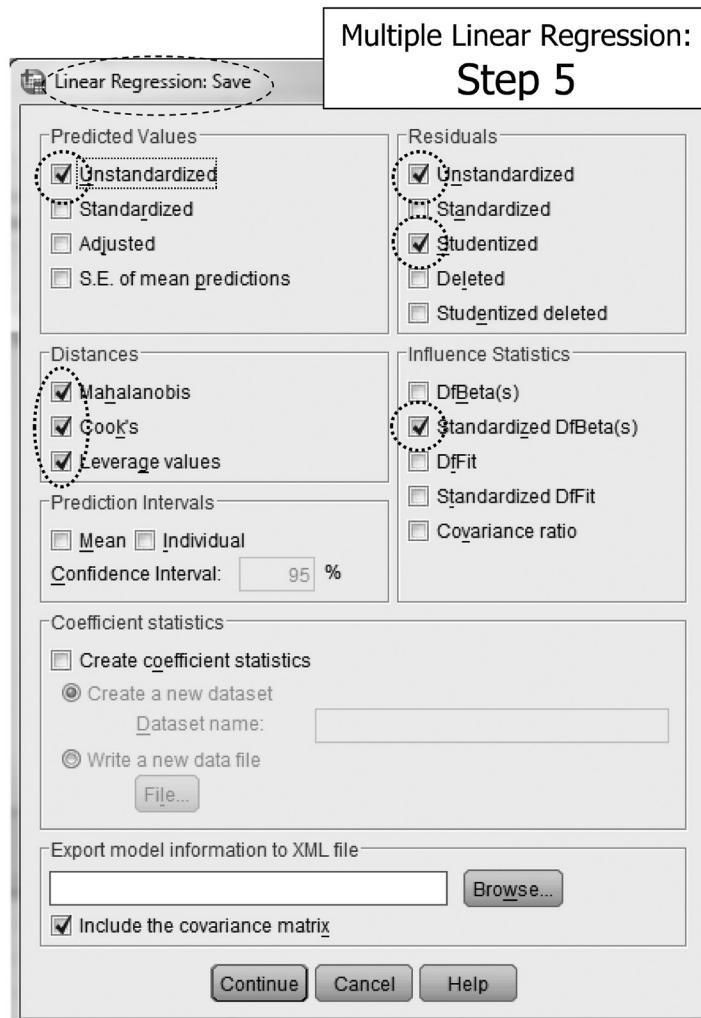
Step 3. From the Linear Regression dialog box (see screenshot Step 2), clicking on “Statistics” will provide the option to select various regression coefficients and residuals. From the Statistics dialog box (see screenshot Step 3), place a checkmark in the box next to the following: (1) estimates, (2) confidence intervals, (3) model fit, (4) R squared change, (5) descriptives, (6) part and partial correlations, (7) collinearity diagnostics, (8) Durbin-Watson, and (9) Casewise diagnostics. For this example, we apply an alpha level of .05, thus we will leave the default confidence interval percentage at 95. If we were using a different alpha, the confidence interval would be the complement of alpha (e.g., $\alpha = .01$ then $CI = 1 - .01 = 99$). We will also leave the default of ‘3 standard deviations’ for defining outliers for the casewise diagnostics. Click on “Continue” to return to the original dialog box.



Step 4. From the Linear Regression dialog box (see screenshot Step 2), clicking on "Plots" will provide the option to select various residual plots. From the Plots dialog box, place a checkmark in the box next to the following: (1) histogram, (2) normal probability plot, and (3) produce all partial plots. Click on "Continue" to return to the original dialog box.



Step 5. From the Linear Regression dialog box (see screenshot Step 2), clicking on "Save" will provide the option to save various predicted values, residuals, and statistics that can be used for diagnostic examination. From the Save dialog box under the heading of **Predicted Values**, place a checkmark in the box next to the following: (1) unstandardized.



Under the heading of **Residuals**, place a checkmark in the box next to the following: (1) unstandardized and (2) studentized. Under the heading of **Distances**, place a checkmark in the box next to the following: (1) Mahalanobis, (2) Cook's, and (3) Leverage values. Under the heading of **Influence Statistics**, place a checkmark in the box next to the following: (1) Standardized DFBETA(s). Click on "Continue" to return to the original dialog box. From the "Linear Regression" dialog box, click on "OK" to return and generate the output.

Interpreting the output. Annotated results are shown in Table 4.3.

■ TABLE 4.3

SPSS Results for the Multiple Regression GRE-GPA Example

Descriptive Statistics			
	Mean	Std. Deviation	N
Graduate grade point average	3.5000	.33166	11
GRE Total Score	112.7273	16.33457	11
Undergraduate grade point average	3.1091	.40113	11

The table labeled "Descriptive Statistics" provides basic descriptive statistics (means, standard deviations, and sample sizes) for the independent and dependent variables.

The table labeled "Correlations" provides the Pearson correlation coefficient values, *p* values, and sample size for the simple bivariate Pearson correlation between the independent and dependent variables.

The correlation between graduate GPA and GRE-Total (*p* = .002) and the correlation between graduate GPA and undergraduate GPA (*p* = .004) are statistically significant.

		Graduate grade point average	GRE Total Score	Undergraduate grade point average
Pearson Correlation	Graduate grade point average	1.000	.784	.752
	GRE Total Score	.784	1.000	.301
	Undergraduate grade point average	.752	.301	1.000
Sig. (1-tailed)	Graduate grade point average		.002	.004
	GRE Total Score	.002		.184
	Undergraduate grade point average	.004	.184	
N	Graduate grade point average	11	11	11
	GRE Total Score	11	11	11
	Undergraduate grade point average	11	11	11

TABLE 4.3 (continued)

SPSS Results for the Multiple Regression GRE-GPA Example

Variables Entered/Removed ^b			
Model	Variables Entered	Variables Removed	Method
1	Undergraduate grade point average, GRE Total Score		Enter

"Variables Entered/Removed" lists the independent variables included in the model and the method they were entered (i.e., Enter).

a. All requested variables entered.
b. Dependent Variable: Graduate grade point average

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	Change Statistics				Sig. F Change	Durbin-Watson
					R Square Change	F Change	df1	df2		
1	.953 ^a	.908	.885	.11272	.908	39.291	2	8	.000	2.116

a. Predictors: (Constant), Undergraduate grade point average, GRE Total Score
b. Dependent Variable: Graduate grade point average

R is the multiple correlation coefficient.

*R*² is the squared multiple correlation coefficient (aka, coefficient of determination). It represents the proportion of variance in the dependent variable that is explained by the independent variables.

Adjusted *R*² is interpreted as the percentage of variation in the dependent variable that is explained after adjusting for sample size and the number of predictors.

ANOVA ^b					
Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression	.998	2	.499	39.291
	Residual	.102	8	.013	
	Total	1.100	10		

The *p* value (.000) indicates we reject the null hypothesis. The probability of finding a sample value of multiple *R*² of .908 or larger when the true population multiple correlation coefficient is zero is less than 1%.

Total *SS* is partitioned into *SS* regression and *SS* residual. Regression sum of squares indicates variability explained by the regression model. Residual sum of squares indicates variability *not* explained by the regression model.

The *F* statistic tests the overall regression model (i.e., that the population multiple correlation coefficient is zero).

a. Predictors: (Constant), Undergraduate grade point average, GRE Total Score
b. Dependent Variable: Graduate grade point average

TABLE 4.3 (continued)

SPSS Results for the Multiple Regression GRE-GPA Example

The 'constant' is the intercept and the unstandardized coefficient tells us that if the predictors were zero, graduate GPA (the dependent variable) would be .638. The 'GRE-Total' and 'UGPA' are the slopes. For every one point increase in GRE-Total, graduate GPA will increase by about 1/10 of one point (holding constant undergraduate GPA). For every one point increase in undergraduate GPA, graduate GPA will increase by about 1/2 of one point (holding constant GRE-Total).

The test statistic, t , is calculated as the unstandardized coefficient divided by its standard error. Thus the slope for undergraduate GPA is calculated as (difference due to rounding):

$$t = \frac{.469}{.093} = 5.043$$

The p value for the intercept ('the constant') ($p = .087$) indicates that the intercept is *not* statistically significantly different from zero (this finding is usually of less interest than the slopes). The p values for GRE-Total and undergraduate GPA (the independent variables) ($p = .001$) indicate that the slopes are statistically significantly different from zero.

Model	Coefficients ^a										Collinearity Statistics	
	Unstandardized Coefficients		Standardized Coefficients		t	Sig.	95.0% Confidence Interval for B		Correlations			
	B	Std. Error	Beta				Lower Bound	Upper Bound	Zero-order	Partial	Part	Tolerance
1 (Constant)	.638	.327			1.954	.087	-.115	1.391				
GRE Total Score	.012	.002	.614	5.447	.001		.007	.018	.784	.887	.585	.909
UGPA	.469	.093	.567	5.030	.001		.254	.684	.752	.872	.541	.909
												1.100

a. Dependent Variable: Graduate grade point average

Zero-order correlations are the simple bivariate Pearson correlations between the dependent variable and the independent variables.

The partial correlation of .887 is the correlation between GRE-Total and graduate GPA (dependent variable) when the linear effect of undergraduate GPA has been removed from both GRE-Total and graduate GPA. Squaring this indicates that 78.7% of the variation in graduate GPA that is not explained by undergraduate GPA is explained by GRE-Total.

The part correlation of .585, when squared (i.e., .342) indicates that GRE-Total explains an additional 34% of the variance in graduate GPA over and above the variance in graduate GPA which is explained by undergraduate GPA.

Cloud bubble: Collinearity statistics to be reviewed under assumptions.

'Collinearity diagnostics' will be examined in our discussion of assumptions.

Collinearity Diagnostics^a

Model	Dimension	Eigenvalue	Condition Index	Variance Proportions		
				(Constant)	GRE Total Score	Undergraduate grade point average
1	1	2.981	1.000	.00	.00	.00
	2	.012	15.727	.03	.86	.40
	3	.007	20.537	.97	.13	.60

a. Dependent Variable: Graduate grade point average

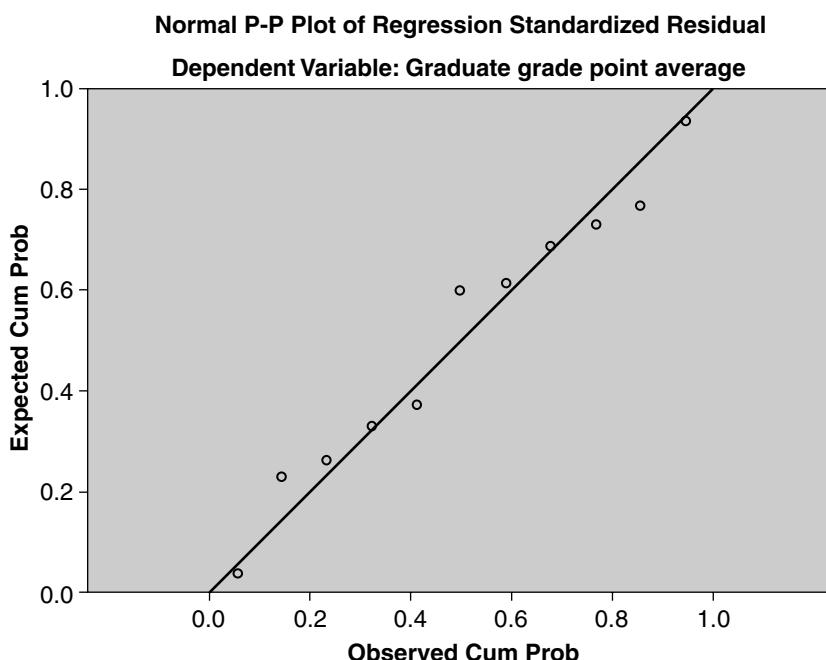
TABLE 4.3 (continued)

SPSS Results for the Multiple Regression GRE-GPA Example

'Residual statistics' and related graphs (histogram and Q-Q plot of standardized residuals, not presented here) will be examined in our discussion of assumptions.

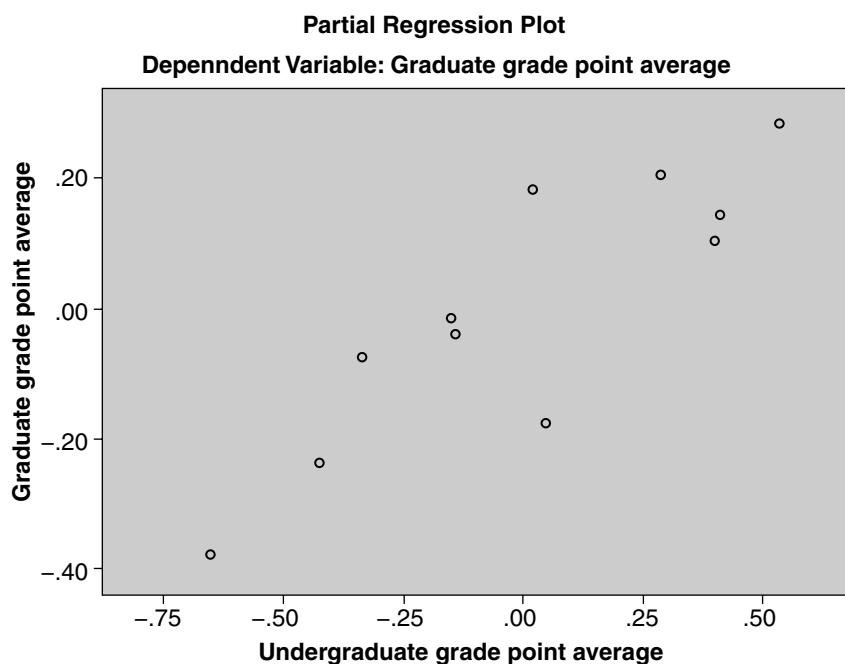
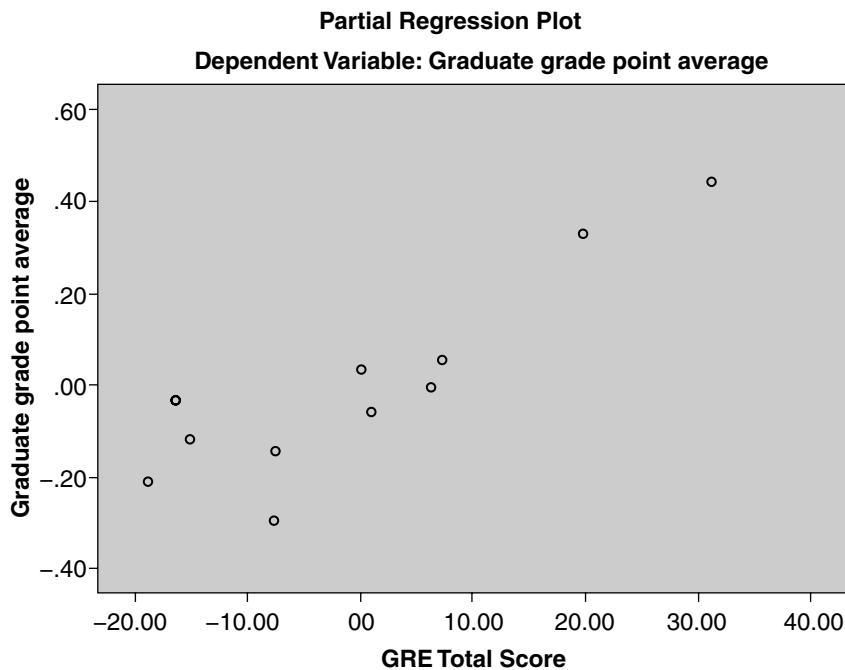
	Residuals Statistics ^a				
	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	3.0714	3.9448	3.5000	.31597	11
Std. Predicted Value	-1.357	1.408	.000	1.000	11
Standard Error of Predicted Value	.038	.079	.058	.011	11
Adjusted Predicted Value	3.0599	3.9117	3.4954	.30917	11
Residual	-.19943	.17207	.00000	.10082	11
Std. Residual	-1.769	1.527	.000	.894	11
Stud. Residual	-1.881	1.716	.017	1.008	11
Deleted Residual	-.22531	.21754	.00458	.12935	11
Stud. Deleted Residual	-2.355	2.020	.000	1.145	11
Mahal. Distance	.240	4.053	1.818	1.048	11
Cook's Distance	.012	.260	.092	.081	11
Centered Leverage Value	.024	.405	.182	.105	11

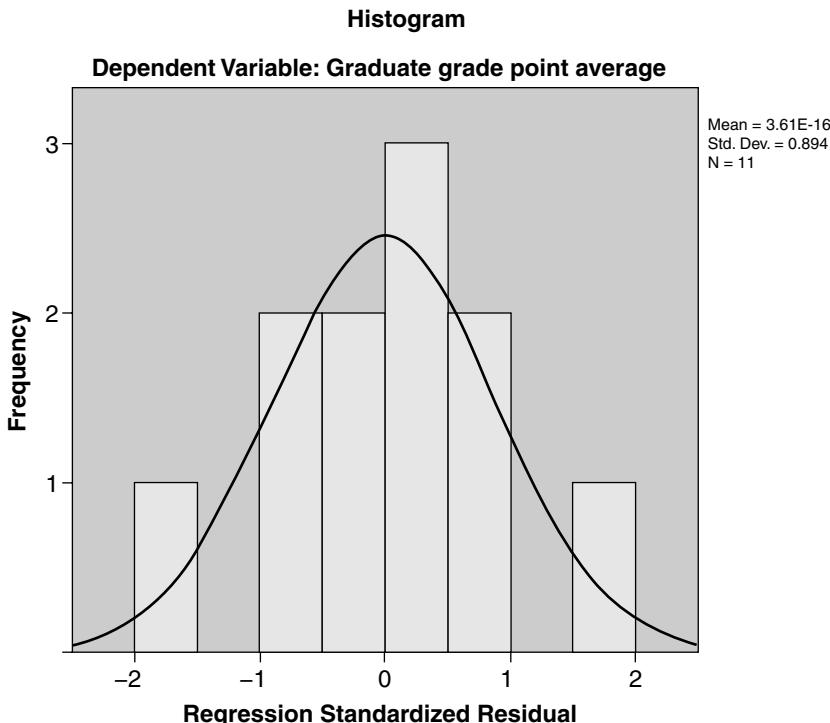
a. Dependent Variable: Graduate grade point average



■ TABLE 4.3 (continued)

SPSS Results for the Multiple Regression GRE-GPA Example





4.4 DATA SCREENING

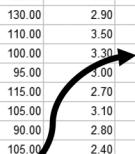
As you may recall, there were a number of assumptions associated with multiple linear regression. These included (a) independence, (b) homoscedasticity, (c) linearity, (d) normality, and (e) noncollinearity. Although fixed values of X were discussed in assumptions, this is not an assumption that will be tested, but is instead related to the use of the results (i.e., extrapolation and interpolation).

Before we begin to examine assumptions, let us review the values that we requested to be saved to our dataset (see dataset screenshot that follows).

1. **PRE_1** are the unstandardized predicted values (i.e., \hat{Y}).
2. **RES_1** are the unstandardized residuals, simply the difference between the observed and predicted values. For student 1, for example, the observed value for the graduate GPA (i.e., the dependent variable) was 4 and the predicted value was 3.94483. Thus, the unstandardized residual is simply $4 - 3.94483$ or .05517.
3. **SRE_1** are the studentized residuals, a type of standardized residual that is more sensitive to outliers as compared to standardized residuals. Studentized residuals are computed as the unstandardized residual divided by an estimate of the standard deviation with that case removed. As a guideline, studentized residuals with an absolute value greater than 3 are considered outliers (Stevens, 1984).
4. **MAH_1** are Mahalanobis distance values, which measure how far that particular case is from the average of the independent variable and thus can be helpful

in detecting outliers. These values can be reviewed to determine cases that are exerting leverage. Barnett and Lewis (1994) produced a table of critical values for evaluating Mahalanobis distance. Squared Mahalanobis distances divided by the number of variables (D^2/df) which are greater than 2.5 (for small samples) or 3 to 4 (for large samples) are suggestive of outliers (Hair et al., 2006). Later, we follow another convention for examining these values using the chi-square distribution.

5. **COO_1** are Cook's distance values and provide an indication of influence of individual cases. As a rule of thumb, Cook's values greater than one suggest that case is potentially problematic.
6. **LEV_1** are leverage values, a measure of distance from a respective case to the average of the predictor.
7. **SDB0_1**, **SDB1_1** and **SDB2_1** are standardized DFBETA values for the intercept and slopes, respectively, and are easier to interpret as compared to their unstandardized counterparts. Standardized DFBETA values greater than an absolute value of two suggest that the case may be exerting undue influence on the calculation of the parameters in the model (i.e., the slopes and intercept).



	GRE_Total	UGPA	GGPA	PRE_1	RES_1	SRE_1	MAH_1	COO_1	LEV_1	SDB0_1	SDB1_1	SDB2_1
1	145.00	3.20	4.00	3.94483	.05517	68954	4.05261	.15608	.40526	-.33730	.59269	-.11447
2	120.00	3.70	3.90	3.86758	.03242	34570	2.17001	.01772	.21700	-.14218	.00022	.17391
3	125.00	3.60	3.80	3.88303	-.08303	-85451	1.65890	.08410	.16589	.35251	-.12349	-.32176
4	130.00	2.90	3.70	3.61728	.08272	86503	1.89338	.09712	.18934	-.03978	.40341	-.27892
5	110.00	3.50	3.60	3.64922	-.04922	-49101	1.18272	.02126	.11827	.07842	.08006	-.17822
6	100.00	3.30	3.50	3.43085	.06915	68899	1.16223	.04134	.11622	.04775	-.22828	.17583
7	95.00	3.00	3.40	3.22793	.17207	1.71646	1.18109	.25952	.11811	.63170	-.75577	.04127
8	115.00	2.70	3.30	3.33660	-.03660	-36688	1.25898	.01242	.12590	-.08873	-.05787	.13770
9	105.00	3.10	3.20	3.39943	-.19943	-1.88064	2.3956	.15299	.02396	.29610	.38705	-.09942
10	90.00	2.80	3.10	3.07188	.02812	29777	2.07186	.01255	.20719	.14662	-.12853	.03898
11	105.00	2.40	3.00	3.07136	-.07136	-81994	3.12865	.15177	.31287	-.51278	-.02036	.55938

As we look at the raw data, we see nine new variables have been added to our dataset. These are our predicted values, residuals, and other diagnostic statistics. The residuals will be used for diagnostics to review the extent to which our data meet the assumptions of multiple linear regression.

4.4.1 Independence

Here we will plot (1) studentized residuals (which were requested and created through the 'Save' option when generating our model) against unstandardized predicted values and (2) studentized residuals against each independent variable to examine the extent to which independence was met. If you need a refresher on generating this plot, please review the data screening chapter. If the assumption of independence is met, the points should fall randomly within a band of -2.0 to +2.0.

4.4.2 Homoscedasticity

Homoscedasticity is evident when a plot of residuals appears fairly constant over the range of unstandardized predicted values (i.e., a random display of points) and observed values of the independent variables. If the display of residuals increases

or decreases across the plot, then there may be an indication that the assumption of homoscedasticity has been violated. The plot used to examine independence can also be used for homoscedasticity.

4.4.3 Linearity

Since we have more than one independent variable, we have to take a different approach to examining linearity than what was done with simple linear regression. However, we can use the same information gleaned from our examination of independence and homoscedasticity for reviewing the assumption of linearity. We can also review the partial regression plots that we asked for when generating the regression model. A separate partial regression plot is provided for each independent variable, where we are looking for linearity (rather than some type of polynomial). Even with a small sample size, the partial regression plots suggest evidence of linearity.

4.4.4 Normality

Normality can be understood by examining residuals as well as various diagnostics to examine our data for influential cases. Let us begin by examining the unstandardized residuals for normality. Just as we saw with simple linear regression, for multiple linear regression the distributional shape of the unstandardized residuals should be normal. A refresher on the steps for generating normality evidence was presented in the data screening chapter, and thus they will not be repeated here.

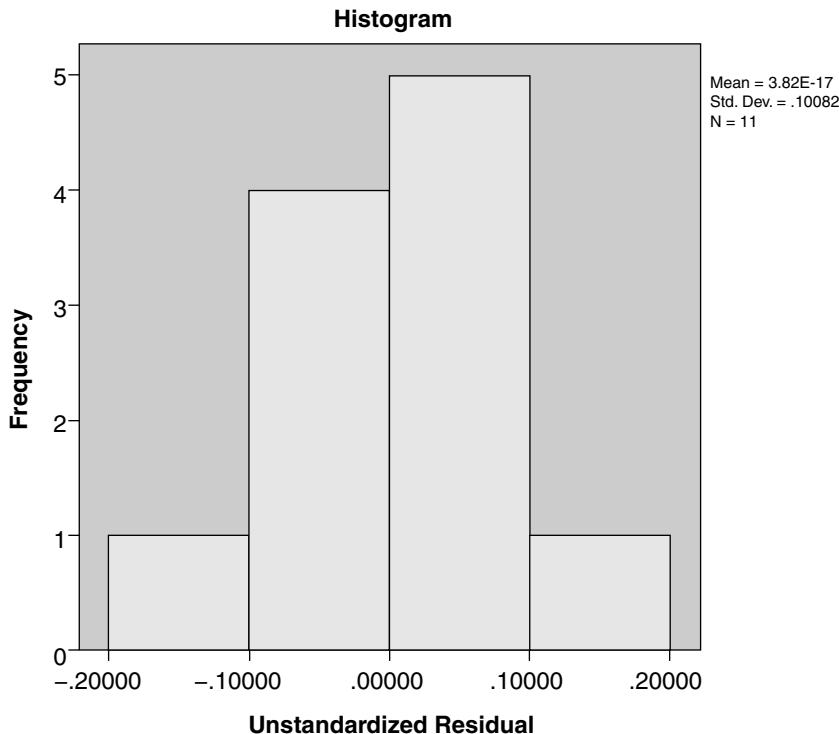
4.4.4.1 Interpreting Normality Evidence

By this point, we are well versed in interpreting quite a range of normality statistics and will do the same for multiple linear regression.

Descriptives

		Statistic	Std. Error
Unstandardized Residual	Mean	.0000000	.03039717
	95% Confidence Interval for Mean	Lower Bound Upper Bound	-.0677291 .0677291
	5% Trimmed Mean		.0015202
	Median		.0281190
	Variance		.010
	Std. Deviation		.10081601
	Minimum		-.19943
	Maximum		.17207
	Range		.37150
	Interquartile Range		.14051
	Skewness		-.336
	Kurtosis		.484

The skewness statistic of the residuals is $-.336$ and kurtosis is $.484$ —both being within the range of what would be considered normal (approximately an absolute value of 2.0 for skew and 7.0 for kurtosis), suggesting some evidence of normality. Given the very small sample size, the histogram below reflects as normal a distribution as might be expected.



There are a few other statistics that can be used to gauge normality. The results for the formal test of normality, the Shapiro-Wilk test (*SW*) (Shapiro & Wilk, 1965), is presented below and suggests that our sample distribution for the residual is *not* statistically significantly different than what would be expected from a normal distribution, as the *p* value is greater than α ($p = .918$).

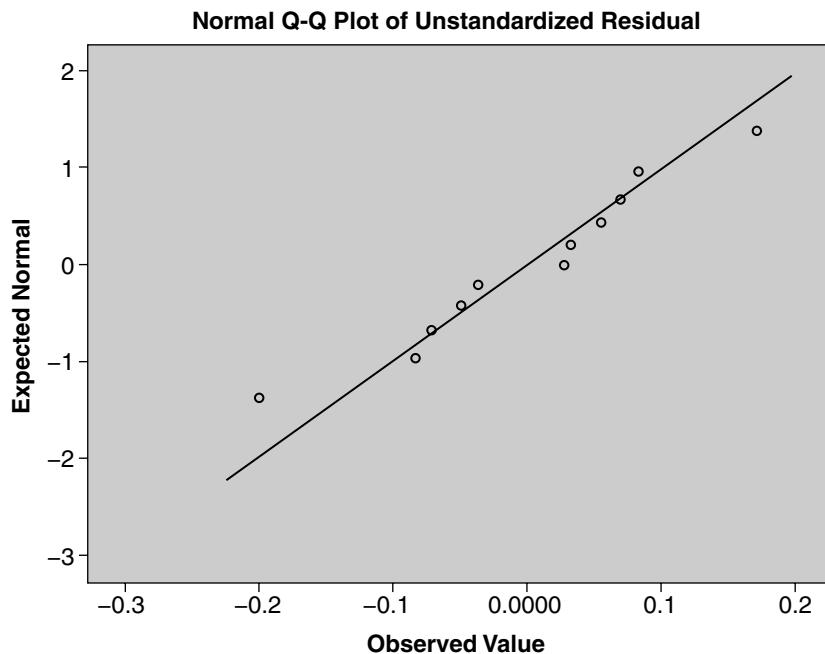
Tests of Normality

Unstandardized Residual	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
	.155	11	.200*	.973	11	.918

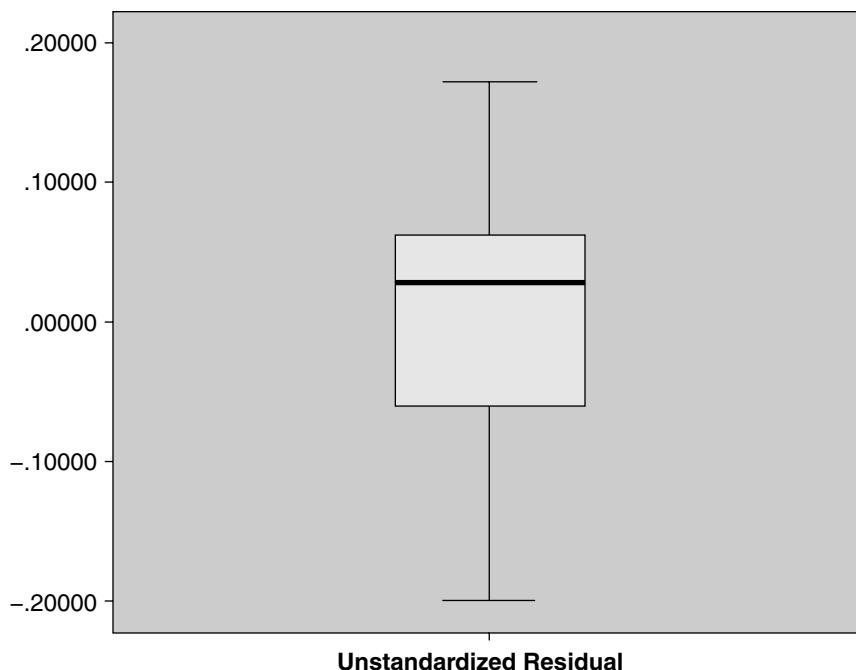
a. Lilliefors Significance Correction

*. This is a lower bound of the true significance.

Quantile-quantile (Q-Q) plots are also often examined to determine evidence of normality. The Q-Q plot of residuals suggests relative normality with points that fall on or close to the diagonal line, suggesting evidence of normality.



The boxplot below also suggests a relatively normal distribution of residuals with no outliers.



Considering the forms of evidence we have examined, skewness and kurtosis statistics, the Shapiro-Wilk test, histogram, the Q-Q plot, and the boxplot, all suggest normality is a reasonable assumption.

4.4.4.2 Screening Data for Influential Points

Casewise Diagnostics

Recall that we requested a number of statistics to help in diagnostics. One that we requested was for ‘Casewise diagnostics.’ If we had any cases with large values for the standardized residual (outside three standard deviations), information would have been included in our output to indicate the case number, value of the standardized residual, predicted value, and unstandardized residual. This information can be used to more closely examine case(s) with the extreme values on the standardized residuals.

Cook’s Distance

Cook’s distance provides an overall measure for the influence of individual cases. Values greater than one suggest that the case may be problematic in terms of undue influence on the model. Examining the residual statistics in our output (see following table), we see that the maximum value for Cook’s distance is .260, well under the point at which we should be concerned.

Residuals Statistics^a

	Minimum	Maximum	Mean	Std. Deviation	N
Predicted Value	3.0714	3.9448	3.5000	.31597	11
Std. Predicted Value	-1.357	1.408	.000	1.000	11
Standard Error of Predicted Value	.038	.079	.058	.011	11
Adjusted Predicted Value	3.0599	3.9117	3.4954	.30917	11
Residual	-.19943	.17207	.00000	.10082	11
Std. Residual	-1.769	1.527	.000	.894	11
Stud. Residual	-1.881	1.716	.017	1.008	11
Deleted Residual	-.22531	.21754	.00458	.12935	11
Stud. Deleted Residual	-2.355	2.020	.000	1.145	11
Mahal. Distance	.240	4.053	1.818	1.048	11
Cook’s Distance	.012	.260	.092	.081	11
Centered Leverage Value	.024	.405	.182	.105	11

a. Dependent Variable: Graduate grade point average

Mahalanobis Distances

Mahalanobis distances are measures of the distance from each case to the mean of the independent variable for the remaining cases. We can use the value of Mahalanobis distance as a test statistic value with the chi-square distribution. With two independent variables and one dependent variable, we have three degrees of freedom. Given an alpha level of .05 (alpha of .001 if you want to be a bit more liberal), the chi-square critical value is 7.82. Thus, any Mahalanobis distance greater than 7.82 suggests that case is an outlier. With a maximum of 4.053 (see previous table), there is no evidence to suggest there are outliers in our data.

Centered Leverage Values

Centered leverage values less than .20 suggest there are no problems with cases that are exerting undue influence. Values greater than .5 indicate problems.

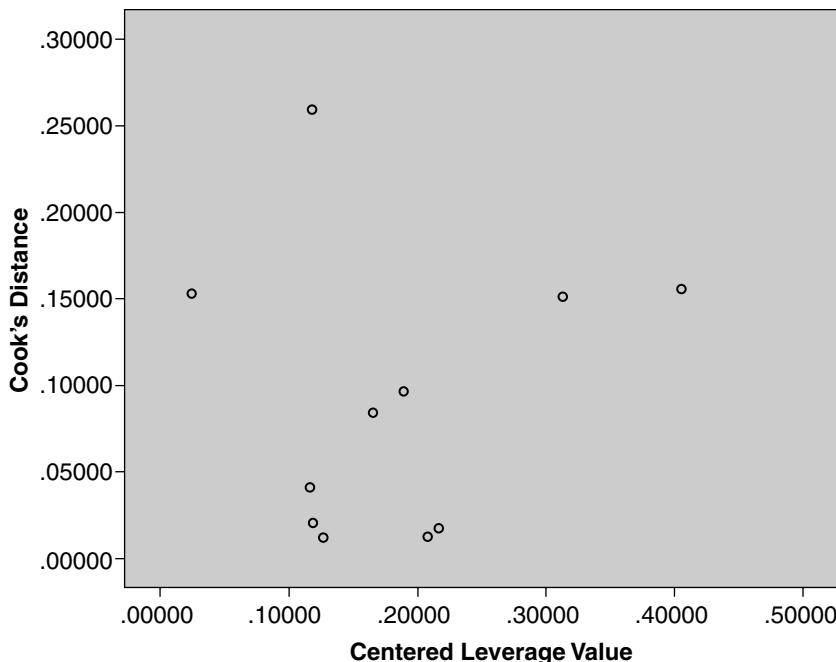
DFBETA

We also asked to save DFBETA values. These values provide another indication of the influence of cases. DFBETA provides information on the change in the predicted value when the case is deleted from the model. For standardized DFBETA values, values greater than an absolute value of 2.0 should be examined more closely. Looking at the minimum and maximum DFBETA values, there are no cases suggestive of undue influence.

Descriptive Statistics			
	N	Minimum	Maximum
Standardized DFBETA Intercept	11	-.51278	.63170
Standardized DFBETA GRE_Total	11	-.75577	.59269
Standardized DFBETA UGPA	11	-.32176	.55938
Valid N (listwise)	11		

Diagnostic Plots

There are a number of diagnostic plots that can be generated from the values we saved. For example, a plot of Cook's distance against centered leverage values provides a way to identify influential cases (i.e., cases with leverage of .50 or above and Cook's distance of 1.0 or greater). Here there are no cases that suggest undue influence.



4.4.5 Noncollinearity

Detecting multicollinearity can be done by reviewing the VIF and tolerance statistics. From the table below, we see tolerance and VIF values. Tolerance is calculated as

$(1 - R^2)$ and values close to zero (a suggested cut point is .10 or less) suggest potential multicollinearity problems. Why? A tolerance of .10 suggests that 90% (or more) of the variance in one of the independent variables can be explained by another independent variable. VIF is the ‘variance inflation factor’ and is the reciprocal of tolerance, where

$$VIF = \frac{1}{\text{tolerance}}. \text{ VIF values greater than 10 (which correspond to a tolerance of .10)}$$

suggest potential multicollinearity.

Collinearity Statistics	
Tolerance	VIF
.909	1.100
.909	1.100

Collinearity diagnostics (see SPSS output below) can also be reviewed. Multiple eigenvalues close to zero indicate independent variables that have strong intercorrelations. The condition index is calculated as the square root of the ratio of the largest

eigenvalue to each preceding eigenvalue (e.g., $\sqrt{\frac{2.981}{.012}} = 15.76$). A general recommendation for interpreting condition indices is that values in the range of 10 to 30

should be of concern, greater than 30 indicates trouble, and greater than 100 indicates disaster (Belsley, 1991). In this case, both the eigenvalues and condition indices suggest possible problems with multicollinearity.

Collinearity Diagnostics^a

Model	Dimension	Eigenvalue	Condition Index	Variance Proportions		
				(Constant)	GRE Total Score	Undergraduate grade point average
1	1	2.981	15.727	.00	.00	.00
	2	.012		.03	.86	.40
	3	.007		.97	.13	.60

a. Dependent Variable: Graduate grade point average

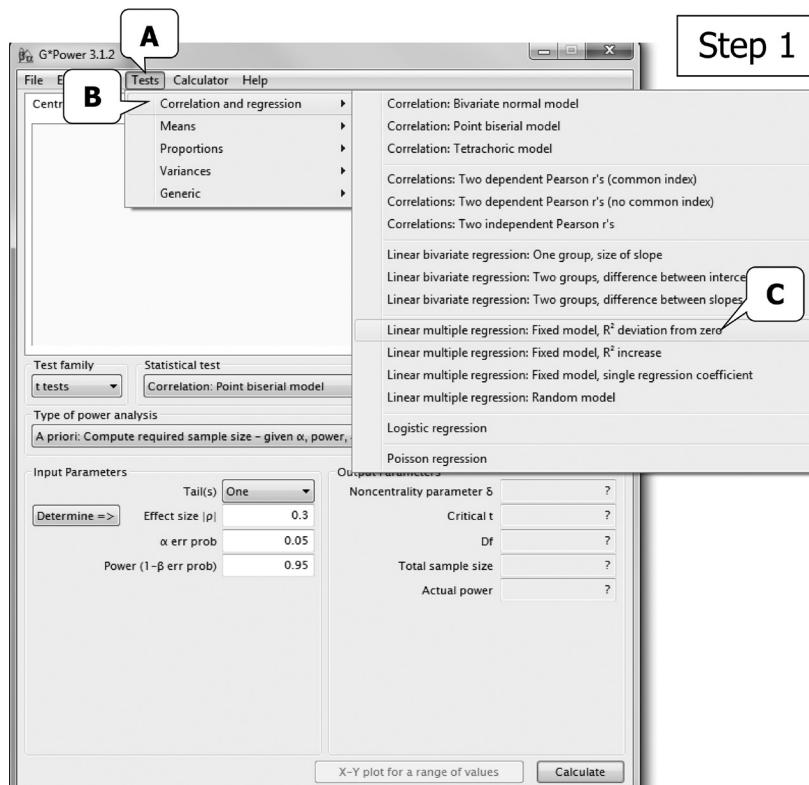
Noncollinearity can also be examined by computing regression models, where each independent variable is considered the outcome and is predicted by the remaining independent variables (the dependent variable is not included in these models). If any of the resultant R_k^2 values are close to one (greater than .9 is a good cut point), then there may be a collinearity problem. For the example data, $R_{12}^2 = .091$ and therefore collinearity is not a concern. Note that in multiple regression situations where there are two independent variables (as in this example with GRE-Total and undergraduate GPA), only one regression needs to be conducted to check for multicollinearity, as the results for regressing undergraduate GPA on GRE-Total are the same as regressing GRE-Total on undergraduate GPA.

4.5 POWER USING G*POWER

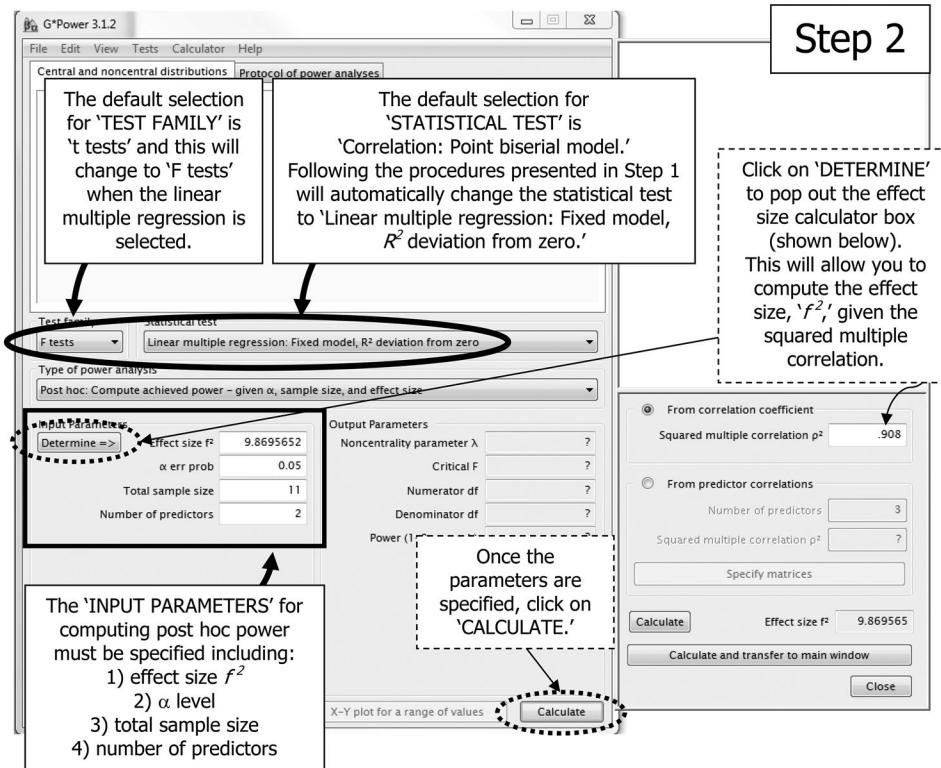
A priori and post hoc power can be determined using the specialized software described previously in this text (e.g., G*Power), or you can consult a priori power tables (e.g., Cohen, 1988). As an illustration, we use G*Power to compute the post hoc power of our test.

4.5.1 Post Hoc Power for Multiple Linear Regression Using G*Power

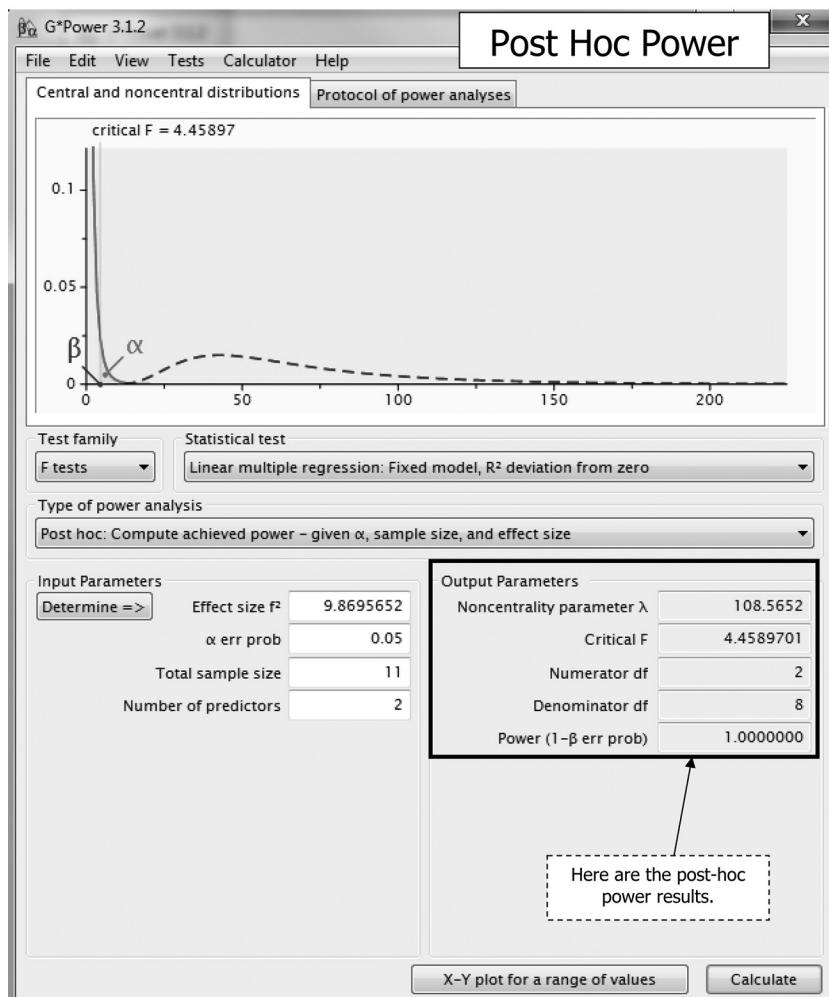
The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family. In our case, we conducted a multiple linear regression. To find regression, we select 'Tests' in the top pull-down menu, then 'Correlation and regression,' and then 'Linear multiple regression: Fixed model, R^2 deviation from zero.' This will allow us to determine power for the hypothesis that the overall multiple R^2 is equal to zero (i.e., power for the overall regression model). Once that selection is made, the 'Test family' automatically changes to 'F test.'



The 'Type of power analysis' desired needs to be selected. To compute post hoc power, select 'Post hoc: Compute achieved power—given α , sample size, and effect size.'



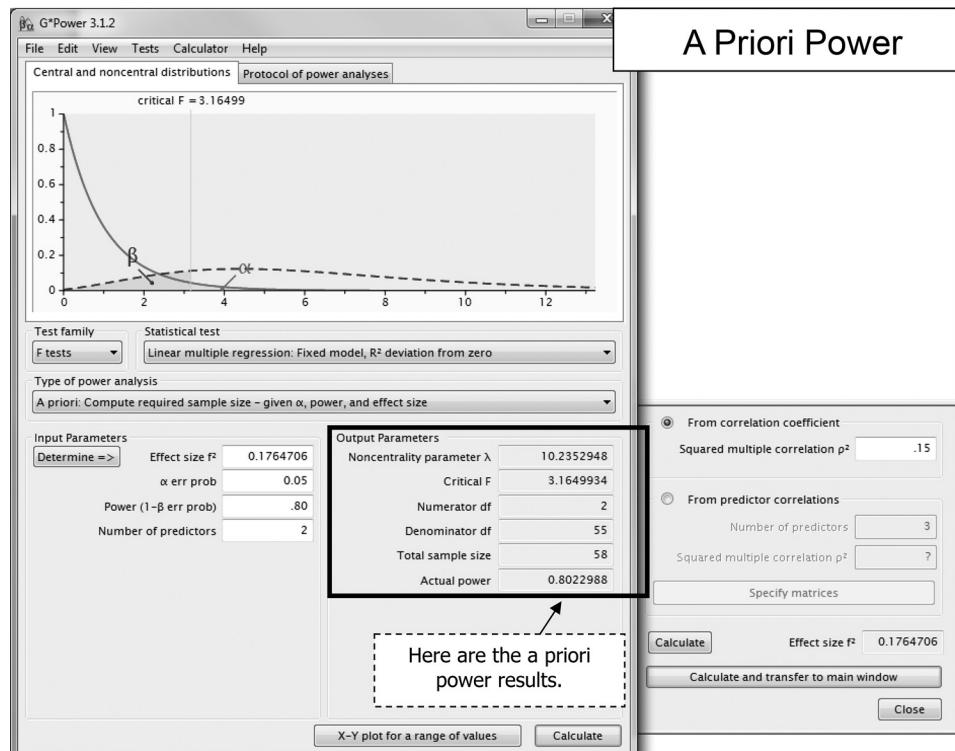
The 'INPUT PARAMETERS' must then be specified. We compute the effect size, f^2 , last and so we skip that for the moment. The α level we used was .05, the total sample size was 11, and there were two independent variables. Next we use the pop-out effect size calculator in G*Power to compute the effect size f^2 . To do this, click on 'DETERMINE,' which is displayed under 'INPUT PARAMETERS.' In the pop-out effect size calculator, input the value for the squared multiple correlation. Click on 'CALCULATE' to compute the effect size f^2 . Then click on 'CALCULATE AND TRANSFER TO MAIN WINDOW' to transfer the calculated effect size (i.e., 9.8695652) to the 'INPUT PARAMETERS.' Once the parameters are specified, click on 'CALCULATE' to find the power statistics.



The 'OUTPUT PARAMETERS' provide the relevant statistics given the input just specified. Here we were interested in determining post hoc power for a multiple linear regression with a computed effect size f^2 of 9.8695652, an alpha level of .05, total sample size of 11, and two predictors. Based on those criteria, the post hoc power for the overall multiple linear regression model was 1.0000. In other words, given the input parameters, the probability of rejecting the null hypothesis when it is really false (in this case, the probability that the multiple correlation coefficient is zero) was at the maximum (i.e., 1.00) (sufficient power is often .80 or above). Do not forget that conducting power analysis a priori is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters). Conducting power for change in R^2 and for the slopes can be conducted similarly by selecting the test family of 'Linear multiple regression: Fixed model, R^2 increase' or 'Linear multiple regression: Fixed model, single regression coefficient,' respectively.

4.5.2 A Priori Power for Multiple Linear Regression Using G*Power

For a priori power, we can determine the total sample size needed for multiple linear regression given the estimated effect size f^2 , α level, desired power, and number of predictors. We follow Cohen's (1988) conventions for effect size (i.e., small $R^2 = .02$; moderate $R^2 = .13$; large $R^2 = .26$). If we had estimated a moderate effect R^2 of .15, alpha of .05, observed power of .80, and two independent variables, we would need a total sample size of 58.



4.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example paragraph for the results of the multiple linear regression analysis. Recall that our graduate research assistant, Addie Venture, was assisting the Assistant Dean in Graduate Student Services, Dr. Golly. Dr. Golly wanted to know if graduate GPA could be predicted by the total score on the required graduate entrance exam (GRE-Total) and by undergraduate GPA. The research question presented to Dr. Golly by Addie included the following: Can graduate GPA be predicted from the GRE-Total and undergraduate GPA?

Addie then assisted Dr. Golly in generating a multiple linear regression as the test of inference, and a template for writing the research question for this design is presented below.

Can [dependent variable] be predicted from [list independent variables]?

It may be helpful to preface the results of the multiple linear regression with information on an examination of the extent to which the assumptions were met. The assumptions include (a) independence, (b) homoscedasticity, (c) normality, (d) linearity, (e) noncollinearity, and (f) values of X are fixed. Because the last assumption (fixed X) is based on interpretation, it will not be discussed here.

A multiple linear regression model was conducted to determine if graduate GPA (dependent variable) could be predicted from GRE-Total scores and undergraduate GPA (independent variables). The null hypotheses tested were that the multiple R^2 was equal to zero and that the regression coefficients (i.e., the slopes) were equal to zero. The data were screened for missingness and violation of assumptions prior to analysis. There was no missing data.

Linearity. Review of the partial scatterplot of the independent variables (GRE-Total and undergraduate GPA) and the dependent variable (midterm exam scores) indicate linearity is a reasonable assumption. Additionally, with a random display of points falling within an absolute value of two, a scatterplot of unstandardized residuals to predicted values provided further evidence of linearity.

Normality. The assumption of normality was tested via examination of the unstandardized residuals. Review of the Shapiro-Wilk test for normality ($SW = .973$, $df = 11$, $p = .918$) and skewness ($-.336$) and kurtosis ($.484$) statistics suggested that normality was a reasonable assumption. The boxplot suggested a relatively normal distributional shape (with no outliers) of the residuals. The Q-Q plot and histogram suggested normality was reasonable. Examination of casewise diagnostics, including Mahalanobis distance, Cook's distance, DFBETA values, and centered leverage values, suggested there were no cases exerting undue influence on the model.

Independence. A relatively random display of points in the scatterplots of studentized residuals against values of the independent variables and studentized residuals against predicted values provided evidence of independence. The Durbin-Watson statistic was computed to evaluate independence of errors and was 2.116, which is considered acceptable. This suggests that the assumption of independent errors has been met.

Homoscedasticity. A relatively random display of points, where the spread of residuals appears fairly constant over the range of values of the independent variables (in the scatterplots of studentized residuals against predicted values and studentized residuals against values of the independent variables) provided evidence of homoscedasticity.

Noncollinearity. Tolerance was greater than .10 (.909) and the variance inflation factor was less than 10 (1.100), suggesting that multicollinearity was not an issue. However, the eigenvalues for the predictors were close to zero (.012 and .007) and the respective condition indices were in the range of concern (between 10–30, 15.727, and 20.537, respectively). A review of GRE-Total regressed on undergraduate GPA, however, produced a multiple R squared of .091, which

suggests noncollinearity. Thus, though there is some isolated cause for concern, the evidence in aggregate suggests that multicollinearity is not an issue.

Here is an example paragraph of results for the multiple linear regression (remember that this will be prefaced by the previous paragraph reporting the extent to which the assumptions of the test were met).

The results of the multiple linear regression suggest that a significant proportion of the total variation in graduate GPA was predicted by GRE-Totals and undergraduate GPA, $F(2, 8) = 39.291, p < .001$. Additionally, we find:

- (a) For GRE-Totals, the unstandardized partial slope (.012) and standardized partial slope (.614) are statistically significantly different from zero ($t = 5.447, df = 8, p < .001$); with every one point increase in the GRE-Totals, graduate GPA will increase by approximately 1/100 of one point when controlling for undergraduate GPA.
- (b) For undergraduate GPA, the unstandardized partial slope (.469) and standardized partial slope (.567) are statistically significantly different from zero ($t = 5.030, df = 8, p < .001$); with every one point increase in undergraduate GPA, graduate GPA will increase by approximately ½ of one point when controlling for GRE-Totals.
- (c) The confidence interval around the unstandardized partial slopes do not include zero (GRE-Totals, .007, .018; undergraduate GPA, .254, .684) further confirming that these variables are statistically significant predictors of graduate GPA. Thus, GRETOT and UGPA were shown to be statistically significant predictors of GGPA, both individually and collectively.
- (d) The intercept (or average graduate GPA when GRE-Totals and undergraduate GPA is zero) was .638, not statistically significantly different from zero ($t = 1.954, df = 8, p = .087$).
- (e) Multiple R^2 indicates that approximately 91% of the variation in graduate GPA was predicted by GRE-Totals scores and undergraduate GPA. Interpreted according to Cohen (1988), this suggests a large effect.
- (f) Estimated power to predict multiple R^2 is at the maximum, 1.00.

We note that the more advanced regression models described in this chapter can all be conducted using SPSS. For further information on regression analysis with SPSS, see Morgan and Griego (1998), Weinberg and Abramowitz (2002), and Meyers, Gamst, and Guarino (2006).

PROBLEMS

Conceptual Problems

1. The correlation of salary and cumulative grade point average controlling for socio-economic status is an example of which one of the following?
 - a. Bivariate correlation
 - b. Partial correlation
 - c. Regression correlation
 - d. Semi-partial correlation

2. Variable 1 is to be predicted from a combination of variable 2 and one of variables 3, 4, 5, or 6. The correlations of importance are as follows:

$$r_{13} = .8 \quad r_{23} = .2$$

$$r_{14} = .6 \quad r_{24} = .5$$

$$r_{15} = .6 \quad r_{25} = .2$$

$$r_{16} = .8 \quad r_{26} = .5$$

Which of the following multiple correlation coefficients will have the largest value?

- a. $r_{1,23}$
- b. $r_{1,24}$
- c. $r_{1,25}$
- d. $r_{1,26}$

3. The most accurate predictions are made when the standard error of estimate equals which one of the following?

- a. \bar{Y}
- b. s_y
- c. 0
- d. 1

4. The intercept can take on a positive value only. True or false?

5. Adding an additional predictor to a regression equation will necessarily result in an increase in R^2 . True or false?

6. The best prediction in multiple regression analysis will result when each predictor has a high correlation with the other predictor variables and a high correlation with the dependent variable. True or false?

7. Consider the following two situations:

Situation 1 $r_{y1} = .6$ $r_{y2} = .5$ $r_{12} = .0$

Situation 2 $r_{y1} = .6$ $r_{y2} = .5$ $r_{12} = .2$

I assert that the value of R^2 will be greater in Situation 2. Am I correct?

8. Values of variables X_1, X_2, X_3 are available for a sample of 50 students. The value of $r_{12} = .6$. I assert that if the partial correlation $r_{12,3}$ were calculated it would be larger than .6. Am I correct?

9. A researcher is building a regression model. There is theory to suggest that science ability can be predicted by literacy skills when controlling for child characteristics (e.g., age and socioeconomic status). Which one of the following variable selection procedures is suggested?

- a. Backward elimination
- b. Forward selection
- c. Hierarchical regression
- d. Stepwise selection

10. I assert that the forward selection, backward elimination, and stepwise regression methods will always arrive at the same final model, given the same dataset and level of significance. Am I correct?

11. I assert the R^2_{adj} will always be larger for the model with the most predictors. Am I correct?
12. In a two-predictor regression model, if the correlation among the predictors is .95 and VIF is 20, then we should be concerned about collinearity. True or false?

Computational Problems

1. You are given the following data, where X_1 (hours of professional development) and X_2 (aptitude test scores) are used to predict Y (annual salary in thousands):

Y	X_1	X_2
40	100	10
50	200	20
50	300	10
70	400	30
65	500	20
65	600	20
80	700	30

Determine the following values: intercept; b^1 ; b^2 ; SS_{res} ; SS_{reg} ; F ; s^2_{res} ; $s(b^1)$; $s(b^2)$; t^1 ; t^2 .

2. You are given the following data, where X_1 (final percentage in science class) and X_2 (number of absences) are used to predict Y (standardized science test score in third grade):

Y	X_1	X_2
300	65	7
480	98	0
350	70	3
420	80	2
400	82	0
335	70	3
370	75	4
390	80	1
485	99	0
415	95	2
375	88	3

Determine the following values: intercept; b^1 ; b^2 ; SS_{res} ; SS_{reg} ; F ; s^2_{res} ; $s(b^1)$; $s(b^2)$; t^1 ; t^2 .

3. Complete the missing information for this regression model ($df = 23$).

Y'	=	25.1	+	1.2 X_1	+	1.0 X_2	-	.50 X_3	
(2.1)		(1.5)		(1.3)		(.06)		standard errors	
(11.9)		()		()		()		t ratios	
()		()		()		()		Significant at .05?	

4. Consider a sample of elementary school children. Given that $r(\text{strength, weight}) = .6$, $r(\text{strength, age}) = .7$, and $r(\text{weight, age}) = .8$, what is the first-order partial correlation coefficient between strength and weight holding age constant?
5. For a sample of 100 adults, you are given that $r_{12} = .55$, $r_{13} = .80$, and $r_{23} = .70$. What is the value of $r_{1(2-3)}$?
6. A researcher would like to predict salary from a set of four predictor variables for a sample of 45 subjects. Multiple linear regression analysis was utilized. Complete the following summary table ($\alpha = .05$) for the test of significance of the overall regression model:

Source	SS	df	MS	F	Critical Value and Decision
Regression	()	()	20	()	
Residual	400	()	()		
Total	()	()			

7. Calculate the partial correlation $r_{12,3}$ and the part correlation $r_{1(2-3)}$ from the following bivariate correlations: $r_{12} = .5$, $r_{13} = .8$, $r_{23} = .9$.
8. Calculate the partial correlation $r_{13,2}$ and the part correlation $r_{1(3-2)}$ from the following bivariate correlations: $r_{12} = .21$, $r_{13} = .40$, $r_{23} = -.38$.
9. You are given the following data, where X_1 (verbal aptitude) and X_2 (prior reading achievement) are to be used to predict Y (reading achievement):

Y	X_1	X_2
2	2	5
1	2	4
1	1	5
1	1	3
5	3	6
4	4	4
7	5	6
6	5	4
7	7	3
8	6	3
3	4	3
3	3	6
6	6	9
6	6	8
10	8	9
9	9	6
6	10	4
6	9	5
9	4	8
10	4	9

Determine the following values: intercept; b^1 ; b^2 ; SS_{res} ; SS_{reg} ; F ; s^2_{res} ; $s(b^1)$; $s(b^2)$; t^1 ; t^2 .

10. You are given the following data, where X_1 (years of teaching experience) and X_2 (salary in thousands) are to be used to predict Y (morale):

Y	X_1	X_2
125	1	24
130	2	30
145	3	32
115	2	28
170	6	40
180	7	38
165	5	48
150	4	42
195	9	56
180	10	52
120	2	33
190	8	50
170	7	49
175	9	53
160	6	49

Determine the following values: intercept; b^1 ; b^2 ; SS_{res} ; SS_{reg} ; F ; s^2_{res} ; $s(b^1)$; $s(b^2)$; t^1 ; t^2 .

Interpretive Problem

1. Use SPSS to develop a multiple regression model with data supplied to illustrate concepts in other chapters in this textbook. Write up your results, including interpretation of effect size and testing of assumptions.

REFERENCES

- Algina, J., & Olejnik, S. (2000). Determining sample size for accurate estimation of the squared multiple correlation coefficient. *Multivariate Behavioral Research*, 35, 119–136.
- Andrews, D. F. (1971). Significance tests based on residuals. *Biometrika*, 58, 139–148.
- Barnett, V., & Lewis, T. (1994). *Outliers in statistical data* (3rd ed.). Chichester, UK: Wiley.
- Bates, D. M., & Watts, D. G. (1988). *Nonlinear regression analysis and its applications*. New York, NY: Wiley.
- Belsley, D. A. (1991). *Conditioning diagnostics: Collinearity and weak data in regression*. New York, NY: Wiley.
- Belsley, D. A., Kuh, E., & Welsch, R. E. (1980). *Regression diagnostics*. New York, NY: Wiley.
- Bernstein, I. H. (1988). *Applied multivariate analysis*. New York, NY: Springer-Verlag.
- Berry, W. D., & Feldman, S. (1985). *Multiple regression in practice*. Beverly Hills, CA: Sage.
- Chatterjee, S., & Price, B. (1977). *Regression analysis by example*. New York, NY: Wiley.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.

- Cohen, J., & Cohen, P. (1983). *Applied multiple regression/correlation analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Erlbaum.
- D'Agostino, R. B. (1971). An omnibus test of normality for moderate and large size samples. *Biometrika*, 58, 341–348.
- DerkSEN, S., & Keselman, H. J. (1992). Backward, forward and stepwise automated subset selection algorithms: Frequency of obtaining authentic and noise variables. *British Journal of Mathematical and Statistical Psychology*, 45, 265–282.
- Erdfelder, E., Faul, F., & Buchner, A. (1996). GPOWER: A general power analysis program. *Behavior Research Methods, Instruments & Computers*, 28(1), 1–11.
- Glass, G. V., & Hopkins, K. D. (1996). *Statistical methods in education and psychology* (3rd ed.). Boston, MA: Allyn and Bacon.
- Hair, J. F., Black, W. C., Babin, B. J., Anderson, R. E., & Tatham, R. L. (2006). *Multivariate data analysis* (6th ed.). Upper Saddle River, NJ: Pearson Prentice Hall.
- Hays, W. L. (1988). *Statistics* (4th ed.). New York, NY: Holt, Rinehart and Winston.
- Hocking, R. R. (1976). The analysis and selection of variables in linear regression. *Biometrics*, 32(1), 1–49.
- Hoerl, A. E., & Kennard, R. W. (1970a). Ridge regression: Application to non-orthogonal models. *Technometrics*, 12, 591–612.
- Hoerl, A. E., & Kennard, R. W. (1970b). Ridge regression: Biased estimation for non-orthogonal models. *Technometrics*, 12, 55–67.
- Huberty, C. J. (1989). Problems with stepwise methods—Better alternatives. In B. Thompson (Ed.), *Advances in social science methodology* (Vol. 1, pp. 43–70). Greenwich, CT: JAI Press.
- Kleinbaum, D. G., Kupper, L. L., Muller, K. E., & Nizam, A. (1998). *Applied regression analysis and other multivariable models* (3rd ed.). Pacific Grove, CA: Duxbury.
- Knofszynski, G. T. (2008). Sample sizes when using multiple linear regression for prediction. *Educational and Psychological Measurement*, 68(3), 431–442.
- Larsen, W. A., & McCleary, S. J. (1972). The use of partial residual plots in regression analysis. *Technometrics*, 14, 781–790.
- Mansfield, E. R., & Conerly, M. D. (1987). Diagnostic value of residual and partial residual plots. *The American Statistician*, 41, 107–116.
- Marquardt, D. W., & Snee, R. D. (1975). Ridge regression in practice. *The American Statistician*, 29, 3–19.
- Maxwell, S. E. (2000). Sample size and multiple regression analysis. *Psychological Methods*, 5, 434–458.
- Mendoza, J. L., & Stafford, K. L. (2001). Confidence intervals, power calculation, and sample size estimation for the squared multiple correlation coefficient under the fixed and random regression models: A computer program and useful standard tables. *Educational and Psychological Measurement*, 61, 650–667.
- Meyers, L. S., Gamst, G., & Guarino, A. J. (2006). *Applied multivariate research: Design and interpretation*. Thousand Oaks, CA: Sage.
- Mickey, R. M., Dunn, O. J., & Clark, V. A. (2004). *Applied statistics: Analysis of variance and regression* (3rd ed.). Hoboken, NJ: Wiley.
- Miller, A. J. (1984). Selection of subsets of regression variables (with discussion). *Journal of the Royal Statistical Society, A*(147), 389–425.

- Miller, A. J. (1990). *Subset selection in regression*. New York, NY: Chapman & Hall.
- Miller, R. G. (1997). *Beyond ANOVA: Basics of applied statistics*. Boca Raton, FL: CRC Press.
- Morgan, G. A., & Griego, O. V. (1998). *Easy use and interpretation of SPSS for Windows: Answering research questions with statistics*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Murphy, K. R., & Myors, B. (2004). *Statistical power analysis: A simple and general model for traditional and modern hypothesis tests* (2nd ed.). Mahwah, NJ: Lawrence Erlbaum Associates Publishers.
- Myers, J. L., & Well, A. D. (1995). *Research design and statistical analysis*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Myers, R. H. (1986). *Classical and modern regression with applications*. Boston, MA: Duxbury.
- Myers, R. H. (1990). *Classical and modern regression with applications* (2nd ed.). Boston, MA: Duxbury.
- Pedhazur, E. J. (1997). *Multiple regression in behavioral research* (3rd ed.). Fort Worth, TX: Harcourt Brace.
- Rousseeuw, P. J., & Leroy, A. M. (1987). *Robust regression and outlier detection*. New York, NY: Wiley.
- Ruppert, D., & Carroll, R. J. (1980). Trimmed least squares estimation in the linear model. *Journal of the American Statistical Association*, 75, 828–838.
- Schafer, W. D. (1991). Reporting hierarchical regression results. *Measurement and Evaluation in Counseling and Development*, 24, 98–100.
- Seber, G. A. F., & Wild, C. J. (1989). *Nonlinear regression*. New York, NY: Wiley.
- Shapiro, S. S., & Wilk, M. B. (1965). An analysis of variance test for normality (complete samples). *Biometrika*, 52(3 and 4), 591–611.
- Smithson, M. (2001). Correct confidence intervals for various regression effect sizes and parameters: The importance of noncentral distributions in computing intervals. *Educational and Psychological Measurement*, 61(4), 605–632.
- Steiger, J. H., & Fouladi, R. T. (1992). R2: A computer program in interval estimation power calculation, and hypothesis testing for the squared multiple correlation. *Behavior Research Methods, Instruments & Computers*, 4, 581–582.
- Stevens, J. P. (1984). Outliers and influential data points in regression analysis. *Psychological Bulletin*, 95(2), 334–344.
- Stevens, J. P. (2009). *Applied multivariate statistics for the social sciences* (5th ed.). New York: Psychology Press.
- Tabachnick, B. G., & Fidell, L. S. (2007). *Using multivariate statistics* (5th ed.). Boston: Pearson.
- Thompson, M. L. (1978). Selection of variables in multiple regression. Part I: A review and evaluation. Part II: Chosen procedures, computations and examples. *International Statistical Review*, 46, 1–19 and 129–146.
- Weinberg, S. L., & Abramowitz, S. K. (2002). *Data analysis for the behavioral sciences using SPSS*. Cambridge, UK: Cambridge University Press.
- Weisberg, S. (1985). *Applied linear regression* (2nd ed.). New York: Wiley.
- Wetherill, G. B. (1986). *Regression analysis with applications*. London: Chapman & Hall.
- Wilcox, R. R. (1996). *Statistics for the social sciences*. San Diego, CA: Academic.

- Wilcox, R. R. (2003). *Applying contemporary statistical procedures*. San Diego, CA: Academic.
- Wonnacott, T. H., & Wonnacott, R. J. (1981). *Regression: A second course in statistics*. New York, NY: Wiley.
- Wu, L. L. (1985). Robust M-estimation of location and regression. In N. B. Tuma (Ed.), *Sociological Methodology* (pp. 316–388). San Francisco, CA: Jossey-Bass.

Chapter 5

LOGISTIC REGRESSION

CHAPTER OUTLINE

5.1 What Logistic Regression Is and How It Works	118
5.1.1 Characteristics	121
5.1.2 Sample Size	133
5.1.3 Power	133
5.1.4 Effect Size	133
5.1.5 Assumptions	134
5.2 Mathematical Introduction Snapshot	138
5.3 Computing Logistic Regression Using SPSS	139
5.4 Data Screening	150
5.4.1 Noncollinearity	151
5.4.2 Linearity	152
5.4.3 Independence	155
5.4.4 Absence of Outliers	156
5.4.5 Assessing Classification Accuracy	157
5.5 Power Using G*Power	160
5.5.1 Post Hoc Power for Logistic Regression Using G*Power	160
5.5.2 A Priori Power for Logistic Regression Using G*Power	162
5.6 Research Question Template and Example Write-Up	163

KEY CONCEPTS

1. Logit
2. Odds
3. Odds ratio

In the previous chapter, we examined the use of the least squares criterion in multiple regression models that allow us to examine the relationship between one or more predictors when the outcome is continuous. In this chapter, we are introduced to logistic regression, which can also be used when the outcome is categorical and that allows model prediction. Logistic regression and discriminant analysis (which is discussed in an upcoming chapter) share similarities, and there can be confusion on when one is more appropriate than the other. Understanding that you may not be fully familiar with discriminant analysis, we'll offer the condensed version of how the two procedures contrast. The assumptions of multivariate normality and equal variance-covariance matrices, which are required in discriminant analysis, do not hold for logistic regression. Thus, logistic regression is more robust than discriminant analysis when these assumptions are not met. Additionally, logistic regression is oftentimes less interpretatively challenging than discriminant analysis given that it falls within the regression family, more common to many researchers as compared to discriminant analysis.

For the purposes of this chapter, we will concentrate on binary logistic regression, which is used when the outcome has only two categories (i.e., dichotomous, binary, or sometimes referred to as a Bernoulli outcome). The logistic regression procedure appropriate for more than two categories is called multinomial (or polytomous) logistic regression. Readers interested in learning more about multinomial logistic regression will be provided some additional references later in this chapter. Also in this chapter, we discuss methods that can be used to enter predictors in logistic regression models. Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying logistic regression, (b) determine and interpret the results of logistic regression, (c) understand and evaluate the assumptions of logistic regression, and (d) have a basic understanding of methods of entering the covariates.

5.1 WHAT LOGISTIC REGRESSION IS AND HOW IT WORKS

Oso Wyse, one of the four amazingly talented statistical gurus in the statistics and research lab, has just had a conversation with his faculty advisor. He finds himself embarking on a challenging statistical project.

Oso Wyse finds himself on his first statistical expedition as a graduate research assistant in the stats and research lab. After an introduction from his faculty advisor, Oso meets with Dr. Malani, a faculty member in the early childhood department. Dr. Malani has collected data on children who will be entering kindergarten in the fall. Interested in kindergarten readiness

issues, Dr. Malani wants to know if a teacher observation scale for social development and family structure (single-parent versus two-parent home) can predict whether children are prepared or unprepared to enter kindergarten. Oso suggests the following research question to Dr. Malani: *Can kindergarten readiness (prepared vs. unprepared) be predicted by social development and family structure (single-parent vs. two-parent home)?* Given that the outcome is dichotomous, Oso determines that binary logistic regression is the appropriate statistical procedure to use to answer Dr. Malani's question. Oso then proceeds with assisting Dr. Malani in analyzing the data.

If the dependent variable is binary (i.e., dichotomous or having only two categories), then ordinary least squares (OLS) regression described earlier in this text is inappropriate. Although OLS regression can easily accommodate dichotomous independent variables through dummy coding (i.e., assignment of 1 and 0 to the categories where '1' is traditionally coded as the category of interest, i.e., case outcome; '0' is traditionally coded as the noncase outcome or reference category), it is an entirely different case when the *outcome* is dichotomous. Applying OLS regression to a binary outcome creates problems. For example, a dichotomous outcome violates normality and homogeneity assumptions in OLS regression. In addition, OLS estimates are based on linear relationships between the independent and dependent variables, and forcing a linear relationship (as seen in Figure 5.1) in the case of a binary outcome is erroneous

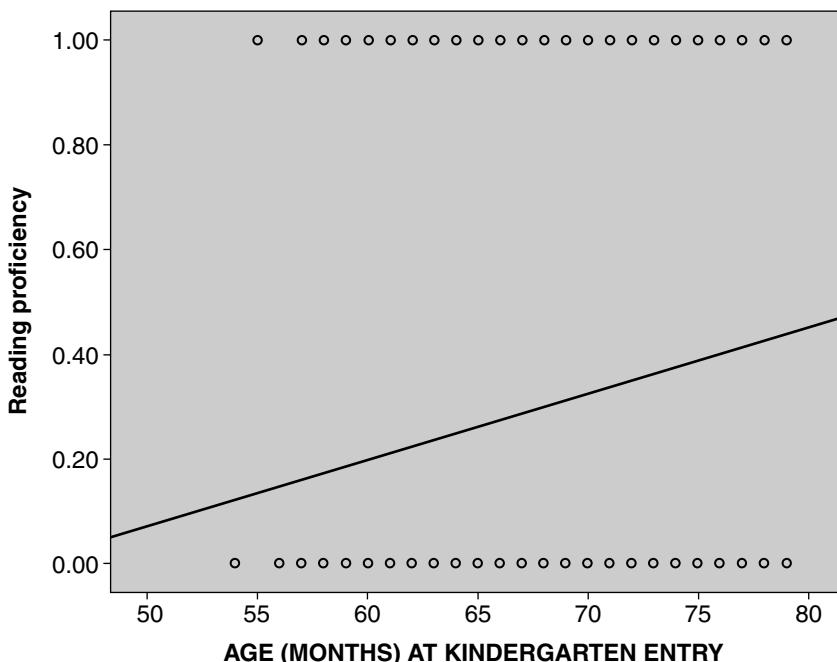


FIGURE 5.1
Nonlinearity of Binary Outcome

[although we found at least one author (Hellevik, 2009) who argues that OLS regression can be used with dichotomous outcomes]. As seen in this figure, there is obviously not a linear relationship between age at kindergarten entry and reading proficiency (i.e., proficient or not proficient).

As part of the regression family, logistic regression still allows a prediction to be made; however, now the prediction is whether or not the unit under investigation falls into one of the two categories of the dependent variable. Initially used mostly in the hard sciences, this method has become more broadly popular in recent years, as there are many situations where researchers want to examine outcomes that are discrete, rather than continuous, in nature. Some examples of dichotomous dependent variables are pass/fail, surviving surgery/not, admit/reject, vote for/against, employ/not, win/lose, or purchase/not. The idea of using a dichotomous variable was introduced in the multiple regression chapter as the concept of a *dummy variable*, where the first condition is indicated by a value of 1 (e.g., prepared for kindergarten), whereas a value of 0 indicates the opposite condition (e.g., unprepared for kindergarten). Understanding the coding of 0 and 1 is very important for interpretation purposes. Again, ‘1’ is traditionally the case outcome (with results interpreted in terms of cases) and ‘0’ the noncase or reference category. For the purposes of this text, our discussion will concentrate on dichotomous outcomes where logistic regression is appropriate (i.e., binary logistic regression, referred to throughout this chapter simply as logistic regression). For conditions for which there are more than two possible categories for the dependent variable (e.g., three categories, such as ‘above satisfactory performance,’ ‘satisfactory performance,’ and ‘below satisfactory performance’), multinomial logistic regression may be appropriate. An example of the data structure for a logistic regression model with a binary outcome (prepared vs. unprepared for kindergarten), one continuous predictor (social development) and one dichotomous dummy-coded predictor (family structure: single-parent vs. two-parent home) is presented in Table 5.1.

■ TABLE 5.1

Kindergarten Readiness Example Data

Child	Social Development (X_1)	Family Structure (X_2)	Kindergarten Readiness (Y)
1	15	Single-parent home (0)	Unprepared (0)
2	12	Single-parent home (0)	Unprepared (0)
3	18	Single-parent home (0)	Prepared (1)
4	20	Single-parent home (0)	Prepared (1)
5	11	Single-parent home (0)	Unprepared (0)
6	17	Single-parent home (0)	Prepared (1)
7	14	Single-parent home (0)	Unprepared (0)
8	18	Single-parent home (0)	Prepared (1)
9	13	Single-parent home (0)	Unprepared (0)
10	10	Single-parent home (0)	Unprepared (0)

TABLE 5.1 (continued)

Kindergarten Readiness Example Data

Child	Social Development (X_1)	Family Structure (X_2)	Kindergarten Readiness (Y)
11	22	Two-parent home (1)	Unprepared (0)
12	25	Two-parent home (1)	Prepared (1)
13	23	Two-parent home (1)	Prepared (1)
14	21	Two-parent home (1)	Prepared (1)
15	30	Two-parent home (1)	Prepared (1)
16	27	Two-parent home (1)	Prepared (1)
17	26	Two-parent home (1)	Prepared (1)
18	28	Two-parent home (1)	Prepared (1)
19	24	Two-parent home (1)	Unprepared (0)
20	30	Two-parent home (1)	Prepared (1)

5.1.1 Characteristics

5.1.1.1 Logistic Regression Equation

As we learned previously with ordinary least squares regression, knowledge of the independent variable(s) provides the information necessary to be able to estimate a precise numerical value of the dependent variable, a predicted value. The following formula recaps the sample multiple regression equation where Y is the predicted outcome for individual i based on (a) the Y intercept, a , the value of Y when all predictor values are zero; (b) the product of the value of the independent variables, X 's, and the regression coefficients, b_k ; and (c) the residual, ε_i :

$$Y_i = a + b_1 X_1 + \dots + b_m X_m + \varepsilon_i$$

As we see, the logistic regression equation is similar in concept to simple and multiple linear regression, but operates much differently. In logistic regression, the binary dependent variable is transformed into a logit variable (which is the natural log of the odds of the dependent variable occurring or not occurring), and the parameters are then estimated using maximum likelihood. The end result is that the odds of an event occurring are estimated through the logistic regression model (whereas OLS estimates a precise numerical value of the dependent variable).

To understand how the logistic regression equation operates, there are three primary computational concepts that must be understood: probability, odds, and the logit. These express the same thing, only in different ways (Menard, 2000). Let us first consider probability.

5.1.1.2 Probability

The overarching difference between OLS regression (i.e., simple and multiple linear regression) and logistic regression is the measurement scale of the outcome. With

OLS regression, our outcome is continuous in scale (i.e., interval or ratio measurement scale). In binary logistic regression, our outcome is dichotomous—one of two categories. Let us use kindergarten readiness ('prepared for kindergarten' coded as '1' vs. 'unprepared' coded as '0') as an example of our logistic regression outcome. Therefore, what the regression equation allows us to predict is substantially different for OLS as compared to logistic regression. In comparison to OLS, which allows us to compute a precise numerical value (e.g., a specific predicted score for the dependent variable), the logistic regression equation allows us to compute a *probability*—more specifically, the *probability* that the dependent variable will occur. The logistic regression equation, therefore, generates predicted probabilities that fall between values of 0 and 1. The probability of a case or unit being classified into the lowest numerical category [i.e., $P(Y = 0)$] or in the case of our example, the probability that a child will be 'unprepared' for kindergarten] is the equal to 1 minus the probability that it falls within the highest numerical category [i.e., $P(Y = 1)$] or the probability that a child will be 'prepared' for kindergarten]. This equates to $P(Y = 0) = 1 - P(Y = 1)$. Applied to our example, the probability that a child will be unprepared for kindergarten is equal to one minus the probability that a child will be prepared for kindergarten. In other words, the knowledge of the probability of one category occurring (e.g., unprepared for kindergarten) allows us to easily determine the probability that the other category will occur (e.g., prepared) as the total probability must equal 1.0. Remember, however, that probabilities have to fall within the range of 0 to 1. As we know, it is not possible to have a negative probability, nor is it possible to have a probability greater than 1 (i.e., greater than 100%). If we try to model the probability as the dependent variable in our OLS equation, it is mathematically possible that the predicted values would be negative or greater than 1—values that are outside the range of what is feasible when considering probabilities. Therefore, this is where our logistic regression equation takes a turn from what we learned with linear regression.

5.1.1.3 Odds and Logit (or Log Odds)

So far, we have talked about the outcome of our logistic regression equation as being a probability, and we also know that predicted probabilities must be between 0 and 1. As we think about how to estimate probabilities, we will see that this takes a few steps to achieve. Rather than the dependent variable being a probability, if it were an *odds value*, then values greater than 1 would be possible and appropriate. **Odds** are simply the ratio of the probability of the dependent variable's two outcomes. The odds that the outcome of a binary variable is 1 (i.e., public school attendance) rather than 0 (or private school attendance), is simply the ratio of the odds that Y is equal 1 to the odds that Y does not equal 1. In mathematical terms, this can be written as follows:

$$Odds(Y = 1) = \frac{P(Y = 1)}{1 - P(Y = 1)}$$

As we see in Table 5.2, when the probability that $Y = 1$ (e.g., prepared for kindergarten) equals .50 (column 1 in Table 5.2), then $1 - P(Y = 1)$ (or unprepared for kindergarten)

TABLE 5.2

Illustration of Logged Odds

$P(Y = 1)$	$1 - P(Y = 1)$	$Odds(Y = 1) = \frac{P(Y = 1)}{1 - P(Y = 1)}$	$\ln[Odds(Y = 1)] = \ln\left[\frac{P(Y = 1)}{1 - P(Y = 1)}\right]$
.001	.999	.001/.999 = .001	$\ln(.001) = -6.908$
.100	.900	.100/.900 = .111	$\ln(.111) = -2.198$
.300	.700	.300/.700 = .429	$\ln(.429) = -.846$
.500	.500	.500/.500 = 1.000	$\ln(1.000) = 0.000$
.700	.300	.700/.300 = 2.333	$\ln(2.333) = .847$
.900	.100	.900/.100 = 9.000	$\ln(9.000) = 2.197$
.999	.001	.999/.001 = 999.000	$\ln(999) = 6.907$

is .50 (column 2) and the odds are equal to 1.00 (column 3). When the probability of $Y = 1$ (e.g., prepared) is very small (say, .100 or less), then the odds for being prepared for kindergarten are also very small and approach 0 (i.e., the smaller the probability that a child is prepared for kindergarten). However, as the probability of $Y = 1$ (e.g., being prepared for kindergarten) increases, the odds (column 3) increase tremendously. Thus, the issue that we are faced with when using odds is that while odds can be infinitely large, we are still limited in that the minimum value is 0 and we still do not have data that can be modeled linearly. When $P(Y = 1) < .5$, the slope below an odds of 1.0 is very steep; yet when $P(Y = 1) > .5$, the slope above odds of 1.0 is much more gradual. It might also be worth noting at this point that the reciprocal odds have the same magnitude of effect but are asymmetrical, and the natural log functions to create a symmetrical outcome variable.

Changing the scale of the odds by taking the natural logarithm of the odds (also called *logit Y* or *log odds*) provides us with a value of the dependent variable that can theoretically range from negative infinity to positive infinity. Thus, taking the log odds of Y creates a linear relationship between X and the probability of Y (Pampel, 2000). The natural log of the odds is calculated as follows, with the residual being the difference between the predicted probability and the actual value of the dependent variable (0 or 1):

$$\ln \frac{P(Y = 1)}{1 - P(Y = 1)} = Logit(Y)$$

In column 4 of Table 5.2, we see what happens when the logit transformation is made. As the odds increase from 1 to positive infinity, the logit (or log odds) of Y becomes larger and larger (and remains positive). As the odds decrease from 1 to 0, the logit (or log odds) of Y is negative and grows larger and larger (in absolute value).

The logit of Y equation is interpreted very similarly to that of OLS. For each one-unit change in the independent variable, the logistic regression coefficients represent the change in the predicted log odds of being in a category. In comparison to OLS

regression, the regression coefficients have the exact same interpretation. The difference in interpretation with logistic regression is that the outcome now represents a *log odds* rather than a precise numerical value as we saw with OLS regression. Linking the logit back to probabilities, a one-unit change in the logit equals a bigger change in probabilities that are near the center as compared to the extreme values. This happens because of the linearization once we take the natural log. Taking the natural log stretches the S-shaped curve into a linear form, thus the values at the extreme are stretched less, so to speak, as compared to the values in the middle (Pampel, 2000). By working with log odds, our familiar additive regression equation is applicable:

$$\ln \frac{P(Y=1)}{1-P(Y=1)} = \text{Logit}(Y) = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m$$

It is important to note that although we were accustomed to examining standardized regression coefficients in OLS regression, it is not the norm that standardized coefficients are computed for logistic regression models by statistical software. Standardization is ordinarily accomplished by taking the product of the unstandardized regression coefficient and the ratio of the standard deviation of X to the standard deviation of Y . The interpretation of a standard deviation change in a continuous variable thus makes sense; however, this is not the case for a dichotomous variable nor is it the case for the log odds (which is the predicted outcome and which does not have a standard deviation).

While interpretation of the logistic equation is relatively straightforward, as it holds many similarities to OLS regression, log odds are not a metric that we use often. Therefore understanding what it means when a predictor, X , has some effect on the log odds, Y , can be difficult. This is where odds come back into the picture.

If we exponentiate the logit (Y) (i.e., the outcome of our logistic regression equation), then it converts back to the odds (see equation below). Now we can interpret the independent variables as affecting the odds (rather than log odds) of the outcome:

$$\text{Odds}(Y=1) = e^{\text{logit}(Y)} = e^{\ln[\text{Odds}(Y=1)]} = e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m} = (e^\alpha)(e^{\beta_1 X_1})(e^{\beta_2 X_2}) \dots (e^{\beta_m X_m})$$

As can be seen here, the exponentiation creates an equation that is multiplicative rather than additive, and this then changes the interpretation of the exponentiated coefficients. In previous regression equations we have studied, when the product of the regression coefficient and its predictor is 0, that variable adds nothing to the prediction of the dependent variable. In a multiplicative environment, a value of 0 corresponds to a coefficient of 1. In other words, a coefficient of 1 will not change the value of the odds (i.e., the outcome). Coefficients greater than 1 increase the odds, and coefficients less than 1 decrease the odds. In addition, the odds will change more the greater the distance the value is from 1.

Converting the odds back to a probability can be done through the following formula:

$$P(Y = 1) = \frac{Odds(Y = 1)}{1 + Odds(Y = 1)} = \frac{e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m}}{1 + e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m}}$$

Probability values close to 1 indicate increased likelihood of occurrence. In our example, since ‘1’ indicates public school attendance, a probability close to 1 would indicate a child was more likely to attend public school. Children with probabilities close to 0 suggest a decreased probability of attending public school (and increased probability of attending private school).

5.1.1.4 Estimation and Model Fit

Now that we understand the logistic regression process and resulting equations a bit better, it is time to turn our attention to how the equation is estimated and how we can determine how well the model fits. We previously learned with multiple regression that the data from the observed values of the independent variables in the sample were used to estimate or predict the values of the dependent variable. In logistic regression, we are also using the knowledge of the values of our predictor(s) to estimate the outcome (i.e., log odds). Now we are using a method called maximum likelihood estimation to estimate the values of the parameters (i.e., the logistic coefficients). As we just learned, the dependent variable in a logistic regression model is transformed into a logit value, which is the natural log of the odds of the dependent variable occurring or not occurring. Maximum likelihood (ML) estimation is then applied to the model and estimates the odds of occurrence after transformation into the logit. The ‘likelihood’ in maximum likelihood refers to the likelihood of the data occurring given a specific value for population parameters that have been assumed. It is the probability of the data contingent upon a parameter-estimate that is being maximized. Whereas in OLS the sum of squared distance of the observed data to the regression line was minimized, in maximum likelihood the log likelihood is maximized.

The log of the likelihood function (sometimes abbreviated as *LL*) that results from ML estimation then reflects the likelihood of observing the sample statistics given the population parameters. The log likelihood provides an index of how much has not been explained in the model after the parameters have been estimated, and as such, the *LL* can be used as an indicator of model fit. The values of the log likelihood function vary from zero to negative infinity, with values closer to zero suggesting better model fit and larger values (in absolute value terms) indicating poorer fit. The log likelihood value will approach zero the closer the likelihood value is to one. When this happens, this suggests the observed data could be generated from these population parameters. In other words, the smaller the log likelihood, the better the model fit. It follows therefore, that the log likelihood value will grow more negative the closer the likelihood function is to zero. This suggests that the observed data are *less* likely to be generated from these population parameters.

Maximum likelihood estimation performed by statistical software usually begins the estimation process with all regression coefficients equal to the most conservative

estimate (i.e., the least squares estimates). Better model fit is accomplished through the use of an algorithm, which generates new sets of regression coefficients that produce larger log likelihoods. This is an iterative process that stops when the selection of new parameters creates very little change in the regression coefficients and very small increases in the log likelihood—so small that there is little value in any further estimation.

5.1.1.5 Significance Tests

As with multiple regression, there are two tests of significance in logistic regression. Specifically, these involve testing the significance of the overall logistic regression model and testing the significance of each of the logistic regression coefficients.

Test of Significance of the Overall Regression Model

The first test is the test of statistical significance to determine overall model fit and provides evidence of the extent to which the predicted values accurately represent the observed values (Xie, Pendergast & Clarke, 2008). We consider several overall model tests including (a) change in log likelihood, (b) Hosmer and Lemeshow goodness-of-fit test, (c) pseudo-variance explained, and (d) predicted group membership. Additional work (e.g., Xie et al., 2008) has recently been conducted on new methods to assess model fit, but these are not currently available in statistical software or easily computed. Also in this section, we briefly address sensitivity, specificity, false positive, false negative, and cross-validation.

Change in Log Likelihood

One way to test overall model fit is the likelihood ratio test. This test is based on the change in the log likelihood function from a smaller model (often the baseline or intercept-only model) to a larger model that includes one or more predictors (sometimes referred to as the fitted model). Although we indicate that the smaller model is often the intercept-only model, this test can also be used to examine changes in model fit from one fitted model to another fitted model and we will discuss this in a bit. This likelihood ratio test is similar to the overall F test in OLS regression and tests the null hypothesis that all the regression coefficients are equal to zero. Using statistical notation, we can denote the null and alternative hypotheses for the regression coefficients as follows:

$$\begin{aligned} H_0 &: \beta_1 = \beta_2 = \dots = \beta_m = 0 \\ H_1 &: H_0 \text{ is false} \end{aligned}$$

For explanation purposes, we assume the smaller model is the baseline or intercept-only model. The baseline log likelihood is estimated from a logistic regression model that includes only the constant (i.e., intercept) term. The model log likelihood is estimated from the logistic regression model that includes the constant and the relevant

predictor(s). By multiplying the difference in these log likelihood functions by -2 , a chi-square test is produced with degrees of freedom equal to the difference in the degrees of freedom of the models ($df = df_{model} - df_{baseline}$) (where ‘model’ refers to the fitted model that includes one or more predictors). In the case of the constant only model, there is only one parameter estimated (i.e., the intercept), so there is only one degree of freedom. In models that include independent variables, the degrees of freedom are equal to the number of independent variables in the model plus one for the constant. The larger the difference between the baseline and model LL values, the better the model fit. It is important to note that the log likelihood difference test assumes nested models. In other words, all elements that are included in the baseline or smallest model must also be included in the fitted model. As alluded to previously, the change in log likelihood test can be used for more than just comparing the intercept-only model to a fitted model. Researchers often use this test in the model building process to determine if adding predictors (or sets of predictors) aids in model fit by comparing one fitted model to another fitted model. In general, the change in log likelihood is computed as follows:

$$\chi^2 = 2(LL_{model} - LL_{baseline})$$

Hosmer-Lemeshow Goodness-of-Fit Test

The Hosmer-Lemeshow goodness-of-fit test is another tool that can be used to examine overall model fit. The Hosmer-Lemeshow statistic is computed by dividing cases into deciles (i.e., 10 groups) based on their predicted probabilities. Then a chi-square value is computed based on the observed and expected frequencies. This is a chi-square test for which the researcher does *not* want to find statistical significance. Nonstatistically significant results for the Hosmer-Lemeshow test indicate the model has acceptable fit. In other words, the predicted or estimated model is not statistically significantly different from the observed values. Although the Hosmer-Lemeshow test can easily be requested in SPSS, it has been criticized for being conservative (i.e., lacking sufficient power to detect lack of fit in instances such as nonlinearity of an independent variable), too likely to indicate model fit when five or fewer groups (based on the decile groups created in computing the statistic) are used to calculate the statistic, and offers little diagnostics to assist the researcher when the test indicates poor model fit (Hosmer, Hosmer, LeCessie & Lemeshow, 1997). Additionally, this test can be overly conservative unless one has very large sample sizes.

Pseudo-Variance Explained

Another overall model fit index for logistic regression is pseudo-variance explained. This index is akin to multiple R^2 (or the coefficient of determination) in OLS regression and can also be considered an effect size measure for the model. The reason these values are considered pseudo-variance explained in logistic regression is that the variance in a dichotomous outcome, as evident in logistic regression, differs as compared to the variance of a continuous outcome, as present in OLS regression.

There are a number of multiple R^2 pseudo-variance explained values that can be computed in logistic regression. We discuss the following: (a) Cox and Snell (1989), (b) Nagelkerke (1991), (c) Hosmer and Lemeshow (1989), (d) Aldrich and Nelson (1984), (e) Harrell (1986), and (f) traditional R^2 . Of these, SPSS automatically computes the Cox and Snell and Nagelkerke indices. There is, however, no consensus on which (if any) of the pseudo-variance explained indices are best, and many researchers choose not to report any of them in their published results. If you do choose to use and/or report one or more of these values, they should be used only as a guide “without attributing great importance to a precise figure” (Pampel, p. 50).

The Cox and Snell R^2 (1989) is computed as the ratio of the likelihood values raised to the power of $2/n$ (where n is sample size). A problem is that the computation is such that the theoretical maximum of one cannot be obtained, even when there is perfect prediction:

$$R_{CS}^2 = 1 - \left(\frac{LL_{baseline}}{LL_{model}} \right)^{\frac{2}{n}}$$

Nagelkerke (1991) adjusts the Cox and Snell value so that the maximum value of one can be achieved, and it is computed as follows:

$$R_N^2 = \frac{R_{CS}^2}{1 - (LL_{baseline})^{\frac{2}{n}}}$$

Hosmer and Lemeshow's (1989) R^2 is the proportional reduction in the log likelihood (in absolute value terms). Although not provided by SPSS, it can easily be computed by the ratio of the model to baseline $-2LL$. Ranging from zero to one, this value provides an indication of how much the badness of fit of the baseline model is improved by the inclusion of the predictors in the fitted model. Hosmer and Lemeshow's (1989) R^2 is computed as:

$$R_L^2 = \frac{-2LL_{model}}{-2LL_{baseline}}$$

Harrell (1986) proposed that Hosmer and Lemeshow's R^2 be adjusted for the number of parameters (i.e., independent variables) in the model. This adjustment (where m equals the number of independent variables in the model) to the computation makes this R^2 value akin to the adjusted R^2 in OLS regression. It is computed as:

$$R_{LA}^2 = \frac{(-2LL_{model}) - 2m}{-2LL_{baseline}}$$

Aldrich and Nelson (1984) provided an alternative to the R_L^2 that is equivalent to the squared contingency coefficient. This measure has the same problem as the Cox and

Snell R^2 ; the theoretical maximum of one cannot be obtained even when the independent variable(s) perfectly predict the outcome. It is computed as:

$$pseudo\ R^2 = \frac{-2LL_{model}}{-2LL_{model} + n}$$

The traditional R^2 , the coefficient of determination as used in simple and multiple regression, can also be used in logistic regression (only with binary logistic regression, as the mean and variance of a dichotomous variable make sense; however, the mean, for example, in a dummy-coded variable situation, is equal to the proportion of cases in the category labeled as 1). R^2 can be computed by correlating the observed values of the binary dependent variable with the predicted values (i.e., predicted probabilities) obtained from the logistic regression model and then squaring the correlated value. Predicted probability values can easily be saved when generating logistic regression models in SPSS.

Predicted Group Membership

Another test of model fit for logistic regression can be accomplished by evaluating predicted to observed group membership. Assuming a cut value of .50, cases with predicted probabilities at .5 or above are predicted as 1 and predicted probabilities below .5 are predicted as 0. A crosstab table of predicted to observed predicted probabilities provides the frequency and percentage of cases correctly classified. Correct classification would be seen in cases that have the same value for both the predicted and observed values. A perfect model produces 100% correctly classified cases. A model that classifies no better than chance would provide 50% correctly classified cases. Press's Q is a chi-square statistic with one degree of freedom and can be used as a formal test of classification accuracy. It is computed as:

$$Q = \frac{[N - (nK)]^2}{N(K - 1)}$$

where N is the total sample size, n represents the number of cases that were correctly classified, and K equals the number of groups. As with other chi-square statistics we have examined, this test is sensitive to sample size. Also, it is important to note that focusing solely on the correct classification overall (as is done with Press's Q) may result in overlooking one or more groups that have unacceptable classification. The researcher should evaluate the classification of each group in addition to the overall classification.

Sensitivity is the probability that a case coded as 1 for the dependent variable (aka ‘positive’) is classified correctly. In other words, sensitivity is the percentage of correct predictions of the cases that are coded as 1 for the dependent variable. In the kindergarten readiness example that we will review later, of those 12 children who were prepared for kindergarten (i.e., coded as 1 for the dependent variable), 11 were correctly classified. Thus, the sensitivity is 11/12 or about 92%.

Specificity is the probability that a case coded as 0 for the dependent variable (aka ‘negative’) is classified correctly. In other words, specificity is the percentage of correct predictions of the cases that are coded as 0 for the dependent variable. In the kindergarten readiness example that we will review later, of those 8 children who were unprepared for kindergarten (i.e., coded as 0 for the dependent variable), 7 were correctly classified. Thus, the specificity is 7/8 or 87.5%.

False positive rate is the probability that a case coded as 0 for the dependent variable (aka ‘negative’) is classified *incorrectly*. In other words, this is the percentage of cases in error where the dependent variable is predicted to be 1 (i.e., prepared), but in fact the observed value is 0 (i.e., unprepared). In the kindergarten readiness example that we will review later, of those 8 children who were unprepared for kindergarten (i.e., coded as 0 for the dependent variable), 1 was incorrectly classified. Thus, the false positive rate is 1/8 or 12.5%. The false positive rate is also computed as one minus specificity.

False negative rate is the probability that a case coded as 1 for the dependent variable (aka ‘positive’) is classified *incorrectly*. In other words, this is the percentage of cases in error where the dependent variable is predicted to be 0 (i.e., unprepared), but in fact the observed value is 1 (i.e., prepared). In the kindergarten readiness example that we will review later, of those 12 children who were prepared for kindergarten (i.e., coded as 1 for the dependent variable), 1 was incorrectly classified. Thus, the false negative rate is 1/12 or about 8%. The false negative rate is also computed as one minus sensitivity.

Cross-validation

A recommended best practice in logistic regression is to cross-validate the results. If the sample size is sufficient, this can be accomplished by using 75%–80% of the sample to derive the model and then use the remaining cases (the holdout sample) to determine its accuracy. With cross-validation, you are in essence testing the model on two samples—a primary sample (which represents the largest percentage of the sample size) and a holdout sample (that which remains). If classification accuracy of the holdout sample is within 10% of the primary sample, this provides evidence of the utility of the logistic regression model.

5.1.1.6 Test of Significance of the Logistic Regression Coefficients

The second test in logistic regression is the test of the statistical significance of each regression coefficient, b_k . This test allows us to determine if the individual coefficients are statistically significantly different from zero. The null and alternative hypothesis can be illustrated in the same mathematical notation as we used with OLS regression:

$$H_0 : \beta_k = 0$$

$$H_1 : \beta_k \neq 0$$

Interpreting the test provides evidence of the probability of obtaining the observed sample coefficient by chance if the null hypothesis was true (i.e., if the population regression coefficient value was zero). The Wald statistic, which follows a chi-square distribution, is used as the test statistic for regression coefficients in SPSS. This is calculated by squaring the ratio of the regression coefficient divided by its standard error:

$$W = \frac{\beta_k^2}{SE_{\beta_k^2}}$$

When the logistic regression coefficients are large (in absolute value), rounding error can create imprecision in estimation of the standard errors. This can result in inaccuracies in testing the null hypothesis, and more specifically, increased Type II errors (i.e., failing to reject the null hypothesis when the null hypothesis is false). An alternative to the Wald test, in situations such as this, is the difference in log likelihood test previously described to compare models with and without the variable of interest (Pampel, 2000).

Raftery (1995) proposed a Bayesian information criterion (BIC), computed as the difference between the chi-square value and the natural log of the sample size, that could also be applied to testing logistic regression coefficients:

$$BIC = \chi^2 - \ln n$$

To reject the null hypothesis, the BIC should be positive (i.e., greater than zero). That is, the chi-square value must be greater than the natural log of the sample size. BIC values below zero suggest that the variable contributes little to the model. BIC values between 0 and +2 are considered weak; between 2 and 6, positive; between 6 and 10, strong; and more than 10, very strong.

Beyond determining statistical significance of the individual predictors, you may also want to assess which predictors are adding the most to the model. In OLS regression, we examined the standardized regression coefficients. There are no traditional standardized regression coefficients provided in SPSS for logistic regression, but they are easy to calculate. Simply standardize the predictors before generating the logistic regression model, and then run the model as desired. You can then interpret the logistic regression coefficients as standardized regression coefficients (if necessary, review the multiple regression chapter).

We can also form a confidence interval around the logistic regression coefficient, b_k . The confidence interval formula is the same as in OLS regression, that is, the logistic regression coefficient plus or minus the product of the tabled critical value and the standard error:

$$CI(b_k) = b_k \pm t_{(\alpha/2)} s_b$$

The null hypothesis that we tested was $H_0 : \beta_k = 0$. It follows that if our confidence interval contains zero, then the logistic regression coefficient (b_k) is not statistically

significantly different from zero at the specified significance level. We can interpret this to say that β_k will be included in $(1-\alpha)\%$ of the sample confidence intervals formed from multiple samples.

5.1.1.7 Methods of Predictor Entry

The three categories of model building that will be discussed include (a) simultaneous logistic regression, (b) stepwise logistic regression, and (c) hierarchical regression.

Simultaneous Logistic Regression

With simultaneous logistic regression, all the independent variables of interest are included in the model in one set. This method of model building is usually used when the researcher does not hypothesize that some predictors are more important than others are. This method of entry allows you to evaluate the contribution of an independent variable over and above that of all other predictors in the model (i.e., each independent variable is evaluated as if it was the last one to enter the equation). One problem that may be encountered with this method of entry is related to strong correlations between the predictor and the outcome. An independent variable that has a strong bivariate correlation with the dependent variable may indicate a weak correlation when entered simultaneously with other predictors. In SPSS, this method of entry is referred to as “ENTER.”

Stepwise Logistic Regression

Stepwise logistic regression is a data-driven model building technique where the computer algorithms drive variable entry rather than theory. Issues with this type of technique have previously been outlined in the discussion associated with this method in multiple regression and thus are not rehashed here. If stepwise logistic regression is determined to be the most appropriate strategy to build your model, Hosmer and Lemeshow (2000) suggest setting a more liberal criteria for variable inclusion (e.g., $\alpha = .15$ to $.20$). They also provide specific recommendations on dealing with interaction terms and scales of variables. Because it is only in unusual instances that this method of model building is appropriate (e.g., exploratory research), additional coverage of the suggestions by Hosmer and Lemeshow is not presented.

SPSS offers forward and backward stepwise methods. For both forward and backward methods, options include conditional, LR, and Wald. The differences between these options are mathematically driven. The LR method of entry uses the $-2LL$ for estimating entry of independent variables. The conditional method also uses the likelihood ratio test, but one that is considered to be computationally quicker. The Wald method applies the Wald test to determining entry of the independent variables. With forward stepwise methods, the model begins with a constant only and, based on some criterion, independent variables are added one at a time until a specified cutoff is achieved (e.g., all independent variables included in the model are statistically significant and

any additional variables not included in the model are not statistically significant). Backward stepwise methods work in the reverse fashion where initially all independent variables (and the constant) are included. Independent variables are then removed until only those that are statistically significant remain in the model, and including an omitted independent variable would not improve the model.

Hierarchical Regression

In hierarchical regression, the researcher specifies *a priori* a sequence for the individual predictor variables (not to be confused with hierarchical linear models, which is a regression approach for analyzing nested data collected at multiple levels, such as child, classroom, and school). The analysis proceeds in a forward selection, backward elimination, or stepwise selection mode according to a researcher specified, theoretically based sequence, rather than an unspecified statistically based sequence. In SPSS, this is conducted by entering predictors in **blocks** and selecting their desired method of entering variables in each block (e.g., simultaneously, forward, backward, stepwise). Because this method was explained in detail in reference to multiple regression and operation of this method of variable selection is the same in logistic regression, additional information will not be presented.

5.1.2 Sample Size

Simulation research suggests that logistic regression is best used with large samples. Samples of size 100 or greater are needed to accurately conduct tests of significance for logistic regression coefficients (Long, 1997). Note that for illustrative purposes, the example in this chapter uses a sample size of 20. We recognize this is insufficient in practice, but have used it for greater ease in presenting the data.

5.1.3 Power

Power in logistic regression can be computed *a priori* (which is ideal) to determine requisite sample size, as well as post hoc. It is important to note the relationship between goodness-of-fit and power in the context of logistic regression, for example, the power of the Hosmer-Lemeshow goodness-of-fit statistic in detecting ill fit (e.g., nonlinearity in predictor variables) (Xie, Pendergast, & Clarke, 2008).

5.1.4 Effect Size

We have already talked about multiple R^2 pseudo-variance explained values, which can be used not only to gauge model fit, but also as measures of effect size. Another important statistic in logistic regression is the **odds ratio** (*OR*), also an effect size index that is similar to R^2 . The odds ratio is computed by exponentiating the logistic regression coefficient e^{b_k} . Conceptually, this is the odds for one category (e.g., prepared for kindergarten) divided by the odds for the other category (e.g., unprepared for kindergarten). The null hypothesis to be tested is that $OR = 1$, which indicates that

there is no relationship between a predictor variable and the dependent variable. If an odds ratio of 1 indicates no effect, then an odds ratio greater than 1 indicates higher odds of the outcome occurring. An odds ratio of less than 1 indicates lower odds of the outcome occurring. Thus, we want to find *OR* to be significantly different from 1.

When the independent variable is continuous, the odds ratio represents the amount by which the odds change for a one-unit increase in the independent variable. When the odds ratio is greater than 1, the independent variable increases the odds of occurrence. When the odds ratio is less than 1, the independent variable decreases the odds of occurrence. The odds ratio is provided in SPSS output as “*Exp(B)*” in the table labeled “Variables in the Equation.” In predicting kindergarten readiness, social development is a continuous covariate with a resulting odds ratio of 2.631. We can interpret this odds ratio to be that for every one unit increase in social development, the odds of being ready for kindergarten (i.e., prepared) increase by 263%, controlling for the other variables in the model.

In the case of categorical variables, including dichotomous, multinomial, and ordinal variables, odds ratios are often interpreted in terms of their relative size or the change in odds ratios in comparing models. Consider first the case of a dichotomous variable. In the model predicting kindergarten readiness, type of household is one independent variable included in the model where a two-parent home is coded as ‘1’ and a single-parent home as ‘0.’ An odds ratio of .002 indicates that the odds of being prepared for kindergarten (compared to unprepared for kindergarten) are decreased by a factor of .002 by being in a single-parent home (as opposed to living in a two-parent home). We could also state that the odds that a child from a single-parent home will be prepared for kindergarten are .998 (i.e., $1 - .002$).

In the case of a categorical variable with more than two categories, the odds ratio is interpreted relative to the reference (or left out) category. For example, say we have a predictor in our model that is mother’s education level with categories that include (1) less than high school diploma, (2) high school diploma or GED, and (3) at least some college. Say we set the last category (‘at least some college’) as the reference category. An odds ratio of .86 for the category of ‘high school diploma or GED’ for mother’s education level suggests that the odds of being prepared for kindergarten (as compared to unprepared) decrease by a factor of .86 when the child’s mother has a high school diploma or GED, relative to when the child’s mother has at least some college, when the other variables in the model are controlled.

Odds ratio values can also be converted to Cohen’s *d* using the following equation:

$$d = \frac{\ln(OR)}{1.81}$$

5.1.5 Assumptions

Compared to OLS regression, the assumptions of logistic regression are somewhat relaxed; however, four primary assumptions must still be considered: (a) noncollinearity,

(b) linearity, (c) independence of errors, and (d) fixed values of X . In this section, we also discuss conditions that are needed in logistic regression, as well as diagnostics, that can be performed to more closely examine the data.

5.1.5.1 Noncollinearity

Noncollinearity is applicable to logistic regression models with multiple predictors just as it was in multiple regression (but is not applicable when there is only one predictor in any regression model). This assumption has already been explained in detail in the multiple regression chapter and thus will not be reiterated other than to explain tools that can be used to detect multicollinearity. Although SPSS does not provide an option to easily generate collinearity statistics in logistic regression, you can generate an OLS regression model (i.e., a traditional multiple linear regression) with the same variables used in the logistic regression model and request collinearity statistics there. Because it is only the collinearity statistics that are of interest, do not be concerned in generating an OLS regression model that violates some OLS basic assumptions (e.g., normality). We have previously discussed tolerance and the variance inflation factor as two collinearity diagnostics (where tolerance is computed as $1 - R_k^2$ where R_k^2 is the variance in each independent variable, X , explained by the other independent variables and VIF is $\frac{1}{1 - R_k^2}$). In reviewing these statistics, tolerance values less than .20 suggest that multicollinearity exists, and values less than .10 suggest serious multicollinearity. VIF values greater than 10 indicate a violation of noncollinearity.

The effects of a violation of noncollinearity in logistic regression are the same as that in multiple regression. First, it will lead to instability of the regression coefficients across samples, where the estimates will bounce around quite a bit in terms of magnitude, and even occasionally result in changes in sign (perhaps opposite of expectation). This occurs because the standard errors of the regression coefficients become larger, thus making it more difficult to achieve statistical significance. Another result that may occur involves an overall regression that is significant, but none of the individual predictors is significant. Violation will also restrict the utility and generalizability of the estimated regression model.

5.1.5.2 Linearity

In OLS regression, the dependent variable is assumed to have a linear relationship with the continuous independent variable(s), but this does not hold in logistic regression. Because the outcome in logistic regression is a logit, the assumption of linearity in logistic regression refers to linearity between the *logit of the dependent variable* and the continuous independent variable(s). Hosmer and Lemeshow (1989) suggest several strategies for detecting nonlinearity, the easiest of which to apply is likely the Box-Tidwell transformation. This strategy is also valuable as it is not overly sensitive to minor violations of linearity. This involves generating a logistic regression model that includes all independent variables of interest along with an interaction term for

each—the interaction term being the product of the continuous independent variable and its natural log [i.e., $X * \ln(X)$]. Statistically significant interaction terms suggest nonlinearity. It is important to note that the assumption of linearity is applicable only for continuous predictors. A violation of linearity can result in biased parameters estimates, as well as the expected change in the logit of Y not being constant across the values of X . As noted later, the Hosmer-Lemeshow test has decreased power in detecting lack of fit in situations where linearity is violated (Xie et al., 2008).

5.1.5.3 Independence of Errors

Independence of errors is applicable to logistic regression models just as it was with OLS regression, and a violation of this assumption can result in underestimated standard errors (and thus overestimated test statistic values and perhaps finding statistical significance more often than is really viable, as well as affecting confidence intervals). This assumption has already been explained in detail during the discussion of multiple regression assumptions and thus additional information will not be provided here.

5.1.5.4 Fixed X

The last assumption is that the values of X_k are **fixed**, where the independent variables X_k are fixed variables rather than random variables. Because this assumption was discussed in detail in relation to multiple regression, we only summarize the main points. When X is fixed, the regression model is only valid for those particular values of X_k that were actually observed and used in the analysis. Thus, the same values of X_k would be used in replications or repeated samples. As discussed in the previous regression chapter, generally we may not want to make predictions about individuals having combinations of X_k scores outside of the range of values used in developing the prediction model; this is defined as *extrapolating* beyond the sample predictor data. On the other hand, we may not be quite as concerned in making predictions about individuals having combinations of X_k scores within the range of values used in developing the prediction model; this is defined as *interpolating* within the range of the sample predictor data. Table 5.3 summarizes the assumptions of logistic regression and the impact of their violation.

■ TABLE 5.3

Assumptions and Violation of Assumptions: Logistic Regression Analysis

Assumption	Effect of Assumption Violation
Noncollinearity of X 's	<ul style="list-style-type: none"> Regression coefficients can be quite unstable across samples (as standard errors are larger) Restricted generalizability of the model
Linearity	<ul style="list-style-type: none"> Bias in slopes and intercept Expected change in logit of Y is not a constant and depends on value of X
Independence	Influences standard errors of the model and thus hypothesis tests and confidence intervals
Values of X 's are fixed	<ul style="list-style-type: none"> Extrapolating beyond the range of X combinations: prediction errors larger, may also bias slopes and intercept Interpolating within the range of X combinations: smaller effects than when extrapolating; if other assumptions met, negligible effect

5.1.5.5 Conditions

Although not assumptions, the following conditions should be met with logistic regression: nonzero cell counts, nonseparation of data, lack of influential points, and sufficient sample size.

Nonzero Cell Counts

The first condition is related to nonzero cell counts in the case of nominal independent variables. A zero cell count occurs when the outcome is constant for one or more categories of a nominal variable (e.g., all females pass the course). This results in high standard errors because entire groups of individuals have odds of 0 or 1. Strategies to remove zero cell counts include recoding the categories (e.g., collapsing categories) or adding a constant to each cell of the crosstab table. If the overall model fit is what is of primary interest, then you may choose not to do anything about zero cell counts. The overall relationship between the set of predictors and the dependent variable is not generally impacted by zero cell counts. However, if zero cell counts are retained and the results of the individual predictors are what is of interest, it would be wise to provide a limitation to your results recognizing higher standard errors that are produced due to zero cell counts, as well as caution that the values of the individual regression coefficients may be affected. Careful review of the data prior to computing the logistic regression model can help thwart potential problems with zero cell counts.

Nonseparation of Data

Another condition that should be examined is that of complete or quasi-complete separation. Complete separation arises when the dependent variable is perfectly predicted and results in an inability to estimate the model. Quasi-complete separation occurs when there is less than complete separation and results in extremely large coefficients and standard errors. These conditions may occur when the number of variables equals (or nearly equals) the number of cases in the dataset, such that large coefficients and standard errors result.

Lack of Influential Points

Outliers and influential cases are problematic in logistic regression analysis just as with OLS regression. Severe outliers can cause the maximum likelihood estimator to reduce to zero (Croux, Flandre & Haesbroeck, 2002). Residual analysis and other diagnostic tests are equally beneficial for detecting miscoded data and unusual (and potentially influential) cases in logistic regression as they are in OLS regression. SPSS provides the option for saving a number of values including predicted values, residuals, and influence statistics. Both probabilities and group membership predicted values can be saved. Residuals that can be saved include (a) unstandardized, (b) logit, (c) studentized, (d) standardized, and (e) deviance. The three types of influence values that can be saved include Cook's, leverage values, and DFBETAs.

The wide variety of values that can be saved suggests that there are many types of diagnostics that can be performed. Review should be conducted when standardized or

studentized residuals are greater than an absolute value of 3.0 and DFBETA values are greater than one. Leverage values greater than $(m + 1)/N$ (where m equals the number of independent variables) indicate an influential case (values closer to 1 suggest problems, while those closer to 0 suggest little influence). If outliers or influential cases are found, it is up to you to decide if removal of the case is warranted. It may be that they, while uncommon, are completely plausible so that they are retained in the model. If they are removed from the model, it is important to report the number of cases that were removed prior to analysis (and evidence to suggest what caused you to remove them). A review of the multiple regression chapter provides further details on diagnostic analysis of outliers and influential cases.

5.2 MATHEMATICAL INTRODUCTION SNAPSHOT

To summarize the mathematics that underlie logistic regression, odds are simply the ratio of the probability of the dependent variable's two outcomes and computed as:

$$Odds(Y = 1) = \frac{P(Y = 1)}{1 - P(Y = 1)}$$

Changing the scale of the odds by taking the natural logarithm of the odds (aka *logit Y* or *log odds*) provides us with a value of the dependent variable that can theoretically range from negative infinity to positive infinity and thus creates a linear relationship between X and the probability of Y (Pampel, 2000). The natural log of the odds is calculated as follows:

$$\ln \frac{P(Y = 1)}{1 - P(Y = 1)} = Logit(Y)$$

and working with log odds, our familiar additive regression equation is applicable:

$$\ln \frac{P(Y = 1)}{1 - P(Y = 1)} = Logit(Y) = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m$$

If we exponentiate the logit (Y) (i.e., the outcome of our logistic regression equation), then it converts back to the odds (as noted by the calculation here) which allows us to interpret the independent variables as affecting the odds (rather than log odds) of the outcome:

$$Odds(Y = 1) = e^{\text{logit}(Y)} = e^{\ln[Odds(Y = 1)]} = e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m} = (e^\alpha)(e^{\beta_1 X_1})(e^{\beta_2 X_2}) \dots (e^{\beta_m X_m})$$

Converting the odds back to a probability can be done through the following formula:

$$P(Y = 1) = \frac{Odds(Y = 1)}{1 + Odds(Y = 1)} = \frac{e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m}}{1 + e^{\alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m}}$$

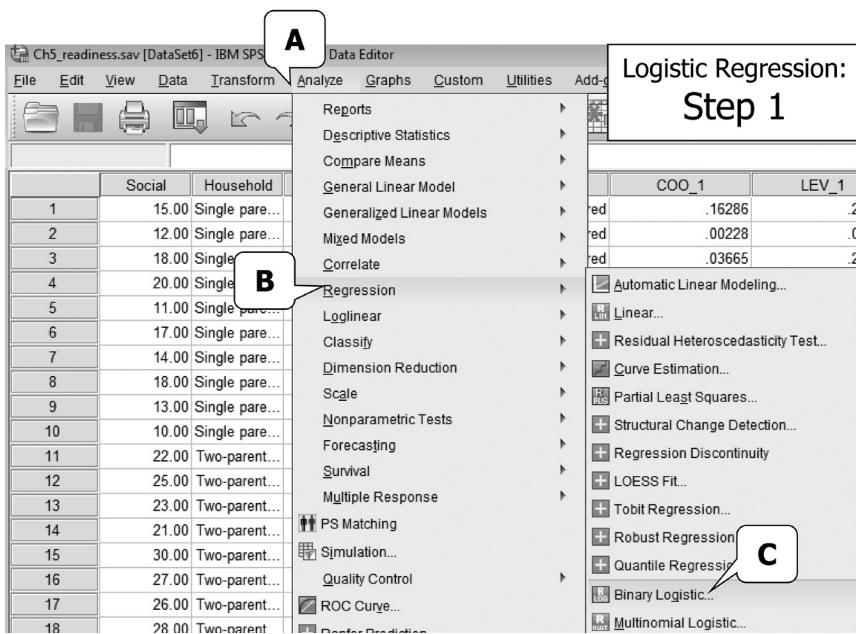
Probability values close to one indicate increased likelihood of occurrence.

5.3 COMPUTING LOGISTIC REGRESSION USING SPSS

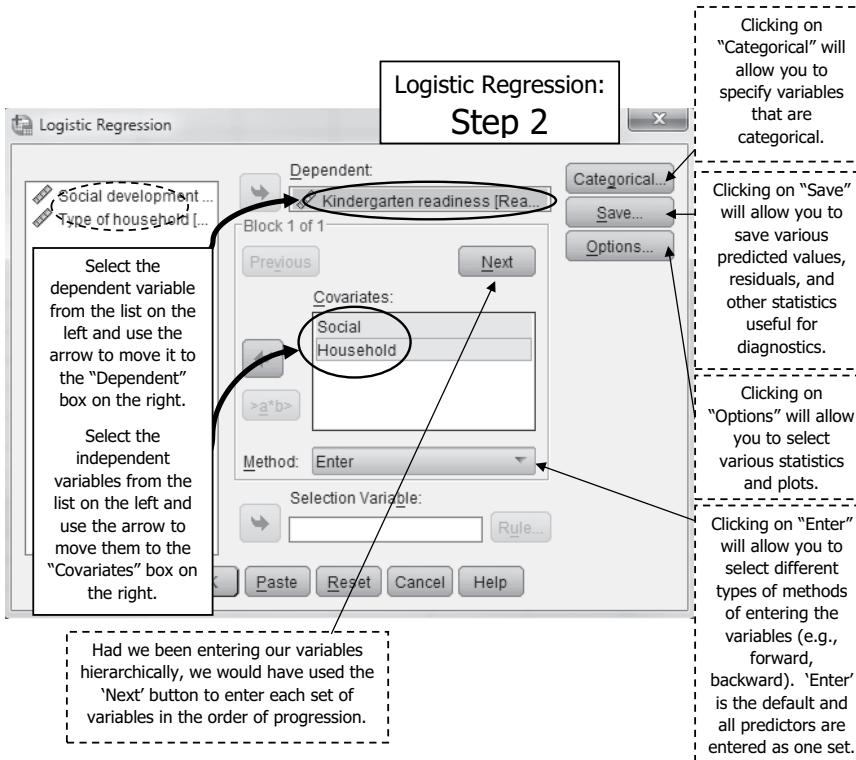
Next, we consider SPSS for the logistic regression model. Before we conduct the analysis, let us review the data (*Ch5.readiness.sav*) (note that we recognize the sample size of 20 does not meet minimum sample size criteria previously specified; however for illustrative purposes we felt it important that we be able to show the entire dataset and this would have been more difficult with the recommended sample size for logistic regression). With one dependent variable and two independent variables, the dataset must consist of three variables or columns, one for each independent variable and one for the dependent variable. Each row still represents one individual. As seen in the screenshot below, the SPSS data is in the form of three columns that represent the two independent variables (a continuous teacher-administered social development scale and household—a dichotomous variable, single- vs. two-adult household) and one binary dependent variable (kindergarten readiness screening test—prepared vs. not prepared). As our dependent variable is dichotomous, we will conduct binary logistic regression. When the dependent variable consists of more than two categories, multinomial logistic regression is appropriate (although not illustrated here).

	Social	Household	Readiness
1	15.00	.00	.00
2	12.00	.00	.00
3	18.00	.00	1.00
4	20.00	.00	1.00
5	11.00	.00	.00
6	17.00	.00	1.00
7	14.00	.00	.00
8	18.00	.00	1.00
9	13.00	.00	.00
10	10.00	.00	.00
11	22.00	1.00	.00
12	25.00	1.00	1.00
13	23.00	1.00	1.00
14	21.00	1.00	1.00
15	30.00	1.00	1.00
16	27.00	1.00	1.00
17	26.00	1.00	1.00
18	28.00	1.00	1.00
19	24.00	1.00	.00
20	30.00	1.00	1.00

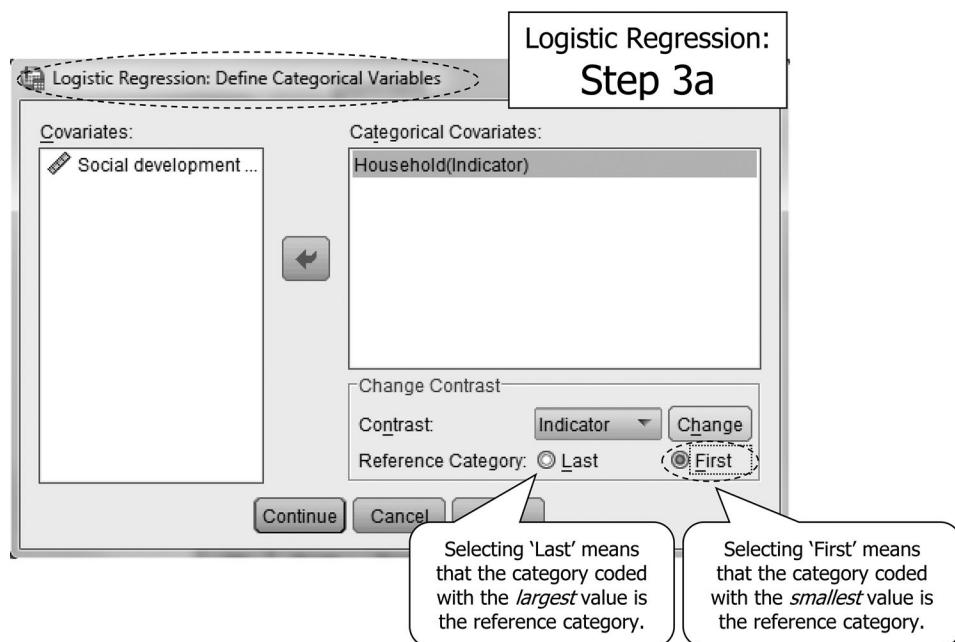
Step 1. To conduct a binary logistic regression, go to “Analyze” in the top pull-down menu, then select “Regression,” and then select “Binary Logistic.” Following the screenshot for Step 1 produces the “Logistic Regression” dialog box.



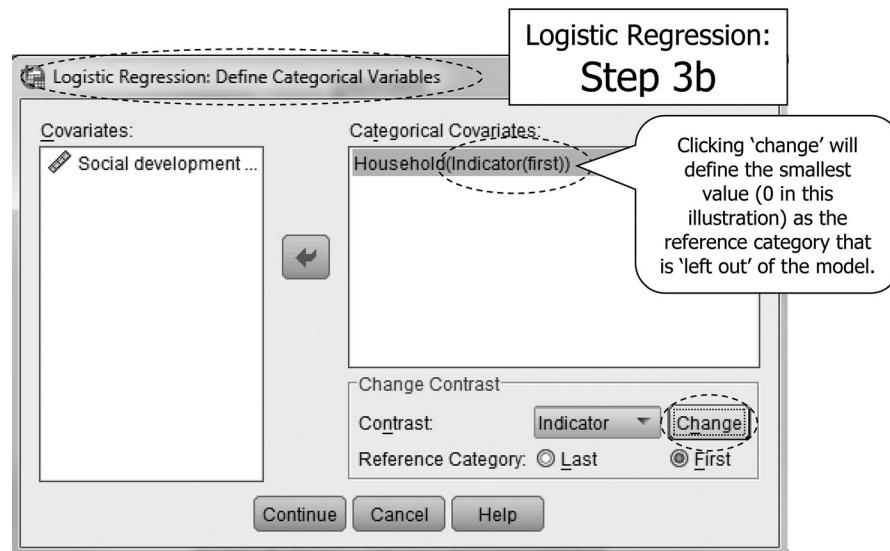
Step 2. Click the dependent variable (e.g., ‘Readiness’) and move it into the “Dependent” box by clicking the arrow button. Click the independent variables and move them into the “Covariate(s)” box by clicking the arrow button (see screenshot Step 2).



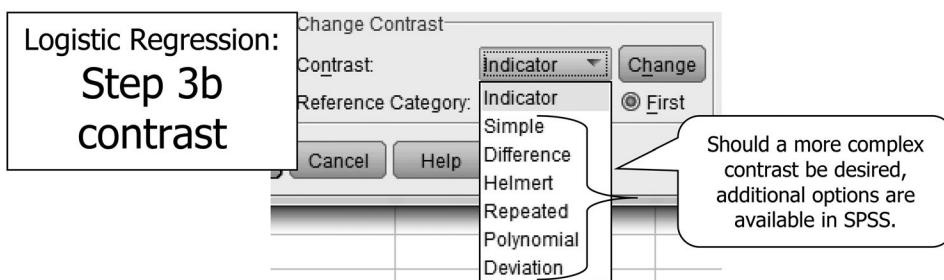
Step 3. From the Logistics Regression dialog box (see screenshot Step 2), clicking on “Categorical” will provide the option to define as categorical those variables that are nominal or ordinal in scale, as well as to select which category of the variable is the reference category through the Define Categorical Variables dialog box (see screenshot Step 3a). From the list of covariates on the left, click the categorical covariate(s) (e.g., ‘Household’) and move it into the “Categorical Covariates” box by clicking the arrow button. By default, ‘(Indicator)’ will appear next to the variable name. Indicator refers to traditional dummy coding and you have the option of selecting which value is the reference category. For binary variables (only two categories), using the ‘Last’ value as the reference category means that the category coded with the largest value will be the category ‘left out’ of the model (or referent), and using the ‘First’ value as the reference category means that the category coded with the smallest value will be the category ‘left out’ of the model. Here, two-parent households were coded as 1 and single-parent households as 0. We use single-parent households (coded as 0) as the reference category. Thus, we select the radio button for ‘First’ (see screenshot Step 3a) to define single-parent households as the reference category.



Next, we need to click the button labeled ‘Change’ (see screenshot Step 3b) to define the first value (i.e., zero or single-parent household) as the reference (or ‘left out’) category. By doing that, the name of our categorical covariate will now read Household(Indicator(first)). Had we had a categorical variable with more than two categories, we could just define the variable as categorical within logistic regression and select either the first or the last value as the reference category. If neither the first or last were what you wanted as the reference category, then some recoding of the data is necessary.

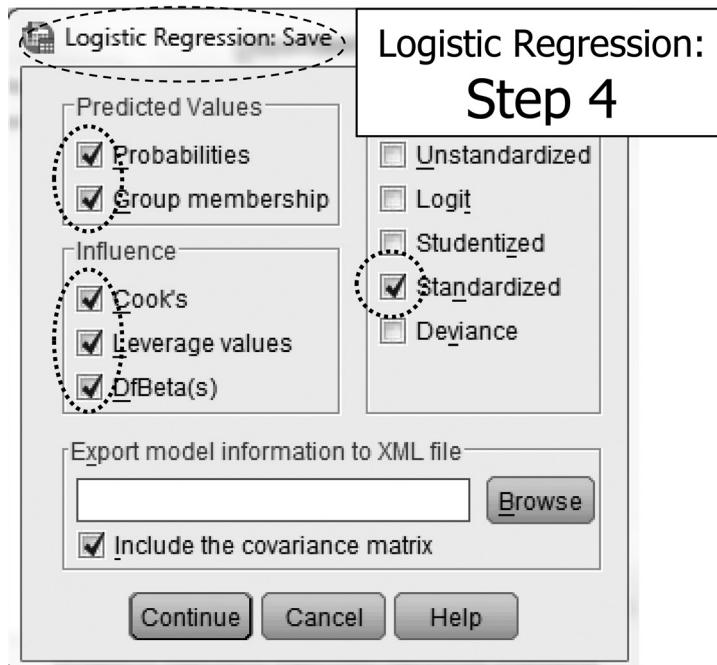


Before we move on, notice that the button for **Contrast** is a toggle menu with **Indicator** as the default option. Selecting the toggle menu allows you to select other types of contrasts often discussed in relation to ANOVA contrasts (e.g., Simple, Difference, Helmert) (see screenshot Step 3b contrast). These will not be reviewed here. Click on "Continue" to return to the Logistic Regression dialog box.

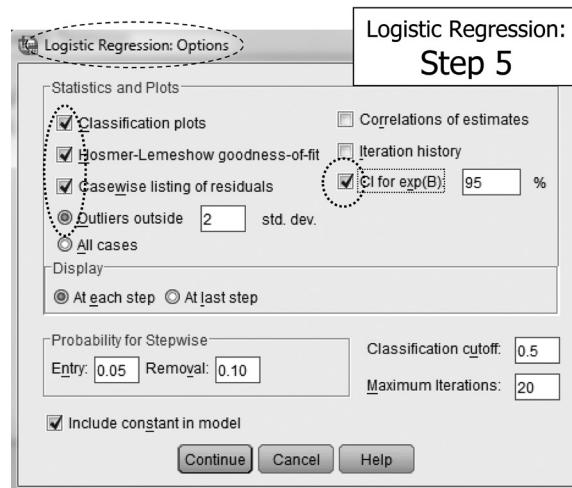


Step 4. From the Logistic Regression dialog box (see screenshot Step 2), clicking on "Save" will provide the option to save various predicted values, residuals, and statistics that can be used for diagnostic examination (see screenshot Step 4). From the Save dialog box under the heading of **Predicted Values**, place a checkmark in the box next to the following: (1) probabilities and (2) group membership. Under the heading of **Residuals**, place a checkmark in the box next to the following: (1) standardized.

Under the heading of **Influences**, place a checkmark in the box next to the following: (1) Cook's, (2) Leverage values, and (3) DFBETA(s). Click on "Continue" to return to the original dialog box.



Step 5. From the Logistic Regression dialog box (see screenshot Step 2), clicking on "Options" will allow you to generate various statistics and plots. From the Options dialog box (see screenshot Step 5) under the heading of **Statistics and Plots**, place a checkmark in the box next to the following: (1) Classification plots, (2) Hosmer-Lemeshow goodness-of-fit, (3) Casewise listing of residuals, (4) Outliers outside, and (5) CI for exp(B). For Outliers outside, you must specify a numeric value of standard deviations to define what you consider to be an outlier. Common values may be 2 (in a normal distribution, 95% of cases will be within ± 2 standard deviations), 3 (in a normal distribution, about 99% of cases will be within ± 3 standard deviations), or 3.29 (in a normal distribution, about 99.9% of cases will be within ± 3.29 standard deviations). For this illustration, we will use a value of 2. For CI for exp(B), you must specify a confidence interval. This should be the complement of the alpha being tested. If you are using an alpha of .05, then the CI will be $1 - .05$ or 95. All the remaining options in the "Options" dialog box will be left as the default settings. Click on "Continue" to return to the original dialog box. From the "Logistic Regression" dialog box, click on "OK" to generate the output.



Interpreting the output. Annotated results are presented in Table 5.4.

■ TABLE 5.4

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example

Case Processing Summary		
	N	Percent
Selected Cases Included in Analysis	20	100.0
Missing Cases	0	.0
Total	20	100.0
Unselected Cases	0	.0
Total	20	100.0

a. If weight is in effect, see classification table for the total number of cases.

This table provides information on sample size and missing data. The sample size is 20 and we have no missing data.

Dependent Variable Encoding	
Original Value	Internal Value
Unprepared	0
Prepared	1

Information on how the values of the dependent variable are coded is provided under 'internal value.' 'Unprepared' is coded as 0 and 'prepared' is coded as 1.

Categorical Variables Codings		
	Frequency	Parameter coding
		(1)
Type of household	Single parent household	10
	Two-parent household	10

Information on how the values of the categorical variable(s) are coded is provided as 'parameter coding.' 'Single parent household' is coded as 0 and 'two-parent household' is coded as 1. The sample size per group is presented in the 'frequency' column.

TABLE 5.4 (continued)

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example

Block 0: Beginning Block

Block 0 is a summary of the model with the constant only (i.e., none of the predictors are included). The **classification table** provides the percentage of cases correctly predicted given the constant only. Without including covariates, we can correctly predict children who are prepared for kindergarten 100% of the time but fail to predict any children (0%) who are unprepared. Here all children are predicted to be prepared.

Observed		Predicted		Percentage Correct	
		Kindergarten readiness			
		Unprepared	Prepared		
Step 0	Kindergarten readiness	Unprepared	0	8	
		Prepared	0	12	
	Overall Percentage			100.0 60.0	

a. Constant is included in the model.

b. The cut value is .500

Variables in the Equation

	B	S.E.	Wald	df	Sig.	Exp(B)
Step 0	Constant	.405	.456	.789	1	.374

Variables not in the Equation

		Score	df	Sig.
Step 0	Variables	8.860	1	.003
	Social development			
	Household(1)	3.333	1	.068
	Overall Statistics	11.168	2	.004

Variables not in the equation provides an indication of whether each covariate will statistically significantly contribute to predicting the outcome. Only social development ($p = .003$) is of value in the logistic model. The value of 11.168 for **overall statistics** is a residual chi-square statistic. Since the p value for it indicates statistical significance ($p = .004$), this indicates that including the two covariates improves the model as compared to the constant only model.

■ TABLE 5.4 (continued)

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example

Omnibus Tests of Model Coefficients			
	Chi-square	df	Sig.
Step 1	Step	15.793	.000
	Block	15.793	.000
	Model	15.793	.000

Model Summary			
	-2 Log likelihood	Cox & Snell R Square	Nagelkerke R Square
1	11.128 ^a	.546	.738

Hosmer and Lemeshow Test			
Step	Chi-square	df	Sig.
1	4.691	7	.698

Method = Enter indicates that the method of entering the predictors was simultaneous entry (recall this is the default method in SPSS and is called "Enter").

The -2LL for the constant only model is computed as the sum of chi-square for the constant only model and -2LL for the full model:

$$\chi^2_{\text{Model}} + -2\text{LL} = 15.793 + 11.128 = 26.921$$

Model summary statistics provide overall model fit. For good model fit, the value of -2LL for the full model (11.128) should be less than -2LL for the constant only model (26.921). This is a chi-square value with degrees of freedom equal to the number of parameters in the full model (i.e., 2 predictors plus one constant) minus the number of parameters in the baseline model (i.e., 1).

Thus there are 2 *df*. Using the chi-square table, with an alpha of .05 and 2 *df*, the critical value is 5.99. Since 11.128 is larger than the critical value, we reject the null hypothesis that the best prediction model is the constant only model. In other words, the full model (with predictors) is better at predicting kindergarten readiness than the constant only model.

a. Estimation terminated at iteration number 7 because parameter estimates changed by less than .001.

The two *R*² values are pseudo *R*² and are interpreted similarly to multiple *R*². These can be used as effect size indices for logistic regression and Cohen's interpretations for correlation can be used to interpret. Both values indicate a large effect.

As a measure of classification accuracy, non-statistical significance (*p* = .698) indicates good model fit for the Hosmer and Lemeshow test. This test is affected by small sample size, however; caution should be used when interpreting the results of this test when sample size is less than 50.

■ TABLE 5.4 (continued)

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example

Contingency Table for Hosmer and Lemeshow Test						
		Kindergarten readiness = Unprepared		Kindergarten readiness = Prepared		Total
		Observed	Expected	Observed	Expected	
Step 1	1	2	1.988	0	.012	2
	2	2	1.922	0	.078	2
	3	1	1.651	1	.349	2
	4	2	1.292	0	.708	2
	5	0	.607	2	1.393	2
	6	1	.404	2	2.596	3
	7	0	.100	2	1.900	2
	8	0	.030	2	1.970	2
	9	0	.005	3	2.995	3

The **classification table** provides information on how well group membership was predicted. Cells on the diagonal indicate correct classification. For example, children who were prepared for kindergarten were accurately classified 91.7% of the time as compared to unprepared children (87.5%). Overall, 90% of children were correctly classified. This is computed as the number of correctly classified cases divided by total sample size:

$$\frac{7+11}{20} = .90$$

		Classification Table ^a		Percentage Correct	
		Predicted			
		Kindergarten readiness			
		Unprepared	Prepared		
Step 1		7	1	87.5	
		1	11	91.7	
Overall Percentage				90.0	

a. The cut value is .500

Using Press's *Q* and given the chi-square critical value of 3.841 (*df* = 1), we find:

$$Q = \frac{[N - (nK)]^2}{N(K-1)} = \frac{[20 - (18)(2)]^2}{20(2-1)} = 12.8$$

We reject the null hypothesis. There is evidence to suggest that the predictions are statistically significantly better than chance.

■ TABLE 5.4 (continued)

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example

NOTE!
Interpretations of B coefficients are usually done via odds ratios.

The Wald statistic is used to test the statistical significance of each covariate.

The p value for 'social' ($p = .030$) indicates that the slope is statistically significantly different from zero. This tells us that the independent variable is contributing to predicting kindergarten preparedness. The intercept ($p = .032$) is also statistically significantly different from zero.

Exp(B) values are the odds ratios.
The odds ratio of 2.631 for social development indicates that the odds for being prepared for kindergarten are over 2-1/2 times greater (or 263%) for every one point increase in social development. The odds for household are nearly zero. This indicates that the odds for being prepared for kindergarten are about the same regardless of the child's household structure (single- versus two-parent home).

Variables in the Equation							
	B	S.E.	Wald	df	Sig.	Exp(B)	95% C.I. for EXP(B)
Step 1 ^a	.967	.446	4.696	1	.030	2.631	1.097 6.313
Social development							
Household(1)	-6.216	3.440	3.265	1	.071	.002	.000 1.693
Constant	-15.404	7.195	4.584	1	.032	.000	

a. Variable(s) entered on step 1: Social development, Household.

The B coefficient is interpreted as the change in the logit of the dependent variable given a one-unit change in the independent variable. Recall that the logit is the natural log of the dependent variable occurring. With B equal to .967, this tells us that a one-unit change in social development will result in nearly a one-unit change in the logit of kindergarten preparedness. The constant is the expected value of the logit of kindergarten readiness for children of single parents (recall this was coded as 0) and when social development is zero.

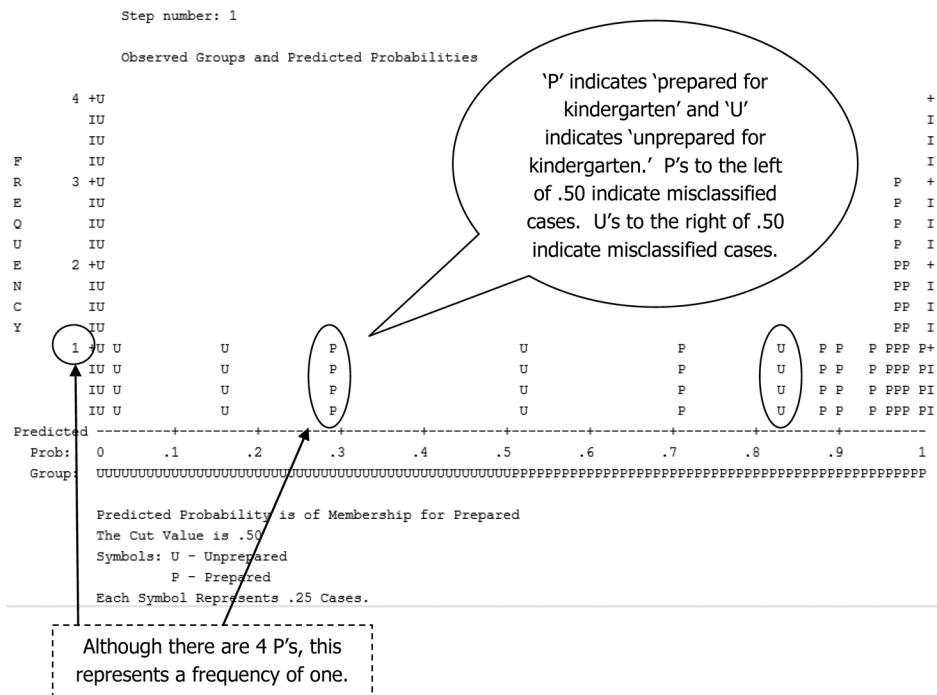
Since the odds of 1.00 (which indicates similar odds for falling into either category of the outcome) are not contained within the interval for social development, this suggests the odds ratio is statistically significantly different from zero. Note that the odds ratio is only computed for the predictors and not for the intercept (i.e., constant).

A negative B indicates that an increase in value of that independent variable will result in a *decrease* in the predicted probability of the dependent variable.

A positive B indicates that an increase in value of that independent variable will result in an *increase* in the predicted probability of the dependent variable.

TABLE 5.4 (continued)

SPSS Results for the Binary Logistic Regression Kindergarten Readiness Example



Case	Selected Status ^a	Observed		Predicted Group	Temporary Variable		
		Kindergarten readiness			Resid	ZResid	
		U	P				
8	S	U**		P	-.832	-2.226	
15	S	P**	.832	U	.786	1.918	

a. S = Selected, U = Unselected cases, and ** = Misclassified cases.

b. Cases with studentized residuals greater than 2.000 are listed

Recall we told SPSS to identify residuals that were outside 2 standard deviations. Based on that decision, cases 8 and 15 were identified as potential outliers. We review this output in the discussion on outliers.

5.4 DATA SCREENING

Previously we described a number of assumptions used in logistic regression. These included (a) noncollinearity, (b) linearity between the predictors and logit of the dependent variable, and (c) independence of errors. We also review the data to ensure there are no outliers.

Before we begin to examine assumptions, let us review the values that we requested to be saved to our data file (see dataset screenshot that follows).

1. **PRE_1** are the predicted probabilities.
2. **PGR_1** is the predicted group membership (here group membership is either prepared or unprepared for kindergarten).
3. **COO_1** are Cook's influence statistics. As a rule of thumb, Cook's values greater than one suggest that case is potentially problematic.
4. **LEV_1** are leverage values. As a general guide, leverage values less than .20 suggest there are no problems with cases exerting undue influence. Values greater than .5 indicate problems.
5. **ZRE_1** are standardized residuals computed as the residual divided by an estimate of the standard deviation of the residual. Standardized residuals have a mean of zero and standard deviation of one.
6. **DFB0_1, DFB1_1, and DFB2_1** are DFBETA values and indicate the difference in a beta coefficient if that particular case were excluded from the model.

As we look at the raw data, we see eight new variables have been added to our dataset. These are predicted values, residuals, and other diagnostic statistics.

	Social	Household	Readiness	PRE_1	PGR_1	COO_1	LEV_1	ZRE_1	DFB0_1	DFB1_1	DFB2_1
1	15.00	.00	.00	.29087	.00	.16286	.28420	-.64046	-.168367	.07492	-.02694
2	12.00	.00	.00	.02202	.00	.00228	.09212	-.15005	-.33145	.01897	-.10172
3	18.00	.00	.00	.88198	1.00	.03665	.21502	.36580	-.80889	.06219	-.61089
4	20.00	.00	1.00	.98104	1.00	.00177	.08403	.13902	-.24278	.01681	-.14108
5	11.00	.00	.00	.00848	.00	.00046	.05082	-.09250	-.15052	.00877	-.04979
6	17.00	.00	1.00	.73959	1.00	.13483	.27690	.59338	-.79435	.07718	-.96766
7	14.00	.00	.00	.13486	.00	.04579	.22703	-.39482	-.127676	.06695	-.25156
8	18.00	.00	1.00	.88198	1.00	.03665	.21502	.36580	-.80889	.06219	-.61089
9	13.00	.00	.00	.05593	.00	.01077	.15379	-.24340	-.69346	.03854	-.18626
10	10.00	.00	.00	.00324	.00	.00009	.02651	-.05702	-.06664	.00393	-.02313
				41706	.00	.31732	.30726	-.84584	-.158416	.09948	-.136519
				92875	1.00	.01215	.13675	.27698	-.56887	.03572	-.15362
				65309	1.00	.18337	.25662	.72883	-.25348	.01592	4.1597
				21377	.00	1.58721	.30146	1.91780	6.53464	-.41034	4.10130
				99939	1.00	.00000	.00691	.02466	-.01393	.00087	-.00535
				98904	1.00	.00058	.04980	.10526	-.15271	.00959	-.05321
				97167	1.00	.00275	.08620	.17075	-.31209	.01960	-.10037
18	28.00	1.00	1.00	.99581	1.00	.00012	.02696	.06489	-.07074	.00444	-.02581
19	24.00	1.00	.00	.83204	1.00	.120520	.19569	-.222568	3.84582	-.24150	.50163
20	30.00	1.00	1.00	.99939	1.00	.00000	.00691	.02466	-.01393	.00087	-.00535

5.4.1 Noncollinearity

It is not possible to request multicollinearity statistics, such as tolerance and VIF, using logistic regression in SPSS. We can, however, estimate those values by running the same variables in a multiple regression model and requesting only the collinearity statistics. We are not interested in the parameter estimates of the model—only the collinearity statistics. Tolerance values less than .10 and VIF values greater than 10 indicate multicollinearity (Menard, 1995). Because the steps for generating multiple regression were presented previously in the text, we will not reiterate them here. Rather, we will merely present the applicable portion of the output of this model. From the output that follows with a tolerance of .248 and VIF of 4.037, we have evidence that we do not have multicollinearity. In examining collinearity diagnostics, a general recommendation for interpreting condition indices is that values in the range of 10 to 30 should be of concern, greater than 30 indicates trouble, and greater than 100 indicates disaster (Belsley, 1991). Here the condition index of dimension three (14.259) is within the range of cause for concern. The last three columns refer to variance proportions. Multiplying these values by 100 provides a percentage of the variance of the regression coefficient that is related to a particular eigenvalue. Multicollinearity is suggested when covariates have high percentages associated with a small eigenvalue (and large condition index). Thus, for purposes of reviewing for multicollinearity, concentrate only on the rows with small eigenvalues. In this example, 100% of the variance of the regression coefficient for social development and 73% for type of household are related to eigenvalue 3 (the dimension with the smallest eigenvalue and largest condition index). This suggests some concern for multicollinearity. In summary, we have met the assumption of noncollinearity with the tolerance and VIF values, but there is some concern for multicollinearity with the condition index and variance proportion values.

Coefficients ^a			
Model		Collinearity Statistics	
		Tolerance	VIF
1	Social development	.248	4.037
	Type of household	.248	4.037

a. Dependent Variable: Kindergarten readiness

Collinearity Diagnostics^a

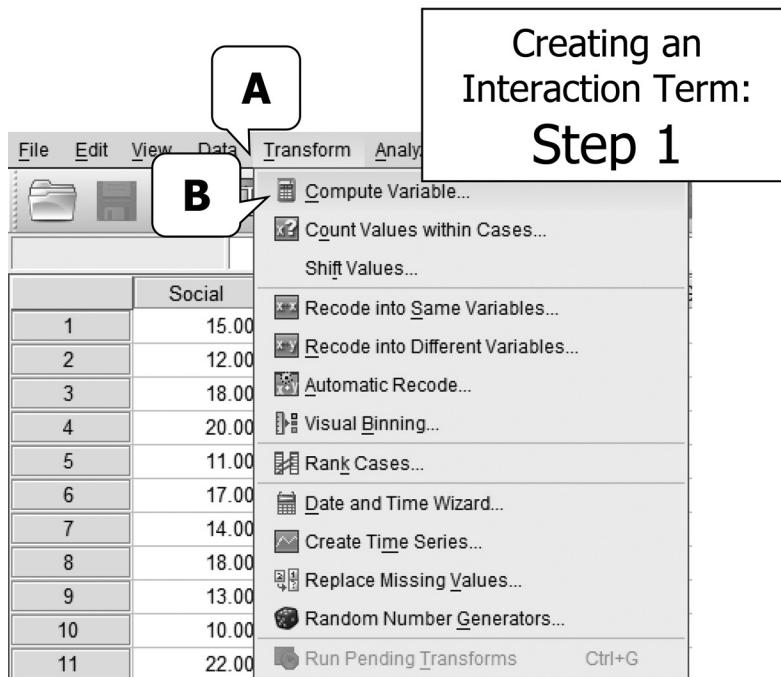
Model	Dimension	Eigenvalue	Condition Index	Variance Proportions		
				(Constant)	Social development	Type of household
1	1	2.683	1.000	.00	.00	.01
	2	.303	2.974	.05	.00	.25
	3	.013	14.259	.95	1.00	.73

a. Dependent Variable: Kindergarten readiness

5.4.2 Linearity

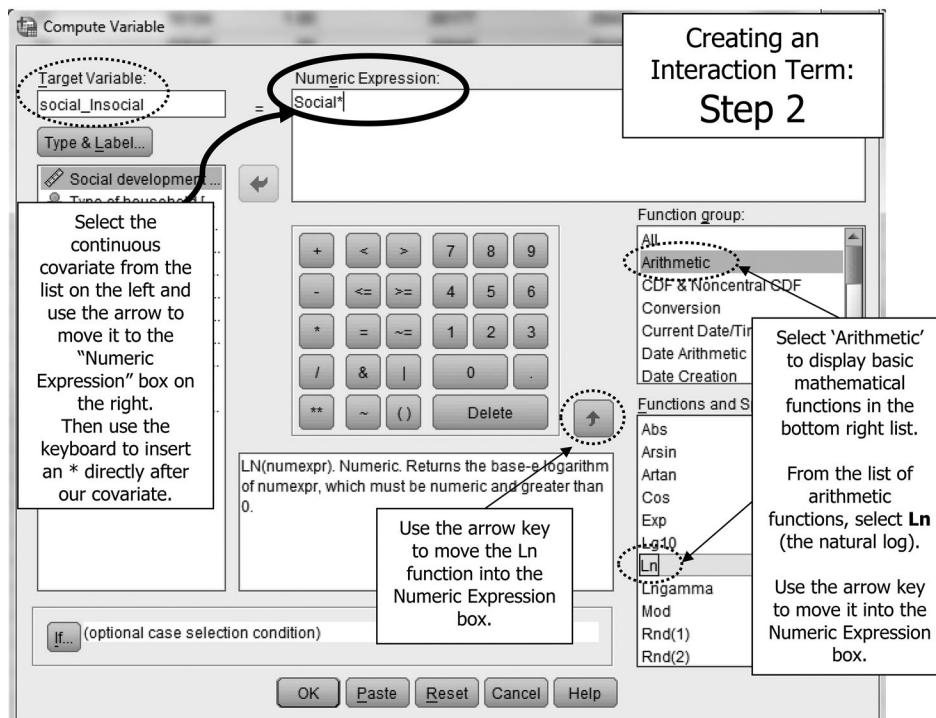
Recall that the linearity assumption is applicable only to continuous variables. Thus, we will test this assumption only for social development. The Tidwell transformation test can be used to test that the assumption of linearity has been met. To generate this test, for each *continuous* independent variable we must first create an interaction term that is the product of the independent variable and its natural log (*ln*). Here we have only one continuous independent variable—social development. Thus, only one interaction term will be created.

Step 1. To create an interaction term of our continuous variable and the natural log of this variable, go to “Transform” in the top pull-down menu, then select “Compute Variable.” Following the screenshot below (see screenshot Creating an Interaction Term: Step 1) produces the “Compute Variable” dialog box.

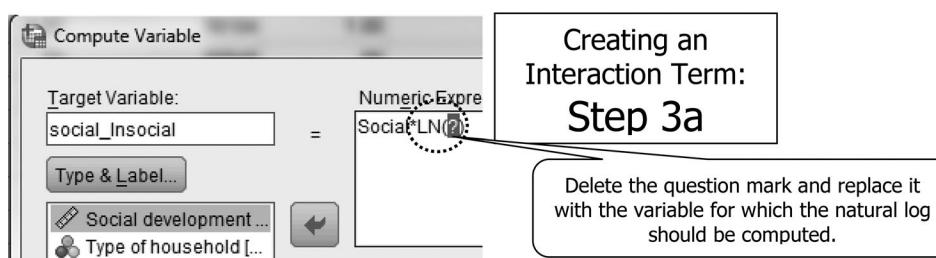


Step 2. In the Target Variable box in the upper left corner, enter the variable name that you want to appear as the column header (see screenshot Creating an Interaction Term: Step 2). Since this is the column header name, this name cannot begin with special characters or numbers and cannot have any spaces. If you wish to define the label for this variable (i.e., what will appear on the output; this *can* include special characters, spaces, and numbers), then click on the Type & Label box directly underneath Target Variable where additional text to define the name of the variable can be included. Next, click on the continuous covariate (i.e., social development) and move it into the Numeric Expression box by clicking on the arrow in the middle of the screen. Using

either the keyboard on screen or your keyboard, click on the asterisk key (i.e., *). This will be used as the multiplication sign. Next, under Function group, click on arithmetic to display all of the basic mathematical functions. From this alphabetized list click on Ln (natural log). To move this function into the Numeric Expression box, click on the arrow key in the right central part of the dialog box.

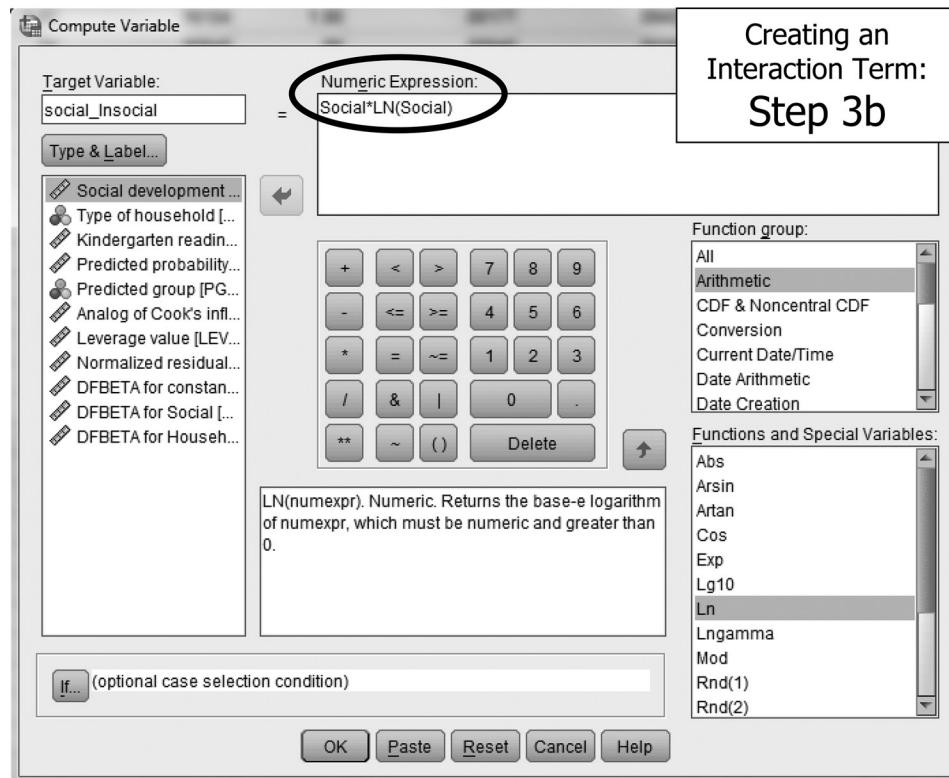


Step 3. Once the natural log function is displayed in the Numeric Expression box, a question mark enclosed inside parentheses will appear (see screenshot Creating an Interaction Term: Step 3a). This is SPSS's way of asking which variable you want the natural log computed for. Here it is the continuous covariate, social development.



Here we want to compute the natural log for the continuous covariate, social development. To move this variable into the parentheses, use the backspace or delete key

to remove the question mark. Then, click on the continuous covariate, social development, and move it into the parentheses next to LN in the Numeric Expression box by clicking on the arrow in the middle of the screen (see screenshot Creating an Interaction Term: Step 3b). The numeric expression should then read: Social*LN(Social). Click OK to compute and create the new variable in the dataset.



Step 4. The next step is to include the newly created variable (i.e., the interaction of the continuous variable with its natural log) into the logistic regression model, along with the other predictors. As those steps have been presented previously, they will not be reiterated here. The output indicates that the interaction term is not statistically significant ($p = .300$), which suggests we have met the assumption of linearity.

Variables in the Equation

		B	S.E.	Wald	df	Sig.	Exp(B)	95% C.I. for EXP(B)	
								Lower	Upper
Step 1 ^a	Social	12.953	11.897	1.185	1	.276	421981.259	.000	5.647E15
	Household(1)	-8.208	5.264	2.432	1	.119	.000	.000	8.236
	social_Insocial	-2.948	2.845	1.074	1	.300	.052	.000	13.845
	Constant	-76.228	64.345	1.403	1	.236	.000		

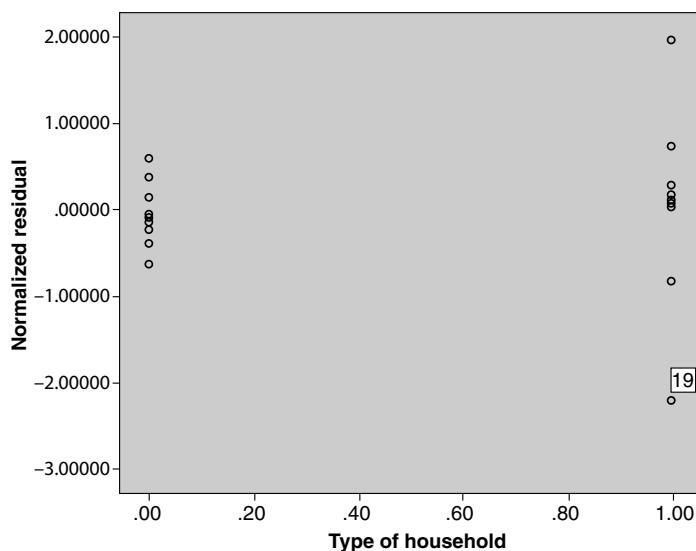
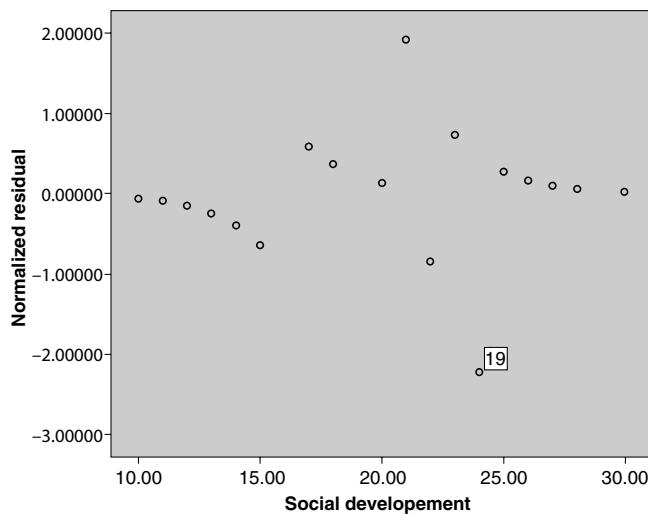
a. Variable(s) entered on Step 1: Social development, Household.

5.4.3 Independence

We plot the standardized residuals (which were requested and created through the 'Save' option) against the values of X to examine the extent to which independence was met. If you need a refresher on generating this plot, please review the data screening chapter.

5.4.3.1 Interpreting Independence Evidence

If the assumption of independence is met, the points should fall randomly within a band of -2.0 to $+2.0$. Here we have pretty good evidence of independence, especially given the small sample size relative to logistic regression, as all but one point (case 19) is within an absolute value of 2.0.



5.4.4 Absence of Outliers

Just as we saw in multiple regression, there are a number of diagnostics that can be used to examine the data for outliers.

5.4.4.1 Cook's Distance

Cook's distance provides an overall measure for the influence of individual cases. Values greater than one suggest that a case may be problematic in terms of undue influence on the model. Examining the residual statistics provided in the binary logistic regression output (see following table), we see that the maximum value for Cook's distance is 1.58, which indicates at least one influential point.

5.4.4.2 Leverage Values

These values range from 0 to 1, with values close to 1 indicating greater leverage. As a general rule, leverage values greater than $(m + 1)/n$ (where m equals the number of independent variables; here $(2 + 1)/20 = .15$ indicates an influential case. With a maximum of .307, there is evidence to suggest one or more cases are exerting leverage.

5.4.4.3 DFBETA

We saved the DFBETA values as another indication of the influence of a case. The DFBETA provide information on the change in the predicted value when the case is deleted from the model. For logistic regression, the DFBETA values should be smaller than one. Looking at the minimum and maximum DFBETA values for the intercept (labeled 'constant') and for household, we have at least one case that is suggestive of undue influence.

	N	Minimum	Maximum
Analog of Cook's influence statistics	20	.00000	1.58721
Leverage value	20	.00691	.30726
Normalized residual	20	-2.22568	1.91780
DFBETA for constant	20	-1.68367	6.53464
DFBETA for Social	20	-.41034	.09948
DFBETA for Household(1)	20	-1.36519	4.10130
Valid N (listwise)	20		

From our logistic regression output, we can review the Casewise List to determine cases with studentized residuals larger than two standard deviations (recall from the Options dialog box that we told SPSS to identify residuals outside two standard deviations). Here there were two cases (cases 8 and 15) that were identified as outliers and the relevant statistics (e.g., observed group, predicted value, predicted group, residual, and standardized residual) are provided. We examine these cases to make sure there

was not a data entry error. If the data are correct, then we determine whether to keep or filter out the case(s).

Casewise List^b

Case	Selected Status ^a	Observed Kindergarten readiness	Predicted	Predicted Group	Temporary Variable	
					Resid	ZResid
8	S	U**	.832	P	-.832	-2.226
15	S	P**	.214	U	.786	1.918

a. S = Selected, U = Unselected cases, and ** = Misclassified cases.

b. Cases with studentized residuals greater than 2.000 are listed.

Since we have a small dataset, we can easily review the values of our diagnostics and see which cases are problematic in terms of exerting undue influence and/or outliers. Those that are circled are values that fall outside of the recommended guidelines and thus are suggestive of outlying or influential cases. Due to the already small sample size, we will not filter out any of these potentially problematic cases. However, in this situation (i.e., with diagnostics that suggest one or more influential cases), you may want to consider filtering out those cases or, at a minimum, reviewing the data to be sure that there was not a data entry error for that case.

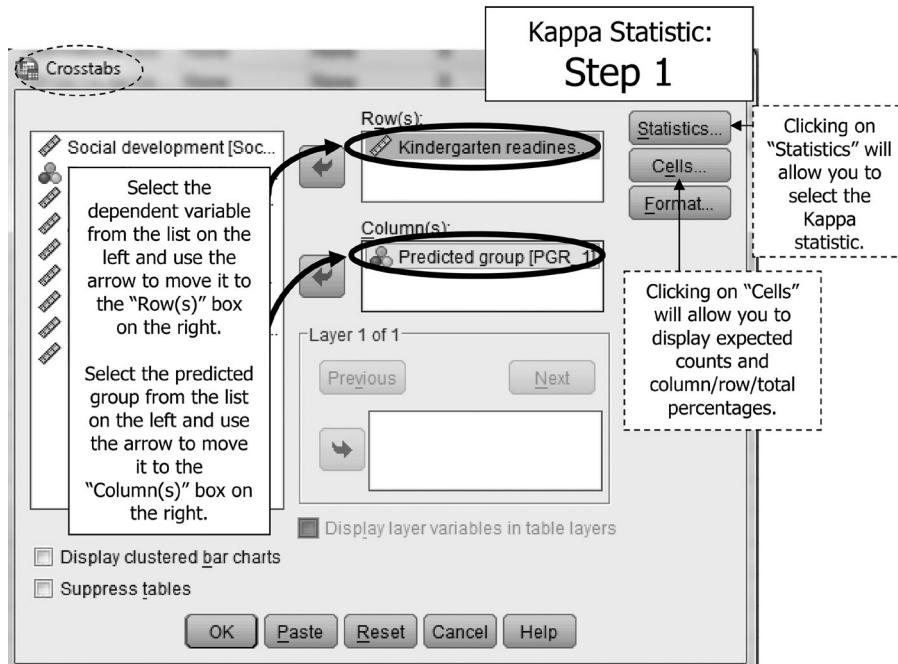
	Social	Household	Readiness	PRE_1	PGR_1	COO_1	LEV_1	ZRE_1	DFB0_1	DFB1_1	DFB2_1
1	15.00	.00	.00	.29087	.00	.16286	.28420	-.64046	-1.68367	.07492	-.02664
2	12.00	.00	.00	.02202	.00	.00228	.09212	-.15005	.33145	.01897	-.10172
3	18.00	.00	1.00	.88198	1.00	.03665	.21502	.36580	-.80889	.06219	-.61089
4	20.00	.00	1.00	.98104	1.00	.00177	.08403	.13902	-.24278	.01681	-.14108
5	11.00	.00	.00	.00848	.00	.00046	.05082	-.09250	-.15052	.00877	-.04979
6	17.00	.00	1.00	.73959	1.00	.13483	.27690	.59338	-.79435	.07718	-.96766
7	14.00	.00	.00	.13486	.00	.04579	.22703	-.39482	-1.27676	.06695	-.25156
8	18.00	.00	1.00	.88198	1.00	.03665	.21502	.36580	-.80889	.06219	-.61089
9	13.00	.00	.00	.05593	.00	.01077	.15379	-.24340	-.69346	.03854	-.18626
10	10.00	.00	.00	.00324	.00	.00009	.02651	-.05702	-.06664	.00393	-.02313
11	22.00	1.00	.00	.41706	.00	.31732	.30726	-.84584	-1.58416	.09948	-1.36519
12	25.00	1.00	1.00	.92875	1.00	.01215	.13675	.27698	-.56887	.03572	-.15362
13	23.00	1.00	1.00	.65309	1.00	.18337	.25662	.72883	-.25348	.01592	.41597
14	21.00	1.00	1.00	.21377	.00	1.58721	.30146	.191780	6.53464	-.41034	4.10130
15	30.00	1.00	1.00	.99939	1.00	.00000	.00691	.02466	-.01393	.00087	-.00535
16	27.00	1.00	1.00	.98904	1.00	.00058	.04980	.10526	-.15271	.00959	-.05321
17	26.00	1.00	1.00	.97167	1.00	.00275	.08620	.17075	-.31209	.01960	-.10037
18	28.00	1.00	1.00	.99581	1.00	.00012	.02696	.06489	-.07074	.00444	-.02581
19	24.00	1.00	.00	.83204	1.00	1.20520	.19569	-.222568	3.84582	-.24150	.50163
20	30.00	1.00	1.00	.99939	1.00	.00000	.00091	.02466	-.01393	.00087	-.00535

5.4.5 Assessing Classification Accuracy

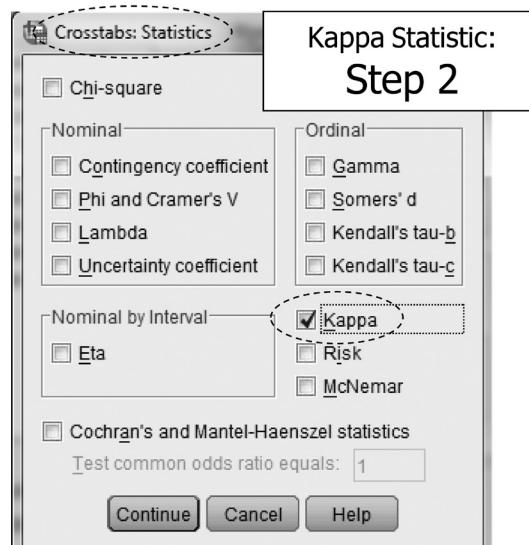
In addition to examining Press's Q for classification accuracy, we can generate a kappa statistic. Kappa is the proportion of agreement above that expected by chance. A kappa statistic of 1.0 indicates perfect agreement, whereas a kappa of 0 indicates chance agreement. Negative values can occur and indicate weaker than chance agreement. General rules of interpretation for kappa are small $< .30$, moderate $.30$ to $.50$, and large $> .50$.

Step 1. Kappa statistics are generated through the Crosstab procedure (go to "Analyze" in the top pull-down menu, then "Descriptive Statistics," and then "Crosstabs"). Once the

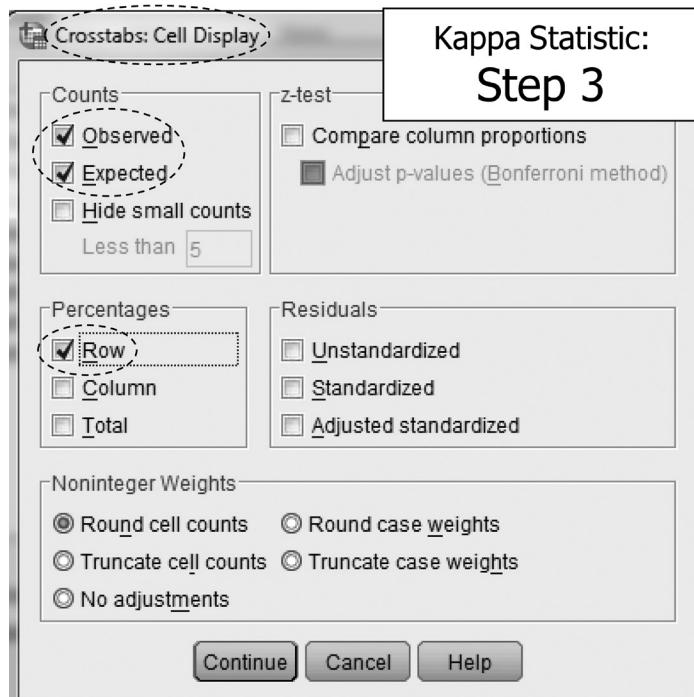
Crosstabs dialog box is open, select the dependent variable from the list on the left and use the arrow key to move it to Row(s). Select the predicted group (PGR_1) from the list on the left and use the arrow key to move it to Column(s) (see screenshot Kappa Statistic: Step 1).



Step 2. Click on the Statistics option button. Place a checkmark in the box next to Kappa (see screenshot Kappa Statistic: Step 2). Then click on Continue to return to the main dialog box.



Step 3. Click on the Cells option button. In the Cell Display dialog box, place a checkmark in the box next to observed, expected, and row (see screenshot Kappa Statistic: Step 3). Then click on Continue to return to the main dialog box. Then click OK to generate the output.



The crosstab table is interpreted as we have seen in the past. The columns represent the predicted group membership and the rows represent the observed group membership. This table should look familiar to the one that was provided to us with the logistic regression results.

Kindergarten readiness * Predicted group Crosstabulation

		Predicted group		Total	
		Unprepared	Prepared		
Kindergarten readiness	Unprepared	Count	7	8	
		Expected Count	3.2	8.0	
		% within Kindergarten readiness	87.5%	100.0%	
	Prepared	Count	1	12	
		Expected Count	4.8	12.0	
		% within Kindergarten readiness	8.3%	91.7%	
Total		Count	8	20	
		Expected Count	8.0	20.0	
		% within Kindergarten readiness	40.0%	100.0%	

What is of most interest is the table labeled Symmetric Measures, as this table contains the Kappa statistic. With a Kappa statistic of .792, and using our rules for interpretation, this is considered to be a large value, which suggests strong agreement.

Symmetric Measures

	Value	Asymp. Std. Error ^a	Approx. T ^b	Approx. Sig.
Measure of Agreement Kappa	.792		.140	.3540
N of Valid Cases	20			.000

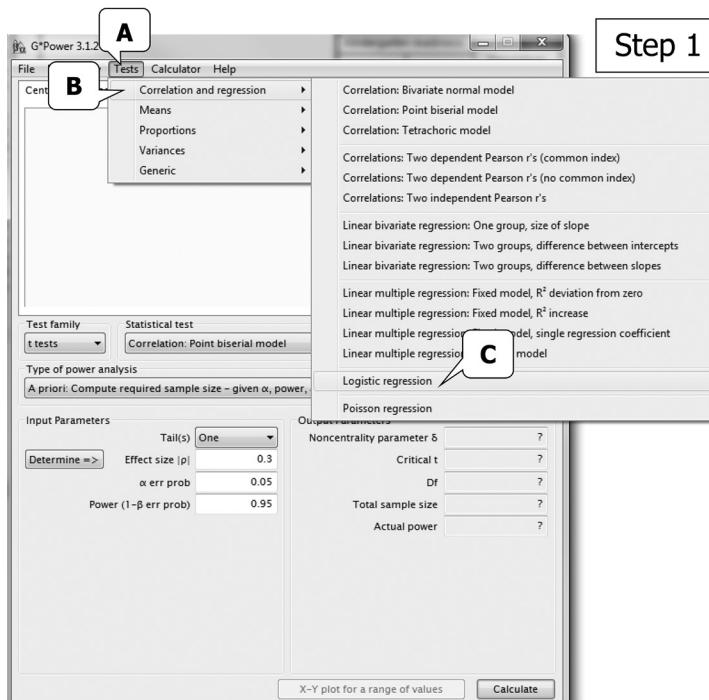
- a. Not assuming the null hypothesis.
 b. Using the asymptotic standard error assuming the null hypothesis.

5.5 POWER USING G*POWER

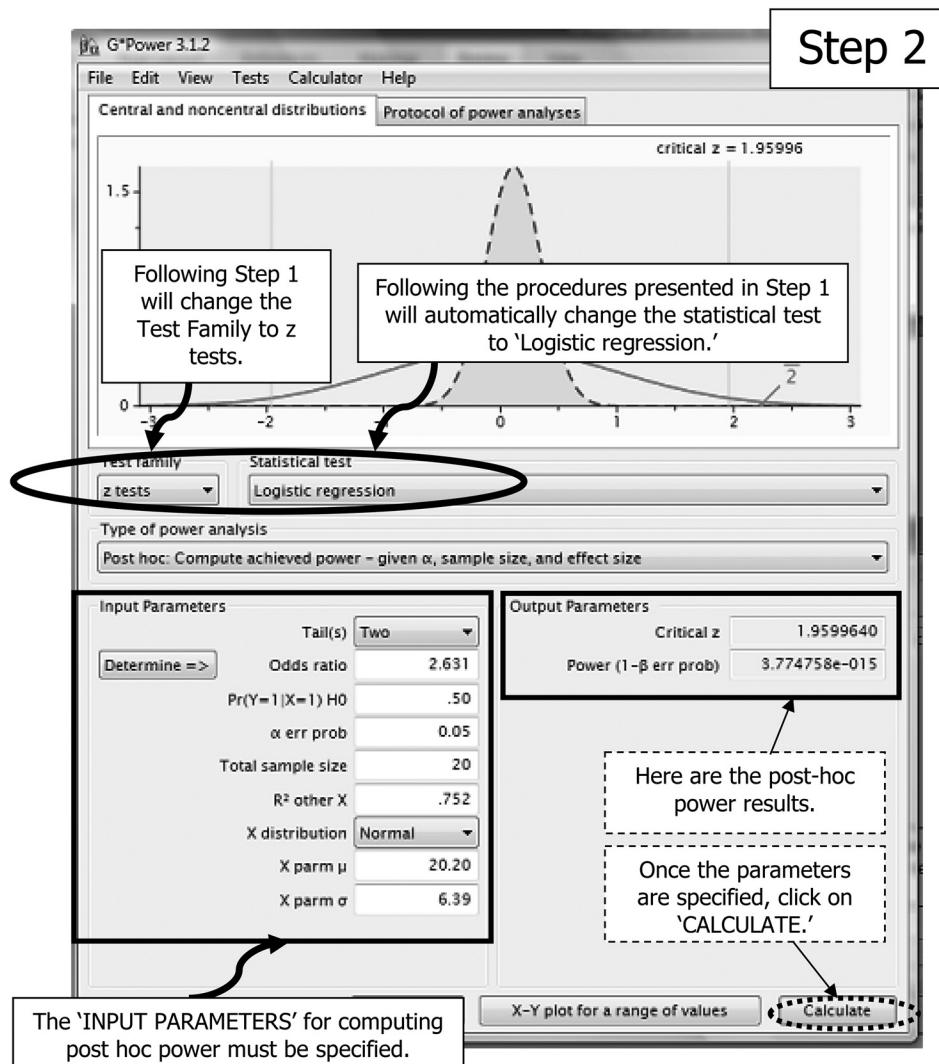
A priori and post hoc power can again be determined using the specialized software described previously in this text (e.g., G*Power), or you can consult a priori power tables (e.g., Cohen, 1988). As an illustration, we use G*Power to first compute post hoc power of our example.

5.5.1 Post Hoc Power for Logistic Regression Using G*Power

The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family. For logistic regression, we select 'Tests' in the top pull-down menu, then 'Correlation and regression,' and finally 'Logistic regression.' Once that selection is made, the 'Test family' automatically changes to 'z tests.'



The 'Type of power analysis' desired then needs to be selected. To compute post hoc power, select 'Post hoc: Compute achieved power—given α , sample size, and effect size.' For this illustration, we will compute power for the continuous covariate.



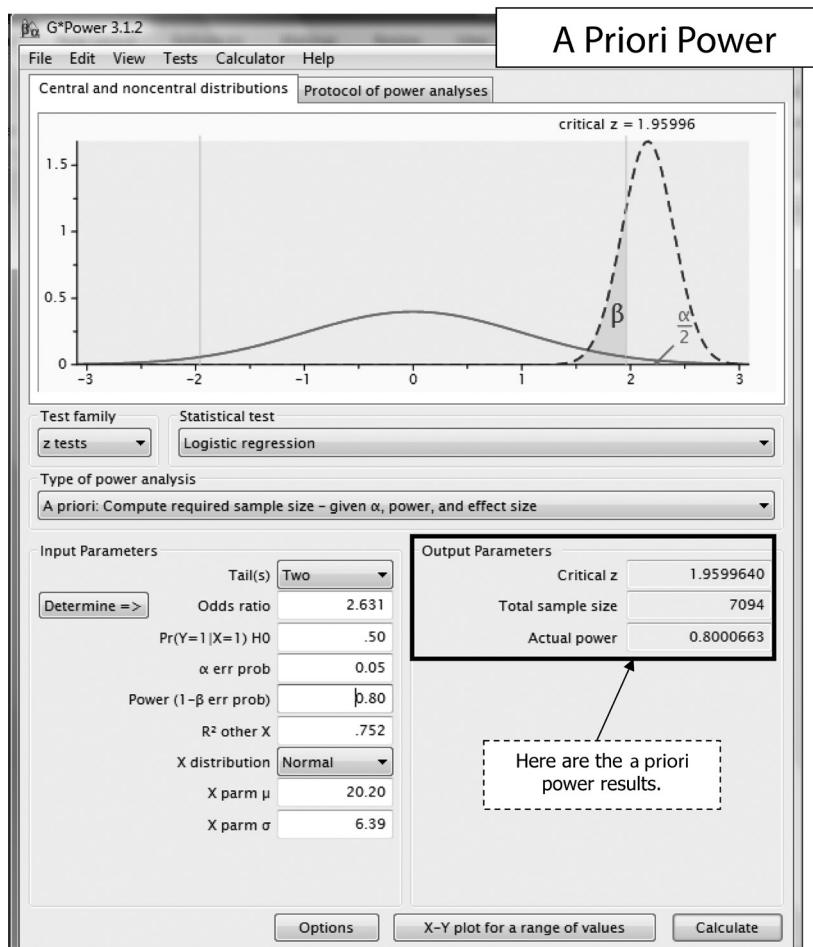
The 'INPUT PARAMETERS' must then be specified. In our example, we conducted a two-tailed test. The odds ratio for our continuous variable social development was 2.631. The probability that $Y = 1$ given that $X = 1$ under the null hypothesis is set to .50. The alpha level we used was .05 and the total sample size was 20. R^2 other X refers to the squared correlation between social development and our other covariate. In this case, the simple bivariate correlation between these variables is .867 and the squared correlation is .752. Social development is a continuous variable, thus it follows a normal distribution. The last two parameters to be specified are for the mean and standard deviation of our covariate. In this case, the mean of social development was 20.20 and

the standard deviation was 6.39. Once the parameters are specified, click on 'CALCULATE' to find the power statistics.

The 'OUTPUT PARAMETERS' provide the relevant statistics for the input just specified. In this example, we were interested in determining post hoc power for a logistic regression model. Based on the criteria specified, the post hoc power was substantially less than 1. In other words, the probability of rejecting the null hypothesis when it is really false was significantly less than 1% (sufficient power is often .80 or above). This finding is not surprising given the very small sample size. Keep in mind that conducting power analysis *a priori* is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters).

5.5.2 A Priori Power for Logistic Regression Using G*Power

For a *priori* power, we can determine the total sample size needed for logistic regression given the same parameters just discussed. In this example, had we wanted an *a priori* power of .80 given the same parameters just defined, we would need a total sample size of 7094.



5.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example paragraph for the results of the logistic regression analysis. Recall that our graduate research assistant, Oso Wyse, was assisting Dr. Malani, a faculty member in the early childhood department. Dr. Malani wanted to know if kindergarten readiness (prepared vs. unprepared) could be predicted by social development (a continuous variable) and type of household (single- vs. two-parent home). The research question presented to Dr. Malani from Oso included the following: Can kindergarten readiness be predicted from social development and type of household?

Oso then assisted Dr. Malani in generating a logistic regression as the test of inference, and a template for writing the research question for this design is presented below.

Can [dependent variable] be predicted from [list independent variables]?

It may be helpful to preface the results of the logistic regression with information on an examination of the extent to which the assumptions were met. The assumptions include (a) independence, (b) linearity, and (c) noncollinearity. We will also examine the data for outliers and influential points.

Logistic regression was conducted to determine whether social development and type of household (single-parent vs. two-parent home) could predict kindergarten readiness.

The assumptions of logistic regression were tested. Specifically, these include (a) noncollinearity, (b) linearity, and (c) independence of errors.

In terms of **noncollinearity**, a VIF value of 4.037 (below the value of 10.0, which indicates the point of concern) and tolerance of .248 (above the value of .10, which suggests multicollinearity) provided evidence of noncollinearity. However, there was some concern for multicollinearity. In examining the collinearity diagnostics, a condition index value of 14.259 was observed, which falls within the range of concern (specifically 10–30). Review of the variance proportions suggested that 100% of the variance of the regression coefficient for social development and 73% for type of household were related to the smallest eigenvalue. This also suggests concern for multicollinearity. Thus, while we met the assumption of noncollinearity with the tolerance and VIF values, there is some concern for multicollinearity with the condition index and variance proportion values.

Linearity was assessed by reestimating the model and including, along with the original predictors, an interaction term that was the product of the continuous independent variable (i.e., social development) and its natural logarithm. The interaction term was not statistically significant, thus providing evidence of linearity ($\text{social}^*\ln(\text{social})$, $B = -2.948$, $SE = 2.845$, $\text{Wald} = 1.074$, $df = 1$, $p = .300$).

Independence was assessed by examining a plot of the standardized residuals against values of each independent variable. With the exception of one case, which was slightly outside the band, all cases were within an absolute value of 2.0, thus indicating the assumption of independence has been met.

In reviewing for **outliers and influential points**, Cook's distance values were generally within the recommended range of less than 1.0, although the

maximum value was 1.587. Leverage values ranged from .007 to .307, well under the recommended .50, suggesting outliers were not problematic. DFBETA values beyond 1.0 also suggested cases that may be exerting influence on the model. Based on the evidence reviewed, there are some cases that are suggestive of outlying and influential points. Due to the small sample size, however, these cases were retained. Readers are urged to interpret the results with caution given the possible influence of outliers.

Here is an example paragraph of results for the logistic regression (remember that this will be prefaced by the previous paragraph reporting the extent to which the assumptions of the test were met).

Logistic regression analysis was then conducted to determine whether kindergarten readiness (prepared vs. unprepared) could be predicted from social development and type of household (single- versus two-parent home). Good model fit was evidenced by nonstatistically significant results on the Hosmer and Lemeshow test, $\chi^2 (n = 20) = 4.691$, $df = 7$, $p = .698$, and large effect size indices when interpreted using Cohen (1988) (Cox and Snell $R^2 = .546$; Nagelkerke $R^2 = .738$). These results suggest that the predictors, as a set, reliably distinguished between children who are ready for kindergarten (i.e., prepared) versus unprepared. Of the two predictors in the model, only social development was a statistically significant predictor of kindergarten readiness ($\text{Wald} = 4.696$, $df = 1$, $p = .030$). The odds ratio for social development suggests that for every one point increase in social development, the odds are about 2 and 2/3 greater for being prepared for kindergarten as compared to unprepared. Type of household was not statistically significant, which suggests that the odds for being prepared for kindergarten (relative to unprepared) are similar regardless of being raised in a single-parent versus a two-parent household. The table below presents the results for the model including the regression coefficients, Wald statistics, odds ratios, and 95% confidence intervals for the odds ratios. This is followed by a table that presents the group means and standard deviations of each predictor for both children who are prepared and unprepared for kindergarten.

Logistic Regression Results

	95% CI for Exp(B)						
	B	SE	Wald	p	Exp(B)	Lower	Upper
Intercept (constant)	-15.404	7.195	4.584	.032	NA		
Social development	.967	.446	4.696	.030	2.631	1.097	6.313
Type of household (two-parent home)	-6.216	3.440	3.265	.071	.002	.000	1.693

Group Means (and Standard Deviations) of Predictors

Predictor	Prepared for Kindergarten	Unprepared for Kindergarten
Social development	23.58 (4.74)	15.13 (5.14)
Type of household (two-parent home)	.67 (.49)	.25 (.46)

Overall, the logistic regression model accurately predicted 90% of the children in our sample, with children who are prepared for kindergarten slightly more likely to be classified correctly (91.7% of children prepared for kindergarten and 87.5% of children unprepared correctly classified). To account for chance agreement in classification, the Kappa coefficient was computed and found to be .792, a large value. Additionally, Press's Q was calculated to be 12.8, providing evidence that the predictions based on the logistic regression model are statistically significantly better than chance. Post hoc power, calculated using G*Power (v. 3.1), was less than .01, indicating very weak power.

PROBLEMS

Conceptual Problems

1. Which one of the following represents the primary difference between OLS regression and logistic regression?
 - a. Computer processing time to estimate the model
 - b. The measurement scales of the independent variables that can be included in the model
 - c. The measurement scale of the dependent variable
 - d. The statistical software that must be used to estimate the model
2. Which one of the following is NOT an appropriate dependent variable for binary logistic regression?
 - a. Bernoulli
 - b. Dichotomous
 - c. Multinomial
 - d. One variable with two categories
3. Which of the following would NOT be appropriate outcomes to examine with binary logistic regression?
 - a. Employment status (employed; unemployed not looking for work; unemployed looking for work)
 - b. Enlisted member of the military (member vs. nonmember)
 - c. Marital status (married vs. not married)
 - d. Recreational athlete (athlete vs. nonathlete)
4. Which of the following represents what is being predicted in binary logistic regression?
 - a. Mean difference between two groups
 - b. Odds that the unit of analysis belongs to one of two groups
 - c. Precise numerical value
 - d. Relationship between one group compared to the other group
5. While probability, odds, and log odds may be computationally different, they all relay the same basic information.
 - a. True
 - b. False

6. A researcher is studying diet soda drinking habits and has coded ‘diet soda drinker’ as ‘1’ and ‘non diet soda drinker’ as ‘0.’ Which of the following is a correct interpretation given a probability value of .52?
 - a. The odds of being a diet soda drinker are about equal to those of not being a diet soda drinker.
 - b. The odds of being a diet soda drinker are substantially greater than not being a diet soda drinker.
 - c. The odds of being a diet soda drinker are substantially less than not being a diet soda drinker.
 - d. Cannot be determined from the information provided.
7. Which of the following is a correct interpretation of the logit?
 - a. The log odds become larger as the odds increase from 1 to 100.
 - b. The log odds become smaller as the odds increase from 1 to 100.
 - c. The log odds stay relatively stable as the odds decrease from 1 to 0.
 - d. The change in log odds becomes larger when the independent variables are categorical rather than continuous.
8. Which of the following correctly contrasts the estimation of OLS regression as compared to logistic regression?
 - a. The sum of the squared distance of the observed data to the regression line is minimized in logistic regression. The log likelihood function is maximized in OLS regression.
 - b. The sum of the squared distance of the observed data to the regression line is maximized in logistic regression. The log likelihood function is minimized in OLS regression.
 - c. The sum of the squared distance of the observed data to the regression line is maximized in OLS regression. The log likelihood function is minimized in logistic regression.
 - d. The sum of the squared distance of the observed data to the regression line is minimized in OLS regression. The log likelihood function is maximized in logistic regression.
9. Which of the following is NOT a test that can be used to evaluate overall model fit for logistic regression models?
 - a. Change in log likelihood
 - b. Hosmer-Lemeshow goodness-of-fit
 - c. Cox and Snell R^2 squared
 - d. Wald test
10. A researcher is studying diet soda drinking habits and has coded ‘diet soda drinker’ as ‘1’ and ‘non diet soda drinker’ as ‘0.’ She has predicted drinking habits based on the individual’s weight (measured in pounds). Given this scenario, which of the following is a correct interpretation of an odds ratio of 1.75?
 - a. For every one unit increase in being a diet soda drinker, the odds of putting on an additional pound increase by 75%.
 - b. For every one unit increase in being a diet soda drinker, the odds of putting on an additional pound decrease by 75%.

- c. For every one pound increase in weight, the odds of being a diet soda drinker decrease by 75%.
- d. For every one pound increase in weight, the odds of being a diet soda drinker increase by 75%.

Computational Problems

1. You are given the following data, where X_1 (high school cumulative grade point average) and X_2 (participation in school-sponsored athletics; 0 = nonathlete and 1 = athlete; use 0 as the reference category) are used to predict Y (college enrollment immediately after high school, '1,' versus delayed college enrollment or no enrollment, '0').

X_1	X_2	Y
4.15	1	1
2.72	0	1
3.16	0	0
3.89	1	1
4.02	1	1
1.89	0	0
2.10	0	1
2.36	1	1
3.55	0	0
1.70	0	0

Determine the following values based on simultaneous entry of independent variables: intercept; $-2LL$; constant; b_1 ; b_2 ; $se(b_1)$; $se(b_2)$; odds ratios; $Wald_1$; $Wald_2$.

2. You are given the following data, where X_1 (participation in high school honors classes; yes = 1, no = 0; use 0 as the reference category) and X_2 (participation in co-op program in college; yes = 1; no = 0; use 0 as the reference category) are used to predict Y (baccalaureate graduation with honors = 1 versus graduation without honors = 0).

X_1	X_2	Y
0	1	1
0	0	1
1	0	0
1	1	1
1	1	1
0	0	0
1	0	1
0	1	1
1	0	0
0	0	0

Determine the following values based on simultaneous entry of independent variables: intercept; $-2LL$; constant; b_1 ; b_2 ; $se(b_1)$; $se(b_2)$; odds ratios; $Wald_1$; $Wald_2$.

Interpretive Problem

1. Use SPSS to develop a logistic regression model with the example survey data on the website. Utilize ‘do you smoke’ as the dependent (binary) variable to find at least two strong predictors from among the continuous and/or categorical variables in the dataset. Write up the results in APA style, including testing for the assumptions. Determine and interpret a measure of effect size.

REFERENCES

- Aldrich, J. H., & Nelson, F. D. (1984). *Linear probability, logit, and probit models*. Beverly Hills, CA: Sage.
- Belsley, D. A. (1991). *Conditioning diagnostics: Collinearity and weak data in regression*. New York: Wiley.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.
- Cox, D. R., & Snell, E. J. (1989). *Analysis of binary data* (2nd ed.). London: Chapman and Hall.
- Croux, C., Flandre, C., & Haesbroeck, G. (2002). The breakdown behavior of the maximum likelihood estimator in the logistic regression model. *Statistics and Probability Letters*, 60, 377–386.
- Harrell, F. E. J. (1986). The LOGIST procedure. In SAS Institute, Inc. (Ed.), *SUGI supplemental library user's guide* (5th ed., pp. 269–293). Cary, NC: SAS Institute, Inc.
- Hellevik, O. (2009). Linear versus logistic regression when the dependent variable is a dichotomy. *Quality and Quantity*, 43(1), 59–74.
- Hosmer, D. W., Hosmer, T., LeCessie, S., & Lemeshow, S. (1997). A comparison of goodness-of-fit tests for the logistic regression model. *Statistics in Medicine*, 16, 965–980.
- Hosmer, D. W., & Lemeshow, S. (1989). *Applied logistic regression*. New York: Wiley.
- Hosmer, D. W., Lemeshow, S., & Sturdivant, R. X. (2000). *Applied logistic regression* (2nd ed.). Hoboken, NJ: John Wiley & Sons, Inc.
- Long, J. S. (1997). *Regression models for categorical and limited dependent variables*. Thousand Oaks, CA: Sage.
- Menard, S. (1995). *Applied logistic regression analysis*. Thousand Oaks, CA: Sage.
- Menard, S. (2000). *Applied logistic regression analysis* (2nd ed.). Thousand Oaks, CA: Sage.
- Nagelkerke, N. J. D. (1991). A note on a general division of the coefficient of determination. *Biometrika*, 78, 691–692.
- Pampel, F. C. (2000). *Logistic regression: A primer*. Thousand Oaks, CA: Sage.
- Raftery, A. E. (1995). Bayesian model selection in social research. In P. V. Marsden (Ed.), *Sociological methodology 1995* (pp. 111–163). London: Tavistock.
- Xie, X.-J., Pendergast, J., & Clarke, W. (2008). Increasing the power: A practical approach to goodness-of-fit test for logistic regression models with continuous predictors. *Computational Statistics and Data Analysis*, 52, 2703–2713. doi:10.1016/j.csda.2007.09.027

Chapter 6

MULTIVARIATE ANALYSIS OF VARIANCE: SINGLE FACTOR, FACTORIAL, AND REPEATED MEASURES DESIGNS

CHAPTER OUTLINE

6.1	What Multivariate Analysis of Variance Is and How It Works	170
6.1.1	Characteristics	172
6.1.2	Sample Size	180
6.1.3	Power	181
6.1.4	Effect Size	181
6.1.5	Assumptions	183
6.1.6	Conditions	187
6.2	Mathematical Introduction Snapshot	188
6.2.1	Mathematical Introduction Snapshot for One-Way and <i>k</i> -Way MANOVA Models	188
6.2.2	Mathematical Introduction Snapshot for Repeated Measures MANOVA	191
6.3	Computing MANOVA Using SPSS	191
6.3.1	Computing Factorial MANOVA	191
6.3.2	Computing Repeated Measures MANOVA Using SPSS	207
6.4	Data Screening	227
6.4.1	Data Screening for One-Way and <i>k</i> -Way MANOVA Models	227
6.4.2	Data Screening for Repeated Measures MANOVA	239



6.5 Power Using G*Power	251
6.5.1 Power for One-Way and k -Way MANOVA Models	251
6.5.2 Power for Repeated Measures MANOVA	256
 6.6 Research Question Template and Example Write-Up	 260
6.6.1 Research Question Template and Example Write-Up for One-Way and k -Way MANOVA Models	260
6.6.2 Research Question Template and Example Write-Up for Repeated Measures MANOVA	262

KEY CONCEPTS

1. Contrast
2. Mean vector
3. Sum of squares and cross-products
4. Variate

The previous chapters have allowed us to examine procedures that provide for the analyses of multiple independent variables, which are very useful tools as it is rarely the case that a single variable can account for all the variation of the dependent variable. This chapter brings in a bit of the ‘old’ as it is an extension of the analysis of variances procedures with which we are already familiar, and it also brings in the ‘new’ as we are now examining more than one dependent variable—thus the name, multivariate (i.e., multiple dependent variables) analysis of variance (MANOVA). There are many different variations of MANOVA, just as there were with ANOVA. This chapter will focus on single and factorial designs, as well as repeated measures MANOVA.

Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying multivariate analysis of variance, both single and factorial designs, as well as repeated measures designs; (b) determine and interpret the results of MANOVA, including repeated measures designs; and (c) understand and evaluate how to screen data prior to conducting MANOVA (single, factorial, or repeated measures).

6.1 WHAT MULTIVARIATE ANALYSIS OF VARIANCE IS AND HOW IT WORKS

Challie Lenge and Ott Lier, two of the graduate student researchers, have a new project waiting for them today.

As we enter the stats lab today, Ott is sharing the details of their new project, on which they’ll be working with Dr. Childs, a reading researcher who recently began analyzing data from the publicly available Early Childhood Longitudinal

Study-Kindergarten Class (ECLS-K) of 1998–1999. Ott had met with Dr. Childs the previous day when Dr. Childs shared a bit about her research interest in using the ECLS-K. At that time, Dr. Childs indicated her interest in exploring how family structure and gender interact in relation to differences in reading, mathematics, and science outcomes collectively, and she came to the stats lab seeking to collaborate on the statistical analyses. After summarizing the conversation with Challie, Ott and Challie brainstorm about the project and decide that a factorial multivariate analysis of variance (MANOVA) is a perfect fit.

At this point in your statistical career, you should already be quite familiar with tests of means such as the independent and dependent *t* test and various ANOVA models (e.g., one-way, factorial, repeated measures). This chapter is an extension of these familiar tests of means. More specifically, we are now covering procedures that allow for analysis of means of multiple dependent outcomes. The dependent variable is now actually the combination of multiple outcomes (where the vector of multiple means is a centroid), and you'll see that some of the terms we learned in ANOVA have multivariate counterparts (see Box 6.1). The current chapter focuses first on single-factor (i.e., one independent variable) and factorial designs (i.e., two or more independent variables, also known as *k*-way MANOVA) and then on repeated measures MANOVA. Before we go any further, it is important to establish notation. Notation μ_{mj} refers to the population mean of dependent variable m for independent variable group j . Dependent variables, m , range from 1 to p , and categories or groups of the independent variable, j , range from 1 to k . An interesting side note, MANOVA and discriminant analysis (which is covered in a later chapter) are mathematically equal although they answer different questions, and we'll actually use discriminant analysis as a follow-up to statistically significant omnibus MANOVA results.

BOX 6.1 ANOVA-MANOVA COUNTERPART TERMS

ANOVA	MANOVA
Homogeneity of variance	Homogeneity of variance-covariance matrices
Sums of squares	Sums of squares and cross products
<i>F</i> test	Hotelling's trace, Pillai's trace, Roy's largest root, Wilks's lambda
Univariate normality	Multivariate normality

With one-way ANOVA, our research question template was very simply: Is there a mean difference in [dependent variable] between [independent variable]? The questions posed by MANOVA are quite similar, but we are now looking at a vector of means (i.e., a centroid). Thus, in MANOVA we are able to answer questions that ask about k group differences in the *centroid*, the vector of means of p dependent variables. In MANOVA factorial designs, we are able to look at not only main effects

but also interaction effects and determine which (if any) dependent variables differ as a result of the main effect(s) and/or interactions. Now, the question you are likely asking yourself is, ‘Why in the world do we need to complicate our statistical lives any more by learning MANOVA when we can simply conduct multiple ANOVAS?’ Ahhh . . . Excellent question. As with any situation where multiple tests are conducted, we increase the chance of a Type I error. For example, if we have three dependent variables and are testing at an alpha of .05, one option is to conduct three ANOVAs. By doing so, however, our alpha increases from .05 to .14 (i.e., $1 - .95^3$) which is unacceptable. Additionally, multiple ANOVAs do not allow you to answer the *really* cool question in which you are most interested, which is, what are the combined effects of the dependent variables as a result of the independent variable(s)? MANOVA not only solves the increased Type I error problem but also gives us much more sophisticated analyses by allowing for the dependent variables to be examined collectively. In this process, we may be able to find group differences that were not evident when the dependent variables were examined in isolation. Should the effect of the independent variable(s) on the combined dependent variables be statistically significant, then these effects are examined separately.

6.1.1 Characteristics

This chapter examines multiple types of MANOVA models. Our discussion will generally separate one-way and factorial designs from repeated measures MANOVA designs.

6.1.1.1 Characteristics of One-Way and k -Way MANOVA Models

In a one-way MANOVA, we consider the effect of one factor or independent variable on two or more dependent variables. In a two-factor (i.e., two-way) MANOVA, we examine the effect of two factors or independent variables on two or more dependent variables. Each factor consists of two or more levels (or categories). This yields what we call a **factorial MANOVA design**, because more than a single factor is included. As with ANOVA designs, we see then that the two-factor MANOVA is an extension of the one-factor MANOVA. Again, we ask the question of why a researcher would want to complicate things more, this time by considering a second factor. Additional factors are considered for the same reasons we did so in factorial ANOVA. First, the researcher may have a genuine interest in studying the second factor. Rather than studying each factor separately in two analyses, the researcher includes both factors in the same analysis. This allows a test of the effect not only of each individual factor, known as **main effects**, but also of both factors collectively. This latter effect is known as an **interaction effect** and provides information about whether the two factors are operating independent of one another (i.e., no interaction exists) or whether the two factors are operating together to produce some additional impact (i.e., an interaction exists). If two separate analyses were conducted, one for each independent variable, no information would be obtained about the interaction effect. As becomes evident, assuming a factorial MANOVA with two independent variables, the researcher will test three

hypotheses: one for each factor or main effect individually and a third for the interaction between the factors. Factorial MANOVA models with more than two independent variables will, accordingly, test for additional main effects and interactions.

A second reason for including an additional factor is an attempt to reduce the error (or within-groups) variation, which is variation that is unexplained by the first factor. The use of a second factor provides a more precise estimate of error variance. For this reason, a two-factor design is generally more powerful than two one-factor designs, as the second factor and the interaction serve to control for additional extraneous variability (and testing two one-factor designs is also problematic in that we increase the Type I error rate). A third reason for considering two factors simultaneously is to provide greater generalizability of the results and to provide a more efficient and economical use of observations and resources. Thus, the results can be generalized to more situations, and the study will be more cost efficient in terms of time and money.

In addition, for the two-factor MANOVA every level of the first factor (hereafter known as factor A) is paired with every level of the second factor (hereafter known as factor B). In other words, every combination of factors A and B is included in the design of the study, yielding what is referred to as a **fully crossed design**. If some combinations are not included, then the design is not fully crossed and may form some sort of a nested design (we will talk about nesting in our discussion of multilevel models, discussed in a later chapter). Each individual responds to only one combination of the factors. If individuals or cases respond to more than one combination of the factors, this would be some sort of repeated measures design, which we examine in the following chapter. In this chapter, we only consider models where all factors are fixed. Thus, the overall design is known as a fixed-effects model. If one or both factors are random, then the design is not a fixed-effects model, again an avenue of MANOVA that we will not explore in this text.

6.1.1.2 Hypotheses of One-Way and k -Way MANOVA Models

As a refresher, the null and alternative hypotheses, respectively, for testing the main effect of factor A in a *univariate* ANOVA are as follows:

$$H_0: \mu_{1\cdot} = \mu_{2\cdot} = \dots = \mu_{J\cdot}$$

$$H_1: \text{not all the } \mu_{\cdot j} \text{ are equal}$$

Testing the main effect of factor B in a univariate ANOVA follows similarly. For the two-factor MANOVA model, there are three sets of hypotheses, one for each of the main effects and one for the interaction effect. This is no different from what we saw with ANOVA. What differs with MANOVA is that we are no longer testing simple population means but rather now we are testing population mean vectors. The MANOVA null hypothesis tests for each dependent variable that all k groups of the independent variable have the same mean.

$$H_0 : \begin{bmatrix} \mu_{11} = \mu_{12} = \dots = \mu_{1k} \\ \mu_{21} = \mu_{22} = \dots = \mu_{2k} \\ \vdots \\ \mu_{p1} = \mu_{p2} = \dots = \mu_{pk} \end{bmatrix}$$

In the *multivariate* case, the null and alternative hypotheses, respectively, for testing the main effect of factor A can best be demonstrated using matrix notation, as the notation above can get a bit cumbersome. The following notation represents the column vector of population means for the dependent variables p for group j of the independent variable:

$$\boldsymbol{\mu}_j = \begin{pmatrix} \mu_{1j} \\ \mu_{2j} \\ \vdots \\ \mu_{pj} \end{pmatrix}$$

Given this, the multivariate null and alternative hypotheses are testing vector equality and can be noted much more compactly as follows:

$$H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = \dots = \boldsymbol{\mu}_k$$

$$H_1 : \text{not all the } \mu_k \text{ are equal}$$

In lay terms, the null hypothesis is that all groups have equal means for all dependent variables, and the alternative hypothesis is that for at least one group of the independent variable there is a mean difference on at least one dependent variable.

6.1.1.3 Omnibus Multivariate Tests of One-Way and k -Way MANOVA Models

As a reminder, the null hypothesis in MANOVA tests equality of vectors of means across groups. Determining differences across the multiple outcomes can be done by examining a number of different omnibus tests, analogous to the omnibus ANOVA F test, and with the exception of data screening and testing assumptions, this is the first step in MANOVA. Unlike ANOVA where the F test is the only test reviewed, in MANOVA multiple omnibus tests that differ mathematically may be used to evaluate the overall hypothesis. The four omnibus tests that we will review include Wilks's lambda, Pillai's trace (also known as Pillai-Bartlett trace), Hotelling's trace (also known as Hotelling-Lawley trace), and Roy's largest root (also known as Roy's greatest characteristic root). Wilks's lambda is the product of the unexplained variances for each discriminant variate [where *variante* is the multivariate equivalent of a variable, representing the statistic that is computed from the linear combination of dependent variables; these are also sometimes referred to as *discriminant function variates*]. Pillai's, the sum of explained variances on the discriminant variates, and Hotelling's, the

sum of the ratio of sum of squares between to sum of squares within (i.e., SSB/SSW) for each discriminant variate (Bray & Maxwell, 1985), are similar and can be approximated by the F statistic.

With so many different omnibus tests from which to select, how is a researcher to decide? The conditions that are being tested should guide the test to select. Many times the results will be similar, but there are some exceptions where power is increased. In ideal or nearly ideal conditions (e.g., assumptions are met, balanced, or nearly balanced designs), Pillai's and Wilk's are preferred. In ideal conditions *and* when the dependent variable represents a single construct, Roy's is preferred. In less balanced designs *or* when homogeneity of variance-covariance is violated, Pillai's is recommended. Should the omnibus MANOVA test be statistically significant, additional analysis is required.

There are different approaches that can be taken when the overall MANOVA is statistically significant. One approach is to follow up with univariate ANOVAs and post hoc tests to determine group differences. However, this approach is criticized because, unlike MANOVA, it does not account for the correlation between dependent variables and is thus less powerful. It is possible to have an overall statistically significant MANOVA but no statistically significant univariate ANOVAs. A more advantageous approach when the omnibus MANOVA is statistically significant, one to which we will adhere, is to conduct discriminant analysis. As we'll learn in a later chapter, discriminant analysis examines the combination of independent variables that can separate groups or categories. Thus, when using discriminant analysis as the follow-up to MANOVA, the dependent variables in our MANOVA model will become the independent variables in discriminant analysis, and the independent variables in MANOVA will become the dependent variables in discriminant analysis (note that discriminant analysis is not a multivariate outcome procedure; in other words, each factor in MANOVA will be tested as a univariate discriminant analysis). An entire chapter of this textbook is devoted to discriminant analysis, so only a cursory introduction is provided here.

Discriminant analysis is MANOVA backwards. MANOVA answers the question of whether mean differences on a combination of continuous dependent variables can be determined based on group membership. Discriminant analysis answers the question of whether group membership can be determined from a combination of continuous predictor variables. If we are working with human subjects, for example, discriminant analysis allows prediction or identification of individuals within groups based on a set of 'discriminating variables,' which we will term independent variables for simplicity and ease in understanding their purpose. The emphasis in discriminant analysis is on interpreting the pattern of differences based on the predictors as a whole, so that the differences in dimensions along which the groups differ can be better understood. Group membership in discriminant analysis can be binary (i.e., dichotomous or two categories, termed two-group discriminant analysis) or multi-category (i.e., multinomial or more than two categories, termed multiple discriminant analysis), and the independent variables should be at least interval in scale as means and variances will be computed from them.

6.1.1.4 Planned and Post Hoc Comparison Procedures of One-Way and k -Way MANOVA Models

Planned and post hoc comparisons can be conducted in MANOVA in a similar fashion, as we learned in ANOVA models. This section examines specific types of contrasts or comparisons (terms that are used interchangeably). A contrast in ANOVA is a weighted combination of the means. A **contrast** in MANOVA is a weighted linear combination of group mean vectors. Statistically, a contrast in MANOVA is the linear combination of group mean vectors, denoted as follows:

$$\Psi = c_1\mu_1 + c_2\mu_2 + \dots + c_k\mu_k$$

where the c_j are known as contrast coefficients (or weights), which are positive, zero, and negative values used to define a particular contrast Ψ_i (the contrast estimated using sample mean vectors is denoted as ψ) and the μ_k are population group mean vectors. In other words, a contrast is simply a particular combination of the group mean vectors, depending on which means the researcher is interested in comparing. A set of contrasts is *orthogonal* if they represent nonredundant and independent sources of variation. A set of two contrasts is *not* orthogonal if the set of contrasts does not sum to zero.

One way of classifying contrasts is whether the contrasts are formulated prior to the research or following a significant omnibus test. **Planned contrasts** (also known as specific or *a priori* contrasts) involve particular comparisons that the researcher is interested in examining *prior* to data collection. These planned contrasts are generally based on theory, previous research, and/or specific hypotheses. Here the researcher is interested in certain specific contrasts *a priori*, where the number of such contrasts is usually small. Planned contrasts are done without regard to the result of the omnibus test (i.e., whether or not the overall test is statistically significant). In other words, the researcher is interested in certain specific contrasts, but not in the omnibus test that examines all possible contrasts. In this situation, the researcher could care less about the multitude of possible contrasts and need not even examine the overall test; rather the concern is only with a few contrasts of substantive interest. In addition, the researcher may not be as concerned with the family-wise error rate for planned comparisons because only a few of them will actually be carried out. Fewer planned comparisons are usually conducted (due to their specificity) than post hoc comparisons (due to their generality), so planned contrasts generally yield narrower confidence intervals, are more powerful, and have a higher likelihood of a Type I error than post hoc comparisons.

Post hoc contrasts are formulated such that the researcher provides no advance specification of the actual contrasts to be tested. This type of contrast is done *only* following a statistically significant omnibus test. Post hoc is Latin for “after the fact,” referring to contrasts tested after a statistically significant omnibus test in MANOVA. Here the researcher may want to take the family-wise error rate into account somehow to achieve better overall Type I error protection. Post hoc contrasts are also known as unplanned, *a posteriori*, or postmortem contrasts. You likely remember studying

a number of multiple comparison procedures (e.g., Tukey's HSD, Dunn's procedure, and more) in ANOVA that could be applied following a statistically significant post hoc test. These will not be covered in the context of MANOVA, as we will adhere to the convention that a statistically significant omnibus MANOVA is proceeded by discriminant analysis (rather than multiple comparison procedures).

6.1.1.5 Characteristics of Repeated Measures MANOVA

By now we understand that the ‘dependent variable’ for MANOVA is actually the *collective examination* of multiple outcomes. With repeated measures MANOVA, the outcome is still multiple dependent variables—but with a twist . . . the units of analysis have been measured multiple times on the multiple dependent variables. For example, an exercise scientist might measure body mass index as well as waist circumference before *and* after participants complete an exercise program. Thus, there are two outcomes (body mass index as well as waist circumference) that have been measured at two points in time (before and after the program). Because the functional relationship between these variables is tenable, repeated measures MANOVA is appropriate. As another example, it would be reasonable to assume a functional relationship between mathematics self-efficacy and mathematics self-concept. When measured at two or more points in time, repeated measures MANOVA is again a viable option for analysis of this data. As another illustration, children are rated by multiple raters on externalizing and internalizing behaviors. While they are not rated at different points in time, they are rated on the same outcomes—which appear to be functionally related—by multiple raters, again making repeated measures MANOVA appropriate. Thus in repeated measures MANOVA, we have multiple within-subjects variables that are functionally related and thus appropriate to collectively analyze. It is also possible to add between-groups factors (i.e., independent variables). The between-groups factor, if/when applicable, can be a single factor (i.e., one independent variable) or a factorial design (i.e., two or more independent variables). When there are *both* within-subjects variables (e.g., multiple outcomes at two or more time points) and between-groups variables (e.g., one or more independent variables) in the model, this is referred to as a *mixed model*, as well as a *doubly multivariate model*, depending on the analytic approach—more on this later!

In MANOVA, we are able to answer questions that ask about k group differences in the *centroid*, the vector of means of p dependent variables. In MANOVA factorial designs, we are able to look at not only main effects but also interaction effects and determine which (if any) dependent variables differ as a result of the main effect(s) and/or interactions. Repeated measures MANOVA extends this process in that now we have multiple within (i.e., over time) factors and potentially also between-groups factors (i.e., independent variables) to examine, with the capacity to examine main effects and interactions still present.

In repeated measures MANOVA, we consider the effect of two or more dependent variables that have been repeatedly measured (e.g., measured at two or more points in

time or multiple raters rating with the same rating scales) (i.e., within-subjects factor). Between subjects, factors can also be added to repeated measures MANOVA, such as one or more factors or independent variables (between-subjects factor) on two or more dependent variables that are measured multiple times can be considered. There are different analytic approaches for repeated measures MANOVA. Assuming the researcher wishes to follow a multivariate approach in at least some capacity, repeated measures MANOVA models can be approached as either a **multivariate mixed model (MMM)** or a **doubly multivariate model (DMM)** (Boik, 1988). A multivariate approach to repeated measures should be followed if it is reasonable to assume there is some functional relation among the multiple dependent variables. In the later discussion of assumptions, it will be seen that the DMM is the recommended approach, because the MMM assumption of sphericity is often violated.

6.1.1.6 Hypothesis of Repeated Measures MANOVA

With the doubly multivariate repeated measures design, we are actually able to address three research questions and their related hypotheses. The first research question relates to the collective multiple dependent measures across time, and the remaining address the individual repeated measures (with an additional question for each dependent variable; this example lists only two dependent measures—A and B). [Note: The most common application of repeated measures in the social sciences is multiple measures over time (e.g., pre, post, and follow-up). Thus, throughout this chapter the ‘repeated measures’ will imply repeated measures over time. However, situations such as multiple raters using the same instrument can also be applicable to repeated measures.]

1. Is there a difference in the linear combination of the multiple dependent measures across time?
2. Is there a difference in [dependent measure A] across time?
3. Is there a difference in [dependent measure B] across time?

Stated another way, the research question could be phrased: Is there a mean difference over time in the noncommensurate dependent variables jointly between the independent variable(s)?

6.1.1.7 Omnibus Multivariate Tests for Repeated Measures MANOVA

The same four omnibus tests that we learned about in MANOVA are again the omnibus multivariate tests in repeated measures MANOVA when approached as a doubly multivariate model. These omnibus tests differ mathematically and are used to evaluate the overall or omnibus hypothesis—that being the difference in the linear combination of the multiple dependent measures across time. These include Wilks’s lambda, Pillai’s trace (also known as Pillai-Bartlett trace), Hotelling’s trace (also known as Hotelling-Lawley trace), and Roy’s largest root (also known as Roy’s

greatest characteristic root). Wilks's lambda is the product of the unexplained variances for each discriminant variate [where *variate* is the multivariate equivalent of a variable and represents the “statistical criterion which results from the linear combination of variables” (1985, p. 10); these are also sometimes referred to as *discriminant function variates*]. Pillai's, the sum of explained variances on the discriminant variates (Bray & Maxwell, 1985), and Hotelling's, the sum of *SSB/SSW* for each discriminant variate (Bray & Maxwell, 1985), are similar and can be approximated by the *F* statistic.

With so many different omnibus tests from which to select, how is a researcher to decide? The conditions that are being tested should guide the test to select. Many times the results will be similar, but there are some exceptions where power is increased. In ideal or nearly ideal conditions (e.g., assumptions are met, balanced, or nearly balanced design), Pillai's and Wilk's are preferred. In ideal conditions *and* when the dependent variable represents a single construct, Roy's is preferred. In less balanced designs *or* when homogeneity of variance-covariance is violated, Pillai's is recommended. Should the omnibus repeated measures MANOVA test be statistically significant, the predictors must then be examined.

There are different approaches that can be taken when the overall MANOVA is statistically significant. One approach is to follow up with univariate ANOVAs and post hoc tests to determine group differences. However, this approach is criticized because, unlike MANOVA, it does not account for the correlation between dependent variables and is thus less powerful. It is possible to have an overall statistically significant MANOVA but no statistically significant univariate ANOVAs. A more advantageous approach when the omnibus MANOVA is statistically significant, one to which we will adhere, is to conduct discriminant analysis. As we'll learn in a later chapter, discriminant analysis examines the combination of independent variables that can separate groups or categories. Thus, when using discriminant analysis as the follow-up to MANOVA, the dependent variables in our MANOVA model will become the independent variables in discriminant analysis, and the independent variables in MANOVA will become the dependent variables in discriminant analysis (note that discriminant analysis is not a multivariate outcome procedure; in other words, each factor in MANOVA will be tested as a univariate discriminant analysis). An entire chapter of this textbook is devoted to discriminant analysis, so only a cursory introduction is provided here.

6.1.1.8 Planned and Post Hoc Comparison Procedures for Repeated Measures MANOVA

Planned and post hoc comparison procedures are applicable in repeated measures MANOVA when there is at least one between-subjects factor that has more than two levels. These procedures, which were detailed in our discussion of MANOVA, are applicable to repeated measures MANOVA as well. Should you need a refresher, please refer to the section in the previous chapter.

6.1.2 Sample Size

We will examine sample size from the perspective of one-way and factorial designs, as well as from repeated measures MANOVA designs.

6.1.2.1 Sample Size for One-Way and k -Way MANOVA Models

Considerations of sample size in one-way and k -way MANOVA are both overall and within cell. As an absolute minimum criterion, the sample size per cell must be larger than the number of dependent variables. There are a number of suggested guidelines for sample size. One recommendation is that there be, at minimum, 20 cases per cell. Thus, if you have a one-way MANOVA and the independent variable has only two categories, the requisite sample size is 40 (i.e., 2 categories x 20). A 2 x 3 MANOVA (i.e., two independent variables with two and three categories, respectively) requires a sample size of 120 (i.e., 6 total categories or cells x 20). (As a refresher from ANOVA . . . The number of cells in factorial MANOVA equals the product of the number of categories of the independent variables. For example, a MANOVA with three independent variables having 2, 3, and 4 categories respectively, has $2 \times 3 \times 4$ or 24 cells.) A better approach to determining sample size is to conduct power analysis (an illustration of using G*Power is provided later in the chapter).

An issue related to overall sample size is cell size—which segues into a discussion of *balanced* and *unbalanced* designs. When the sample sizes within cell are unequal, an unbalanced design exists. Unbalanced designs are not a deal killer in MANOVA, per se, but it can be problematic particularly if there are fewer cases in a cell than there are dependent variables—i.e., singularity. In cases of singularity, the assumption of equal variance-covariance cannot be tested. Another general guideline: The recommended ratio of the largest to smallest cell size should be within a ratio of 1.5:1. As we know, sample size and power are intricately related. In the case of MANOVA, as the number of dependent variables increase, the sample size also has to increase in order to sustain sufficient power.

6.1.2.2 Sample Size for Repeated Measures MANOVA

There exist relatively few standing rules for sample size for repeated measures MANOVA. Thus, applying the suggestions for MANOVA seem prudent. Considerations of sample size in repeated measures MANOVA are both overall and within cell. As an absolute minimum criteria, the product of the number of dependent variables and number of measurement occasions (i.e., the number of within-subjects factors) should be less than the sample size. For example, if there are 3 dependent measures, and each has been measured on 4 occasions, then the total sample size should be more than 12 (i.e., 3×4). There are also the recommendations for sample size that were discussed in relation to MANOVA—such as at minimum, 20 cases per cell. A better approach to determining sample size is to conduct power analysis (an illustration of using G*Power is provided later in the chapter). Let us briefly broach missing data on the dependent

measures. Should there be missing data on one or more dependent measures (i.e., one or more cases is not measured on all dependent variables at all measurement occasions), the researcher must generally apply a missing data technique if the common default of listwise deletion is not desirable. Applying missing value techniques to outcome data is still controversial; however, there is some research on replacement of missing outcomes in randomized controlled trials (Puma, Olsen, Bell, & Price, 2009).

The earlier discussion of *balanced* and *unbalanced* designs for one-way and factorial MANOVA is applicable to repeated measures designs as well.

6.1.3 Power

As for power (the probability of correctly rejecting a false null hypothesis), one can consider either planned power (a priori) or observed power (post hoc), as discussed in previous chapters. In MANOVA as with ANOVA, we know that power is primarily a function of α , sample size, and effect size. For planned power, one inputs each of these components either into a statistical table or power chart or, more commonly nowadays, into statistical software (such as G*Power). Planned power is most often used by researchers to determine adequate sample sizes in MANOVA models (regardless of one-way, k -way, or repeated measures), which is highly recommended. Many disciplines recommend a minimum power value, such as .80. Thus, these methods are a useful way to determine the sample size that would generate a desired level of power. Observed power can be determined by some statistics software, such as SPSS, and indicates the power that was actually observed in a completed study.

6.1.4 Effect Size

We will examine effect size from the perspective of one-way and factorial designs, as well as from repeated measures MANOVA designs.

6.1.4.1 Effect Size for One-Way and k -Way MANOVA Models

Wilks's generalized correlation ratio, also known as multivariate eta squared as it is the multivariate extension of eta squared for univariate outcomes, is a measure of effect in MANOVA. Wilks's generalized correlation ratio represents the multivariate proportion of variance accounted for by the best linear combination of dependent variables. Multivariate eta squared shared properties of eta squared. This effect ranges from zero to +1.00, and as a generalized correlation ratio, is a generalization of R^2 . Wilks's generalized correlation ratio is interpreted as the proportion (or percentage) of variance in the discriminant function that is explained by groups (i.e., between-subjects factor(s)). A multivariate eta squared of zero suggests that *none* of the total variance in the total data can be explained. An eta squared of 1.00 indicates that *all* the variance in the total data can be explained. In a two-group MANOVA, multivariate eta squared represents the proportion of variance in the discriminant function that is explained by group membership (Kline, 2004b). Multivariate eta squared can be larger in value than eta

squared for any of the individual outcomes in MANOVA, creating a positively biased statistic (i.e., overestimates the association) (Kline, 2004b). This problem occurs in MANOVA because the dependent variables are considered in aggregate in Wilks's lambda. Additionally, only a portion of the discriminant function variance (estimated as multivariate eta squared) is explained by group membership, and the discriminant function explains only a portion of the variance of the outcomes. The bias is most evident for n 's (i.e., group sample sizes) less than 30. The sum of eta squared may be greater than 1.0 for all effects, and the interpretation can thus be difficult. Multivariate eta squared is calculated as:

$$\text{Multivariate } \eta^2 = 1 - \Lambda$$

Partial eta squared is the effect size reported in standard statistical software such as SPSS. It is calculated as:

$$\text{partial } \eta^2 = \frac{df_1 F}{df_1 F + df_2}$$

The multivariate counterpart of Hedge's g , a univariate standardized mean difference measure of effect, is the square root of Mahalanobis generalized distance, D_M^2 . Mahalanobis generalized distance measures the distance between group centroids relative to the pooled within-groups variance-covariance matrix. Mahalanobis generalized distance can be computed in multiple ways, one of which is using Wilks's lambda. The calculation for a two-group design with any number of dependent variables follows (Kline, 2004b):

$$D_M^2 = \frac{df_w (1 - \Lambda)}{\Lambda} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

Where df_w is the ANOVA within-groups degrees of freedom (e.g., $n_1 + n_2 - 2 = N - 2$ for a two-group design), Λ is the multivariate Wilks's lambda from the MANOVA, and n_1 and n_2 are the sample sizes of the groups of the independent variable. Again, the square root of this value represents a multivariate standardized mean difference effect.

The square root of Mahalanobis generalized distance can be extended to scenarios with more than two groups:

$$D_{M_\Psi}^2 = \frac{df_w (1 - \Lambda_\Psi)}{\Lambda_\Psi} \left(\sum_{i=1}^a \frac{c_i^2}{n_i} \right)$$

Where df_w is the ANOVA total within-groups degrees of freedom (i.e., $N - a$), Λ_Ψ is the multivariate Wilks's lambda from the MANOVA *contrast* (not the omnibus multivariate Wilks's lambda), a equals the number of groups in the between-subjects factor, c_i^2 is the squared contrast weight, and n_i equals the sample size of the i th group.

Cohen's (1988) subjective standards can be used as follows to interpret these effect sizes: small effect, $\sqrt{D_M^2} = .20$, multivariate $\eta^2 = .01$; medium effect, $\sqrt{D_M^2} = .50$, multivariate $\eta^2 = .06$; large effect, $\sqrt{D_M^2} = .80$, multivariate $\eta^2 = .14$. Note that these are subjective standards developed for the behavioral sciences; your discipline may use other standards. Should there be little to no intercorrelations among the dependent variables, any multivariate effect will likely result in a negligible effect, as the variables are measuring independent aspects of group differences. Thus, examining linear combination of the variables is really pointless (Kline, 2004). Wilks's generalized correlation ratio and the square root of Mahalanobis generalized distance both assume homogeneity of variance-covariance matrices (Kline, 2004b). For further discussion on these measures of effect, as well as additional multivariate effect size indices, see Kline (Kline, 2004a, 2004b).

6.1.4.2 Effect Size for Repeated Measures MANOVA

Measures of multivariate effect size, specifically Wilks's generalized correlation ratio (also known as multivariate eta squared) and the square root of Mahalanobis generalized distance (notated as $\sqrt{D_M^2}$), were presented previously in the discussion of MANOVA and do not change appreciably with the doubly multivariate model. For that reason, they are not repeated here. Rather, what does bear repeating is the summary of the interpretations. Cohen's (1988) subjective standards can be used as follows to interpret Wilks's generalized correlation ratio and the square root of Mahalanobis generalized distance: small effect, $\sqrt{D_M^2} = .20$, multivariate $\eta^2 = .01$; medium effect, $\sqrt{D_M^2} = .50$, multivariate $\eta^2 = .06$; large effect, $\sqrt{D_M^2} = .80$, multivariate $\eta^2 = .14$.

6.1.5 Assumptions

We will examine assumptions from the perspective of one-way and factorial designs, as well as from repeated measures MANOVA designs.

6.1.5.1 Assumptions for One-Way and k -Way MANOVA Models

The assumptions for MANOVA share commonality with ANOVA but are transcended to the multivariable level and thus are a bit more complex. These assumptions include (a) independence, (b) multivariate normality of the dependent variables, (c) linearity (a component of multivariate normality), and (d) homogeneity of variance-covariance matrices for the dependent variables.

Independence

The first assumption is concerned with independence of the observations. Violations of this assumption can detrimentally affect standard error values and thus any resulting

hypothesis tests. In particular, even small violations of this assumption can result in a quite dramatically increased actual alpha level as compared to nominal alpha level (Barcikowski, 1981; Scariano & Davenport, 1987). At minimum when there is dependence, tests should be conducted at a decreased alpha level (e.g., .01) given that the actual alpha will likely be higher. Multilevel analysis is also an appropriate approach to consider if this assumption is violated. Researchers may also consider the use of groups as the unit of analysis and analyze group means, rather than individual cases.

Multivariate Normality for the Dependent Variables

For ANOVA, it was assumed that observations were normally distributed in each group. In MANOVA, we are now looking for multivariate normality—the linear combination of variables needs to be normally distributed. It is common practice to examine univariate normality first before assessing multivariate normality. The following can be used to detect univariate normality violations: frequency distributions, normal probability (Q-Q) plots, and skewness statistics. The simplest procedure involves checking for symmetry in a histogram, frequency distribution, boxplot, or skewness and kurtosis statistics. Although **nonzero kurtosis** (i.e., a distribution that is either flat, platykurtic, or has a sharp peak, leptokurtic) will have minimal effect on the parameter estimates, **nonzero skewness** (i.e., a distribution that is not symmetric with either a positive or a negative skew) will have much more impact on these estimates. Thus, finding asymmetrical distributions is a must. One rule of thumb is to be concerned if the skewness value is larger than 1.5 or 2.0 in magnitude.

Another useful graphical technique for examining univariate normality is the normal probability plot (or Q-Q plot). With normally distributed residuals, the points on the normal probability plot will fall along a straight diagonal line, whereas nonnormal data will not. There is a difficulty with this plot because there is no criterion with which to judge deviation from linearity. It is recommended that skewness and/or the normal probability plot be considered at a minimum when determining normality evidence. Plots of the standardized residuals to values expected from the model can also be reviewed as evidence of the extent to which multivariate normality is achieved. Systematic patterns may suggest violations.

The formal test of normality, the Shapiro-Wilk test (*SW*) (Shapiro & Wilk, 1965), provides evidence of the extent to which our sample distribution is statistically different from a normal distribution. Reviewing multiple indices of univariate normality is recommended. Should it be the case that there are conflicting interpretations, review the skewness and Shapiro-Wilk's test. If either skewness or Shapiro-Wilk suggests substantial departure from normality, then normality is rejected (D'Agostino, 1986).

In terms of multivariate normality, we will employ the use of a macro in SPSS (DeCarlo, 1997) to examine a number of multivariate normality indices that include (a) multivariate kurtosis (Mardia, 1970), (b) multivariate skewness and kurtosis based on Small's (1980) multivariate extension of univariate skewness and kurtosis (Looney,

1995), (c) multivariate normality omnibus test (Looney, 1995), (d) largest squared and plot of squared Mahalanobis distance, and (e) critical values for hypothesis test for a single multivariate outlier using Mahalanobis distance (Penny, 1996).

Because univariate normality is a necessary condition for multivariate normality, rejection of the assumption of univariate normality means that multivariate normality is also rejected. Univariate normality is not a sufficient condition for multivariate normality, however, and meeting univariate normality should not imply tenable multivariate normality. Transformations can be used to normalize nonnormal data. However, again there is the problem of dealing with transformed variables measured along some other scale than that of the original variables, which complicates the interpretation. MANOVA is robust to violations of multivariate normality (Huberty & Olejnik, 2006), even with an unbalanced design, when there are 20 or more cases per cell. Other studies suggest an overall sample size of 40, with 10 or more per group of the independent variable to ensure robustness to nonnormality (Seo, Kanda, & Fujikoshi, 1995). Outliers are a different story. MANOVA is actually quite sensitive to outliers, and screening for univariate and multivariate outliers within cells is essential.

Linearity

When all pairs of dependent variables are bivariate normally distributed (i.e., linear), the assumption of linearity is met (thus making this a component of multivariate normality). This assumption can be tested by examining scatterplots of dependent variables as well as Bartlett's test of sphericity, which assesses collective intercorrelation of the dependent variables. If violated, transformations of the dependent variable may be possible. In cases where the variable is not theoretically essential to the model, ill-performing dependent variables may be considered for removal as well.

Homogeneity of Variance-Covariance Matrices for the Dependent Variables

This assumption is also referred to as homoscedasticity. In plain language, when this assumption is met, this means that variation or dispersion between groups on the collective dependent variables is equal (or at least not statistically different). Box's M , which simultaneously examines the KJ group variances and covariances, can be used to test this assumption, and nonstatistically significant Box's M indicates the assumption has been met. Box's test is very sensitive to nonnormality. MANOVA is not robust to violations of this assumption, and this worsens under the following conditions: as the number of dependent variables increases and as the imbalance in the sample sizes per cell increases (particularly when the largest group size is two or more times the size of the smallest group) (Huberty & Olejnik, 2006). If and/or when this assumption is violated, use Levene's test for univariate analysis to determine the dependent variable that has the heterogeneous variance and apply a variance stabilizing transformation,

such as natural log or square root. Because violations of this assumption and that of multivariate normality often go hand in hand, then a transformation to correct for non-normality may also correct for unequal variance-covariance matrices. To replace the omnibus test, the Yao test is an option for examining specific contrasts (Huberty & Olejnik, 2006).

Violations of this assumption are most problematic when the observed probability values are within close range of the alpha level. This is why: The alpha level is too conservative in cases where the cells that have the larger sample sizes also have the larger variance (thus results that are not statistically significant *may* have been if there had been equal variance-covariance matrices—therefore, if you *do* find statistical significance in this situation, there is little reason for concern for violation of the assumption). The alpha level is too liberal in cases where the cells that have the smaller sample sizes have the larger variances (thus results that are statistically significant *may not* have been if there had been equal variance-covariance matrices—this situation is extremely problematic; therefore, if you *do not* find statistical significance in this situation, there is little reason for concern for violation of the assumption).

Concluding Thoughts on Assumptions

As mentioned in previous chapters, there is no rule stating that research that violates assumptions must be scrapped. However, researchers who face violations of one or more assumptions must handle these situations on a case-by-case basis, considering both the goal of the analyses and the extent to which the assumptions were violated and the resulting effect of violation. It is also important that researchers present the evidence found, along with justification for decisions that were made and limitations that may result, as applicable. A summary of the assumptions and the effects of their violation is presented in Table 6.1.

TABLE 6.1

Assumptions and Violation of Assumptions: One-Way, *k*-Way, and Repeated Measures Multivariate Analysis of Variance

Assumption	Effect of Assumption Violation
Independence	<ul style="list-style-type: none">Increased likelihood of a Type I error in the <i>F</i> statistic (inflates actual alpha level to substantially greater than the nominal alpha level); influences standard errors of means and thus inferences about those means
Multivariate normality of dependent variables	<ul style="list-style-type: none">Minimal effect with moderate violationMinimal effect on Type I errors
Homogeneity of variance-covariance matrices of dependent variables	<ul style="list-style-type: none">Increased likelihood of a Type I and/or Type II errorLess effect with equal or nearly equal <i>n</i>'s (ratio of smallest to largest within 1:1.5); when outside this ratio, larger variance in the group with the <i>larger n</i> results in a conservative test and larger variance in the group with the <i>smaller n</i> results in a liberal test; in unbalanced designs with ratios of 1:2 or greater, dramatically increased Type I error rates can occur even with mild heterogeneity

6.1.5.2 Assumptions for Repeated Measures MANOVA

The assumptions for repeated measures MANOVA share commonality with both MANOVA and repeated measures ANOVA and are also somewhat dependent on how you approach your analyses. Regardless of how you approach your repeated measures MANOVA, the assumptions of independence, multivariate normality, linearity (a component of multivariate normality), and homogeneity of variance-covariance hold (see Table 6.1 for a summary). If testing the effects of the multivariate repeated measures design is approached as a multivariate mixed model, then the additional assumption of multivariate sphericity must hold. If the analysis is approached as a doubly multivariate model, then multivariate sphericity is *not* an assumption that must be tested. For this reason, evaluating the repeated measures MANOVA as a doubly multivariate model is recommended (Lix & Hinds, 2004).

To reiterate, the assumptions of the doubly multivariate repeated measures MANOVA include (a) independence, (b) multivariate normality of the dependent variables, (c) linearity (a component of multivariate normality), and (d) homogeneity of variance/covariance. Each of these were presented previously in the discussion of MANOVA and do not change appreciably with the doubly multivariate model. For that reason, they are not repeated here.

6.1.6 Conditions

In multivariate analysis of variance, because this is a test of means, a condition of the test is that the scale of measurement on the dependent variables is at the interval or ratio level. In the case, specifically, of repeated measures MANOVA, each dependent measure must also be measured on more than one occasion. The independent variable in MANOVA is a grouping or discrete variable, so it usually conducted with independent variables, which are categorical—nominal or ordinal in scale. However, there is one caveat to the measurement scale of the independent variable. The independent variable(s) in MANOVAs can also be used in the case of interval or ratio values that are *discrete*. Recall that discrete variables are variables that can only take on certain values and that arise from the counting process. An example of a discrete variable that may be a good candidate for being an independent variable in a MANOVA model is number of children. What would make this a good candidate? The responses to this variable would likely be relatively limited (in the general population it may be anticipated that the range would be from zero children to five or six—although outliers may be a possibility) and each discrete value would likely have multiple cases (with fewer cases having larger numbers of children). Applying this is obviously at the researcher's discretion; at some point, the number of discrete values can become so numerous as to be unwieldy in a MANOVA model. Thus, while at first glance we may not consider it appropriate to use interval or ratio variables as independent variables in MANOVA models, there are situations where it is feasible and may be appropriate.

In summary, the conditions of MANOVA are as follows: (a) one or more independent variables, each with two or more levels; (b) the levels of the independent

variable are fixed by the researcher; (c) subjects are exposed to only one level of the independent variable; and (d) the dependent variables are measured at least at the interval level.

6.2 MATHEMATICAL INTRODUCTION SNAPSHOT

6.2.1 Mathematical Introduction Snapshot for One-Way and k -Way MANOVA Models

Similar to ANOVA, the variance in MANOVA is partitioned into variances due to differences among scores within groups (error or unsystematic variability) and variances due to differences between groups (systematic variability) (see Figure 6.1 for this concept illustrated with a two-way MANOVA design). These squared differences are sums of squares. Dividing the sums of squares by degrees of freedom provides an estimate of variance attributed to source (i.e., main effect, interaction, error). Summing the squared differences creates a sum-of-squares-and-cross-products matrix (the counterpart for which is ANOVA's sum of squares). Given multiple dependent variables, matrices of the dependent variables are used to calculate the 'sum of square and cross-products' (aka 'cross-products' or 'sum of products') (see Figure 6.1). The total cross-products matrix is partitioned into cross-products matrices for differences associated with the first independent variable, with the second independent variable, and with the interaction between the independent variables, as well as for the error (subjects within groups).

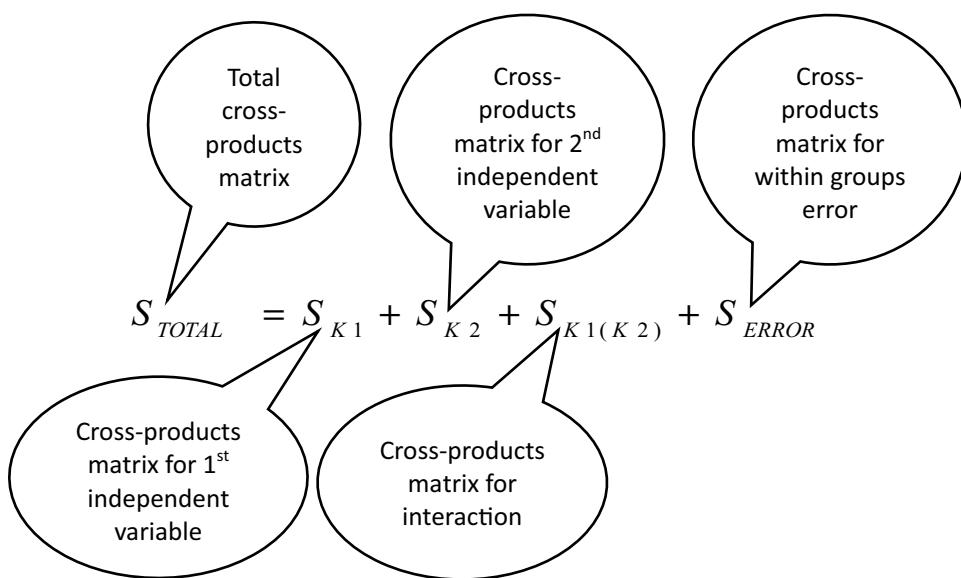


FIGURE 6.1

Partitioning Variation in Two-Way MANOVA

6.2.1.1 Partitioning the Variation

You are already familiar with partitioning the sums of squares in ANOVA. The total sum of squares in Y , denoted as SS_{total} , represents the amount of total variation in Y . The total variation is then partitioned into variation between the groups (i.e., the categories or levels of the independent variable), denoted by SS_{betw} , and variation within the groups (i.e., units or cases within each category or level of the independent variable), denoted by SS_{within} . In the one-factor analysis of variance, we therefore partition SS_{total} as follows:

$$SS_{\text{total}} = SS_{\text{between}} + SS_{\text{within}}$$

or

$$\sum_{i=1}^n \sum_{j=1}^J (Y_{ij} - \bar{Y}_{..})^2 = \sum_{i=1}^n \sum_{j=1}^J (\bar{Y}_{.j} - \bar{Y}_{..})^2 + \sum_{i=1}^n \sum_{j=1}^J (Y_{ij} - \bar{Y}_{.j})^2$$

where SS_{total} is the total sum of squares due to variation among all of the observations without regard to group membership, SS_{between} is the between-groups sum of squares due to the variation between the groups, and SS_{within} is the within-groups sum of squares due to the variation within the groups combined across groups.

The goal in MANOVA is still to determine how much of that variation can be explained, but now the concept of partitioning the variance is working with partitioning the *covariance matrix*. Now working with multiple dependent variables, partitioning the variance is accomplished by partitioning the *multivariate* variance (i.e., the total covariance matrix) into variation explained by the model and variation unexplained (i.e., the residual), and this is done via **cross-products**. Through cross-products, the relationship between the dependent variables can be examined by partitioning the variation into model cross-products and residual cross-products. We still have an ANOVA summary table, but it is formed with matrices rather than scalars, with the matrix consisting of p^2 elements. The total sum of squares and cross-products in MANOVA is the multivariate equivalent of ANOVA's total sums of squares and consists of a $p \times p$ matrix. The total sum of squares in MANOVA can be denoted as follows:

$$T = \sum_{i=1}^g \sum_{j=1}^{n_i} (Y_{ij} - \bar{y}_{..})(Y_{ij} - \bar{y}_{..})'$$

Where g equals the group mean vectors and j equals the groups of the independent variables. The first part of the equation, $(Y_{ij} - \bar{y}_{..})$, denotes the difference between the vector of observations (Y_{ij}) and grand mean vector $(\bar{y}_{..})$. This is multiplied by its transpose. This total sum of squares and cross-products can be partitioned into

error sum of squares and cross-products (**E**) and hypothesis sum of squares and cross-products (**H**) as follows:

$$\begin{aligned} \mathbf{T} &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{..}) (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{..})' \\ &= \sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.}) (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.})' + \sum_{i=1}^g n_i (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..}) (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})' \end{aligned}$$

Where

$$\sum_{i=1}^g \sum_{j=1}^{n_i} (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.}) (\mathbf{Y}_{ij} - \bar{\mathbf{y}}_{i.})' = \text{error sum of squares and cross-products (aka within-groups), } \mathbf{E}$$

$$\sum_{i=1}^g n_i (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..}) (\bar{\mathbf{y}}_{i.} - \bar{\mathbf{y}}_{..})' = \text{hypothesis sum of squares and cross-products (aka between-groups), } \mathbf{H}$$

The MANOVA summary table is provided in Table 6.2.

In ANOVA, we talked in terms of mean squares (i.e., variances). The concept is applicable in MANOVA, although the terminology differs a bit. In ANOVA, sums of squares are divided by degrees of freedom to produce variances or mean squares. In MANOVA, a ‘determinant’ (the variance in matrix terms) is found for each cross-product matrix. In ANOVA, ratios of variances (i.e., mean squares) are formed to test main effects and interactions. In MANOVA, ratios of determinants are formed to test main effects and interactions.

Multivariate *F* degrees of freedom are computed as:

$$\begin{aligned} df_1 &= p(df_{\text{effect}}) \\ df_2 &= s \left[df_{\text{error}} - \frac{p - df_{\text{effect}} + 1}{2} \right] - \left[\frac{p(df_{\text{effect}}) - 2}{2} \right] \end{aligned}$$

Where:

p = number of dependent variables

df_{effect} = number of levels of first independent variable minus 1 multiplied by the number of levels of the second independent variable minus 1

df_{error} = number of levels of first independent variable multiplied by the number of levels of the second independent variable multiplied by the number of scores per cell for each dependent variable

s = minimum of (*p*, *df_{effect}*)

■ TABLE 6.2

Multivariate Analysis of Variance Summary Table

Source	SSCP	<i>df</i>
Hypothesis (aka between groups)	H	<i>g</i> – 1
Error (aka within groups)	E	<i>N</i> – <i>g</i>
Total	T	<i>N</i> – 1

6.2.2 Mathematical Introduction Snapshot for Repeated Measures MANOVA

MANOVA constructs the weighted linear combination of the dependent variables with group separation that best accounts for intercorrelation between the variables (Kline, 2004b). In repeated measures MANOVA, following the doubly multivariate approach, the multivariate general linear model is employed, testing linear hypotheses using p -dimensional observed vectors (Timm, 1980). The hypothesis of interest is most often that the vector means are equal. For further reading on the mathematics underlying repeated measures MANOVA, readers are referred to Timm (1980), which provides a much more technical view of the subject. Researchers interested in an extension of repeated measures MANOVA are encouraged to consider growth curve modeling.

6.3 COMPUTING MANOVA USING SPSS

Our applied examples will first illustrate factorial MANOVA and then repeated measures MANOVA.

6.3.1 Computing Factorial MANOVA

We will use SPSS to compute MANOVA. However, before we conduct the analysis, let us talk about the data. The data we are using is the Early Childhood Longitudinal Study Kindergarten Class of 1998–1999 (ECLS-K) (<http://nces.ed.gov/ecls/kindergarten.asp>), available through the U.S. Department of Education National Center for Education Statistics Institute of Education Sciences (NCES IES). Thank you to NCES for making this data publicly available. The ECLS-K is a longitudinal study, following a nationally representative sample of public and private school children from kindergarten through eighth grade. Data was collected in the fall and spring of kindergarten and first grade (1998–1999 and 1999–2000, waves 1 through 4 respectively), spring of third (2002, wave 5) and fifth (2004, wave 6) grades, and spring of eighth grade (2007, wave 7). Direct (administered to the child) and indirect (administered to the parents and/or teachers) assessments were conducted including measures of physical development, social and emotional skills, and cognitive abilities.

As with many of the data we will work with, the ECLS-K is a complex sample (i.e., not a simple random sample), employing a multistage probability design to select a nationally representative sample of children who attended kindergarten during 1998–1999. During the base year, geographic areas (counties or groups of counties) were the primary sampling units (PSUs), schools within sampled PSUs were the second-stage sampling units, and children within schools were the third- and last-stage sampling units. For each successive wave of data collection, the sampling design was modified. For example, in first grade, a 30% subsample of children in the base year were included, and this sample was freshened with children who were not enrolled in kindergarten during the base year of data collection. You get the picture. . . . Additional details on the technical aspects of the ECLS-K can be accessed from the technical report (Tourangeau, Nord, Le, & Sorongon, 2009). If you access the full dataset (<http://nces.ed.gov/ecls/dataproducts.asp>), you will find a large number of weights as well as

stratum and unit variables that can be used to address the complex sampling design. We won't get into the technical aspects of this, but when the data are analyzed to adjust for the sampling design (including nonsimple random sampling procedure and disproportionate sampling), the end results are then representative of the intended population (in this case, nationally representative of children who attended public or private kindergarten in 1998–1999). The purpose of this text is not to serve as a primer for understanding complex samples, and thus readers interested in learning more about complex survey designs are referred to any number of excellent resources (Hahs-Vaughn, 2005; Hahs-Vaughn, McWayne, Bulotskey-Shearer, Wen, & Faria, 2011a, 2011b; Lee, Forthofer, & Lorimor, 1989; Skinner, Holt, & Smith, 1989). Additionally, MANOVA is complex enough, without adding the complex sampling design in the mix! Thus, the applications in the textbook do not illustrate how to adjust for the complex sample design. As such, the results that we see should not be interpreted to represent any larger population but only that select sample of individuals who actually completed the survey. I want to stress that the reason why the sampling design has not been illustrated in the textbook applications is that the point of this portion of the textbook is to illustrate how to use statistical software to generate various procedures and how to interpret the output, not to ensure the results are representative of the intended population. Please do not let this discount or diminish the need to apply this critical step in your own analyses when using complex survey data, as quite a large body of research exists that describes the importance of effectively analyzing complex samples and provides evidence of biased results when the complex sample design is not addressed in the analyses (Hahs-Vaughn, 2005, 2006a, 2006b; Hahs-Vaughn et al., 2011a, 2011b; Kish & Frankel, 1973, 1974; Korn & Graubard, 1995; Lee et al., 1989; Lumley, 2004; Pfeffermann, 1993; Skinner et al., 1989).

Getting back to our specific data file . . . We are using the **ECLSK_MANOVA_N1393.sav** file. This is data from the U.S., and the data file has been delimited to include only children who meet the following criteria:

- a) nonzero third-grade panel weight [$C5CW0 > 0$];
- b) child did not receive special education services in waves 2, 4, or 5 [$F2SPECS = 2$ and $F4SPECS = 2$ and $F5SPECS = 2$];
- c) child did not change schools during the first five waves of data collection [$FKCHGSCH = 0$ and $R3R2SCHG = 1$ and $R4R2SCHG = 1$ and $R4R3SCHG = 1$ and $R5R4SCHG = 1$];
- d) child attended a public school [$S5SCTYP = 4$];
- e) mother's language to child is English ('never speaks non-English') [$P1HMLANG = 1$];
- f) home language during kindergarten and first grade was English [$WKLNGST = 2$ and $W1LANGST = 2$] ($n = 1,396$).

(Note: Because of the massive size of the full 1998–1999 ECKLS-K kindergarten through grade 8 child file, we haven't made it available from the textbook's companion

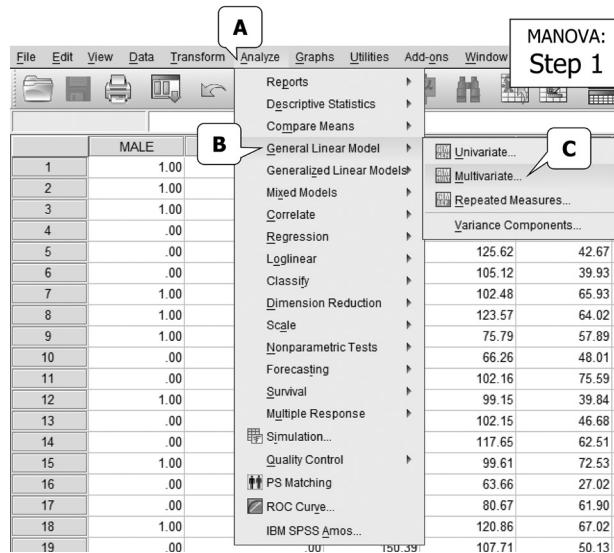
website. However, you can access the full file, along with teacher and school data files, from NCES at <http://nces.ed.gov/ecls/dataproducts.asp>.)

Before we run the data, it's always important to examine frequency distributions of the variables that will be used in the model to assess missing data, potential data entry problems, and similar. With this data, we have very little missing data after we've made the delimitations (reading IRT missing = 1 and science IRT missing = 2), and thus I've taken the liberty to perform listwise deletion on the missing items (resulting in $n = 1,393$).

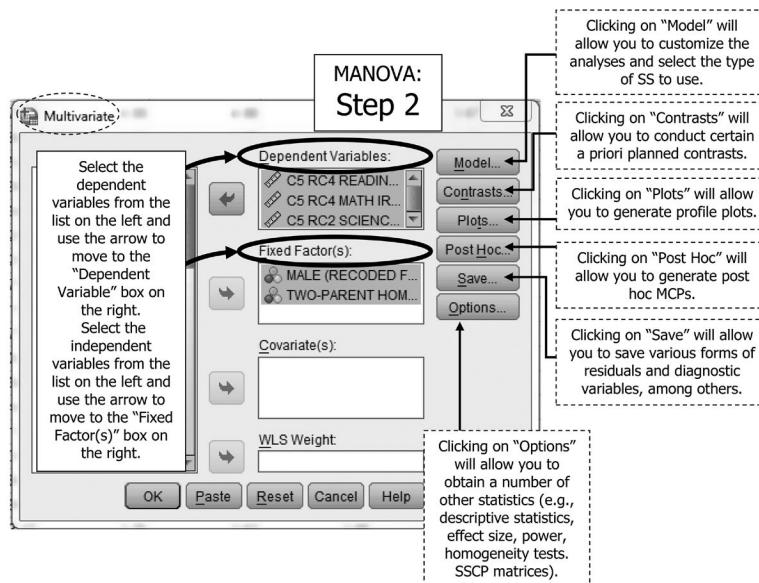
Let's look at the data. For the MANOVA illustration, we'll be working with the first five variables (variables 1–5 in your SPSS file). A few other variables are retained in this data file, however, in the event you want to test other models. The first five variables are the variables we'll use to illustrate MANOVA. This is followed by other variables (both categorical and continuous, variables 6–36) which have been left so you can test other models. The remaining variables in the dataset are those that were used to delimit the sample (variables 37–50), recode gender and family type (variables 51–52) and finally, the child ID variable in case you are interested in merging variables from the full dataset with this smaller, delimited file. Each row in the dataset still represents one child. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the children were measured. For the MANOVA illustration, we will work with the first six variables in the dataset.

	MALE	TWOPARENT_GR3	C5R4RSCL	C5R4MSCL	C5R2SSCL
1	1.00	1.00	144.29	116.95	68.73
2	1.00	1.00	82.91	61.47	42.58
3	1.00	1.00	123.92	104.03	61.07
4	.00	.00	105.29	71.69	35.68
5	.00	1.00	119.42	125.62	42.67
6	.00	1.00	135.10	105.12	39.93
7	1.00	.00	127.38	102.48	65.93
8	The first independent variable is Male (where males were coded as 1 and females as 0)		16	The dependent variables are spring 3rd grade reading, mathematics, and science IRT scale scores	
9	The second independent variable is an indicator of family type (where two-parent home in 3rd grade = 1 and non two-parent home = 0)		11	\$4.02	
10			12	57.89	
11			13	48.01	
12			14	75.59	
13			15	127.84	
14			16	99.15	
15			17	39.84	
16			18	126.42	
17			19	102.15	
18			20	46.68	
19			21	163.31	
20			22	117.65	
21			23	62.51	
22			24	99.61	
23			25	72.53	
24			26		
25			27		
26			28		
27			29		
28			30		
29			31		
30			32		
31			33		
32			34		
33			35		
34			36		
35			37		
36			38		
37			39		
38			40		
39			41		
40			42		
41			43		
42			44		
43			45		
44			46		
45			47		
46			48		
47			49		
48			50		

Step 1. To conduct a MANOVA, go to “Analyze” in the top pull-down menu, then select “General Linear Model,” and then select “Multivariate.” Following the screenshot below (Step 1) produces the “Multivariate” dialog box.

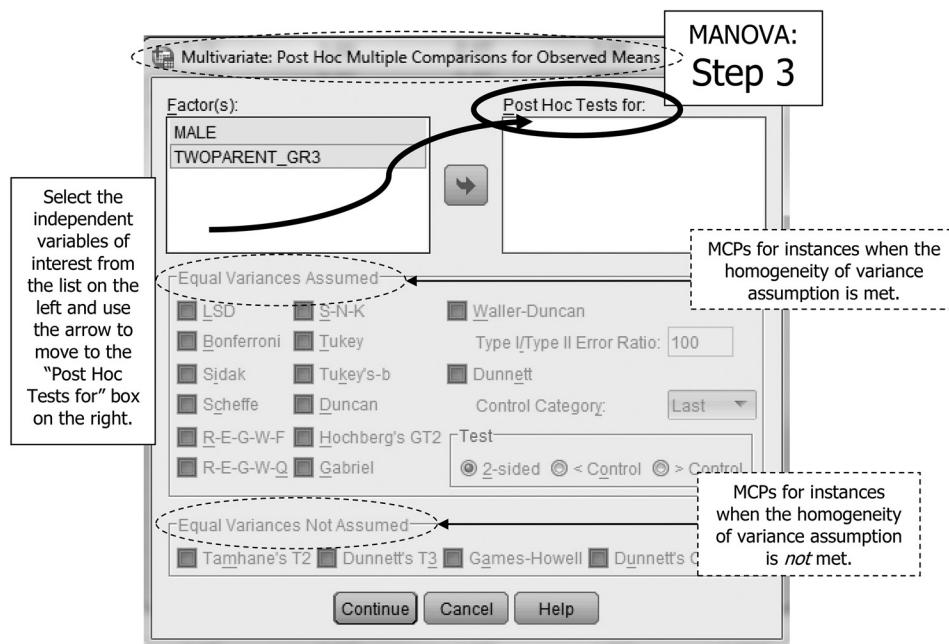


Step 2. Click the dependent variables (e.g., reading, math, science scores) and move them into the “Dependent Variables” box by clicking the arrow button (see screenshot Step 2). Click the first independent variable (e.g., male) and move it into the “Fixed Factor(s)” box by clicking the arrow button. Follow this same step to move the second independent variable into the “Fixed Factors” box.

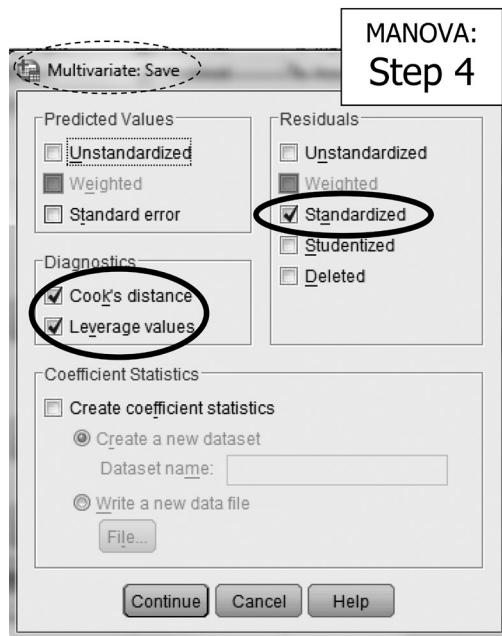


Step 3. From the “Multivariate” dialog box (see screenshot Step 2), click on “Post Hoc” to display the Multivariate: Post Hoc Multiple Comparisons for Observed Means dialog box (see screenshot Step 3). Various post hoc multiple comparison procedures

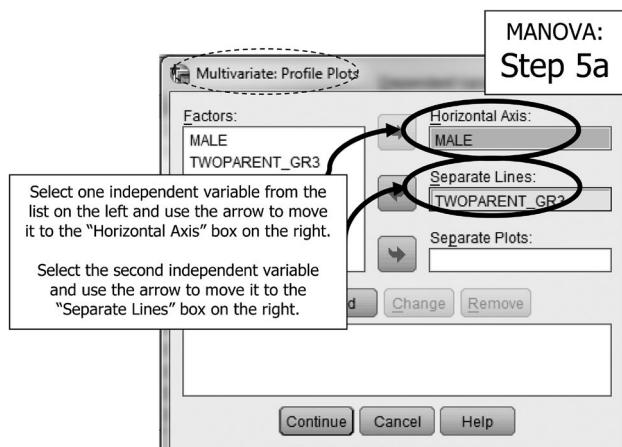
(MCPs) can be selected from this site. Because both of our variables are binary, even if we select post hoc tests they won't be displayed, as there are fewer than three groups in each of our independent variables. We are interested in following up any statistically significant findings with discriminant analysis; thus, this is presented simply to see the options available. If you were to make a selection, you would click on the names of the independent variables in the 'Factor(s)' list box in the top left (e.g., 'MALE' and 'TWOPARENT_GR3') and move them to the 'Post Hoc Tests for' box in the top right by clicking on the arrow key. Check an appropriate MCP for your situation by placing a checkmark in the box next to the desired MCP. Click on "Continue" to return to the original dialog box.

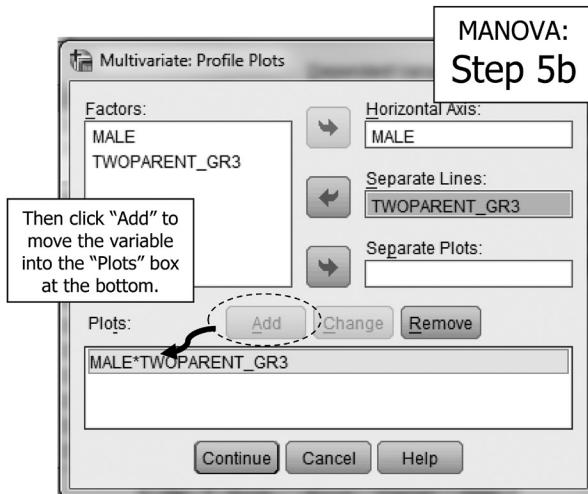


Step 4. From the "Multivariate" dialog box (see screenshot Step 2), click on "Save" to select those elements that you want to save (see screenshots Steps 2a through 2d). In our case, we want to save the Standardized residuals, Cook's distance, and Leverage values, which will be used later to examine the extent to which normality and independence are met. Cook's distance provides a measure of the change in the residuals with the exclusion of a case. Small Cook's distance values are desirable. Uncentered leverage values provide a measure of influence of each case on model fit. Click on "Continue" to return to the original dialog box. From the "Multivariate" dialog box, click on "OK" to return to generate the output.

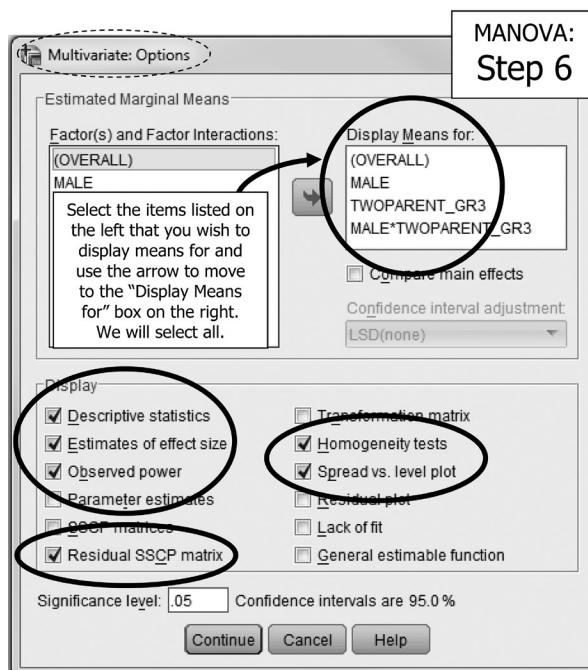


Steps 5a and 5b. From the “Multivariate” dialog box (see screenshot Step 2), click on “Plots” to display the Multivariate: Profile Plots dialog box (see screenshot Step 5a). Click the first independent variable (e.g., ‘MALE’) and move it into the “Horizontal Axis” box by clicking the arrow button (see screenshot Step 5a). [*Tip: Placing the independent variable that has the most categories or levels on the horizontal axis of the profile plots will make for easier interpretation of the graph. In our case, both are binary so this is a moot point.*] Then click the second independent variable (e.g., ‘TWOPARENT_GR3’) and move it into the “Separate Lines” box by clicking the arrow button (see screenshot Step 5a). Then click on “Add” to move the variable into the “Plots” box at the bottom of the dialog box (see screenshot Step 5b). Click on “Continue” to return to the original dialog box.





Step 6. From the “Multivariate” dialog box (see screenshot Step 2), click on “Options” to display the Multivariate: Options dialog box (see screenshot Step 6). Check the following selections: “Descriptive Statistics,” “Estimates of effect size,” “Observed power,” “Residual SSCP matrix,” “Homogeneity tests,” and “Spread versus level plot.” We also want to display means for all the factors and interactions by clicking the variable names in the box on the top left and moving to the Display Means for box. Click on “Continue” to return to the original dialog box. From the “Multivariate” dialog box, click on “OK” to return to generate the output.



Interpreting the output. Annotated results are presented in Table 6.3.

TABLE 6.3

SPSS Results for the Two-Factor MANOVA

Between-Subjects Factors				
		Value Label	N	
MALE (RECODED FROM 'GENDER')	.00 1.00	FEMALE MALE	714 679	
TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')	.00 1.00	NON-TWO PARENT HOME TWO-PARENT HOME	415 978	

The table labeled "Between-Subjects Factors" provides sample sizes for each of the categories of the independent variables (recall that the independent variables are the 'between subjects factors').

Descriptive Statistics					
	MALE (RECODED FROM 'GENDER')	TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')	Mean	Std. Deviation	N
C5 RC4		NON-TWO PARENT HOME	124.4289	25.32221	235
READING IRT SCALE SCORE	MALE	TWO-PARENT HOME	136.5994	25.82007	479
		Total	132.5937	26.27045	714
C5 RC4 MATH IRT SCALE SCORE	MALE	NON-TWO PARENT HOME	120.2862	25.62189	180
		TWO-PARENT HOME	134.6827	25.79636	499
		Total	130.8662	26.50547	679
		NON-TWO PARENT HOME	122.6320	25.50473	415
	Total	TWO-PARENT HOME	135.6214	25.81257	978
		Total	131.7517	26.38992	1393
		NON-TWO PARENT HOME	93.1106	22.04189	235
	FEMALE	TWO-PARENT HOME	103.6039	22.83838	479
		Total	100.1502	23.09708	714
C5 RC2 SCIENCE IRT SCALE SCORE	MALE	NON-TWO PARENT HOME	98.0022	23.49151	180
		TWO-PARENT HOME	109.8919	23.44908	499
		Total	106.7400	24.02406	679
		NON-TWO PARENT HOME	95.2323	22.78380	415
	Total	TWO-PARENT HOME	106.8122	23.35289	978
		Total	103.3623	23.77445	1393
		NON-TWO PARENT HOME	46.5627	13.41577	235
	FEMALE	TWO-PARENT HOME	53.9529	13.85336	479
		Total	51.5205	14.13528	714
C5 RC2 SCIENCE IRT SCALE SCORE	MALE	NON-TWO PARENT HOME	49.3297	14.45817	180
		TWO-PARENT HOME	57.5127	14.05864	499
		Total	55.3434	14.60913	679
		NON-TWO PARENT HOME	47.7628	13.92824	415
	Total	TWO-PARENT HOME	55.7692	14.06448	978
		Total	53.3839	14.48967	1393

■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA

Box's Test of Equality of Covariance Matrices ^a	
Box's M	14.019
F	.775
df1	18
df2	2220072, 108
Sig.	.732

Tests the null hypothesis that the observed covariance matrices of the dependent variables are equal across groups.

a. Design: Intercept + MALE +
TWOPARENT_GR3 + MALE *
TWOPARENT_GR3

The *F* test (and associated p value) for Box's Test for Equality of Covariance Matrices is reviewed to determine if equal covariances can be assumed. In this case, we meet the assumption (as p is greater than α). The variance-covariance matrices of the dependent variables are relatively similar across all levels of the independent variables.

Bartlett's Test of Sphericity ^a	
Likelihood Ratio	.000
Approx. Chi-Square	2622.961
df	5
Sig.	.000

Tests the null hypothesis that the residual covariance matrix is proportional to an identity matrix.

a. Design: Intercept + MALE +
TWOPARENT_GR3 + MALE *
TWOPARENT_GR3

The chi-square test (and associated p value) for Bartlett's Test of Sphericity is reviewed to determine if there are significant correlations among the dependent variables after controlling for the independent variables. In this case, we do not meet the assumption (as p is less than α). The null hypothesis is that the residual covariance matrix is equal to an identity matrix. Good models have residual correlations near zero indicating random residuals. In this case, statistical significance indicates that there are uncorrelated dependent variables after controlling for the independent variables.

The caveat: Bartlett's Test of Sphericity is sensitive to large samples (i.e., may provide erroneous results) and thus best used with small sample sizes.

■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA

Multivariate Tests ^a									
Effect		Value	F	Hypo. df	Error df	Sig.	Partial Eta Squared	Noncent. Parameter	Observed Power ^c
MALE	Pillai's Trace	.957	10208.596 ^b	3.000	1387.000	.000	.957	30625.788	1.000
	Intercept	.043	10208.596 ^b	3.000	1387.000	.000	.957	30625.788	1.000
	Wilks' Lambda								
	Hotelling's Trace	22.081	10208.596 ^b	3.000	1387.000	.000	.957	30625.788	1.000
	Roy's Largest Root	22.081	10208.596 ^b	3.000	1387.000	.000	.957	30625.788	1.000
	Pillai's Trace	.066	32.561 ^b	3.000	1387.000	.000	.066	97.683	1.000
	Wilks' Lambda	.934	32.561 ^b	3.000	1387.000	.000	.066	97.683	1.000
	Hotelling's Trace	.070	32.561 ^b	3.000	1387.000	.000	.066	97.683	1.000
	Roy's Largest Root	.070	32.561 ^b	3.000	1387.000	.000	.066	97.683	1.000
	Pillai's Trace	.067	33.382 ^b	3.000	1387.000	.000	.067	100.146	1.000
	TWOPAR	Wilks' Lambda	.933	33.382 ^b	3.000	1387.000	.000	.067	100.146
	ENT_GR3	Hotelling's Trace	.072	33.382 ^b	3.000	1387.000	.000	.067	100.146
MALE * TWOPAR	Roy's Largest Root	.072	33.382 ^b	3.000	1387.000	.000	.067	100.146	1.000
	Pillai's Trace	.000	.181 ^b	3.000	1387.000	.909	.000	.544	.084
	Wilks' Lambda	1.000	.181 ^b	3.000	1387.000	.909	.000	.544	.084
	Hotelling's Trace	.000	.181 ^b	3.000	1387.000	.909	.000	.544	.084
ENT_GR3	Roy's Largest Root	.000	.181 ^b	3.000	1387.000	.909	.000	.544	.084

a. Design: Intercept + MALE + TWOPARENT_GR3 + MALE * TWOPARENT_GR3

b. Exact statistic

c. Computed using alpha = .05

Pillai's Trace represents the sum of the proportion of explained variance on all the discriminant functions similar to the ratio of SS between to SS total.

Wilks' Lambda represents the ratio of error variance to total variance for each variate.

Hotellings's Trace compares directly to the *F* ratio in ANOVA, and is the sum of the ratio of between variance to error variance for each variate (i.e., the sum of eigenvalues for each variate).

Roy's Largest Root represents the maximum possible between-group differences and is conceptually the same as the *F* ratio in ANOVA, representing the proportion of explained to unexplained variance for the first discriminant function. In many cases, this test is the most powerful.

In terms of power, all four have similar power even with small sample sizes. All four test statistics are relatively robust to violations of multivariate normality (with some exceptions: Roy's root is less robust if there is a platykurtic distribution AND Roy's root is less robust when the assumption of homogeneity of covariances is *not* met).

When there is a balanced design, Pillai's is most robust to assumption violations. With an unbalanced design, Pillai's is less robust when homogeneity of variance-covariances is not met.

Multivariate test results examine the extent that each independent variable and interaction are statistically significant for at least one dependent variable. The *F* test takes into account the covariances (i.e., sum of cross products) and the group means (i.e., sum of squares between and within groups)

Multivariate Tests^a

Regardless of which test we select, there are statistically significant main effects for MALE and TWOPARENT, but no statistically significant interaction.

This means that there are simultaneous mean differences in the outcomes when comparing girls to boys and when comparing children from two-parent homes to other home environments. But there are not differential simultaneous mean differences for girls who come from two parent vs. other types of homes (or likewise boys from two parent vs. other types of homes).

$$\begin{aligned} \text{partial } \eta^2 &= \frac{df_1 F}{df_1 F + df_2} \\ &= \frac{3(32.561)}{4(32.561) + 1387} = \frac{97.683}{1484.683} = .066 \end{aligned}$$

Which to report? Wilk's Lambda if assumptions are met. Pillai's trace if assumptions are violated.

■ TABLE 6.3 (continued)

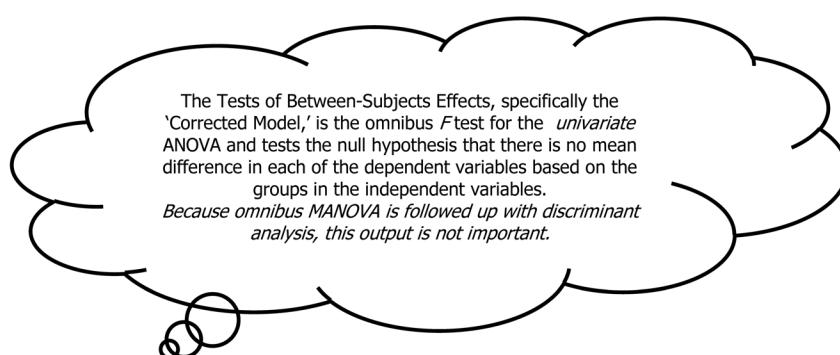
SPSS Results for the Two-Factor MANOVA

Levene's Test of Equality of Error Variances ^a				
	F	df1	df2	Sig.
C5 RC4 READING IRT SCALE SCORE	.503	3	1389	.680
C5 RC4 MATH IRT SCALE SCORE	.482	3	1389	.695
C5 RC2 SCIENCE IRT SCALE SCORE	1.397	3	1389	.242

Tests the null hypothesis that the error variance of the dependent variable is equal across groups.

a. Design: Intercept + MALE + TWOPARENT_GR3 + MALE * TWOPARENT_GR3

The *F* test (and associated *p* value) for Levene's Test for the **univariate** Equality of Error Variances. These results are reviewed to determine if equal variances can be assumed *only* if univariate ANOVAs are being used as a follow-up to statistically significant omnibus MANOVA. In this case, we would have met the assumption (as *p* is greater than α). Note that *df1* is calculated as $(JK - 1)$ and *df2* is calculated as $(N - JK)$.



Tests of Between-Subjects Effects									
Source	Dependent Variable	Type III Sum of Squares	df	Mean Square	F	Sig.	Partial Eta Squared	Noncent Parameter	Observed Power ^d
Corrected Model	C5 RC4 READING IRT SCALE SCORE	51807.314 ^a	3	17269.105	26.140	.000	.053	78.421	1.000
	C5 RC4 MATH IRT SCALE SCORE	51172.293 ^b	3	17057.431	32.208	.000	.065	96.624	1.000
	C5 RC2 SCIENCE IRT SCALE SCORE	22554.259 ^c	3	7518.086	38.720	.000	.077	116.160	1.000
Intercept	C5 RC4 READING IRT SCALE SCORE	19151362.03	1	19151362.03	28989.369	.000	.954	28989.369	1.000
	C5 RC4 MATH IRT SCALE SCORE	11775385.99	1	11775385.99	22234.321	.000	.941	22234.321	1.000
	C5 RC2 SCIENCE IRT SCALE SCORE	3092761.100	1	3092761.100	15928.417	.000	.920	15928.417	1.000
MALE	C5 RC4 READING IRT SCALE SCORE	2640.957	1	2640.957	3.998	.046	.003	3.998	.515
	C5 RC4 MATH IRT SCALE SCORE	8990.001	1	8990.001	16.975	.000	.012	16.975	.985
	C5 RC2 SCIENCE IRT SCALE SCORE	2879.181	1	2879.181	14.828	.000	.011	14.828	.970
TWOPARENT_GR3	C5 RC4 READING IRT SCALE SCORE	50768.032	1	50768.032	76.847	.000	.052	76.847	1.000
	C5 RC4 MATH IRT SCALE SCORE	36036.124	1	36036.124	68.044	.000	.047	68.044	1.000
	C5 RC2 SCIENCE IRT SCALE SCORE	17444.490	1	17444.490	89.843	.000	.061	89.843	1.000
MALE * TWOPARENT_GR3	C5 RC4 READING IRT SCALE SCORE	356.376	1	356.376	.539	.463	.000	.539	.114
	C5 RC4 MATH IRT SCALE SCORE	140.252	1	140.252	.265	.607	.000	.265	.081
	C5 RC2 SCIENCE IRT SCALE SCORE	45.219	1	45.219	.233	.629	.000	.233	.077
Error	C5 RC4 READING IRT SCALE SCORE	917620.603	1389	660.634					
	C5 RC4 MATH IRT SCALE SCORE	735619.991	1389	529.604					
	C5 RC2 SCIENCE IRT SCALE SCORE	269696.926	1389	194.166					

■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA

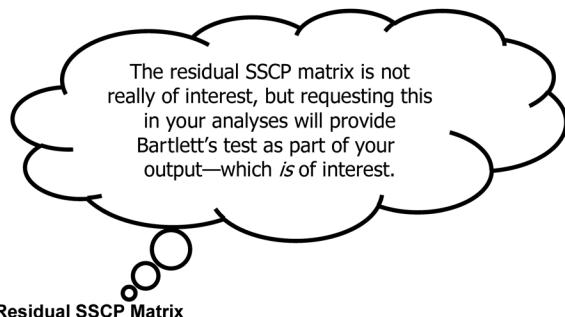
Total	C5 RC4 READING IRT SCALE SCORE	25149820.30	1393						
	C5 RC4 MATH IRT SCALE SCORE	15669276.21	1393						
	C5 RC2 SCIENCE IRT SCALE SCORE	4262085.508	1393						
Corrected Total	C5 RC4 READING IRT SCALE SCORE	969427.917	1392						
	C5 RC4 MATH IRT SCALE SCORE	786792.284	1392						
	C5 RC2 SCIENCE IRT SCALE SCORE	292251.185	1392						

a. R Squared = .053 (Adjusted R Squared = .051)

b. R Squared = .065 (Adjusted R Squared = .063)

c. R Squared = .077 (Adjusted R Squared = .075)

d. Computed using alpha = .05



Residual SSCP Matrix

		C5 RC4 READING IRT SCALE SCORE	C5 RC4 MATH IRT SCALE SCORE	C5 RC2 SCIENCE IRT SCALE SCORE
Sum-of-Squares and Cross-Products	C5 RC4 READING IRT SCALE SCORE	917620.603	582145.544	359751.635
	C5 RC4 MATH IRT SCALE SCORE	582145.544	735619.991	289680.258
	C5 RC2 SCIENCE IRT SCALE SCORE	359751.635	289680.258	269696.926
Covariance	C5 RC4 READING IRT SCALE SCORE	660.634	419.111	259.000
	C5 RC4 MATH IRT SCALE SCORE	419.111	529.604	208.553
	C5 RC2 SCIENCE IRT SCALE SCORE	259.000	208.553	194.166
Correlation	C5 RC4 READING IRT SCALE SCORE	1.000	.709	.723
	C5 RC4 MATH IRT SCALE SCORE	.709	1.000	.650
	C5 RC2 SCIENCE IRT SCALE SCORE	.723	.650	1.000

Based on Type III Sum of Squares

■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA

Estimated Marginal Means

The 'Grand Mean' represents the overall mean, regardless of group membership, for each of the dependent variables. The 95% CI represents the CI of the grand mean.

Dependent Variable	Mean	Std. Error	95% Confidence Interval	
			Lower Bound	Upper Bound
			1. Grand Mean	
C5 RC4 READING IRT SCALE SCORE	128.999	.758	127.513	130.486
C5 RC4 MATH IRT SCALE SCORE	101.152	.678	99.821	102.483
C5 RC2 SCIENCE IRT SCALE SCORE	51.839	.411	51.034	52.645

2. MALE (RECODED FROM 'GENDER')

Dependent Variable	MALE (RECODED FROM 'GENDER')	Mean	Std. Error	95% Confidence Interval	
				Lower Bound	Upper Bound
C5 RC4 READING IRT SCALE SCORE	FEMALE	130.514	1.024	128.506	132.522
C5 RC4 MATH IRT SCALE SCORE	MALE	127.484	1.117	125.293	129.676
C5 RC2 SCIENCE IRT SCALE SCORE	FEMALE	98.357	.916	96.560	100.155
C5 RC2 SCIENCE IRT SCALE SCORE	MALE	103.947	1.000	101.984	105.910
C5 RC4 READING IRT SCALE SCORE	FEMALE	50.258	.555	49.169	51.346
C5 RC4 READING IRT SCALE SCORE	MALE	53.421	.606	52.233	54.609

The table labeled 'TWO-PARENT' provides descriptive statistics for each of the categories of the second independent variable for each dependent variable. In addition to means, the *SE* and 95% CI of the means are reported.

The table labeled 'MALE' provides descriptive statistics for each of the categories of the first independent variable by each dependent variable. In addition to means, the *SE* and 95% CI of the means are reported.

3. TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')

Dependent Variable	TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')	Mean	Std. Error	95% Confidence Interval	
				Lower Bound	Upper Bound
C5 RC4 READING IRT SCALE SCORE	NON-TWO PARENT HOME	122.358	1.273	119.860	124.855
C5 RC4 MATH IRT SCALE SCORE	TWO-PARENT HOME	135.641	.822	134.028	137.254
C5 RC2 SCIENCE IRT SCALE SCORE	NON-TWO PARENT HOME	95.556	1.140	93.321	97.792
C5 RC2 SCIENCE IRT SCALE SCORE	TWO-PARENT HOME	106.748	.736	105.304	108.192
C5 RC4 READING IRT SCALE SCORE	NON-TWO PARENT HOME	47.946	.690	46.592	49.300
C5 RC4 READING IRT SCALE SCORE	TWO-PARENT HOME	55.733	.446	54.859	56.607

■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA

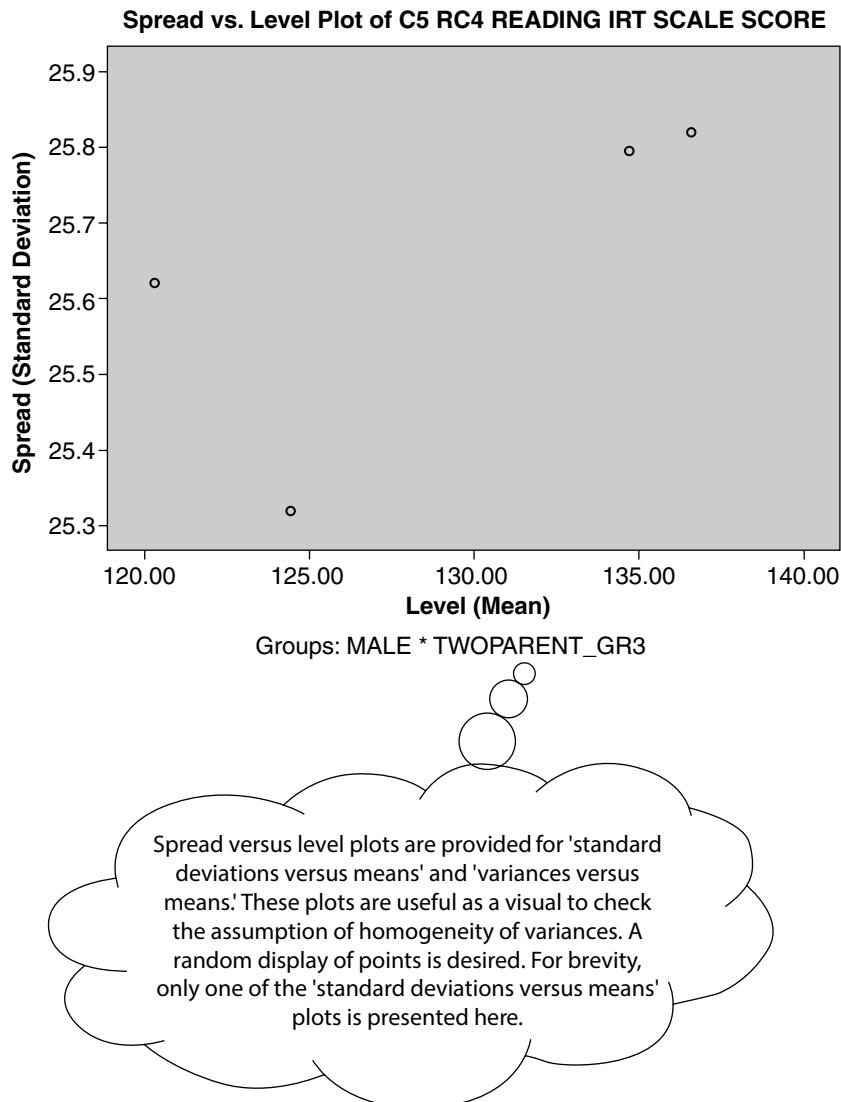
4. MALE (RECODED FROM 'GENDER') * TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')

Dependent Variable	MALE (RECODED FROM 'GENDER')	TWO-PARENT HOME DURING 3RD GRADE (RECODED FROM 'P5HFAMIL')	Mean	Std. Error	95% Confidence Interval	
					Lower Bound	Upper Bound
C5 RC4 READING	FEMALE	NON-TWO PARENT HOME	124.429	1.677	121.140	127.718
		TWO-PARENT HOME	136.599	1.174	134.296	138.903
IRT SCALE SCORE	MALE	NON-TWO PARENT HOME	120.286	1.916	116.528	124.044
		TWO-PARENT HOME	134.683	1.151	132.426	136.940
C5 RC4 MATH IRT	FEMALE	NON-TWO PARENT HOME	93.111	1.501	90.166	96.055
		TWO-PARENT HOME	103.604	1.051	101.541	105.667
SCALE SCORE	MALE	NON-TWO PARENT HOME	98.002	1.715	94.637	101.367
		TWO-PARENT HOME	109.892	1.030	107.871	111.913
C5 RC2 SCIENCE	FEMALE	NON-TWO PARENT HOME	46.563	.909	44.780	48.346
		TWO-PARENT HOME	53.953	.637	52.704	55.202
IRT SCALE SCORE	MALE	NON-TWO PARENT HOME	49.330	1.039	47.292	51.367
		TWO-PARENT HOME	57.513	.624	56.289	58.736

The table labeled '**MALE*TWO-PARENT**' provides descriptive statistics for each dependent variable for each of the categories of the first independent variable by the second independent variable (i.e., cell means representing the interaction of the independent variables). In addition to means, the *SE* and 95% CI of the means are reported.

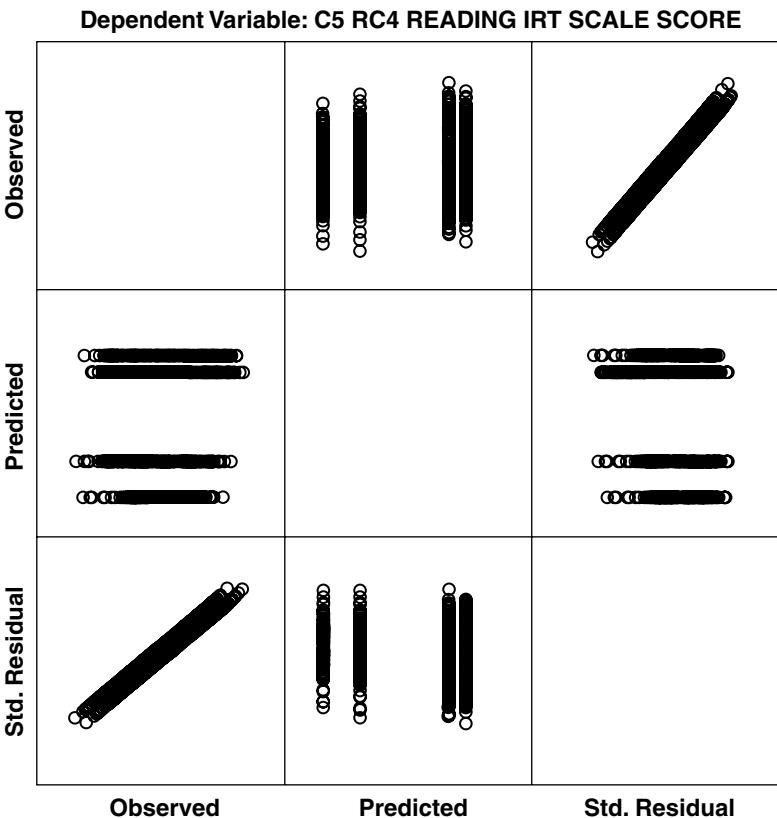
■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA



■ TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA



Model: Intercept + MALE + TWOPARENT_GR3 + MALE * TWOPARENT_GR3

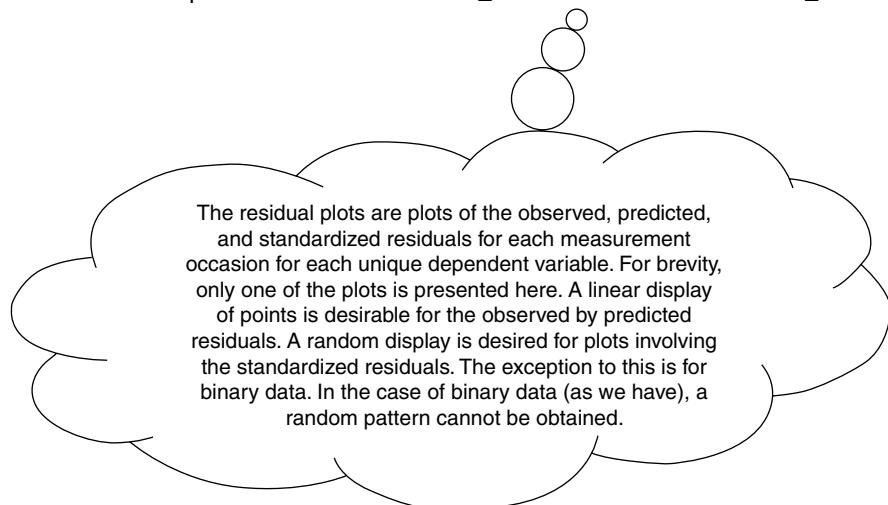
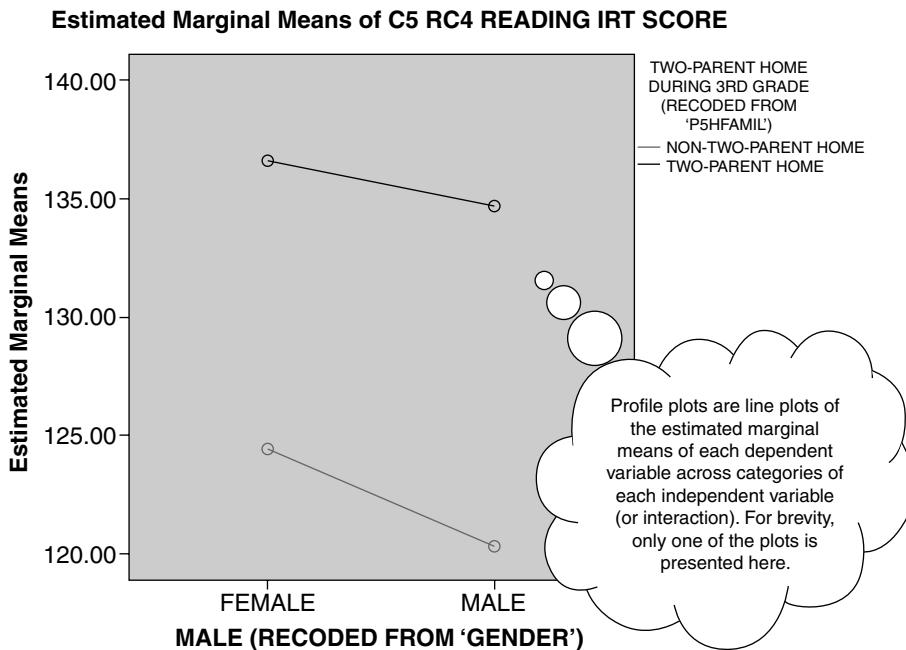


TABLE 6.3 (continued)

SPSS Results for the Two-Factor MANOVA



6.3.2 Computing Repeated Measures MANOVA Using SPSS

We will use SPSS to compute repeated measures MANOVA using data from the same secondary data source as applied with the factorial MANOVA: the Early Childhood Longitudinal Study Kindergarten Class of 1998–1999 (ECLS-K) (<http://nces.ed.gov/ecls/kindergarten.asp>), available through the U.S. Department of Education National Center for Education Statistics Institute of Education Sciences (NCES IES). Specifically, we are using the **ECLSK_REPEATMANOVA_N1344.sav** file. This is data from the U.S., and the data file has been delimited to include only children who meet the following criteria:

- nonzero third-grade panel weight [$C5CW0 > 0$];
- child did not receive special education services in waves 2, 4, or 5 [$F2SPECS = 2$ and $F4SPECS = 2$ and $F5SPECS = 2$];
- child did not change schools during the first five waves of data collection [$FKCHGSCH = 0$ and $R3R2SCHG = 1$ and $R4R2SCHG = 1$ and $R4R3SCHG = 1$ and $R5R4SCHG = 1$];

- d) child attended a public school [S5SCTYP = 4];
- e) mother's language to child is English ('never speaks non-English') [P1HMLANG = 1];
- f) home language during kindergarten and first grade was English [WKLANGST = 2 and W1LANGST = 2] ($n = 1,396$).

(Note: Because of the massive size of the full 1998–1999 ECKLS-K kindergarten through grade 8 child file, we haven't made it available from the textbook's companion website. However, you can access the full file, along with teacher and school data files, from NCES at <http://nces.ed.gov/ecls/dataproducts.asp>.)

Before we run the data, it's always important to examine frequency distributions of the variables that will be used in the model to assess missing data, potential data entry problems, and similar. With this example, we are now looking at five measurement occasions (fall and spring kindergarten; spring first grade; spring third grade; and spring fifth grade). There is very little missing data after we've made the delimitations (maximum of about 2% in fall kindergarten reading), and thus I've taken the liberty to perform listwise deletion on the missing items (resulting in $n = 1,344$).

Let's look at the data. For the repeated measures MANOVA illustration, we'll be working with the first set of variables. The repeated measures reading IRT scale scores are variables 1–5 in the SPSS file. The repeated measures mathematics IRT scale scores are variables 6–10 in the SPSS file. The between-subjects factor of 'male' is variable 16 in the SPSS dataset. A few other variables are retained in this data file, however, if you want to test other models (both within and between). The remaining variables in the dataset are those that were used to delimit the sample (variables 21–34), recode gender and family type (variables 35–36) and finally, the child ID variable in case you are interested in merging variables from the full dataset with this smaller, delimited file. Each row in the dataset still represents one child. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the children were measured. For the repeated measures MANOVA illustration, we will work with the repeated variables in the first section of the dataset. These are the within factors.

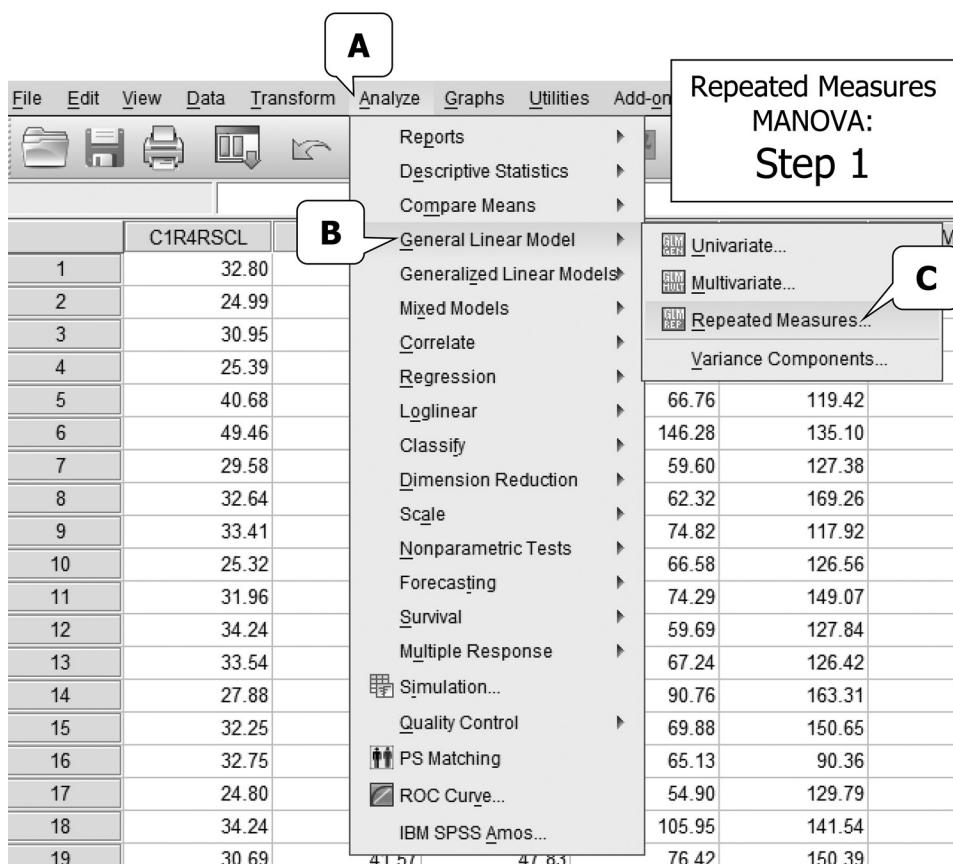
	C1R5CL	C2R5CL	C3R5CL	C4R5CL	C5R5CL	C1R4MSCL	C2R4MSCL	C3R4MSCL	C4R4MSCL	C5R4MSCL	C1R3SCAL	C2R3SCAL	C3R3SCAL	C4R3SCAL	C5R3SCAL
1	32.80	35.00	37.53	63.55	144.29	26.73	28.47	42.85	68.64	116.95	28.120	26.651	33.346	41.928	61.73
2	24.99	30.04	38.10	56.72	82.91	17.45	21.07	22.33	42.28	61.47	14.763	24.065	31.603	27.179	42.58
3	30.95	38.93	51.53	60.86	123.92	29.58	34.86	39.24	60.43	104.03	32.209	35.449	33.271	38.988	61.07
4	25.39	36.22	41.06	59.43	105.29	17.83	32.67	30.14	37.60	71.69	22.709	24.117	29.776	32.338	35.68
5	40.68	44.63	47.11	66.76	119.42	28.63	34.76	44.71	63.39	125.62	22.995	29.668	28.65	36.978	42.67
6	49.46	76.16	93.10	145.28	135.10	37.81	37.67	49.37	74.38	105.12	24.432	31.410	29.51	35.687	39.93
7	29.58	45.33	48	59.60	127.38	28.78	44.55	47.38	63.25	102.48	26.778	31.592	36.21	35.654	65.93
8	32.64	36.09	52	62.32	169.26	36.24	51.04	51.82	66.40	123.57	28.542	41.508	36.68	44.520	64.02
9	33.41	48.25	56.0	74.82	117.92	24.67	31.57	35.04	50.28	75.79	26.851	28.133	27.65	39.732	57.89
10	25.32	38.40	41.0	66.58	126.56	19.36	28.09	30.30	40.64	66.26	22.947	26.200	28.10	31.803	48.01

The **within factors** are reading and mathematics IRT scale scores measured in waves 1, 2, 3, 4, and 5.

In the event you want to explore some additional repeated measures MANOVA models, there is an additional set of repeated measures, general knowledge IRT scale scores, also measured in waves 1, 2, 3, 4, and 5.

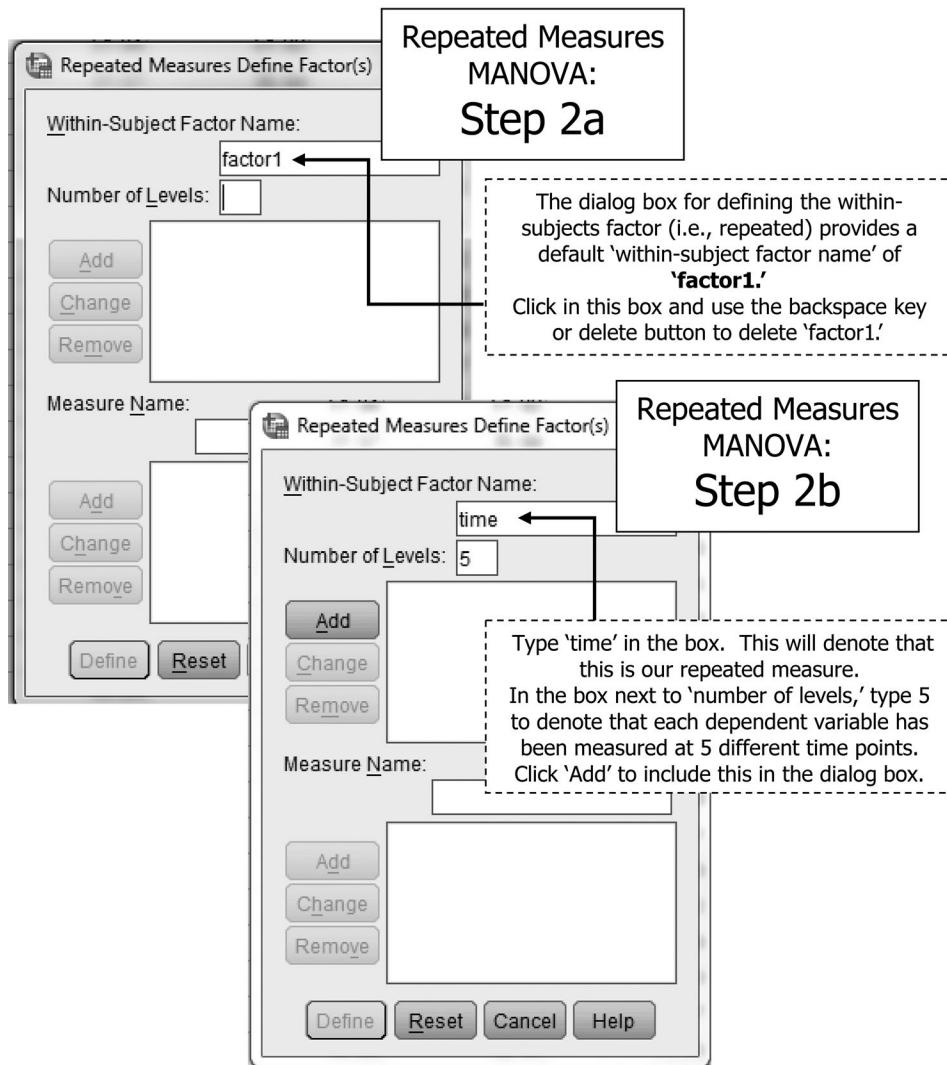
We will also work with ‘male’ as the between-subjects factor. With this variable, ‘males’ were coded as 1 and ‘females’ were coded as 0.

Step 1. To conduct a repeated measures MANOVA, go to “Analyze” in the top pull-down menu, then select “General Linear Model,” and then select “Repeated Measures.” Following the screenshot below (see repeated measures MANOVA screenshot Step 1) produces the “Repeated Measures” dialog box.

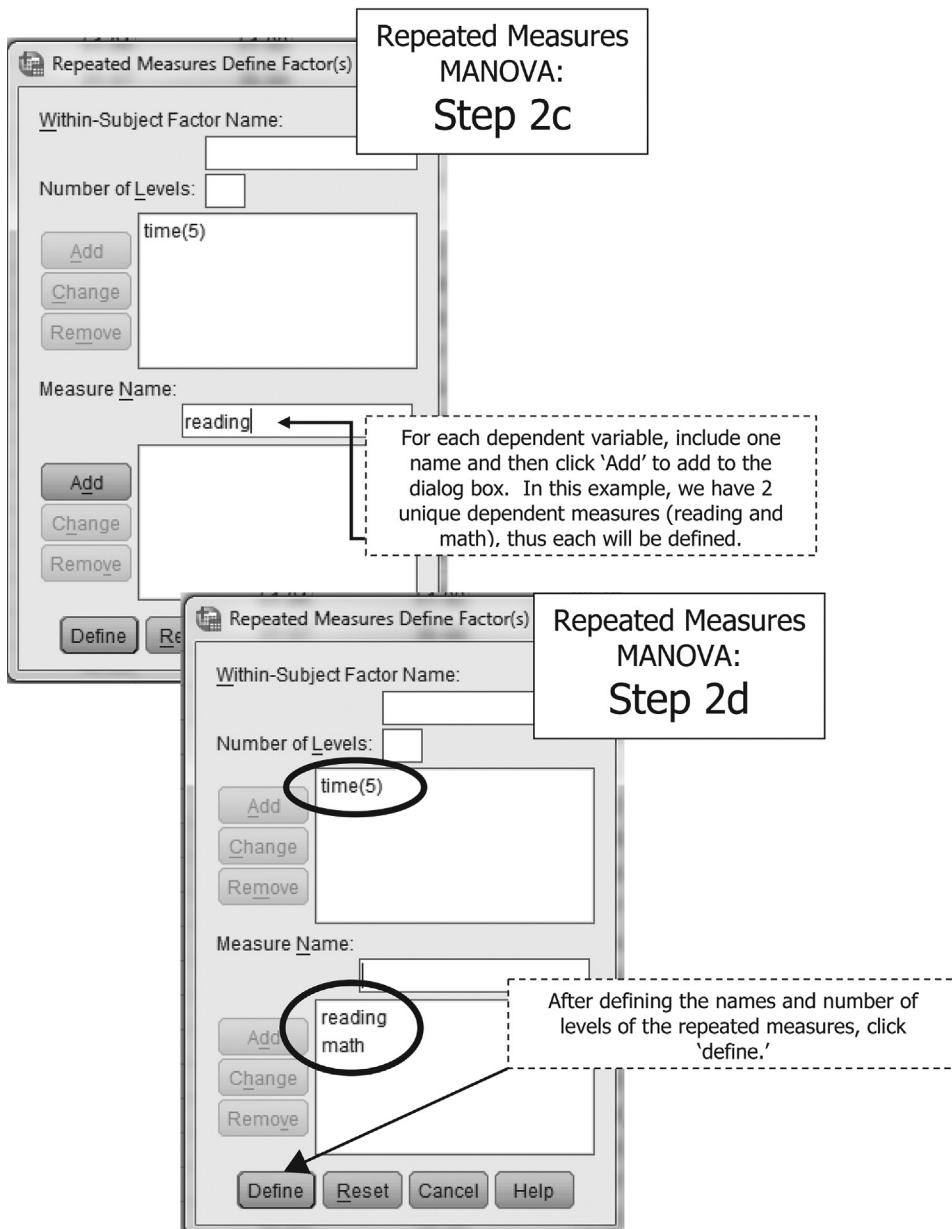


Step 2. From the ‘repeated measures define factor(s)’ dialog box, delete ‘factor1’ that appears as the ‘within-subject factor name’ and replace it with a term to denote the repeated measures aspect of your dependent variables (see repeated measures MANOVA screenshot Step 2a). In this illustration, we’ll call it ‘time’; however, if the measures were taken at regular intervals (e.g., week, month, or year), you may want to denote your within-subjects factor with a term that more descriptively defines the interval at which the measure was taken. Next, define the ‘number of levels’ of the within-subjects factor. This is essentially defining how many time points of the repeated measures are included in the analyses. In this illustration, we have measures

at five points in time, thus we enter 5 for 'number of levels' (see repeated measures MANOVA screenshot Step 2b).

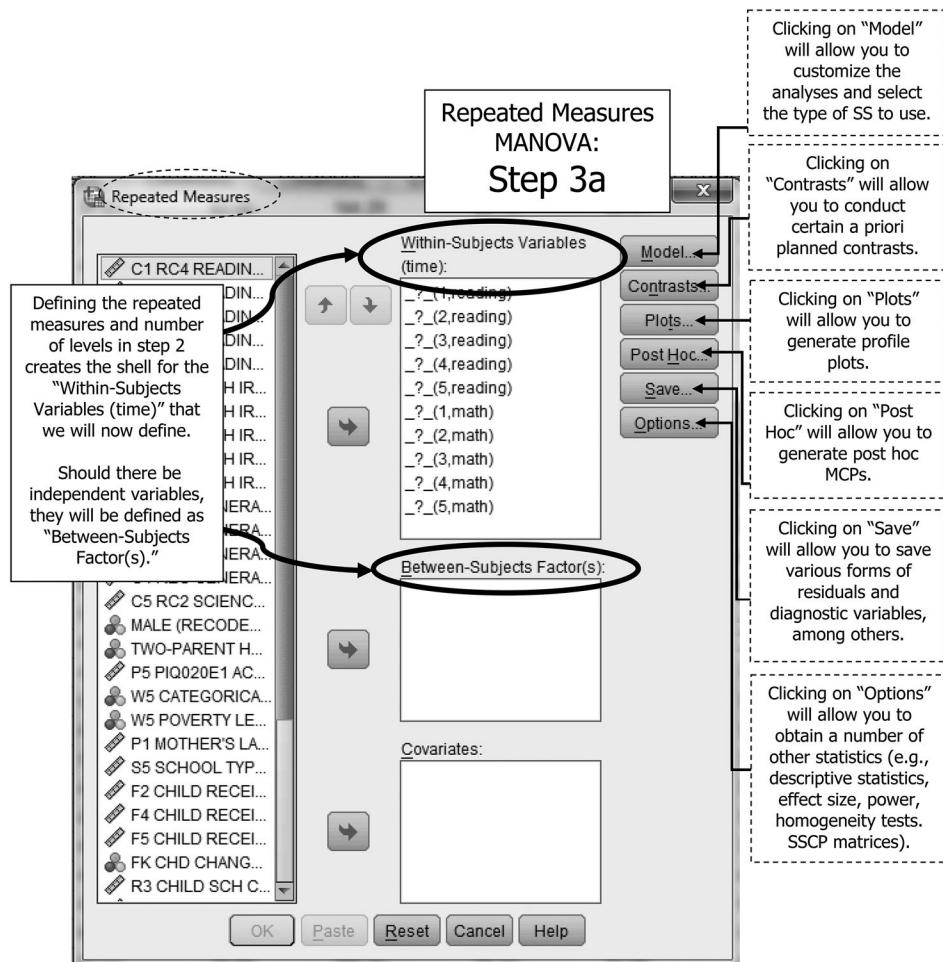


In repeated measures MANOVA screenshot Step 2c, we see the next step is to define the measure name. In this example, we have two unique dependent variables: reading and mathematics. Thus, we will type 'reading' in the 'Measure Name' box and then click 'Add' to add this to the dialog box. We repeat this step to add 'math' as a 'Measure Name' (see repeated measures MANOVA screenshot Step 2d). Follow this same step to define any additional names of dependent variables. Click on "Define" to bring up the repeated measures dialog box.

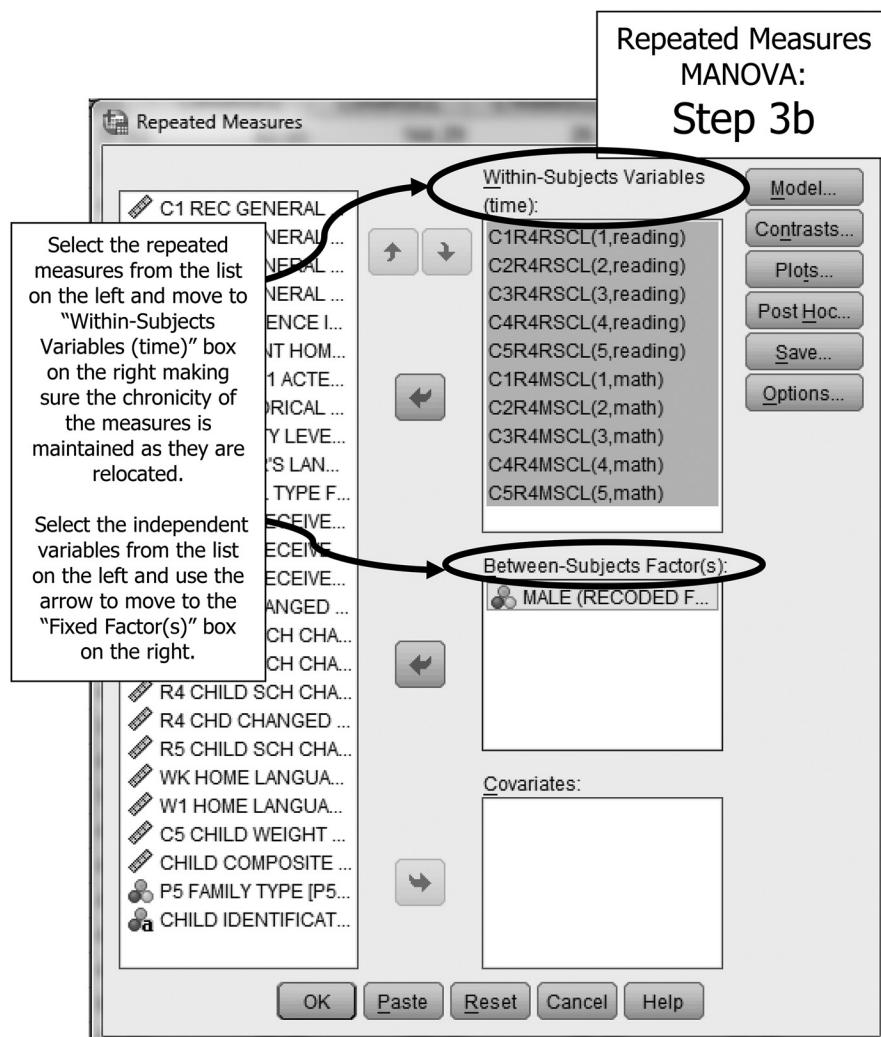


Step 3. From the main Repeated Measures dialog box, we see that the previous step (see repeated measures MANOVA screenshot Steps 2a through 2d) defined each of our dependent variables (reading and math) for each measurement occasion (in this example, five different time points) (see repeated measures MANOVA screenshot Step 3a). Now, we must define the variables that represent each of those time points. It is very important that the chronological order of the measures is replicated as the variables are

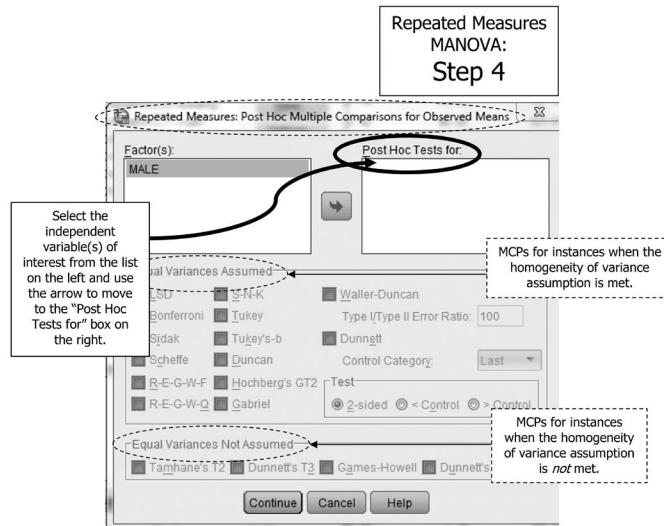
relocated from the variable list on the left to the "Within-subjects Variables (time)" box on the right. In other words, the first measurement occasion (fall kindergarten represented by C1) must be aligned with '(1,reading)' to denote the first data collection for this variable. Likewise for the remaining variables. If the variables are chronologically aligned in the variable list, they can be moved as a set by highlighting all the variables of interest. If this is not the case, each dependent variable should be moved individually.



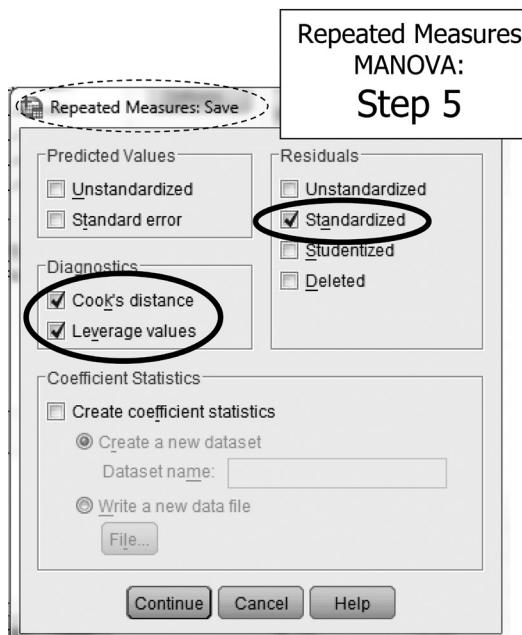
Click the repeated measures for each dependent variable (e.g., reading, math) and move them into the "Within-subjects Variables (time)" box by clicking the arrow button (see repeated measures MANOVA screenshot Step 3b). If there are also between-subjects factors (i.e., independent variables), move them to the "Between-Subjects Factor(s)" box in a similar fashion. For this example, we are including one only independent variable. Thus, click 'male' and move it into the "Between-Subjects Factor(s)" box by clicking the arrow button. Follow this same step to move any other independent variables into the "Between-Subjects Factor(s)" box.



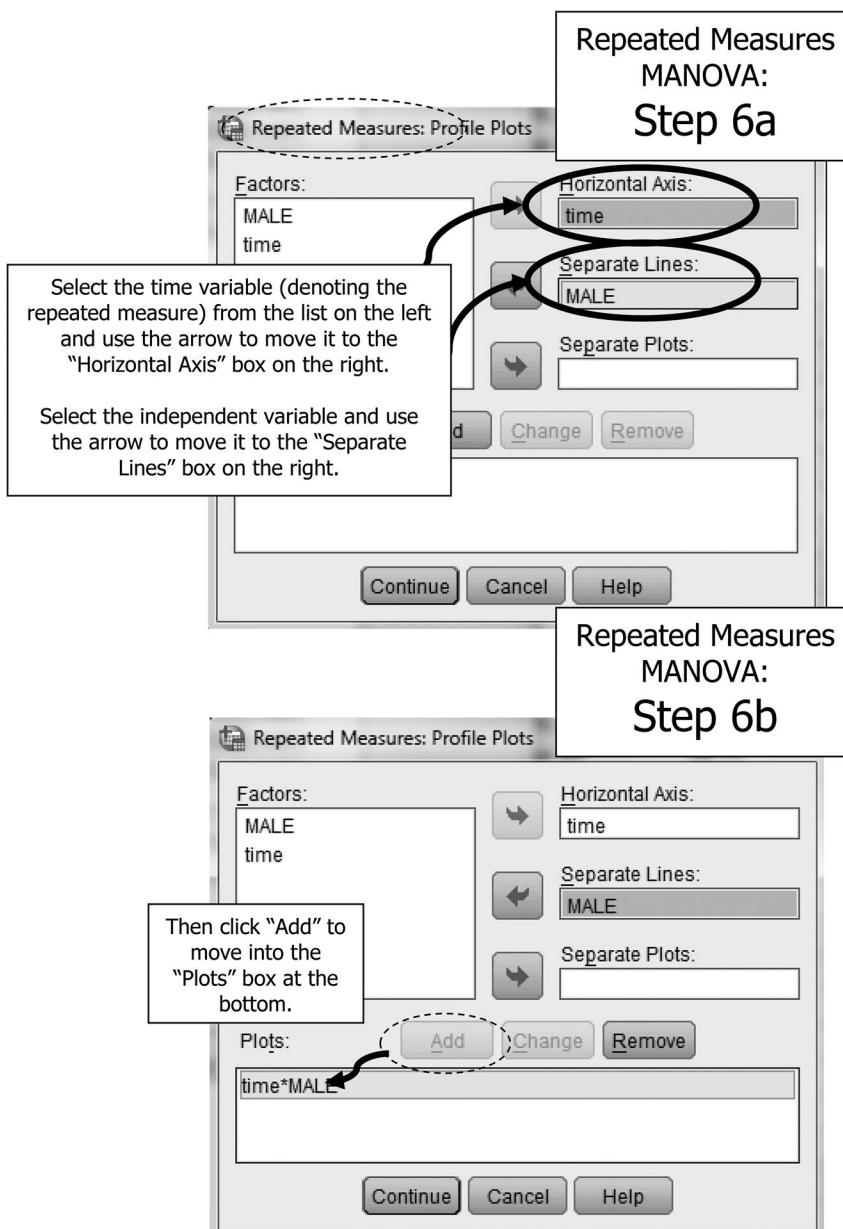
Step 4. From the “Repeated Measures” dialog box (see repeated measures MANOVA screenshot Step 3a), click on “Post Hoc” to display the Repeated Measures: Post Hoc Multiple Comparisons for Observed Means dialog box (see repeated measures MANOVA screenshot Step 4). Various post hoc multiple comparison procedures (MCPs) can be selected from this site. Because our variable is binary, even if we select post hoc tests the results will not be displayed as there are fewer than three groups in the independent variable. If you were to make a selection, you would click on the names of the independent variable in the ‘Factor(s)’ list box in the top left (e.g., ‘MALE’) and move to the ‘Post Hoc Tests for’ box in the top right by clicking on the arrow key. Check an appropriate MCP for your situation by placing a checkmark in the box next to the desired MCP. Click on “Continue” to return to the original dialog box.



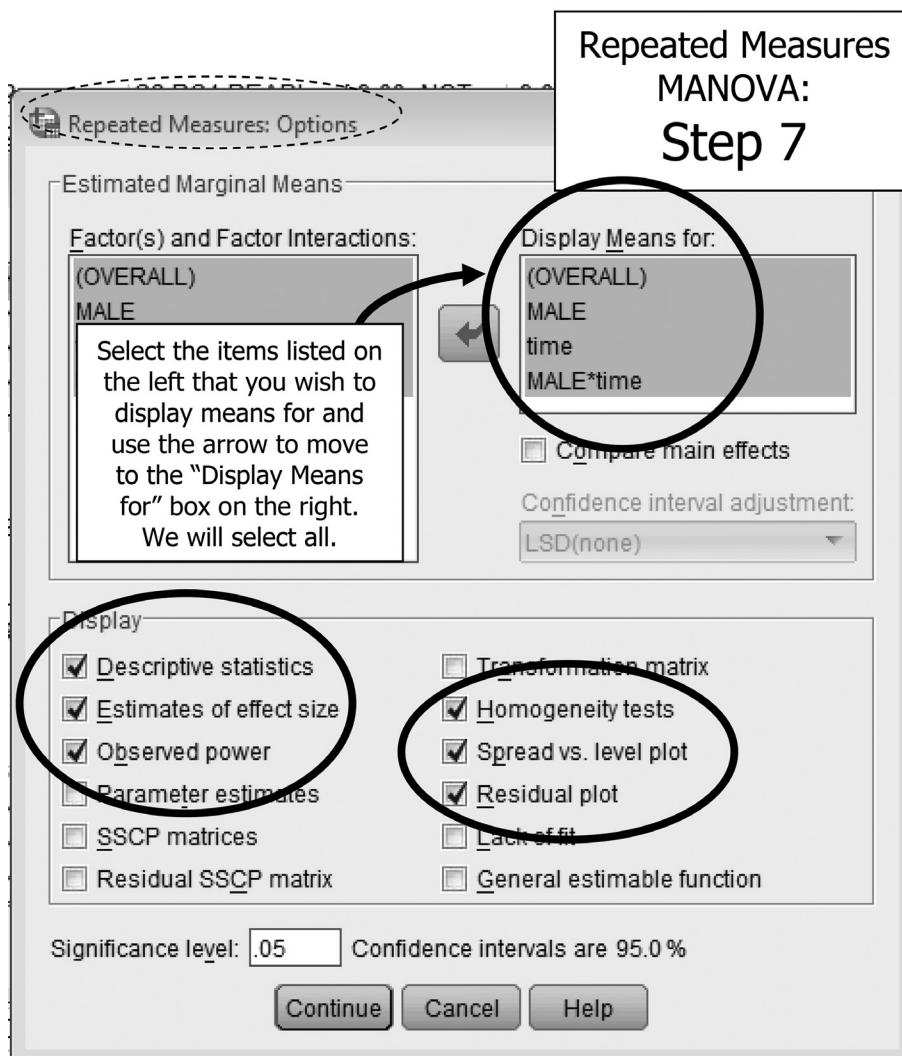
Step 5. From the “Repeated Measures” dialog box (see repeated measures MANOVA screenshot Step 3a), click on “Save” to select those elements that you want to save (see repeated measures MANOVA screenshot Step 5). In our case, we want to save the Standardized residuals, Cook’s distance, and Leverage values, which will be used later to examine the extent to which normality and independence are met. Cook’s distance provides a measure of the change in the residuals with the exclusion of a case. Small Cook’s distance values are desirable. Uncentered leverage values provide a measure of influence of each case on model fit. Click on “Continue” to return to the original dialog box.



Step 6. From the “Repeated Measures” dialog box (see repeated measures MANOVA screenshot Step 3a), click on “Plots” to display the Repeated Measures: Profile Plots dialog box (see repeated measures MANOVA screenshot Step 6a). Click the independent variable (i.e., ‘MALE’) and move it into the “Separate Lines” box by clicking the arrow button (see repeated measures MANOVA screenshot Step 6a). Then click on “Add” to move this into the “Plots” box at the bottom of the dialog box (see repeated measures MANOVA screenshot Step 6b). Click on “Continue” to return to the original dialog box.



Step 7. From the “Repeated Measures” dialog box (see repeated measures MANOVA screenshot Step 3a), click on “Options” to display the Repeated Measures: Options dialog box (see repeated measures MANOVA screenshot Step 7). Check the following selections: “Descriptive Statistics,” “Estimates of effect size,” “Observed power,” “Homogeneity tests,” “Spread vs. Level Plot,” and “Residual Plot.” We also want to display means for all the factors and interactions by clicking the variable names in the box on the top left and moving to the “Display Means for” box. Click on “Continue” to return to the original dialog box. From the “Repeated Measures” dialog box (see repeated measures MANOVA screenshot Step 3a), click on “OK” to return to generate the output.



Interpreting the output. Annotated results are presented in Tables 6.4.

TABLE 6.4

SPSS Results for the Repeated Measures MANOVA

Within-Subjects Factors		
Measure	time	Dependent Variable
reading	1	C1R4RSCL
	2	C2R4RSCL
	3	C3R4RSCL
	4	C4R4RSCL
	5	C5R4RSCL
	1	C1R4MSCL
	2	C2R4MSCL
	3	C3R4MSCL
	4	C4R4MSCL
	5	C5R4MSCL
Between-Subjects Factors		
		Value Label
MALE (RECODED FROM 'GENDER')		Value Label
.00		FEMALE
1.00		MALE
		N
		681
		663

Box's Test of Equality of Covariance Matrices ^a		
Box's M	209.764	
F	3.784	
df1	55	
df2	5807401.089	
Sig.	.000	

Tests the null hypothesis that the observed covariance matrices of the dependent variables are equal across groups.

a. Design: Intercept + MALE

Within Subjects Design: time

The table labeled "Within-Subjects Factors" lists each unique dependent variable and the specific repeated measures associated with each (recall that the repeated measures are the 'between subjects factors').

The table labeled "Between-Subjects Factors" provides sample sizes for each of the categories of the independent variables (recall that the independent variables are the 'between subjects factors').

The *F* test (and associated *p* value) for Box's Test for Equality of Covariance Matrices is reviewed to determine if equal covariances can be assumed. In this case, we do not meet the assumption (as *p* is less than α). The covariance matrices of the dependent variables cannot be assumed to be relatively similar across all levels of the independent variable. In this example, since the sample sizes of boys to girls are very similar (boys outnumbering girls by only about 3%), we can disregard Box's test and will use Pillai's Trace as the multivariate test.

■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA

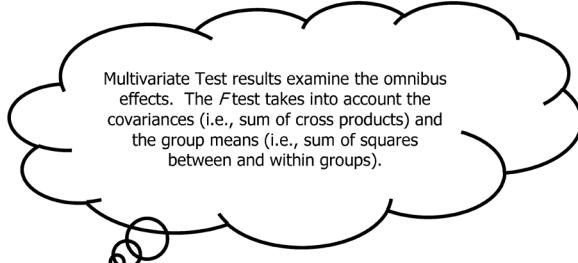
The table labeled "Descriptive Statistics" provides basic descriptive statistics (means, standard deviations, and sample sizes) for each cell of the design.

Descriptive Statistics

	MALE	Mean	Std. Deviation	N
C1 RC4 READING IRT SCALE SCORE	FEMALE	36.4623	9.09909	681
	MALE	35.7253	10.28692	663
	Total	36.0987	9.70662	1344
C2 RC4 READING IRT SCALE SCORE	FEMALE	49.2737	13.33427	681
	MALE	47.5793	13.62304	663
	Total	48.4378	13.49910	1344
C3 RC4 READING IRT SCALE SCORE	FEMALE	55.6704	17.04792	681
	MALE	53.6108	17.04154	663
	Total	54.6544	17.06953	1344
C4 RC4 READING IRT SCALE SCORE	FEMALE	81.7899	22.47575	681
	MALE	79.1078	23.31056	663
	Total	80.4668	22.92212	1344
C5 RC4 READING IRT SCALE SCORE	FEMALE	133.1565	26.04651	681
	MALE	130.8594	26.62739	663
	Total	132.0234	26.34991	1344
C1 RC4 MATH IRT SCALE SCORE	FEMALE	27.1188	8.16079	681
	MALE	28.0011	9.90379	663
	Total	27.5540	9.06996	1344
C2 RC4 MATH IRT SCALE SCORE	FEMALE	38.2300	10.50029	681
	MALE	39.4090	12.43731	663
	Total	38.8116	11.50751	1344
C3 RC4 MATH IRT SCALE SCORE	FEMALE	44.7730	12.73346	681
	MALE	47.0357	14.71841	663
	Total	45.8892	13.78988	1344
C4 RC4 MATH IRT SCALE SCORE	FEMALE	63.4322	16.13855	681
	MALE	66.9349	18.70499	663
	Total	65.1601	17.53303	1344
C5 RC4 MATH IRT SCALE SCORE	FEMALE	100.4706	23.12008	681
	MALE	106.8578	24.05597	663
	Total	103.6214	23.79303	1344

■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA



The test of the intercept is of little interest. It tests if the grand mean differs from zero.

Multivariate Tests ^a									
Effect		Value	F	Hypo. df	Error df	Sig.	Partial Eta Squared	Noncent. Parameter	Obs. Power ^c
Between	Pillai's Trace	.956	14508.948 ^b	2.000	1341.000	.000	.956	29017.896	1.000
	Wilks' Lambda	.044	14508.948 ^b	2.000	1341.000	.000	.956	29017.896	1.000
	Hotelling's Trace	21.639	14508.948 ^b	2.000	1341.000	.000	.956	29017.896	1.000
	Roy's Largest Root	21.639	14508.948 ^b	2.000	1341.000	.000	.956	29017.896	1.000
Subjects	Pillai's Trace	.056	40.147 ^b	2.000	1341.000	.000	.056	80.293	1.000
	Wilks' Lambda	.944	40.147 ^b	2.000	1341.000	.000	.056	80.293	1.000
	Hotelling's Trace	.060	40.147 ^b	2.000	1341.000	.000	.056	80.293	1.000
	Roy's Largest Root	.060	40.147 ^b	2.000	1341.000	.000	.056	80.293	1.000
Time	Pillai's Trace	.956	3657.966 ^b	8.000	1335.000	.000	.956	29263.725	1.000
	Wilks' Lambda	.044	3657.966 ^b	8.000	1335.000	.000	.956	29263.725	1.000
	Hotelling's Trace	21.920	3657.966 ^b	8.000	1335.000	.000	.956	29263.725	1.000
	Roy's Largest Root	21.920	3657.966 ^b	8.000	1335.000	.000	.956	29263.725	1.000
Within	Pillai's Trace	.051	8.990 ^b	8.000	1335.000	.000	.051	71.918	1.000
	Wilks' Lambda	.949	8.990 ^b	8.000	1335.000	.000	.051	71.918	1.000
	Hotelling's Trace	.054	8.990 ^b	8.000	1335.000	.000	.051	71.918	1.000
	Roy's Largest Root	.054	8.990 ^b	8.000	1335.000	.000	.051	71.918	1.000

a. Design: Intercept + MALE
Within Subjects Design: time
b. Exact statistic
c. Computed using alpha = .05

Pillai's Trace represents the sum of the proportion of explained variance on all the discriminant functions similar to the ratio of SS between to SS total.

Wilks' Lambda represents the ratio of error variance to total variance for each variate.

Hotellings's Trace compares directly to the *F* ratio in ANOVA, and is the sum of the ratio of between variance to error variance for each variate (i.e., the sum of eigenvalues for each variate).

Roy's Largest Root represents the maximum possible between-group differences and is conceptually the same as the *F* ratio in ANOVA, representing the proportion of explained to unexplained variance for the first discriminant function. In many cases, this is the most powerful.

In terms of power, all four have similar power even with small sample sizes. All four test statistics are relatively robust to violations of multivariate normality (with some exceptions: Roy's root is less robust if there is a platykurtic distribution AND Roy's root is less robust when the assumption of homogeneity of covariances is *not* met).

When there is a balanced design, Pillai's is most robust to assumption violations. With an unbalanced design, Pillai's is less robust when homogeneity of variance-covariances is not met.

Regardless of which test we select, there are statistically significant multivariate main effects of male and time. Of primary consideration is the statistically significant interaction of male by time as this qualifies any main effects and suggests that the difference between boys and girls on the linear combination of reading and math outcomes is different over time. In other words, there are simultaneous mean differences in the outcomes when comparing girls to boys.

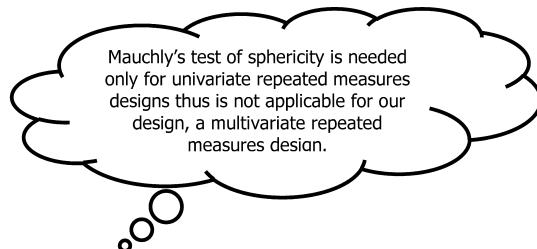
$$\text{partial } \eta^2 = \frac{df_1 F}{df_1 F + df_2}$$

$$= \frac{8(8.990)}{8(8.990) + 1335} = \frac{71.92}{1406.92} = .051$$

Which to report? Wilks' Lambda if assumptions are met. Pillai's trace if assumptions are violated.

■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA



Mauchly's Test of Sphericity ^a								
Within Subjects Effect	Measure	Mauchly's W	Approx. Chi-Square	df	Sig.	Epsilon ^b		
						Greenhouse-Geisser	Huynh-Feldt	Lower-bound
time	reading	.116	2881.885	9	.000	.544	.546	.250
	math	.216	2054.879	9	.000	.555	.556	.250

Tests the null hypothesis that the error covariance matrix of the orthonormalized transformed dependent variables is proportional to an identity matrix.

a. Design: Intercept + MALE

Within Subjects Design: time

b. May be used to adjust the degrees of freedom for the averaged tests of significance. Corrected tests are displayed in the Tests of Within-Subjects Effects table.

Tests of Within-Subjects Effects

Multivariate ^{a,b}								
Within Subjects Effect	Value	F	Hypothesis df	Error df	Sig.	Partial Eta Squared	Noncent. Parameter	Observed Power ^c
time	Pillai's Trace	.956	1229.357	8.000	10736.000	.000	.478	9834.856
	Wilks' Lambda	.065	3936.088 ^d	8.000	10734.000	.000	.746	31488.701
	Hotelling's Trace	14.151	9491.523	8.000	10732.000	.000	.876	75932.182
	Roy's Largest Root	14.128	18959.516 ^d	4.000	5368.000	.000	.934	75838.065
MALE	Pillai's Trace	.021	13.948	8.000	10736.000	.000	.010	111.582
	Wilks' Lambda	.979	14.012 ^c	8.000	10734.000	.000	.010	112.098
	Hotelling's Trace	.021	14.077	8.000	10732.000	.000	.010	112.614
	Roy's Largest Root	.021	27.575 ^d	4.000	5368.000	.000	.020	110.301

a. Design: Intercept + MALE

Within Subjects Design: time

b. Tests are based on averaged variables.

c. Exact statistic

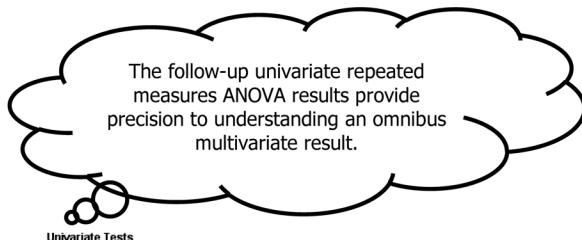
d. The statistic is an upper bound on F that yields a lower bound on the significance level.

e. Computed using alpha = .05

You'll notice some similarities between this multivariate within-subjects effects table and the table that presented previously. The difference is that *all* the calculations in the previous table are 'exact statistics' whereas some of the values in this table are based on 'averaged' values. For this reason, disregard the values in this table and use those from the previous multivariate table.

■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA



Univariate Tests									
Source	Measure	Type III Sum of Squares	df	Mean Square	F	Sig.	Partial Eta Squared	Noncent. Parameter	Observed Power ^a
time	reading	Sphericity Assumed	7799987.234	4	1949996.808	13655.299	.000	.911	54621.197
	reading	Greenhouse-Geisser	7799987.234	2.178	3581704.479	13655.299	.000	.911	29737.562
	reading	Huynh-Feldt	7799987.234	2.183	3573012.524	13655.299	.000	.911	29809.904
	reading	Lower-bound	7799987.234	1.000	7799987.234	13655.299	.000	.911	13655.299
	math	Sphericity Assumed	4786316.696	4	1196579.174	13457.649	.000	.909	53830.597
time * MALE	reading	Greenhouse-Geisser	4786316.696	2.219	2157309.306	13457.649	.000	.909	29857.828
	reading	Huynh-Feldt	4786316.696	2.224	2151981.155	13457.649	.000	.909	29931.754
	reading	Lower-bound	4786316.696	1.000	4786316.696	13457.649	.000	.909	13457.649
	math	Sphericity Assumed	735.522	4	183.881	1.288	.272	.001	5.151
	math	Greenhouse-Geisser	735.522	2.178	337.747	1.288	.277	.001	2.804
Error(time)	reading	Huynh-Feldt	735.522	2.183	336.928	1.288	.277	.001	.281
	reading	Lower-bound	735.522	1.000	735.522	1.288	.257	.001	.1.288
	math	Sphericity Assumed	6700.844	4	1675.211	18.841	.000	.014	75.363
	math	Greenhouse-Geisser	6700.844	2.219	3020.233	18.841	.000	.014	41.801
	math	Huynh-Feldt	6700.844	2.224	3012.774	18.841	.000	.014	41.904
	math	Lower-bound	6700.844	1.000	6700.844	18.841	.000	.014	18.841
	reading	Sphericity Assumed	766558.293	5368	142.801				
	reading	Greenhouse-Geisser	766558.293	2922.514	262.294				
	reading	Huynh-Feldt	766558.293	2929.624	261.658				
	reading	Lower-bound	766558.293	1342.000	571.206				
	math	Sphericity Assumed	477292.647	5368	88.914				
	math	Greenhouse-Geisser	477292.647	2977.430	160.304				
	math	Huynh-Feldt	477292.647	2984.802	159.908				
	math	Lower-bound	477292.647	1342.000	355.658				

a. Computed using alpha = .05

The statistically significant interaction supersedes the main effects of gender and time, thus follow-up examination focuses only on the interaction.

For the interaction of gender by time, the results are statistically significant only for math. In other words, there is differential growth for boys and girls over time but only for math, not reading.

■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA

Tests of Within-Subjects Contrasts										
Source	Measure	time	Type III SS	df	Mean Square	F	Sig.	Partial Eta Squared	Noncent. Parameter	Obs. Power ^a
time	reading	Linear	6733463.403	1	6733463.403	21820.830	.000	.942	21820.830	1.000
		Quadratic	922702.322	1	922702.322	6411.537	.000	.827	6411.537	1.000
		Cubic	136479.930	1	136479.930	1740.875	.000	.565	1740.875	1.000
		Order 4	7341.579	1	7341.579	182.100	.000	.119	182.100	1.000
	math	Linear	4284999.880	1	4284999.880	22273.759	.000	.943	22273.759	1.000
		Quadratic	426205.225	1	426205.225	5997.538	.000	.817	5997.538	1.000
		Cubic	73429.112	1	73429.112	1369.760	.000	.505	1369.760	1.000
		Order 4	1682.479	1	1682.479	43.578	.000	.031	43.578	1.000
	reading	Linear	566.871	1	566.871	1.837	.176	.001	1.837	.273
		Quadratic	141.400	1	141.400	.983	.322	.001	.983	.168
		Cubic	5.790	1	5.790	.074	.786	.000	.074	.058
		Order 4	21.461	1	21.461	.532	.466	.000	.532	.113
MALE	math	Linear	5972.377	1	5972.377	31.045	.000	.023	31.045	1.000
		Quadratic	682.201	1	682.201	9.600	.002	.007	9.600	.872
		Cubic	24.717	1	24.717	.461	.497	.000	.461	.104
		Order 4	21.549	1	21.549	.558	.455	.000	.558	.116
	reading	Linear	414113.845	1342	308.580					
		Quadratic	193131.000	1342	143.913					
		Cubic	105209.201	1342	78.397					
		Order 4	54104.247	1342	40.316					
Error(time)	math	Linear	258172.399	1342	192.379					
		Quadratic	95367.039	1342	71.063					
		Cubic	71940.960	1342	53.607					
		Order 4	51812.249	1342	38.608					

a. Computed using alpha = .05

The precision and pattern of the interaction between gender and time can be assessed from the within-subjects contrasts. For the interaction of gender and time, there is a statistically significant linear and quadratic trend only for math outcomes, not reading. This means that boys and girls differ in their tendency to show linear as well as quadratic growth over time. Quadratic growth indicates a 'bend' in the data (e.g., U or upside down U shape).

TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA

The results in this table are based on the 'average.' In other words, if the within-subjects variables were averaged across the multiple time points, these are any between-subjects effects that would result. The results from this table can be disregarded.

Tests of Between-Subjects Effects

Transformed Variable: Average

Source	Measure	Type III SS	df	Mean Square	F	Sig.	Partial Eta Squared	Noncent. Parameter	Obs. Power ^a
MALE	reading	33227135.849	1	33227135.849	27400.603	.000	.953	27400.603	1.000
	math	21240781.523	1	21240781.523	23378.773	.000	.946	23378.773	1.000
	reading	6025.723	1	6025.723	4.969	.026	.004	4.969	.606
	math	13574.390	1	13574.390	14.941	.000	.011	14.941	.971
	reading	1627366.217	1342	1212.642					
	math	1219273.928	1342	908.550					

a. Computed using alpha = .05

This indicates that if the reading and math outcomes were averaged across all 5 time points, there would be a statistically significant main effect for males.

Estimated Marginal Means**1. Grand Mean**

Measure	Mean	Std. Error	95% Confidence Interval	
			Lower Bound	Upper Bound
reading	70.324	.425	69.490	71.157
math	56.226	.368	55.505	56.948

The 'Grand Mean' represents the overall mean, regardless of group membership and averaged across all measurement occasions, for each of the unique dependent variables. The 95% CI represents the CI of the grand mean.

2. MALE

Measure	MALE (RECODED FROM 'GENDER')	Mean	Std. Error	95% Confidence Interval	
				Lower Bound	Upper Bound
reading	FEMALE	71.271	.597	70.100	72.441
	MALE	69.377	.605	68.190	70.563
	FEMALE	54.805	.517	53.792	55.818
	MALE	57.648	.524	56.621	58.675

The table labeled 'MALE' provides descriptive statistics for each of the categories of the independent variable by each dependent variable averaged across measurement occasions. In addition to means, the SE and 95% CI of the means are reported.

TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA

Measure	time	Mean	Std. Error	3. time ←	
				95% Confidence Interval	
				Lower Bound	Upper Bound
reading	1	36.094	.265	35.575	36.613
	2	48.426	.368	47.705	49.148
	3	54.641	.465	53.728	55.553
	4	80.449	.624	79.224	81.674
	5	132.008	.718	130.599	133.417
math	1	27.560	.247	27.075	28.045
	2	38.819	.314	38.204	39.435
	3	45.904	.375	45.169	46.640
	4	65.184	.476	64.250	66.117
	5	103.664	.643	102.402	104.926

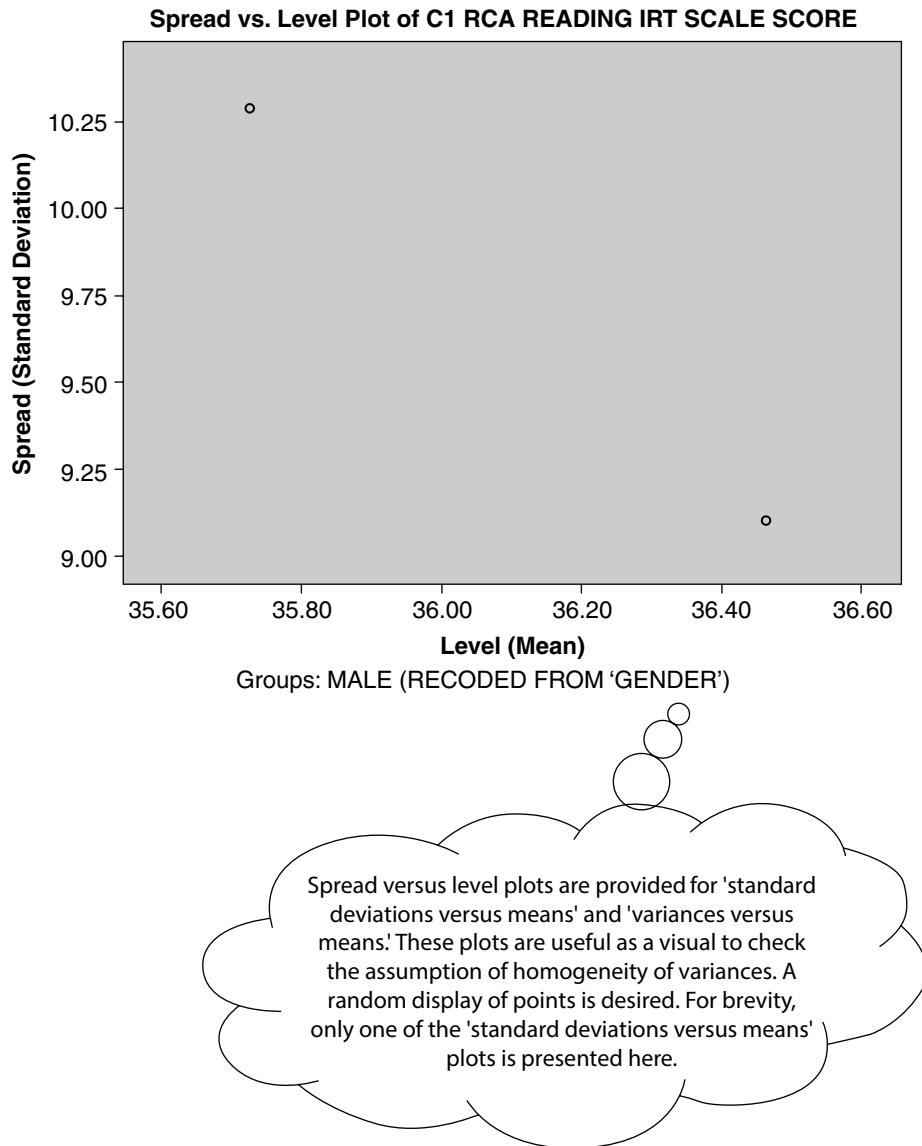
Measure	MALE	time	Mean	Std. Error	4. MALE * time ←	
					95% Confidence Interval	
					Lower Bound	Upper Bound
reading	FEMALE	1	36.462	.372	35.733	37.192
		2	49.274	.516	48.261	50.287
		3	55.670	.653	54.389	56.952
		4	81.790	.877	80.069	83.511
		5	133.157	1.009	131.177	135.136
	MALE	1	35.725	.377	34.986	36.465
		2	47.579	.523	46.552	48.606
		3	53.611	.662	52.312	54.909
		4	79.108	.889	77.364	80.852
		5	130.859	1.023	128.853	132.866
math	FEMALE	1	27.119	.347	26.438	27.800
		2	38.230	.441	37.366	39.094
		3	44.773	.527	43.739	45.807
		4	63.432	.669	62.120	64.744
		5	100.471	.904	98.698	102.244
	MALE	1	28.001	.352	27.311	28.692
		2	39.409	.446	38.533	40.285
		3	47.036	.534	45.988	48.083
		4	66.935	.678	65.605	68.264
		5	106.858	.916	105.061	108.655

'Time' represents the means by measurement occasion, regardless of group membership, across all measurement occasions, for each of the unique dependent variables. The 95% CI represents the CI of by time point.

The table labeled 'MALE*time' provides descriptive statistics for each dependent variable for each of the categories of the independent variable by measurement occasion. In addition to means, the *SE* and 95% CI of the means are reported.

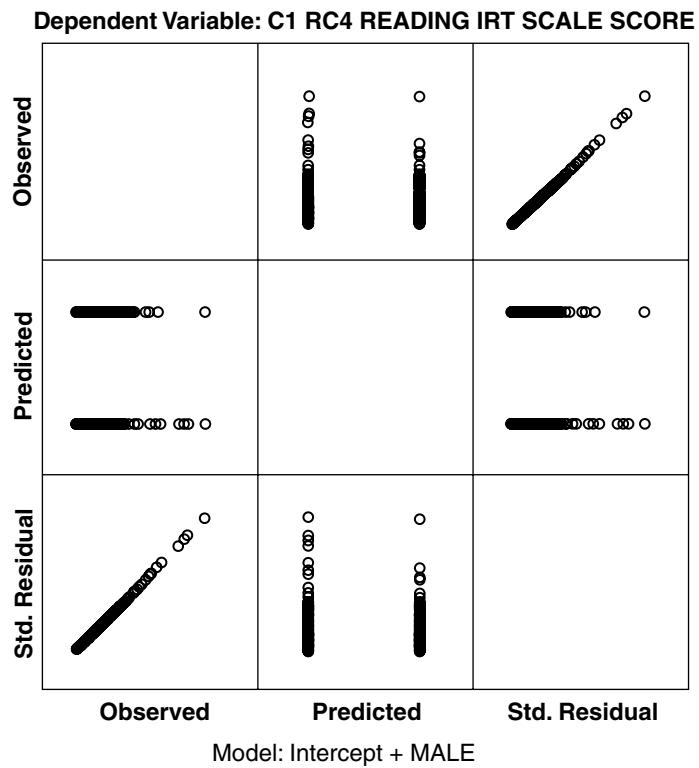
■ TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA



■ TABLE 6.4 (continued)

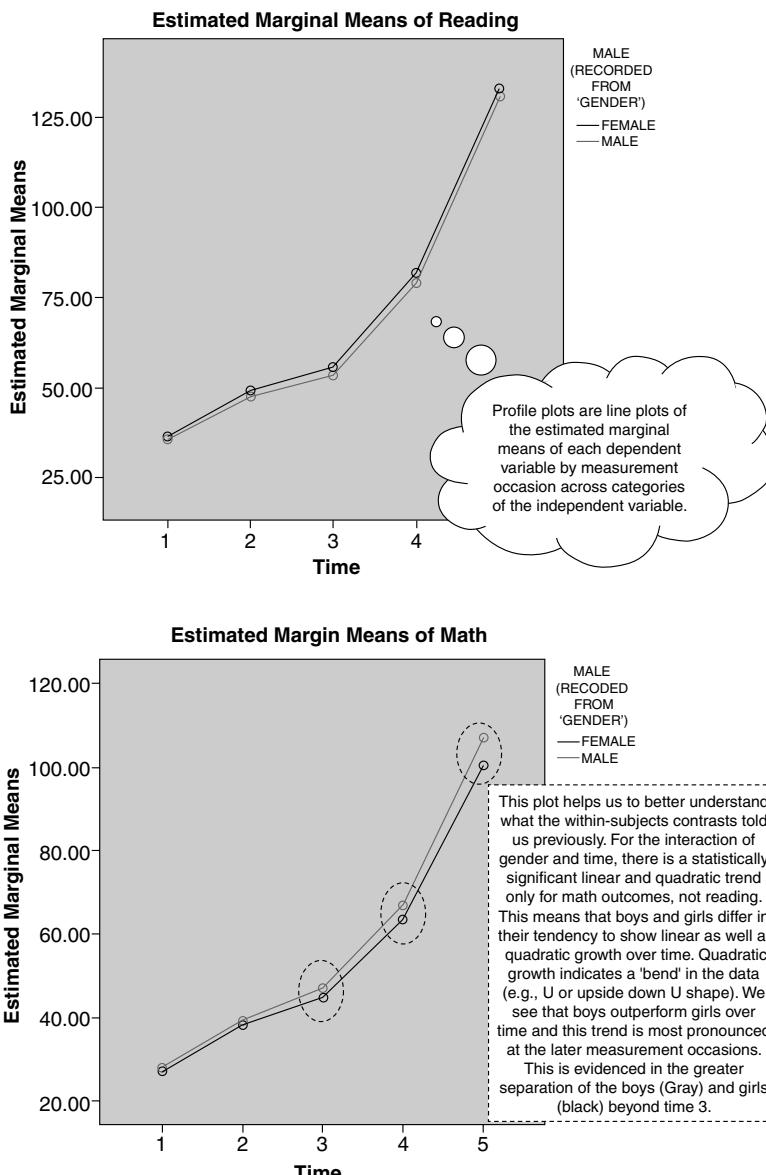
SPSS Results for the Repeated Measures MANOVA



The residual plots are plots of the observed, predicted, and standardized residuals for each measurement occasion for each unique dependent variable. For brevity, only one of the plots is presented here. A linear display of points is desirable for the observed by predicted residuals. A random display is desired for plots involving the standardized residuals. The exception to this is for binary data. In the case of binary data (as we have), a random pattern cannot be obtained.

TABLE 6.4 (continued)

SPSS Results for the Repeated Measures MANOVA



6.4 DATA SCREENING

6.4.1 Data Screening for One-Way and k -Way MANOVA Models

The assumptions for MANOVA include (a) independence, (b) multivariate normality of the dependent variables, (c) linearity (a component of multivariate normality), and (d) homogeneity of variance-covariance matrices for the dependent variables.

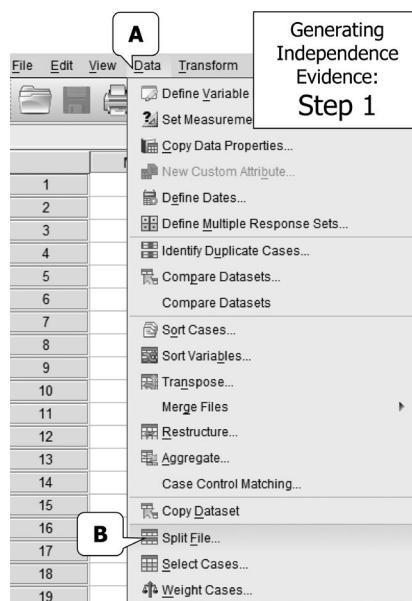
6.4.1.1 Independence

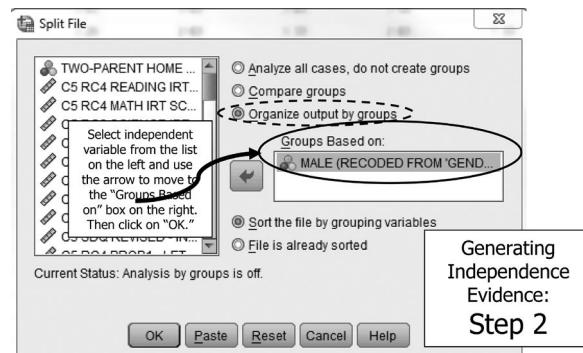
If subjects have been randomly assigned to conditions (or to the different combinations of the levels of the independent variables in the MANOVA), the assumption of independence has been met. In many cases, we use independent variables that do not allow random assignment, such as preexisting or self-selected characteristics, such as those in this sample. We can plot residuals against levels of our independent variables in a scatterplot to get an idea of whether or not there are patterns in the data and thereby provide an indication of whether we have met this assumption. Given we have multiple independent variables in the factorial MANOVA, we will split the scatterplot by levels of one independent variable ('MALE') and then generate a bivariate scatterplot for 'TWOPARENT' by residual. Remember that the residual was added to the dataset by saving it when we generated the factorial MANOVA model.

Please note that some researchers do not believe that the assumption of independence can be tested. If there is not random assignment to groups, then these researchers believe this assumption has been violated—period. The plot that we generate will give us a general idea of patterns, however, in situations where random assignment was not performed or not possible.

Splitting the File

The first step is to split our file by the levels of one of our independent variables (e.g., 'MALE'). To do that, go to 'Data' in the top pull-down menu and then select 'Split File' (see generating independence evidence screenshot Step 1). Next, move the variable on which you want to split the file into the 'groups based on' box, and click the radio button for 'organize output by groups' (see generating independence evidence screenshot Step 2).



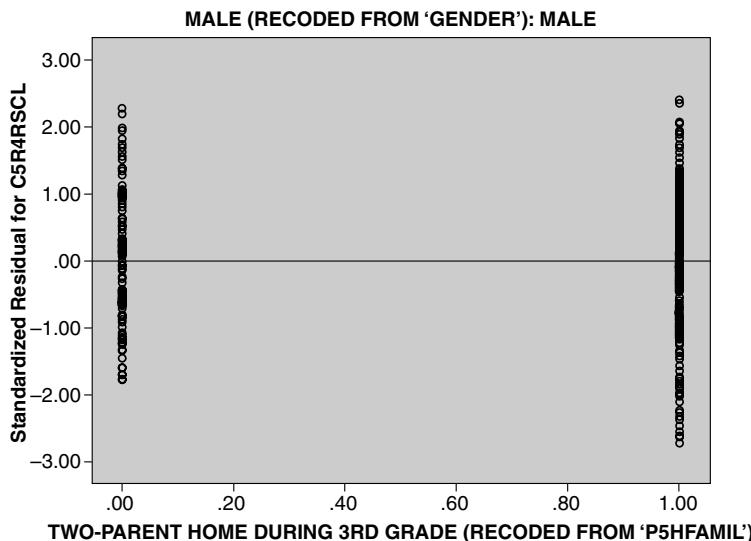


Generating the Scatterplot

Instructions for generating a scatterplot (if you need a refresher) are presented in the data screening chapter.

Interpreting Independence Evidence

In examining the scatterplots for evidence of independence, the points should fall relatively randomly above and below a horizontal line at zero and within the band of +2.0 to -2.0. In this example, our scatterplot for each level of two-parent home generally suggests evidence of independence with a relatively random display of residuals above and below the horizontal line at zero for both males and females (i.e., each category of the independent variable that was used to split the file) and generally



within a band of an absolute value of ± 2.0 . We repeat this process for each standardized residual. Additional graphs are not presented here, but all are similar to the one provided here, suggesting a relatively random display of residuals above and below zero. Thus, since we have not met the assumption of independence through random assignment of cases to groups, this gives us some assurance that independence is a reasonable assumption.

6.4.1.2 Multivariate Normality of the Dependent Variables

Our examination of multivariate normality will commence with review of univariate normality. This will be followed by examination of multivariate normality.

Univariate Normality Evidence

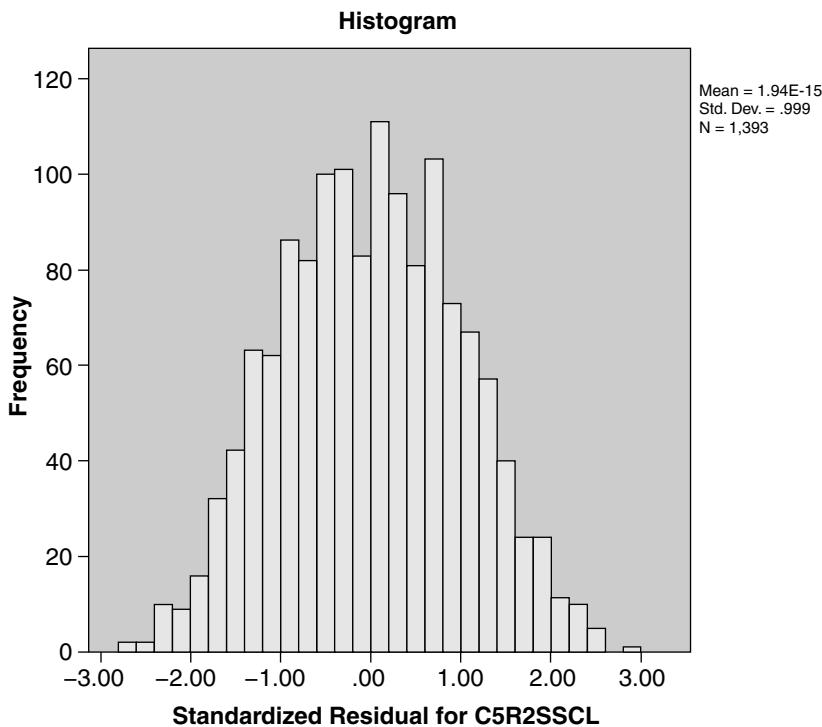
Testing for multivariate normality is difficult using statistical software. As a work around, the next best solution is to test for univariate normality of residuals, as univariate normality is a necessary (but insufficient) condition for multivariate normality. We will request skewness and kurtosis statistics, histograms, formal tests of normality, Q-Q plots, and boxplots to explore this assumption. The steps for generating these were presented in the data screening chapter, should you need a refresher.

Interpreting Univariate Normality Evidence

We have quite a bit of experience in interpreting univariate normality evidence, including skewness and kurtosis, histograms, and boxplots.

	Statistics		
	Standardized Residual for C5R4RSCL	Standardized Residual for C5R4MSCL	Standardized Residual for C5R2SSCL
N Valid	1393	1393	1393
Missing	0	0	0
Skewness	-.230	-.209	.037
Std. Error of Skewness	.066	.066	.066
Kurtosis	-.417	-.556	-.485
Std. Error of Kurtosis	.131	.131	.131

The skewness statistics of the residuals are within the range of $-.230$ to $.037$ and kurtosis are within the range of $-.556$ to $-.417$, all within the range of an absolute value of 2.0 and 7.0 , respectively, suggesting some evidence of normality. The histograms of residuals are also relatively normal and there is little to suggest that normality may not be an unreasonable assumption. For brevity, only one graph is presented here.



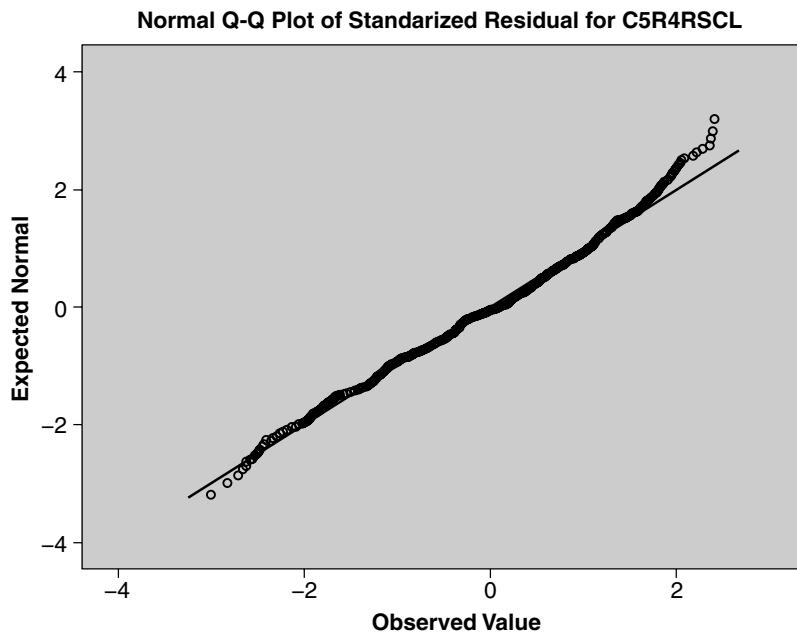
The output for the Shapiro-Wilk test is presented below and suggests that our sample distributions for residuals are statistically significantly different from what would be expected from a normal distribution ($p < .001$ for all).

Tests of Normality

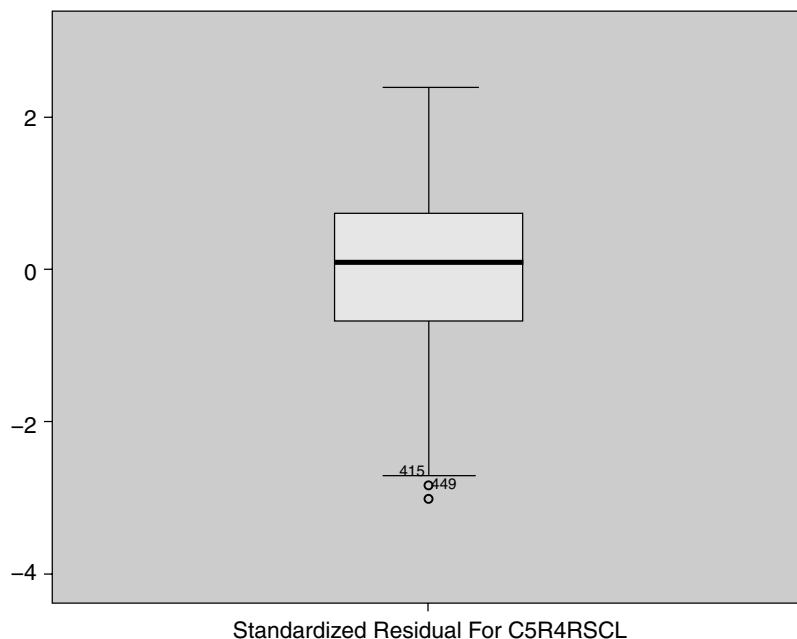
	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Standardized Residual for C5R4RSCL	.042	1393	.000	.991	1393	.000
Standardized Residual for C5R4MSCL	.036	1393	.000	.990	1393	.000
Standardized Residual for C5R2SSCL	.025	1393	.037	.996	1393	.001

a. Lilliefors Significance Correction

The Q-Q plot of residuals shown on the next page suggests relative normality with the only nonnormality suggested in the tails of the distribution. For brevity, only one graph is presented here.

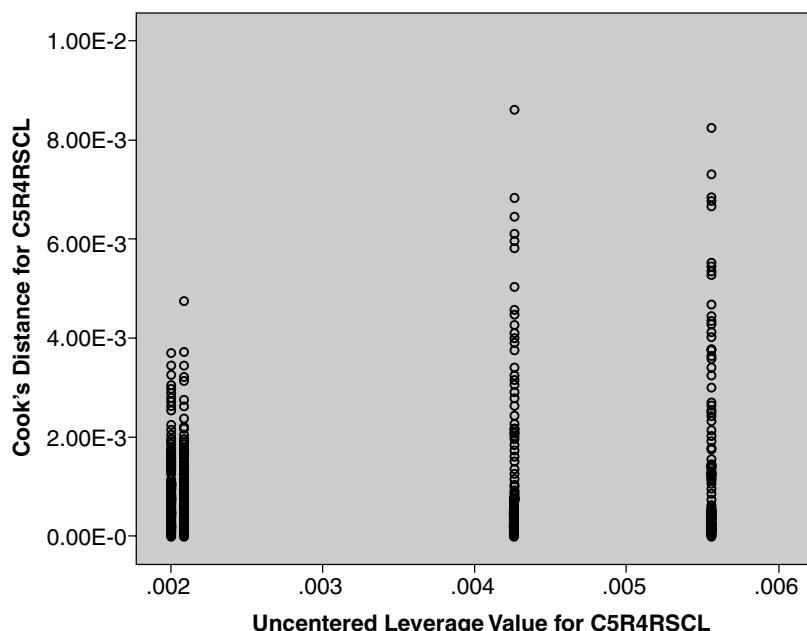


Examination of the boxplot of standardized residuals for reading, presented here, suggests a relatively normal distributional shape of residuals with two potential outliers. The boxplots of math and science residuals (not presented for brevity) are also relatively normal in shape with no outliers.



With the exception of the Shapiro-Wilk test, the other forms of evidence examined, skewness and kurtosis statistics, the Q-Q plot, and the boxplot, suggest normality is a reasonable assumption. Even though there is some nonnormality suggested by the Shapiro-Wilk test, we can be reasonably assured that we have met the assumption of univariate normality of the dependent variables.

There are other diagnostic tools that can be used to screen our data for influential points that may impact univariate or multivariate normality and that can be generated from the values we saved when conducting MANOVA. For example, a plot of Cook's distance against uncentered leverage values provides a way to identify influential cases (i.e., cases with leverage of .50 or above and Cook's distance of 1.0 or greater). Here there are no cases that suggest undue influence.



Cook's distance provides an overall measure for the influence of individual cases. Values greater than one suggest that the case may be problematic in terms of undue influence on the model. Examining the residual statistics in our output (see following table), we see that the maximum value for Cook's distance is .01, well under the point at which we should be concerned.

Statistics

		Cook's Distance for C5R4RSCL	Cook's Distance for C5R4MSCL	Cook's Distance for C5R2SSCL
N	Valid	1393	1393	1393
	Missing	0	0	0
Minimum		.00	.00	.00
Maximum		.01	.01	.01

In summary, data screening of *univariate normality* and diagnostics suggest that normality is reasonable.

Multivariate Normality Evidence

Although univariate normality is a necessary condition of multivariate normality, it is not a sufficient condition. SPSS does not have built-in capability to examine *multivariate normality*. However, we will enlist the use of a macro that will do just that. The macro is based on the work of DeCarlo (1997) and is extremely simple to use.

First, copy the macro language (the original language is accessible from <http://www.columbia.edu/~ld208/normtest.sps> and printed in the 1997 article as well; there is a bug in that macro that is fixed in the language accessible from this site: https://dl.dropboxusercontent.com/u/3385251/Mult_Macro_fixAndy.sps) and paste it into a SPSS syntax file. To open a new syntax file, in SPSS go to File in the top pull-down menu, then New then Syntax. Your syntax file with the macro language should look something like this:

```

normtest_rev.sps - IBM SPSS Statistics Syntax Editor
File Edit View Data Transform Analyze Direct Marketing Graphs Utilities Add-ons Run Tools Window Help
Active: temp

1  *Some edits to this ...
2  *For original and th...
3  *http://spssx-discus...
4  ****
5  *
6  * Psychologi...
7  *
8  * \* normtestvars=x...
9  * (variable names ca...
10 * (1983). Algorithms...
11 * Fisher's g statist...
12 * Mardia's p-value f...
13 define
14 tendiffdefine.

Line   Command           Information
19  execute  Unrecognized text appears on the EXECUTE command. This command allows n
                                subcommands.

IBM SPSS Statistics Processor is ready | Unicode:ON In 225 Col 11 | NUM

```

Save this syntax file by going to File then Save as in the top pull-down menu in SPSS.

Second, open the SPSS data file that includes the variables you want to examine. In this case, it's the ECLS-K data. The dataset you will be working with needs to be open when running the syntax.

Third, open another new syntax file. To open a new syntax file, in SPSS go to File in the top pull-down menu, then New then Syntax. Type the following two-line command:

```

INCLUDE file='C:\normtest.sps'.
normtest vars=x1,x2,x3,x4 .

```

The first line of this syntax (i.e., INCLUDE file='C:\normtest.sps'.) tells SPSS where to find the macro file (the file which was saved to syntax in Step 1). Thus, in the first line of this new syntax file, change 'normtest' to the name of *your* syntax file that holds the macro. The second line (i.e., 'normtest vars=x1,x2,x3,x4 /.') specifies the variables to examine. Thus, change the 'x1' to the name of the first variable *in your dataset* to examine, 'x2' to the name of the second variable in your dataset, and so forth. In this illustration, we are reviewing the standardized residuals (ZRE_1, ZRE_2, ZRE_3) so the second line of our syntax specifies: normtest vars= ZRE_1, ZRE_2, ZRE_3 /.

Fourth, save this two-line syntax file by going to **File** then **Save as** in the top pull-down menu in SPSS. Then click on **Run** in the top horizontal menu in SPSS to generate the output (which follows). In summary, you will have two syntax (i.e., sps) files: one which contains the macro and a second which contains only two lines of syntax that is used to run the macro on the dataset that is opened in SPSS. [Note: If you run into problems in running the syntax, such as getting errors in generating the graphs, run SPSS as an administrator (if you are able) and that should solve the problem.]

Interpreting Multivariate Normality Evidence

What follows is the output from the multivariate normality macro. Generally, there is evidence of a violation of multivariate normality based on the multivariate skewness and kurtosis tests, as well as on the omnibus multivariate normality test.

Output from DeCarlo (1997) SPSS Macro for Multivariate Normality

Run MATRIX procedure:	Double check the sample size ('number of observations')		
Number of observations:	1393		
Number of variables:	3		
Measures and tests of skew:			
ZRE_1	-.2301	-.2298	-3.4764 .0005
ZRE_2	-.2092	-.2090	-3.1679 .0015
ZRE_3	.0370	.0370	.5659 .5715
Measures and tests of kurtosis:			
ZRE_1	-.4173	-.4201	-4.0089 .0001
ZRE_2	-.5561	-.5585	-5.9430 .0000
ZRE_3	-.4846	-.4871	-4.8948 .0000
Omnibus tests of normality (both chisq, 2 df):			
D'Agostino & Pearson K sq	K sq	p-value	Jarque & Bera LM test
ZRE_1	28.1562	.0000	22.5059 .0000
ZRE_2	45.3556	.0000	28.2421 .0000
ZRE_3	24.2788	.0000	14.0897 .0009

Statistically significant univariate skew for reading (ZRE_1) and math (ZRE_2) and kurtosis for all three dependent variables indicate violation of univariate normality.

Statistically significant omnibus tests of univariate normality also indicate violation of multivariate normality.

***** Multivariate Statistics *****

<< Tests of multivariate skew: >>

Small's test (chisq)		
Q1	df	p-value
22.1106	3.0000	.0001

Srivastava's test		
chi(b1p)	df	p-value
10.3201	3.0000	.0160

<< Tests of multivariate kurtosis: >>

A variant of Small's test (chisq)		
VQ2	df	p-value
53.4114	3.0000	.0000

Srivastava's test		
b2p	N(b2p)	p-value
2.8906	-1.4436	.1488

Mardia's test		
b2p	N(b2p)	p-value
13.8359	-3.9662	.0001

Omnibus test of multivariate normality:

(based on Small's test, chisq)		
VQ3	df	p-value
75.5220	6.0000	.0000

----- END MATRIX -----

Statistically significant *multivariate* skew and kurtosis indicate violation of multivariate normality. In this case, there is a violation of multivariate skew for both Small's and Srivastava's tests and a violation of multivariate kurtosis for the variant of Small's and Mardia's test. The only evidence of multivariate normality is based on Srivastava's multivariate kurtosis.

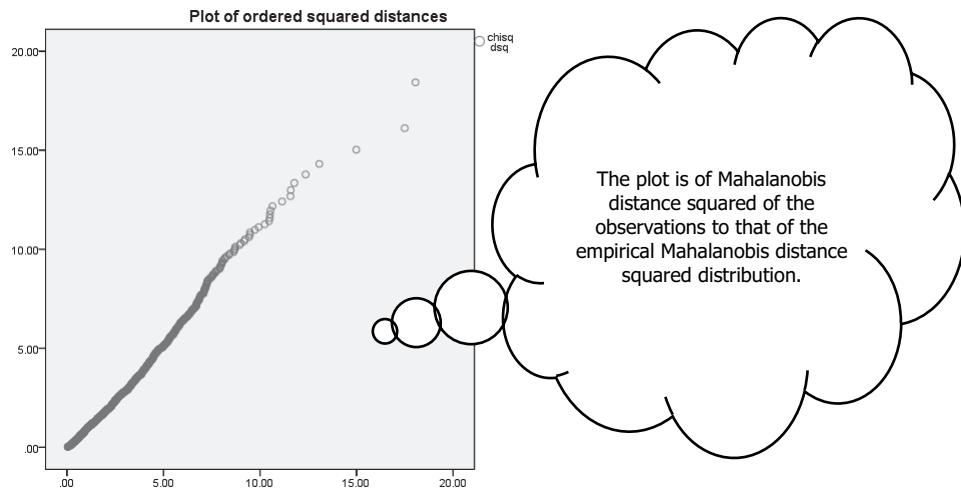
The omnibus test of *multivariate* normality also indicates a violation.

Critical values (Bonferroni) for a single multivar. outlier:

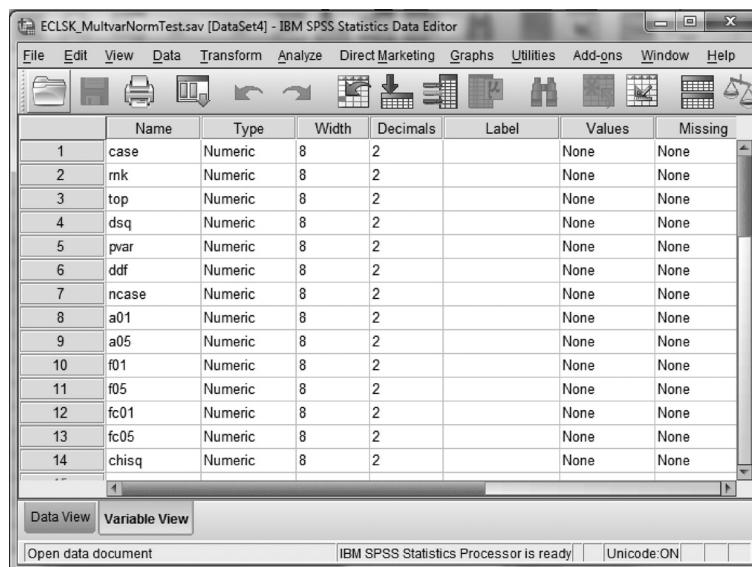
```
critical F(.05/n) =23.08  df = 3,1389
critical F(.01/n) =26.36  df = 3,1389
```

5 observations with largest Mahalanobis distances:

rank = 1	case# =1063	Mahal D sq =	18.05
rank = 2	case# =1153	Mahal D sq =	17.49
rank = 3	case# = 231	Mahal D sq =	14.99
rank = 4	case# = 221	Mahal D sq =	13.07
rank = 5	case# =1165	Mahal D sq =	12.36

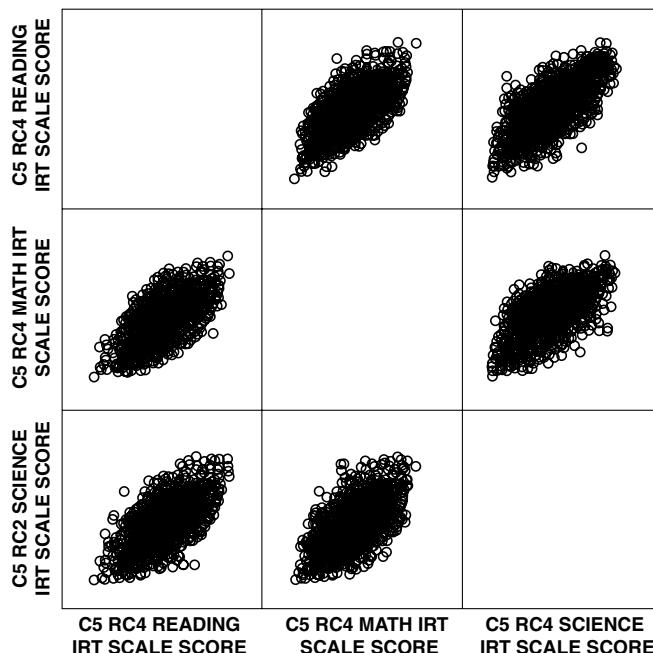


In addition to this output, a new SPSS data file is created that includes the requisite variables generated in the event additional exploration is desired.



6.4.1.3 Linearity

Linearity of the dependent variables, a component of multivariate normality, can be examined by scatterplots of all pairs of dependent variables. The linearity assumption means that a straight line provides a reasonable fit to the data. If the relationship is not a linear one, then the linearity assumption is violated. In our sample, linearity is reasonable given the matrix scatterplot with points suggesting at least moderately strong positive linear relations.



6.4.1.4 Homogeneity of Variance-Covariance Matrices

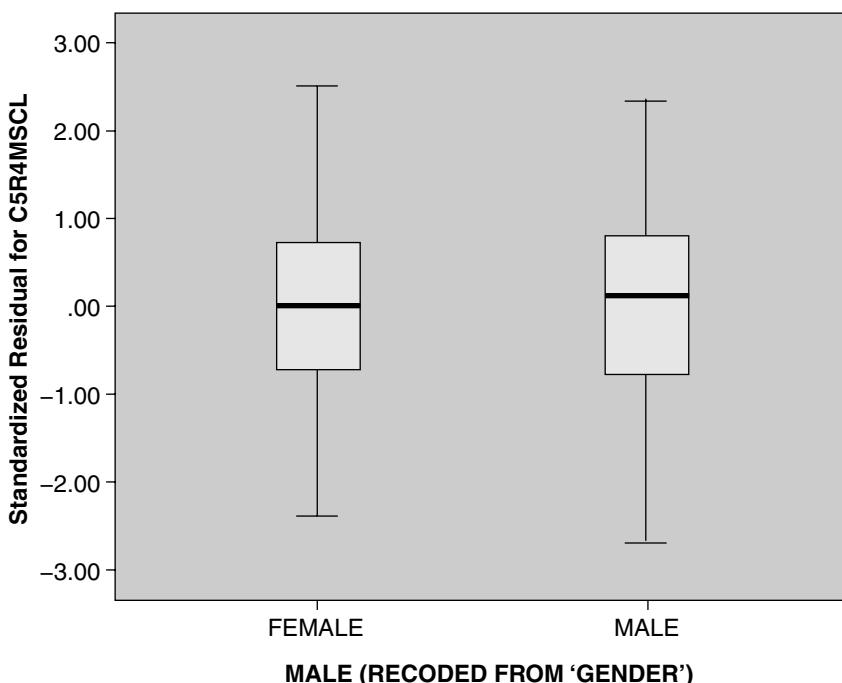
Homogeneity of variance-covariance can be tested using Box's M test, a test that can be generated when computing MANOVA. Statistically significant results indicate a violation of the assumption. As seen in Table 6.3 (SPSS output) and presented here again, Box's M is not statistically significant (Box's $M = 14.019, p = .732$), thus we have evidence of equal variance-covariance matrices.

Box's Test of Equality of Covariance Matrices^a

Box's M	14.019
F	.775
df1	18
df2	2220072.108
Sig.	.732

Tests the null hypothesis that the observed covariance matrices of the dependent variables are equal across groups.

a. Design: Intercept + MALE + TWOPARENT_GR3 + MALE * TWOPARENT_GR3



A more subjective and visual examination of the homogeneity of variance-covariance assumption can be accomplished via spread-versus-level plots and boxplots of the standardized residuals to the factors. Spread-versus-level plots were requested and reviewed in the output (see Table 6.3) and suggest reasonable homogeneity. Boxplots of standardized residuals to the factors suggest evidence of homogeneity, as there are no substantial differences in the box lengths or whisker lengths for the predictors by group. The graph presented here, the standardized residual for math to male, suggests very similar box lengths or whisker lengths for the predictors by group, providing another form of evidence for homogeneity of variance-covariance. Additional boxplots are not presented but suggest similar homogeneity.

6.4.2 Data Screening for Repeated Measures MANOVA

The assumptions for repeated measures MANOVA include (a) independence, (b) multivariate normality of the dependent variables, (c) linearity (a component of multivariate normality), and (d) homogeneity of variance-covariance matrices for the dependent variables.

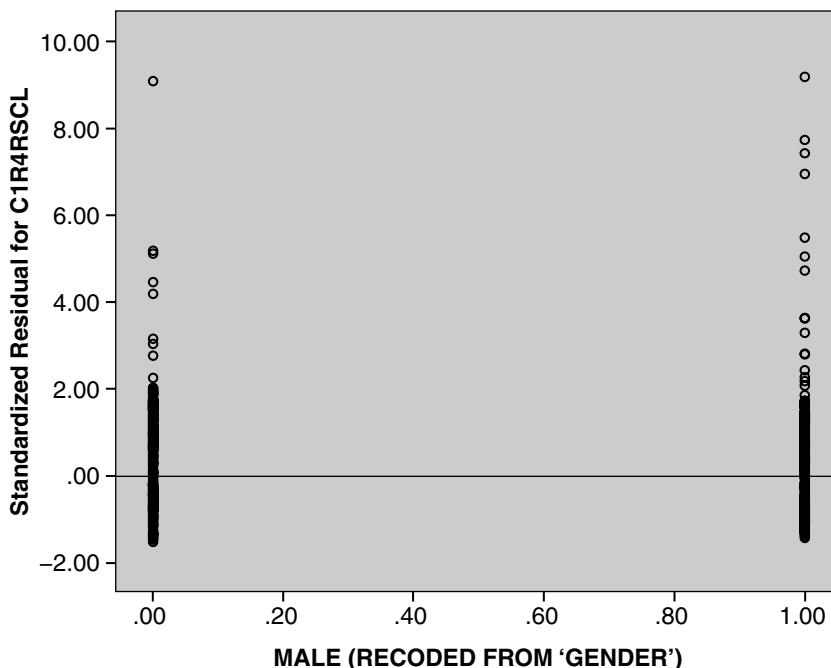
6.4.2.1 Independence

With repeated measures MANOVA, the assumption of independence must consider both the within-subjects factor and the between-subjects factor, if applicable. Just as with the one-way and k -way MANOVA, we can plot residuals against levels of the independent variable(s) in a scatterplot to get an idea of whether or not there

are patterns in the data and thereby provide an indication of whether we have met this assumption. In this example, we have only one between-subjects factor. In cases where there are multiple independent variables, the scatterplot can be split by levels of one independent variable and then a bivariate scatterplot for the other independent variable by residual can be generated. Remember that the residual was added to the dataset by saving it when we generated the repeated measures MANOVA model. Evidence of meeting the assumption is met when the points fall relatively randomly within an absolute value of 2.0. The plots that we generate will give us a general idea of patterns, however, in situations where random assignment was not performed or not possible.

Interpreting Independence Evidence

In examining the scatterplots for evidence of independence, the points should fall relatively randomly above and below a horizontal line within an absolute value of 2.0. In this example, our scatterplot generally suggests evidence of independence with a relatively random display of residuals above and below the horizontal line at zero for both males and females (i.e., each category of the independent variable) with most points falling within an absolute value of 2.0. We repeat this process for each standardized residual. Additional graphs are not presented here, but all are similar to the one provided, suggesting a relatively random display of residuals within a band of -2.0 to +2.0. Thus, since we have not met the assumption of independence through random assignment of cases to groups, this gives us some assurance that independence is a reasonable assumption.



6.4.2.2 Multivariate Normality of the Dependent Variables

Our examination of multivariate normality will commence with review of univariate normality. This will be followed by examination of multivariate normality. (See the data screening chapter for a refresher on generating univariate normality, and one-way and k -way MANOVA for generating multivariate normality evidence.) The skewness statistics of the residuals for reading are within the range of $-.210$ to 3.042 and for math within the range of $-.209$ to 1.438 , all suggesting that normality is relatively reasonable. Kurtosis statistics for reading are within the range of $-.507$ to 18.775 and for math within the range of $-.676$ to 4.845 . While most of the residuals are within the range of an absolute value of 2.0, suggesting some evidence of normality, the residuals for reading measured at the first wave (fall kindergarten) reflect a quite peaked distribution based on the kurtosis statistics (18.775).

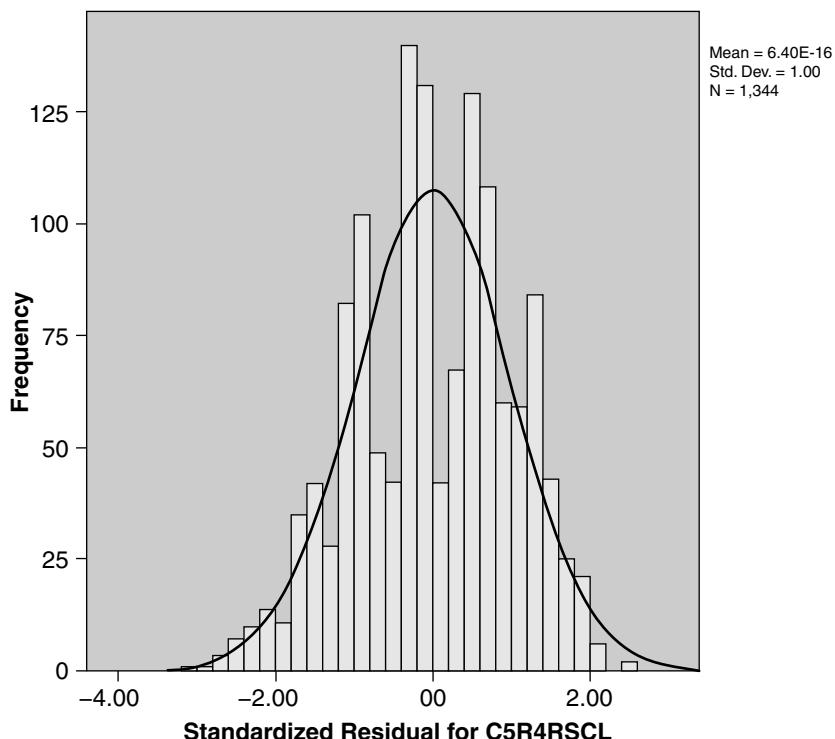
Statistics [READING]

		Standardized Residual for C1R4RSCL	Standardized Residual for C2R4RSCL	Standardized Residual for C3R4RSCL	Standardized Residual for C4R4RSCL	Standardized Residual for C5R4RSCL
N	Valid	1344	1344	1344	1344	1344
	Missing	0	0	0	0	0
Skewness		(-3.042)	(-2.204)	(-2.208)	(-.840)	(-.210)
Std. Error of Skewness		.067	.067	.067	.067	.067
Kurtosis		(18.775)	(7.837)	(6.664)	(.633)	(-.507)
Std. Error of Kurtosis		.133	.133	.133	.133	.133

Statistics [MATH]

		Standardized Residual for C1R4MSCL	Standardized Residual for C2R4MSCL	Standardized Residual for C3R4MSCL	Standardized Residual for C4R4MSCL	Standardized Residual for C5R4MSCL
N	Valid	1344	1344	1344	1344	1344
	Missing	0	0	0	0	0
Skewness		(-1.438)	(-.891)	(.962)	(.524)	(-.209)
Std. Error of Skewness		.067	.067	.067	.067	.067
Kurtosis		(4.845)	(1.484)	(1.853)	(.154)	(-.676)
Std. Error of Kurtosis		.133	.133	.133	.133	.133

As we learned with MANOVA, nonzero kurtosis has minimal effect on the parameter estimates. Thus, the peaked distributions found in this example are not call for alarm. For brevity, only one graph (last measurement occasion for reading) is presented here.



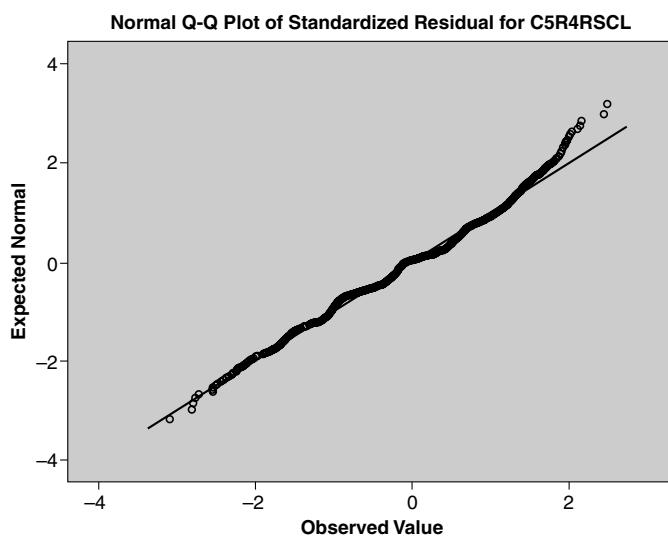
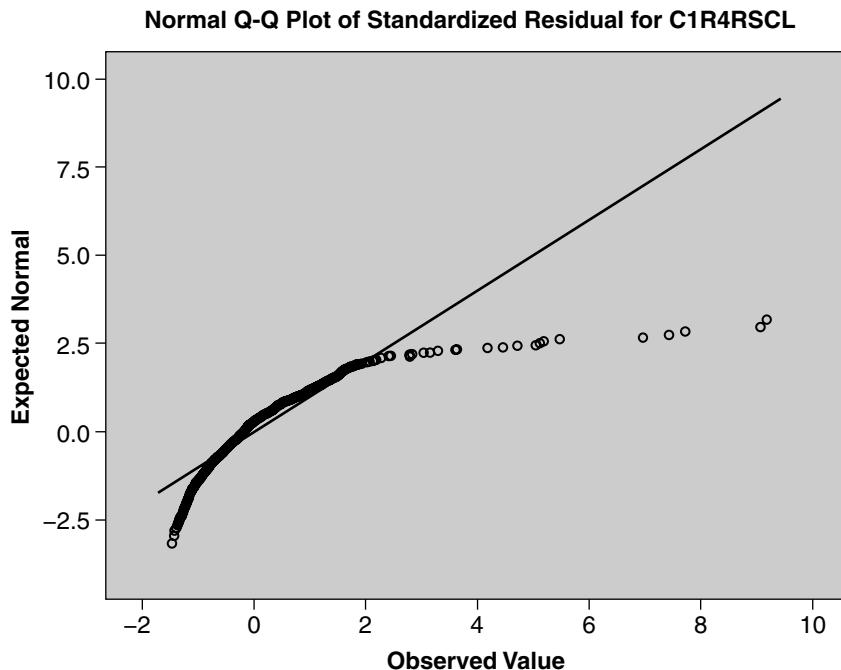
The output for the Shapiro-Wilk test is presented below and suggests that our sample distributions for all residuals are statistically significantly different than what would be expected from a normal distribution ($p < .001$ for all).

Tests of Normality

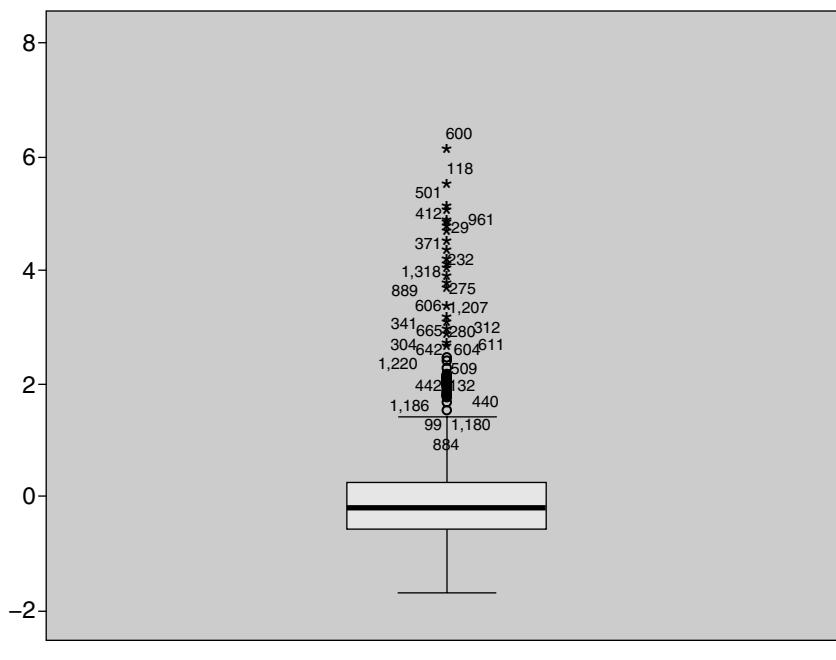
	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Standardized Residual for C1R4RSCL	.115	1344	.000	.795	1344	.000
Standardized Residual for C2R4RSCL	.145	1344	.000	.829	1344	.000
Standardized Residual for C3R4RSCL	.166	1344	.000	.805	1344	.000
Standardized Residual for C4R4RSCL	.101	1344	.000	.954	1344	.000
Standardized Residual for C5R4RSCL	.060	1344	.000	.988	1344	.000
Standardized Residual for C1R4MSCL	.079	1344	.000	.921	1344	.000
Standardized Residual for C2R4MSCL	.053	1344	.000	.959	1344	.000
Standardized Residual for C3R4MSCL	.069	1344	.000	.954	1344	.000
Standardized Residual for C4R4MSCL	.063	1344	.000	.982	1344	.000
Standardized Residual for C5R4MSCL	.042	1344	.000	.986	1344	.000

a. Lilliefors Significance Correction

Two Q-Q plots of residuals are shown here. For brevity, only two graphs are presented, selected as they illustrate the disparity in normality evidence. The first is reading at the first measurement occasion. The lack of adherence to the diagonal line suggests evidence of nonnormality, and this is not surprising, in particular, given the kurtosis statistic for this residual. The second residual plot, also for reading, was conducted at the last measurement occasion, and is more suggestive of normality.



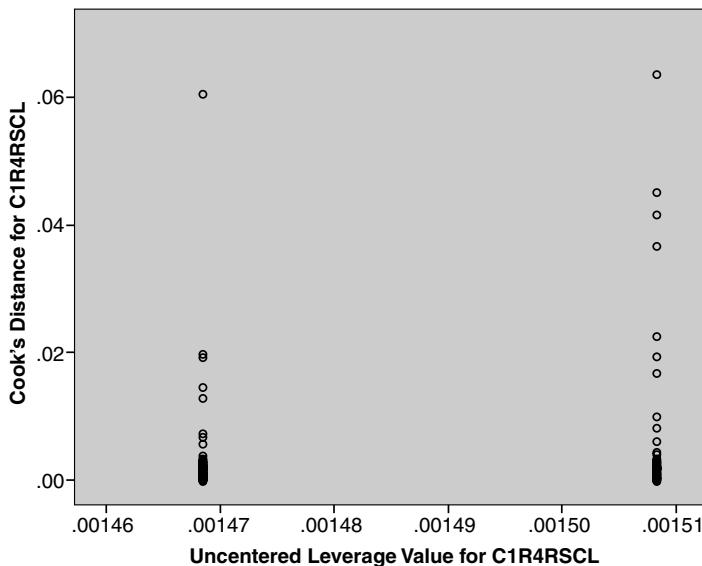
Examination of the boxplot of standardized residuals for reading at the third measurement occasion, presented here, suggests a number of potential outliers. The boxplots of the other residuals (not presented for brevity) are similar, with the first and second measurement occasions generally having a few more potential outliers, the fourth having a few less, and the last measurement occasion having none or very few. Given the number of cases, finding outliers at various measurement occasions is quite common, and it would actually be quite surprising if this were not the case. Because our skewness statistics are generally within an acceptable range, we will not concern ourselves with the potential outliers suggested by the boxplots.



Standardized Residual for C3R4RSCL

Considering the forms of evidence we have examined, skewness and kurtosis statistics, the Shapiro-Wilk test, the Q-Q plot, and the boxplot, there is quite a bit of evidence to suggest nonnormality. It does not, however, appear to be especially problematic given that it centers more on nonzero kurtosis rather than nonzero skew.

There are other diagnostic tools that can be used to screen our data for influential points that may impact univariate or multivariate normality and that can be generated from the values we saved when conducting MANOVA. For example, a plot of Cook's distance against uncentered leverage values provides a way to identify influential cases (i.e., cases with leverage of .50 or above and Cook's distance of 1.0 or greater). Here there are no cases that suggest undue influence. For brevity, the additional graphs are not presented but likewise suggest no evidence of influential points.



Cook's distance provides an overall measure for the influence of individual cases. Values greater than one suggest that the case may be problematic in terms of undue influence on the model. Examining the residual statistics in our output (see following table), we see that the maximum value for Cook's distance is .06, well under the point at which we should be concerned.

Statistics [READING]

	Cook's Distance for C1R4RSCL	Cook's Distance for C2R4RSCL	Cook's Distance for C3R4RSCL	Cook's Distance for C4R4RSCL	Cook's Distance for C5R4RSCL
Minimum	.00	.00	.00	.00	.00
Maximum	.06	.04	.03	.02	.01

Statistics [MATH]

	Cook's Distance for C1R4MSCL	Cook's Distance for C2R4MSCL	Cook's Distance for C3R4MSCL	Cook's Distance for C4R4MSCL	Cook's Distance for C5R4MSCL
Minimum	.00	.00	.00	.00	.00
Maximum	.04	.02	.02	.01	.00

In summary, data screening of *univariate normality* and diagnostics suggests some evidence of nonnormality.

Multivariate Normality Evidence

Although univariate normality is a necessary condition of multivariate normality, it is not a sufficient condition. We will again use DeCarlo's macro (1997) for examining multivariate normality.

Interpreting Multivariate Normality Evidence

What follows is the output from the multivariate normality macro. Generally, there is evidence of a violation of multivariate normality based on the multivariate skewness and kurtosis tests, as well as on the omnibus multivariate normality test.

Output from DeCarlo (1997) SPSS Macro for Multivariate Normality

Run MATRIX procedure:

Number of observations:
1344

Double check the sample size ('number of observations')

Number of variables:
10

Measures and tests of skew:

	g1	sqrt(b1)	z(b1)	p-value
ZRE_1	3.0422	3.0388	24.9894	.0000
ZRE_2	2.2037	2.2013	21.1854	.0000
ZRE_3	2.2078	2.2054	21.2070	.0000
ZRE_4	.8397	.8388	11.0842	.0000
ZRE_5	-.2096	-.2093	-3.1172	.0018
ZRE_6	1.4377	1.4360	16.3925	.0000
ZRE_7	.8906	.8896	11.6121	.0000
ZRE_8	.9623	.9612	12.3282	.0000
ZRE_9	.5244	.5238	7.4297	.0000
ZRE_10	-.2086	-.2084	-3.1030	.0019

Statistically significant univariate skew for all residuals indicates violation of univariate normality.

Measures and tests of kurtosis:

	g2	b2-3	z(b2)	p-value
ZRE_1	18.7751	18.7008	18.2953	.0000
ZRE_2	7.8367	7.8031	14.5584	.0000
ZRE_3	6.6644	6.6352	13.8025	.0000
ZRE_4	.6331	.6263	3.7467	.0002
ZRE_5	-.5067	-.5092	-5.1114	.0000
ZRE_6	4.8453	4.8228	12.2801	.0000
ZRE_7	1.4845	1.4745	6.7747	.0000
ZRE_8	1.8531	1.8417	7.7320	.0000
ZRE_9	.1535	.1485	1.1432	.2530
ZRE_10	-.6756	-.6776	-7.8590	.0000

Statistically significant univariate kurtosis for all but one residual indicates violation of univariate normality.

Omnibus tests of normality (both chisq, 2 df):

	D'Agostino & Pearson K sq		Jarque & Bera LM test	
	K sq	p-value	LM	p-value
ZRE_1	959.1901	.0000	21652.8612	.0000
ZRE_2	660.7678	.0000	4495.1127	.0000
ZRE_3	640.2454	.0000	3554.9035	.0000
ZRE_4	136.8960	.0000	179.5628	.0000
ZRE_5	35.8428	.0000	24.3370	.0000
ZRE_6	419.5157	.0000	1764.4675	.0000
ZRE_7	180.7364	.0000	299.0404	.0000
ZRE_8	211.7670	.0000	396.9173	.0000
ZRE_9	56.5077	.0000	62.6985	.0000
ZRE_10	71.3926	.0000	35.4329	.0000

Statistically significant omnibus tests of univariate normality also indicate violation of normality.

***** Multivariate Statistics *****

< Tests of multivariate skew: >

Small's test (chisq)		p-value	
Q1	df	.0000	
1012.5589	10.0000		
Srivastava's test		p-value	
chi(blp)	df	.0000	
406.0041	10.0000		

Statistically significant *multivariate* skew and kurtosis indicate violation of multivariate normality. In this case, there is a violation of multivariate skew and kurtosis for both Small's and Srivastava's tests.

< Tests of multivariate kurtosis: >

A variant of Small's test (tchisq)		p-value	
VQ2	df	.0000	
658.7969	10.0000		
Srivastava's test		p-value	
b2p	N(b2p)	.0000	
5.1530	50.9487		
Mardia's test		p-value	
b2p	N(b2p)	.0000	
222.8219	121.6605		

Omnibus test of multivariate normality:

(based on Small's test, chisq)		p-value	
VQ3	df	.0000	
1671.3558	20.0000		

The omnibus test of *multivariate* normality also indicates a violation.

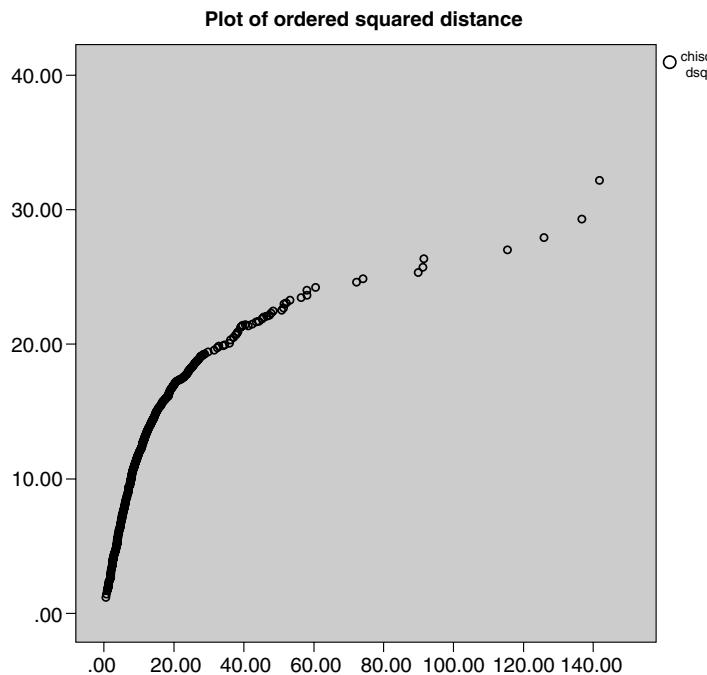
----- END MATRIX -----

Critical values (Bonferroni) for a single multivar. outlier:

```
critical F(.05/n) =37.65  df = 10,1333
critical F(.01/n) =41.52  df = 10,1333
```

5 observations with largest Mahalanobis distances:

rank = 1	case# = 494	Mahal D sq =	141.79
rank = 2	case# = 118	Mahal D sq =	136.75
rank = 3	case# = 1282	Mahal D sq =	125.88
rank = 4	case# = 39	Mahal D sq =	115.55
rank = 5	case# = 1072	Mahal D sq =	91.59



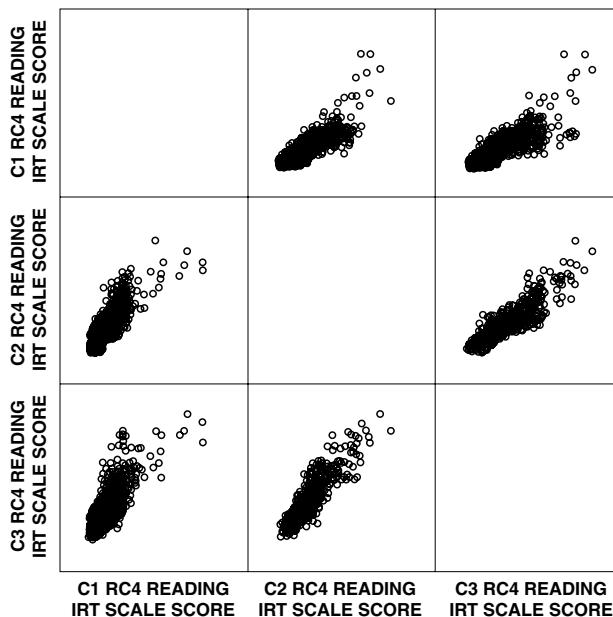
In addition to this output, a new SPSS data file is created that includes the requisite variables generated in the event additional exploration is desired.

	Name
1	case
2	rnk
3	top
4	dsq
5	pvar
6	ddf
7	ncase
8	a01
9	a05
10	f01
11	f05
12	fc01
13	fc05
14	chisq
15	

6.4.2.3 Linearity

Linearity of the dependent variables can be examined by scatterplots of all pairs of dependent variables. The linearity assumption means that a straight line provides a reasonable fit to the data. If the relationship is not a linear one, then the linearity assumption is violated. Given the large number of variables in this illustration, a matrix scatterplot produces indistinguishable blobs. Rather, what is presented here is a

matrix scatterplot of the first three measurement occasions for reading, the points suggesting at least moderate to strong positive linear relations. The remaining scatterplots are similarly interpreted.



6.4.2.4 Homogeneity of Variance-Covariance Matrices

Homogeneity of variance-covariance can be tested using Box's M test, a test that can be generated when computing repeated measures MANOVA. Statistically significant results indicate a violation of the assumption. As seen in Table 6.4 (SPSS output) and presented here again, Box's M is statistically significant (Box's $M = 209.764, p < .001$), thus we have evidence of unequal variance-covariance matrices—a violation of the assumption, in other words.

Box's Test of Equality of Covariance Matrices ^a	
Box's M	209.764
F	3.784
df1	55
df2	5807401.089
Sig.	.000

Tests the null hypothesis that the observed covariance matrices of the dependent variables are equal across groups.

a. Design: Intercept + MALE
Within Subjects Design: time

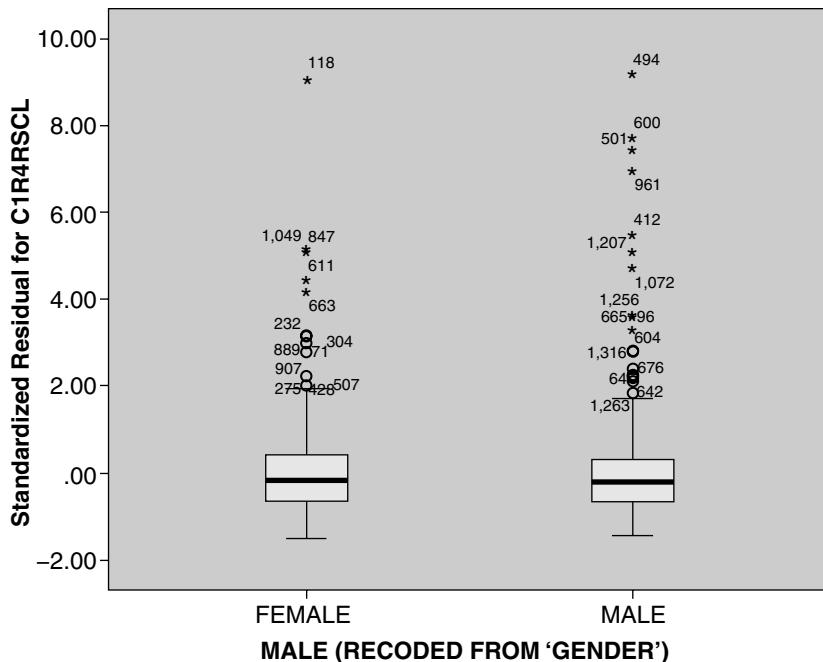
As we learned in our discussion of MANOVA, when the assumption of equal variance-covariance matrices is met, this means that variation or dispersion between groups on the collective dependent variables is equal (or at least not statistically different). Box's

M simultaneously examines the KJ group variances and covariances, and nonstatistically significant Box's M indicates the assumption has been met. Box's test is very sensitive to nonnormality. MANOVA is not robust to violations of this assumption, and this worsens under the following conditions: as the number of dependent variables increases and as the imbalance in the sample sizes per cell increases (particularly when the largest group size is two or more times the size of the smallest group) (Huberty & Olejnik, 2006). If and/or when this assumption is violated, we can use Levene's test for univariate analysis to determine the dependent variable(s) that has/have the heterogeneous variance. Variance stabilizing transformations such as natural log or square root may be considered. For this example, we generated Levene's test via independent t tests. The results for Levene's suggest that four of the five math outcomes have unequal variances. At this point, we will not consider a transformation, and this is due to the following: (a) the ratio of males to females is nearly equal, (b) the results of Box's test may have been influenced by nonnormality of some of the variables, and (c) alpha is too conservative in cases where the cells that have the larger sample sizes also have the larger variance (thus results that are not statistically significant *may* have been had there been equal variance-covariance matrices); in this example, males were the larger n and also had the larger variance for all measurement occasions in math—thus if we *do* find statistical significance in this situation, there is little reason for concern for violation of the assumption.

		Levene's Test for Equality of Variances	
		F	Sig.
C1 RC4 READING IRT SCALE SCORE	Equal variances assumed	.076	.783
C2 RC4 READING IRT SCALE SCORE	Equal variances assumed	.261	.609
C3 RC4 READING IRT SCALE SCORE	Equal variances assumed	1.841	.175
C4 RC4 READING IRT SCALE SCORE	Equal variances assumed	.785	.376
C5 RC4 READING IRT SCALE SCORE	Equal variances assumed	.046	.831
C1 RC4 MATH IRT SCALE SCORE	Equal variances assumed	8.085	.005
C2 RC4 MATH IRT SCALE SCORE	Equal variances assumed	12.545	.000
C3 RC4 MATH IRT SCALE SCORE	Equal variances assumed	9.890	.002
C4 RC4 MATH IRT SCALE SCORE	Equal variances assumed	18.577	.000
C5 RC4 MATH IRT SCALE SCORE	Equal variances assumed	.259	.611

A more subjective and visual examination of the homogeneity of variance-covariance assumption can be accomplished via spread-versus-level plots and boxplots of the standardized residuals to the factors. Spread-versus-level plots were requested and reviewed in the output (see Table 6.4) and suggest reasonable homogeneity. Evidence of homogeneity is present when there are no substantial differences in the box lengths or whisker lengths for the residuals by group. The graph presented here, the standardized

residual for reading at the first measurement occasion to male, suggests very similar box lengths or whisker lengths for the residual by group, providing another form of evidence for homogeneity of variance-covariance. Additional boxplots are not presented but suggest similar homogeneity.



6.5 POWER USING G*POWER

We will examine power from the perspective of one-way and factorial designs as well as from repeated measures MANOVA designs.

6.5.1 Power for One-Way and *k*-Way MANOVA Models

With five or fewer dependent variables, the power of MANOVA equals or exceeds that of ANOVA. In other words, the probability of rejecting the null hypothesis if it is true is higher when conducting MANOVA as compared to ANOVA when there are five or fewer outcomes. We will use G*Power to illustrate how to compute post hoc and a priori power.

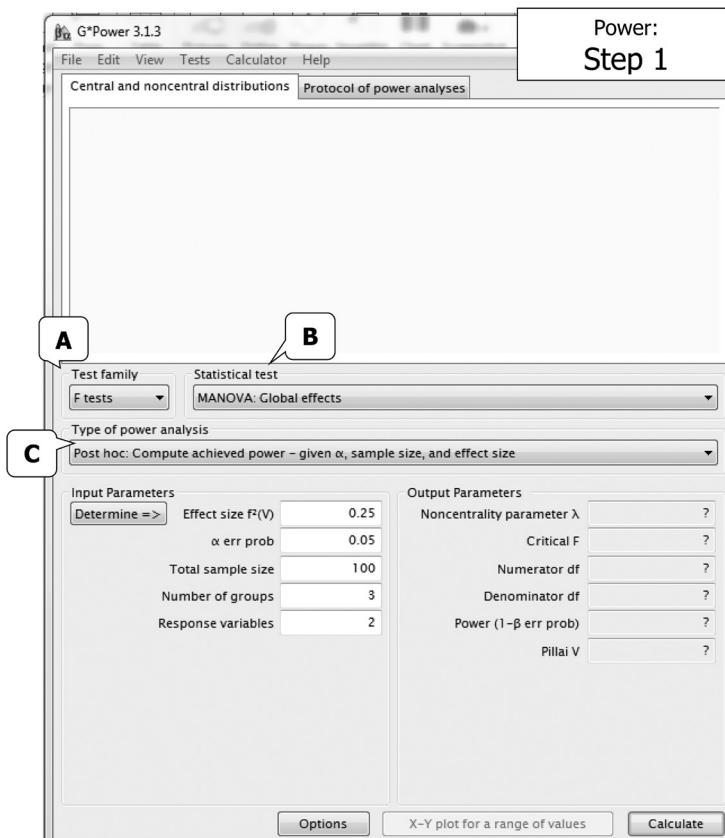
6.5.1.1 Post Hoc Power for Factorial MANOVA Using G*Power

Global Effects

With MANOVA, G*Power provides the option to calculate power both for 'global effects' as well as for 'special effects and interactions.' These are separate calculations, thus researchers who are interested in power for both aspects must conduct two separate power analyses. We will illustrate computing power for global effects, but

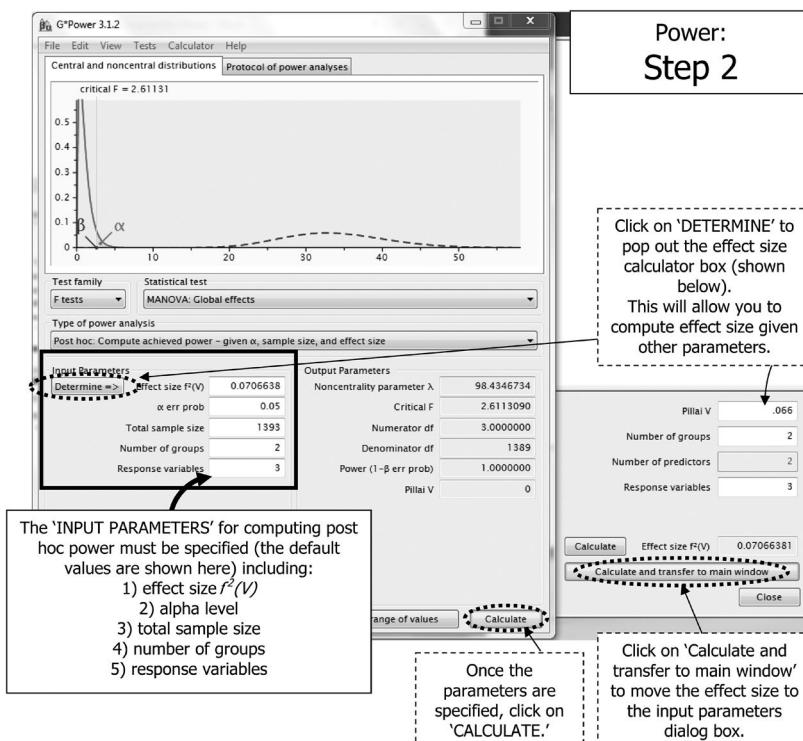
note that computing post hoc power for the other main effect(s) and interaction(s) are similarly obtained.

The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family (see Power: Step 1). In our case, we conducted a factorial MANOVA. To find the MANOVA, we change the 'Test Family' in the drop down menu to 'F tests' and then change the 'Statistical test' to 'MANOVA: Global effects.' The 'Type of power analysis' desired then needs to be selected. To compute post hoc power, we need to select 'Post hoc: Compute achieved power—given α , sample size, and effect size.'

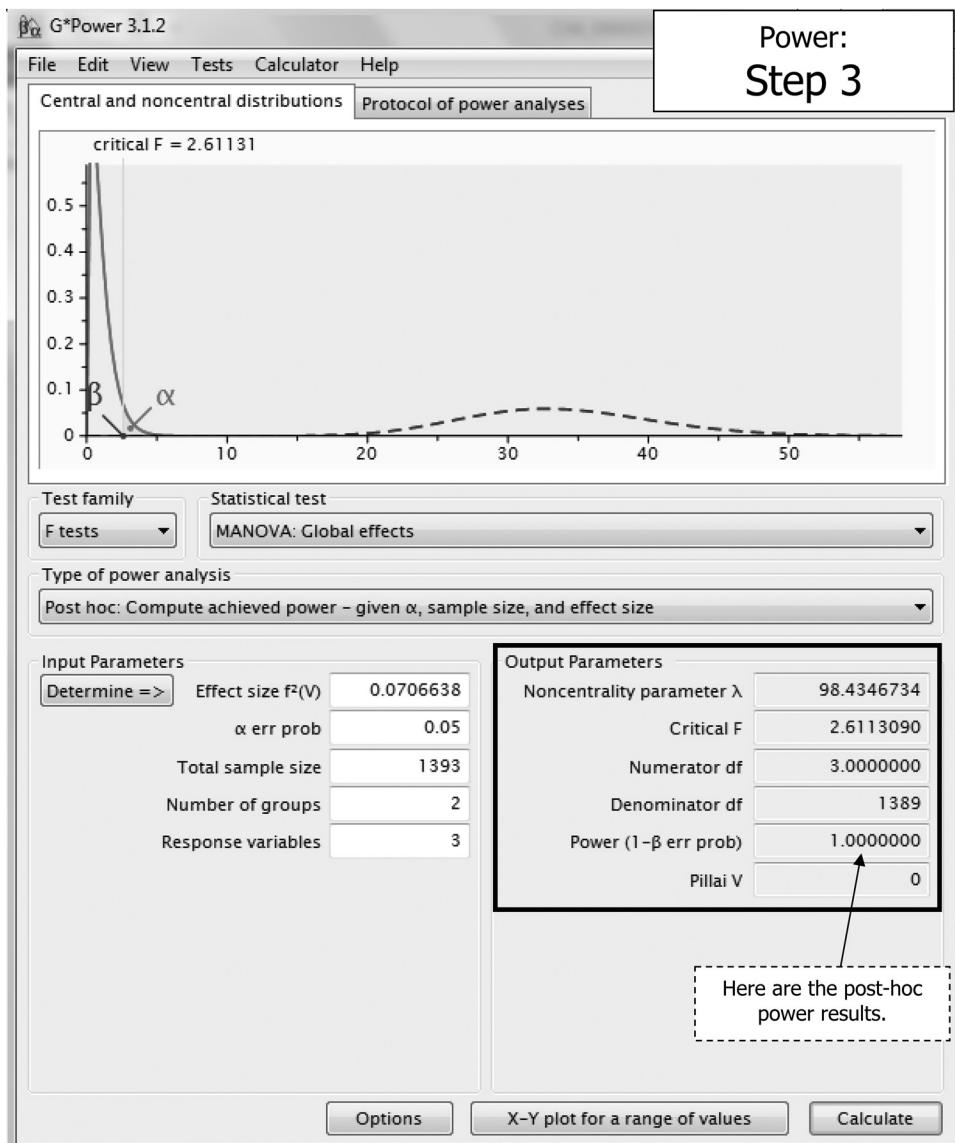


The 'INPUT PARAMETERS' must then be specified (see Power: Step 2). We compute the effect size $f^2(V)$ last, so skip that for the moment. In our example, the alpha level we used was .05 and the total sample size was 1393. Remember that we are calculating 'global effects,' thus we are computing post hoc power for each omnibus main effect. In this example, both factors had two categories (thus the entry for 'number of groups'). Because the effect size differs for each, there is still a need to calculate global main effects for both gender and two-parent home. There are three dependent variables, so we enter '3' for 'response variables.'

We skipped filling in the first parameter, the effect size, for a reason. SPSS only provides a partial eta squared effect size. Thus we will use the pop-out effect size calculator in G*Power to compute the effect size $f^2(V)$. To pop out the effect size calculator, click on 'DETERMINE' which is displayed under 'INPUT PARAMETERS.' In the pop-out effect size calculator, enter Pillai's V for 'male' that was calculated in SPSS (i.e., .066). The 'NUMBER OF GROUPS' for male is 2, and the 'RESPONSE VARIABLES' is 2. Clicking on 'CALCULATE AND TRANSFER TO MAIN WINDOW' in the pop-out effect size calculator will calculate the effect size and transfer the calculated effect size (i.e., .07066) to the 'INPUT PARAMETERS.' Once the parameters are specified, click on 'CALCULATE' to find the power statistics.

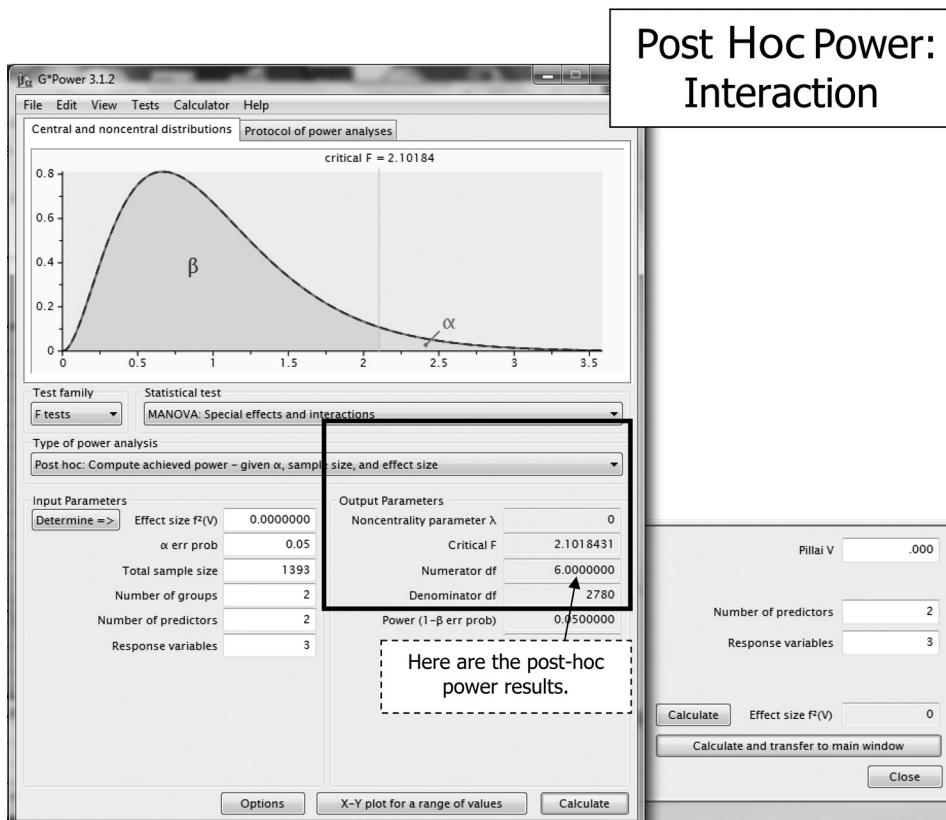


The 'OUTPUT PARAMETERS' provide the relevant statistics given the input just specified (see Power: Step 3). In this example, we were interested in determining post hoc power for the global effects for one of the factors in a factorial MANOVA with a computed effect size $f^2(V)$ of .07, alpha level of .05, total sample size of 1393, two groups in our factor, and three dependent variables. Based on those criteria, the post hoc power for the main effect of attractiveness was 1.00. In other words, the probability of rejecting the null hypothesis when it is really false was 1.00, which would be considered maximum power (sufficient power is often .80 or above). Note that this value is the same as that reported in SPSS. Keep in mind that conducting power analysis a priori is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters).



Interactions

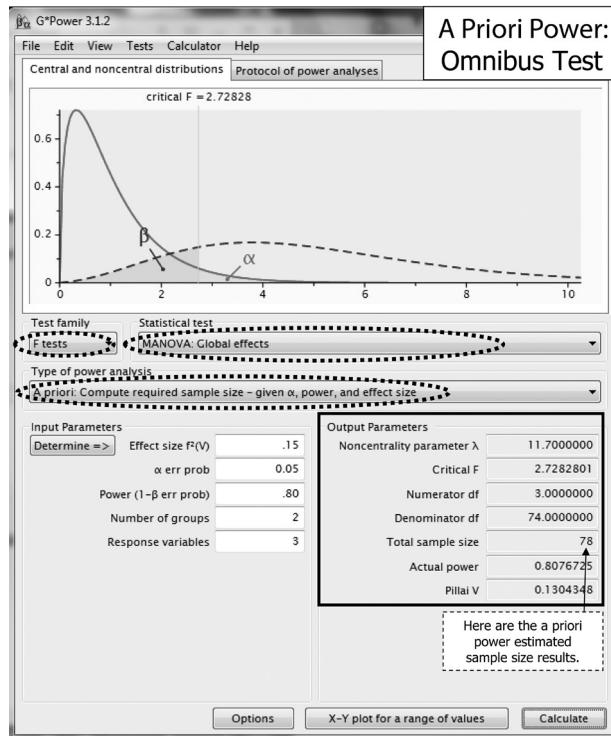
Calculation of power for interactions is conducted similarly (see post hoc power interaction screenshot on the next page) but using Pillai's value for the interaction omnibus test instead. The input of .000 for Pillai's omnibus interaction results in the following output for interaction power. The post hoc power of the interaction effect for this test was .05—the probability of rejecting the null hypothesis when it is really



false (in this case, the probability that the means of the dependent variable would be equal for each cell) was about 5%, which would be considered very low power (sufficient power is often .80 or above). Note that this value is very close to that reported in SPSS (.08).

6.5.1.2 A Priori Power for Factorial MANOVA Using G*Power

For a priori power, we can determine the total sample size needed for the main effects and/or interactions given an estimated effect size $f^2(V)$, alpha level, desired power, number of groups (i.e., number of categories of the independent variable or interaction, depending on which a priori power is of interest), and number of response variables (i.e., dependent variables) (see a priori power omnibus test screenshot on the next page). We follow Cohen's (1988) conventions for effect size (i.e., small $f^2 = .02$; moderate $f^2 = .15$; large $f^2 = .35$). In this example, had we estimated a moderate effect of .15, alpha of .05, desired power of .80, two groups and three dependent variables, we would need a total sample size of 78 (or about 39 individuals per cell) to detect global main effects.



6.5.2 Power for Repeated Measures MANOVA

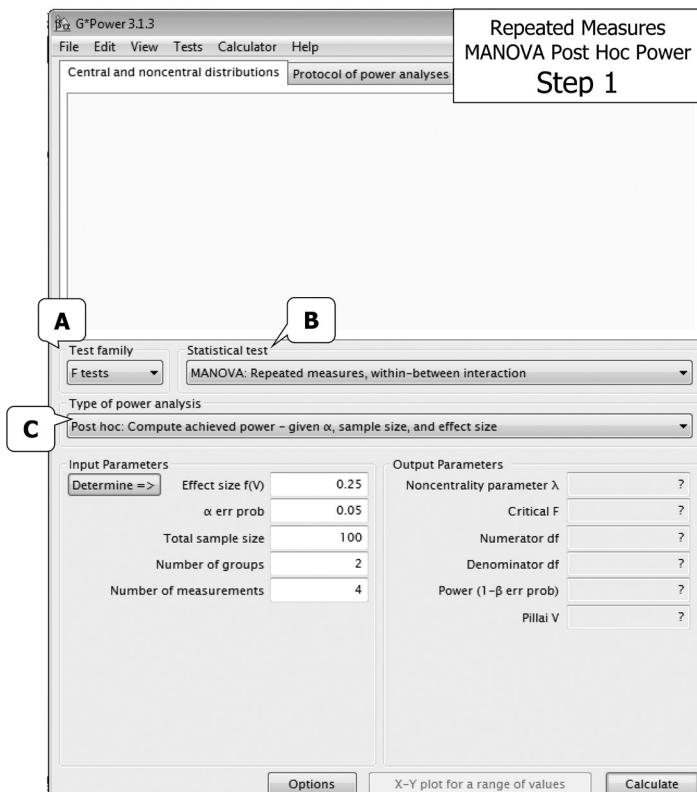
We will use G*Power to illustrate how to compute post hoc and a priori power for repeated measures MANOVA.

6.5.2.1 Post Hoc Power for Repeated Measures MANOVA Using G*Power

Within-Between Interaction

With repeated measures MANOVA, G*Power provides the option to calculate power for within and/or between factors as well as the within-between interaction. These are separate calculations, thus researchers who are interested in power for the other aspects must conduct separate power analyses. For illustrative purposes, we will generate power for the within-between interaction, but note that computing post hoc power for the other effects are similarly obtained.

The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family. In our case, we conducted a repeated measures MANOVA. To find the MANOVA, we change the 'Test Family' in the drop down menu to 'F tests' (see Step 1 part A) and then change the 'Statistical test' to 'MANOVA: Repeated measures, within-between interaction' (see Step 1 part B). The 'Type of power analysis' desired then needs to be selected. To compute post hoc power, we need to select 'Post hoc: Compute achieved power—given α , sample size, and effect size' (see Repeated Measures MANOVA Post Hoc Power: Step 1).



The 'INPUT PARAMETERS' must then be specified (see Repeated Measures MANOVA Post Hoc Power: Step 2). We compute the effect size $f(V)$ last, so skip that for the moment. In our example, the alpha level we used was .05 and the total sample size was 1344. Remember that we are calculating the 'within-between interaction' thus we must specify the number of groups of the between-subjects factor. In this example, we had two categories (thus the entry for 'number of groups'). There are five measurement occasions, so we enter '5' for 'number of measurements.'

We skipped filling in the first parameter, the effect size, for a reason. SPSS only provides a partial eta squared effect size. Thus we will use the pop-out effect size calculator in G*Power to compute the effect size $f(V)$. To pop out the effect size calculator, click on 'DETERMINE' which is displayed under 'INPUT PARAMETERS.' In the pop-out effect size calculator, enter Pillai's V for 'Male' that was calculated in SPSS (i.e., .056). The 'NUMBER OF GROUPS' for male is 2, and the 'NUMBER OF MEASUREMENTS' is 5. Clicking on 'CALCULATE AND TRANSFER TO MAIN WINDOW' in the pop-out effect size calculator will calculate the effect size and transfer the calculated effect size (i.e., .2435612) to the 'INPUT PARAMETERS.' Once the parameters are specified, click on 'CALCULATE' to find the power statistics (see Repeated Measures MANOVA Post Hoc Power: Step 3).

**Repeated Measures
MANOVA Post Hoc Power
Step 2**

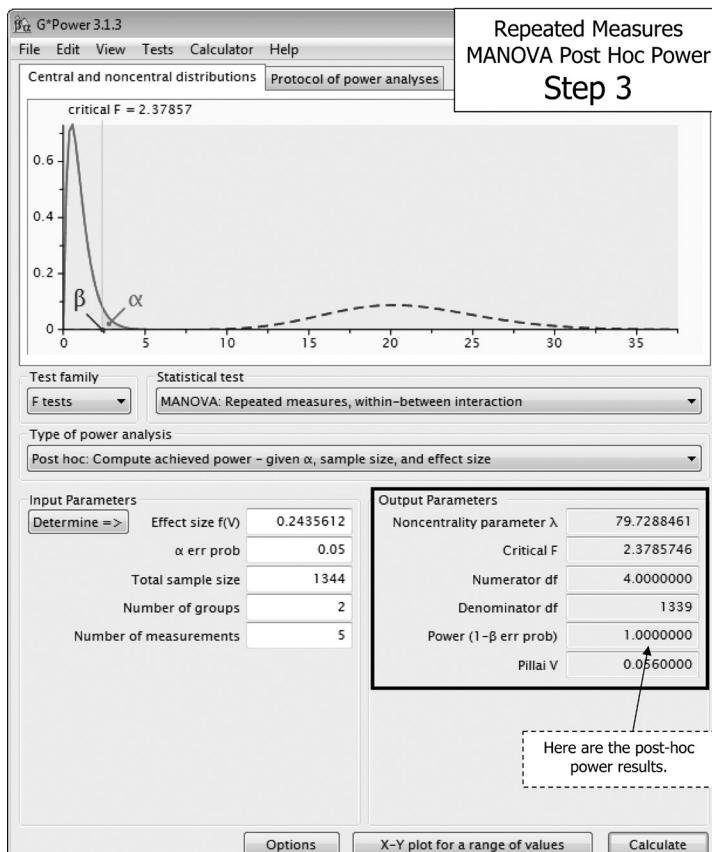
The 'INPUT PARAMETERS' for computing post hoc power must be specified (the default values are shown here) including:

- 1) effect size $f^2(V)$
- 2) alpha level
- 3) total sample size
- 4) number of groups
- 5) number of measurements

Click on 'DETERMINE' to pop out the effect size calculator box (shown below). This will allow you to compute effect size given other parameters.

Once the parameters are specified, click on 'CALCULATE.'

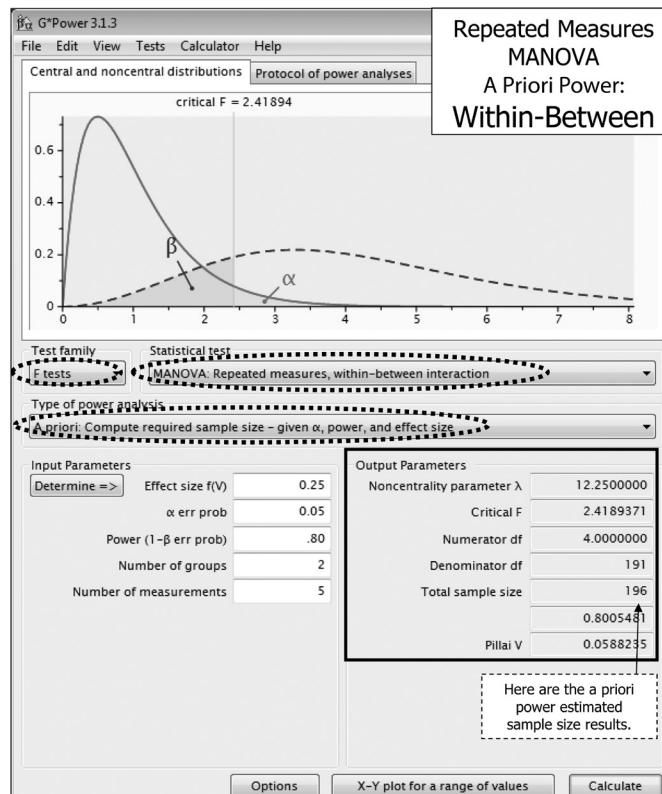
Click on 'Calculate and transfer to main window' to move the effect size to the input parameters dialog box.



The 'OUTPUT PARAMETERS' provide the relevant statistics given the input just specified. In this example, we were interested in determining post hoc power for the within-between interactions in a repeated measures MANOVA with a computed effect size $f(V)$ of approximately .24, alpha level of .05, total sample size of 1344, two groups in our between-subjects factor, and five measurement occasions. Based on those criteria, the post hoc power for the main effect of attractiveness was 1.00. In other words, the probability of rejecting the null hypothesis when it is really false was 1.00, which would be considered maximum power (sufficient power is often .80 or above). Note that this value is the same as that reported in SPSS. Keep in mind that conducting power analysis a priori is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters).

6.5.2.2 A Priori Power for Repeated Measures MANOVA Using G*Power

For a priori power, we can determine the total sample size needed for the within and/or between factors and within-between interaction. Let's again illustrate using the within-between interaction. We can generate a priori power for the interaction given an estimated effect size $f(V)$, alpha level, desired power, number of groups (i.e., number of categories of the between-subjects factor), and number of measurement occasions. In this example, we will leave the estimated effect as the default $f(V)$ as .25. We will define an alpha of .05, desired power of .80, two groups, and five measurement occasions. Given these parameters, we would need a total sample size of 196 to detect the within-between interaction.



6.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

6.6.1 Research Question Template and Example Write-Up for One-Way and *k*-Way MANOVA Models

Recall that Challie and Ott were assisting Dr. Childs in examining group differences in collective academic outcomes using data from the ECLS-K. The research question presented to Dr. Childs was the following: Is there a difference in the reading, mathematics, and science outcomes jointly between boys and girls and between children from two-parent versus other family structures?

Challie and Ott then assisted Dr. Childs in conducting MANOVA, and a template for writing the research question for MANOVA is presented below.

Is there a mean difference in the [list dependent variables] jointly between [list independent variables]?

It may be helpful to preface the results of the MANOVA with information on an examination of the extent to which the data were thoroughly screened.

Prior to analysis, reading, mathematics, and science IRT scale scores were examined for accuracy in data entry, missing values, and the extent to which MANOVA assumptions were met. Frequency distributions of the independent variables suggested that the range of values were within what was to be expected. An examination of means and standard deviations also suggested data accuracy and the absence of values that may have not been coded properly (e.g., missing data codes).

In terms of missing data, after delimitations there was very little missing data ($n = 3$ cases with missing data). These cases were excluded from the analyses, resulting in a final sample size of 1,393.

The assumptions of MANOVA were examined. These include independence, multivariate normality of dependent variables, linearity (a component of multivariate normality), and homogeneity of variances-covariances. Each of the assumptions were thoroughly examined prior to data analysis.

Independence. The assumption of independence was reviewed by plotting standardized residuals against levels of the independent variables in a scatterplot. The scatterplots generally suggest evidence of independence with a relatively random display of residuals above and below the horizontal line at zero for each category of the independent variables that were used to split the file. Thus, since we have not met the assumption of independence through random assignment of cases to groups, this gives us some assurance that independence is a reasonable assumption.

Multivariate normality. Standardized residuals were examined to determine univariate normality, a necessary condition for multivariate normality. The skewness statistics of the standardized residuals were within the range of $-.230$ to $.037$ and kurtosis within the range of $-.556$ to $-.417$, all within the range of an absolute value of 2.0 and 7.0, respectively, suggesting evidence of normality. However, only the skewness of the standardized residual for science was not statistically significant based on DeCarlo's SPSS macro (1997). The histograms of standardized residuals were also relatively normal. The Q-Q plot of residuals

suggest relative normality with the only nonnormality suggested in the tails of the distribution. Examination of the boxplot of standardized residuals for reading suggests a relatively normal distributional shape of residuals with two potential outliers. The boxplots of math and science standardized residuals were also relatively normal in shape with no outliers. Shapiro-Wilk's formal test for normality suggests that the sample distributions for standardized residuals are statistically significantly different from what would be expected from a normal distribution ($p < .001$ for all). Influential points were examined by plotting Cook's distance against uncentered leverage values, and the plot did not suggest any cases exerting undue influence. Cook's distance provides an overall measure for the influence of individual cases. Values greater than one suggest that the case may be problematic in terms of undue influence on the model. Examining the residual statistics, the maximum value for Cook's distance is .01, well under the point at which we should be concerned. Generally, most forms of evidence suggest normality is a reasonable assumption. In addition, although testing for univariate normality doesn't guarantee multivariate normality, departures from multivariate normality are usually negligible when univariate normality is met for each variable. In this situation, multivariate normality was examined using DeCarlo's (1997) SPSS macro for multivariate normality and did not appear to be met. Multivariate skewness (Small's $\chi^2 = 22.11$, $df = 3$, $p = .0001$; Srivastava's $\chi = 10.32$, $df = 3$, $p = .02$) and multivariate kurtosis (Small's variant $\chi^2 = 53.41$, $df = 3$, $p < .001$; Mardia's = 13.86, $N(b2p) = -3.97$, $p < .001$) were suggestive of a violation of multivariate normality as was the omnibus test of multivariate normality (Small's test variant $\chi^2 = 75.52$, $df = 6$, $p < .001$). The only exception was one test of multivariate kurtosis which suggested multivariate normality (Srivastava's = 2.89, $N(b2p) = -1.44$, $p = .15$). Violations of multivariate normality have minimal effect on Type I errors (i.e., rejecting the null hypothesis when it is true).

Linearity. Linearity of the dependent variables was examined by scatterplots of all pairs of dependent variables. All scatterplots suggested at least a moderately strong positive linear relation.

Homogeneity of variances-covariances. In terms of homogeneity of variances-covariances, boxplots of the dependent variables to predictors were examined as a visual means to determine the extent to which equal variances can be assumed. There were no substantial differences in the box lengths or whisker lengths for the predictors by group suggesting evidence of homogeneity. Spread-versus-level plots were also examined and a relatively random display of points provided another form of visual evidence that this assumption was met. As a formal test, Box's M ($M = 14.019$) provides evidence of equal variance-covariance matrices, $F(18, 2220072) = .775$, $p = .732$.

Here is an example summary of results for MANOVA (remember that this will be prefaced by the previous information reporting the extent to which the data were thoroughly screened).

Using data from the ECLS-K, a 2×2 multivariate analysis of variance (MANOVA) was conducted to determine whether there were simultaneous mean differences in academic outcomes (specifically reading, mathematics, and science IRT scale scores) based on two predictors (gender with 'male' coded as 1 and 'female' as 0; family structure with 'two-parent home' coded as 1 and 'other family structure' coded as 0). The omnibus Wilks's lambda was statistically significant

for the main effect for males indicating that the combined dependent variables differed, on average, between males and females (i.e., there is a difference based on gender, collapsing across time), $\Lambda = .934$, $F(3, 1387) = 32.561$, $p < .001$, partial $\eta^2 = .066$. Partial eta squared suggests a moderate main effect for males. The overall Wilks's lambda was statistically significant for the main effect for two-parent home indicating that the combined dependent variables differed, on average, between children from two-parent as compared to other family structures, $\Lambda = .933$, $F(3, 1387) = 33.3821$, $p < .001$, partial $\eta^2 = .067$. Partial eta squared suggests a moderate main effect for two-parent home. The omnibus Wilks's lambda was not statistically significant for the interaction between male and two-parent home, $\Lambda = 1.00$, $F(3, 1387) = .909$, partial $\eta^2 < .001$. Partial eta squared for the interaction suggests an inconsequential effect size for the interaction of gender and parental home environment.

Note: If there are statistically significant main effects, interactions, or a priori tests, include a section that details the follow-up using discriminant analysis. A preface to this section is provided here, and a full write-up is included in the discriminant analysis chapter.

Follow-up discriminant analysis. Discriminant analyses were conducted as follow-ups to the statistically significant main effects to determine if [statistically significant main effects/interactions/a priori contrasts] could be differentiated based on [dependent variables]. A discriminant analysis was conducted to determine . . .

6.6.2 Research Question Template and Example Write-Up for Repeated Measures MANOVA

Recall that the graduate research students were assisting Dr. Childs in examining repeated noncommensurate outcomes for boys and girls using data from the ECLS-K. The research question presented to Dr. Childs was the following: Is there a difference over time in reading and mathematics outcomes between boys and girls?

Challie, Ott, and Addie then assisted Dr. Childs in conducting a repeated measures MANOVA, and a template for writing the research question for repeated measures MANOVA is presented below. Note that this addresses both the within- and the between-factors.

Is there a mean difference over time in the [list noncommensurate dependent variables] jointly between [list independent variable(s)]?

It may be helpful to preface the results of the repeated measures MANOVA with information on an examination of the extent to which the data were thoroughly screened.

Prior to analysis, reading and mathematics IRT scale scores were examined for accuracy in data entry, missing values, and the extent to which repeated measures MANOVA assumptions were met. Frequency distributions of the independent variables suggested that the range of values were within what was to be expected. An examination of means and standard deviations also

suggested data accuracy and the absence of values that may have not been coded properly (e.g., missing data codes).

In terms of missing data, after delimitations there was very little missing data (approximately 2%). These cases were excluded from the analyses, resulting in a final sample size of 1,344.

The assumptions of repeated measures MANOVA were examined. These include independence, multivariate normality of dependent variables, linearity (a component of multivariate normality), and homogeneity of variances/covariances. Each of the assumptions were thoroughly examined prior to data analysis.

Independence. The assumption of independence was reviewed by plotting standardized residuals against levels of the independent variable in a scatterplot. The scatterplots generally suggest evidence of independence with a relatively random display of residuals above and below the horizontal line at zero and within a band of an absolute value of 2.0 for each category of the independent variable. Thus, since we have not met the assumption of independence through random assignment of cases to groups, this gives us some assurance that independence is a reasonable assumption.

Multivariate normality. Standardized residuals were examined to determine univariate normality, a necessary condition for multivariate normality. The skewness statistics of the standardized residuals ranged from $-.210$ to 3.042 (only one of which was above 3.0 and only two above 2.0), generally suggesting normality. Eight of the ten residuals had kurtosis statistics ranging from $-.676$ to 6.664 , smaller than the recommended 7.0 suggested by West (1996), again suggesting evidence of normality. Two residuals, however, had kurtosis values suggesting nonnormality (7.837 and 18.775). Based on DeCarlo's SPSS macro (1997) univariate skewness and kurtosis were statistically significant, indicating nonnormality ($p < .001$). The histograms and Q-Q plots of standardized residuals were commensurate with the skew and kurtosis statistics, suggesting some nonnormality for a few residuals. Examination of the boxplots of standardized residuals suggest some potential outliers for most measurement occasions. Shapiro-Wilk's formal test for normality suggests that the sample distributions for standardized residuals are statistically significantly different from what would be expected from a normal distribution ($p < .001$ for all). Influential points were examined by plotting Cook's distance against uncentered leverage values, and the plots did not suggest any cases exerting undue influence. Cook's distance provides an overall measure for the influence of individual cases. Values greater than one suggest that the case may be problematic in terms of undue influence on the model. Examining the residual statistics, the maximum value for Cook's distance is $.06$, well under the point at which we should be concerned. In aggregate, the indices suggest some nonnormality, particularly for the early measurement occasions. The departures from normality are not particularly concerning as it appears to be confined to nonzero kurtosis. Although testing for univariate normality doesn't guarantee multivariate normality, departures from multivariate normality are usually negligible when univariate normality is met for each variable. In this situation, multivariate normality was examined using DeCarlo's (1997) SPSS macro for multivariate normality and did not appear to be met. Multivariate skewness (Small's $\chi^2 = 1012.56$, $df = 10$, $p < .001$; Srivastava's $\chi = 406.00$, $df = 10$, $p < .001$) and multivariate kurtosis (Small's

variant $\chi^2 = 658.80$, $df = 10$, $p < .001$; Srivastava's = 5.15, $N(b2p) = 50.95$, $p < .001$). Mardia's = 222.82, $N(b2p) = 121.66$, $p < .001$) were suggestive of a violation of multivariate normality as was the omnibus test of multivariate normality (Small's test variant $\chi^2 = 1671.3558$, $df = 20$, $p < .001$). Violations of multivariate normality have minimal effect on Type I errors (i.e., rejecting the null hypothesis when it is true).

Linearity. Linearity of the dependent variables was examined by scatterplots of all pairs of dependent variables. All scatterplots suggested at least a moderate to strong positive linear relation.

Homogeneity of variances-covariances. In terms of homogeneity of variances-covariances, boxplots of each measurement occasion of the dependent variables to the independent variable were examined as a visual means to determine the extent to which equal variances can be assumed. There were no substantial differences in the box lengths or whisker lengths for the predictors by group suggesting evidence of homogeneity. Spread-versus-level plots were also examined and a relatively random display of points provided another form of visual evidence that this assumption was met. As a formal test, Box's M ($M = 209.76$) suggested a violation of the assumption of equal variance/covariance matrices, $F(55, 5807401.09) = 3.784$, $p < .001$. Levene's test for univariate analysis was computed to determine the dependent variable(s) that had heterogeneous variance. The results for Levene's suggest that four of the five math outcomes had unequal variances. Variance stabilizing transformations such as natural log or square root are sometimes considered in this situation, however this was not a consideration in this study due to the following: (a) the ratio of males to females is nearly equal, (b) the results of Box's test may have been influenced by nonnormality of some of the variables, and (c) alpha is too conservative in cases where the cells that have the larger sample sizes also have the larger variance (thus results that are not statistically significant *may* have been had there been equal variance-covariance matrices); in this example, males were the larger n and also had the larger variance for all measurement occasions in math—thus if we *do* find statistical significance in this situation, there is little reason for concern for violation of the assumption.

Here is an example summary of results for repeated measures MANOVA (remember that this will be prefaced by the previous information reporting the extent to which the data were thoroughly screened).

Using data from the ECLS-K, a repeated measures multivariate analysis of variance (MANOVA) was conducted to determine whether there were simultaneous mean differences over time (five measurement occasions beginning in fall kindergarten and through spring third grade) in academic outcomes (specifically reading and mathematics IRT scale scores) based on gender (with 'male' coded as 1 and 'female' as 0). The omnibus Wilks's lambda was statistically significant for the main effect for males, indicating that the combined dependent variables differed over time, on average, between males and females (i.e., there is a different between gender, collapsing across time), Pillai's trace = .056, $F(2, 1341) = 40.147$, $p < .001$, partial $\eta^2 = .056$. Partial eta squared suggests a moderate main effect for males. The overall Pillai's trace

was statistically significant for the main effect of time indicating differences over time on the combination of dependent variables, Pillai's trace = .956, $F(8, 1335) = 3657.97, p < .001$, partial $\eta^2 = .956$. Partial eta squared suggests a large main effect for time. The omnibus test was also statistically significant for the interaction between gender and time, Pillai's trace = .051, $F(8, 1335) = 8.990$, partial $\eta^2 = .051$. Partial eta squared for the interaction suggests a moderate effect size for the interaction of gender and time. Table χ presents the means and standard deviations by measurement occasion for reading and math by gender. [Table χ has not been presented here; however, you will want to do so in your study.] The within-subjects contrast indicates that for the interaction of gender by time, the results are statistically significant only for mathematics. In other words, there is differential linear ($F = 31.045, p < .001$, partial $\eta^2 = .023$) and quadratic ($F = 682.201, p = .002$, partial $\eta^2 = .007$) growth for boys and girls over time but only in mathematics, not reading. In both cases, there is a small effect.

PROBLEMS

Conceptual Problems: One-Way and k -Way MANOVA Models

1. Marie wants to know which type of counseling intervention produces the biggest mean difference in a composite measure of marriage satisfaction. She collects the following data on 250 couples who have completed marriage counseling: (a) type of intervention (three categories: couple; small group; large group), (b) Marital Satisfaction Questionnaire (interval measurement scale), (c) Couples' Happiness Scale (interval measurement scale), and (d) Marriage Inventory (interval measurement scale). Is computing a MANOVA appropriate given this scenario?
 - a. Yes
 - b. No
2. In univariate data screening, Oscar finds that the standardized residuals for the three continuous outcome variables for his MANOVA have skewness values ranging from $-.72$ to $+1.13$ and kurtosis values ranging from $.89$ to $+1.20$. He conducts the Shapiro-Wilk test and finds p values ranging from $.32$ to $.76$. Oscar runs DeCarlo's (1997) macro and finds statistically significant multivariate skew and kurtosis and omnibus multivariate tests for standardized residuals ($p < .001$). Which of the following is a reasonable assessment based on this information?
 - a. There is some evidence to suggest the assumption of homogeneity of variance-covariance is met.
 - b. There is some evidence to suggest the assumption of homogeneity of variance-covariance has been violated.
 - c. There is some evidence to suggest the assumption of linearity is met.
 - d. There is some evidence to suggest the assumption of linearity has been violated.
 - e. There is some evidence to suggest the assumption of multivariate normality is met.
 - f. There is some evidence to suggest the assumption of multivariate normality has been violated.

3. Your friend has conducted a 3×4 MANOVA and found a statistically significant interaction effect. What do you suggest to your friend to do now?
 - a. Compute a priori power
 - b. Conduct discriminant analysis
 - c. Generate univariate ANOVA
 - d. Have a party and celebrate
4. In generating MANOVA, you find a p value of .002 for Box's test. What is implied by this finding?
 - a. The assumption of equal variance-covariance has not been met.
 - b. The assumption of linearity has been violated.
 - c. The correlations among the dependent variables are not zero.
 - d. The omnibus F test for main effects is statistically significant.
5. A researcher conducts a 2 (factor A) \times 3 (factor B) MANOVA with four dependent variables. Level of significance is .05, and the p value for Roy's largest root is .03 for factor A, .14 for factor B, and .58 for A \times B. Which one of the following is suggested?
 - a. Nonstatistically significant main effect for factor A. Statistically significant main effect for factor B and interaction effect.
 - b. A statistically significant main effect for factor A. Nonstatistically significant main effects for factor B and interaction effect.
 - c. The assumption of equal variance-covariance has been met.
 - d. The assumption of equal variance-covariance has not been met.
 - e. There is some evidence to suggest the assumption of multivariate normality has been violated.
 - f. There is some evidence to suggest the assumption of multivariate normality has been met.
6. Which test is commonly used to assess the assumption of linearity?
 - a. Box's M
 - b. Bartlett's
 - c. Shapiro
 - d. Wilks
7. Which test is commonly used to assess the omnibus MANOVA?
 - a. Box's M
 - b. Bartlett's
 - c. Shapiro
 - d. Wilks
8. The recommended follow-up to a statistically significant omnibus MANOVA is which one of the following?
 - a. Discriminant analysis
 - b. Power analysis
 - c. Profile plots
 - d. Univariate ANOVA

9. Partial eta squared for a statistically significant main effect is .20. Which of the following is a correct interpretation of this value?
 - a. The assumption of equal variance-covariance has been met.
 - b. The assumption of multivariate normality has been met.
 - c. The results are not statistically significant.
 - d. There is a large effect.
 - e. There is a small effect.
10. A researcher reports the following: Multivariate skewness (Small's $\chi^2 = 286.34$, $p = .27$), multivariate kurtosis (Small's variant $\chi^2 = 572.65$, $p = .34$), omnibus test (Small's test variant $\chi^2 = 386.96$, $p = .73$). Which one of the following is suggested?
 - a. The assumption of equal variance-covariance has been met.
 - b. The assumption of multivariate normality has been met.
 - c. The results are not statistically significant.
 - d. There is a large effect.
 - e. There is a small effect.

Conceptual Problems: Repeated Measures MANOVA

1. Malani collects data on 200 couples before they enter marriage counseling and again after they have completed marriage counseling. She collects the following data: (a) type of intervention (three categories: couple; small group; large group), (b) Marital Satisfaction Questionnaire (interval measurement scale measured pre and post), (c) Couples' Happiness Scale (interval measurement scale measured pre and post), and (d) Marriage Inventory (interval measurement scale measured pre and post). Is computing a repeated measures MANOVA appropriate given this data scenario?
 - a. Yes
 - b. No
2. Which one of the following data scenarios illustrates a good candidate for repeated measures MANOVA?
 - a. 40 meter sprint speed is measured before and after a regimented training program, hypothesizing differences by athletic level
 - b. Children are measured in the fall and spring on reading and science
 - c. Ratings of disruptive behavior and social competence are collected on children with autism spectrum disorder
 - d. The number of days absent are collected on employees before and after a policy change on personal leave
3. Which one of the following data scenarios illustrates, at least initially, a good candidate for repeated measures MANOVA?
 - a. Continuous measures, including caregiver quality of life, burden, and social support were collected prior to study implementation. Client symptoms and

- institutionalization were collected at the conclusion of the study on families with a relative with dementia.
- b. Internalizing and externalizing behavior (both continuous measures) collected at the beginning of the academic year on boys and girls by grade level
 - c. Data on group cohesion was collected on employees who participated in a face-to-face team building exercise as well as employees who completed an online module on team building
 - d. Risk factors and treadmill test measures (both continuous measures) collected on a group of elderly before and after engagement in a strenuous supervised exercise program
4. The assumption that is not required when pursuing a doubly multivariate repeated MANOVA is which one of the following?
- a. Independence
 - b. Multivariate normality
 - c. Homogeneity of variance/covariance
 - d. Multivariate sphericity
5. Which one of the following plots are helpful as a visual inspection of homogeneity of variance-covariance?
- a. Boxplots
 - b. Histograms
 - c. Profile plots
 - d. Spread-versus-level plots
6. Univariate normality is a necessary but insufficient condition for multivariate normality.
- a. True
 - b. False
7. In generating repeated measures MANOVA, you find a p value of .16 for Box's test. What is implied by this finding?
- a. The assumption of equal variance-covariance has been met.
 - b. The assumption of linearity can be assumed.
 - c. The correlations among the dependent variables are not zero.
8. In univariate data screening, Wesley finds that the standardized residuals for the two repeated continuous outcome variables for his repeated measures MANOVA have skewness values ranging from .06 to .85 and kurtosis values ranging from -.63 to +2.10. He conducts the Shapiro-Wilk test and finds p values ranging from .22 to .81. Wesley runs DeCarlo's (1997) macro and finds nonstatistically significant multivariate skew and kurtosis and omnibus multivariate tests for standardized residuals ($p > .47$). Which of the following is a reasonable assessment based on this information?
- a. There is some evidence to suggest the assumption of homogeneity of variance-covariance is met.
 - b. There is some evidence to suggest the assumption of homogeneity of variance-covariance has been violated.

- c. There is some evidence to suggest the assumption of linearity is met.
 - d. There is some evidence to suggest the assumption of linearity has been violated.
 - e. There is some evidence to suggest the assumption of multivariate normality is met.
 - f. There is some evidence to suggest the assumption of multivariate normality has been violated.
9. Which test is commonly used to assess the omnibus repeated measures MANOVA?
- a. Box's M
 - b. Bartlett's
 - c. Shapiro
 - d. Wilks
10. A friend comes to you seeking advice on conducting their repeated measures MANOVA. They are particularly interested in achieving maximum power. They have data on 40 cases including two within-subjects factors and one between-subjects factor, with equal cases in each of the two categories of the independent variable. The assumptions of the test are met. Which one of the following multivariate tests do you suggest given this scenario?
- a. Hotelling's trace
 - b. Pillai's trace
 - c. Roy's largest root
 - d. Wilks's lambda
 - e. Any. All have similar power in this scenario.

Computational Problems: One-Way and *k*-Way MANOVA Models

1. Use SPSS to conduct MANOVA (ECLSK_MANOVA_N1393.sav). The dependent variables are externalizing problems [C5SDQEXR] and internalizing problems [C5SDQINR]. The independent variables are male [MALE] and two-parent family structure [TWOPARENT_GR3]. Report the results for testing the assumption of homogeneity of variance-covariance matrices as well as the omnibus multivariate tests for the main effects and interaction.
2. Use SPSS to conduct MANOVA (ECLSK_MANOVA_N1393.sav). The dependent variables are perceived interest/competence in reading [C5SDQRDS] and perceived interest/competence in math [C5SDQMTR]. The independent variables are male [MALE] and two-parent family structure [TWOPARENT_GR3]. Report the results for testing the assumption of homogeneity of variance-covariance matrices as well as the omnibus multivariate tests for the main effects and interaction.

Computational Problems: Repeated Measures MANOVA

1. Use SPSS to conduct repeated measures MANOVA with only within-factors (ECLSK_REPEATEDMANOVA_N1344.sav). There are two repeated dependent variables, *reading* and *general knowledge*. More specifically, reading IRT scale

score in fall kindergarten [C1R4RSCL] and reading IRT scale score in spring kindergarten [C2R4RSCL] and general knowledge IRT scale score in fall kindergarten [C1RGSCAL] and general knowledge IRT scale score in spring kindergarten [C2RGSCAL]. Report the results for the omnibus multivariate tests for the main effect of time.

2. Use SPSS to conduct repeated measures MANOVA with only within-factors (ECLSK_REPEATEDMANOVA_N1344.sav). There are two repeated dependent variables, *mathematics* and *general knowledge*. More specifically, math IRT scale score in fall kindergarten [C1R4MSCL] and math IRT scale score in spring kindergarten [C2R4MSCL] and general knowledge IRT scale score in fall kindergarten [C1RGSCAL] and general knowledge IRT scale score in spring kindergarten [C2RGSCAL]. Report the results for the omnibus multivariate tests for the main effect of time.

Interpretive Problem: One-Way and *k*-Way MANOVA Models

1. Use SPSS to develop a MANOVA model with other variables in the ECLS-K dataset (ECLSK_MANOVA_N1393.sav). Write up the results in APA style, including testing for the assumptions. Report effect size values and post hoc power for the omnibus test.

Interpretive Problem: Repeated Measures MANOVA

1. Use SPSS to develop a repeated MANOVA model with other variables in the ECLS-K dataset (ECLSK_REPEATEDMANOVA_N1344.sav). Write up the results in APA style, including testing for the assumptions. Report effect size values and post hoc power for the omnibus test.

REFERENCES

- Barcikowski, R. S. (1981). Statistical power with group mean as the unit of analysis. *Journal of Educational Statistics*, 6(3), 267–285.
- Boik, R. J. (1988). The mixed model for multivariate repeated measures: Validity conditions and an approximate test. *Psychometrika*, 53(4), 469–486.
- Bray, J. H., & Maxwell, S. E. (1985). *Multivariate analysis of variance*. Newbury Park, CA: SAGE.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.
- D'Agostino, R. B. (1986). Tests for the normal distribution. In R. B. D'Agostino & M. A. Stephens (Eds.), *Goodness-of-fit techniques* (pp. 367–419). New York, NY: Marcel Dekker.
- DeCarlo, L. T. (1997). On the meaning and use of kurtosis. *Psychological Methods*, 2, 292–307.
- Hahs-Vaughn, D. L. (2005). A primer for using and understanding weights with national datasets. *Journal of Experimental Education*, 73(3), 221–248.
- Hahs-Vaughn, D. L. (2006a). Analysis of data from complex samples. *International Journal of Research & Method in Education*, 29(2), 163–181.

- Hahs-Vaughn, D. L. (2006b). Weighting omissions and best practices when using large-scale data in educational research. *Association for Institutional Research Professional File*, 101, 1–9.
- Hahs-Vaughn, D. L., McWayne, C. M., Bulotskey-Shearer, R. J., Wen, X., & Faria, A. (2011a). Complex sample data recommendations and troubleshooting. *Evaluation Review*, 35(3), 304–313. doi: 10.1177/0193841X11412070
- Hahs-Vaughn, D. L., McWayne, C. M., Bulotskey-Shearer, R. J., Wen, X., & Faria, A. (2011b). Methodological considerations in using complex survey data: An applied example with the head start family and child experiences survey. *Evaluation Review*, 35(3), 269–303.
- Huberty, C. J., & Olejnik, S. (2006). *Applied MANOVA and discriminant analysis* (2nd ed.). Hoboken, NJ: Wiley.
- Kish, L., & Frankel, M. R. (1973, October 17). *Inference from complex samples*. Paper presented at the annual meeting of the Royal Statistical Society.
- Kish, L., & Frankel, M. R. (1974). Inference from complex samples. *Journal of the Royal Statistical Society, Series B*, 36, 1–37.
- Kline, R. B. (2004a). *Beyond significance testing: Reforming data analysis methods in behavioral research*. Washington, DC: American Psychological Association.
- Kline, R. B. (2004b). Supplemental chapter on multivariate effect size estimation. Retrieved January 20, 2015, from <http://www.apa.org/books/resources/kline>
- Korn, E. L., & Graubard, B. I. (1995). Examples of differing weighted and unweighted estimates from a sample survey. *American Statistician*, 49, 291–305.
- Lee, E. S., Forthofer, R. N., & Lorimor, R. J. (1989). *Analyzing complex survey data*. Newbury Park, CA: Sage.
- Lix, L. M., & Hinds, A. M. (2004). Multivariate contrasts for repeated measures designs under assumption violations. *Journal of Modern Applied Statistical Methods*, 3(2), 333–344.
- Looney, S. W. (1995). How to use tests for univariate normality to assess multivariate normality. *American Statistician*, 49, 64–70.
- Lumley, T. (2004). Analysis of complex survey samples. *Journal of Statistical Software*, 9(8), 1–19.
- Mardia, K. V. (1970). Measures of multivariate skewness and kurtosis with applications. *Biometrika*, 57, 519–530.
- Penny, K. I. (1996). Appropriate critical values when testing for a single multivariate outlier by using the Mahalanobis distance. *Applied Statistics*, 45, 73–81.
- Pfeffermann, D. (1993). The role of sampling weights when modeling survey data. *International Statistical Review*, 61(2), 317–337.
- Puma, M. J., Olsen, R. B., Bell, S. H., & Price, C. (2009). *What to do when data are missing in group randomized controlled trials (NCES 2009–0049)*. Washington, DC: National Center for Education Evaluation and Regional Assistance, Institute of Education Sciences, U.S. Department of Education.
- Scariano, S. M., & Davenport, J. M. (1987). The effects of violation of independence assumptions in the one-way ANOVA. *The American Statistician*, 41(2), 123–129.
- Seo, T., Kanda, T., & Fujikoshi, Y. (1995). The effects of nonnormality on tests for dimensionality in canonical correlation and MANOVA tests. *Journal of Multivariate Analysis*, 52, 325–337.
- Shapiro, S. S., & Wilk, M. B. (1965). An analysis of variance test for normality (complete samples). *Biometrika*, 52(3 and 4), 591–611.

- Skinner, C. J., Holt, D., & Smith, T. M. F. (Eds.). (1989). *Analysis of complex samples*. New York: Wiley.
- Small, N. J. H. (1980). Marginal skewness and kurtosis in testing multivariate normality. *Applied Statistics*, 29, 85–87.
- Timm, N. H. (1980). Multivariate analysis of variance of repeated measurements. In P. R. Krishnaiah (Ed.), *Handbook of statistics, analysis of variance* (Vol. 1, pp. 41–87). New York: North Holland.
- Tourangeau, K., Nord, C., Le, T., & Sorongon, A. (2009). *Early childhood longitudinal study, kindergarten class of 1998–99 (ECLS-K), combined user's manual for the ECLS-K eighth-grade and K-8 full sample data files and electronic codebooks*. Washington, DC: National Center for Education Statistics, Institute of Education Sciences, U.S. Department of Education.

Chapter 7

DISCRIMINANT ANALYSIS

CHAPTER OUTLINE

7.1	What Discriminant Analysis Is and How It Works	275
7.1.1	Characteristics	277
7.1.2	Sample Size	286
7.1.3	Power	287
7.1.4	Effect Size	287
7.1.5	Assumptions	289
7.2	Mathematical Introduction Snapshot	292
7.3	Computing Discriminant Analysis Using SPSS	293
7.3.1	Generating Kappa Statistic for Classification Accuracy	314
7.4	Data Screening	315
7.4.1	Independence	316
7.4.2	Linearity	316
7.4.3	Noncollinearity	320
7.4.4	Normality of Independent Variables	321
7.4.5	Homogeneity of Variance-Covariance Matrices	324
7.5	Power Using G*Power	324
7.5.1	Post Hoc Power for Discriminant Analysis Using G*Power	324
7.5.2	A Priori Power for Discriminant Analysis Using G*Power	327
7.6	Research Question Template and Example Write-Up	328

KEY CONCEPTS

1. Discrimination
2. Classification
3. Discriminant function
4. Cut score

In recent chapters we examined various regression models—these allow us to examine the relationship between one or more predictors when the outcome is continuous (ordinary least squares regression) or categorical (logistic regression). This was followed by a discussion of multivariate analysis of variance procedures that allow for multivariate examination of mean differences. In this chapter, we are introduced to discriminant analysis, a procedure developed by Fisher (1936) that provides for classification of groups when the outcome is categorical, the same goal as that of logistic regression but, as we'll see, goes a bit beyond what logistic regression allows us to do. This procedure can be helpful in situations where classification into groups for purposes of intervention, training, receipt of something (e.g., monetary loan), and similar scenarios is needed or wanted. Cluster analysis, discussed in the next chapter, is similar in that it allows cases to be classified based on a set of variables. However, with discriminant analysis, the researcher knows *a priori* the mutually exclusive groups within which he or she wants to classify (i.e., the categorical dependent variable). In comparison, with cluster analysis, there are no such *a priori* groups into which the independent variables will be classified.

As we'll see in this chapter, discriminant analysis and analysis of variance share computational similarities. Likewise, logistic regression and discriminant analysis also share similarities, and there can be confusion in determining when one is more appropriate than the other. We'll see that discriminant analysis has stricter assumptions than logistic regression. The assumptions of multivariate normality and equal variance-covariance matrices, which are required in discriminant analysis, do not hold for logistic regression. Thus, logistic regression is more robust than discriminant analysis when these assumptions are not met. On the flip side, discriminant analysis is more powerful than logistic regression when these assumptions *do* hold. While discriminant analysis may be more interpretatively challenging as compared to logistic regression because it is not within the regression family, it is a very important statistical tool for discriminating between groups. Determining when discriminant analysis rather than logistic regression should be used can be confusing, and Box 7.1 may be helpful in understanding how the procedures compare.

For the purposes of this chapter, we will examine two-group discriminant analysis (similar to binary logistic regression where the dependent variable or grouping variable has only two categories or levels) as well as the three-group discriminant analysis (i.e., the dependent variable has three categories), although discriminant analysis can be applied in situations with more than three groups. Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying discriminant analysis and the discriminant function, (b) determine and interpret the results

BOX 7.1 DISCRIMINANT ANALYSIS COMPARISON TO LOGISTIC REGRESSION

Element	Discriminant Analysis	Logistic Regression
Scale of Variable	<ul style="list-style-type: none">• <i>Dependent variable</i> = one categorical with mutually exclusive and exhaustive groups known a priori• <i>Independent variable(s)</i> = one or more continuous <p><i>Consideration for decision:</i> If the model includes one or more categorical independent variables, logistic regression should be used.</p>	<ul style="list-style-type: none">• <i>Dependent variable</i> = one categorical with mutually exclusive and exhaustive groups known a priori• <i>Independent variable(s)</i> = one or more continuous and/or categorical
Assumptions	<ul style="list-style-type: none">• Independence• Linearity• Noncollinearity• Multivariate normality• Homogeneity of variance-covariance matrices <p><i>Consideration for decision:</i> If the assumptions of multivariate normality and/or homogeneity of variance-covariance matrices do not hold, logistic regression should be used.</p>	<ul style="list-style-type: none">• Independence• Linearity• Noncollinearity

of discriminant analysis, and (c) understand and evaluate the assumptions of discriminant analysis.

7.1 WHAT DISCRIMINANT ANALYSIS IS AND HOW IT WORKS

You may recall that Oso recently spearheaded the effort to analyze data using logistic regression. We revisit the stats lab today and find Oso excited to take on a similar analytic challenge.

Oso accomplished his previous statistical endeavor with logistic regression with such fortitude that the faculty advisor for the stats lab has asked him to consult with a faculty member in psychology on another challenging but exciting project. Dr. Isealine is a faculty member in biological sciences and has collected data on middle school students before and after their participation in a university-sponsored science summer camp. Dr. Isealine wants to know if science content knowledge, self-efficacy, and self-regulation prior to attending the science summer camp can classify students who identify as either wanting or not wanting to be a scientist when they enter the camp (i.e., at baseline). After looking at the data and speaking with Dr. Isealine, Oso has suggested that they explore group membership using discriminant analysis given the interest in classification of persons. Oso suggests the following research question to Dr. Isealine: *To what extent can identification of wanting or not wanting to be a scientist be identified from science content knowledge, self-efficacy, and*

self-regulation? With a dichotomous outcome and the goal of identifying or predicting group membership, Oso will use discriminant analysis as the statistical procedure to answer the question. Oso then proceeds with assisting Dr. Isealine in analyzing the data.

We have already learned about one statistical procedure—logistic regression—that can be used when the dependent variable is categorical (e.g., binary, dichotomous, or having only two categories). Logistic regression falls within the family of regression models, allowing a prediction to be made (specifically predicting group membership in the categorical outcome). Discriminant analysis shares similarity to logistic regression in that it allows researchers to use independent variables (aka predictor variables or ‘discriminating variables’) to identify group membership in a categorical outcome (which we will refer to again as the dependent variable), a process that is termed *classification*. (Discriminant analysis has a second purpose of *discrimination*, determining patterns of predictors in forming groups, that will be discussed and that goes beyond what logistic regression can do.) It is important to note that the categories of the dependent variable must be defined *a priori*. If your goal is to use a set of independent variables to create subgroups that are not already predefined, then you likely need a technique such as cluster analysis, which is covered in a different chapter.

Discriminant analysis differs from logistic regression in that computationally it is comparable to MANOVA rather than regression. Discriminant analysis is MANOVA backwards. MANOVA answers the question of whether mean differences on a combination of continuous dependent variables can be determined based on group membership. Discriminant analysis answers the question of whether group membership can be determined from a combination of continuous predictor variables. If we are working with human subjects, for example, discriminant analysis allows prediction or identification of individuals within groups based on a set of ‘discriminating variables,’ which we will term independent variables for simplicity and ease in understanding their purpose. The emphasis in discriminant analysis is on interpreting the pattern of differences based on the predictors as a whole, so that the differences in dimensions along which the groups differ can be better understood.

Group membership in discriminant analysis can be binary (i.e., dichotomous or two categories, termed two-group discriminant analysis) or multi-category (i.e., multinomial or more than two categories, termed multiple discriminant analysis), and the independent variables should be at least interval in scale as means and variances will be computed from them. The data structure for a discriminant analysis that examined a binary outcome (want vs. do not want to be a scientist) and three continuous predictors (science content knowledge, self-efficacy, and self-regulation) will look similar to that for the data structure of a binary logistic regression model (see example presented in chapter 5). The exception in these illustrations is that all predictor variables are now continuous in scale. The first 20 cases of the data file are presented in Table 7.1 (please note that for brevity the entire data file is not printed but can be accessed from the online materials).

TABLE 7.1

Identification as Scientist Example Data

Child	Science Content Knowledge (X_1)	Self-Efficacy (X_2)	Self-Regulation (X_3)	"I want to be a scientist" (Y)
1	1	3.44	3.71	No
2	0	3.89	4.14	No
3	2	4.56	4.43	Yes
4	1	3.78	3.86	Yes
5	0	4.67	5.00	No
6	0	3.33	3.71	No
7	3	3.67	3.71	Yes
8	0	4.89	5.00	Yes
9	2	4.33	5.00	No
10	3	4.00	3.43	No
11	5	4.67	4.00	Yes
12	2	4.00	4.29	Yes
13	0	4.67	4.57	No
14	0	3.78	3.86	Yes
15	0	1.11	1.43	Yes
16	0	5.00	5.00	No
17	1	4.22	3.74	Yes
18	0	3.89	3.43	Yes
19	2	4.86	4.81	No
20	1	4.20	4.32	Yes
...

7.1.1 Characteristics

7.1.1.1 Discriminant Function

Discriminant analysis is a tool that can be used for two primary and quite distinct purposes: (1) to determine or describe differences between or among the groups of the dependent variable by examining patterns from the predictors, termed *discrimination*; and (2) to determine or predict group membership, termed *classification*. The first purpose has also been referred to as descriptive discriminant analysis and the latter, predictive discriminant analysis (Huberty & Hussein, 2003).

Discrimination

The first goal, discrimination, is determining the dimensions along which groups differ. In other words, determining the features or attributes that best separate cases into the categories of the dependent variable. This can be likened to how we think of latent variables. This goal of discrimination is completed by examining patterns in

the canonical discriminant function, referred to from this point forward simply as the ‘discriminant function,’ to see if there is differentiation of units in the groups of our dependent variable. This goal is computationally reflective of MANOVA. The calculation is presented later in the mathematical introduction snapshot in Section 7.2. As we’ll see, the discriminant function is visually and functionally similar to a regression equation. The discriminant function score is predicted from the set of discriminating or independent variables, which are each weighted by a coefficient and then summed with the constant. The discriminant scores are then coordinates that are defined by the discriminant function. The unstandardized coefficient represents the amount of change on the discriminant function if the score on the respective variable changed by one unit. As we learned with regression, the unstandardized coefficients offer information on absolute contribution of a variable. Should relative importance of a variable in the discriminant function be of value, examination of standardized coefficients (forthcoming in the discussion) is needed.

The coefficients for the first function are derived such that there are maximum group differences on the function. When you are trying to classify into two groups (i.e., a binary dependent variable), only one discriminant function will be computed. If your dependent variable has more than two categories, then the number of discriminant functions that can be computed is equal to the *smaller* of (a) one less than the number of groups in the dependent variable ($k - 1$) or (b) the number of independent variables. In the case of three or more groups, the first function that is computed creates maximum separation. The first example that is illustrated has two groups in the dependent variable and three predictors. Thus, the number of discriminant functions equals one as $k - 1 = 2 - 1 = 1$, and this is less than the number of independent variables (which is $p = 3$). As another example, if a researcher has a model with five groups and three independent variables, there would be three discriminant functions calculated as the number of independent variables ($p = 3$) is less than $k - 1$ (which is $5 - 1 = 4$).

The second discriminant function again derives maximal difference between group means; in other words, it creates as much additional separation as possible, under the condition that the second function is uncorrelated to the first function. This process continues until the maximum number of discriminant functions, each uncorrelated with the previous, is reached. Again, the maximum number of discriminant functions that can be computed equals the fewer of either the number of independent variables or the number of groups of the dependent variable minus one. The most typical position in functional space is the location of the group centroid, which will be explained in an upcoming section.

Standardized Coefficients

As stated previously, standardized coefficients can be helpful when we want to know the relative importance of the predictor variables in the discriminant model. In other words, standardized coefficients can be reviewed to answer the question:

Which independent variables are the strongest contributors to the discriminant function scores? Standardized coefficients are coefficients that would have resulted if all the independent variables in our original data had standard deviations of one. The standardized coefficients should be reviewed in absolute value terms. The larger the standardized coefficient value (regardless of sign), the greater the contribution of that variable to the discriminant function score. It is important to note that standardized coefficients simultaneously consider the contribution of all independent variables. Because of this, a limitation of the standardized coefficient arises when two independent variables are strongly correlated. In those cases, the variables share in the contribution to the discriminant score and the standardized coefficients may be smaller than if only one of the highly correlated variables is used in the model. Standardized coefficients may also be larger (assuming a positive relationship between variables) than if only one of the highly correlated variables is used in the model but with opposite signs (and likewise smaller but with opposite signs when there is a negative relationship) and thus the contributions cancel each other out. Should the data be properly screened for collinearity, this should not be an issue.

Classification

The second goal of discriminant analysis, *classification*, predicting or identifying group membership from a set of predictors, is usually the primary question of interest and is akin to examining main effects in MANOVA. With this goal, we examine the classification rule or model that best separates cases into the categories of the dependent variable (i.e., ‘groups’), the goal of which is to maximize separation of the groups. Classification is performed by comparing the position or distance of a case to the group centroids to determine the closest centroid, and there are a number of different ways that classification can be achieved.

Classification based on a linear combination of predictor variables was first proposed by Fisher (1936), and classification function coefficients that are based on Fisher’s linear discriminant function can be requested when generating discriminant analysis. Fisher’s classification is based on a linear combination of predictor variables that achieves both of the following: (a) maximizes the differences *between* the groups of the dependent variable and (b) minimizes the variation *within* the groups. The classification function coefficients themselves are not usually interpreted, given that they are not standardized and each group has a different function, making interpretation difficult.

Mahalanobis distance is another statistic that can be used as a measure of distance for classification purposes. In general terms, Mahalanobis distance is a measure of spatial distance between two or more correlated variables. Mahalanobis distance, when squared, provides a squared measure of generalized distance from each case to each group centroid with cases classified into the group with the smallest squared distance value. This measure of distance assumes equal covariance matrices. In cases where

this assumption is violated, a modification can be made; however, the modification will not be covered in this text.

Classification can also be performed using the canonical discriminant functions rather than the predictor variables (as proposed by Fisher). Although standard statistical software alleviates this issue, classification in this fashion is less mathematically intense than what was proposed by Fisher, particularly as the size of the sample increases, and the results are generally the same if the assumption of homogeneity of variance-covariance is met. Classification results are still relatively robust to even moderate violations of homogeneity of variance-covariance but not so with substantial violations (Klecka, 1980).

Classification Matrix

The classification matrix is a cross-tabulation of the distribution of observed group membership to predicted group membership. The percentage of cases correctly classified provides another, albeit more indirect, measure of accuracy of the predictor variables in discriminating to groups. By chance, cases would be predicted to the correct group 50% of the time. Forthcoming in the chapter are a number of statistics that can be used to determine classification accuracy.

7.1.1.2 Interpreting the Discriminant Functions

The discriminant function(s) can be interpreted by studying (a) the relative position of the case to the group centroid and (b) the relationship between the independent variables and the discriminant functions. When more than one function is computed, the researcher must also determine how many functions are required. Recall that the maximum number of discriminant functions equals the lesser of the number of groups of the dependent variable *or* the number of independent variables. If the number of independent variables is less than the number of groups of the dependent variable, then the maximum number of discriminant functions equals the number of independent variables. The discriminant scores, which are computed in standard deviation units, are coordinates in our p -dimensional discriminant function space, and thus plots can be informative.

With two groups in the dependent variable, one discriminant function is produced. When there are more than two groups, the researcher must determine how many discriminant functions are statistically significant. There are a number of indices that are computed in discriminant analysis that can be useful in interpreting the number of important discriminant functions. These include (a) eigenvalues, (b) canonical correlations, and (c) Wilks's lambda.

Eigenvalues

An eigenvalue is the variance of the linear combination of variables that results from the decomposition of the correlation matrix. The largest variance (i.e., eigenvalue) is

associated with the first combination. The next largest variance, which is uncorrelated with the first, is second. This pattern continues for all combinations. Eigenvalues are a type of squared correlation (specifically, squared canonical correlation coefficients), and thus can be interpreted as shared variance. In discriminant analysis, eigenvalues reflect the degree of power to discriminate, with larger eigenvalues suggesting greater discrimination. Discriminant analysis models that examine more than two groups in the dependent variable will produce multiple eigenvalues whose magnitudes can be compared. Unfortunately, there are no standard criteria for interpreting how large the percent of variance accounted for must be in order to be recognized as useful in interpreting the model.

Canonical Correlations

Canonical correlations measure the relationship between the groups of the dependent variable and the discriminant function. Canonical correlations are interpreted similarly to other types of correlations, with values closer to zero representing little to no relationship and values closer to an absolute value of one representing stronger relationships between groups and the discriminant function. The canonical correlation coefficient squared represents the proportion of variance in the discriminant function that is explained by the groups of the dependent variable. It is important to note that the first canonical correlation will not always have the largest value. Because the canonical correlation measures the relationship between the groups and the discriminant function, if there is a weak relationship between the groups and first discriminant function, the first canonical correlation coefficient will be weak, even though the first discriminant function, by default, represents maximum group differences. When the groups of the dependent variable are similar in respect to the independent variables, the canonical correlations will be weak.

Wilks's Lambda

Wilks's lambda, Λ , is the statistic that measures discrimination via multivariate group differences. Lambda values can range between 0 and +1.0. Smaller values of lambda indicate that the independent variables are more effective in discriminating cases into the groups of the dependent variable. The smaller and closer to zero the lambda value is, the more separation there is *between* groups relative to within groups. The larger and closer to one the lambda value is, the less separation or discrimination there is between groups. Lambda values of +1.0 are found when group centroids are the same, reflecting absolutely no differences between groups. Wilks's lambda can be interpreted as a test of inference using the chi-square distribution. Statistically significant Wilks's lambda values indicate the discriminant function is statistically significant, and this can be interpreted as evidence of discrimination between groups.

When there are only two groups in the dependent variable, only one Wilks's lambda value will be computed. More than two groups will produce multiple Wilks's lambda statistics. In those cases, the researcher should examine the first Wilks's lambda to

determine if the first discriminant function is statistically significant. If the first Wilks's lambda is statistically significant, review of the next Wilks's lambda is warranted. If the second lambda is statistically significant, this suggests that a second dimension would add to our understanding of the group differences. This review of lambdas continues until all are reviewed and found statistically significant or not statistically significant. Nonstatistically significant results on the second lambda provide evidence to suggest that one dimension is sufficiently representing all the observed between group differences. When there are multiple lambdas that are statistically significant, this provides evidence that those statistically significant discriminant functions are working together as a set (and they will be used as a set rather than individually). In some cases, a researcher may find that none of the lambdas are statistically significant, and this indicates that the groups of the dependent variable are not different.

After determining which, if any, lambdas are statistically significant, the researcher would proceed to examine the structure coefficients and group centroids for those discriminant functions to better interpret them.

Structure Coefficients

Another type of coefficient produced when computing discriminant analysis is a structure coefficient. The structure coefficients are pooled within group correlations. These structure coefficients represent correlations between a single independent variable and the discriminant function within the groups of the dependent variable. Just as with other bivariate correlations, the closer the structure coefficient value is to an absolute value of one, the stronger the relationship between the respective independent variable and the discriminant function. As the structure coefficient approaches a value of 0 (again in absolute value terms), this suggests very little relationship between the respective independent variable and the discriminant function within group.

Centroids

Centroids were mentioned previously in our discussion of discriminant function space. As a reminder, a group centroid is the most typical position in discriminant function space. In studying correlations, we became acquainted with scatterplots, graphs that represent coordinates in two-dimensional space. With discriminant analysis, we are studying p -dimensional space with the independent variables as axes defining this space, where p represents the number of discriminating or independent variables and the size of the dimensional space equaling, at most, one less than the number of groups of the dependent variable. Each measure in our dataset (i.e., each value or score) represents one point in this p -dimensional space, and the coordinates are the unit's score on each of the variables. Groups that differ will be reflective of points that clump together in the p -dimensional space. Just as we use a mean to summarize a 'typical' score or value for a variable, a *centroid* is used with discriminant analysis to summarize group position. Thus, centroids can be examined to better understand how groups differ. The

point in this p -dimensional space where each axis has a value of zero, where the mean on each axis for all cases in aggregate is found, and where the origin is located is termed the *grand centroid*.

In the case where a centroid does not define a new dimension, the number of discriminant functions needed is less than the number that is mathematically calculated (i.e., the groups of the dependent variable can be identified by fewer functions than are computed), the researcher can examine statistical tests to determine the number of dimensions needed.

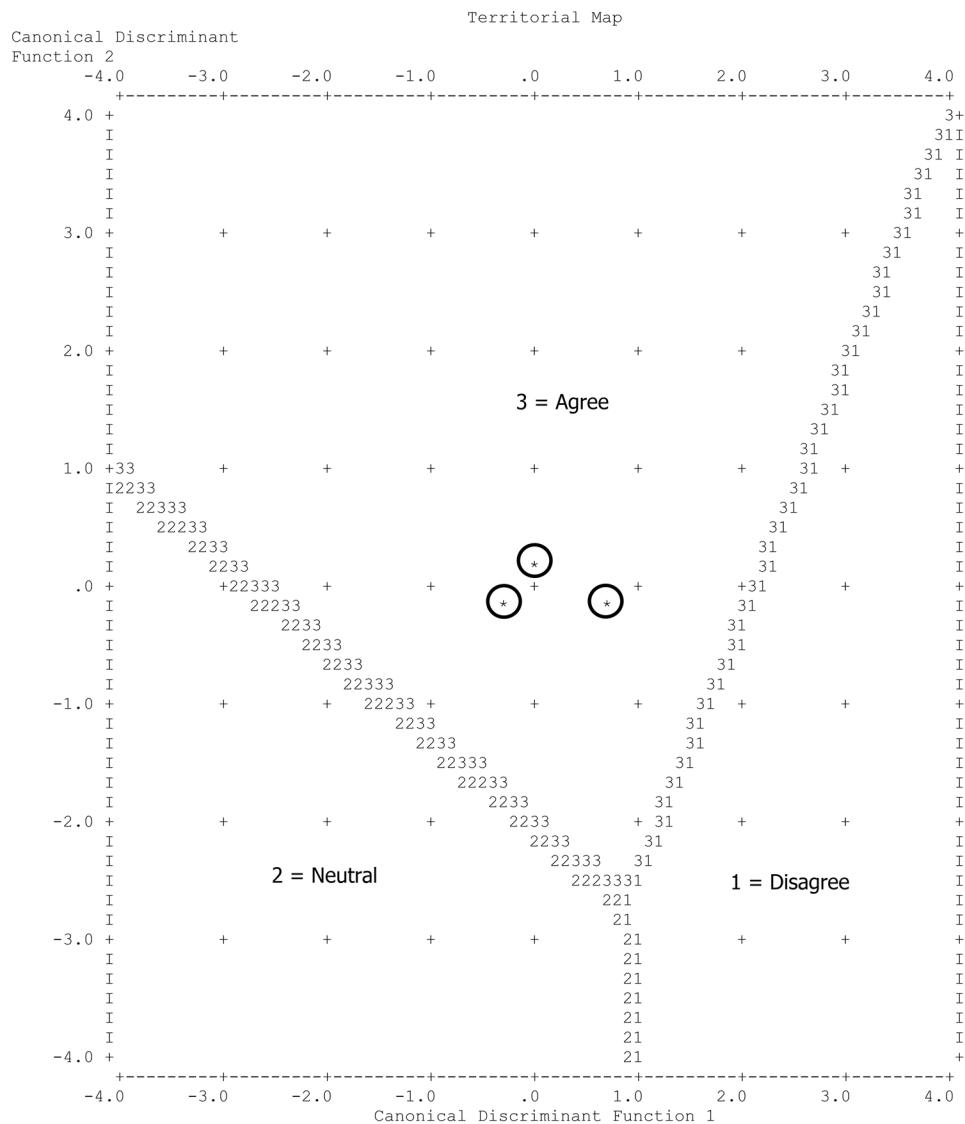
7.1.1.3 Discriminant Function Plots

Discriminant function plots and graphs can be helpful for visually interpreting the degree of discrimination. Centroids and coordinates of points for cases can easily be plotted when there are only two discriminant functions, as this is a two-dimensional space and visually quite easy to interpret. As the number of discriminant functions increase, the p -dimensional discriminant function space also increases, and visual interpretations of the entire space can become unwieldy. However, researchers may want to consider viewing plots of the first two discriminant functions considering that they are most important.

Figure 7.1 is a *territorial map*. In this illustration, the territorial map is a scatterplot or map of how groups differ on the two discriminant functions (function 1 on the horizontal and function 2 on the vertical) resulting from a three-group discriminant analysis (the results of which are presented in Table 7.5). As the number of groups in the dependent variable increase, the number of distinct regions in the territorial map will thus also increase.

The boundaries of the planes within the space of the territorial plot are created by the pairs of scores with territories representing those regions that have the strongest association with each group. In other words, the territories define the areas where units of a group are most likely to be located. In Figure 7.1, for example, students classified as group 3 ('agree' to the statement, "I want to be a scientist") have a range of negative to positive scores on discriminant function 1 (horizontal axis), but have only positive scores on discriminant function 2 (vertical axis). Thus, we see the top 'V' portion of the territorial map is the region where students classified as group 3 will be located. In comparison, students classified as group 1 ('disagree' to the statement, "I want to be a scientist") have only positive scores on discriminant function 1 but a range of positive and negative scores on discriminant function 2.

The group centroid plots can also be helpful in visualizing classification of cases into groups. Separate plots, representing each group individually (an example is provided in the SPSS output interpretation), as well as combined plots, representing all groups of the dependent variable, can be created. In the combined plot (Figure 7.2), the location of each case in discriminant function space can be examined.



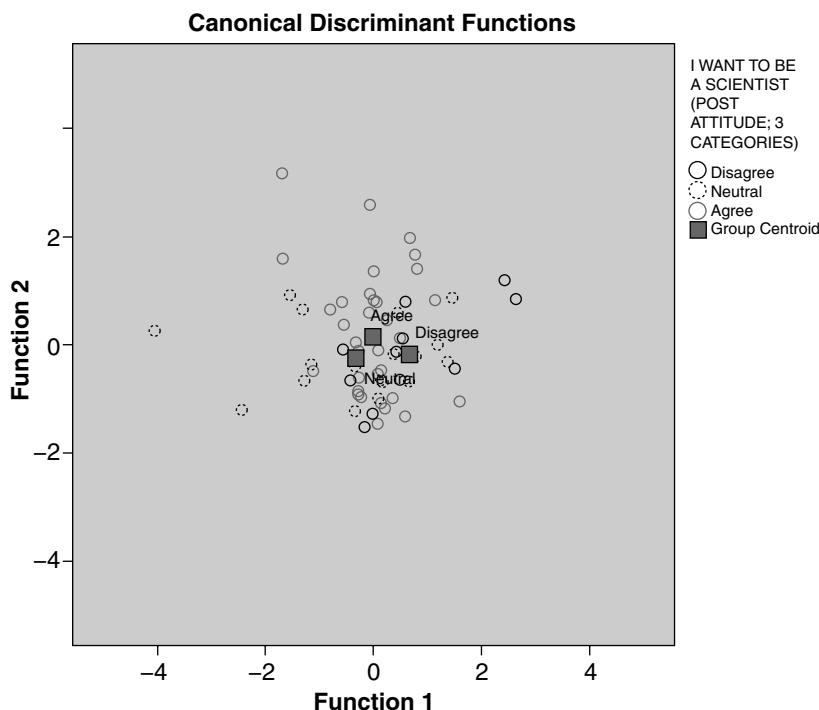
Symbols used in territorial map

Symbol	Group	Label
1	1	Disagree
2	2	Neutral
3	3	Agree

* Indicates a group centroid

FIGURE 7.1

Territorial Plot

**FIGURE 7.2**

Combined Canonical Discriminant Function Plot

7.1.1.4 Cut Score

Classifications in discriminant analysis are based on a cut score. In a nutshell, discriminant scores less than the cut score are classified into group '0,' and discriminant scores greater than the cut scores are classified into group '1.' When calculating the cut score in discriminant analysis, the researcher must determine how they want to treat the size of the samples of the groups of the dependent variable. Prior probabilities can be computed assuming equal sizes of groups or they can be computed proportional to the size of the group. Equal prior probabilities assume that, regardless of the observed sample sizes of the groups of the dependent variable, each group has an equal chance of occurring in the population. If group sizes are unbalanced and relative group size observed in the sample is assumed to be representative of what is seen in the population, then prior probabilities should be computed from the observed group sizes. If unsure whether to assume equal probabilities or other probabilities, the conservative approach is to select equal probabilities when defining the model.

The cut score is calculated as follows, and this equation is appropriate to use when *prior probabilities are calculated from group sizes*:

$$Z_{CS} = \frac{n_A Z_B + n_B Z_A}{n_A + n_B}$$

Where Z_{CS} equals the optimal cut score between groups 0 and 1

n_A equals the number of observations in group A

n_B equals the number of observations in group B

Z_A equals the centroid for group A

Z_B equals the centroid for group B

When prior probabilities are assumed to be equal for groups, the cut score is calculated as follows:

$$Z_{C_Equal} = \frac{Z_A + Z_B}{2}$$

Where Z_{C_Equal} equals the optimal cut score for equal group sizes

Z_A equals the centroid for group A

Z_B equals the centroid for group B

7.1.1.5 Cross-Validation

Cross-validation is a way that researchers can gauge the ability of their data to correctly predict cases. One way to cross-validate results is to have a hold-out sample, one subset being used to derive the functions and the other to test classification into groups. SPSS includes an option to generate leave-one-out cross-validation results. This is a jackknife cross-validation procedure where the discriminant function is computed using all cases but one (i.e., ‘leave-one-out’) and then the one excluded case is classified. This process is repeated for all cases, with each case having a turn at being excluded. Reasoning that a case that is being predicted should not be used to create the function, jackknife cross-validation produces more reliable results (but often lower rates of correct classification) in comparison to the original classification.

7.1.1.6 Putting the Pieces Together

Partly due to the fact that there can be two different objectives for generating discriminant analysis (i.e., discrimination and classification), there may appear to be more decisions to make when computing discriminant analysis as compared to other procedures. Table 7.2 may be helpful in understanding how the pieces fit together.

7.1.2 Sample Size

There are various recommendations for sample size requirements of discriminant analysis, both within group as well as overall. A guideline for *overall* sample size suggests that there should be a sample size of at least 20 for every one predictor, with a more liberal ratio of a sample of size 5 for every predictor.

One suggestion for the *within-group* sample size is that the number of units in the dependent variable’s smallest group should be larger than the number of independent

TABLE 7.2

How the Discriminant Analysis Pieces Fit Together

Goal	Statistics to Review
Primary: Discrimination (interpreting patterns of predictors to understand how the groups are separated)	<ul style="list-style-type: none"> • Discriminant function score • Unstandardized and standardized discriminant function coefficients • Eigenvalues • Canonical correlations • Wilks's lambda • Discriminant function plots
Primary: Classification (decision rule for classifying into groups where the number and meaning of the underlying dimensions are irrelevant)	<ul style="list-style-type: none"> • Fisher's linear discriminant function • Mahalanobis distance • Cut score • Classification cross-tabulation and classification accuracy • Acceptable classification indices <ul style="list-style-type: none"> ◦ Standards of comparison ◦ Press's Q ◦ Kappa
Effect: Overall discriminant function	<ul style="list-style-type: none"> • Eta squared or partial eta squared
Effect: Individual discriminant function	<ul style="list-style-type: none"> • Squared canonical correlation coefficients

variables. This is a very liberal standard. A more conservative standard for within-group sample size is that each group should have at least 20 units. When classifying unbalanced groups in discriminant analysis, groups that have the larger size will have an increased probability of classification. In extremely unbalanced cases and to increase the balance within groups, researchers may want to consider drawing a random sample of units from the larger size group.

7.1.3 Power

As noted earlier, discriminant analysis is more powerful than logistic regression when the assumptions of the test hold. G*Power will be applied later to test a priori and post hoc power.

7.1.4 Effect Size

Effect size in discriminant analysis can be computed for the *overall discriminant analysis* (via partial eta squared) and for the *individual discriminant functions* (via squared canonical correlation). Although not technically a measure of effect, we will also consider indices to determine acceptable classification.

7.1.4.1 Overall Discriminant Analysis Effect Size

The effect size for the overall discriminant analysis is the same as that determined for the overall MANOVA, that being eta squared or partial eta squared. Using Cohen's (1988) guidelines for interpretation, eta squared values of .01 are interpreted to be small, .06 are moderate, and .14 are large. If the software used to generate discriminant

analysis does not provide either of these values, they can easily be calculated from the Wilks's lambda provided when computing discriminant analysis:

$$\text{partial } \eta^2 = 1 - \Lambda^{1/3}$$

7.1.4.2 Individual Discriminant Function Effect Size

When squared, the canonical correlation for each discriminant function provides the proportion of shared variance between groups and predictors for that individual discriminant function. Cohen's guidelines for the interpretation of correlation coefficients can be applied to interpret canonical correlation coefficient values. More specifically, canonical correlation values of .10 are interpreted to be small, .30 are moderate, and .50 are large. Thus, squared canonical correlation values of .01 are small, .09 are moderate, and .25 are large. An additional measure of effect size is the squared pooled within-groups correlation between the predictor variable and the standardized canonical discriminant functions. These values can be found in the output labeled 'structure matrix.'

7.1.4.3 Acceptable Classification

After computing discriminant analysis, it is important to determine acceptable classification. This can be done by computing standards of comparison, Press's Q , and kappa.

Standards of Comparison

One way to determine the accuracy of the discriminant function in predicting is to compute standards of comparison, a chance criterion measure of overall acceptable classification. When sample sizes of the groups are equal, standards of comparison are computed as follows, where k equals the number of groups or categories of the dependent variable. Probability of chance classification for a two-group function, for example, is .50.

$$C_{\text{equal}} = \frac{1}{k}$$

When group sample sizes are unequal, the proportional chance criterion standards of comparison is computed as follows, where p equals the proportion of units in group 1, and $(1 - p)$ equals the proportion of units in group 2:

$$C_{\text{unequal}} = p^2 + (1 - p)^2$$

In the absence of absolute criterion for evaluating classification accuracy, a general recommendation is that classification accuracy that is at least $\frac{1}{4}$ greater than the chance criterion standard provides justification for using the results as a way to discriminate groups.

Press's *Q*

Press's *Q* is another index that can be used to gauge acceptable overall classification. Press's *Q* is a statistical test for discriminating power of the classification matrix when compared with a model of chance. As this is a chi-square test with one degree of freedom (i.e., $\chi^2 = 3.841$), this test is sensitive to sample size. Press's *Q* is calculated as follows:

$$Q = \frac{[N - (nk)^2]}{N(k - 1)}$$

Where *N* equals the total sample size, *n* equals the number of observations correctly classified, and *k* equals the number of groups of the dependent variable. Statistically significant Press's *Q* statistics provide evidence to suggest that the predictions are statistically significantly better than chance.

Kappa

The kappa statistic is another measure that can be used to assess classification accuracy. Kappa is the proportion of agreement after chance agreement has been removed, and values can range from -1.0 to +1.0. A kappa value of one indicates perfect agreement, while a kappa of zero indicates the agreement is the same as expected by chance. Thus, values closer to one suggest stronger agreement. Negative values, while rare, can occur and suggest that agreement is poorer than expected by chance. Cohen's guidelines for interpretation of correlation coefficients can be applied to interpret kappa values. More specifically, kappa values of .10 are interpreted to have weak prediction, .30 have moderate prediction, and .50 or larger suggest strong prediction.

The predicted group membership values, which can be saved when computing discriminant analysis, are used to compute the kappa. The kappa statistic is computed using the predicted group membership values along with the dependent variable for the discriminant analysis.

7.1.5 Assumptions

Many of the same assumptions that have been considered for previous statistical procedures are applicable for discriminant analysis as well: (a) independence, (b) linearity, (c) noncollinearity, (d) multivariate normality, and (e) homogeneity of variance-covariance matrices. Conditions for discriminant analysis are that, as stated previously: (a) there must be at least two a priori mutually exclusive groups for the dependent variable and (b) discriminating (i.e., independent or predictor) variables must be measured on at least an interval scale and the number of independent variables should be less than the total number of cases minus two.

7.1.5.1 Independence

We should conceptually be familiar with the assumption of independence from previous chapters (e.g., ANOVA). Independence in discriminant analysis is similar to what

has been seen with other statistical procedures—scores or values of the predictors are independently and randomly sampled from the population. Lack of independence affects the estimated standard errors of the model. Unfortunately, there are no satisfactory means by which to examine independence when generating discriminant analysis. If the data were collected using a simple random sample, the researcher could feel relatively confident that this assumption has been met. Data collected by nonsimple random sampling methods run the risk of increased homogeneity, which often translates to increased Type I errors (i.e., too often finding statistical significance when the null hypothesis is really true). Researchers whose data do not conform to a simple random sample may wish to caution readers to this risk.

7.1.5.2 Linearity

Linearity among all pairs of predictors within groups is required. Violation of the assumption of linearity reduces power but does not appreciably increase the chance of a Type I error. Transformations may be considered for nonlinear distributions. Linearity can be assessed by reviewing simple bivariate scatterplots of all pairs of predictors within groups or levels of the dependent variable.

7.1.5.3 Collinearity

Collinearity is a concept that was introduced with multiple regression. Noncollinearity occurs when the predictor variables are *not* highly correlated, and this is desired. Multicollinearity, on the other hand, occurs when multiple predictor variables are highly correlated. Singularity is the most extreme case of multicollinearity, occurring when there is a perfect relationship between predictors. When multicollinearity is present, inverting the discriminant matrix becomes unstable, and when singularity occurs, the discriminant matrix cannot be inverted. Tolerance tests, such as those computed with multiple regression, can be used to test for collinearity. Independent variables that do not meet tolerance should be excluded.

7.1.5.4 Multivariate Normality

Multivariate normality is met when there is normality of the linear combinations of the sampling distributions of the predictor means. A cursory assessment of multivariate normality can be done by examining univariate normality of the independent variables using statistics such as skewness, kurtosis, formal tests of normality, and graphs for visual inspection. Normality violations due to outliers are more problematic than violations due to skewness, for which discriminant analysis is fairly robust. Multivariate outliers can be detected by examining Mahalanobis distance, which is evaluated as a chi-square statistic with degrees of freedom equal to the number of independent variables in the discriminant model. Mahalanobis distance values that are statistically significant at $p < .001$ are considered multivariate outliers and should be considered for removal. In the event nonnormality due to situations other than outliers is detected, discriminant analysis results will be relatively robust if the samples sizes of the groups

of the dependent variable are equal and a nondirectional (i.e., two-tailed) test is used. In the case of unequal sample sizes of groups, the results are generally still robust when there are five or less independent variables and the sample size is at least 20 in the smallest group of the dependent variable. Transformations of predictor variables may be used to improve adherence to this assumption.

7.1.5.5 Homogeneity of Variance-Covariance Matrices

Homogeneity of variance-covariance matrices is the multivariate version of the univariate equal variances assumption. When homogeneity of variance-covariance matrices is met, the within-group variances are similar. When the goal of discriminant analysis is testing inferences, discriminant analysis is relatively robust to violations of homogeneity of variance-covariance with large sample sizes or when there is a balanced design (i.e., equal sample sizes in the groups of the dependent variable). In comparison, when the goal of discriminant analysis is classification, the results are not robust to violations as groups with greater dispersion (i.e., heterogeneity) will generally have cases that are overclassified. Additionally, the smaller the sample size, the less robust classification results become.

When there are more than two groups (and thus two or more discriminant functions), this assumption can be evaluated by examining scatterplots of the scores of the first two discriminant functions. Homogeneity is evident when the scatterplots produce scores that are fairly equal in size (i.e., there is similarity in the overall size of the scatterplots). Box's M test can also be used, although this test has been criticized for being overly sensitive with large samples (i.e., tends to be statistically significant even when the differences in the variance-covariance matrices are not that large when the sample size of the discriminant analysis model is large), insensitive with small samples (i.e., not powerful enough to detect heterogeneity in variance-covariance matrices when the sample size is small), and can also be quite sensitive where there is multivariate nonnormality. Thus, researchers who use Box's M as evidence for homogeneity may wish to use a more liberal alpha level, such as .001, for evaluating Box's M .

There are several options that can be pursued if this assumption is violated. Transformations of predictor variables may be used to improve adherence to this assumption. Separate covariance matrices can be used during classification (an option available when generating discriminant analysis) or nonparametric classification (e.g., logistic regression) can be used. It is important to note that if separate covariance matrices are used, the sample size must be sufficient to cross-validate results, as this option can lead to an overfit model.

7.1.5.6 Concluding Thoughts on Assumptions

There is no rule that research that violates assumptions must be scrapped. However, researchers who face violations of assumptions must handle these situations on a

TABLE 7.3

Assumptions and Violation of Assumptions: Discriminant Analysis

Assumption	Effect of Assumption Violation
Independence	<ul style="list-style-type: none"> Influences standard errors of the model and thus hypothesis tests and confidence intervals
Linearity	<ul style="list-style-type: none"> Reduces power but does not appreciably increase Type I error rates
Noncollinearity of X 's	<ul style="list-style-type: none"> Instability in discriminant matrix inversion
Multivariate normality	<ul style="list-style-type: none"> Reduces efficiency and accuracy in tests of inference Increases misclassifications
Homogeneity of variance-covariance matrices	<ul style="list-style-type: none"> Overclassifies to groups with greater dispersion

case-by-case basis, considering both the goal of the analyses and the extent to which the assumptions were violated and the resulting effect of violation. It is also important that researchers present the evidence found, along with justification for decisions that were made.

When examining assumptions related to discriminant analysis, if violations are found, the researcher must consider the primary goal of their analyses. If the primary goal is classification and the percentage of cases correctly classified is high, then violations of the assumptions are likely minimal, and efforts to decrease violations (e.g., data transformations) may result in only marginal improvements that are more trouble than they are worth. However, if the percentage of correctly classified cases is low, misclassifications may be due to either violations of assumptions or a weak model. In this case, the researcher may want to consider steps that will improve evidence of meeting assumptions. Violations of assumptions are also generally more problematic with smaller sample sizes. Thus, being more particular to satisfying assumptions with small n 's is warranted. The assumptions are summarized in Table 7.3.

7.2 MATHEMATICAL INTRODUCTION SNAPSHOT

The discriminant function is calculated as follows:

$$D_{ik} = d_{0k} + d_{1k}z_{i1} + d_{2k}z_{i2} + \dots + d_{jk}z_{ip}$$

Where D_{ik} = discriminant function score for the i th case on the k th discriminant function (where k is the minimum of either one less than the number of groups in the dependent variable or the number of predictors)

d_{0k} = constant

d_{jk} = value of the j th coefficient for the k th discriminant function

z_{ip} = value on discriminating variable z for the i th case of the p th predictor (1 through p , where p equals the number of independent variables)

The discriminant function likely looks familiar, as it is visually and functionally similar to a regression equation. Examining the equation, we see that the discriminant score is computed for each case by summing a constant to the products of the value for the case on each predictor variable and the coefficient for that variable. In other words, the discriminant function score is predicted from the set of discriminating or independent variables, which are each weighted by a coefficient and then summed with the constant. The discriminant scores are then coordinates that are defined by the discriminant function. The unstandardized coefficient represents the amount of change on the discriminant function if the score on the respective variable changed by one unit. As we learned with regression, the unstandardized coefficients offer information on absolute contribution of a variable. Should relative importance of a variable in the discriminant function be of value, examination of standardized coefficients (forthcoming in the discussion) is needed.

7.3 COMPUTING DISCRIMINANT ANALYSIS USING SPSS

Now we consider SPSS for computing discriminant analysis. Our dataset has 66 cases [Ch7_sciencecamp.sav]. The number of cases within each category of our binary dependent variable [PRE_WANT_SCI_BI] is approximately equal ($n = 31$, no; $n = 35$, yes), and surpasses the general criteria for at least 20 units in each group. With slight imbalance in the size of the groups, however, we recognize there may be an increased probability of classification into the larger group (i.e., ‘yes’). We are working with four columns of data, the first representing the binary dependent variable (‘yes’ or ‘no’ response to the statement, “I want to be a scientist,” which was measured prior to participating in the science summer camp) and the following three columns representing the continuous independent variables that were also measured prior to participating in the science summer camp (science content knowledge, self-efficacy, and self-regulation). (Note: If using the dataset from the website, you’ll notice that there are additional variables in the file. These are provided so that you have the opportunity to explore and practice additional models beyond what is illustrated in the textbook.) Again, our dependent variable is categorical, thus making discriminant analysis appropriate.

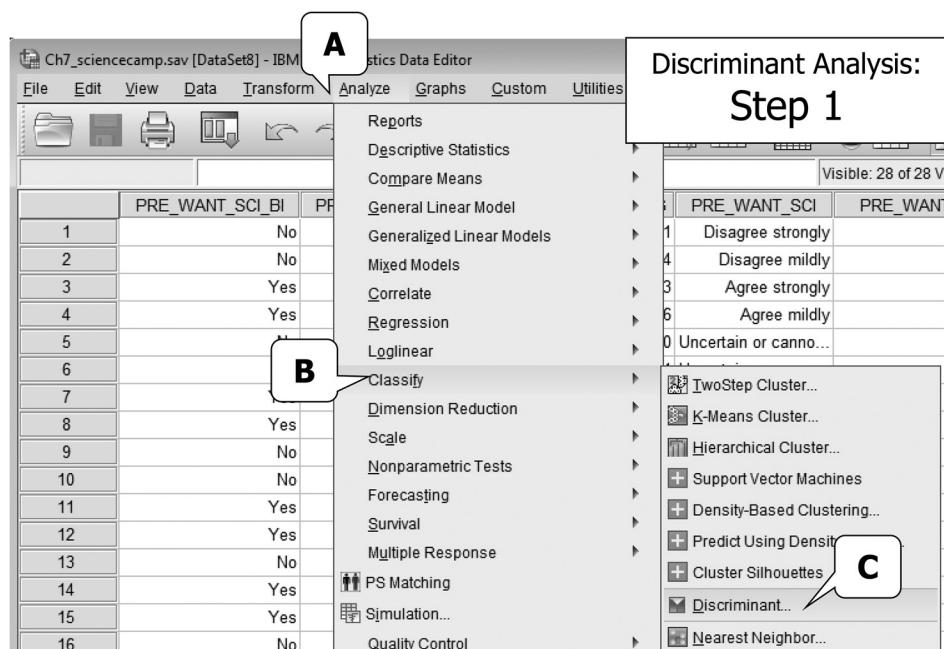
(The three-group discriminant analysis results, which use these same three predictor variables but the three-category dependent variable measured *after* participation in the summer camp [POST_WANT_SCI_3CAT], are also interpreted and presented in Table 7.4. The research question examined for the three-group analysis would therefore be slightly different than that posed for the two-group, given that the outcomes reflect different points in time—one prior and one after participation in the summer camp. However, for illustrative purposes, the screenshots reflect the use of the two-group dependent variable.)

The dependent variable is 'Pre_Want_Sci_Bi' and represents the student's answer ('yes' or 'no') to the statement, "I want to be a scientist." This is a binary variable where '1' represents 'Yes' and '0' represents 'No.'

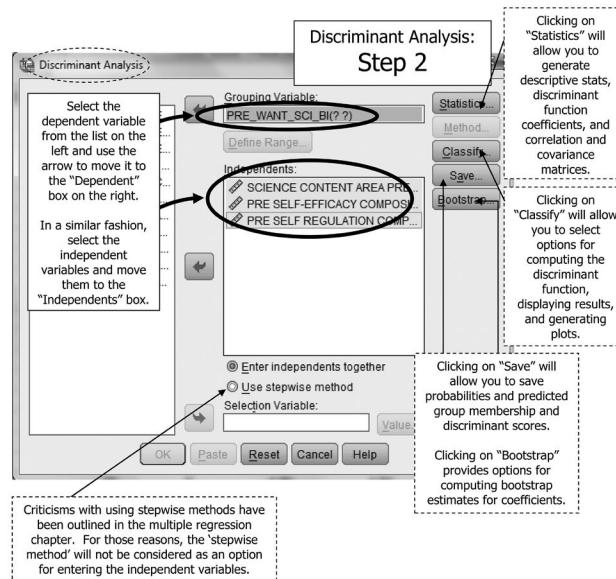
The independent variables are labeled 'Pretest,' 'Pre_Selfeff,' and 'Pre_Selfreg' where each value represents the child's score on a pretest measure including a science content knowledge test and self-report scores on self-efficacy and self-regulation.

	PRE_WANT_SCI_BI	PRETEST	PRE_SELFEFF	PRE_SELFREG
1	No	1	3.44	3.71
2	No	0	3.89	4.14
	Yes	2	4.56	4.43
	Yes	1	3.78	3.86
	No	0	4.67	5.00
	No	0	3.33	3.71
	Yes	3	3.67	3.71
	Yes	0	4.89	5.00
	No	2	4.33	5.00
	No	3	4.00	3.43
	Yes	5	4.67	4.00
	Yes	2	4.00	4.29
	No	0	4.67	4.57
	Yes	0	3.78	3.86
14				
15	Yes	0	1.11	1.43
16		No	0	5.00
17		Yes	1	4.22
18		Yes	0	3.89
19		No	2	4.86
20		Yes	1	4.20
				4.32

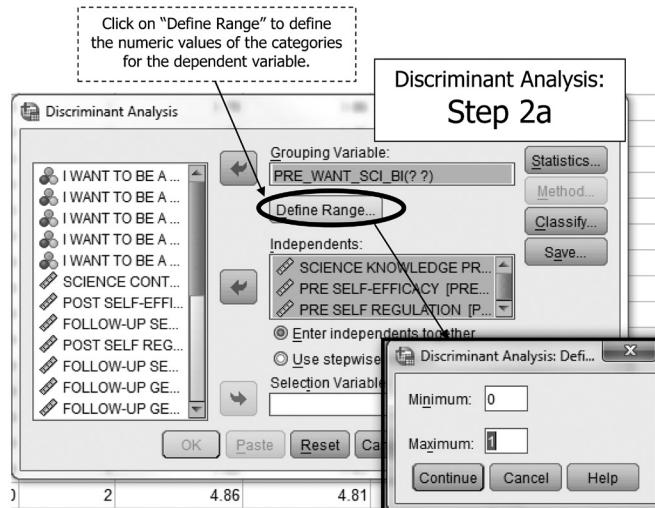
Step 1. To conduct discriminant analysis, go to "Analyze" in the top pull-down menu, then select "Classify," and then select "Discriminant." Following the screenshot below (see screenshot Step 1) produces the "Discriminant Analysis" dialog box.



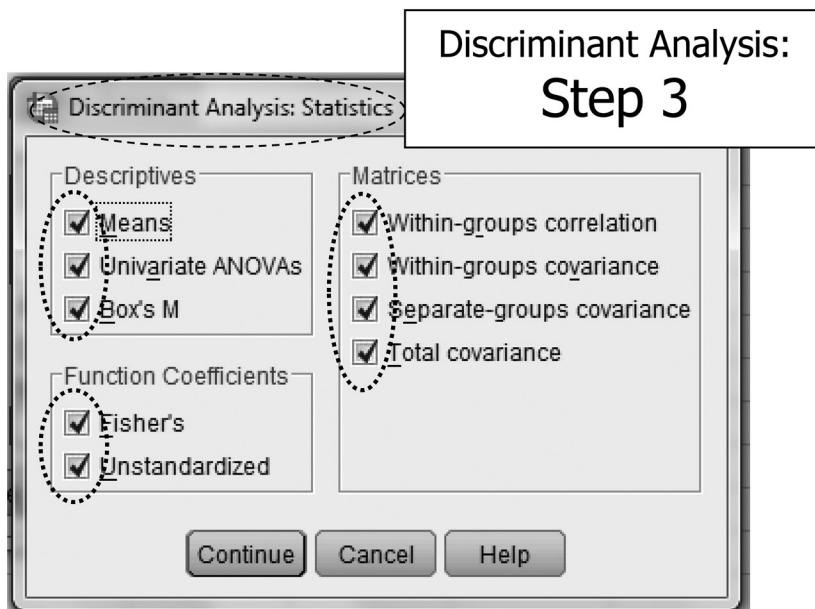
Step 2. Click the dependent variable (e.g., ‘Pre_Want_Sci_Bi’) and move it into the “Grouping Variable” box by clicking the arrow button. Click the independent variables and move them into the “Independents” box by clicking the arrow button (see screenshot Step 2).



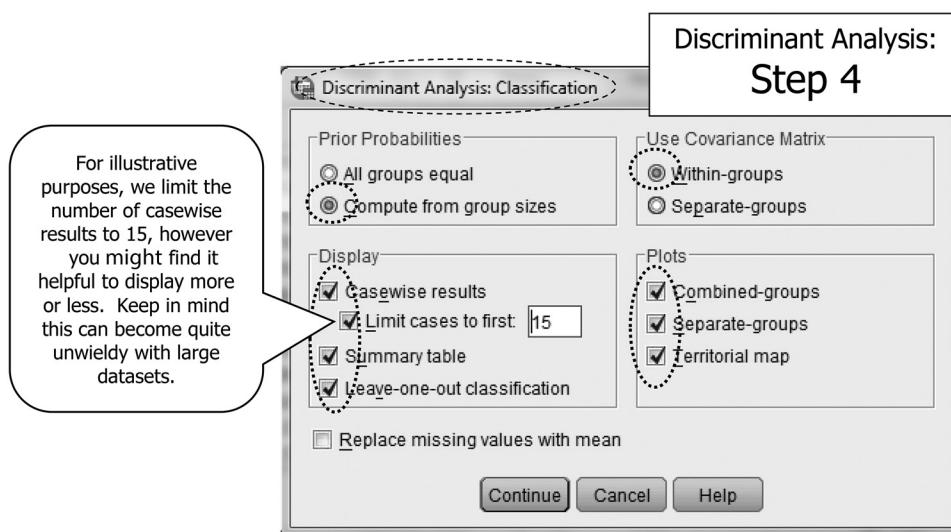
Step 2a. After selecting the dependent and independent variables, you must indicate the range of values of the dependent variable for which the discriminant analysis will be computed. Where there are only two groups, the two numbers that represent those categories are entered. If the dependent variable has more than two groups, the range of categories for which you want to analyze must be defined. Click “Define Range” to open the dialog box that allows you to identify the smallest and largest numeric category for the dependent variable (see screenshot Step 2a). In this example, there were two categories defined with ‘0’ and ‘1.’ Click on “Continue” to return to the main dialog box.



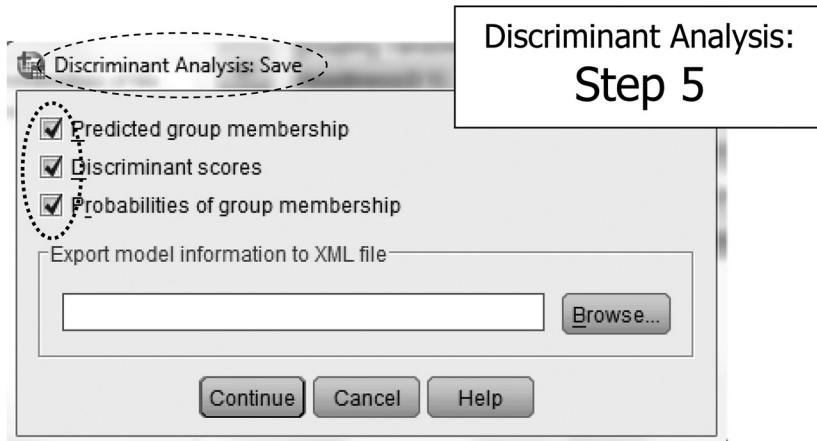
Step 3. From the “Discriminant Analysis” dialog box (see screenshot Step 2), click on “Statistics” to select options for descriptive statistics, function coefficients, and matrices (see screenshot Step 3). Place a checkmark in all boxes, and click on “Continue” to return to the “Discriminant Analysis” dialog box.



Step 4. From the “Discriminant Analysis” dialog box (see screenshot Step 2), click on “Classify” to select options for classification. From the “Classification” dialog box, under the heading of **Prior Probabilities** (see screenshot Step 4), place a checkmark in the box next to the following: (1) compute from group sizes. Under the heading of **Display**, place a checkmark in the box next to the following: (1) casewise results, (2) summary table, and (3) leave-one-out classification. When you select casewise results, you will then have the option to ‘limit cases to first,’ which tells SPSS to limit the casewise results to only the first so many cases in your dataset. For illustrative purposes, we will limit to the first 15 cases. You do not have to set a limit; however, the output can become quite voluminous even with a relatively small sample size such as what we are working with. If you truly want to understand your data, however, it may be helpful to output all casewise results, perhaps transferring them to a spreadsheet where aggregated results can be computed for review. Under the heading of **Use Covariance Matrices**, place a checkmark in the box next to the following: (1) within-groups. Under the heading of **Plots**, place a checkmark in the box next to the following: (1) combined groups, (2) separate groups, and (3) territorial maps. Click on “Continue” to return to the original dialog box.



Step 5. From the “Discriminant Analysis” dialog box (see screenshot Step 2), clicking on “Save” will provide the option to save predicted group membership, discriminant scores, and probabilities of group membership (see screenshot Step 5). These can be used for diagnostic and other examination. From the “Save” dialog box, place a checkmark in all boxes. Click on “Continue” to return to the original dialog box.



Interpreting the output. Annotated results for the two-group discriminant analysis are presented in Table 7.4, and annotated results for the three-group discriminant analysis are presented in Table 7.5. Please note that the two-group output includes much more detailed annotation. The annotation for the three-group output primarily presents interpretation of key results and highlights appreciable differences when generating two- versus three-group discriminant analysis.

TABLE 7.4

SPSS Results for Two-Group Discriminant Analysis

SYNTAX

```

GET
FILE='G:\filename.sav'.
DATASET NAME DataSet1 WINDOW=FRONT.
DISCRIMINANT
/GROUPS=PRE_WANT_SCI_BI(0 1)
/VARIABLES=PRETEST PRE_SELFEFF PRE_SELFREG
/ANALYSIS ALL
/SAVE=CLASS SCORES PROBS
/PRIORS SIZE
/STATISTICS=MEAN STDDEV UNIVF BOXM COEFF RAW CORR COV GCOV TCOV TABLE
CROSSVALID
/PLOT=COMBINED SEPARATE MAP
/PLOT=CASES(15)
/CLASSIFY=NOMISSING POOLED.

```

Analysis Case Processing Summary		
Unweighted Cases	N	Percent
Valid	66	100.0
Excluded	Missing or out-of-range group codes	0 .0
	At least one missing discriminating variable	0 .0
	Both missing or out-of-range group codes and at least one missing discriminating variable	0 .0
Total	0	.0
Total	66	100.0

This table simply provides descriptive data on the number of cases included in the analyses and any cases that were excluded (and the number excluded by reason).

I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	Mean	Std. Deviation	Valid N (listwise)	
			Unweighted	Weighted
No	SCIENCE KNOWLEDGE PRETEST	1.6774	1.77740	31 31.000
	PRE SELF-EFFICACY	4.1604	.51809	31 31.000
	PRE SELF REGULATION	4.3189	.57013	31 31.000
Yes	SCIENCE KNOWLEDGE PRETEST	2.4286	2.63780	35 35.000
	PRE SELF-EFFICACY	4.1772	.69762	35 35.000
	PRE SELF REGULATION	4.1208	.69920	35 35.000
Total	SCIENCE KNOWLEDGE PRETEST	2.0758	2.28918	66 66.000
	PRE SELF-EFFICACY	4.1693	.61524	66 66.000
	PRE SELF REGULATION	4.2134	.64479	66 66.000

'Group Statistics' provide basic statistics for the independent variables for each category of the dependent variable.

For students indicating 'yes' in response to wanting to be a scientist, the larger mean is associated with higher science content knowledge and higher self-efficacy but slightly lower self-regulation.

The 'Tests of Equality of Group Means' are one-way ANOVA F test results, using the discriminant analysis dependent variable as the one-way ANOVA independent variable and each independent variable of the DA as the one-way ANOVA outcome. The results give an indication on the potential of each predictor in the discriminant model. In this case, none of the predictors appear to contribute to the discriminant model.

	Wilks' Lambda	F	df1	df2	Sig. -
SCIENCE KNOWLEDGE PRETEST	.973	1.792	1	64	.185
PRE SELF-EFFICACY	1.000	.012	1	64	.913
PRE SELF REGULATION	.976	1.579	1	64	.214

Wilks' lambda values can give some indication of a variable's potential contribution to the discriminant model. Smaller Wilks' lambda values indicate more potential for the variable to be able to discriminate.

TABLE 7.4 (continued)

SPSS Results for Two-Group Discriminant Analysis

		SCIENCE KNOWLEDGE PRETEST	PRE SELF-EFFICACY	PRE SELF REGULATION
Covariance	SCIENCE KNOWLEDGE PRETEST	5.177	.171	.108
	PRE SELF-EFFICACY	.171	.384	.293
	PRE SELF REGULATION	.108	.293	.412
Correlation	SCIENCE KNOWLEDGE PRETEST	1.000	.121	.074
	PRE SELF-EFFICACY	.121	1.000	.737
	PRE SELF REGULATION	.074	.737	1.000

'Pooled within-group matrices' provides both covariance and correlation coefficients between the independent variables for all cases combined. Notice the degrees of freedom that are used to calculate the coefficients in these matrices differ by one as compared to the following covariance matrices.

a. The covariance matrix has 64 degrees of freedom.

The pooled within-group correlation coefficients can be used to examine the extent to which the assumption of non-collinearity has been met. Coefficients over .90 (in absolute value terms) suggest multi-collinearity. These coefficients suggest evidence this assumption has been met.

		SCIENCE KNOWLEDGE PRETEST	PRE SELF-EFFICACY	PRE SELF REGULATION
No	I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	3.159	-.069	.106
	PRE SELF-EFFICACY	-.069	.268	.182
	PRE SELF REGULATION	.106	.182	.325
Yes	SCIENCE KNOWLEDGE PRETEST	6.958	.383	.109
	PRE SELF-EFFICACY	.383	.487	.391
	PRE SELF REGULATION	.109	.391	.489
Total	SCIENCE KNOWLEDGE PRETEST	5.240	.172	.068
	PRE SELF-EFFICACY	.172	.379	.288
	PRE SELF REGULATION	.068	.288	.416

Covariance matrices for each group or category of the dependent variable can be found from the 'Covariance Matrices' table. The difference in coefficient values for the 'total' in this table as compared to the previous is due to differences in degrees of freedom (note 65 df used for these calculations).

a. The total covariance matrix has 65 degrees of freedom.

Analysis 1

Box's Test of Equality of Covariance Matrices

Log Determinants		
I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	Rank	Log Determinant
No	3	-1.810
Yes	3	-.608
Pooled within-groups	3	-.998

Log determinants provide a measure of variability of the group, with larger values indicating more variable groups. Thus, large differences in the log determinant values may be indicative of covariance matrices that differ between groups.

The ranks and natural logarithms of determinants printed are those of the group covariance matrices.

TABLE 7.4 (continued)

SPSS Results for Two-Group Discriminant Analysis

Test Results	
Box's M	11.140
Approx.	1.761
df1	6
F	28476.747
df2	
Sig.	.103

Tests null hypothesis of equal population covariance matrices.

Box's M can be used to test the assumption of equal covariances across groups. In this example, this test is not statistically significant ($p = .103$), thus there is evidence of equal covariances across groups. Had this test been statistically significant, we may have wanted to request 'separate matrices' to see if different classification results were computed. Recall, however, that Box's M is quite sensitive to both large and small sample sizes.

Summary of Canonical Discriminant Functions

Eigenvalues				
Function	Eigenvalue	% of Variance	Cumulative %	Canonical Correlation
1	.087 ^a	100.0	100.0	.283

a. First 1 canonical discriminant functions were used in the analysis.

The larger the eigenvalue, the greater the discrimination between groups. When there are more than two groups, multiple eigenvalues will be presented.

The squared canonical correlation is an effect size for the individual discriminant function. In this example, the squared canonical correlation is .08. Using Cohen's standards, this would be interpreted as a small effect.

Wilks' Lambda				
Test of Function(s)	Wilks' Lambda	Chi-square	df	Sig.
1	.920	5.232	3	.156

The overall effect size index can be computed from Wilks' lambda. Using Cohen's standards, this is a small effect.

$$\text{partial } \eta^2 = 1 - \Lambda^{1/3} = 1 - (.92)^{1/3} = 1 - .97 = .03$$

The p value for Wilks' lambda provides evidence of the extent to which the set of predictors statistically significantly discriminated between the groups. In this case, the set was not statistically significant in discriminating children who did or did not want to be a scientist.

Standardized Canonical Discriminant Function Coefficients	
	Function
	1
SCIENCE KNOWLEDGE	.549
PRETEST	
PRE SELF-EFFICACY	.878
PRE SELF REGULATION	-1.220

Standardized canonical coefficients allow comparison of variables that have different initial scales. Larger coefficients (in absolute value terms) indicate more discriminatory ability. In this case, pre self-regulation does a better job in the discriminant model.

Structure Matrix	
	Function
	1
SCIENCE KNOWLEDGE	.566
PRETEST	
PRE SELF REGULATION	-.532
PRE SELF-EFFICACY	.046

Correlations of each predictor with the discriminant function are presented in the 'Structure Matrix.' When squared, these correlations are effect size values indicating the proportion of shared variance between the independent variable and the discriminant function.

Pooled within-groups correlations between discriminating variables and standardized canonical discriminant functions

Variables ordered by absolute size of correlation within function.

TABLE 7.4 (continued)

SPSS Results for Two-Group Discriminant Analysis

Canonical Discriminant Function

Coefficients

	Function
	1
SCIENCE KNOWLEDGE	.241
PRETEST	
PRE SELF-EFFICACY	1.417
PRE SELF REGULATION	-1.900
(Constant)	1.595

Unstandardized coefficients

Based on these values, we can write out the discriminant function equation:

$$D = 1.595 + .241z_1 + 1.417z_2 + (-1.900)z_3$$

Functions at Group Centroids

I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	Function
	1
No	-.309
Yes	.274

Unstandardized canonical discriminant functions evaluated at group means

These values represent the group means of the discriminant function scores (by function if multiple functions are computed). The aggregated mean of discriminant function scores equals zero. One-half between the centroids (i.e., adding these values together and dividing by two) we find the cut score.

$$CutScore = \frac{-.374 + .274}{2} = -.05$$

A student that has a discriminant function score of greater than -.05 would likely state 'yes' they want to be a scientist (and less than -.05, 'no').

Classification Statistics

Classification Processing Summary

Processed	66
Missing or out-of-range group codes	0
Excluded	0
At least one missing discriminating variable	0
Used in Output	66

In the event that any cases were not classified in the discriminant model, the number of cases that failed classification for reason of exclusion would be presented here. In this illustration, all 66 cases were classified.

Prior Probabilities for Groups

I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	Prior	Cases Used in Analysis	
		Unweighted	Weighted
No	.470	31	31.000
Yes	.530	35	35.000
Total	1.000	66	66.000

Classification in discriminant analysis can be based on prior probabilities that represent equal distribution into groups (i.e., 'all groups equal' this is the default) or proportional to group size ('compute from group sizes'). We selected proportional to group size, thus the prior probabilities are calculated proportionally (i.e., 31/66 = .47).

Classification Function Coefficients

	I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	
	No	Yes
SCIENCE KNOWLEDGE	-.007	.134
PRETEST		
PRE SELF-EFFICACY	6.190	7.016
PRE SELF REGULATION	6.076	4.968
(Constant)	-26.747	-25.686

Classification function coefficients can be used to create discriminant function equations by group with cases classified into the group with the highest G .

$$G_{no} = -26.747 - .007z_1 + 6.190z_2 + 6.076z_3$$

$$G_{yes} = -25.686 + .134z_1 + 7.016z_2 + 4.968z_3$$

Fisher's linear discriminant functions

TABLE 7.4 (continued)

SPSS Results for Two-Group Discriminant Analysis

Casewise statistics can be quite useful to better understand discrimination on a case-by-case basis.

Each unit/case in the data is presented in the table, noted by **case #** thus the largest case # will reflect your sample size. The **actual group** represents the observed category of the dependent variable to which the case is identified.

Regardless of the number of groups in your dependent variable, the 'highest' and 'second highest' will be presented in the casewise statistics. The **highest group** represents the group of the dependent variable that was assigned the highest probability. For case 1, this student was *observed* to be group 0 which was 'no' (did not want to be a scientist) and *predicted* to be in group 0. Thus, 'no' was the group that the discriminant model assigned the highest probability. Given there were only two groups, the **second highest group** is therefore '1' which represented 'yes.'

	Case #	Actual Group	Highest Group						Second Highest Group			Discriminant Scores
			Predicted Group	P(D>d G=g)		P(G=g D=d)	Squared Mahalanobis Distance to Centroid	Group	P(G=g D=d)		Squared Mahalanobis Distance to Centroid	
				p	df							
Orig.	1	0	0	.976	1	.516	.001	1	.484	.375	.339	
The p represents the typicality probability (Huberty & Olejnik, 2006), the probability this case is from the assigned group. Very small typicality values (< .05) suggest large distance between the case and the centroid, and may be considered outliers.	2	0	0	.647	1	.578	.210	1	.422	1.083	.767	
	3	1	1	.878								
	4	1	1	.681								
	5	0	0	.326								
	6	0	0	.668								
	7	1	1	.853								
	8	1	0 ^{**}	.505								
	9	0	0	.332	1	.649	.943	1	.351	2.415	-1.280	
	10	0	1 ^{**}	.230	1	.729	1.441	0	.271	3.181	1.474	
	11	1	1	.123	1	.767	2.379	0	.233	4.518	1.816	
	12	1	0 ^{**}	.931	1	.525	.007	1	.475	.448	-.396	
	13	0	0	.867								
	14	1	0 ^{**}	.944								
	15	1	1	.856								
	1	0	0	.687	3	.504	1.481	1	.496	1.756		
Cross-validated ^b	2	0	0	.870	3	.572	.714	1	.428	1.536		
	3	1	1	.926	3	.547	.466	0	.453	.597		
	4	1	1	.857	3	.507	.769	0	.493	.583		
	5	0	0	.580	3	.638	1.962	1	.362	3.338		
	6	0	0	.512	3	.557	2.305	1	.443	3.003		
	7	1	1	.838	3	.592	.850	0	.408	1.352		
	8	1	0 ^{**}	.509	3	.639	2.319	1	.361	3.700		
	9	0	0	.607	3	.637	1.837	1	.363	3.206		
	10	0	1 ^{**}	.600	3	.760	1.868	0	.240	3.936		
	11	1	1	.358	3	.757	3.229	0	.243	5.254		
	12	1	0 ^{**}	.987	3	.529	.140	1	.471	.613		
	13	0	0	.660	3	.523	1.599	1	.477	2.029		
	14	1	0 ^{**}	.803	3	.534	.994	1	.466	1.509		
	15	1	0 ^{**}	.000	3	.688	42.831	1	.312	44.659		

For the original data, squared Mahalanobis distance is based on canonical functions.

For the cross-validated data, squared Mahalanobis distance is based on observations.

**. Misclassified case

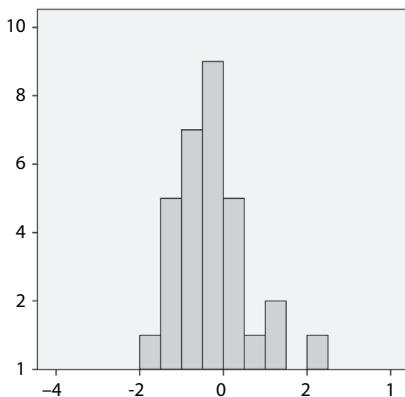
b. Cross validation is done only for those cases in the analysis. In cross validation, each case is classified by the functions derived from all cases other than that case.

TABLE 7.4 (continued)

SPSS Results for Two-Group Discriminant Analysis

Canonical Discriminant Function 1

I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY = Yes

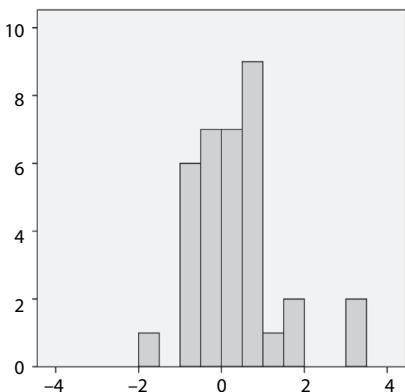


Mean = -0.27
Std. Dev. = 0.878
N = 31

Histograms provide a visual of the distributional shapes of the discriminant function scores by group. For this example, we can see there is a lot of overlap between the groups. Had there been more separation, one group would have been distributed more to the left and the other group distributed more to the right. This provides visual confirmation of the results already interpreted—this function does not discriminate well.

Canonical Discriminant Function 1

I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY = Yes



Mean = 0.27
Std. Dev. = 1.096
N = 35

		I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY	Predicted Group Membership		Total
			No	Yes	
Original	Count	No	20	11	31
		Yes	13	22	35
	%	No	64.5	35.5	100.0
		Yes	37.1	62.9	100.0
Cross- validated ^b	Count	No	19	12	31
		Yes	14	21	35
	%	No	61.3	38.7	100.0
		Yes	40.0	60.0	100.0

'Classification results' provides information on the percentage of cases correctly classified in both the initial analyses as well as the cross-validated analyses. In this example, about 64% of the original and 61% of the cross-validated cases were correctly classified. The cross-validated results are computed using the leave-one-out approach. All the data are used to calculate the discriminant function, then each case is classified using the discriminant function computed from all cases except for the one being classified, and the probability of misclassification is calculated from these results.

a. 63.6% of original grouped cases correctly classified.

b. Cross validation is done only for those cases in the analysis. In cross validation, each case is classified by the functions derived from all cases other than that case.

c. 60.6% of cross-validated grouped cases correctly classified.

■ TABLE 7.5

SPSS Results for Three-Group Discriminant Analysis

SYNTAX

```

GET
  FILE='G:\filename.sav'.
DATASET NAME DataSet1 WINDOW=FRONT.
DISCRIMINANT
  /GROUPS=POST_WANT_SCI_3CAT(1 3)
  /VARIABLES=PRETEST PRE_SELEFF PRE_SELFREG
  /ANALYSIS ALL
  /SAVE=CLASS SCORES PROBS
  /PRIORS SIZE
  /STATISTICS=MEAN STDDEV UNIVF BOXM COEFF RAW CORR COV GCOV TCOV TABLE
CROSSVALID
  /PLOT=COMBINED SEPARATE MAP
  /PLOT=CASES(15)
  /CLASSIFY=NONMISSING POOLED.

```

Analysis Case Processing Summary

Unweighted Cases		N	Percent
Valid		66	100.0
	Missing or out-of-range	0	.0
	group codes	0	.0
	At least one missing discriminating variable	0	.0
Excluded	Both missing or out-of-range group codes and at least one missing discriminating variable	0	.0
	Total	0	.0
Total		66	100.0

Group Statistics

I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)	Mean	Std. Deviation	Valid N (listwise)	
			Unweighted	Weighted
Disagree	SCIENCE KNOWLEDGE	1.6364	1.80404	11
	PRETEST			11.000
	PRE SELF-EFFICACY	4.0705	.62913	11
Neutral	PRE SELF REGULATION	4.4156	.48003	11
	SCIENCE KNOWLEDGE	1.6667	1.68034	18
	PRETEST			18.000
Agree	PRE SELF-EFFICACY	4.2903	.41226	18
	PRE SELF REGULATION	4.1884	.55478	18
	SCIENCE KNOWLEDGE	2.4054	2.64007	37
Total	PRETEST			37.000
	PRE SELF-EFFICACY	4.1398	.69429	37
	PRE SELF REGULATION	4.1655	.72536	37
	SCIENCE KNOWLEDGE	2.0758	2.28918	66
	PRETEST			66.000
	PRE SELF-EFFICACY	4.1693	.61524	66
	PRE SELF REGULATION	4.2134	.64479	66
	SCIENCE KNOWLEDGE			66.000
	PRETEST			66.000

'Group Statistics' provide basic statistics for the independent variables for each category of the dependent variable.

■ TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

The 'Tests of Equality of Group Means' are one-way ANOVA *F* test results, using the discriminant analysis dependent variable as the one-way ANOVA independent variable and each independent variable of the DA as the one-way ANOVA outcome. The results give an indication on the potential of each predictor in the discriminant model. In this case, no predictor appears to be contribute to the discriminant model.

Tests of Equality of Group Means					
	Wilks' Lambda	F	df1	df2	Sig.
SCIENCE KNOWLEDGE	.973	.870	2	63	.424
PRETEST					
PRE SELF-EFFICACY	.984	.525	2	63	.594
PRE SELF REGULATION	.980	.650	2	63	.526

Wilks' lambda values can give some indication of a variable's potential contribution to the discriminant model. Smaller Wilks' lambda values indicate more potential for the variable to be able to discriminate.

Pooled Within-Groups Matrices ^a				
		SCIENCE KNOWLEDGE PRETEST	PRE SELF- EFFICACY	PRE SELF REGULATION
Covariance	SCIENCE KNOWLEDGE	5.261		
	PRETEST		.190	.092
	PRE SELF-EFFICACY		.190	.301
Correlation	PRE SELF		.092	.420
	REGULATION			
	SCIENCE KNOWLEDGE	1.000		
Correlation	PRETEST		.133	
	PRE SELF-EFFICACY		.133	1.000
	PRE SELF		.062	.748
	REGULATION			1.000

a. The covariance matrix has 63 degrees of freedom.

The pooled within-group correlation coefficients can be used to examine the extent to which the assumption of non-collinearity has been met. Coefficients over .90 (in absolute value terms) suggest multicollinearity. These coefficients suggest evidence this assumption has been met.

■ TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

Covariance Matrices ^a				
I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)		SCIENCE KNOWLEDGE PRETEST	PRE SELF-EFFICACY	PRE SELF REGULATION
Disagree	SCIENCE KNOWLEDGE PRETEST	3.255	-.116	.295
	PRE SELF-EFFICACY	-.116	.396	.199
	PRE SELF REGULATION	.295	.199	.230
Neutral	SCIENCE KNOWLEDGE PRETEST	2.824	.167	.106
	PRE SELF-EFFICACY	.167	.170	.033
	PRE SELF REGULATION	.106	.033	.308
Agree	SCIENCE KNOWLEDGE PRETEST	6.970	.285	.029
	PRE SELF-EFFICACY	.285	.482	.455
	PRE SELF REGULATION	.029	.455	.526
Total	SCIENCE KNOWLEDGE PRETEST	5.240	.172	.068
	PRE SELF-EFFICACY	.172	.379	.288
	PRE SELF REGULATION	.068	.288	.416

a. The total covariance matrix has 65 degrees of freedom.

Analysis 1

Box's Test of Equality of Covariance Matrices

Log Determinants		
I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)	Rank	Log Determinant
Disagree	3	-2.152
Neutral	3	-2.001
Agree	3	-1.239
Pooled within-groups	3	-1.005

The ranks and natural logarithms of determinants printed

are those of the group covariance matrices.

TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

Test Results	
Box's M	36.815
Approx.	2.776
F	
df1	12
df2	4519.431
Sig.	.001

Tests null hypothesis of equal population covariance matrices.

Box's *M* can be used to test the assumption of equal covariances across groups. In this example, this test is statistically significant ($p = .001$), thus there is evidence of heterogeneity of covariances across groups. Given this, we can re-run the discriminant model and request 'separate matrices' to see if different classification results will be computed using separate covariance matrices rather than within-group covariance matrices. The classification results using separate matrices are presented at the very end of the output.

Summary of Canonical Discriminant Functions

Eigenvalues				
Function	Eigenvalue	% of Variance	Cumulative %	Canonical Correlation
1	.109 ^a	75.3	75.3	.313
2	.036 ^a	24.7	100.0	.186

a. First 2 canonical discriminant functions were used in the analysis.

The squared canonical correlation is an effect size for the individual discriminant function. In this example, the squared canonical correlation for the first function is .10 and .03 for the second function. Using Cohen's standards, these are interpreted as small effects.

Wilks' Lambda				
Test of Function(s)	Wilks' Lambda	Chi-square	df	Sig.
1 through 2	.871	8.591	6	.198
2	.965	2.179	2	.336

The overall effect size index can be computed from Wilks' lambda. Using Cohen's standards, these are small effects.

$$\text{partial } \eta^2 = 1 - \Lambda^{1/3} = 1 - (.87)^{1/3} = 1 - .95 = .05$$

$$\text{partial } \eta^2 = 1 - \Lambda^{1/3} = 1 - (.97)^{1/3} = 1 - .99 = .01$$

The *p* values for Wilks' lambda provide evidence of the extent to which the set of predictors statistically significantly discriminated between the groups. In this case, the set was not statistically significant in discriminating children who did or did not want to be a scientist for either function.

Standardized Canonical Discriminant Function

	Coefficients	
	1	2
SCIENCE KNOWLEDGE	.023	.928
PRETEST	-1.399	-.380
PRE SELF-EFFICACY	1.424	-.146
PRE SELF REGULATION		

Standardized canonical coefficients allow comparison of variables that have different initial scales. Larger coefficients (in absolute value terms) indicate more discriminatory ability. In this case, pre self-regulation does a better job in the discriminant model in the first function.

TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

Structure Matrix		
	Function	
	1	2
PRE SELF REGULATION	.379*	-.373
SCIENCE KNOWLEDGE	-.076	.869*
PRETEST		
PRE SELF-EFFICACY	-.330	-.366*

Correlations of each predictor with the discriminant function are presented in the 'Structure Matrix.' When squared, these correlations are effect size values indicating the proportion of shared variance between the independent variable and the discriminant function.

Pooled within-groups correlations between discriminating variables and standardized canonical discriminant functions

Variables ordered by absolute size of correlation within function.

*. Largest absolute correlation between each variable and any discriminant function

Canonical Discriminant Function Coefficients

	Function	
	1	2
SCIENCE KNOWLEDGE	.010	.405
PRETEST		
PRE SELF-EFFICACY	-2.257	-.613
PRE SELF REGULATION	2.197	-.225
(Constant)	.133	2.666

Based on these values, we can write out the discriminant function equations:

$$D_1 = .133 + .010z_1 + (-2.257)z_2 + 2.197z_3$$

$$D_2 = 2.67 + .405z_1 + (-.613)z_2 + (-.225)z_3$$

Unstandardized coefficients

Functions at Group Centroids

I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)	Function	
	1	2
Disagree	.663	-.163
Neutral	-.332	-.234
Agree	-.035	.162

These values represent the group means of the discriminant function scores by function.

Unstandardized canonical discriminant functions evaluated at group means

Classification Statistics

Classification Processing Summary

Processed	66		
Missing or out-of-range group codes	0		
Excluded At least one missing discriminating variable	0		
Used in Output	66		

In the event that any cases were not classified in the discriminant model, the number of cases that failed classification for reason of exclusion would be presented here. In this illustration, all 66 cases were classified.

TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

I WANT TO BE A SCIENTISIT (POST ATTITUDE; 3 CATEGORIES)	Prior	Cases Used in Analysis	
		Unweighted	Weighted
Disagree	.167	11	11.000
Neutral	.273	18	18.000
Agree	.561	37	37.000
Total	1.000	66	66.000

Classification in discriminant analysis can be based on prior probabilities that represent equal distribution into groups (i.e., 'all groups equal'; this is the default) or proportional to group size ('compute from group sizes'). We selected proportional to group size, thus the prior probabilities are calculated proportionally (i.e., $11/66 = .167$).

	I WANT TO BE A SCIENTISIT (POST ATTITUDE; 3 CATEGORIES)		
	CATEGORIES)		
	Disagree	Neutral	Agree
SCIENCE KNOWLEDGE PRETEST	.000	-.039	.125
PRE SELF-EFFICACY	5.393	7.683	6.770
PRE SELF REGULATION	6.648	4.478	5.041
(Constant)	-27.446	-27.126	-25.241

Fisher's linear discriminant functions

Classification function coefficients can be used to create discriminant function equations by group with cases classified into the group with the highest G .

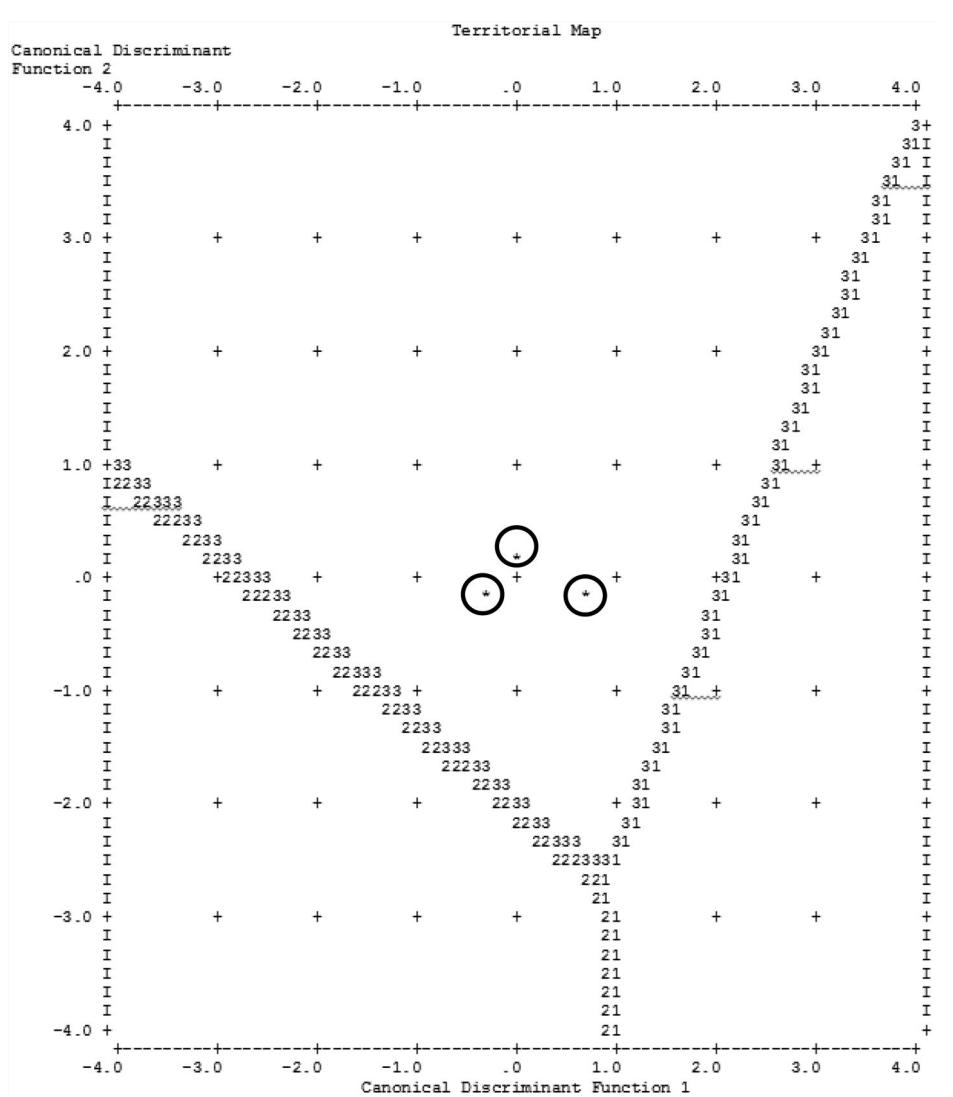
$$G_{\text{Disagree}} = -27.446 + .000z_1 + 5.393z_2 + 6.648z_3$$

$$G_{\text{Neutral}} = -27.126 - .039z_1 + 7.683z_2 + 4.478z_3$$

$$G_{\text{Agree}} = -25.241 + .125z_1 + 6.770z_2 + 5.041z_3$$

■ TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis



Symbols used in territorial map

Symbol Group Label

1	1	Disagree
2	2	Neutral
3	3	Agree
*		Indicates a group centroid

TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

	Case Number	Actual Group	Casewise Statistics									
			Predicted Group	Highest Group			Second Highest Group			Discriminant Scores		
				P(D>d G=g)	p	df	P(G=g D=d)	Group	P(G=g D=d)	Squared Mahalanobis Distance to Centroid	Function 1	Function 2
Original	1	1	3**	.852	2	.588	.320	2	.217	.868	.529	.122
	2	1	3**	.635	2	.521	.909	2	.267	.804	.461	-.651
	3	3	3	.835	2	.562	.362	2	.326	.011	-.400	-.316
	4	3	3	.955	2	.579	.093	2	.268	.193	.090	-.115
	5	3	3	.274	2	.455	2.590	2	.294	2.025	.585	-1.322
	6	2	3**	.673	2	.546	.791	1	.239	.014	.770	-.215
	7	3	3	.816	2	.659	.407	2	.215	1.203	.047	.795
	8	3	3	.267	2	.455	2.642	2	.360	1.672	.084	-1.459
	9	2	3**	.339	2	.494	2.161	1	.336	.503	1.357	-.308
	10	2	3**	.382	2	.626	1.926	2	.326	1.792	-1.333	.655
	11	2	3**	.233	2	.641	2.917	2	.321	2.662	-1.563	.927
	12	3	3	.843	2	.581	.343	2	.220	.846	.540	.057
	13	2	3**	.362	2	.477	2.030	2	.391	.984	-.356	-1.226
	14	3	3	.787	2	.542	.479	2	.296	.252	.081	-.520
	15	3	3	.236	2	.702	2.890	1	.166	3.342	.764	1.663
Cross-validated ^b	1	1	3**	.680	3	.603	1.510	2	.224	2.052		
	2	1	3**	.733	3	.529	1.283	2	.272	1.174		
	3	3	3	.902	3	.558	.575	2	.328	.195		
	4	3	3	.880	3	.574	.673	2	.271	.734		
	5	3	3	.361	3	.429	3.206	2	.306	2.442		
	6	2	3**	.441	3	.565	2.694	1	.254	1.868		
	7	3	3	.854	3	.654	.782	2	.219	1.532		
	8	3	3	.318	3	.426	3.520	2	.378	2.318		
	9	2	3**	.434	3	.506	2.738	1	.367	1.012		
	10	2	3**	.521	3	.653	2.258	2	.298	2.388		
	11	2	3**	.312	3	.686	3.567	2	.275	3.952		
	12	3	3	.947	3	.579	.366	2	.222	.843		
	13	2	3**	.551	3	.488	2.107	2	.376	1.188		
	14	3	3	.707	3	.532	1.396	2	.302	1.089		
	15	3	3	.000	3	.469	43.322	1	.371	41.363		

For the original data, squared Mahalanobis distance is based on canonical functions.

For the cross-validated data, squared Mahalanobis distance is based on observations.

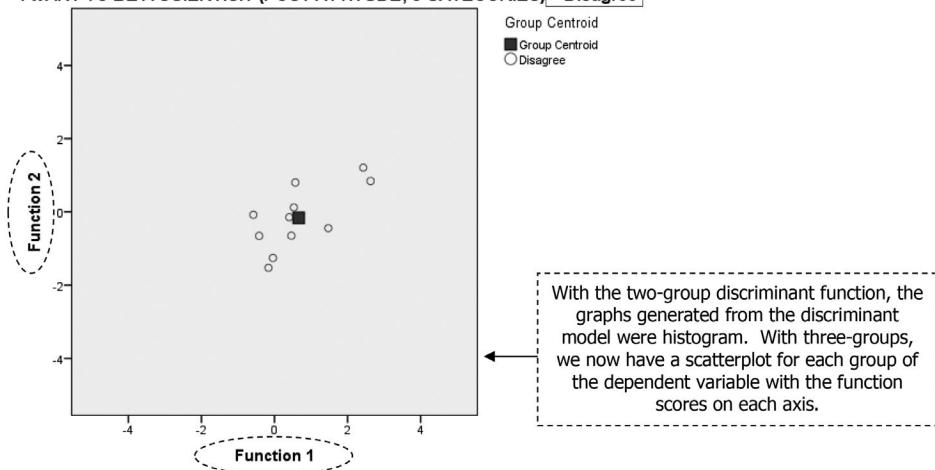
** Misclassified case

b. Cross validation is done only for those cases in the analysis. In cross validation, each case is classified by the functions derived from all cases other than that case.

Separate-Groups Graphs

Canonical Discriminant Functions

I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES) = Disagree



■ TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

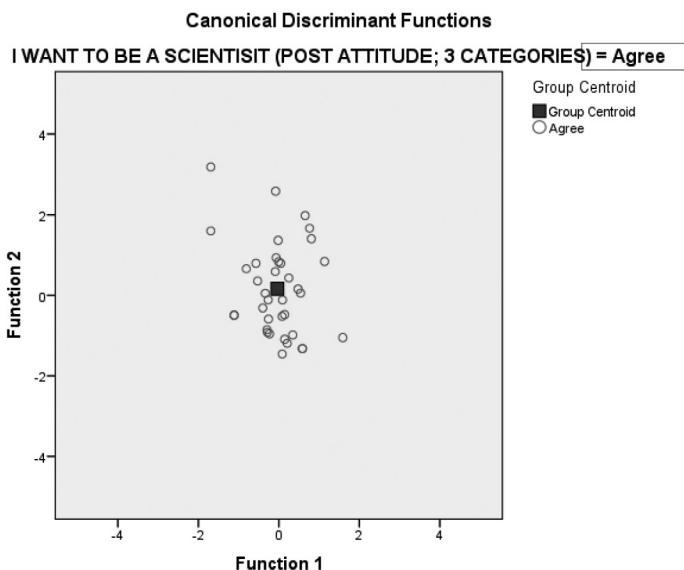
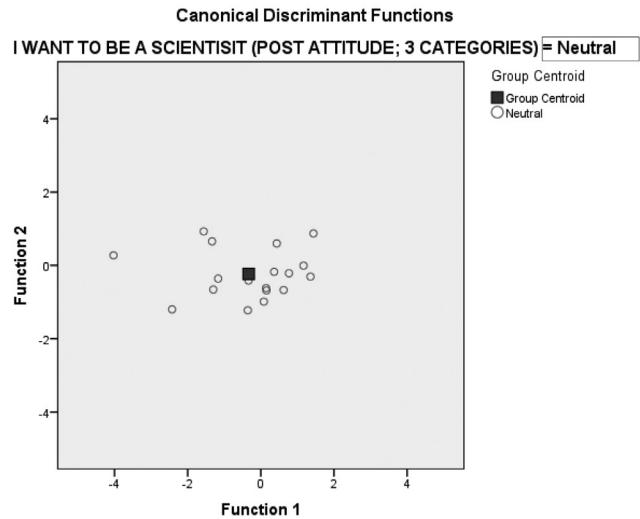
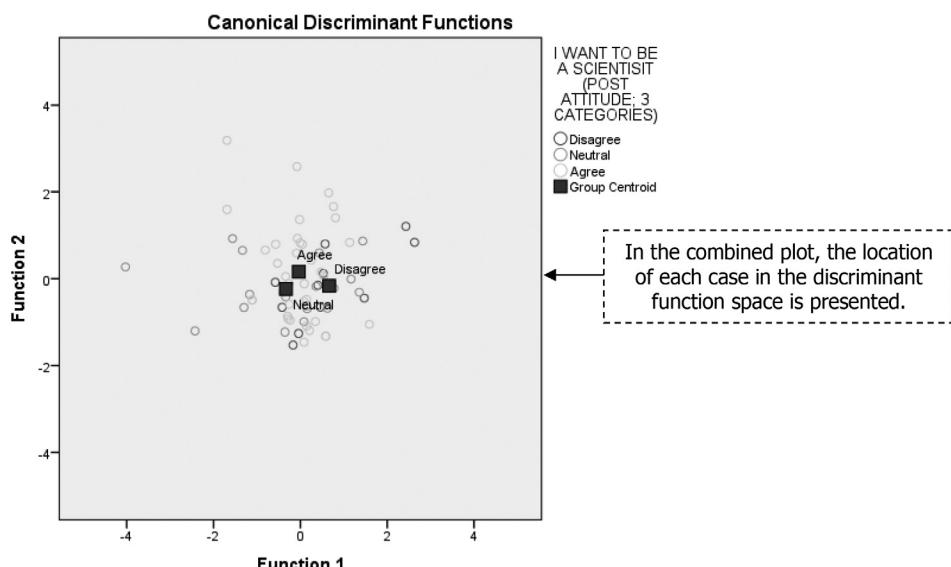


TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis



		I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)	Predicted Group Membership			Total
			Disagree	Neutral	Agree	
Original	Count	Disagree	1	0	10	11
	Count	Neutral	0	2	16	18
	Count	Agree	1	0	36	37
	%	Disagree	9.1	.0	90.9	100.0
	%	Neutral	.0	11.1	88.9	100.0
	%	Agree	2.7	.0	97.3	100.0
Cross-validated ^b	Count	Disagree	0	0	11	11
	Count	Neutral	0	1	17	18
	Count	Agree	1	0	36	37
	%	Disagree	.0	.0	100.0	100.0
	%	Neutral	.0	5.6	94.4	100.0
	%	Agree	2.7	.0	97.3	100.0

a. 59.1% of original grouped cases correctly classified.

b. Cross validation is done only for those cases in the analysis. In cross validation, each case is classified by the functions derived from all cases other than that case.

c. 56.1% of cross-validated grouped cases correctly classified.

'Classification results' provides information on the percentage of cases correctly classified in both the initial analyses as well as the cross-validated analyses.

In this example, about 59% of the original and 56% of the cross-validated cases were correctly classified

TABLE 7.5 (continued)

SPSS Results for Three-Group Discriminant Analysis

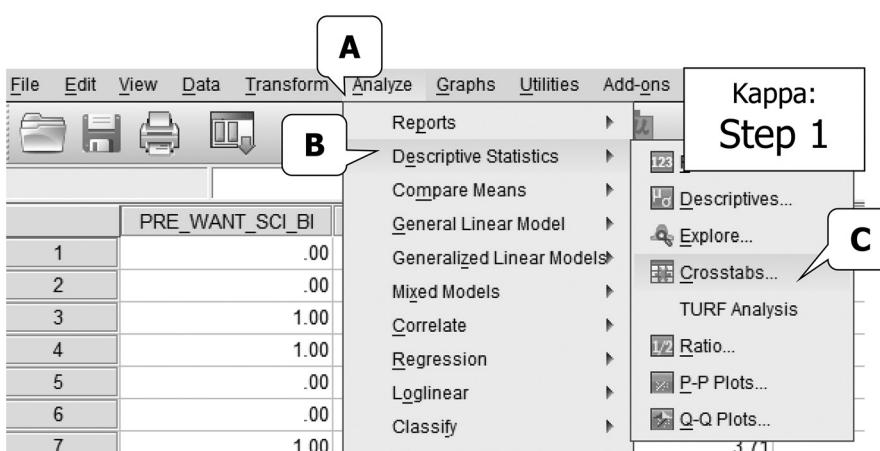
		Classification Results ^a			
		Predicted Group Membership			Total
Original Count	I WANT TO BE A SCIENTIST (POST ATTITUDE; 3 CATEGORIES)	Disagree	Neutral	Agree	
		2	1	8	11
	Disagree	2	5	11	18
	Neutral	1	3	33	37
	Agree	18.2	9.1	72.7	100.0
	% Disagree	11.1	27.8	61.1	100.0
	Neutral	2.7	8.1	89.2	100.0
	Agree				

^a a. 60.6% of original grouped cases correctly classified.

Given that Box's M suggested we may have violated the assumption of homogeneity of variance-covariance matrices, we can re-run the discriminant analysis and request that the covariance matrix used be based on 'separate groups' rather than 'within-groups.' Doing that, these are the classification results that are computed. They are not appreciably different than the model that used within-group covariance matrices.

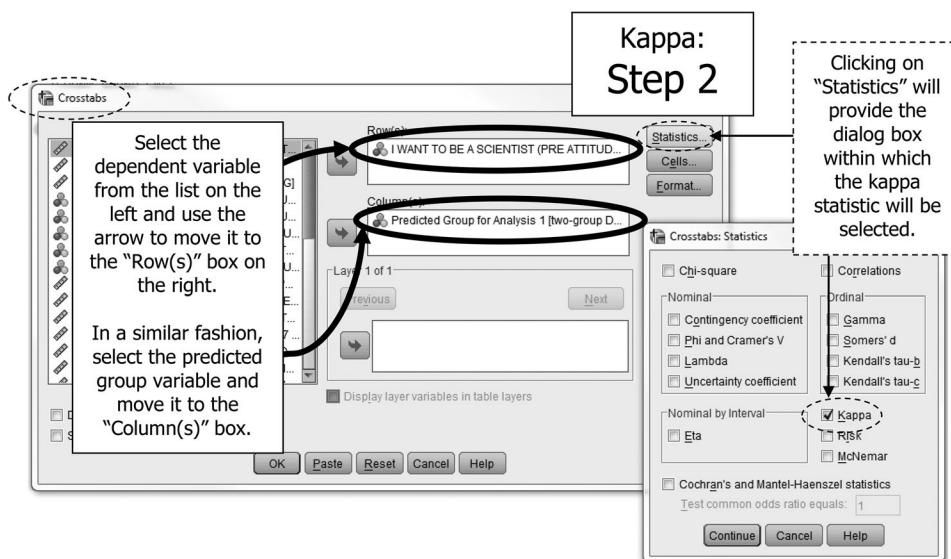
7.3.1 Generating Kappa Statistic for Classification Accuracy

Kappa Step 1. To generate a cross-tabulation table and the kappa statistic, go to "Analyze" in the top pull-down menu, then select "Descriptive Statistics," and then select "Crosstabs." Following the screenshot below (see screenshot Kappa: Step 1) produces the "Crosstabs" dialog box.



Kappa Step 2. Click the dependent variable (e.g., 'Pre_Want_Sci_Bi') and move it into the "Row(s)" box by clicking the arrow button. Click the predicted group membership variable (Dis_1) that was saved when generating the discriminant analysis and move it into the "Column(s)" box by clicking the arrow button (see screenshot Kappa: Step 2). Click on "Continue" to return to the original dialog box, and then click "Ok" to generate the output.

Interpreting the Kappa Statistic. Kappa can range from -1.0 to +1.0. A kappa value of one indicates perfect agreement, while a kappa of zero indicates the agreement is the



same as expected by chance. Negative values are rare but can occur and suggest that agreement is poorer than expected by chance. Cohen's guidelines for interpreting correlation coefficients can be applied to interpret kappa values, and thus kappa values of .10 are interpreted to have weak prediction, .30 have moderate prediction, and .50 or larger suggest strong prediction. In this example, a kappa value of .273 suggests moderate prediction.

Symmetric Measures

	Value	Asymp. Std. Error ^a	Approx. T ^b	Approx. Sig.
Measure of Agreement	Kappa	.273	.118	.220
N of Valid Cases		66		.026

a. Not assuming the null hypothesis.

b. Using the asymptotic standard error assuming the null hypothesis.

7.4 DATA SCREENING

Previously we described a number of assumptions of discriminant analysis. These included (a) independence, (b) linearity, (c) noncollinearity, (d) normality of independent variables, and (e) homogeneity of variance-covariance matrices. We also review the data to ensure there are no outliers.

Before we begin to examine assumptions, let us review the values that we requested to be saved to our data file (see following screenshot that shows values for the first 20 cases).

1. **DIS_1** are predicted group membership values
2. **DIS1_1** are discriminant scores (the average of these by group are the group centroids)
3. **DIS1_2** are probabilities of membership in group 0 for analysis 1
4. **DIS2_2** are probabilities of membership in group 1 for analysis 1

Note that additional variables will be created depending on the number of categories in your dependent variable. For example, if you run the *three-group* discriminant analysis example presented in Table 7.4, you will see that two columns of discriminant scores (rather than just one) will be created (one for discriminant function 1 and one for discriminant function 2), and three probability variables will be created (one for each category of the dependent variable).

As we look at the raw data, we see four new variables have been added to our dataset.

	Dis_1	Dis1_1	Dis1_2	Dis2_1
1	.00	-.33866	.51644	.48356
2	.00	-.76701	.57823	.42177
3	1.00	.12026	.44973	.55027
4	1.00	-.13773	.48716	.51284
	.00	-1.29082	.65042	.34958
	.00	-.73757	.57404	.42596
	1.00	.45916	.40147	.59853
8	.00	-.97593	.60762	.39238
9	.00	-1.28021	.64902	.35098
10	1.00	1.47428	.27068	.72932
11	1.00	1.81629	.23316	.76684
12	.00	-.39556	.52472	.47528
13	.00	-.47663	.53649	.46351
14	.00	-.37920	.52234	.47766
15	1.00	.45591	.40192	.59808
16	.00	-.81849	.58553	.41447
17	1.00	.71245	.36655	.63345
18	1.00	.59243	.38295	.61705
19	1.00	-.16635	.49133	.50867
20	.00	-.40825	.52656	.47344
..

7.4.1 Independence

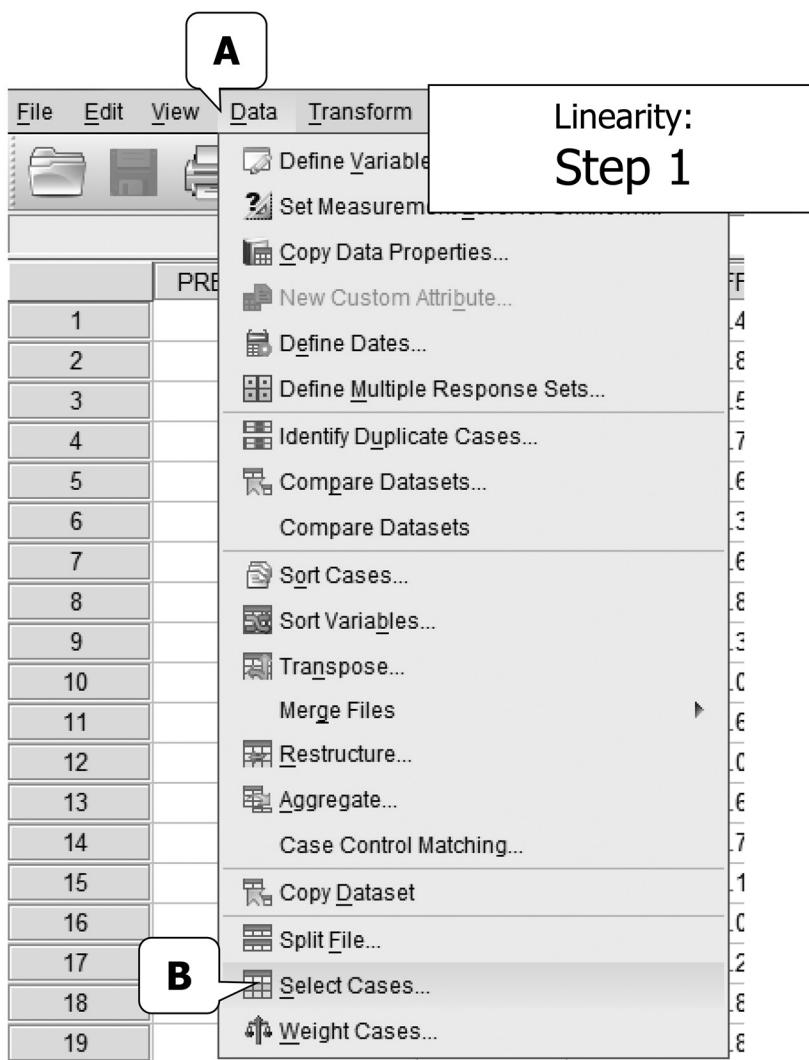
The first assumption is concerned with independence of the observations. As stated earlier, unfortunately, there are no satisfactory means by which to examine independence in discriminant analysis. In this case, data were not collected using a simple random sampling method, and thus our results run the risk of increased homogeneity. This may translate to increased Type I errors, and we may wish to caution readers to this risk as we present the results.

7.4.2 Linearity

The assumption of linearity in discriminant analysis is met when there are linear relationships between all pairs of independent variables within each group of the dependent variable. Violations of linearity may result in decreased power. Researchers that find nonlinear relationships may want to consider transforming or dropping one or more independent variables.

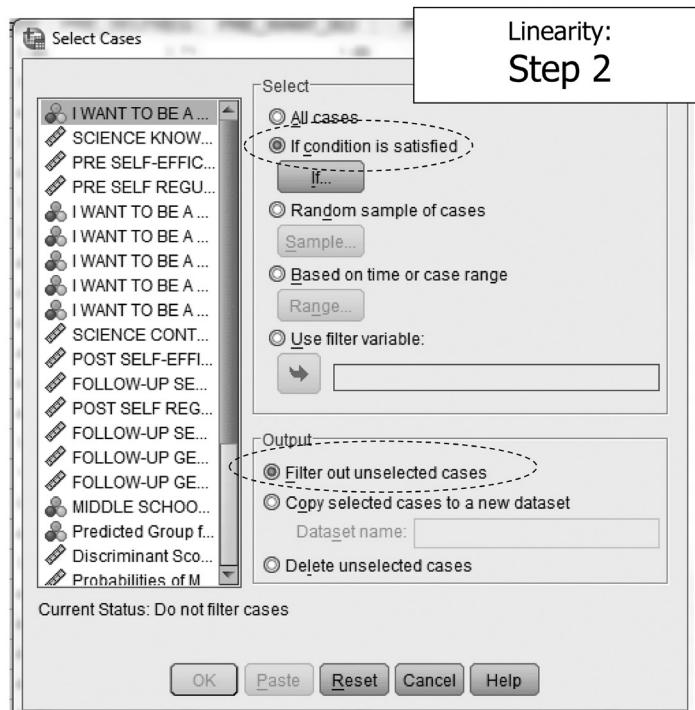
With a small number of predictors and a limited number of groups in the dependent variable, checking for evidence of linearity is relatively manageable. As one or both of the number of predictors or number of groups increases, examining evidence of linearity can be quite laborious but is an important data screening check. To check for linearity, we will first select one group of the dependent variable, and then generate a scatterplot of all pairs of the independent variables. That process is then repeated (but not illustrated here) for each remaining group of the dependent variable.

Linearity Step 1 (selecting group). To examine linearity, go to “Data” then “Select cases” (see screenshot Linearity: Step 1) to open the “Select Cases” dialog box.

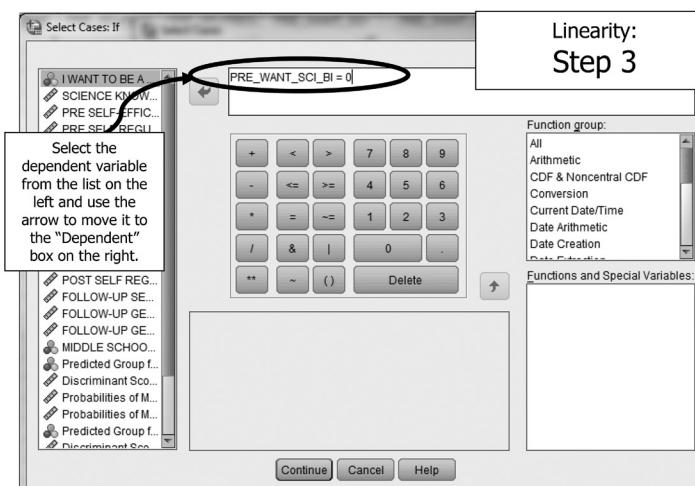


Linearity Step 2 (selecting group). From the “Select Cases” dialog box, click the radio button for “If condition is satisfied,” and then click “If . . .” (see screenshot Linearity: Step 2) to open the dialog box that allows you to define the selection

condition. The default Output option is to "Filter out unselected cases," and this is the desirable option in this case.



Linearity Step 3 (selecting group). To define the conditional statement to filter out one group of our dependent variable, click on the dependent variable from the list on the left and move to the box on the right by clicking the arrow (see screenshot Linearity: Step 3). Using your keyboard or the buttons within the dialog box, enter the numeric value for the group of the dependent variable that you want select. For this illustration, we will first select cases with a value of '0' on the dependent variable. Click on "Continue" to return to the original dialog box, and then click "OK" to run the analysis.



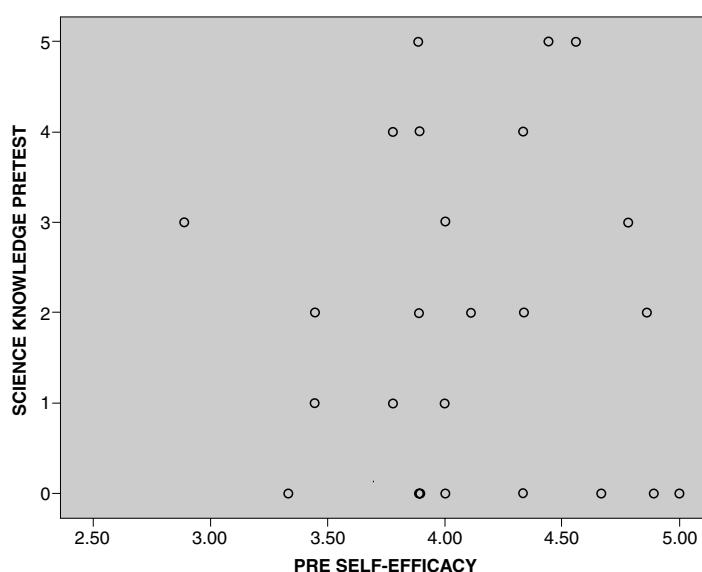
This process will filter out cases in group 1 of the dependent variable and allow you to examine linearity only for cases in group 0 (i.e., ‘no’ to the statement, “I want to be a scientist”). In the data view in SPSS, the cases filtered out will have slashes across their case number. Only those cases that do not have slashes will be included in any further analyses until the selection of cases is modified (or until the data file is closed in SPSS).

7.4.2.1 Using Scatterplots to Examine Linearity by Group

Linearity Step 4 (creating scatterplot). The previous steps allowed us to select one group of our dependent variable. We now need to generate a scatterplot within that group. Recall that linearity is reviewed for all possible pairs of the independent variables. Thus, generate scatterplots for all pairs of independent variables. When those are completed, return to Linearity: Step 1 to select a different group of the dependent variable by filtering out cases.

7.4.2.2 Interpreting Linearity Evidence

As you likely recall from your introductory statistics course, scatterplots are also often examined to determine visual evidence of linearity. Scatterplots are graphs that depict coordinate values of X and Y . Linearity is suggested by points that fall in a straight line. This line may suggest a positive relation (as scores on X increase, scores on Y increase and vice versa), a negative relation (as scores on X increase, scores on Y decrease, and vice versa), little or no relation (relatively random display of points), or a polynomial relation (e.g., curvilinear). In this example, our scatterplot suggests little evidence of any type of relationship. The scatterplots by group for the remaining independent variable pairs are not presented here. However, in examining them, there is evidence to suggest that linearity is reasonable (even though in some cases there is very little relationship of any sort, as seen here).



7.4.3 Noncollinearity

From the discriminant analysis output, we can check collinearity by reviewing the pooled within-groups correlation coefficients. These correlations are just as they sound—averaged (i.e., ‘pooled’) across groups of the dependent variable. Thus, these coefficients allow us to examine collinearity between predictors within group. Pooled within-group correlations that are .90 are stronger (in absolute value terms) suggest multicollinearity of those variables, and removal of one of the strongly correlated predictors is warranted as they are not uniquely contributing to the model.

Should the pooled within-groups correlations suggest potential multicollinearity, you may wish to conduct further checks of collinearity. As we did with logistic regression, we can estimate collinearity diagnostic values by running the same variables in a multiple regression model and requesting only the collinearity statistics. We are not interested in the parameter estimates of the model—only the collinearity statistics. Tolerance values less than .10 and VIF values greater than 10 indicate multicollinearity (recall that tolerance is simply the reciprocal of VIF) (Menard, 1995). Because the steps for generating collinearity statistics via multiple regression were presented previously in the text, we will not reiterate them here. Rather, we will merely present the applicable portion of the output for this model. With a range of tolerance values between .465 and .981 and VIF values ranging between 1.019 and 2.150, we have evidence of noncollinearity.

Coefficients^a

Model	Unstandardized Coefficients			Standardized Coefficients		t	Sig.	Collinearity Statistics	
	B	Std. Error	Beta					Tolerance	VIF
(Constant)	.750	.444				1.690	.096		
SCIENCE KNOWLEDGE	.033	.027		.151		1.231	.223		
1 PRETEST								.981	1.019
PRE SELF-EFFICACY	.195	.146		.239		1.337	.186		
PRE SELF REGULATION	-.262	.138		-.335		-1.891	.063		
								.465	2.150
								.471	2.122

a. Dependent Variable: I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY

The square root of the ratio of the largest eigenvalue to each successive eigenvalue is the ‘condition index,’ and this statistic is presented in the table labeled ‘Collinearity Diagnostics.’ The ‘dimension’ refers to a linear combination of the predictor variables, and the eigenvalue is the variance for that particular combination. A general guideline for interpreting condition indices is that values in the range of 10 to 30 should be of concern, greater than 30 indicates trouble, and greater than 100 indicates disaster (Belsley, 1991). Using these standards, there is some concern for multicollinearity.

The last four columns in this table refer to variance proportions. Multiplying these values by 100 provides a percentage of the variance of the regression coefficient that is related to a particular eigenvalue. Multicollinearity is suggested when covariates

have high percentages associated with a small eigenvalue (and large condition index). Thus for purposes of reviewing for multicollinearity, concentrate only on the rows with small eigenvalues. In this illustration, for example, the last dimension has a large condition index (24.48) with pre-self-efficacy and pre-self-regulation having large corresponding variance proportions. This suggests there may be issues with the estimation of those regression weights and thus there may be some multicollinearity.

Collinearity Diagnostics*

Model	Dimension	Eigen-value	Condition Index	Variance Proportions			
				(Constant)	SCIENCE KNOWLEDGE PRETEST	PRE SELF-EFFICACY	PRE SELF REGULATION
1	1	3.524	1.000	.00	.03	.00	.00
	2	.457	2.776	.00	.96	.00	.00
	3	.013	16.704	.99	.00	.10	.17
	4	.006	24.480	.01	.01	.90	.82

a. Dependent Variable: I WANT TO BE A SCIENTIST (PRE ATTITUDE) BINARY

In summary, there is some suggestion of multicollinearity based on examination of the condition indices and variance proportion values. However, we have met the assumption of noncollinearity with the tolerance and VIF values. Given that the pooled within-groups correlation coefficients were in the acceptable range, it seems reasonable to assume that the assumption of noncollinearity has been met.

7.4.4 Normality of Independent Variables

Univariate normality is a necessary condition for multivariate normality, thus we will examine univariate normality of independent variables. There are no residuals created when computing discriminant analysis, as there are in regression, thus we will examine the raw scores of the independent variables.

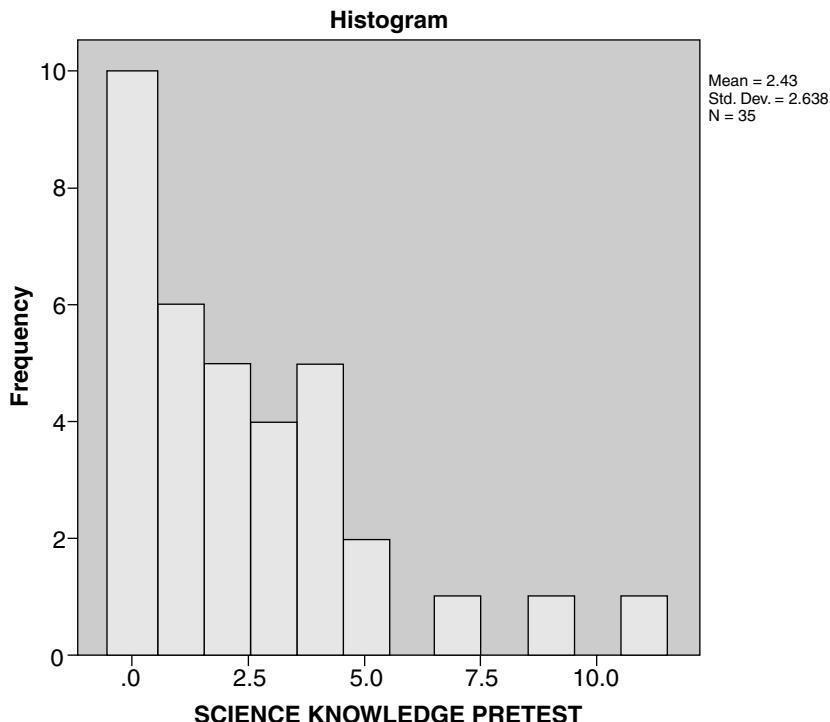
Descriptives

		Statistic	Std. Error
SCIENCE KNOWLEDGE PRETEST	Skewness	1.533	.398
	Kurtosis	2.653	.778
PRE SELF-EFFICACY	Skewness	-2.666	.398
	Kurtosis	10.412	.778
PRE SELF REGULATION	Skewness	-1.802	.398
	Kurtosis	5.328	.778

A general guideline is that skewness and kurtosis should be within the range of an absolute value of 2.0 and 7.0, respectively, and when that happens, there is evidence

to suggest normality. In this case, skewness and kurtosis are generally acceptable (although slightly high for pre-self-efficacy). When considering normality, skewness is more problematic than kurtosis, so we are less troubled by the large kurtosis values, although we will note that.

In reviewing the histograms, all reflect nonnormality with the skew quite obvious in all three (for illustrative purposes, only one histogram is presented).



There are a few other statistics that can be used to gauge normality, including the formal test of normality, the Shapiro-Wilk test (*SW*) (Shapiro & Wilk, 1965). The output for the Shapiro-Wilk test is presented and suggests that our sample distributions are all statistically significantly different from what would be expected from a normal distribution (in all cases, $p < .001$).

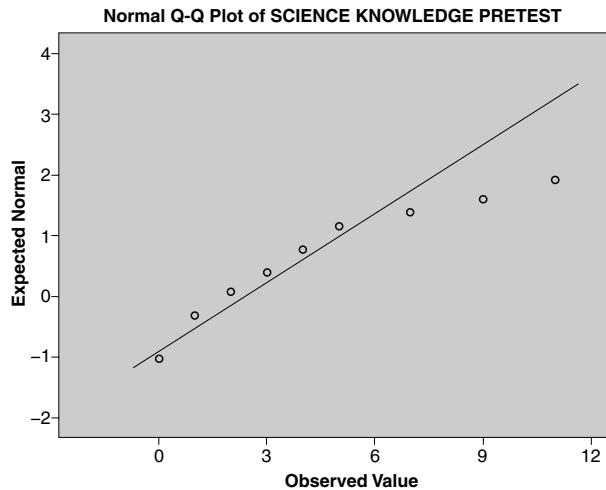
Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
SCIENCE KNOWLEDGE PRETEST	.179	35	.006	.832	35	.000
PRE SELF-EFFICACY	.171	35	.011	.761	35	.000
PRE SELF REGULATION	.118	35	.200*	.861	35	.000

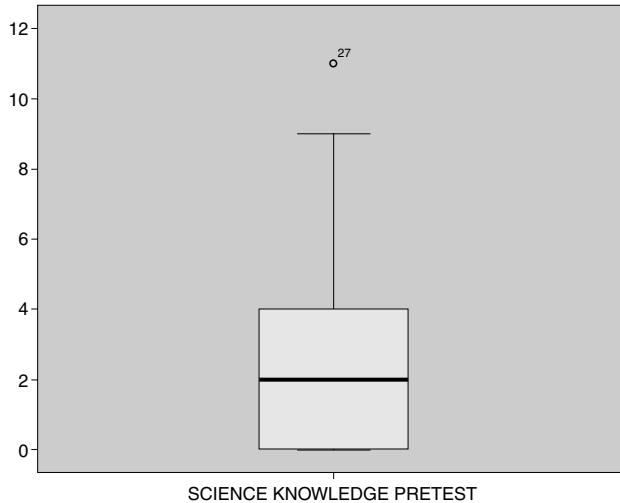
*. This is a lower bound of the true significance.

a. Lilliefors Significance Correction

The Q-Q plot of science content knowledge suggests some nonnormality, and nonnormality is evident in the other distributions as well (although the remaining graphs are not presented). For a refresher on interpreting Q-Q plots, review the data screening chapter.



Examination of the boxplot of science content knowledge also suggests some nonnormality, as well as the potential of an outlier. The boxplots of the other independent variables were very similar, also suggesting nonnormality with the presence of potential outliers.



Considering the forms of evidence we have examined, skewness and kurtosis statistics, the Shapiro-Wilk test, histogram, the Q-Q plot, and the boxplot, all suggest some degree of univariate nonnormality. Recall that in the case of unequal sample sizes of groups, such as we have, the discriminant analysis results will still be relatively robust given that we have only three independent variables (the standard being five or less independent variables) and the sample size of the smallest group is greater than 30 (the standard being at least 20 in the smallest group of the dependent variable).

7.4.5 Homogeneity of Variance-Covariance Matrices

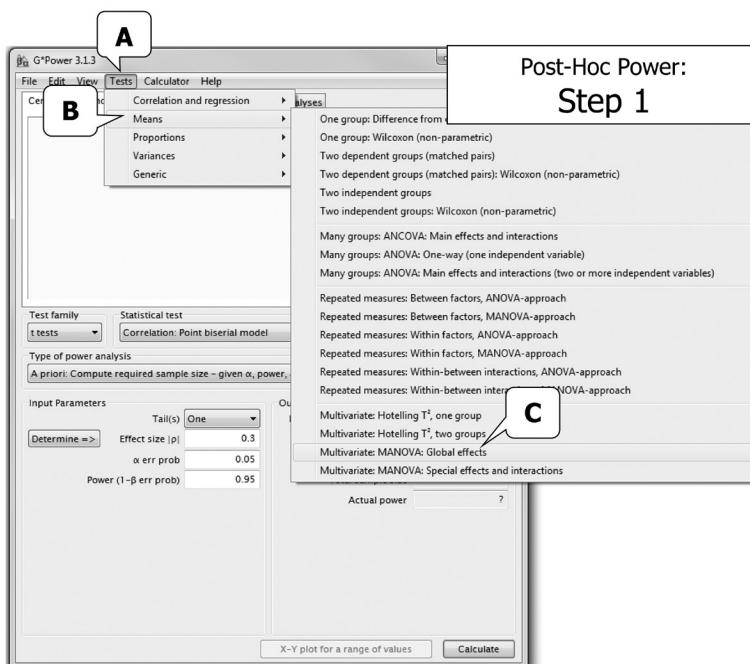
Homogeneity of variance-covariance can be tested using Box's M test, a test that can be generated when computing your discriminant analysis. Statistically significant results indicate a violation of the assumption. As seen in Table 7.3, Box's M is not statistically significant (Box's $M = 11.14, p = .10$), thus we have evidence of equal variance-covariance matrices. As noted earlier, Box's M has been criticized as being sensitive to large samples, and thus researchers may also want to use visual means to explore this assumption.

7.5 POWER USING G*POWER

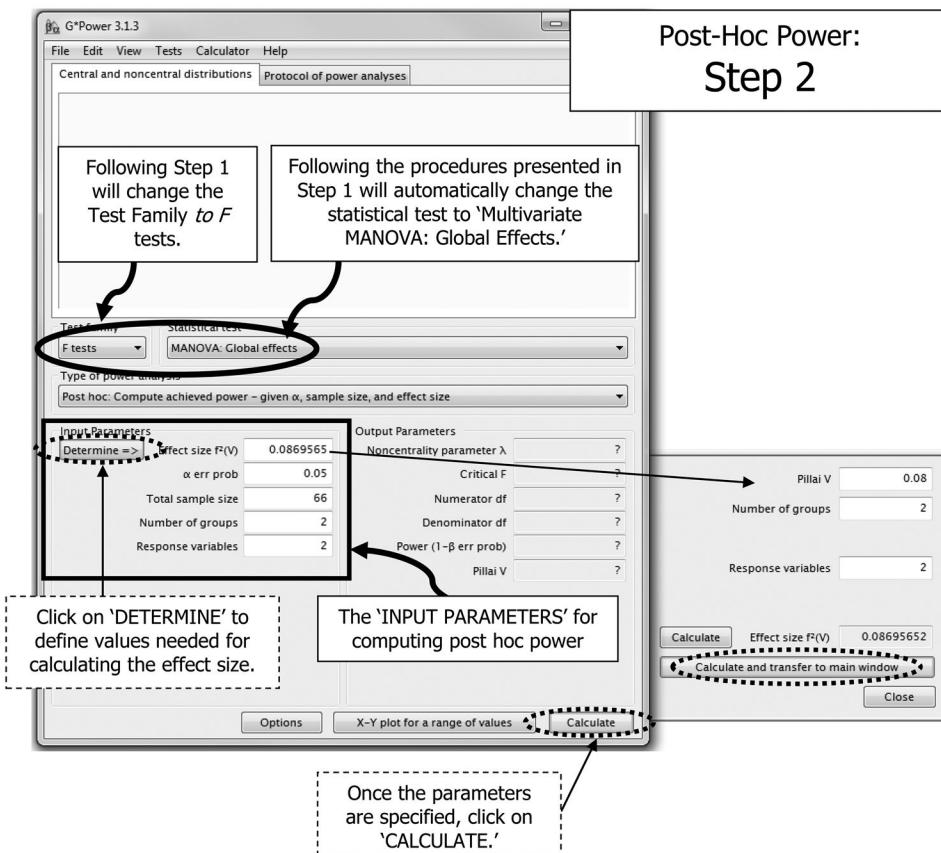
A priori and post hoc power can again be determined using the specialized software described previously in this text (e.g., G*Power), or you can consult a priori power tables (e.g., Cohen, 1988). You'll find that there is not an option for conducting power analysis specifically for discriminant analysis in G*Power. However, because the test of inference discriminant analysis is identical to that of MANOVA, power analysis for the two are the same as well (Cohen, 1988). As an illustration, we use G*Power to first compute post hoc power of our example.

7.5.1 Post Hoc Power for Discriminant Analysis Using G*Power

The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family. Recall that power analysis for discriminant analysis will be done via MANOVA. Thus, we select 'Tests' in the top pull-down menu, then 'Means,' and finally 'Multivariate MANOVA: Global effects' (see screenshot Post-Hoc Power: Step 1). Once that selection is made, the 'Test family' automatically changes to 'F tests' (see screenshot Post-Hoc Power: Step 2).



The 'Type of power analysis' desired then needs to be selected. To compute post hoc power, select 'Post hoc: Compute achieved power—given α , sample size, and effect size' (see screenshot Post-Hoc Power: Step 2).

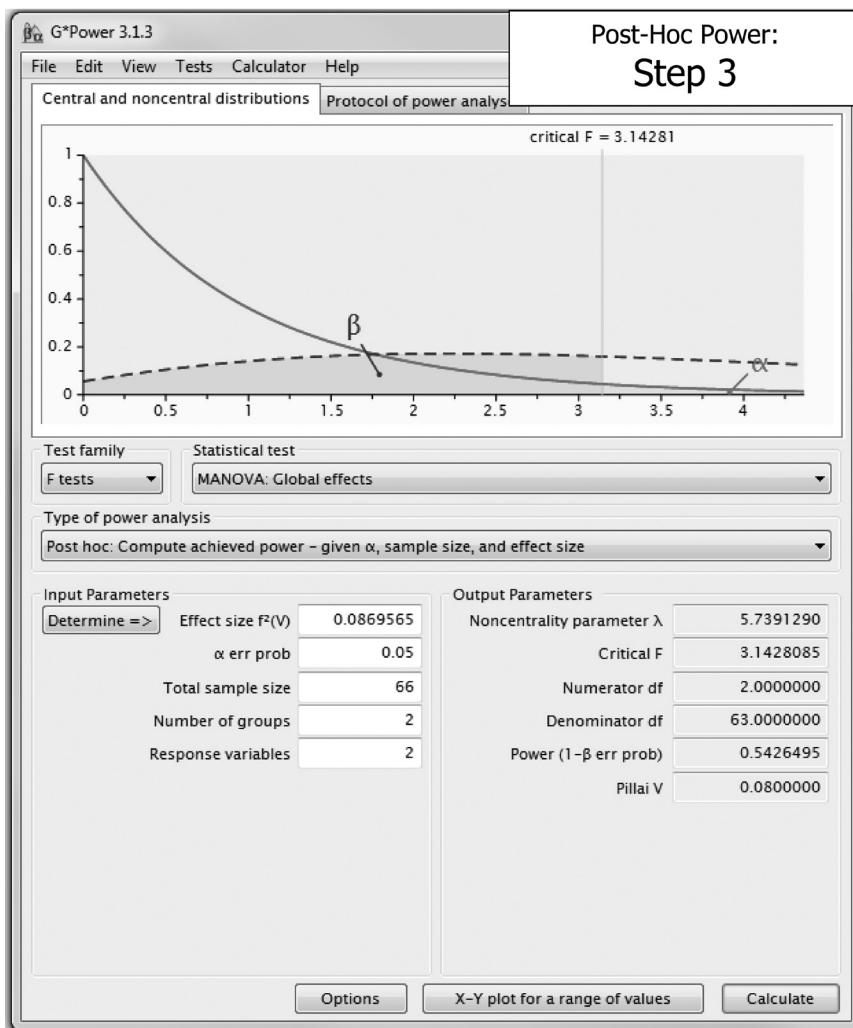


The 'INPUT PARAMETERS' must then be specified. The effect size, f^2 , is the first input parameter that must be defined. Clicking on 'DETERMINE' will provide a dialog box to define various parameters for computing this effect size. Pillai's V can be obtained from computing MANOVA with the discriminant analysis independent variables serving as the dependent variables in the MANOVA and the discriminant analysis dependent variable serving as the independent variable in the MANOVA. Pillai's V in this example is .08. Completing the rest of the statistics to compute the effect size, we have two groups (reflecting the number of groups of the dependent variable) and two response variables (reflecting the number of independent variables). After the values are entered, click on 'CALCULATE AND TRANSFER TO MAIN WINDOW' to compute the effect size f^2 and insert the value in the input parameters dialog box.

In terms of the remaining input parameters, the alpha level used was .05, and the total sample size was 66. The number of groups refers to the number of groups (or categories or levels) in the dependent variable, and in our case is two. The number of response variables refers to the number of independent variables or predictors. Again,

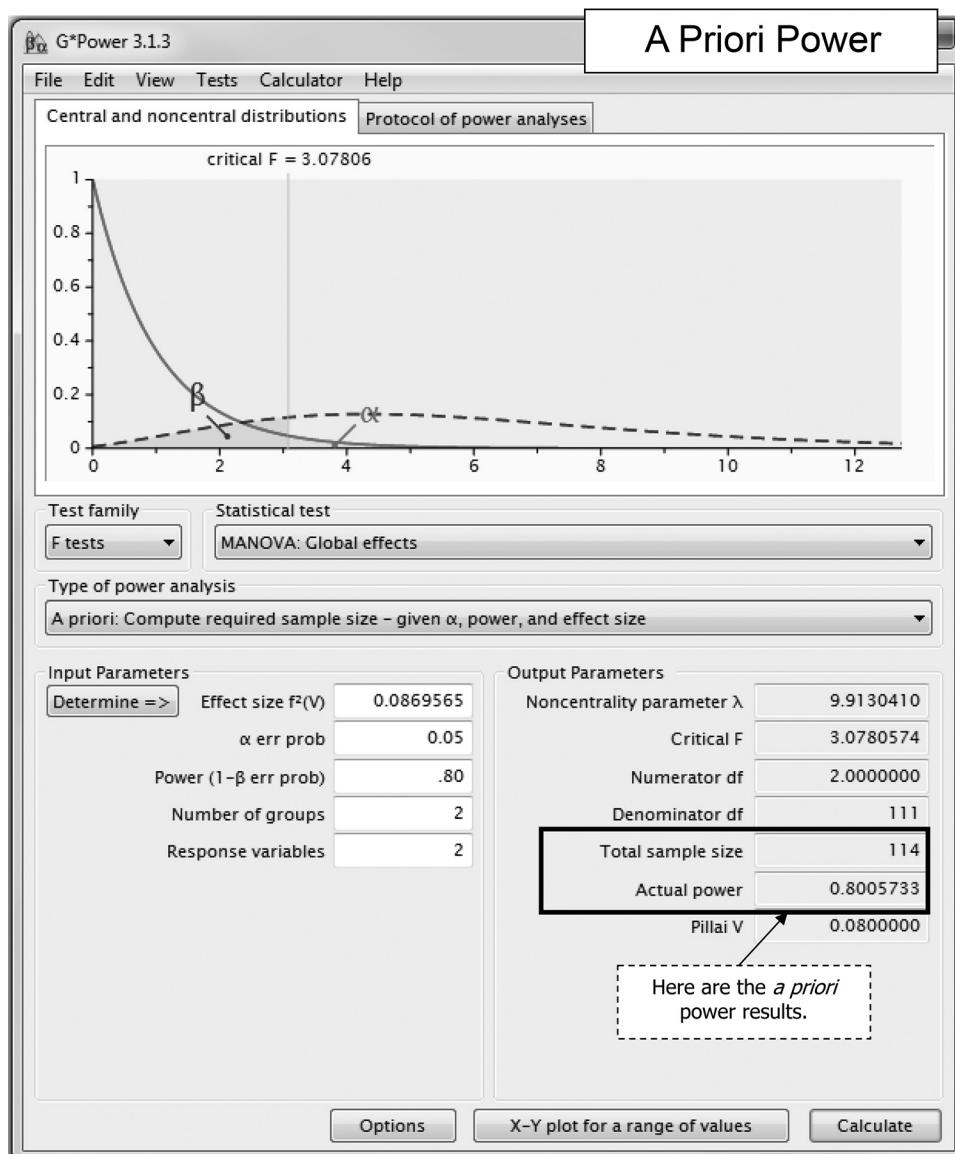
in our case, two. Once the parameters are specified, click on 'CALCULATE' to find the power statistics.

The 'OUTPUT PARAMETERS' provide the relevant statistics for the input just specified (see screenshot Post-Hoc Power: Step 3). In this example, we were interested in determining post hoc power for MANOVA that can be used as a proxy for power for discriminant analysis. Based on the criteria specified, the post hoc power was quite a bit less than what is considered sufficient power of around .80 (in our case, .54). In other words, the probability of rejecting the null hypothesis when it is really false was only about 54%. This finding is not surprising given the relatively small sample size. Keep in mind that conducting power analysis a priori is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters).



7.5.2 A Priori Power for Discriminant Analysis Using G*Power

For a priori power, we can determine the total sample size needed for discriminant analysis again using MANOVA and given the same parameters just discussed. In this example, had we wanted an a priori power of .80 given the same parameters just defined, we would need a total sample size of 114.



7.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example write-up for the results of the discriminant analysis. Recall that our graduate research assistant, Oso, was assisting Dr. Isealine, a faculty member in the psychology department. Dr. Isealine wanted to know if students identifying as ‘wanting to be a scientist’ (yes vs. no) could be predicted by science content knowledge, self-efficacy, and self-regulation (all continuous variables). The research question presented to Dr. Isealine from Oso included the following: Can students identifying as wanting or not wanting to be a scientist be predicted from science content knowledge, self-efficacy, and self-regulation?

Oso then assisted Dr. Isealine in generating discriminant analysis. A template for writing the research question for this design is presented below, following very similarly to that seen with logistic regression.

To what extent can [grouping variable] be identified from [list independent variables]?

It may be helpful to preface the results of the discriminant analysis with information on an examination of the extent to which the assumptions were met. The assumptions include (a) independence, (b) linearity, (c) noncollinearity, (d) normality of independent variables, and (e) homogeneity of variance-covariance matrices. We will also examine the data for outliers and influential points.

Discriminant analysis was conducted to determine whether self-report by middle school students as wanting or not wanting to be a scientist could be identified from science content knowledge, self-efficacy, and self-regulation.

There were four assumptions of discriminant analysis that were tested. Specifically, these assumptions included (a) linearity, (b) noncollinearity, (c) normality of independent variables, and (d) homogeneity of variance-covariance matrices.

Linearity was assessed by examining scatterplots of all pairs of independent variables within groups. The scatterplots did not suggest evidence of non-linear relationships.

In terms of **noncollinearity**, VIF values in the range of .465 to .981 (below the value of 10.0 which indicates the point of concern) and tolerance within the range of 1.019 to 2.150 (acceptable values being greater than .10) provide evidence of noncollinearity. However, in examining other collinearity diagnostics, condition index values between 10–30 were observed (specifically 16.70 and 24.48), and review of the variance proportions suggested that 90% and 82% of the variance of the regression coefficients for pre-self-efficacy and pre-self-regulation were related to the smallest eigenvalue. This suggests some concern for multicollinearity.

Univariate normality of the independent variables was assessed as a necessary condition of multivariate normality. Skewness and kurtosis were examined. While skewness was generally acceptable (less than 3.0 in absolute value), kurtosis

was large for pre-self-efficacy in particular (10.41). Review of formal tests of normality again suggest some evidence of nonnormality. For all three independent variables, Shapiro-Wilk was statistically significant ($p < .001$). These non-normality tendencies were also reflected in visual inspections of histograms and Q-Q plots. Boxplots for all three variables indicate potential outliers. Based on the evidence reviewed, there appears to be some degree of nonnormality. Due to the relatively small sample size however, all cases were retained. Readers are urged to interpret the results with caution given the possible influence of outliers.

As a formal test of **homogeneity of variance-covariance**, Box's M ($M = 11.14$) provides evidence of equal covariance matrices, $F(6, 28476.75) = 1.76, p = .10$.

Here is an example write-up of results for the discriminant analysis (remember that this will be prefaced by the previous paragraph reporting the extent to which the assumptions of the test were met).

Discriminant analysis was then conducted to determine whether middle school students self-report to the statement, "I want to be a scientist" (yes vs. no), could be predicted from three continuous variables (science content knowledge, self-efficacy, and self-regulation). The overall Wilks's lambda was statistically significant, $\Lambda = .92$, $\chi^2 (4, N = 66) = 5.23, p = .156$, partial $\eta^2 = .03$ [*Note to researchers: Partial eta squared is not included in the output but has to be calculated*], indicating that overall the predictors did not differentiate between students reporting that they did or did not want to be a scientist. The canonical R^2 was .08 [*Note to researchers: This can be found in the eigenvalue table—the canonical correlation coefficient in the table must be squared to arrive at this value; the squared value is not included in the SPSS output*], indicating that 8% of the variability of the scores for the discriminant function can be accounted for by differences between students who report that they want versus do not want to be a scientist. Both partial eta squared and the squared canonical correlation coefficient can be interpreted as suggesting a small effect.

In Table 1, we present the structure (loading) matrix of correlations [*Note to researchers: These values come from the structure matrix table*] between the predictors and the discriminant function as well as the standardized weights [*Note to researchers: These values come from the standardized canonical discriminant function coefficients table*]. Based on the coefficients, science content knowledge (measured at baseline; $F(1, 64) = 1.792, p = .185$) demonstrates the strongest relationship with the discriminant function followed by self-regulation (measured at pretest; $F(1, 64) = 1.579, p = .214$). In other words, the best predictors for distinguishing between students who want versus do not want to be a scientist when they enter the science summer camp (i.e., at baseline) are their baseline science knowledge and self-regulation (also at baseline). This is followed by baseline self-efficacy ($F(1, 64) = .012, p = .913$). However, none of the predictors was statistically significant in contributing to the discriminant model.

At baseline, students who state 'yes' in response to the statement, "I want to be a scientist" have higher average science knowledge ($M = 2.43, SD = 2.64$) as compared to students who state 'no' ($M = 1.68, SD = 1.78$). Self-efficacy is

■ TABLE 1

Structure Matrix and Standardized Coefficients

Predictor	Structure (Loading) Matrix Correlations	Standardized Coefficients for Discriminant Function
Science content knowledge	.566	.549
Pre-self-regulation	-.532	-1.220
Pre-self-efficacy	.046	.878

■ TABLE 2

Group Means (SD) of Predictors of Wanting to Be a Scientist

Predictor	'Yes' (<i>n</i> = 35)	'No' (<i>n</i> = 31)
Science content knowledge	2.43 (2.64)	1.68 (1.78)
Pre-self-efficacy	4.18 (.70)	4.16 (.52)
Pre-self-regulation	4.12 (.70)	4.32 (.57)

comparable between the two groups at baseline ('Yes,' $M = 34.18$, $SD = .70$; $M = 4.16$, $SD = .52$). However, self-regulation is slightly lower for students who state 'yes' they want to be a scientist ($M = 4.12$, $SD = .70$) as compared to 'no' ($M = 4.32$, $SD = .57$). On the basis of this relationship, the discriminant function will be labeled 'low self-regulation, high science knowledge and self-efficacy.' Table 2 provides the group means and standard deviations of predictors.

Using the sample proportions as prior probabilities, the discriminant function accurately predicted 64% of the students in our sample with only slightly more students stating 'no' they did not want to be a scientist as classified correctly (65% 'no' and 63% 'yes' correctly classified). To account for chance agreement in classification, the kappa coefficient was computed and found to be .27, a moderate value. Cross-validation, using the leave-one-out (jackknife) procedure, was conducted to provide evidence of how well the discriminant function would predict using a new sample (in other words, to determine the stability of the classification procedure). The percentage of correctly classified cases using the cross-validation was very similar, 61% (62% 'no' and 60% 'yes' of cross-validated cases correctly classified). This indicates a high degree of consistency in the classification scheme.

Post hoc power, calculated using G*Power (v. 3.1), was less than .54, indicating only moderate power.

PROBLEMS

Conceptual Problems

1. Which one of the following represents the primary difference between discriminant analysis and logistic regression?
 - a. The assumptions of the test
 - b. The measurement scales of the independent variables that can be included in the model
 - c. The measurement scale of the dependent variable
 - d. The statistical software that must be used to estimate the model
2. Which one of the following is NOT an appropriate dependent variable for discriminant analysis?
 - a. Binary
 - b. Continuous
 - c. Dichotomous
 - d. Multinomial
3. Which of the following would NOT be appropriate outcomes to examine with discriminant analysis?
 - a. Employment status (employed; unemployed not looking for work; unemployed looking for work)
 - b. Income (measured in whole numbers)
 - c. Marital status (married vs. not married)
 - d. Recreational athlete (athlete vs. nonathlete)
4. Which one of the following is NOT an assumption that must be considered when computing discriminant analysis?
 - a. Collinearity
 - b. Expected frequencies per cell
 - c. Linearity
 - d. Homogeneity of variance-covariance
5. Which one of the following is NOT a statistic that can be used to determine acceptable classification in discriminant analysis?
 - a. Kappa
 - b. Press's Q
 - c. Standards of comparison
 - d. Wilks's lambda
6. Which one of the following is a correct interpretation of Wilks's lambda?
 - a. Lambda values can range from -1.0 to +1.0
 - b. Lambda values near one suggest little to no differences between groups
 - c. Lambda values of zero suggest little separation between groups
 - d. The number of Wilks's lambda values computed will equal the number of groups of the dependent variable

7. Which one of the following represents the proportion of variance in the discriminant function that is explained by groups of the dependent variable?
 - a. Eigenvalue
 - b. Press's Q
 - c. Squared canonical correlation
 - d. Wilks's lambda
8. Cases are assigned to groups of the dependent variable based on which one of the following?
 - a. Canonical correlation
 - b. Cutting score
 - c. Eigenvalue
 - d. Standardized coefficient
9. Which one of the following represents strong prediction?
 - a. Kappa = .56
 - b. Overall $N = 20$
 - c. Structure coefficient = .05
 - d. Wilks's lambda = .90
10. A researcher finds the following standardized coefficients when discriminating to determine prediction for college readiness: (a) standardized entrance exam score, .76; (b) high school cumulative grade point average, .22; (c) college preparedness index, .56; and (d) school district exit exam, .34. Which predictor has the strongest contribution to the discriminant function score?
 - a. Standardized entrance exam score
 - b. High school cumulative grade point average
 - c. College preparedness index
 - d. School district exit exam

Computational Problems

1. You are given the following data, where X_1 (high school cumulative grade point average) and X_2 (standardized entrance exam score) are used to predict Y (college enrollment immediately after high school, '1,' versus delayed college enrollment or no enrollment, '0').

X_1	X_2	Y
4.15	30	1
2.72	17	1
3.16	27	0
3.89	24	1
4.02	25	1
1.89	13	0
2.10	20	1
2.36	21	1
3.55	23	0
1.70	18	0

Determine the following values based on simultaneous entry of independent variables: Wilks's lambda; Box's M ; canonical correlation; squared canonical correlation; cross-validation classification rate.

- You are given the following data, where X_1 (work addiction) and X_2 (intrinsic motivation) are used to predict Y (employed full-time = 1; employed part-time = 2; unemployed or looking for work = 3).

X_1	X_2	Y
76	65	1
62	81	1
59	72	2
83	45	2
77	84	1
93	72	3
80	92	2
92	97	3
37	56	3
56	42	1
72	65	1
45	42	2
98	99	1
86	92	1
69	59	2

Determine the following values based on simultaneous entry of independent variables: Wilks's lambda; Box's M ; canonical correlation; squared canonical correlation; cross-validation classification rate.

Interpretive Problem

- Use SPSS to develop a discriminant analysis model with the example survey data on the website. Utilize 'Do you smoke?' as the dependent (binary) variable to find at least two strong predictors from among the continuous variables in the dataset. Write up the results in APA style, including testing for the assumptions. Determine and interpret a measure of effect size.

REFERENCES

- Belsley, D.A. (1991). *Conditioning diagnostics: Collinearity and weak data in regression*. New York: Wiley.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.
- Fisher, R.A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7, 179–188.

- Huberty, C. J., & Hussein, M. H. (2003). Some problems in reporting use of discriminant analyses. *Journal of Experimental Education*, 71(2), 177–191.
- Klecka, W. R. (1980). *Discriminant analysis*. Newbury Park, CA: Sage.
- Menard, S. (1995). *Applied logistic regression analysis*. Thousand Oaks, CA: Sage.
- Shapiro, S. S., & Wilk, M. B. (1965). An analysis of variance test for normality (complete samples). *Biometrika*, 52(3 and 4), 591–611.

Chapter 8

CLUSTER ANALYSIS

CHAPTER OUTLINE

8.1 What Cluster Analysis Is and How It Works	336
8.1.1 Characteristics	336
8.1.2 Sample Size	344
8.1.3 Power	344
8.1.4 Effect Size	344
8.1.5 Assumptions	344
8.1.6 Conditions	345
8.2 Mathematical Introduction Snapshot	345
8.3 Computing Cluster Analysis Using SPSS	346
8.3.1 Checking Our Results From Step 6	350
8.4 Data Screening	358
8.5 Research Question Template and Example Write-Up	358

KEY CONCEPTS

1. Agglomerative clustering
2. Dendrogram
3. Distance matrix
4. Divisive clustering
5. Hierarchical clustering methods
6. K -means clustering algorithms
7. Partitioning clustering methods

Cluster analysis is a technique that shares some similarity to exploratory factor analysis in that it is not an inferential procedure, but rather a grouping procedure. More specifically, cluster analysis groups homogenous *units* of analysis based on their characteristics or attributes (where units can be defined very broadly such as individuals, employers, locations, organisms, and more) using proximity or distance. Scholars in areas such as psychology refer to this as a person-centered (as compared to variable-centered) approach. In comparison, EFA groups related variables based using relationships. Identifying common groups of units is helpful in many disciplines, and the use of cluster analysis has been applied in many disciplines including archaeology (e.g., Nunes et al., 2013), astronomy (e.g., Carlson, Profumo, & Weninger, 2013), economics (e.g., La Rocca, Stagliano, La Rocca, & Cariola, 2015), marketing (e.g., Aroeana & Michaelidou, 2014), medicine (e.g., Chen et al., 2014), psychology (e.g., Maree, Allott, Killackey, Farhall, & Cotton, 2015), sociology (e.g., Gomez & Parigi, 2015), and many others. Cluster analysis may be used by researchers as an end to itself—simply to identify profiles of persons or things and form a taxonomy. Many times, however, the profiles or classifications created from cluster analysis are then used in an inferential procedure to test a hypothesis to determine, for example, how the profiles relate to some outcome of interest.

Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying cluster analysis, (b) determine and interpret the results of cluster analysis, and (c) understand and evaluate how to screen data prior to conducting cluster analysis.

8.1 WHAT CLUSTER ANALYSIS IS AND HOW IT WORKS

The graduate research assistants in the stats lab have successfully navigated a number of challenging but exhilarating projects. Today, we find them discussing yet another exciting task.

Challie Lenge and Oso Wyse have just been approached by their faculty advisor to consult with Dr. B. Bryant, a faculty member in Sport Sciences who is examining statistics of Division I-A football teams of U.S. colleges and universities. More specifically, Dr. Bryant wants to develop profiles of the 128 Division I-A institutions based on various offensive, defensive, and efficiency statistics. From these profiles, Dr. Bryant intends to assist coaches in preparing for games and predicting outcomes of games. Dr. Bryant also anticipates that different league structures may be proposed based on the cluster analytic profiles. The research question that Dr. Bryant is interested in is the following: *What are the underlying profiles of Division I-A football teams?* Great sports enthusiasts, Challie and Oso quickly rise to the task, and after consulting with Dr. Bryant, recommend cluster analysis.

8.1.1 Characteristics

Cluster analysis uses characteristics or attributes to classify units (e.g., people) based on patterns. Units that are classified similarly will be in closer proximity to each other geometrically (i.e., high within-cluster or internal homogeneity) as compared to units

in different clusters (i.e., high between-cluster or external homogeneity). Cluster analysis is a multistep process that includes (a) variable selection, (b) clustering procedure selection, (c) cluster similarity measure selection, (d) cluster algorithm selection, (e) number of clusters selection, and (f) validation and interpretation of cluster solution. A discussion of each aspect follows, and the general steps in carrying out a cluster analysis are presented in Box 8.1.

8.1.1.1 Variable Selection

Variable selection based on theoretical and conceptual considerations is a common mantra that should be quite familiar, and the application of this thought process to cluster analysis holds. This is not the time to include ‘everything and the kitchen sink,’ as using too many clustering variables increases the chance of collinearity—which reduces the ability to identify distinct profiles. The profiles created from the cluster analysis process are derived from the variables that were used to form them. It follows, then, that great care should be taken in selecting the variables used in cluster analysis so that only theoretically and conceptually relevant variables are included as clustering variables. As part of the interpretation process, researchers should also exclude variables that do not differentiate across the clusters, even if they were initially considered theoretically important.

BOX 8.1 FITTING THE CLUSTER ANALYTIC MODEL

Element	Options
Clustering Procedure	<ul style="list-style-type: none">• Hierarchical: measured in distance<ul style="list-style-type: none">◦ Agglomerative◦ Divisive• Nonhierarchical: measured by within-cluster variation• Two-step clustering
Number of Clusters	Determine the number of clusters to retain and save the cluster solution: <ul style="list-style-type: none">• Agglomeration schedule• Dendrogram: Tree-like diagram; number of clusters to retain is based on clearly delineated breaks• Variance ratio criterion (<i>VRC</i>): more objective, the <i>VRC</i> is a mathematical calculation
Solution Stability	Determine the stability of the cluster solution: <ul style="list-style-type: none">• Rerun the cluster analysis using one or more different criteria (e.g., clustering method, distance measure)• For hierarchical clustering, reorder the units of analysis and rerun the cluster analysis with the same criteria (recall that the order of the units makes a difference when computing cluster analysis using hierarchical cluster)
Cross-Validation	Replicate the cluster analysis: <ul style="list-style-type: none">• Rerun the model with a different sample• Split the sample, using part for the analyses and the remainder for cross-validation
Identify the Clusters	Generate statistics (e.g., means) to compare/contrast the cluster solution and identify the clusters by name

8.1.1.2 Clustering Procedure Selection

In selecting a clustering procedure, we are determining *how* clusters will be formed. The formation of clusters is based on optimization. Optimization could be based on minimizing the variation within clusters, maximizing the distance between clusters, or determining similarity (or dissimilarity) between units in a cluster relative to units in other clusters. Clustering algorithms are the methods by which determination is made on how the clusters are formed. We will discuss hierarchical, nonhierarchical, and two-step clustering methods.

Hierarchical Methods

Hierarchical clustering procedures are not to be confused with hierarchical models (discussed later in the text). Rather, hierarchical clustering procedures are referred to as such because of their tree-like appearance or hierarchical structure that results from this type of method. Agglomerative and divisive are two types of hierarchical techniques that work very similarly but in opposite directions. **Agglomerative clustering** starts with each unit representing one cluster (i.e., the number of clusters equals the sample size), where the unit represents the unit of analysis (e.g., people, organizations, or whatever has been measured). The two most similar clusters, based on the smallest distance between clusters, are then joined or merged. Thus, at this step, the number of clusters is reduced by one. In this first step, therefore, two individual units are joined into a cluster. This process continues in a sequential manner such that a hierarchy of clusters, from bottom to top with each step reducing the number of clusters by one, is formed, with one single member cluster defining the top of the hierarchy. In this way, a tree-like structure is created. Thus, the total number of clusters formed in the process is $n - 1$, and the top cluster encompasses all the clusters from the lower level. **Divisive clustering** is a top-down process, essentially agglomerative in reverse, where all units start as one cluster and the clusters are sequentially divided and stop when each unit is a single cluster. In both agglomerative and divisive clustering, units are not reassigned once they become part of a cluster. This separates hierarchical clustering procedures from other types of clustering methods.

Cluster Similarity Measure Selection

The similarity or dissimilarity between units in hierarchical clustering methods has to be measured in some way. Hierarchical clustering *algorithms* are the mathematical vehicle with which the similarity/dissimilarity in distance between units is calculated, and this can be measured in many ways. Euclidean distance is straight-line distance (i.e., the length of the line that connects two units) and is a common measure of distance for interval and ratio variables. A distance matrix represents, on the off-diagonal, the distance between pairs of objects, and on the diagonal, the distance from each object to itself (i.e., 0). As with correlation matrices, the matrix is symmetrical above and below the diagonal—and thus only one needs to be referred. Chebychev distance is often used with ordinal data and is computed as the maximum of the absolute difference in clustering variables' values. Of the two considerations, the similarity (i.e.,

distance) measure and the clustering algorithm, the clustering algorithm is the most important consideration (Mooi & Sarstedt, 2011).

Cluster Algorithm Selection

There are a number of agglomerative algorithms, the most common of which include (a) single linkage, (b) complete linkage, (c) average linkage, (d) centroid method, and (e) Ward's method. The delineation between the algorithms is based on how the distance from a cluster to a unit or between clusters is defined. Using the same data, different cluster solutions may be derived depending on the algorithm selected.

Single Linkage

Single linkage is also referred to as ***nearest neighbor***. This algorithm is appropriately named as the similarity between clusters is defined as the shortest distance between any unit in one cluster to any unit in another cluster (i.e., the smallest distance between observations in two clusters). This is a versatile algorithm and can work with any number of different cluster patterns. However, clusters that have poor differentiation within cluster may result in joining dissimilar clusters based on the algorithm's objective of joining clusters that have the smallest distance between points. Similarity between clusters is represented by only one pair of units—the closest pair. Even though single linkage works best when clusters are in long chains, single linkage methods are commonly applied first as they are helpful in identifying outliers.

Complete Linkage

Complete linkage is also referred to as ***furthest neighbor***. In comparison to single linkage, complete linkage determines cluster similarity using the maximum or farthest distance between units in a cluster. As such, outliers have great impact on cluster solutions produced using the complete linkage algorithm. Within-group similarity is equal to the diameter of the cluster, as all the units in the cluster are united at the maximum distance. Similarity between clusters is represented by only one pair of units—the farthest pair.

Average Linkage

Similarity of any two clusters, where average linkage has been applied, reflects average similarity or distance of *all* units in one cluster (rather than just the smallest or largest) with all units in another. As compared to single and complete linkage, outliers are less problematic with average linkage given the use of all units. Generally, clusters created using the average linkage algorithm have small and relatively equal within-cluster variation, and the sample sizes within the clusters are more similar than clusters produced with other methods (the latter being a trait this algorithm shares with the centroid method). In SPSS, average linkage is referred to as ***between-groups linkage***.

Centroid Method

The similarity (i.e., distance) between clusters is defined as the distance between the geometric centers (i.e., centroids) of clusters. With each grouping of units, a new centroid is computed, and thus cluster centroids change as clusters are grouped. Clusters created using the centroid method have similar sample sizes within clusters as compared to clusters produced with other methods (this trait is common to both centroid and average linkage algorithms).

Ward's Method

With Ward's method, the similarity between clusters is defined as the sum of squares within clusters summed over all cluster variables. The combined clusters are those that minimize the total sum of squared distance. In other words, units are combined that produce the minimal increase in within-cluster variance. Outliers are especially problematic with this algorithm. The resulting clusters are relatively proportional in units, and thus if one or more clusters are expected to have a small proportion of units, the results from Ward's method may be deceiving.

Nonhierarchical Methods

Nonhierarchical methods are also referred to as *partitioning methods*. There are several nonhierarchical clustering algorithms; most of the most common ones (sequential, parallel, and optimizing) fall under the umbrella of ***k-means clustering algorithms***. These algorithms partition the data into clusters, the number of which is defined by the researcher, and then the units are assigned to a cluster in an iterative fashion, stopping when a criterion that minimizes the within-cluster variation and maximizes the between-cluster distance is reached. Nonhierarchical clustering procedures usually begin by identifying a starting point (i.e., cluster seed). In this first step, the researcher therefore specifies the number of clusters to create—a decision based on previous research, for example. Next, the units are assigned to a cluster based on similarity (e.g., shortest distance using Euclidean distance). This is followed by calculation of the geometric center of the cluster. Finally, the distances from each unit to the centroid are computed, with reassignment conducted such that minimal distance to the centroids is obtained. In an iterative fashion, the geometric centers are again computed with reassignment as appropriate until the solution converges (i.e., there is no additional reassignment of units to clusters).

K-means algorithms use the within-cluster variation, not distance measures, as the means by which clusters are formed. Another difference between hierarchical and nonhierarchical methods is that units can be reassigned to clusters at any point using *k*-means algorithms, and this potential of reassignment is not possible with hierarchical methods. Outliers are less impactful to *k*-means solutions and computationally less demanding for large sample sizes ($n = 500+$) and when there are a large number of clustering variables, and thus often preferred over hierarchical clustering procedures.

Because Euclidean distance is the distance measure with k -means, this method should be applied only when all clustering variables are at least interval in scale. Applying k -means to ordinal data may result in distorted cluster solutions. Lastly, a challenge in using k -means methods is the determination of the initial number of clusters. If determining this initial seed is difficult, VRC can be computed and used as the initial starting point. Additionally, a hierarchical clustering method may be applied first to get an idea of an initial starting point with k -means conducted afterwards.

Sequential Threshold

In this nonhierarchical clustering algorithm, all units within a predefined distance (i.e., ‘threshold’) are included in the first cluster seed. This same process continues for additional cluster seeds. The criticism of this method is that reassignment, even if it improves the cluster, is not possible once a unit has been assigned.

Parallel Threshold

Parallel threshold shares similarity to sequential threshold. Units are assigned to the nearest cluster center based on threshold distance, with all cluster seeds considered at the same time. Units can be reassigned to clusters with adjustments to the threshold levels.

Optimizing Procedure

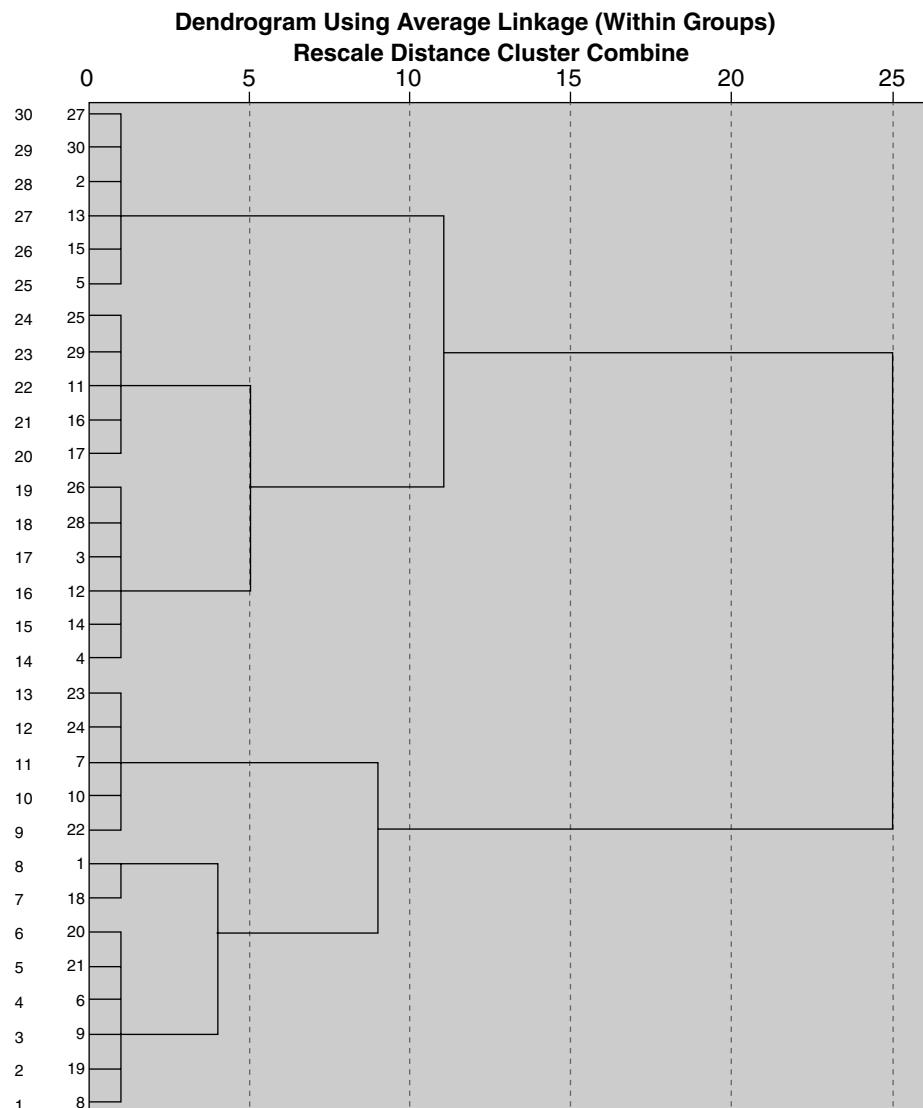
Reassignment of units to a different cluster seed is possible with this method. As implied by the name, therefore, the goal of this method is to create clusters that are optimal in terms of units within a cluster being optimally unique relative to other clusters.

Two-Step Clustering

The two-step clustering procedure handles large data ($n = 500+$) and clustering variables of mixed measurement scales quite well (Chiu, Fang, Chen, Wang, & Jeris, 2001). As the name implies, cluster analysis undertaken using the two-step method is conducted in two stages. In the first step, an algorithm similar to k -means is applied. In the second step, the results from the first step are considered as a modified hierarchical agglomerative clustering procedure is applied in which the units are combined sequentially, forming a cluster tree. The leaves of the cluster tree are the individual units. One of several benefits in the two-step clustering method is that all clustering variables, regardless of measurement scale, are considered simultaneously, and the importance of each clustering variable to the cluster solution is computed. Additionally, the number of clusters can be determined based on the empirical results of the model or the number of clusters can be specified a priori by the researcher. Measures of fit, including Akaike’s Information Criteria (AIC) or Bayes Information Criteria (BIC), are computed to assist in guiding the decision on the number of clusters to retain.

8.1.1.3 Number of Clusters

In our discussion of exploratory factor analysis, the number of factors to retain is a very important decision in the process. Similarly, determining the number of clusters to retain is an important decision that must be made in cluster analysis. A **dendrogram** is similar to a scree plot. It is a graph that can be used to make somewhat subjective interpretations of the number of clusters to retain. On the *X* axis of the dendrogram are the number of clusters, and on the *Y* axis are the units. The lines connect the units within clusters. Similar to a scree plot, the goal is to search for clearly delineated breaks. The challenge with this technique is that it is subjective, and identifying the natural break can be quite difficult. Additionally, the dendrogram becomes quite unwieldy the

**FIGURE 8.1**

Dendrogram Example

larger your sample size. The sample used later in the chapter has a sample size of 128. Although we use the dendrogram to help inform our decision in that example, it's easy to see that samples much larger than that would produce dendograms that are so large that they become useless in informing the decision-making process. Regardless, in some cases, dendograms can be useful aids for visualizing the clustering process.

The agglomeration schedule is another tool that can be used to inform the number of clusters to retain. The agglomeration schedule presents, in progressive order, how the units of analysis were combined into clusters based on distance and clustering method. The number of stages in the agglomeration schedule is one less than the sample size. The coefficients in the agglomeration schedule represent the distance between the units. The smaller the coefficient, the more homogenous the units. Read from the bottom up, large jumps in value in the coefficients in the agglomeration schedule suggest the number of clusters to retain (see agglomeration schedule and interpretation in Table 8.1). In other words, a big difference in coefficient values between two consecutive rows suggests heterogeneity and therefore a stopping point in the cluster solution.

More objective approaches to determining the number of clusters have been proposed. One is the **variance ratio criterion (VRC)** (Calinski & Harabasz, 1974) which is computed as F for a one-way ANOVA as follows:

$$VRC_k = \left(\frac{SS_B / (k-1)}{SS_W / (n-k)} \right)$$

Where n = sample size (i.e., number of units)

k = number of clusters

SS_B = sum of squares between the clusters

SS_w = sum of squares within the clusters

Using the VRC , the number of clusters to retain is then computed as follows, with the number of clusters being that which minimizes ω_k :

$$\omega_k = (VRC_{k+1} - VRC_k) - (VRC_k - VRC_{k-1})$$

While the VRC can be a helpful guide in determining the number of clusters to retain, by default of the calculation, three is the minimum number of clusters that can be selected. Thus, as with EFA, the number of clusters to select should be based on a solution that has meaning and can be interpreted.

8.1.1.4 Cross-Validation of Cluster Solution

Interpretation of the cluster solution goes beyond simply evaluating the cluster solution and also includes examination of the validity of the cluster solution. As is recommended with other procedures such as EFA and CFA, cross-validating the results to

determine cluster analytic stability with one-half of the sample (when sample size provides for that) is ideal. Additionally, sensitivity analyses using different cluster procedures and decisions (e.g., different distance measures) can lend evidence to the validity of the results. Researchers should keep in mind that finding different solutions, when different methods are applied, is common even when the solution is sufficient. There is a fair degree of subjectivity in determining the extent to which solutions may differ but yet judged to be adequate.

8.1.1.5 Interpreting the Cluster Solution

Interpretation of the cluster solution is done by examination of the cluster centroids, the geometric center of the clusters. Clusters that are clearly differentiated from each other in their centroids are desirable. Application of a hypothesis test, such as analysis of variance, can be used at this step.

8.1.2 Sample Size

Cluster analysis is not an inferential technique. Thus, similar to exploratory factor analysis, there is not a power calculation to suggest appropriate sample size for cluster analysis. The most important sample size consideration in cluster analysis is sufficiency in generating representation of the smallest group. In small sample size analyses, outliers may truly be outliers or they may represent a group that includes a small proportion of cases. Only when there is sufficient sample size will there be an increased potential of correctly identifying the cases as one or the other. One general recommendation for cluster sample size is that the number of observations be at least 2^m , where m equals the number of clustering variables. For example, if there were 10 clustering variables, the sample size would need to be at least 2^{10} or 1,024. This suggests that cluster analysis is a large sample size procedure, and reiterates the importance of careful selection of the clustering variables.

8.1.3 Power

As stated previously, there are not power calculations to suggest appropriate sample size for cluster analysis given a priori or post hoc power. What exists are a number of sample size recommendations as presented previously.

8.1.4 Effect Size

Cluster analytic solutions, in and of themselves, do not produce effect size results. Once the cluster solutions are applied in an inferential procedure, effect sizes can then be generated.

8.1.5 Assumptions

The clustering variables are assumed to be independent. If this assumption holds, there should not be high correlations among the clustering variables. Correlations among

the clustering variables that are very strong show evidence of multicollinearity, with which you should be quite familiar. In the case of highly correlated clustering variables, the cluster solution is influenced more strongly by the multicollinear variables rather than the clustering variables that are uncorrelated.

8.1.6 Conditions

Distance measures are appropriate to apply when the clustering variables are interval and/or ratio in scale. Even though ordinal data is not continuous, distance measures are appropriate, as treating them as nominal will disregard the ranking of the categories, allowing proximity to be best understood. Nominal or binary data, on the other hand, must apply a different similarity measure such as a matching coefficient. Nominal data with more than two categories should be converted into dichotomous variables so that matching coefficients can be applied. A viable alternative, when conducting cluster analysis with metric and nominal variables, is to conduct a **two-step clustering procedure**.

More so than not, variables with unlike response scales and/or with different measurement scales are of interest as clustering variables. In these situations, standardization may be the solution. Standardizing to a unit normal distribution (i.e., z , where the mean is zero and the standard deviation is one) is one option, as are other types of standardization (e.g., standardization by range, 0 to +1 or -1 to +1). Researchers may also consider using the correlation, rather than distance, between measures as a way to imply standardization.

8.2 MATHEMATICAL INTRODUCTION SNAPSHOT

As stated previously, similarity (or dissimilarity) is measured by distance in hierarchical methods. One of the most common measures of distance when working with interval or ratio data is Euclidean distance. Euclidean distance is simply the distance of a straight line computed as the square root of the variables' sum of squared differences. The equation for the Euclidean distance follows when the distance between two variables (X and Y) for two units of analysis (1 and 2) is desired:

$$d_{\text{Euclidean}} = \sqrt{(X_1 - X_2)^2 + (Y_1 - Y_2)^2}$$

With more than two units of analysis (e.g., people), a distance matrix is applied. The off-diagonal of the matrix represents the distance between respective units of analysis. Zeros are displayed on the diagonal to indicate that the distance from an object to itself is zero.

Chebychev distance is commonly used when data are ordinal in scale. This is calculated as the maximum of the absolute value of differences and expressed in notation as follows:

$$d_{\text{Chebychev}} = \max(|X_1 - X_2|, |Y_1 - Y_2|,)$$

8.3 COMPUTING CLUSTER ANALYSIS USING SPSS

Next, we consider SPSS software for computing cluster analysis. Before we conduct the analysis, let us talk about the data. The data we are using is team statistics of Division I-A Football Bowl Subdivision (FBS) obtained from ESPN during January 2016 ($n = 128$; FBS_2015.sav) (http://espn.go.com/college-football/statistics/team/_stat/total/sort/totalYards).

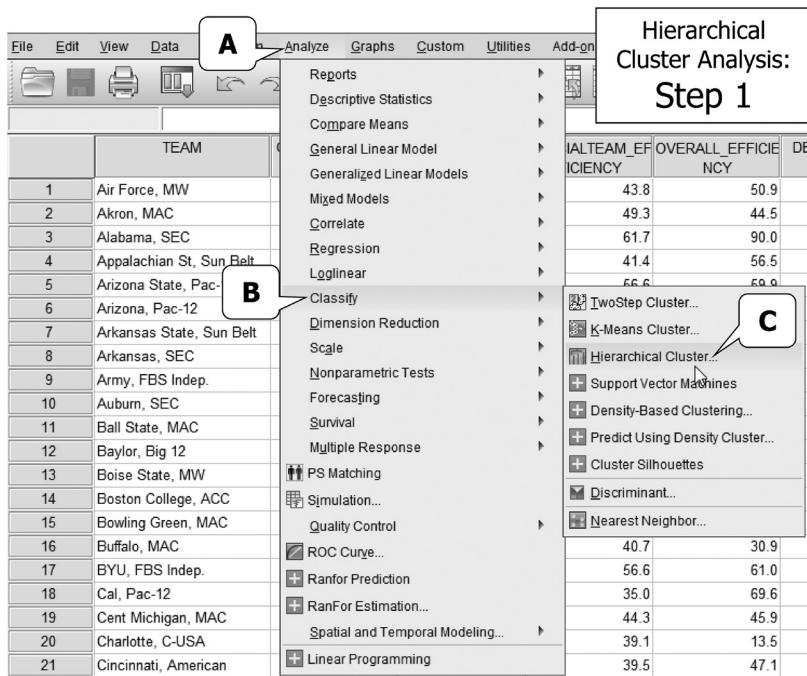
Let's look at the data. For the cluster analysis illustration, we will work with the following variables. These are all measured on a continuous scale. Understand that the inclusion of this number of variables is for illustrative purposes and does not adhere to the recommended condition of sample size of at least 2^m or a sample size of more than 1,000 (where m is the number of variables).

1. Offense efficiency
2. Defense efficiency
3. Special teams efficiency
4. Overall efficiency
5. Defense yards per game
6. Defense passing yards per game
7. Defense rushing yards per game
8. Offense yards per game
9. Offense passing yards per game
10. Offense rushing yards per game

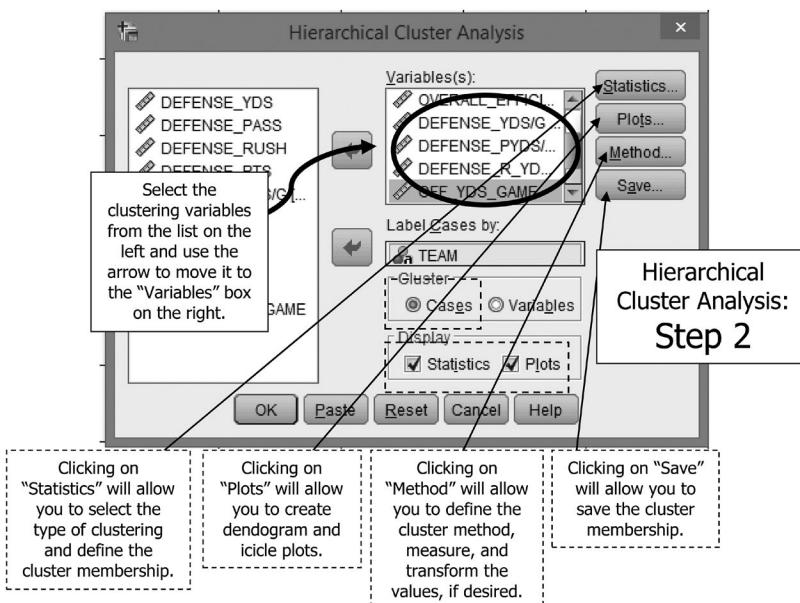
Additional variables are included in the data file just in case you are interested in trying out a slightly different model, as well as testing the extent to which the profiles differ on other variables. Each row in the data set represents one Division I-A football team. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the teams were measured (not all variables are presented in the screenshot).

	TEAM	OFF_EFFICIENCY	DEFENSE_EFFICIENCY	SPECIALTEAM_EFFICIENCY	OVERALL_EFFICIENCY	DEFENSE_YDS	DEFENSE_YDSD	DEFENSE_PASS
1	Air Force, MW	59.0	43.8	43.8	50.9	4382.0	337.0	2475.0
2	Akron, MAC	32.6	59.1	49.3	44.5	3943.0	329.0	2865.0
3	Alabama, SEC	74.7	91.6	61.7	90.0	3356.0	258.0	2394.0
4	Appalachian St, Sun Belt	51.3	62.1	41.4	56.5	3816.0	318.0	2213.0
5	Arizona State, Pac-12	59.6	54.1	56.6	59.9	5353.0	446.0	3860.0
6	Arizona, Pac-12	65.6	35.5	54.4	50.8	5560.0	463.0	3298.0
7	Arkansas State, Sun Belt	42.2	60.1	59.6	53.5	4763.0	397.0	3023.0
8	Arkansas, SEC	88.4	44.3	49.9	74.9	4849.0	404.0	3414.0
9	Army, FBS Indep.	26.7	26.5	41.9	20.9	4507.0	376.0	2525.0
10	Auburn, SEC	59.0	63.7	56.5	65.4	5062.0	422.0	2786.0
11	Ball State, MAC	26.8	20.3	61.3	19.7	6216.0	518.0	3510.0
12	Baylor, Big 12	82.7	67.3	35.3	79.2	4668.0	389.0	2792.0
13	Boise State, MW	52.8	66.6	48.5	60.9	4105.0	342.0	2693.0
14	Boston College, ACC	13.1	84.1	40.2	45.5	3052.0	254.0	2058.0

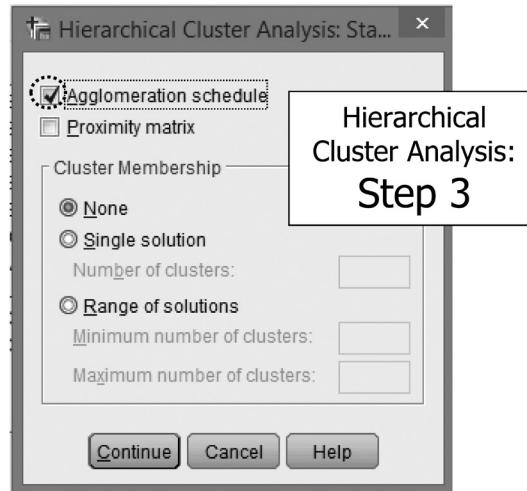
Step 1. To conduct cluster analysis, go to “Analyze” in the top pull-down menu, then select “Classify,” and then select “Hierarchical Cluster.” Following the screenshot below (see screenshot Step 1) produces the “Hierarchical Cluster Analysis” dialog box.



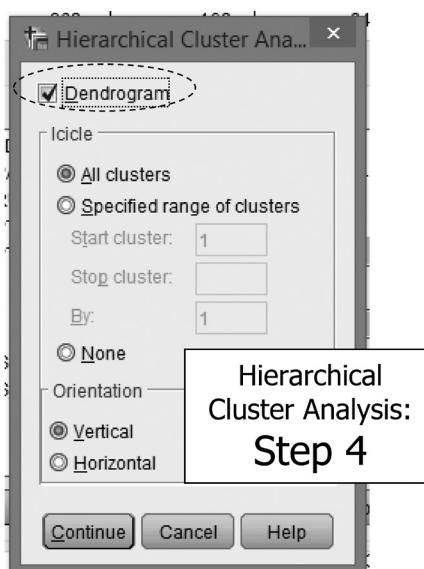
Step 2. Click the clustering variables and move into the “Variables” box by clicking the arrow button (see screenshot Step 2). We will leave the default settings to cluster by “Cases” and to display both “Statistics” and “Plots.” Move the variable “Team” into the dialog box to “Label Cases by,” so that our graphs will be displayed by team name rather than generic ID number.



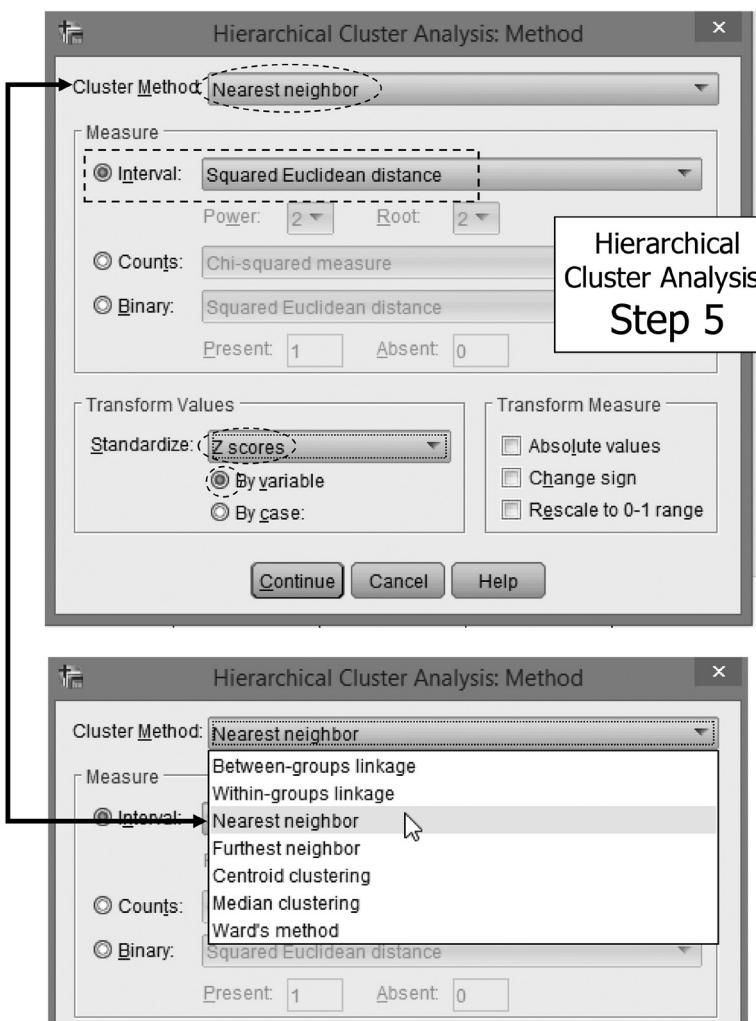
Step 3. From the “Hierarchical Cluster Analysis” dialog box (see screenshot Step 2), clicking on “Statistics” will provide the option to select how the objects are combined, either by the agglomeration schedule or by proximity matrix (i.e., distance matrix) (see screenshot Step 3). From this dialog box, the number or range of clusters to retain can also be defined. At this step, because neither of these are known, we select only “Agglomeration schedule.” Click on “Continue” to return to the “Hierarchical Cluster Analysis” dialog box.



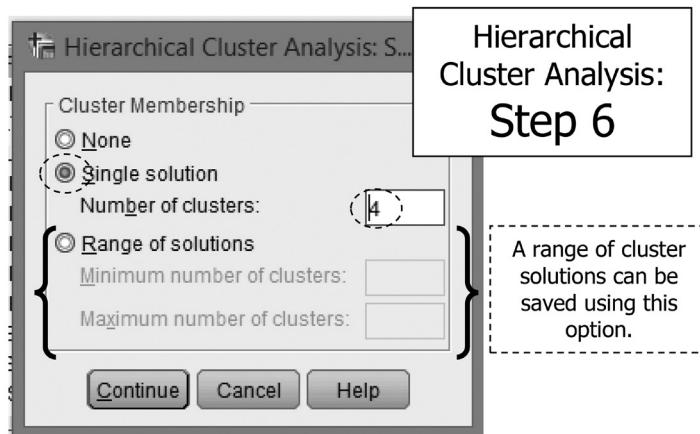
Step 4. From the “Hierarchical Cluster Analysis” dialog box (see screenshot Step 2), clicking on “Plots” will provide the option to create different types of plots (see screenshot Step 4). We will select a dendrogram and an icicle plot with all clusters. Click on “Continue” to return to the “Hierarchical Cluster Analysis” dialog box.



Step 5. From the “Hierarchical Cluster Analysis” dialog box (see screenshot Step 2), click on “Method.” For this illustration, we will select the “Nearest neighbor” clustering algorithm, which is a single linkage method (see screenshot Step 5). Even though single linkage works best when clusters are in long chains, single linkage methods are commonly applied first as they are helpful in identifying outliers. (Later, we will also apply the furthest neighbor clustering algorithm to compare results among the different methods.) Our data can be categorized as ratio in scale. Given that measurement scale, we will select to use squared Euclidean distance, which minimizes the importance of small distances and maximizes the importance of large distances. We will also transform by Z score, by variable. Click on “Continue” to return to the “Hierarchical Cluster Analysis” dialog box. We are not yet ready to save our membership, the last button available from the “Hierarchical Cluster Analysis” dialog box (see screenshot Step 2), so this will conclude selections in SPSS. Click on “OK” to generate the output.



Step 6. This is the last step, and *it should only be considered after the results have been examined. In other words, Step 6 would not normally be selected until after the initial analysis was attempted.* Once the number of clusters is determined, save the cluster solution. From the “Hierarchical Cluster Analysis” dialog box (see screenshot Step 2), click on “Save.” **For this illustration, we will save a single solution of 4 clusters based on the furthest neighbor method** (see screenshot Step 6). (*Thus, if you are following along and generating the cluster solution, you will want to go back to screenshot Step 5 and select ‘furthest neighbor’ as the cluster method.*) Had we desired to examine a range of cluster solutions (e.g., 3–5 clusters), we could have specified that here as well.



8.3.1 Checking Our Results From Step 6

We see that a new variable has been created that specifies the cluster to which each team was assigned.

OFF_POINTS	OFF_POINTS_GA ME	CLU4_1	
473.0	33.8	1	
311.0	23.9	2	
526.0	35.1	1	
477.0	36.7	1	
486.0	37.4	3	
448.0	34.5	3	
520.0	40.0	4	
467.0	35.9	3	
265.0	22.1	4	
357.0	27.5	4	

We can now use this cluster variable to determine differences between teams in each profile using, for example, ANOVA. In this example, of the four clusters, most are in cluster 4 (about 48%) and only about 5% are in cluster 2.

Complete Linkage (furthest neighbor)

	Frequency	Percent	Valid Percent	Cumulative Percent
Valid 1	34	26.6	26.6	26.6
2	6	4.7	4.7	31.3
3	26	20.3	20.3	51.6
4	62	48.4	48.4	100.0
Total	128	100.0	100.0	

The teams in cluster 1, for example, include the following:

Air Force, MW	Notre Dame, FBS Indep.
Alabama, SEC	Oklahoma, Big 12
Appalachian St, Sun Belt	Ole Miss, SEC
Baylor, Big 12	OSU, Big Ten
Clemson, ACC	Pitt, ACC
Florida, SEC	San Diego State, MW
FSU, ACC	Stanford, Pac-12
Ga Southern, Sun Belt	TCU, Big 12
Georgia, SEC	Temple, American
Iowa, Big Ten	Tennessee, SEC
Louisville, ACC	Texas A&M, SEC
LSU, SEC	UNC, ACC
Marshall, C-USA	USF, American
Miami (FL), ACC	Utah, Pac-12
Michigan, Big Ten	VT, ACC
Navy, American	Washington, Pac-12
Northwestern, Big Ten	Wisconsin, Big Ten

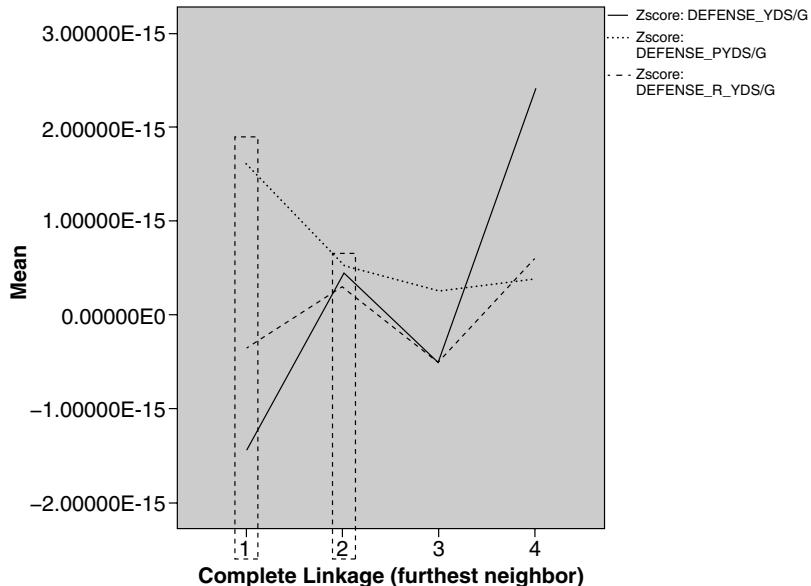
Now, let's standardize the variables that were used to create the cluster analytic solution and create line graphs with them. Standardizing can be done by going to "Analyze" then "Descriptive" then "Descriptive statistics" and placing a checkmark in the 'save standardized values as variables' box. Line graphs can be generated by following line graph screenshots Steps 1 and 2.

The image displays two overlapping SPSS dialog boxes. The top dialog is titled "Line Charts" and shows three chart types: Simple, Multiple, and Drop-line. The "Multiple" option is selected. A title box on the right says "Line Graph: Step 1". Below the chart types, a section labeled "Data in Chart Are" contains three radio button options: "Summaries for groups of cases" (selected), "Summaries of separate variables", and "Values of individual cases". At the bottom are "Define", "Cancel", and "Help" buttons.

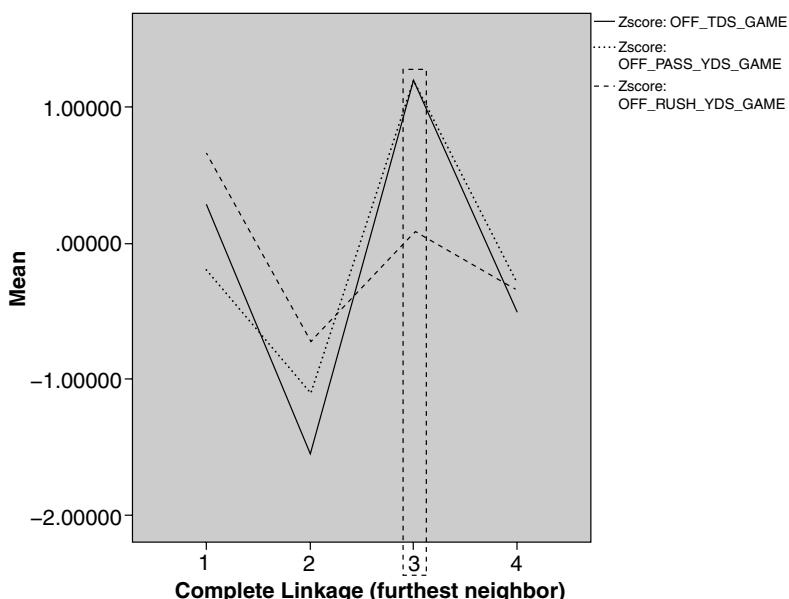
The bottom dialog is titled "Define Multiple Line: Summaries of Separate Variables". It lists variables on the left: ID, TEAM, OFF_EFFICIENCY, DEFENSE_EFFI..., SPECIALTEAM_E..., OVERALL_EFFIC..., DEFENSE_YDS, DEFENSE_YDS/..., DEFENSE_PASS, DEFENSE_PYDS..., DEFENSE_RUSH, DEFENSE_R_YD..., DEFENSE_PTRS, DEFENSE_PTRS/..., OFF_YDS, OFF_YDS_GAME, OFF_PASS, OFF_PASS_YDS..., OFF_RUSH, OFF_RUSH_YDS..., OFF_POINTS, OFF_POINTS_G... . On the right, under "Lines Represent", there are three entries: MEAN(Zscore: DEFENSE_YD..., MEAN(Zscore: DEFENSE_PY..., and MEAN(Zscore: DEFENSE_R_Y...). Under "Category Axis", it says "Complete Linkage (furthest nei...)" and "Panel by". A callout bubble points to the "Category Axis" field with the text: "For the category axis, enter the cluster solution variable that was saved in the cluster screenshot Step 6." At the bottom are "OK", "Paste", "Reset", "Cancel", and "Help" buttons.

Graphing the standardized variables in line graphs, we can better understand how the clusters compare and contrast, and after examining the cluster solution, you will want to identify the profiles (i.e., clusters) by name—similar to naming the constructs identified in factor analysis. As a starting point, let's look at the defense statistics. We see, for

example, that cluster 2 is the most similar on all three defense statistics, and the statistics of those teams hover near the mean. Thus, this cluster could be categorized as ‘average defense.’ In comparison, cluster 1 could be categorized as ‘high mean defense passing yard, moderate mean defense rushing yard, and low mean defense yards/game.’



Similarly, looking at the offense statistics, we can categorize cluster 3 as ‘high offensive and passing yards per game and moderate offensive rushing yards per game.’



Though not presented, a similar graph could be provided for efficiency statistics.

Interpreting the output. Annotated results are presented in Table 8.1.

■ TABLE 8.1

SPSS Results for the Cluster Analysis Example

Agglomeration Schedule						
Stage	Cluster Combined		Coefficients	Stage Cluster First Appears		Next Stage
	Cluster 1	Cluster 2		Cluster 1	Cluster 2	
1	99	122	.614	0	0	52
2	94	118	.620	0	0	21
3	96	123	.730	0	0	4
4	87	96	.758	0	3	10
5	48	59	.990	0	0	63
6	24	27	1.021	0	0	40
7	91	113	1.090	0	0	17
8	92	108	1.128	0	0	62
9	71	121	1.140	0	0	13
10	87	106	1.174	4	0	17
11	10	22	1.231	0	0	16
12	117	127	1.291	0	0	24
13	71	76	1.333	9	0	19
14	93	97	1.367	0	0	15
15	93	103	1.374	14	0	22
110	19	25	3.622	0	0	113
111	1	20	3.649	109	0	112
112	1	37	3.656	111	0	114
113	5	19	3.691	0	110	114
114	1	5	3.760	112	113	115
115	1	83	3.862	114	0	116
116			4.072	115	0	117
117			4.128	116	103	118
118			4.338	117	0	119
119			4.543	118	0	120
120			4.622	119	84	121
121			4.691	120	0	122
122			4.701	121	0	123
123			4.801	122	0	124
124			5.590	123	0	125
125	1	60	5.702	124	0	126
126	1	56	8.693	125	0	127
127	1	4	11.679	126	0	0

This example is based on **nearest neighbor method**.

For this example, a **three cluster solution** is reasonable, as there are two clusters after stage 125.

The column labeled 'cluster combined' in the agglomeration schedule displays the cases that are combined at each stage. The 'coefficient' is the distance at which the merger was made, and this value depends on both the proximity measure and the linkage method applied. For example, cases 99 and 122 are combined at stage 1 at a squared Euclidean distance of .614. This cluster is named 99' as the cluster number is labeled as the smallest case number in the respective cluster.

The column labeled 'next stage' indicates the stage at which the cluster will appear next. We see that case 96, which appears in the third stage, appears again in stage 4.

For brevity, only the first and last portions of the schedule are presented.

Reading the agglomeration schedule, we see that cases 99 and 122 are combined in a cluster at stage 1, and by naming conventions of the clusters where the smallest value represents the name of the cluster, this cluster is named '99.' At stage 2, cases 94 and 118 are combined to become cluster '94.' Jump to step 4, case 87 is added to the previous cluster '96' which now includes cases 96, 123, and 87.

TABLE 8.1 (continued)

PSS Results for the Cluster Analysis Example

A substantial jump in the distance coefficient or drop in the similarity coefficient indicates a good cluster solution has been reached. The solution prior to the jump indicates a reasonable solution. In this illustration, the only large jump is at the final inclusion of cases, specifically between stages 125 and 126 and 127. Thus, the stage *before* the sudden jump, at stage 125 in the table, indicates the optimal stopping point at which to merge clusters. Given there are two stages *after* this optimal point, a three cluster solution is suggested.

The **icicle plot** is read from bottom to top and provides a visual of the agglomeration schedule. The cases are on the horizontal axis, and the number of clusters are on the vertical axis. Cases that have joined to form a cluster are represented by contiguous bars. White space between bars represents boundaries between clusters. *For brevity, only a portion of the plot is presented.*

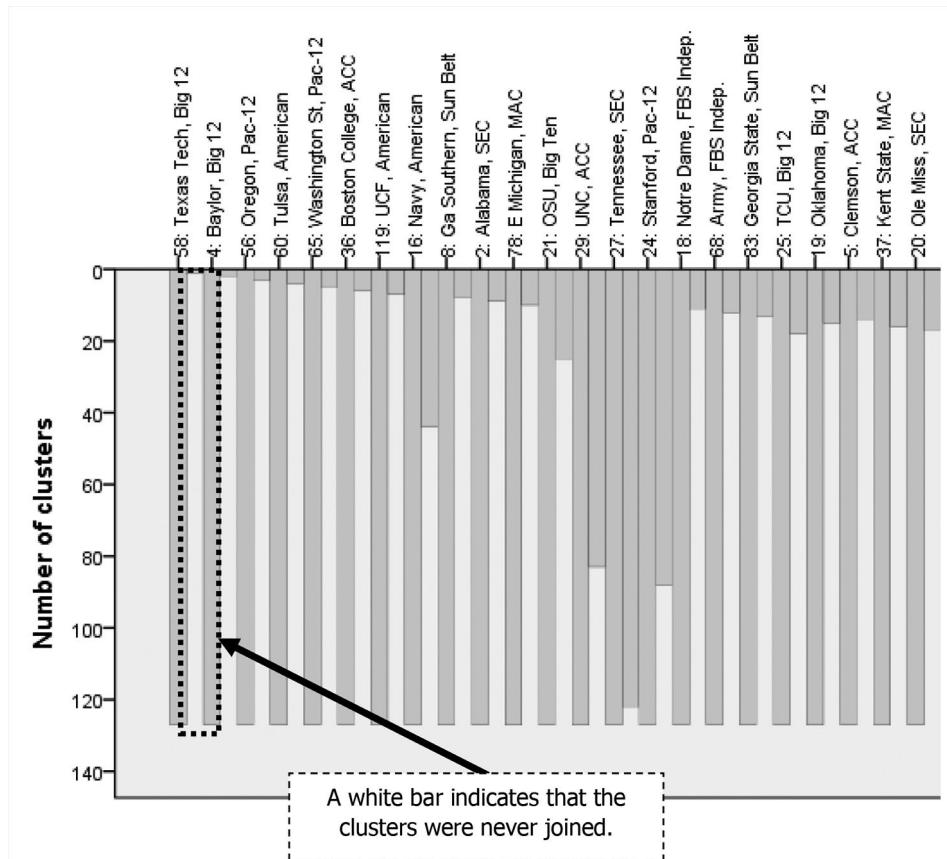


TABLE 8.1 (continued)

PSS Results for the Cluster Analysis Example

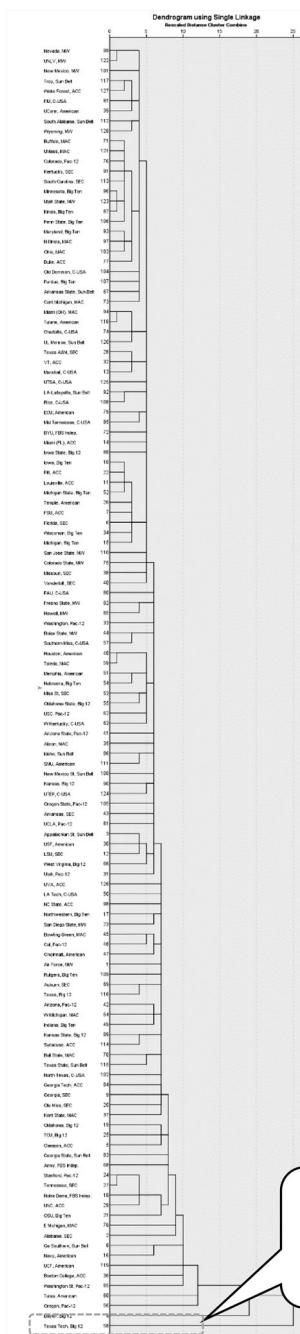
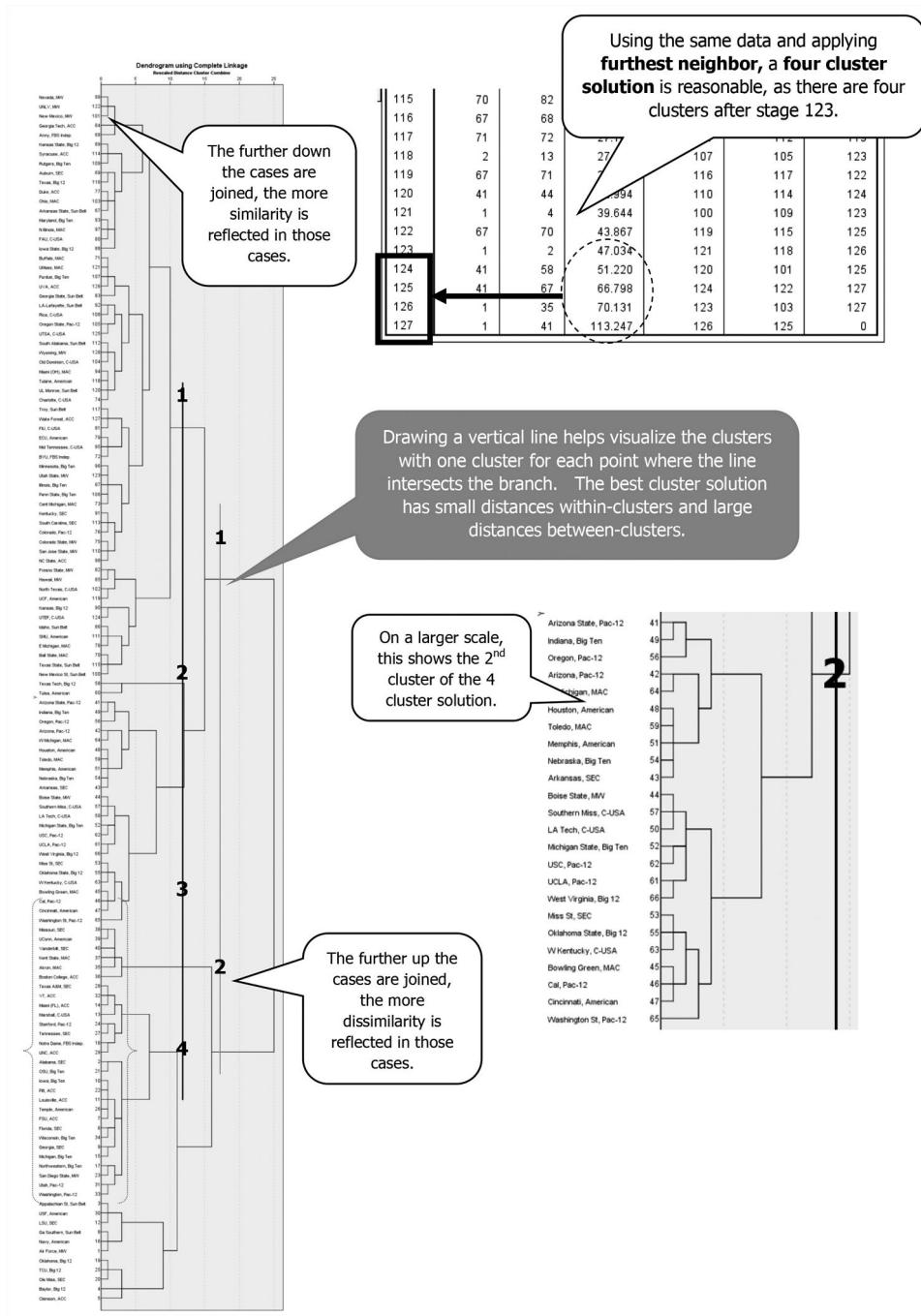


TABLE 8.1 (continued)

PSS Results for the Cluster Analysis Example



8.4 DATA SCREENING

As stated previously, multicollinearity is the main problem that can occur with cluster analysis. Thus, examining bivariate correlations prior to generating the cluster analysis can be used to determine if there are multiple variables with high shared variance (e.g., $r > .90$). In those cases, removing one of the variables that is sharing a substantial portion of the variance is helpful. In this example, the highest correlation is around .86, so we can be comfortable with having met the assumption of noncollinear data.

Correlations										
	OFF_EFFICIENCY	DEFENSE_EFFICIENCY	SPECIALTEAM_EFFICIENCY	OVERALL_EFFICIENT	DEFENSE_YDSIG	DEFENSE_PYDSIG	DEFENSE_RYDSIG	OFF_YDS_GAME	OFF_PASS_YDS_GAME	OFF_RUSH_YDS_GAME
OFF_EFFICIENCY	Pearson Correlation	1	.417**	.174*	.862**	-.192*	.044	-.306**	.797**	.509
	Sig. (2-tailed)		.000	.050	.000	.030	.621	.000	.000	.000
	N	128	128	128	128	128	128	128	128	128
DEFENSE_EFFICIENCY	Pearson Correlation	.417**	1	.088	.809**	-.778**	-.428**	-.747**	.165	.017
	Sig. (2-tailed)		.000	.324	.000	.000	.000	.000	.063	.847
	N	128	128	128	128	128	128	128	128	128
SPECIALTEAM_EFFICIENCY	Pearson Correlation	.174*	.088	1	.254*	-.055	.009	-.085	.011	-.104
	Sig. (2-tailed)		.050	.324		.004	.539	.917	.341	.803
	N	128	128	128	128	128	128	128	128	128
OVERALL_EFFICIENT	Pearson Correlation	.862**	.809**	.254*	1	-.542**	-.197	-.604**	.593**	.312*
	Sig. (2-tailed)		.000	.000	.004		.026	.000	.000	.000
	N	128	128	128	128	128	128	128	128	128
DEFENSE_YDS/G	Pearson Correlation	-.192*	-.778**	-.055	-.542**	1	.720**	.820**	.074	.200
	Sig. (2-tailed)		.030	.000	.539	.000		.000	.408	.024
	N	128	128	128	128	128	128	128	128	128
DEFENSE_PYDSIG	Pearson Correlation	.044	-.428**	.009	-.197	.720**	1	.193*	.236**	.243*
	Sig. (2-tailed)		.621	.000	.917	.026	.000		.029	.006
	N	128	128	128	128	128	128	128	128	128
DEFENSE_RYDS/G	Pearson Correlation	-.306**	-.747**	-.085	-.604**	.820**	.193*	1	-.090	.083
	Sig. (2-tailed)		.000	.000	.341	.000	.029		.311	.352
	N	128	128	128	128	128	128	128	128	128
OFF_YDS_GAME	Pearson Correlation	.797**	.165	.011	.583	.074	.236**	-.080	1	.731*
	Sig. (2-tailed)		.000	.063	.903	.000	.408	.007	.311	.000
	N	128	128	128	128	128	128	128	128	128
OFF_PASS_YDS_GAME	Pearson Correlation	.509**	.017	-.104	.312*	.200	.243**	.083	.731**	1
	Sig. (2-tailed)		.000	.847	.243	.000	.024	.006	.352	.000
	N	128	128	128	128	128	128	128	128	128
OFF_RUSH_YDS_GAME	Pearson Correlation	.453**	.207	.152	.409*	-.157	.015	-.234**	.448**	-.283*
	Sig. (2-tailed)		.000	.019	.087	.000	.077	.863	.008	.000
	N	128	128	128	128	128	128	128	128	128

**. Correlation is significant at the 0.01 level (2-tailed).

*. Correlation is significant at the 0.05 level (2-tailed).

8.5 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example write-up for the results of the cluster analysis. Recall that the stats lab grad assistants, Challie Lenge and Oso Wyse, were assisting Dr. B. Bryant in exploring team statistics of football programs at Division I-A institutions using cluster analysis. The research question examined was: What are the underlying profiles of Division I-A football teams?

The team then assisted Dr. Bryant in conducting cluster analysis, and a template for writing the cluster analysis results follows:

What are the underlying profiles for [variable set]?

It may be helpful to preface the results of the cluster analysis with information on an examination of the extent to which the data were thoroughly screened.

Prior to conducting the analyses, bivariate correlations were generated to determine the extent to which multicollinearity may have been present. The

strongest correlation between variables was approximately .86, below the threshold of .90 which may be an indication of multicollinearity.

Here is an example summary of results for the cluster analysis (remember that this will be prefaced by a section reporting the extent to which the data were thoroughly screened).

A hierarchical cluster analysis was computed to determine the underlying profiles of Division I-A football programs using 10 football team statistics including [list variables here]. Prior to clustering, all variables were standardized with a mean of zero and standard deviation of one. The final solution employed the furthest neighbor method using squared Euclidean distance, although the nearest neighbor method was used to compare cluster solutions. Both the agglomeration schedule and the dendrogram suggested that a four-cluster solution was reasonable. This solution resulted in four distinct profiles categorized as the following: [include the names of the profiles, based on how the means are differentiated by cluster].

At this point, you will want to include descriptive statistics of the clusters, such as a frequency distribution table with counts and percentages of teams per cluster and graphs of means. Then you will want to identify the profiles by name and clearly indicate how they are distinctly different from each other. It is likely that the cluster solution is only the first step. Examining the profile membership, as an independent variable, in further analyses is anticipated (though not presented here).

PROBLEMS

Conceptual Problems

1. Conceptually, cluster analysis is akin to which one of the following statistical procedures but uses units rather than variables for the clustering?
 - a. Analysis of variance
 - b. Exploratory factor analysis
 - c. Multilevel modeling
 - d. Regression
2. Agglomerative clustering algorithms do NOT include which one of the following?
 - a. Average linkage
 - b. Centroid method
 - c. Complete linkage
 - d. k means
3. Nearest neighbor clustering algorithm is also referred to as which one of the following?
 - a. Complete linkage
 - b. Sequential threshold
 - c. Single linkage
 - d. Ward's method

4. The delineation between clustering algorithms is based on a distance measure.
 - a. True
 - b. False
5. Which method begins, rather than ends, with the researcher defining the number of clusters?
 - a. Agglomerative methods
 - b. Divisive methods
 - c. Hierarchical methods
 - d. Nonhierarchical methods
6. When is the application of distance measures in cluster analysis appropriate?
 - a. When a matching coefficient is computed
 - b. When the clustering variables are at least interval in scale
 - c. When the solution will be cross-validated
 - d. When the variance ratio criterion will be computed
7. Identification of a cluster seed is usually the starting point for nonhierarchical clustering procedures.
 - a. True
 - b. False
8. The clustering algorithm that minimizes the total sum of squared distances is which one of the following?
 - a. Complete linkage
 - b. Sequential threshold
 - c. Single linkage
 - d. Ward's method
9. Determination of how clusters are formed is based on which one of the following?
 - a. Agglomerative clustering
 - b. Clustering algorithm
 - c. Dendrogram
 - d. Variable selection
10. Which one of the following methods is also known as a partitioning method?
 - a. Average linkage
 - b. Centroid method
 - c. Complete linkage
 - d. k means

Computational Problems

1. Using the FBS_2015.sav dataset, conduct cluster analysis with the following three variables: (1) defensive passing yards, (2) defensive rushing yards, and (3) defensive points. Use Ward's method and squared Euclidean distance. Standardize the variables using z scores. Determine the best cluster solution using the agglomeration schedule and the dendrogram and categorize the clusters based on means of the clustering variables.

2. Using the FBS_2015.sav dataset, conduct cluster analysis with the following three variables: (1) offensive passing yards, (2) offensive rushing yards, and (3) offensive points. Use Ward's method and squared Euclidean distance. Standardize the variables using z scores. Determine the best cluster solution using the agglomeration schedule and the dendrogram and categorize the clusters based on means of the clustering variables.

Interpretive Problem

1. Use SPSS to conduct cluster analysis with a different set of the Division 1-A football team statistics (FBS.sav). Apply at least two different clustering methods and compare the results. Write up the results.

REFERENCES

- Aroean, L., & Michaelidou, N. (2014). A taxonomy of mobile phone consumers: Insights for marketing managers. *Journal of Strategic Marketing*, 22(1), 73–89.
- Calinski, T., & Harabasz, J. (1974). A dendrite method for cluster analysis. *Common Statistical Theory*, 3(1), 1–27.
- Carlson, E., Profumo, S., & Weninger, C. (2013). Clustering analysis of the morphology of the 130 GeV gamma-ray feature. *Physical Review*, 88(4-A), 1–12.
- Chen, L., Lin, Z., Lin, G., Zhou, C., Chen, Y., Wang, X., & Zheng, Z. (2014). Classification of micro-vascular patterns via cluster analysis reveals their prognostic significance in glioblastoma. *Human Pathology*, 46, 120–128.
- Chiou, T., Fang, D., Chen, J., Wang, Y., & Jeris, C. (2001). *A robust and scalable clustering algorithm for mixed type attributes in large database environment*. Paper presented at the 7th Annual Association for Computing Machinery SIGKDD International Conference in Knowledge Discovery and Data Mining, San Francisco.
- Gomez, C., & Parigi, P. (2015). The regionalization of intergovernmental organization networks: A non-linear process. *Social Networks*, 43, 192–203.
- La Rocca, M., Stagliano, R., La Rocca, T., & Cariola, A. (2015). Investment cash flow sensitivity and financial constraint: A cluster analysis approach. *Applied Economics*, 47(41), 4442–4457.
- Maree, P. R., Allott, K. A., Killackey, E., Farhall, J., & Cotton, S. M. (2015). Exploring cognitive heterogeneity in first-episode psychosis: What cluster analysis can reveal. *Psychiatry Research*, 229(3), 819–827.
- Mooi, E., & Sarstedt, M. (2011). *A concise guide to market research*. Berlin Heidelberg: Springer-Verlag.
- Nunes, K. P., Toyota, R. G., Oliveira, P. M. S., Neves, E. G., Soares, E. A. A., & Munita, C. G. (2013). Preliminary compositional evidence of provenance of ceramics from Hatahara archaeological site, central Amazonia. *Journal of Chemistry*, 2013, 1–6.

Chapter 9

EXPLORATORY FACTOR ANALYSIS

CHAPTER OUTLINE

9.1 What Exploratory Factor Analysis Is and How It Works	363
9.1.1 Characteristics	364
9.1.2 Sample Size	377
9.1.3 Power	378
9.1.4 Effect Size	378
9.1.5 Assumptions	378
9.2 Mathematical Introduction Snapshot	381
9.3 Computing EFA Using SPSS	383
9.3.1 Computing EFA With Continuous Data Using SPSS	383
9.3.2 Computing EFA With Ordinal Data Using SPSS	412
9.4 Data Screening	427
9.4.1 Independence	427
9.4.2 Linearity	427
9.4.3 Absence of Outliers	428
9.4.4 Extreme Multicollinearity and Singularity	432
9.5 Research Question Template and Example Write-Up	433

KEY CONCEPTS

1. Factor
2. Eigenvalue
3. Communality

4. Factor extraction
5. Orthogonal rotation
6. Oblique rotation
7. Factor retention

Up to this point, we have generally concerned our analyses with procedures that have as the goal the examination of one or more a priori outcomes. With this chapter, we begin to deviate from this method of examination in that we are now doing just as the name of this procedure implies—exploring the data. Actually, some would say it is not even that but rather “it is reconnaissance” (Kaiser, 1970, p. 402). Rather than having one or more a priori outcomes, we are using exploratory factor analysis to reduce a large number of variables into identifiable clusters of variables to better understand the structure of the data.

Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying exploratory factor analysis, (b) determine and interpret the results of exploratory factor analysis, and (c) understand and evaluate how to screen data prior to conducting exploratory factor analysis.

9.1 WHAT EXPLORATORY FACTOR ANALYSIS IS AND HOW IT WORKS

As we visit the statistics lab today, we find that Addie Venture and Oso Wyse have been tasked with an exploration analysis of data.

As graduate student researchers in the stats lab, Addie and Oso have become quite accustomed to working with their teammates on data analyses that examine one or more outcomes of interest. Many times, these outcomes have been computed as composite variables from psychological assessments. While Addie and Oso have appreciated the ability to group together individual items to form various constructs, they had never really been concerned with the process underlying that construction—until today, that is. Dr. Wesley, a faculty member from the Higher Education program, is interested in examining the factor structure of measures of perceived use of skills at home and at the workplace for a select group of individuals who participated in the Survey of Adult Skills, a large data collection effort from the Organization for Economic Cooperation and Development’s Programme for the International Assessment of Adult Competencies (PIAAC). Addie and Oso suggest the following research question to Dr. Wesley: *What is the underlying factor structure for perceived use of skills at home and work?* Given that dimension reduction is the goal of the project, the team recommends exploratory factor analysis to answer Dr. Wesley’s question. Always up for adventure and armed with statistical knowledge, Addie and Oso are excited to embark on this task.

Globally, exploratory factor analysis (EFA) is a statistical procedure that allows us to cluster together variables into what we’ll refer to in this chapter as factors (but are also

known as constructs or latent constructs, a term we will use when we discuss confirmatory factor analysis). These variables may be, as just one example, a number of indices designed to measure general skills. Examining each of the variables individually may provide useful information simply by reviewing descriptive statistics of the individual measures. However, even *more* useful information may be provided by examining the underlying constructs from the variables, those variables that group together and make the number of measures parsimonious and more manageable. In essence, what exploratory factor analysis allows us to do is to work with all variables simultaneously, but at the same time know something about their underlying data structure. Exploratory factor analysis is therefore often used to provide evidence of construct validity. The underlying focus of factor analysis deals with finding common variance (distributed among the factors) and eliminating the unique variance that is not of interest (where total variance = common variance + specific variance + error variance).

Although confirmatory factor analysis will be introduced in detail in a later chapter, it is important to broach the topic here so there is a good understanding of when each is most appropriate. Confirmatory factor analysis is a statistical technique that can be used to identify the factor structure of observed variables and to test the hypothesis that a relationship exists between the respective observed variables and one or more underlying latent constructs. Additionally, much of the terminology and concepts that we will discuss in relation to EFA generalize to CFA. The titles of the procedures may give some indication of when one is more appropriate than the other is. By nature of exploration, EFA is appropriate when there is a lack of theory to dictate relationships between the variables. Brown refers to this as a “data-driven approach” (2006, p. 14). In comparison, CFA is appropriate when a strong theoretical base exists such that the relationships between variables are known and can be specified in the modeling process. In fact, it is very common for researchers to first conduct EFA prior to CFA so that there is a better understanding of how the items relate to each other and the underlying constructs or factors.

9.1.1 Characteristics

9.1.1.1 Principal Components Versus Exploratory Factor Analysis

Before we delve into this chapter, it is important to understand the difference between principal components analysis (PCA, sometimes also known as ‘component factor analysis’ or ‘component analysis’) and exploratory factor analysis (EFA, sometimes known as ‘common factor analysis’). There is a difference, although in reading published literature, it seems that many authors understand them to be used interchangeably (and they should not be). If your goal is to estimate underlying factors and attach some meaning to those factors (as a form of construct validity, for example), then EFA is required. PCA, on the other hand, can be used to estimate and understand the contributions of the variables to the linear components within the data, but PCA is simply a method of decomposition—a technique for data reduction only. As stated by Borsboom, “the extraction of a principal components structure, by itself, will not

ordinarily shed much light on the correspondence with a putative latent variable structure” (Borsboom, 2006, p. 426). If the interest is in placing substantive meaning on the factors extracted, EFA is the procedure needed. Throughout the chapter, it will be assumed that EFA is the goal. However, keep in mind that generating PCA or EFA is as simple as a toggle menu option in SPSS. While the results are mathematically different, the solutions you see *may* actually be quite similar. This is, again, one of those times when you must be a responsible researcher and understand the goal of your research (decomposition only, PCA, or extraction of meaning, EFA) so that you can select the appropriate method.

9.1.1.2 Exploratory Factor Analysis Specification Conditions and Decisions

There are a number of conditions that must be understood and decisions that must be made when selecting to use and implement either PCA or EFA. These are related to

- a. determining factorability
- b. fitting the factor model
- c. selecting the factor(s)
- d. rotating the factor solution

As we’ll learn, researchers must first determine if factor analysis is appropriate for both their research question and their data. Within factorability, we will discuss measurement scale, sample size, and sample homogeneity, followed by tools for determining initial factorability. Second, researchers must select procedures to fit the model and estimate the model parameters. Within this realm, factor extraction and factor rotation will be reviewed. Third, the number of common factors to specify when fitting the model has to be determined. Lastly, whether or not to rotate, and how to rotate if needed, must be determined.

9.1.1.3 Factorability

Measurement scale and sample homogeneity are important considerations for determining factorability. Sample size (discussed later in the chapter) is also a consideration. In this section, we will also discuss tools for determining initial factorability.

Measurement Scale of Variables

It is important to remember that factor analysis (PCA and EFA) has as the primary requirement that a correlation matrix (denoted in statistical terms as uppercase bold **R**, the input correlation matrix with unities—or 1.0—in the diagonal, which is also referred to as the unreduced correlation matrix) be calculated from the variables in the model. (Note that a covariance matrix can also be applied in EFA; interpretation tends to be much easier with a correlation matrix. The remainder of the chapter will focus

on a correlation matrix.) With conventional factor analysis, the computed correlation matrix is a Pearson matrix. This, therefore, suggests that the variables applied must be metric (at least interval in scale) so that a linear relationship exists between the variables. (However, this does not guarantee that linearity will be met.) A bit more will be added to this discussion as we talk about factor loadings later in the chapter. Even though one of the conditions of conventional factor analysis is measurement that is at least interval in scale, it is quite common to find factor analysis applied to Likert-type items which are ordinal in scale (e.g., five-point scale ranging from strongly agree to strongly disagree), particularly as the number of levels of the items increases. And should you find that your ordinal items meet the assumption of linearity, then proceeding with the factor analysis is fine (assuming other conditions and assumptions are satisfactorily met). However, items with small numbers of levels (less than seven categories in particular) are often not good candidates for conventional factor analysis, and the factors may be more difficult to interpret. Technically, binary (i.e., dichotomous) items can be factor analyzed with conventional methods, however the interpretation can be problematic as the results can reflect variation in the endorsement rate of the variables rather than the underlying construct (Fabrigar & Wegener, 2012). Categorical variables that have similar splits will tend to correlate even if the context of correlation of the variables doesn't make sense (see Gorsuch, 1983). This problem is augmented with binary data where correlations tend to reflect similar 'difficulty' as evidenced in a testing type of environment. If you do decide to proceed with conventional factor analysis using categorical variables, the factor loadings should be examined with extreme care to determine if they reflect 'difficulty' (where difficulty is defined as approximately the proportion of individuals with a '1' for their item score, as opposed to a '0') as compared to a substantive relationship. The use of binary data in conventional factor analysis can also result in a factor solution with too many factors. In the case of categorical variables, dichotomous in particular, it is highly recommended that a specialized factor analytic program that is designed for that type of data be applied to it. Later in this chapter, SPSS categorical principal components analysis (CATPCA), an add-on in SPSS, will be used to illustrate the application of ordinal data with factor analysis.

Homogeneity of the Sample in Relation to the Underlying Factor Structure

An important condition of factor analysis is that the sample of cases from which the variables were measured must be homogenous in respect to the underlying factor structure. In other words, if your collective sample of cases is known to differ, based on some characteristic, on the set of variables for which you are factor analyzing, then separate factor analysis should be performed for the groups that are anticipated to differ. For example, say that all employees of a company have been surveyed about their perceptions of the work environment, and previous empirical research suggests that those in management positions have different perceptions as compared to nonmanagement positions. The factor analysis should be conducted separately for those groups (i.e., management and nonmanagement) that are expected to differ.

Initial Factorability Assessment

There are a number of indices that should be reviewed prior to conducting the factor analysis that will help you gauge the extent to which the variables and the matrices produced from them are factorable. These include (1) correlation coefficient values, (2) Bartlett's test of sphericity, (3) anti-image correlation matrix, and (4) Kaiser-Meyer-Olkin measure of sampling adequacy.

Correlation coefficient values between the variables being factor analyzed should be .30 (in absolute value terms) or greater. This will ensure sufficient relationships to justify examination of the potential underlying components. Correlations lower than .30 may be due to low variance, which can result when samples are homogenous (but does not necessarily imply homogeneity in the sample). (However, correlations of more complex scores, such as difference scores, may have correlations between .20 and .30 and still have variables that are extremely factorable.) If there are correlation coefficient values that are not satisfactory and that are not theoretically critical, remove the variable with the lowest individual correlation value and rerun (doing so until, collectively, the correlation values reach what you deem acceptable).

Bartlett's test of sphericity is conducted to determine if the observed correlation matrix is statistically significantly different from an identity matrix (i.e., diagonal elements are 1 and off-diagonal elements are 0). Statistically significant results for Bartlett's test are desirable, as they allow you to reject the null hypothesis, which states that the observed correlation matrix equals the identity matrix. We want to see redundant variance, overlapping variance among variables, in order to reduce the variables into a fewer number of latent factors, and this is accomplished with a statistically significant Bartlett's test. Should the null hypothesis not be rejected, this provides evidence that the correlation matrix produced from the variables cannot be factor analyzed.

Kaiser-Meyer-Olkin (KMO) measure of sampling adequacy (MSA) is an index of shared variance in the variables and compares the magnitudes of the observed to those of the partial correlation coefficients. MSA values range from zero to one, and large values are another form of evidence to suggest that the variables are factorable. In its early origination, Kaiser (Kaiser & Rice, 1974, p. 112) proposed the following guidelines for interpreting the index: in the .90s = marvelous; in the .80s = meritorious; in the .70s = middling; in the .60s = mediocre; in the .50s = miserable; below .50 = unacceptable. As we'll see when we compute our factor analysis, an MSA for each individual item and an overall MSA will be generated. If the overall MSA is not satisfactory, remove the variable with the lowest individual MSA value and rerun (doing so until the MSA value reaches what you deem acceptable). The overall KMO-MSA numerator is the sum of squared correlations of all variables, and the denominator is the numerator value (i.e., the sum of squared correlations of all variables) plus the sum of squared partial correlations of each variable i with each variable j , controlling for the other variables. The idea behind the MSA is that the partial correlations (reflected in the denominator) should not be relatively small if one is to expect distinct factors to

emerge from factor analysis (i.e., creating a small denominator that will then provide for a larger MSA value). The size of the MSA can therefore be expected to increase as the following increase—sample size, average correlation, number of variables—and as the number of factors decrease. In SPSS, the MSA values are provided on the diagonal of the anti-image correlation matrix.

9.1.1.4 Fitting the Factor Model

Factor Extraction

Assuming you have made it through the previous examination and have determined that factor analysis is appropriate for your data, the next level of decisions has to deal with implementing or actually computing the factor analysis. Beginning with this section, we will discuss a number of concepts and procedures that should be understood to fit the model appropriately. Fitting the model, in reference to factor analysis, is also known as factor extraction. Although many times the algorithms will produce similar results, this is not always the case. Therefore, understanding how they operate and situations where they are most effective is needed.

Factor analytic models that generate two or more factors will have an infinite number of ways that the factors can be oriented in multidimensional space, each with an equally best-fitting solution (Fabrigar & Wegener, 2012). Let's say that we have a factor model where two factors are suggested. If we think about our items in two-dimensional space, the axes represent the factors and the space between the individual observed variables represents their intercorrelations—variables closer together have stronger relationships with each other. This implies that one single unique best-fitting solution does not exist when the model generates more than one factor. Therefore, this puts the burden on the researchers to select one solution. This decision process of factors to retain is the factor extraction process. Good model fit is achieved when the mathematical model for converting physical distance into predicted correlations between variables is similar to the correlations among observed variables.

A number of different algorithms can be used to fit factor analytic models, all of which calculate orthogonal factors that combine to reproduce the correlation matrix. Our discussion will focus on a few of the most common. Those commonly found in standard statistical software include principal components, unweighted least squares, generalized (weighted) least squares, maximum likelihood, principal axis factoring, alpha factoring, and image factoring. Of these, principal components, principal axis, and maximum likelihood are likely the most common and are those on which our discussion will focus. In addition to its common use in EFA, maximum likelihood is also the most commonly applied estimation method in CFA (Brown, 2006).

Which extraction method selected is the researcher's choice. Generally, any extraction method will require rotation in order for the solution to be interpretable. The

solutions from the different extraction methods will converge in situations where you have a large number of cases and variables and communality estimates that are similar. Evidence of the stability of your factor solution can be seen in cases where there is convergence of factor analytic solutions when using different extraction methods. While applying every estimation procedure to your data would be akin to a fishing expedition, it is quite common to select a small handful of estimation techniques to test the stability of your factor analytic model under different estimation methods—for example, first applying principal axis factoring then proceeding with maximum likelihood and ceasing the analyses when a sound solution is achieved. Generating factor analysis using two different estimation methods has been recommended (Child, 2006). Should the solutions result in discrepancies, an attempt to determine the reason(s) for the discrepancies is appropriate, followed by generation of the factor model with a third estimation technique (Child, 2006).

Principal Components

We have already broached the topic of principal components analysis as compared to common factor analysis, thus we will not delve further into that difference other than to mention a few notables as it relates to how the data is extracted. In a nutshell, the *variance* is analyzed in PCA whereas the *covariance* (communality) is analyzed in common factor analysis. In PCA, the goal is to extract the most variance from the variables with each factor.

Unweighted and Generalized (Weighted) Least Squares

Both unweighted and weighted least squares methods of factor extraction attempt to minimize the squared differences between the observed and reproduced (off-diagonal) correlation matrices. The difference between the two is that variables that share substantial variance with other variables are weighted more heavily, and variables that have more unique variance (i.e., less shared variance) receive less weight. The heavily weighted items thus contribute more to the solution than the items with lesser weight.

Maximum Likelihood (ML)

Maximum likelihood estimation calculates factor loadings that maximize the probability that the observed correlation matrix would be sampled from the population. ML is the most statistically advanced extraction method and one of the most commonly applied.

Principal Axis Factoring

Principal axis factoring has communality estimates, which are estimated through an iterative process, in the diagonal of the correlation matrix. The goal of principal axis factoring is to extract maximum variance from the variables with each factor, and this makes principal axis factoring less desirable in some situations as compared to

other extraction methods that can be more effective in reproducing the correlation matrix. Principal axis factoring is one of the most commonly applied extraction methods (Child, 2006).

Alpha Factoring

Alpha factoring uses an iterative procedure to estimate communalities that then maximize coefficient alpha (i.e., an index of reliability). Unlike score reliability in psychometric research (i.e., consistency of subjects), alpha factoring focuses on determining consistency of *variables*, in other words, extracting factors that are consistently found when repeated samples of *variables* (not subjects) are drawn from a population of variables (not subjects).

Image Factoring

Image factoring uses multiple regression, with each variable serving as the dependent variable and the remaining as the independent variables, to predict image scores that are then used to compute a covariance matrix. The communalities in this extraction method are the variances from the image score covariance matrix. Factor loadings represent covariance values (as compared to correlation values seen in the other estimation procedures) between the factors and variables.

Communalities

The communality, h^2 , interpreted as the *reliability of the variable*, measures the percent of variance (squared multiple correlation) of a given variable explained by all the factors jointly. The total communality is calculated by adding the squares of all the loadings of a variable across the common factors. It is the sum of all the common variance—the proportion of common variance within a variable. Computationally, the communality is the sum of squared factor loadings for a variable across all the factors.

A variable that has a low communality (.20 or below) has low common variance and high specific and error variance. A variable with a low communality may be a candidate for removal from the model, as this suggests that the factor model may not be working well for that variable. Low communalities across the set of variables indicate that the variables have weak relationships with each other. However, please note the following: A low communality can still be meaningful if the variable is contributing to a well-defined factor. The communality coefficient is not the critical element per se, but rather it is the extent to which the variable plays a role in the interpretation of the factor that is key.

It is also possible to have communalities that are too large. A communality that exceeds 1.0 is evidence of a spurious solution and may reflect a sample size that is too small or a factor model that has too few or too many factors. If you find yourself in this situation, and it is unfeasible to collect more data (either more cases and/or more variables), then

remove the variable with the largest communality and rerun, repeating this process until the communality estimates are less than one. It is important to note that communalities are unaffected by rotation but are impacted by extraction method, thus the only communalities provided in standard statistical software such as SPSS are the initial and extracted estimates. Extracted communalities represent the percent of variance in a given variable explained by the extracted factors, which are often fewer in number than all the possible factors, resulting in coefficients less than 1.0 (as a side note, the communalities will be less than one even initially, with exceptions noted previously). Assuming most of the common variance is contained within those extracted factors, then the unique variance can be calculated as $1 - h^2$.

9.1.1.5 Factor Retention

Once variables are factored, the researcher must determine how many factors to retain. While this may hold only a small fraction of this chapter, the number of factors to retain has been characterized as “the crucial decision” in the EFA process as when the optimal number of factors are retained, other EFA results will generally be similar (O’Connor, 2000, p. 396). When too few factors are extracted, important information is lost, potentially important factors are neglected, error in factor loadings increases, and other problematic issues arise (Zwick & Velicer, 1986). When too many factors are retained, factors are unnecessarily split resulting in low loadings and the attribution of importance to factors which really are not (Zwick & Velicer, 1986).

Theoretically, there are as many potential factors as there are variables. For example, in a case where 12 variables are being factor analyzed, theoretically, there are 12 factors. Obviously, a researcher would be ill-guided to retain that many factors, as the goal of factor analysis is parsimony (at least in respect to data reduction)—retention of the smallest number of factors that explains the most variance of the observed variables. Historically, the number of factors to retain from a factor analytic solution have relied more often on visual (and subjective) inspection and subjective rules rather than empirical evidence, and there is not one single tool recommended. Rather, multiple decision rules are recommended and deemed desirable, as is the application of more sophisticated factor retention strategies such as parallel analysis and bootstrapping (Thompson & Daniel, 1996). Despite this recommendation, much published research exists that does not adhere (e.g., Gaskin & Happell, 2014; Henson & Roberts, 2006). Never fear, by the end of the chapter you will have the skills to call yourself a sophisticated researcher!

Scree Plots

Scree plots, where the number of factors to retain is based on where the elbow bends in the plot, are a visual tool that can be used to decide on the number of factors to retain. We see an example of a scree plot in Table 9.4. The factor numbers are plotted on the X axis and the eigenvalues are on the Y axis. In interpreting the scree plot, we look for the clearest delineation where the line goes from being diagonal to being horizontal. Then, to determine the number of factors suggested by the scree plot, we count the number of

straight lines (not dots), stopping at the point where the line becomes more horizontal than diagonal. As with all visual tools, however, there is a certain degree of subjectivity that comes with making this decision. Even among experts, the reliability of scree plot interpretation is low (Streiner, 1998). When used as a decision rule to determine the number of factors to retain, scree plots generally perform better than the eigenvalue greater than one rule but are less accurate than parallel analysis (Zwick & Velicer, 1986).

Kaiser's Rule (Eigenvalues Greater Than One)

This rule is also known as the Unity Rule or the Kaiser-Guttman Criterion as it was proposed by Guttman and modified by Kaiser. Determining the number of factors to retain using Kaiser's Rule is quite simple—only those factors with eigenvalues greater than 1.0 are retained and factors with eigenvalues that are less than 1.0 are dropped. The value of one is the cut point given that the total variance contributed by each variable is one, and the variance of the factors retained should be greater than the contribution of only one variable. Eigenvalues, also known as characteristic roots or latent roots, are a measure of variance that are computed from the input (i.e., unreduced) correlation matrix. More specifically, eigenvalues measure the amount of variance in the total sample that is accounted for by each factor, and eigenvectors summarize this variance for the respective correlation or variance-covariance matrix (Brown, 2006). Factors with small eigenvalues suggest the respective factor is contributing little to explaining the variance in the variables.

Despite its widespread, and often sole, application to determining the number of factors to retain, the application of eigenvalues greater than one consistently misestimates the number of factors (either over- or underestimating) (Zwick & Velicer, 1982, 1986). Other criticisms are that an overestimation of the number of factors occurs when there are low communalities and a large number of variables and an underestimation of the number of factors to retain occurs when there are a small number of variables or when the sample size is very large. Kaiser's Rule tends to work best in conditions of moderate to large communalities, modest sample sizes, and 20–50 variables. Given these limiting conditions within which Kaiser's Rule tends to produce fairly accurate estimates, applying Kaiser's rule should only be done as a starting point (if at all) when generating your factor model. When appropriate, the results should be reviewed and the model recomputed based on a fixed number of factors.

Parallel Analysis

In comparison to the eigenvalue greater than one rule and visual examination of scree plots, there are statistically based procedures that exist for determining the number of factors to retain. Parallel analysis is one such procedure that is considered superior for determining optimal solutions for factor retention, and with 92% accuracy, has been considered the most accurate of the common methods used for retaining factors (including Kaiser's rule, Velicer's minimum average partial—MAP, scree plots, and Bartlett's test) (Zwick & Velicer, 1986). Introduced by Horn (Horn, 1965), parallel analysis is a

method by which the cut-off point for factor retention can be judged, where below the cutoff, the factors possess generally trivial error variance. In simple terms, parallel analysis generates numerous replications of analyses that are drawn from random, normally distributed data with sample size N and number of variables V , concentrating on the number of factors that account for more variance than the factors derived from the random data (O'Connor, 2000). In other words, eigenvalues are extracted from the random data sets that reflect the same number of cases and variables as the observed data (thus the random data parallels the observed data in cases and variables). In the example we will later work with, we have 191 cases and 8 variables. In parallel analysis, there would then be 191 multiplied by 8 random data matrices generated with eigenvalues computed for both the observed correlation matrix and each random data matrix. Decisions on the number of factors to retain are based on comparing the eigenvalues from the original data to the eigenvalues of the random data. Factors are retained when the i th eigenvalue from the observed data is greater than the i th eigenvalue from the random data (O'Connor, 2000). Current practice recommends the use of the eigenvalue which corresponds to the percentile selected by the researcher (e.g., 95th) (Glorfeld, 1995).

Although parallel analysis is not currently available within the point-and-click user interface of popular statistical software such as SPSS, there is user-friendly syntax that has been written that allows users to perform this procedure with their own data within software such as SPSS and SAS (O'Connor, 2000). The syntax can be copied (alleviating potential error in rewriting the code), and the user has to specify only a few simple elements: (a) number of cases, (b) number of variables, (c) location of the data, and (d) the percentile at which the researcher wishes the analysis to be generated.

Number of Variables per Factor

Researchers also need to consider the number of observed variables per factor in their solution. Three variables per factor is the absolute minimum needed to define a factor (Child, 2006). Why at least three variables are needed can be understood by considering a straight line with only two points as estimation of a linear relationship. We can imagine how our line may change if error is introduced by drawing two small circles around each point. Rather than simply two unique points, now these points can be placed anywhere within that circle—this is our margin of error. We can quickly see how different our line may be depending on where the points are placed within the circle. This illustrates that two points are insufficient for estimation of a linear relationship (Child, 2006). Factors that are defined by very few variables (e.g., two or three) may be underdetermined and very unstable when the model is replicated (Brown, 2006).

9.1.1.6 Factor Rotation

Once the data are extracted, it is most always the case that the solution be rotated in order for it to be interpretable. In the world of factor analysis, rotation simply means that the axes (i.e., factor vectors or reference axes) are placed in a different position by turning about the origin (Child, 2006). If the factors are not rotated, axes will lay

between the clusters of variables and the variables will not clearly differentiate to a primary factor. *It is important to note that rotation does nothing to the mathematical fit between the observed and reproduced correlation matrices, as there is mathematical equivalence between solutions prior to rotation and orthogonally rotated solutions.* Rather, rotation serves only to clarify and improve the ability to interpret the solution—there will be clearer differentiation by factor of the factor loadings of the variables. As we discussed with extraction methods, data that has clear correlational patterns will likely produce similar results regardless of rotation method. There are only two types of rotations: orthogonal and oblique, although there are quite a few methods available in standard statistical software that will accomplish the rotation.

Orthogonal Rotation

Variables that are orthogonal are unrelated, and perfect orthogonality is characterized by a correlation value of zero. In orthogonal rotation, axes are rotated at 90-degree angles. Going back to our general understanding of relationships, a correlation of zero means that knowledge of one variable in no way enhances our knowledge of the second. Thus, in the context of factor rotation, orthogonal rotation will produce uncorrelated factors. Considering many situations where factor analysis is applied in the social sciences in particular (and more specifically as we think about human behavior and attributes), however, it is likely the case to anticipate some correlation between factors as (more often than not) the constructs being measured are seldom completely independent of the others. In cases where some correlation between factors does exist, orthogonal rotation will result in a less interpretable solution than oblique rotation. Even if there are substantial correlations between factors, orthogonal rotation will constrain the solution to produce uncorrelated variables, thereby resulting in misleading solutions (Brown, 2006). In cases where there is indeed a lack of relationship between factors, orthogonal and oblique rotations will produce quite similar results.

There are a number of different types of orthogonal rotation techniques available in standard statistical software. These include varimax, quartimax, and equamax, each of which works with a different statistic to maximize or minimize it. Varimax rotation, the most common orthogonal rotation, maximizes the variance of the factor loadings within the factors and across the variables to simplify the factors. Quartimax rotation, on the other hand, maximizes the variance of the factor loadings within the variables and across the factor loadings to simplify the variables. Equamax attempts to bridge varimax and quartimax by simultaneously simplifying both the factors and the variables. Equamax is the least preferred orthogonal rotation, as research suggests it is unstable in situations other than when the number of factors can be specified with confidence.

Oblique Rotation

Factors that are oblique are related, and perfect obliqueness is characterized by a correlation value of one (in absolute value terms). When oblique rotation is applied, the factor axes are rotated independently of each other at different angles (i.e., not just

90 degrees, as is the case with orthogonal rotation). Going back to our general understanding of relationships, a correlation of one means that knowledge of one factor (in this case) perfectly enhances our knowledge of the second factor. Thus, in the context of factor rotation, oblique rotation will produce correlated *factors* (not correlated variables). While oblique seems to be the most defensible option of the two rotations (given that it is reasonable to assume there would be correlation between constructs), be prepared for the possibility that it may increase the difficulty in attaching meaning to your factors. This is because there will likely be an increased number of cross-loading variables in the oblique, as compared to orthogonal, rotated solution. Cross-loading variables are variables that have similar factor loadings for multiple factors. If you are unsure which rotation to select, you may wish to test oblique rotation first and review the factor correlation matrix. Small factor correlations (e.g., less than .30) may warrant orthogonal rotation. There are a few different types of oblique rotation techniques available in standard statistical software, including direct oblimin and promax.

Associated Matrices

The type of rotation selected will alter the matrices generated in your factor solution. In an orthogonal rotation solution, the structure matrix is simply the factor-loading matrix (and is the only matrix that requires review). Oblique rotations will result in generation of both a structure and a pattern matrix. The structure matrix coefficients represent the variance in the observed variables explained by a factor, both a unique (i.e., relationship between the variable and the factor, as with the pattern matrix) and common (i.e., relationship between the variable and the shared variance among the factors) contribution. In oblique rotations, the structure matrix is the product of the pattern and factor correlation matrices, and the loadings in the structure matrix will often be larger than those in the pattern matrix because they reflect overlap in the factors (i.e., are inflated due to this), unless there is a weak relationship between the factors. The pattern matrix coefficients or loadings represent unique contributions only, i.e., unique relationships between the variables and factors. Generally, the larger number of factors, the lower the coefficients in the pattern matrix since there is more common contribution to the variance explained. Because both a structure and pattern matrix are generated with oblique rotation, this requires examination of both when interpreting the meaning of the factors. Of the two matrices, the pattern matrix is the one that is most often reported and interpreted (Brown, 2006).

9.1.1.7 Factor Loadings

In simple terms, the factor loading is the coordinate of a variable along a classification axis. It reflects the relationship between a factor and an observed variable and is the slope of increase (when positive) or decrease (when negative) in the observed variable for each unit of increase or decrease in the factor. The factor-loading value is interpreted in the same units as the measured variables. Now, this is where consideration of the measurement scale of items in the factor analysis come into play. . . . This type of index that measures the relationship between the factor and observed variables is meaningful if, and only if, the observed variables are measured in such a way that the units can be

ordered or ranked and there is equal distance between the units—this implies the variables must be interval or ratio scale. Nominal and ordinal items (even with three or more categories) are usually insufficient to meet this condition (Fabrigar & Wegener, 2012).

The factor loading provides information on the relative contribution that an individual variable makes to a factor, and the researcher must decide which variables load onto which factor. An often-followed, ‘moderately rigorous’ guideline is that a variable should have a factor loading of at least .30 in order to be retained with that factor; however, this rule should be applied only in models with samples of 80 or more (as with samples of this size, a correlation coefficient is statistically significant at an alpha of .01) (Child, 2006, p. 63). A variable with a factor loading of .30, when squared, is interpreted as variance, and this would mean that variable accounts for slightly less than 10% (9% specifically) of the common variance of the factor.

A squared factor loading is a measure of variance accounted for, similar to R^2 . More specifically, it estimates the amount of variance in a factor that is accounted for by the individual variable—the proportion of variance in the item response or variable scores that are explained by a factor. EFA allows the decomposition of observed variance into both common/shared variance and unique variance. In the ideal situation, a variable will have a large coordinate for only one axis and low coordinates for all other axes—providing evidence to suggest that the variable relates to one and only one factor. It is possible to have negative factor loadings. Factors that are defined by variables with both positive and negative factor loadings are called bipolar factors (Child, 2006). The percent of variance in all the variables accounted for by each factor is computed as the sum of the squared factor loadings for that factor divided by the number of variables—which is also the same as dividing the eigenvalue of a factor by the number of variables in the model. Box 9.1 summarizes the process of fitting the factor model.

BOX 9.1 FITTING THE FACTOR MODEL

Element	Options
Factor Extraction	Select an algorithm: <ul style="list-style-type: none"> • Principal components • Unweighted and generalized (weighted) least squares • Maximum likelihood (ML) • Principal axis factoring • Alpha factoring • Image factoring
Communalities	Review communalities: <ul style="list-style-type: none"> • Low communalities (< 2.0): consider removing unless inclusion of the variable is key to interpreting the factor • High communalities (> 1.0): may be evidence of a spurious solution, and may reflect a sample size that is too small or a factor model that has too few or too many factors. Remove the variable with the largest communality and rerun the EFA—repeating this process until the communality estimates are less than one.

Factor Retention	Determine the number of factors to retain: <ul style="list-style-type: none">• Scree plots: Subjective visual tool; can be used as a guide but do not rely on this absent other more objective means• Eigenvalues greater than one: Kaiser's Rule works best and produces fairly accurate results in conditions of moderate to large communalities, modest sample sizes, and 20–50 variables. Given these limiting conditions, applying Kaiser's rule should only be done as a starting point (if at all) when generating your factor model.• Parallel analysis: Most accurate option for determining the number of factors to retain. Decisions on the number of factors to retain are based on statistical analysis, comparing the eigenvalues from the original data to the eigenvalues of randomly generated data.
Number of Variables per Factor	Review the number of variables per factor: <ul style="list-style-type: none">• Minimum: 3 per factor
Factor Rotation	Determine how to rotate the factors: <ul style="list-style-type: none">• Orthogonal (uncorrelated): varimax, quartimax, and equamax• Oblique (correlated): direct oblimin and promax
Factor Loadings	Review factor loadings: <ul style="list-style-type: none">• Ideally, a variable will have a large factor loading for only one factor• A 'moderately rigorous' recommendation: a variable should have a factor loading of at least .30 in order to be retained with that factor<ul style="list-style-type: none">◦ This rule should be applied only in models with samples > 80

9.1.2 Sample Size

Unlike traditional statistical procedures, there is not a power calculation to suggest appropriate sample size for factor analysis. What exists are a number of sample size recommendations for factor analysis that have been made throughout the years, with none reaching consensus as the absolute criterion that must be followed and all later being determined invalid (MacCallum, Widaman, Zhang, & Hong, 1999). These recommendation are generally based on a subject-to-variable ratio (STV) or absolute sample size per number of cases (N).

Case or subject-to-variable ratio (STV) recommendations range from two times the number of cases (Kline, 1979) to five or more times the number of items with a case-to-item ratio greater than or equal to 5 and a minimum of 100 cases, regardless of the case-to-item ratio (Bryant & Yarnold, 1995; Suhr, 2006), more than 5 times the number of items to allow for missing data (Suhr, 2006), 10 times the number of items (Nunally, 1978), and 51 more cases than the number of variables (Lawley & Maxwell, 1971);

Other criterion are based on an absolute number of *cases* (N), with 100 cases being the suggested bare minimum sample size (Gorsuch, 1983; MacCallum et al., 1999), at least 150–300 (tending toward 150 if items are not highly correlated) (Hutcheson & Sofroniou, 1999), at least 200 (Guilford, 1954), at least 250 (Cattell, 1978), and a

sliding scale ranging from 100 to 1,000 (with 100 = poor, 200 = fair, 300 = good, 500 = very good, and 1,000 or greater = excellent) (Comrey & Lee, 1992).

All this to be said, many researchers today would likely agree that these recommendation for STV and absolute number of cases are weak criteria to follow to estimate the sample size for EFA, and there is research to suggest the invalidity of these rules (MacCallum et al., 1999). What is more important is the factorability of the model, as seen through communalities (percent of variance in a variable that is explained jointly by all factors), the degree of overdetermination (ratio of factors to variables), the size of the factor loading, and general model fit. Simulation research suggests that estimating factor structure is achievable, even with small sample sizes (particularly $N > 20$), given the following conditions are met: (a) high communalities (approximately .8 to .9), (b) small number of *expected* factors to be retained (2 to 4), and (c) low model error (which is likely evidenced in situations where communalities are high; RMSR = .00 to .06) (Preacher & MacCallum, 2002). Other simulation research has shown that factors with four or more variables with factor loadings of .60 or greater are interpretable regardless of the sample size (Guadagnoli & Velicer, 1988). Solutions with lower factor loadings (.40) can still be interpreted if the number of cases is at least 150 and the number of variables per factor is larger (> 10) (Guadagnoli & Velicer, 1988).

The take-home message for sample size with EFA is this: Do not adhere to a recommendation criterion for STV or absolute number of cases. Rather, design your study so that you collect the largest sample size that resources will allow. In some cases, this may mean that the sample size will be unnecessarily small. In those instances—and all others, as a matter of fact—be prepared to defend your sample size using previous empirical research, such as the simulation research presented here. And if you are a researcher so inclined to study methodological issues, this is an area ripe for continued examination.

9.1.3 Power

There are no power calculations to suggest appropriate sample size for exploratory factor analysis given a priori or post hoc power. What exists are a number of sample size recommendations as presented previously.

9.1.4 Effect Size

Factor analytic solutions, in and of themselves, do not produce effect size results. Once composite variables are created based on the factor analytic solutions and then those composite variables are applied in an inferential procedure, effect sizes can then be generated.

9.1.5 Assumptions

As with most multivariate statistical procedures, there are a number of assumptions that must be considered with factor analysis, either EFA or PCA. These include

(a) independence, (b) linearity, (c) absence of outliers (both bivariate and multivariate) in cases and variables, and (d) lack of extreme multicollinearity and singularity. As previously discussed, a condition required for conventional factor analysis is continuous data (assuming the factor analytic procedure is computed from a Pearson correlation, as we will assume in this chapter). A large sample size is not necessarily required (as detailed previously) but may be helpful depending on the factor model. Factor analysis is actually robust to violations of the assumption of normality and normality is really not applicable in EFA as it is with many other multivariate procedures. The only exception to this is in the situation where tests of inference are used to determine the number of factors to retain (e.g., when using ML estimation), and in this case, multivariate normality *is* an assumption. Examination of univariate normality, which is not overly sensitive as are multivariate normality tests, can be done through examination of skewness and kurtosis, formal tests of normality, and plots (e.g., Q-Q plots). In terms of multivariate normality, a macro in SPSS (DeCarlo, 1997) (illustrated with MANOVA in chapter 4) can be used to examine a number of multivariate normality indices that include (a) multivariate kurtosis (Mardia, 1970), (b) multivariate skewness and kurtosis based on Small's (1980) multivariate extension of univariate skewness and kurtosis (Looney, 1995), (c) multivariate normality omnibus test (Looney, 1995), (d) largest squared and plot of squared Mahalanobis distance, and (e) critical values for hypothesis test for a single multivariate outlier using Mahalanobis distance (Penny, 1996).

9.1.5.1 Independence

The first assumption is concerned with **independence** of the observations. Violations of this assumption can detrimentally impact standard error values and thus any resulting hypothesis tests. Testing for this assumption is a bit nebulous in exploratory factor analysis, as there are no independent and dependent variables that allow for this type of examination. In the absence of statistical evidence, we will rely on theoretical evidence: If the units have been randomly sampled from a population, there is evidence that the assumption of independence has been met.

9.1.5.2 Linearity

As you recall, factor analysis uses relationships among the variables as the basis for determining factors with conventional factor analysis doing so via a Pearson correlation matrix. Therefore, it is assumed there is a linear relationship among the variables. Bivariate scatterplots can be examined to determine the extent to which this assumption is held.

9.1.5.3 Absence of Outliers in Cases and Variables

Outliers in factor analysis operate in an unfavorable fashion, just as they do in other procedures. One or more outlying cases (either univariate or multivariate) can have undue and unwanted influence on the factor model. In addition to the ways we've screened for outliers in previous procedures (e.g., boxplots), they can also be screened by reviewing standard scores of the variables. Standardized scores with absolute

values of 3.29 or greater (which equates to values more than 3–1/4 standard deviation units from the mean; about .05% of cases are above and below this point in a standardized normal distribution) should be flagged as outliers. Multivariate outliers can be determined by Mahalanobis distance values, which can be calculated using multiple regression, discriminant analysis, or logistic regression (or via simple matrix algebra, without generating other analyses). Multivariate outliers are evidenced by statistically significant Mahalanobis distance scores ($\alpha = .001$ if you tend toward the liberal edge, which is appropriate with EFA), evaluated using a chi-square distribution with degrees of freedom equal to the number of variables. To generate Mahalanobis distance, apply all the variables as independent variables with the dependent variable being a binary variable coded 1 for potential outliers and 0 for all other variables. The process for examining outliers is therefore to look for univariate outliers first. If any are detected, then screen for multivariate outliers.

In factor analysis, it is also possible to have outlying variables, that is, variables that are unrelated to others in the factor model. These outlying variables can be determined by reviewing the following: (a) squared multiple correlations with all other variables and (b) weak correlations with the factors that are identified in the factor analytic model. Outlying variables that are identified can be disregarded.

9.1.5.4 Lack of Extreme Multicollinearity and Singularity

In other procedures, we have discussed how multicollinearity can be problematic because it makes the matrix inversion process unstable. As you recall, multicollinearity is a very strong linear relationship between two or more of the predictors. You may be wondering how it is the case that this can be problematic in factor analysis, as one of the indices we use to determine the ability to factor analyze is the relationship between variables and there is no matrix inversion. In factor analysis, we are concerned with *severe and extreme multicollinearity*, which can be problematic in factor analysis. Singularity is a special case of multicollinearity; it is perfect multicollinearity and occurs when two or more variables perfectly predict and are therefore perfectly redundant. This can occur in factor analysis (just as it did in multiple regression), for example, when a composite variable as well as its component variables are used as predictors in the same factor analytic model.

How do we detect violations of this assumption? Remember that we are looking only for extreme multicollinearity, so we will limit our detection methods quite a bit as compared to our data examination in multiple regression. For EFA, the simplest method is to conduct a series of multiple regression models, one regression model for each variable where that variable is the dependent variable and all remaining variables are the independent variables. If any of the resultant R_k^2 values are close to one (greater than .9 is a good guideline to go by), then there may be an extreme multicollinearity problem. However, large R^2 values may also be due to small sample sizes, so be cautious in interpreting cases where the number of cases is small. If the number of variables is greater than or equal to n , then perfect multicollinearity is a possibility.

TABLE 9.1

Assumptions and Violation of Assumptions: Exploratory Factor Analysis

Assumption	Effect of Assumption Violation
Independence	Influences standard errors of the model and thus hypothesis tests
Linearity	Reduces interpretability of the factor analytic solution
Absence of outlying cases and variables	Exerts undue influence on and distorts the factor analytic solution
Lack of extreme multi-collinearity	Reduces ability to separate effects of variables
Multivariate normality	Minimal effect when violated with exceptions including (a) when hypothesis testing is conducted as part of the EFA, (b) when maximum likelihood is used to estimate the factor model, and (c) with small sample sizes

9.1.5.5 Concluding Thoughts on Assumptions

As mentioned in previous chapters, there is no rule stating that research that violates assumptions must be scrapped. However, researchers who face violations of assumptions must handle these situations on a case-by-case basis, considering both the goal of the analyses and the extent to which the assumptions were violated and the resulting effect of violation. It is also important that researchers present the evidence found, along with justification for decisions that were made. The assumptions are summarized in Table 9.1.

9.2 MATHEMATICAL INTRODUCTION SNAPSHOT

Now that we understand the conditions and decision points, there are a few additional foundational topics related to the underlying mathematics of exploratory factor analysis that may be helpful with which to become acquainted. Note that this is not meant to be a primer on the mathematical proofs nor is it meant to serve as a foundation for which hand calculations can be made. Rather, it is meant to provide a bit more of the mathematical representation for those who are interested in delving deeper into this aspect.

Using matrix algebra, we can express the correlational structure of the common factor model as follows:

$$\mathbf{P} = \Lambda\Phi\Lambda^T + D_\Psi$$

In this equation, \mathbf{P} refers to the population correlation matrix of observed variables. The factor-loading matrix, Λ (lambda), represents the linear influence strength and direction of the latent or component factors on the observed variables. In this matrix, the columns represent the factors and the rows represent the observed variables. Thus, $\Lambda_{3,1}$ refers to the factor loading for (the value of which is the path between) the effect or influence of common factor one on observed variable three.

The transpose of the factor-loading matrix is represented by lambda superscript T , Λ^T . As reviewed in the material on matrix algebra in the appendix, transposing means that

what was originally in the rows now become columns (and what was originally in the columns now become rows).

The covariance matrix among the unique factors is represented by D_ψ (the subscript for which is psi). The diagonals of this matrix are the variances of the unique factors. The off-diagonals are the covariances and are zero when orthogonality is assumed.

The correlation matrix between the factors is represented by Φ (phi). When orthogonality of errors is assumed (i.e., the factors are uncorrelated), the population correlation matrix is simply: $P = \Lambda\Lambda^T + D_\psi$.

Because our interest is in the conceptual understanding of EFA, we'll end our mathematical discussion at this point. The summary of the underlying mathematics of EFA was drawn from Fabrigar and Wegener (2012), which provides a very accessible account. Readers interested in learning more of the mathematics are referred to that source, among others.

■ TABLE 9.2

Factor-Loading Matrix Example

	Factor Matrix ^a	
	[Common Factor 1]	[Common Factor 2]
Index of use of numeracy skills at home	$\Lambda_{11} = .843$	$\Lambda_{12} = -.175$
Index of use of ICT skills at home	$\Lambda_{21} = .673$	$\Lambda_{22} = .066$
Index of use of reading skills at home	$\Lambda_{31} = .528$	$\Lambda_{32} = .153$
Index of use of numeracy skills at work	$\Lambda_{41} = .330$	$\Lambda_{42} = .086$
Index of readiness to learn	$\Lambda_{51} = .412$	$\Lambda_{52} = .504$
Index of use of task discretion at work	$\Lambda_{61} = .059$	$\Lambda_{62} = .311$
Index of learning at work	$\Lambda_{71} = .062$	$\Lambda_{72} = .300$
Index of use of planning skills at work	$\Lambda_{81} = -.183$	$\Lambda_{82} = .296$

Extraction Method: Maximum Likelihood.

a. Two factors extracted. Six iterations required.

■ TABLE 9.3

Example of Correlation Matrix of Common Factors

Factor Correlation Matrix		
Factor	[Common Factor 1]	[Common Factor 2]
1	1.000	
2	$\Phi_{21} = .165$	1.000

Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

9.3 COMPUTING EFA USING SPSS

As we know by now, conventional factor analysis requires continuous data. There are many situations, however, where EFA of ordinal survey (e.g., Likert) items is desirable. Thus, our use of SPSS will first illustrate EFA with continuous data, and this will be followed by an illustration of the use of parallel analysis for factor retention. Next, we will illustrate how to use one of the SPSS add-ons for conducting EFA with ordinal data.

9.3.1 Computing EFA With Continuous Data Using SPSS

Next, we consider SPSS for conducting exploratory factor analysis with data that is continuous in scale (should you have only ordinal items, please see the following SPSS section, “Computing EFA With Ordinal Data”). Before we conduct the analysis, let us talk about the data. The data we are using is the 2013 Survey of Adult Skills (<http://www.oecd.org/site/piaac/surveyofadultskills.htm>), available through the Organisation for Economic Co-operation and Development (OECD). Thank you to OECD for making this data publicly available.

The Survey of Adult Skills, conducted in 33 countries, is part of the Programme for the International Assessment of Adult Competencies (PIAAC), and the first results from the survey were released in 2013. Measured in the survey are “key cognitive and workplace skills needed for individuals to participate in society and for economies to prosper” (see <http://www.oecd.org/site/piaac/surveyofadultskills.htm>). Adults ages 16 to 65 were interviewed in their homes, with 5,000 individuals from each country participating. It is important to note that the Survey of Adult Skills is a complex sample (i.e., not a simple random sample). Although each country was allowed to create their own sampling design and selection plan (for example, some countries oversampled some groups of individuals), it had to adhere to technical standards published by the PIAAC. For example, the U.S. sampling design was a four-stage stratified probability proportional to size design. If you access the full dataset, you will find the last few variables are various weights as well as stratum and unit variables. We won’t get into the technical aspects of this, but when the data are analyzed to adjust for the sampling design (including nonsimple random sampling procedure and disproportionate sampling), the end results are then representative of the intended population. The purpose of the text is not to serve as a primer for understanding complex samples, and thus readers interested in learning more about complex survey designs are referred to any number of excellent resources (Hahs-Vaughn, 2005; Hahs-Vaughn, McWayne, Bulotskey-Shearer, Wen, & Faria, 2011a, 2011b; Lee, Forthofer, & Lorimor, 1989; Skinner, Holt, & Smith, 1989). Additionally, so as to not complicate matters any more than necessary (learning EFA is generally complicated enough!), the applications in this textbook do not illustrate how to adjust for the complex sample design. As such, the results that we see should not be interpreted to represent any larger population but only that select sample of individuals who actually completed the survey. I want to stress that the reason why the sampling design has not been illustrated in the textbook applications is because the point of this section of the textbook is to illustrate how to use statistical software to generate various

procedures and how to interpret the output and not to ensure the results are representative of the intended population. Please do not let this discount or diminish the need to apply this critical step in your own analyses when using complex survey data, as quite a large body of research exists that describes the importance of effectively analyzing complex samples and provides evidence of biased results when the complex sample design is not addressed in the analyses (Hahs-Vaughn, 2005, 2006a, 2006b; Hahs-Vaughn et al., 2011a, 2011b; Kish & Frankel, 1973, 1974; Korn & Graubard, 1995; Lee et al., 1989; Lumley, 2004; Pfeffermann, 1993; Skinner et al., 1989).

Now, let's review the data. We are using the **PIAAC_EFA.sav** file. This is data from the U.S., and the data file has been delimited to include only individuals who were between the ages of 25–29 [AGEG5LFS = 3], who were employed or participated in education or training during the 12 months prior to completing the survey [NEET = 0], and who reported having ‘above high school’ education [B_Q01a_T = 3] ($n = 288$). The size of this sample is more than sufficient to generate EFA, but at the same time small enough to work with for readers who may be using a version of SPSS that limits the number of cases. Additionally, it creates at least an intuitively homogenous sample that would be anticipated to respond similarly on the items. (Note: The complete PIAAC Survey of Adult Skills data file, which includes 5,010 cases, is available from the textbook’s companion website and is titled **PIAAC_SurveyOfAdultSkills.sav**.)

Before we run the data, it's always important to examine frequency distributions of the variables that will be used in the model to assess missing data, potential data entry problems, and similar. With this data, we have some missing data (it has already been coded by the survey collectors as 9996), and thus I've taken the liberty to perform listwise deletion on the missing items (resulting in $n = 191$); however, the remaining variables in the data file have been left as is so that you may practice your data cleaning skills in working with ‘real data.’

Let's look at the data. For the EFA illustration, we'll be working with 13 indices (variables 1–13 in your SPSS file), each of which is measured on a continuous scale.

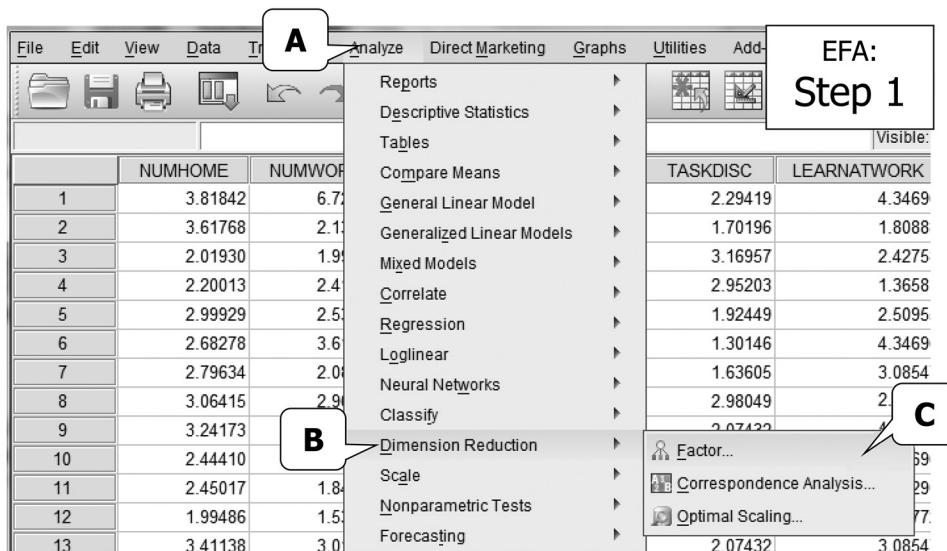
1. Index of use of numeracy skills at home (basic and advanced—derived)
2. Index of use of numeracy skills at work (basic and advanced—derived)
3. Index of use of ICT skills at home (derived)
4. Index of use of reading skills at home (prose and document texts—derived)
5. Index of use of task discretion at work (derived)
6. Index of learning at work (derived)
7. Index of use of planning skills at work (derived)
8. Index of readiness to learn (derived)
9. Index of use of ICT skills at work (derived)
10. Index of use of influencing skills at work (derived)
11. Index of use of reading skills at work (prose and document texts—derived)
12. Index of use of writing skills at work (derived)
13. Index of use of writing skills at home (derived)

The first 13 variables are the indices for EFA. The next three variables in the SPSS dataset were used to delimit the sample. A few variables used for data screening are included (outlier and MAH_1, Mahalanobis distance, which we will discuss as we test assumptions). This is followed by three variables in the dataset that represent the country and participant ID variables. I've left those in the data file just in case you are interested in merging variables from the full dataset with this smaller, delimited file. Each row in the data set still represents one individual. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the respondents were measured. For the EFA illustration, we will work with the 13 continuous index measures.

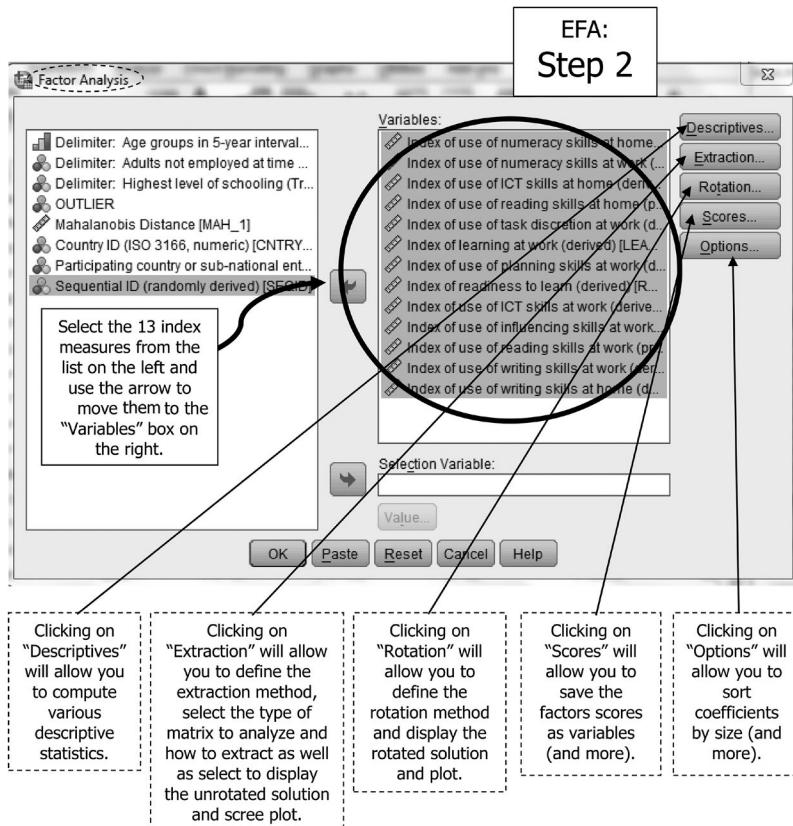
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATWORK	PLANNING	READYTOLEARN	ICTWORK	INFLUENCE	READWORK
1	3.81842	6.72630	3.09433	3.54036	2.29419	4.34696	3.82347	5.00418	4.05051	5.78929	6.21393
2	3.61768	2.13456	3.57861	3.38760	1.70196	1.80888	1.17814	3.21841	4.70515	2.84993	3.61270
3	2.01930	1.99907	2.89841	2.88951	3.16957	2.42758	2.66925	5.00418	2.85857	2.13151	3.56752
4	2.20013	2.41922	2.85982	4.079	2.95203	1.36581	3.82347	2.61279	2.13376	5.78929	2.71635
5	2.99929	2.53842	3.68743	5.9446	1.92449	2.50955	2.66925	5.00418	6.28512	2.79571	2.66331
6	2.68278	3.61287	3.98046	2.0354	1.30146	4.34696	2.22069	2.29040	2.38056	2.71789	4.30735
7	2.79634	2.08571	2.25562	1.93575	1.63605	3.08547	2.22069	1.22896	1.71389	2.93267	2.86236
8	3.06415	2.96927	3.22884	4.5796	2.98049	2.42758	2.22069	2.35854	4.28065	2.63603	3.25087
9	3.24173	3.68996	2.94005	3.0577	2.07432	4.34696	1.63605	3.07297	3.11809	3.05326	4.53684
10	2.44410	1.73260	1.61112	2.88728	2.15631	4.34696	1.43576	2.64003	2.38056	1.76144	2.86718

We will conduct EFA using the 13 index measures (10 are illustrated here).

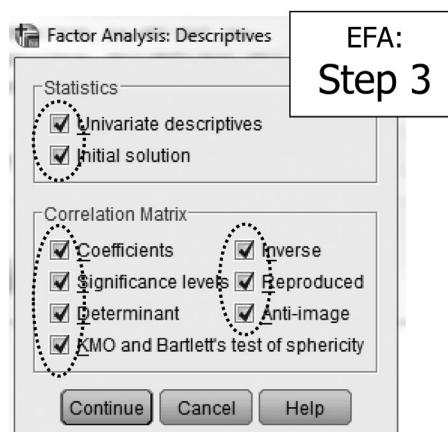
Step 1. To conduct exploratory factor analysis, go to “Analyze” in the top pull-down menu, then select “Dimension Reduction,” and then select “Factor.” Following the screenshot below (Step 1) produces the “Factor Analysis” dialog box.



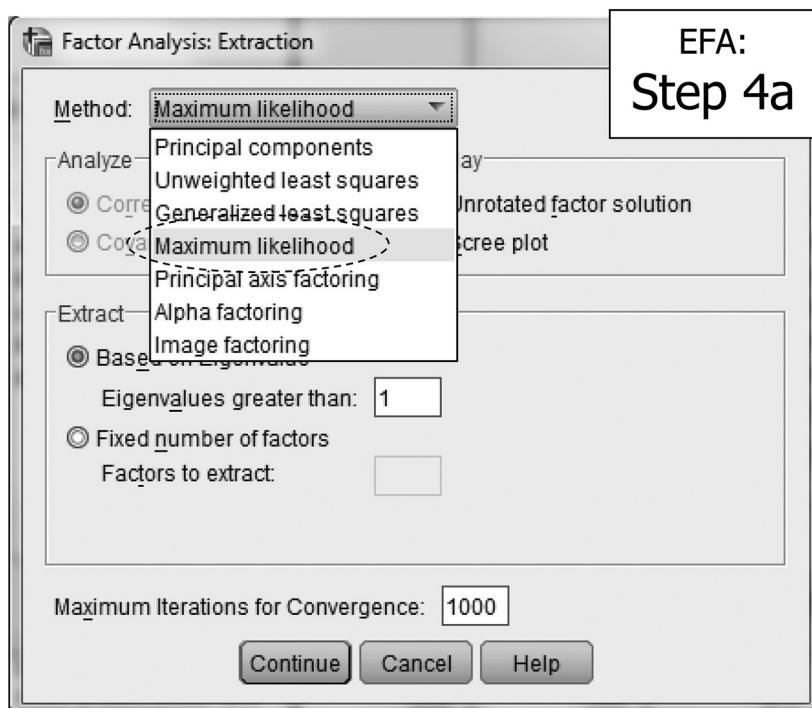
Step 2. Click the 13 index measures and move into the “Variables” box by clicking the arrow button (see screenshot Step 2).



Step 3. From the Factor Analysis dialog box (see screenshot Step 2), clicking on "Descriptives" will provide the option to compute various descriptive statistics (see screenshot Step 3). From the Factor Analysis: Descriptives dialog box, place a checkmark in all the boxes. Click on "Continue" to return to the Factor Analysis dialog box.

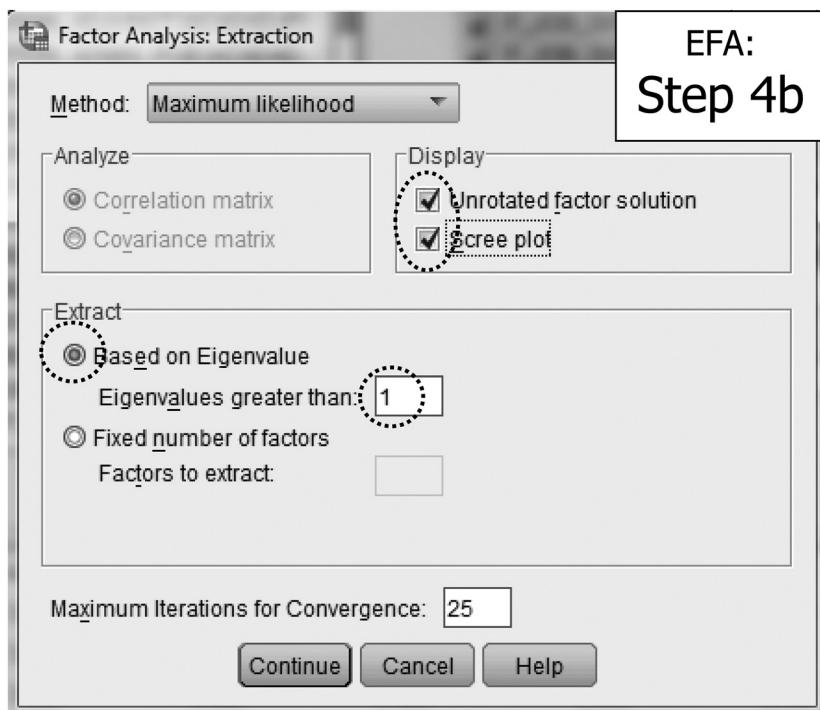


Step 4a. From the Factor Analysis dialog box (see screenshot Step 2), clicking on “Extraction” will provide the option to select various options related extraction methods and what is displayed (see screenshot Step 4a). Using the pull-down menu, click on “Maximum likelihood.” Recall that we discussed how solutions from the different extraction methods will converge in situations where you have a large number of cases and variables and communality estimates that are similar. We also stated that evidence of the stability of the factor solution can be seen in cases where there is convergence of factor analytic solutions when using different extraction methods. Thus, you may want to select a small handful of estimation techniques to test the stability of your factor analytic model under different estimation methods, although for this illustration, we will apply only one.



Step 4b. Also from the Factor Analysis: Extraction dialog box, place a checkmark in the box next to the following: (1) unrotated factor solution and (2) scree plot (see screenshot Step 4b). Under the heading for ‘Extract,’ click the radio button for ‘based on eigenvalue’ and then enter 1 in the box for ‘eigenvalues greater than.’. Recall that the application of Kaiser’s rule consistently (often substantially) overestimates the number of factors (Zwick & Velicer, 1982, 1986), thus we won’t base our factor solution interpretation on it as an important piece of the results. Knowing that it usually overestimates the number of factors to retain, however, it does give us a starting point from which to work. Depending on the solution, we may choose to rerun the model and base the number of factors to extract on a ‘fixed number of factors.’ Leave the default

setting for 'Maximum Iterations for Convergence' at 25. Click on "Continue" to return to the Factor Analysis dialog box.

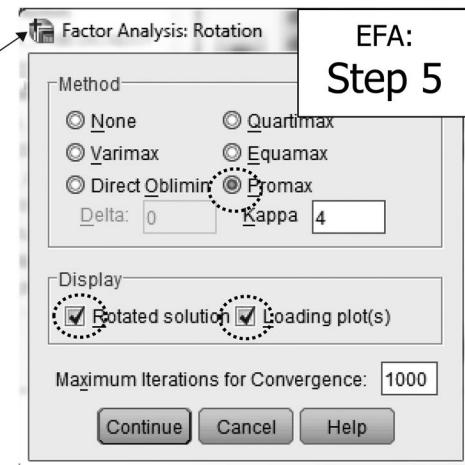


Step 5. From the Factor Analysis dialog box (see screenshot Step 2), clicking on "Rotation" will provide the option to select various options related to rotation methods. Place a checkmark in the box next to the following: (1) rotated solution and (2) loading plot(s) (see screenshot Step 5). In terms of the factor-loading plot, in the event that only one factor is extracted, no plot will be displayed. When two factors are extracted, a two-dimensional plot will be displayed. When three or more factors are extracted, a three-dimensional factor-loading plot of only the first three factors extracted is displayed. Under the heading for 'Method,' click the radio button for 'Promax' and then enter 4 in the box for 'Kappa' (which is the default). (Other values of kappa can be introduced, with the ideal kappa value being one that results in the simplest factor structure with low correlations among the factors; higher kappa values lead to larger correlations among factor and simpler loading structures. The default of 4 is based on previous research which suggests this value produces a generally good solution (Hendrickson & White, 1964).) Change the default setting for 'Maximum Iterations for Convergence' to 1000. The number of iterations to convergence simply defines how many iterations the algorithm can take to perform the rotation. It is likely the case that 1000 is overkill, but it doesn't hurt to set it at a large value just in case it's required. Click on "Continue" to return to the Factor Analysis dialog box.

Clicking on "Rotation" will allow you to define the rotation method, select what to display, and define iterations.

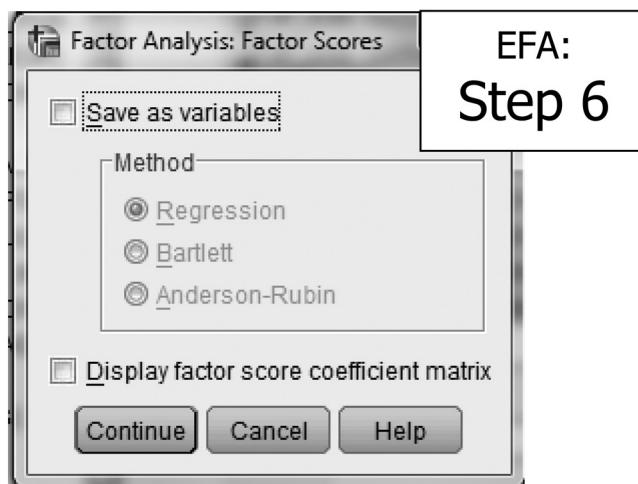
Direct oblimin and promax are oblique rotation methods, allowing the factors to be correlated. Varimax, quartimax, and equamax are orthogonal rotation methods, assuming unrelated factors and maintaining the axes at 90 degrees.

What is displayed in the output is dependent on the method of rotation. The rotated pattern and factor transformation matrices are displayed with orthogonal rotations. The pattern, structure, and factor correlation matrices are displayed with oblique rotations.



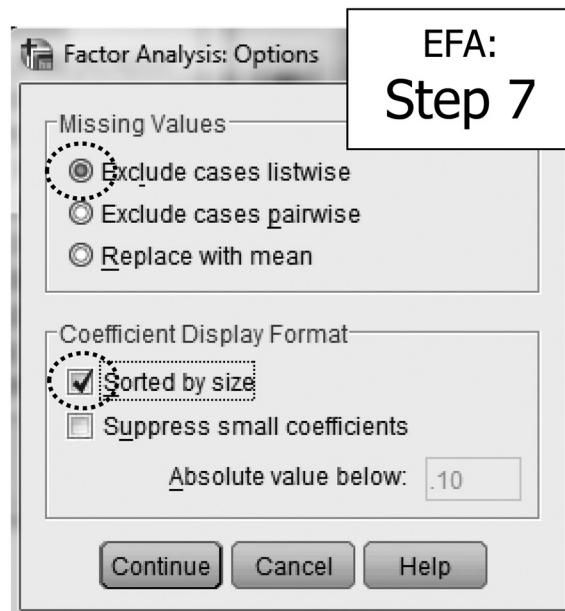
EFA:
Step 5

Step 6. From the Factor Analysis dialog box (see screenshot Step 2), clicking on "Scores" will provide the option to save the variables created as composite scores and to display the factor score coefficient matrix (see screenshot Step 6). Many times, researchers select to skip this step and use methods such as the mean sum (i.e., adding all the items together and dividing by the number of items) as a method to create the composite score. If you do choose to allow the software to create your composite score, there are three methods from which to choose to estimate the factor score coefficients. The *regression method* produces factor scores that have a mean of 0 and a variance that equals the squared multiple correlation between the estimated factor scores and the true factor values. The factor scores estimated from the regression method may be correlated even if the factors are orthogonal. The *Bartlett score* produces factor scores that have a mean of 0. This method minimizes



the sum of squares of the unique factors over the range of variables. The *Anderson-Rubin method* is a modified Bartlett method that produces factor scores with a mean of 0 and standard deviation of 1 and that maintains orthogonality of the estimated factors. Thus, the scores produced are uncorrelated. At this time, do not make any selections on this screen, as we will adhere to the mean sum method for creating a composite score. Click on "Continue" to return to the Factor Analysis dialog box.

Step 7. From the Factor Analysis dialog box (see screenshot Step 2), clicking on "Options" will bring up the dialog box that allows various options for dealing with missing values, as well as options for displaying the coefficients (see screenshot Step 7). We will leave the default setting for the Missing Values as 'exclude cases listwise.' For our purposes, because we have already dealt with missing values, which selection is made for missing values is moot. As you conduct your own research, however, should you have missing values, it should be dealt with prior to generating the factor analysis and not within the EFA, as none of the three options provided are acceptable means for which to address missing values (the exception may be if you have an extremely small percentage of missing, such as 5% or less). Under the heading for Coefficient Display Format, place a checkmark in the box for 'sorted by size.' This will make it much easier to see the clusters of variables produced in the factor solution, as it groups the items by factor in descending order of factor-loading size. Then click on "Continue" to return to the Factor Analysis dialog box. From the "Factor Analysis" dialog box, click on "OK" to generate the output.



Interpreting the output. Annotated results are presented in Table 9.4.

TABLE 9.4

SPSS Results for the Exploratory Factor Analysis Example

Descriptive Statistics										
	Mean	Std. Deviation	Analysis N							
Index of use of numeracy skills at home (basic and advanced - derived)	2.5237439	.87857371	191							
Index of use of numeracy skills at work (basic and advanced - derived)	2.4790569	1.09711244	191							
Index of use of ICT skills at home (derived)	2.6264839	.74799761	191							
Index of use of reading skills at home (prose and document texts - derived)	2.7400939	.70200056	191							
Index of use of task discretion at work (derived)	1.8943278	.77157784	191							
Index of learning at work (derived)	2.5732116	.97107503	191							
Index of use of planning skills at work (derived)	2.1381388	1.10080532	191							
Index of readiness to learn (derived)	2.8349031	1.04219028	191							
Index of use of ICT skills at work (derived)	2.4364373	1.01810011	191							
Index of use of influencing skills at work (derived)	2.5390486	1.08697897	191							
Index of use of reading skills at work (prose and document texts - derived)	2.5950712	.73915941	191							
Index of use of writing skills at work (derived)	2.5990950	1.00602294	191							
Index of use of writing skills at home (derived)	2.4855768	.87936551	191							

Correlation Matrix													
	Index of use of numeracy skills at home (basic and advanced - derived)	Index of use of numeracy skills at work (basic and advanced - derived)	Index of use of reading skills at home (prose and document texts - derived)	Index of use of task discretion at work (derived)	Index of learning at work (derived)	Index of use of planning skills at work (derived)	Index of readiness to learn (derived)	Index of use of ICT skills at work (derived)	Index of use of influencing skills at work (derived)	Index of use of reading skills at work (prose and document texts - derived)	Index of use of writing skills at work (derived)	Index of use of writing skills at home (derived)	
Correlation	1.000	.292	.554	.418	.015	.001	-.202	.248	.101	-.007	.066	.021	.418
	Index of use of numeracy skills at home (basic and advanced - derived)	1.000	.176	.166	.086	.098	.013	.161	.179	.043	.310	.230	.055
	Index of use of numeracy skills at work (basic and advanced - derived)		1.000	.374	.014	.045	-.111	.341	.320	.069	.212	.118	.475
	Index of use of ICT skills at home (derived)			1.000	.049	.107	-.048	.296	.250	.143	.235	.151	.548
	Index of use of reading skills at home (prose and document texts - derived)				1.000	.037	.140	.206	.312	-.096	.074	.140	-.058
	Index of use of task discretion at work (derived)					1.000	.139	.167	.135	.248	.439	.287	.092
	Index of learning at work (derived)						1.000	.026	.113	.451	.178	.207	-.059
	Index of use of planning skills at work (derived)							1.000	.374	.049	.154	.130	.199
	Index of readiness to learn (derived)								1.000	.163	.350	.412	.156
	Index of use of ICT skills at work (derived)									1.000	.330	.297	-.015
	Index of use of influencing skills at work (derived)										1.000	.538	.180
	Index of use of reading skills at work (prose and document texts - derived)											1.000	.145
	Index of use of writing skills at work (derived)												1.000
	Index of use of writing skills at home (derived)												

a Determinant = .037	The off-diagonals of the correlation matrix provide simple bivariate correlations between each of the items in the EFA. Correlations of .30 and above provide evidence of good factorability. Unfortunately, we have quite a few small correlations that will likely cause problems in our factor solution.
	The bottom 1/2 of the matrix presents the <i>p</i> value for each correlation.

The footer provides information on the determinant. Non-zero determinant values will help ensure the factor solution can be computed.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

	Index of use of numeracy skills at home (basic and advanced - derived)	Index of use of numeracy skills at work (basic and advanced - derived)	Index of use of ICT skills at home (basic and advanced - derived)	Index of use of reading skills at home (prose and document - derived)	Index of use of reading skills at work (prose and document - derived)	Index of use of task discretion at work (derived)	Index of use of learning at work (derived)	Index of use of planning skills at work (derived)	Index of readiness to learn (derived)	Index of use of ICT skills at work (derived)	Index of use of influencing skills at work (derived)	Index of use of reading skills at work (prose and document - derived)	Index of use of writing skills at work (derived)	Index of use of writing skills at home (derived)
Index of use of numeracy skills at home (basic and advanced - derived)		.419	-.419	-.737	-.320	-.090	.013	.255	-.078	.210	-.112	.200	.038	-.260
Index of use of numeracy skills at work (basic and advanced - derived)			.036	-.063	.000	.048	-.056	-.068	-.029	.132	-.357	-.137	.219	
Index of use of ICT skills at home (basic and advanced - derived)				.090	.000	.114	.121	-.260	-.387	-.117	.217	.119	.493	
Index of use of reading skills at home (prose and document - derived)					.081	.038	.149	-.204	-.150	-.287	.145	.086	-.760	
Index of use of reading skills at work (prose and document - derived)						.238	-.016	-.265	-.141	-.357	.335	.012	-.078	.197
Index of use of task discretion at work (derived)							.007	-.199	.099	-.156	-.513	-.085	-.076	
Index of use of learning at work (derived)								.007	-.052	.012	.674	-.008	-.057	.167
Index of use of planning skills at work (derived)									.052	-.343	.050	.092	.063	.040
Index of readiness to learn (derived)										.008	-.074	-.177	.420	.060
Index of use of ICT skills at work (derived)											.000	-.232	-.200	.342
Index of use of influencing skills at work (derived)												.055	-.585	-.079
Index of use of reading skills at work (prose and document - derived)													.154	-.191
Index of use of writing skills at work (derived)														.748
Index of use of writing skills at home (derived)														

We won't spend a lot of time examining the inverted correlation but this is another method that can be used to determine factorability. The off-diagonals should be close to zero (Guttman, 1953). Kaiser (1970) extended this work, and thus we'll review the KMO MSA as a technique for factorability.

KMO and Bartlett's Test	
Kaiser-Meyer-Olkin Measure of Sampling Adequacy.	.705
Bartlett's Test of Sphericity	606.968
df	78
Sig.	.000

Measure of Sampling Adequacy (MSA): An overall MSA of .50 or above should be achieved before proceeding with factor analysis. According to Kaiser and Rice (1974), our MSA is 'middling.'

Bartlett's Test of Sphericity: This is a statistical test to determine if the overall correlation matrix is an identity matrix (i.e., the null hypothesis is that our overall correlation matrix is equal to an identity matrix), and thus we want to find statistical significance here—suggesting that we do not have an identity matrix. In practical terms, a statistically significant Bartlett's test indicates at least some of the variables have significant correlations. *A statistically significant Bartlett's test is desirable and suggests factor analysis is appropriate.*

In this example, we've met both criteria suggesting it is appropriate to factor analyze our variables.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Anti-image Matrices													
	Index of use of numeracy skills at home (basic and advanced) - derived	Index of use of numeracy skills at work (basic and advanced) - derived	Index of use of ICT skills at home (basic and advanced) - derived	Index of use of reading skills at home (gross and document texts) - derived	Index of use of task discretion at work (derived)	Index of learning at work (derived)	Index of use of planning skills at work (derived)	Index of readiness to learn (derived)	Index of use of ICT skills at work (derived)	Index of use of influencing skills at work (derived)	Index of use of writing skills at work (derived)	Index of use of writing skills at home (derived)	
Anti-image Covariance	.553	-.183	-.219	-.104	-.043	.005	.099	-.032	.072	-.039	.060	.013	-.080
Index of use of numeracy skills at home (basic and advanced) - derived													
Index of use of numeracy skills at work (basic and advanced) - derived	-.183	.791	.015	-.029	.000	.029	-.031	-.040	-.014	.067	-.152	-.067	.097
Index of use of ICT skills at home (basic and advanced) - derived	-.219	.015	.537	.028	.039	.047	.046	-.104	-.129	-.040	-.063	.040	-.148
Index of use of reading skills at home (gross and document texts) - derived	-.104	-.029	.028	.590	-.039	.017	.062	-.090	-.055	-.108	-.046	.031	-.251
Index of use of task discretion at work (derived)	-.043	.000	.039	-.039	.809	-.010	-.151	-.085	-.180	.173	.005	-.039	.089
Index of learning at work (derived)	.005	.029	.047	.017	-.010	.769	.004	-.115	.047	-.076	-.212	.040	-.033
Index of use of planning skills at work (derived)	.099	-.031	.046	.062	-.151	.004	.706	-.027	.005	-.303	-.003	-.025	-.066
Index of readiness to learn (derived)	-.032	-.040	-.104	-.090	-.085	-.115	-.027	.747	-.159	.024	.037	.029	.017
Index of use of ICT skills at work (derived)	.072	.014	-.129	-.055	-.100	.047	.005	-.159	.622	-.029	-.059	-.161	.021
Index of use of influencing skills at work (derived)	-.039	.067	-.040	-.108	.173	-.076	-.303	.024	-.029	.638	-.080	-.079	.122
Index of use of reading skills at work (gross and document texts) - derived	.060	-.152	-.063	-.046	.005	-.212	-.003	.037	-.059	-.080	.539	-.195	-.024
Index of use of writing skills at work (derived)	.013	-.067	.040	.031	-.039	-.040	-.025	.029	-.161	-.079	-.195	.618	-.066
Index of use of writing skills at home (derived)	-.080	.097	-.148	-.251	.069	-.033	-.066	.017	.021	.122	-.024	-.066	.560
Anti-image Correlations													
Index of use of numeracy skills at home (basic and advanced) - derived	-.707*												
Index of use of numeracy skills at work (basic and advanced) - derived	-.277	-.777*											
Index of use of ICT skills at home (basic and advanced) - derived	.024	.878*											
Index of use of reading skills at home (gross and document texts) - derived	.739*	.050											
Index of use of reading skills at home (gross and document texts) - derived	-.402	.024	.739*										
Index of use of task discretion at work (derived)	-.183	-.043	.050	.746*									
Index of learning at work (derived)	-.064	.000	.060	-.058	.501*								
Index of use of planning skills at work (derived)	.008	.037	.073	.025	-.013	.728*							
Index of readiness to learn (derived)	.159	-.042	.074	.096	-.200	.005	.565*						
Index of use of ICT skills at work (derived)	-.050	-.052	-.165	-.135	-.110	-.151	-.038	.791*					
Index of use of influencing skills at work (derived)	.123	-.020	-.224	-.091	-.253	.068	.008	-.234	.746*				
Index of use of writing skills at work (derived)	-.066	.094	-.068	-.176	.240	-.109	-.452	.035	-.047	.568*			
Index of use of reading skills at work (gross and document texts) - derived	.109	-.233	-.117	-.082	.068	-.330	-.005	.058	-.102	-.136	.746*		
Index of use of writing skills at work (derived)	.022	-.096	.069	.052	-.055	-.068	-.038	.043	-.260	-.125	-.338	.776*	
Index of use of writing skills at home (derived)	-.144	.146	-.270	-.437	.133	-.050	-.105	.026	.036	.204	-.043	-.112	.676*

a. Measures of Sampling Adequacy (MSA)

The anti-image covariance matrix presents the negatives of the partial covariances, and thus the anti-image correlation matrix provides the negatives of the partial correlation coefficients. The measure of sampling adequacy (MSA) for a variable is displayed on the diagonal of the anti-image correlation matrix (denoted by footnote 'a'). Reviewing the anti-image correlations, items with individual MSA values below .50 are considered unacceptable and should be excluded. In this example, all are acceptable.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Communalities^a

	Initial	Extraction
Index of use of numeracy skills at home (basic and advanced - derived)	.447	.547
Index of use of numeracy skills at work (basic and advanced - derived)	.209	.151
Index of use of ICT skills at home (derived)	.463	.545
Index of use of reading skills at home (prose and document texts - derived)	.410	.407
Index of use of task discretion at work (derived)	.191	.232
Index of learning at work (derived)	.231	.232
Index of use of planning skills at work (derived)	.294	.278
Index of readiness to learn (derived)	.253	.288
Index of use of ICT skills at work (derived)	.378	.638
Index of use of influencing skills at work (derived)	.362	.999
Index of use of reading skills at work (prose and document texts - derived)	.461	.849
Index of use of writing skills at work (derived)	.382	.417
Index of use of writing skills at home (derived)	.440	.428

Extraction Method: Maximum Likelihood.

a. One or more communality estimates greater than 1 were encountered during iterations. The resulting solution should be interpreted with caution.

For the sake of brevity, the output from the additional models generated are not presented. However, in addition to removing ‘Index of use of influencing skills at work,’ the following variables also had communalities greater than 1.0 and were removed through the iterative process of running the factor solution, reviewing communalities, and removing the one variable with the largest communality: ‘index of using writing skills at home,’ ‘index of using reading skills at work (prose and document text),’ ‘index of using writing skills at work,’ ‘index of using ICT skills at work.’ *Thus a total of 5 of our original 13 variables were removed from the model as they suggested they were not factorable, leaving 8 variables to factor analyze.*

Because the output prior to the communalities table has been annotated in detail, it will not be presented again. The descriptive statistics and bivariate correlation coefficients do not change because variables are removed from the set that are factor analyzed, and thus they are not presented again. Rather, we’ll pick up with the tables that do reflect new, adjusted values as a result of removing variables from the set.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Inverse of Correlation Matrix								
	Index of use of numeracy skills at home (basic and advanced - derived)	Index of use of numeracy skills at work (basic and advanced - derived)	Index of use of ICT skills at home (derived)	Index of use of reading skills at home (prose and document texts - derived)	Index of use of task discretion at work (derived)	Index of learning at work (derived)	Index of use of planning skills at work (derived)	Index of readiness to learn (derived)
Index of use of numeracy skills at home (basic and advanced - derived)	1.702	-.315	-.712	-.381	.001	.074	.240	-.032
Index of use of numeracy skills at work (basic and advanced - derived)	-.315	1.121	.010	-.028	-.066	-.085	-.055	-.070
Index of use of ICT skills at home (derived)	-.712	.010	1.578	-.201	.060	-.003	.022	-.317
Index of use of reading skills at home (prose and document texts - derived)	-.381	-.028	-.201	1.308	-.006	-.091	-.017	-.203
Index of use of task discretion at work (derived)	.001	-.066	.060	-.006	1.071	.022	-.140	-.229
Index of learning at work (derived)	.074	-.085	-.003	-.091	.022	1.062	-.135	-.155
Index of use of planning skills at work (derived)	.240	-.055	.022	-.017	-.140	-.135	1.090	-.030
Index of readiness to learn (derived)	-.032	-.070	-.317	-.203	-.229	-.155	-.030	1.261

KMO and Bartlett's Test

Kaiser-Meyer-Olkin Measure of Sampling Adequacy.	695
Bartlett's Test of Sphericity	193.696
df	28
Sig.	.000

With an overall MSA of .50 and a statistically significant Bartlett's test, we've again met both criteria suggesting it is appropriate to factor analyze our variables.

Anti-image Matrices

	Index of use of numeracy skills at home (basic and advanced - derived)	Index of use of numeracy skills at work (basic and advanced - derived)	Index of use of ICT skills at home (derived)	Index of use of reading skills at home (prose and document texts - derived)	Index of use of task discretion at work (derived)	Index of learning at work (derived)	Index of use of planning skills at work (derived)	Index of readiness to learn (derived)
Anti-image Covariance								
Index of use of numeracy skills at home (basic and advanced - derived)	.588	-.165	-.265	-.171	.000	.041	.130	-.015
Index of use of numeracy skills at work (basic and advanced - derived)	-.165	.892	.006	-.019	-.055	-.072	-.045	-.049
Index of use of ICT skills at home (derived)	-.265	.006	.634	-.097	.036	-.002	.013	-.160
Index of use of reading skills at home (prose and document texts - derived)	-.171	-.019	-.097	.764	-.004	-.066	-.012	-.123
Index of use of task discretion at work (derived)	.000	-.055	.036	-.004	.934	.020	-.120	-.169
Index of learning at work (derived)	.041	-.072	-.002	-.066	.020	.942	-.117	-.116
Index of use of planning skills at work (derived)	.130	-.045	.013	-.012	-.120	-.117	.918	-.022
Index of readiness to learn (derived)	-.015	-.049	-.160	-.123	-.169	-.116	-.022	.793
Anti-Image Correlation								
Index of use of numeracy skills at home (basic and advanced - derived)	.663 ^a		-.228	-.434	-.255	.000	.055	.177
Index of use of numeracy skills at work (basic and advanced - derived)	-.228	.733 ^b		.008	-.023	-.060	-.078	-.050
Index of use of ICT skills at home (derived)	-.434	.008	.700 ^c		-.140	.046	-.002	.016
Index of use of reading skills at home (prose and document texts - derived)	-.255	-.023	-.140	.793 ^d		-.005	-.077	-.014
Index of use of task discretion at work (derived)	.000	-.060	.046	-.005	.544 ^e	.021	-.129	-.197
Index of learning at work (derived)	.055	-.078	-.002	-.077	.021	.592 ^f	-.125	-.134
Index of use of planning skills at work (derived)	.177	-.050	.016	-.014	-.129	-.125	.585 ^g	-.025
Index of readiness to learn (derived)	-.022	-.059	-.225	-.158	-.197	-.134	-.025	.725 ^h

a. Measures of Sampling Adequacy(MSA)

Reviewing the anti-image correlations, items with individual MSA values below .50 are considered unacceptable and should be excluded. In this example, all are again acceptable.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Communalities	
	Initial Extraction
Index of use of numeracy skills at home (basic and advanced - derived)	.412 .741
Index of use of numeracy skills at work (basic and advanced - derived)	.108 .116
Index of use of ICT skills at home (derived)	.366 .458
Index of use of reading skills at home (prose and document texts - derived)	.236 .303
Index of use of task discretion at work (derived)	.066 .100
Index of learning at work (derived)	.058 .094
Index of use of planning skills at work (derived)	.082 .121
Index of readiness to learn (derived)	.207 .424

Extraction Method: Maximum Likelihood.

Although all our communalities are now under the threshold of 1.0 so as not to generate a warning, we do see that we still have some communalities that are relatively small (under .30), however we will retain them in our set of variables given that we have already removed quite a few variables.

Factor	Total Variance Explained						Rotation Sums of Squared Loadings ^a	
	Initial Eigenvalues			Extraction Sums of Squared Loadings				
	Total	% of Variance	Cumulative %	Total	% of Variance	Cumulative %		
1	2.313	28.909	28.909	1.762	22.028	22.028	1.779	
2	1.344	16.800	45.709	.594	7.430	29.458	.656	
3	.971	12.135	57.844					
4	.894	11.177	69.021					
5	.791	9.887	78.907					
6	.665	8.317	87.224					
7	.620	7.749	94.974					
8	.402	5.026	100.000					

Extraction Method: Maximum Likelihood.

a. When factors are correlated, sums of squared loadings cannot be added to obtain a total variance.

'Initial eigenvalues' presents the variance explained by the initial solution. Only two factors in the initial solution have eigenvalues greater than 1, accounting for about 46% of the variability in the original variables. This suggests much unexplained variation. While we do not suggest applying Kaiser's rule to determine the number of factors to extract, if you do apply that rule to your data, the 'initial eigenvalues' should be eigenvalues reviewed to make the decision as these eigenvalues are derived from the unreduced input correlation matrix.

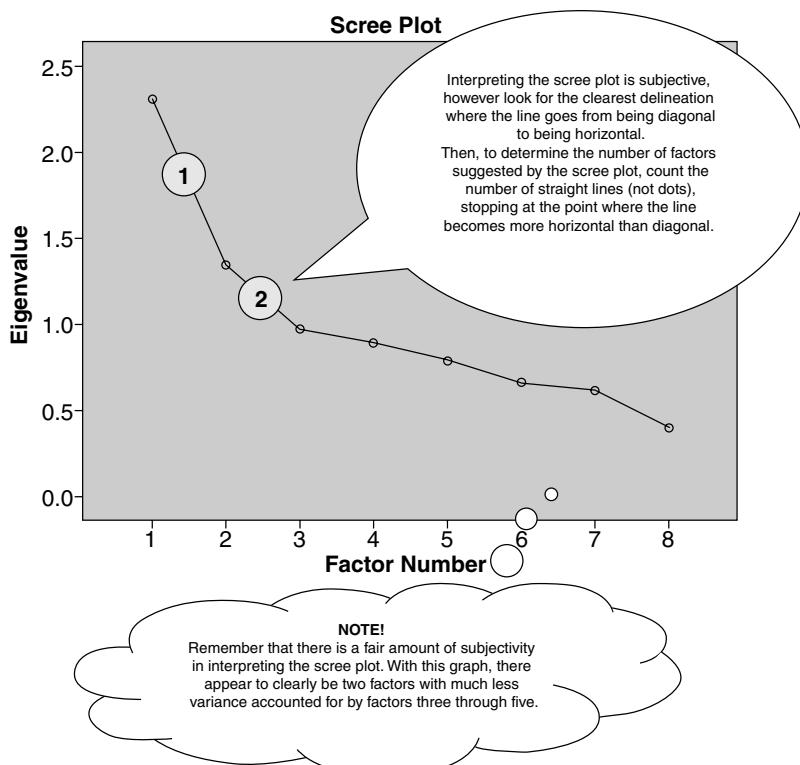
'Extraction sum of squared loadings' presents the variance explained by the extracted factors *before rotation*. The cumulative variability explained by the two factor extracted solution is about 29%, about 17% less than the initial solution. This means about 17% of the variation explained by the initial solution is lost as a result of factors unique to the original variables and variability unexplained by the factor solution.

In the social and behavioral sciences, it is common to find around 60% of the total variance explained in factor analytic models (Child, 2006).

'Rotation sums of squared loadings' provides the variance explained by the extracted factors *after rotation*. Note the footer regarding our oblique rotation (i.e., sums of squared factor loadings) cannot be added together to reflect total variance.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example



	Factor Matrix ^a	
	1	2
Index of use of numeracy skills at home (basic and advanced - derived)	.843	-.175
Index of use of ICT skills at home (derived)	.673	.066
Index of use of reading skills at home (prose and document texts - derived)	.528	.153
Index of use of numeracy skills at work (basic and advanced - derived)	.330	.086
Index of readiness to learn (derived)	.412	.504
Index of use of task discretion at work (derived)	.059	.311
Index of learning at work (derived)	.062	.300
Index of use of planning skills at work (derived)	-.183	.296

The factor matrix presents the unrotated solution. We rotated our solution (and this will nearly always be the case), thus we are not interested in these results.

Extraction Method: Maximum Likelihood.

a. 2 factors extracted. 6 iterations required.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Goodness-of-fit Test

Chi-Square	df	Sig.
9.564	13	.729

The null hypothesis for the goodness-of-fit test is that the factor model sufficiently describes the data. The results of this test provide evidence of the extent to which our factor solution reproduces the variance-covariance matrix. In this example, we fail to reject the null hypothesis providing evidence that the factor model does indeed describe the data—in other words, the relationships among the variables is sufficiently described by the factor model and good fit is suggested.

Reproduced Correlations

		Index of learning at work (derived)	Index of readiness to learn (derived)	Index of use of ICT skills at home (derived)	Index of use of numeracy skills at home (basic and advanced - derived)	Index of use of numeracy skills at work (basic and advanced - derived)	Index of use of planning skills at work (derived)	Index of use of reading skills at home (prose and document texts - derived)	Index of use of task discretion at work (derived)
Reproduced Correlation	Index of learning at work (derived)	.094 ^a	.177	.062	.000	.046	.077	.079	.097
	Index of readiness to learn (derived)	.177	.424 ^a	.311	.259	.179	.074	.295	.181
	Index of use of ICT skills at home (derived)	.062	.311	.458 ^a	.556	.228	-.104	.366	.060
	Index of use of numeracy skills at home (basic and advanced - derived)	.000	.259	.556	.741 ^a	.263	-.206	.419	-.005
	Index of use of numeracy skills at work (basic and advanced - derived)	.046	.179	.228	.263	.116 ^a	-.035	.187	.046
	Index of use of planning skills at work (derived)	.077	.074	-.104	-.206	-.035	.121 ^a	-.051	.081
	Index of use of reading skills at home (prose and document texts - derived)	.079	.295	.366	.419	.187	-.051	.303 ^a	.079
	Index of use of task discretion at work (derived)	.097	.181	.060	-.005	.046	.081	.079	.100 ^a
Residual ^b									
	Index of learning at work (derived)		-.010	-.016	.001	.052	.061	.028	-.060
	Index of readiness to learn (derived)		-.010	.030	-.013	-.018	-.048	.001	.025
	Index of use of ICT skills at home (derived)		-.016	.030	-.002	-.052	-.007	.008	-.046
	Index of use of numeracy skills at home (basic and advanced - derived)		.001	-.013	-.002	.030	.005	.000	.020
	Index of use of numeracy skills at work (basic and advanced - derived)		.052	-.018	-.052	.030	.048	-.021	.040
	Index of use of planning skills at work (derived)		.061	-.048	-.007	.005	.048	.003	.059
	Index of use of reading skills at home (prose and document texts - derived)		.028	.001	.008	.000	-.021	.003	-.030
	Index of use of task discretion at work (derived)		-.060	.025	-.046	.020	.040	.059	-.030

Extraction Method: Maximum Likelihood.

a. Reproduced communalities

b. Residuals are computed between observed and reproduced correlations. There are 5 (17.0%) nonredundant residuals with absolute values greater than 0.05.

The correlation matrix based on the extracted factors is the **reproduced correlation matrix**, and these coefficients should be very close to the values in the original correlation matrix. When that happens (i.e., the reproduced and original coefficients are very close in value), the values in the residual matrix will be close to zero (as the residual values reflect the difference—as simple subtraction—between the original and the reproduced matrix) *and* the extracted factors account for much of the variance in the original correlation matrix—therefore, the extracted factors represent the original data well. The values on the diagonal of the reproduced correlation matrix are also the values presented as extracted communalities.

In this example, our residuals are quite small, the largest being about .50 in absolute value, suggesting the extracted factors are representing the original data well.

TABLE 9.4 (continued)

SPSS Results for the Exploratory Factor Analysis Example

Pattern Matrix ^a		
	Factor	
	1	2
Index of use of numeracy skills at home (basic and advanced - derived)	.869	-.219
Index of use of ICT skills at home (derived)	.670	.034
Index of use of reading skills at home (prose and document texts - derived)	.514	.129
Index of use of numeracy skills at work (basic and advanced - derived)	.322	.071
Index of readiness to learn (derived)	.355	.490
Index of use of task discretion at work (derived)	.022	.312
Index of use of planning skills at work (derived)	-.219	.309
Index of learning at work (derived)	.027	.301

Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

a. Rotation converged in 3 iterations.

The rotation we selected was oblique rotation (assuming correlated factors) using promax rotation. Oblique rotations will produce *both* a factor **pattern matrix** (which are the coefficients for the linear combination of the variables; the factor loadings of each variable onto the factor) and a factor **structure matrix** (which are correlations between the variables and the factors—the product of the pattern and factor correlation matrices—thus taking into account the relationship between factors). It has been suggested that both the pattern and structure matrices be used to interpret the factors (Gorsuch, 1983).

The pattern and structure matrices will be identical with orthogonal rotations (recall that orthogonal means that the factors are not assumed to correlate).

Structure Matrix		
	Factor	
	1	2
Index of use of numeracy skills at home (basic and advanced - derived)	.833	-.076
Index of use of ICT skills at home (derived)	.676	..
Index of use of reading skills at home (prose and document texts - derived)	.535	.214
Index of use of numeracy skills at work (basic and advanced - derived)	.333	.124
Index of readiness to learn (derived)	.436	.549
Index of use of task discretion at work (derived)	.074	.315
Index of learning at work (derived)	.077	.305
Index of use of planning skills at work (derived)	-.168	.273

Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

The factors are identified by reviewing the coefficients. Variables that are associated with factor 1 have high values for factor 1 but not factor 2. (Note that this is where sorting by size is important as the software automatically arranges the variables in descending value, making it easy to see the set of variables that load most strongly on each factor.)

Although the values in the pattern and structure matrices are slightly different, both suggest that the same variables load on the same factors.

Factor Correlation Matrix		
Factor	1	2
1	1.000	.165
2	.165	1.000

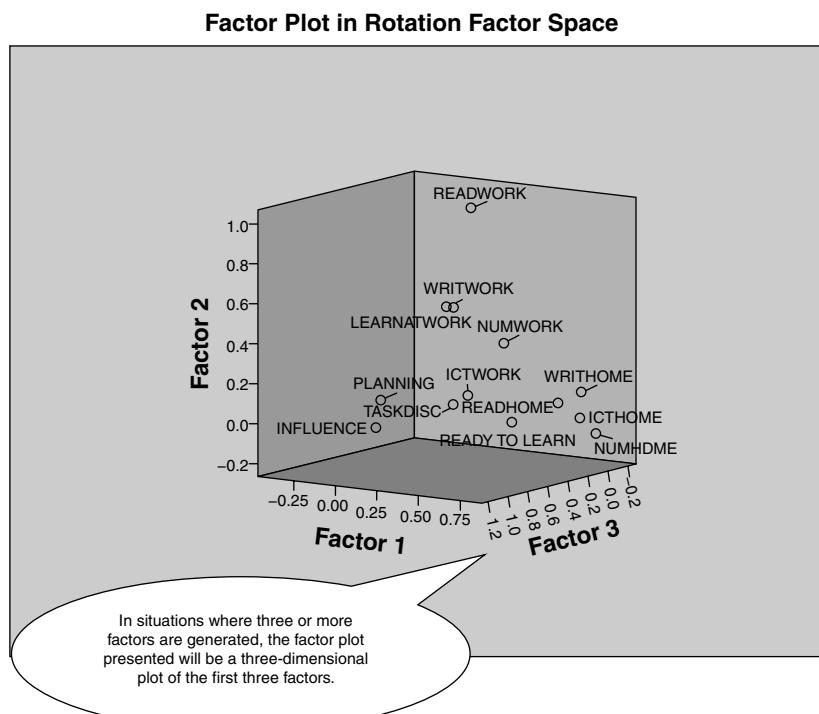
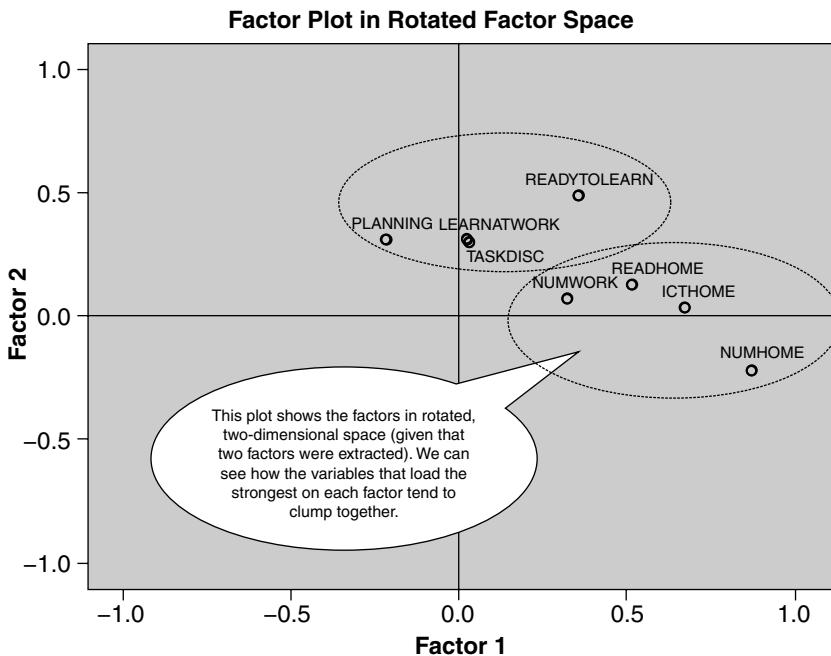
Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

The factor correlations are generated only when oblique rotation is selected (correlations with orthogonal rotations are set to zero). In this example, the correlations between factors are relatively weak which may warrant re-analysis assuming orthogonality.

■ TABLE 9.4 (continued)

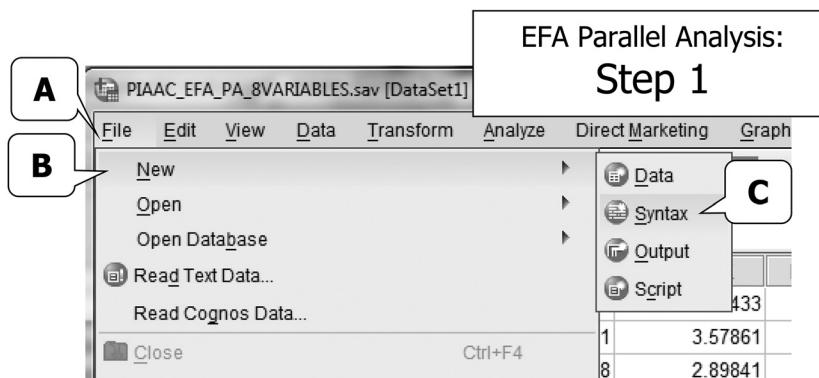
SPSS Results for the Exploratory Factor Analysis Example



9.3.1.1 SPSS Parallel Analysis for Determining Factor Retention

Next, we consider SPSS for conducting parallel analysis (PA). When you run the parallel analysis program, it is important that the data file is open so that the program will recognize that data file as the one with which to generate the PA. We will continue to work with the PIAAC_EFA.sav dataset.

PA Step 1. As mentioned previously, this is not available in the point-and-click user interface but can easily be performed with syntax available from O'Connor (2000). [Additional annotated code, along with syntax to generate artificial raw data that may be helpful for getting a feel for how it works, is accessible online at <https://people.ok.ubc.ca/briocconn/nfactors/rawpar.sps>.] To open a new syntax file, click on "File" then "New" then "Syntax." Following the screenshot below (see screenshot EFA Parallel Analysis: Step 1) produces the "Syntax Editor."



Step 2. It is most helpful to access an electronic copy of the article or the supplementary online material (see <https://people.ok.ubc.ca/briocconn/nfactors/rawpar.sps>) so that the syntax can be copied and pasted directly into the SPSS syntax viewer (however, it has also been provided in Table 9.5). When the syntax is copied into the syntax editor, the code that needs your input will clearly be displayed (see screenshot EFA Parallel Analysis: Step 2). These include the following:

- The GET line tells SPSS that the file currently open is the one that should be used to generate the parallel analysis. For ease, the dataset we are using, PIAAC_EFA .sav, is organized so that the eight variables we will factor analyze are grouped together. In this instance, the GET syntax is GET raw / FILE = * / missing=omit / VAR = NUMHOME to READYTOLEARN. Specifying FILE = * tells the program to read the SPSS data file that is open (thus make sure there is only one dataset open when you run the program, and the one that is open is the one from which you want the parallel analysis generated). The VAR = NUMHOME to READYTOLEARN. tells the program only to generate the parallel analysis on the variables within this range (in this illustration, it happens to be the first eight variables in the data file).

- The number of parallel datasets to compute needs to be defined (100 is the default and is an appropriate starting place): `compute ndatasets = 100`.
- The percentile must be specified (95th is common): `compute percent = 95`.
- The kind of parallel analysis to compute must be specified with 1 referring to PCA and 2 referring to principal axis/common factor analysis (which is what we will generate in this illustration): `compute kind = 2`.
- The type of distribution must be specified with 1 being normally distributed and 2 being permutations of the raw data: `compute randtype = 1`. It is important to note that the distributions of the observed variables remain the same during the parallel analysis procedure. As noted by O'Connor, “Permutations of the raw data set are thus highly accurate and most relevant, especially in cases where the raw data are not normally distributed or when they do not meet the assumption of multivariate normality” (see <https://people.ok.ubc.ca/briocconn/nfactors/raw-par.sps>). O'Connor recommends specifying normally distributed data first (i.e., `compute randtype = 1`) to get a general idea of the number of factors that the parallel analysis suggests retaining. Then specify distributions as permutations of the raw data (i.e., `compute randtype = 2`) with a small number of datasets (e.g., 100) to see how long the program takes to run. Assuming the time for running the program is doable, then run the parallel analysis program with the number of parallel data sets desired for your analyses (with 1,000 generally being sufficient).

**EFA Parallel Analysis:
Step 2**

```

File Edit View Data Transform Analyze Graphs Utilities Add-ons Run
set matrix.
* Enter the name/location of the data file for analyses after "FILE=".
GET
* Enter the desired ...
compute
* Enter the desired ...
compute
* Enter either
compute
* Enter either
compute
***** E...
compute
compute
* principal component...
do if
compute
compute
compute
compute
compute
loop
compute
1  set mxloops=9000 printback=off width=80 seed = 1953125.
2  matrix.
3  * Enter the name/location of the data file for analyses after "FILE=".
4  If you specify "FILE = *", then the program will read the current
5  active SPSS data file; Alternatively, enter the name/location
6  of a previously saved SPSS data file instead of "*".
7  you can use the "/ VAR=" subcommand after "/ missing=omit"
8  subcommand to select variables for the analyses.
9  GET raw / FILE = * / missing=omit / VAR = NUMHOME to READYTOLEARN.
10 * Enter the desired number of parallel data sets here.
11 compute ndatasets = 100.
12 * Enter the desired percentile here.
13 compute percent = 95.
14 * Enter either
15 1 for principal components analysis, or
16 2 for principal axis/common factor analysis.
17 compute kind = 2 .
18 * Enter either
19 1 for normally distributed random data generation parallel anal
20 2 for permutations of the raw data set.
21 compute randtype = 1.
22 **** End of user specifications. *****
23

```

It is important to remove the seed syntax when multiple generations of the code are conducted. If not, every run of the code will produce the same eigenvalues.

TABLE 9.5

SPSS Syntax for Generating Parallel Analysis

```
set mxloops=9000 printback=off width=80 seed = 1953125.  
matrix.  
* Enter the name/location of the data file for analyses after "FILE =";  
If you specify "FILE = *", then the program will read the current,  
active SPSS data file; Alternatively, enter the name/location  
of a previously saved SPSS data file instead of "*";  
you can use the "/ VAR =" subcommand after "/ missing=omit"  
subcommand to select variables for the analyses.
```

This tells SPSS to use the SPSS
data set that is currently open
and to generate the parallel
analysis only on variables
NUMHOME through
READYTOLEARN.

```
GET raw / FILE = * / missing=omit / VAR = NUMHOME to READYTOLEARN.
```

* Enter the desired number of parallel data sets here.

```
compute ndatsets = 100.
```

* Enter the desired percentile here.

```
compute percent = 95.
```

* Enter either

1 for principal components analysis, or

2 for principal axis/common factor analysis.

```
compute kind = 2 .
```

* Enter either

1 for normally distributed random data generation parallel analysis, or

2 for permutations of the raw data set.

```
compute randtype = 1.
```

This tells SPSS to generate 100
parallel datasets (100 is the default
and is a good starting place).

This tells SPSS to compute the 95th percentile.

This tells SPSS to compute principal
axis/common factor analysis.

This tells SPSS to generate normally
distributed random data.

***** End of user specifications. *****

```
compute ncases = nrow(raw).
```

```
compute nvars = ncol(raw).
```

```
* principal components analysis & random normal data generation.  
do if (kind = 1 and randtype = 1).  
compute nm1 = 1 / (ncases-1).  
compute vcv = nm1 * (sscp(raw) - ((t(csum(raw))*csum(raw))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute realeval = eval(d * vcv * d).
```

■ TABLE 9.5 (continued)

SPSS Syntax for Generating Parallel Analysis

```
compute evals = make(nvars,ndatsets,-9999).  
loop #nds = 1 to ndatsets.  
compute x = sqrt(2 * (ln(uniform(ncases,nvars)) * -1) ) &  
          cos(6.283185 * uniform(ncases,nvars) ).  
compute vcv = nm1 * (sscp(x) - ((t(csum(x))*csum(x))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute evals(:,#nds) = eval(d * vcv * d).  
end loop.  
end if.  
  
* principal components analysis & raw data permutation.  
do if (kind = 1 and randtype = 2).  
compute nm1 = 1 / (ncases-1).  
compute vcv = nm1 * (sscp(raw) - ((t(csum(raw))*csum(raw))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute realeval = eval(d * vcv * d).  
compute evals = make(nvars,ndatsets,-9999).  
loop #nds = 1 to ndatsets.  
compute x = raw.  
loop #c = 1 to nvars.  
loop #r = 1 to (ncases -1).  
compute k = trunc( (ncases - #r + 1) * uniform(1,1) + 1 ) + #r - 1.  
compute d = x(#r,#c).  
compute x(#r,#c) = x(k,#c).  
compute x(k,#c) = d.  
end loop.  
end loop.  
compute vcv = nm1 * (sscp(x) - ((t(csum(x))*csum(x))/ncases)).
```

■ TABLE 9.5 (continued)

SPSS Syntax for Generating Parallel Analysis

```
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute evals(:,#nd) = eval(d * vcv * d).  
end loop.  
end if.  
  
* PAF/common factor analysis & random normal data generation.  
do if (kind = 2 and randtype = 1).  
compute nm1 = 1 / (ncases-1).  
compute vcv = nm1 * (sscp(raw) - ((t(csum(raw))*csum(raw))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute cr = (d * vcv * d).  
compute smc = 1 - (1 &/ diag(inv(cr))).  
call setdiag(cr,smc).  
compute realeval = eval(cr).  
compute evals = make(nvars,ndatsets,-9999).  
compute nm1 = 1 / (ncases-1).  
loop #nd = 1 to ndatsets.  
compute x = sqrt(2 * (ln(uniform(ncases,nvars)) * -1) ) &*  
          cos(6.283185 * uniform(ncases,nvars) ).  
compute vcv = nm1 * (sscp(x) - ((t(csum(x))*csum(x))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute r = d * vcv * d.  
compute smc = 1 - (1 &/ diag(inv(r))).  
call setdiag(r,smc).  
compute evals(:,#nd) = eval(r).  
end loop.  
end if.  
  
* PAF/common factor analysis & raw data permutation.  
do if (kind = 2 and randtype = 2).  
compute nm1 = 1 / (ncases-1).  
compute vcv = nm1 * (sscp(raw) - ((t(csum(raw))*csum(raw))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute cr = (d * vcv * d).
```

■ TABLE 9.5 (continued)

SPSS Syntax for Generating Parallel Analysis

```
compute smc = 1 - (1 &/ diag(inv(cr)) ).  
call setdiag(cr,smc).  
compute realeval = eval(cr).  
compute evals = make(nvars,ndatsets,-9999).  
compute nm1 = 1 / (ncases-1).  
loop #nds = 1 to ndatsets.  
compute x = raw.  
loop #c = 1 to nvars.  
loop #r = 1 to (ncases -1).  
compute k = trunc( (ncases - #r + 1) * uniform(1,1) + 1 ) + #r - 1.  
compute d = x(#r,#c).  
compute x(#r,#c) = x(k,#c).  
compute x(k,#c) = d.  
end loop.  
end loop.  
compute vcv = nm1 * (sscp(x) - ((t(csum(x))*csum(x))/ncases)).  
compute d = inv(mdiag(sqrt(diag(vcv)))).  
compute r = d * vcv * d.  
compute smc = 1 - (1 &/ diag(inv(r)) ).  
call setdiag(r,smc).  
compute evals(:,#nds) = eval(r).  
end loop.  
end if.  
  
* identifying the eigenvalues corresponding to the desired percentile.  
compute num = rnd((percent*ndatsets)/100).  
compute results = { t(1:nvars), realeval, t(1:nvars), t(1:nvars) }.  
loop #root = 1 to nvars.  
compute ranks = rnkorder(evals(#root,:)).  
loop #col = 1 to ndatsets.  
do if (ranks(1,#col) = num).  
compute results(#root,4) = evals(#root,#col).  
break.  
end if.
```

■ TABLE 9.5 (continued)

SPSS Syntax for Generating Parallel Analysis

```
end loop.  
end loop.  
compute results(:,3) = rsum(evals) / ndatsets.  
  
print /title="PARALLEL ANALYSIS:".  
do if (kind = 1 and randtype = 1).  
print /title="Principal Components & Random Normal Data Generation".  
else if (kind = 1 and randtype = 2).  
print /title="Principal Components & Raw Data Permutation".  
else if (kind = 2 and randtype = 1).  
print /title="PAF/Common Factor Analysis & Random Normal Data Generation".  
else if (kind = 2 and randtype = 2).  
print /title="PAF/Common Factor Analysis & Raw Data Permutation".  
end if.  
compute specifs = {ncases; nvars; ndatsets; percent}.  
print specifs /title="Specifications for this Run:"  
/rlabels="Ncases" "Nvars" "Ndats" "Percent".  
print results  
/title="Raw Data Eigenvalues, & Mean & Percentile Random Data Eigenvalues"  
/clabels="Root" "Raw Data" "Means" "Prcntyle" /format "f12.6".  
  
do if (kind = 2).  
print / space = 1.  
print /title="Warning: Parallel analyses of adjusted correlation matrices".  
print /title="eg, with SMCs on the diagonal, tend to indicate more factors".  
print /title="than warranted (Buja, A., & Eyuboglu, N., 1992, Remarks on parallel".  
print /title="analysis. Multivariate Behavioral Research, 27, 509-540.)".  
print /title="The eigenvalues for trivial, negligible factors in the real".  
print /title="data commonly surpass corresponding random data eigenvalues".  
print /title="for the same roots. The eigenvalues from parallel analyses".  
print /title="can be used to determine the real data eigenvalues that are".  
print /title="beyond chance, but additional procedures should then be used".  
print /title="to trim trivial factors.".br/>print / space = 2.
```

■ TABLE 9.5 (continued)

SPSS Syntax for Generating Parallel Analysis

```
print /title="Principal components eigenvalues are often used to determine".
print /title="the number of common factors. This is the default in most".
print /title="statistical software packages, and it is the primary practice".
print /title="in the literature. It is also the method used by many factor".
print /title="analysis experts, including Cattell, who often examined".
print /title="principal components eigenvalues in his scree plots to
determine".
print /title="the number of common factors. But others believe this common".
print /title="practice is wrong. Principal components eigenvalues are based".
print /title="on all of the variance in correlation matrices, including both".
print /title="the variance that is shared among variables and the variances".
print /title="that are unique to the variables. In contrast, principal".
print /title="axis eigenvalues are based solely on the shared variance".
print /title="among the variables. The two procedures are qualitatively".
print /title="different. Some therefore claim that the eigenvalues from one".
print /title="extraction method should not be used to determine".
print /title="the number of factors for the other extraction method.".
print /title="The issue remains neglected and unsettled.".
end if.

compute root      = results(:,1).
compute rawdata = results(:,2).
compute percntyl = results(:,4).

save results /outfile= 'screedata.sav' / var=root rawdata means percntyl .

end matrix.
```

Step 3. Now that the syntax is created (see PA_PIAAC_n191.sps), run the program. For this data, we first generate 100 datasets using normally distributed data. Then we generate 1000 datasets using permutations of the raw data.

Interpreting the PA output. Annotated results are presented in Tables 9.6 and 9.7. More specifically, in Table 9.6 the results were generated using 100 datasets (compute ndatsets = 100.) with normally distributed random data (compute randtype = 1.). Table 9.7 results were generated using 1000 datasets (compute ndatsets = 1000.) with normally distributed random data (compute randtype = 2.). In both cases, we arrive at the same conclusion—two factors should be retained.

TABLE 9.6

SPSS Parallel Analysis Results for the Exploratory Factor Analysis Example: 100 Datasets with Normally Distributed Random Data

Run MATRIX procedure:

PARALLEL ANALYSIS:

PAF/Common Factor Analysis & Random Normal Data Generation

Specifications for this Run:
 Ncases 191
 Nvars 8
 Ndsets 100
 Percent 95

To determine the number of factors to retain, compare the eigenvalues derived from the observed raw data to the eigenvalues derived from the parallel analysis computation. When the t th eigenvalue from the observed data is greater than the t th eigenvalue from the random or permuted data in the parallel analysis, then those factors are retained.

Here, we see the first two 'raw data' eigenvalues are greater than the random data **mean** and **percentile** eigenvalues indicating that two factors should be retained.

Raw Data Eigenvalues, & Mean	Root	Raw Data	Means	Random Data Eigenvalues	Prcntile
1.000000	1.613687		.359535	.472506	
2.000000	.463109		.227939	.313412	
3.000000	.041523		.132743	.200257	
4.000000	.040669		.054752	.109795	
5.000000	-.085385		-.014786	.029500	
6.000000	-.088389		-.085607	-.037489	
7.000000	-.193077		-.155585	-.107658	
8.000000	-.256326		-.231419	-.171749	

Warning: Parallel analyses of adjusted correlation matrices eg, with SMCs on the diagonal, tend to indicate more factors than warranted (Buja, A., & Eyuboglu, N., 1992, Remarks on parallel analysis. Multivariate Behavioral Research, 27, 509-540.). The eigenvalues for trivial, negligible factors in the real data commonly surpass corresponding random data eigenvalues for the same roots. The eigenvalues from parallel analyses can be used to determine the real data eigenvalues that are beyond chance, but additional procedures should then be used to trim trivial factors.

■ TABLE 9.6 (continued)

SPSS Parallel Analysis Results for the Exploratory Factor Analysis Example: 100 Datasets with Normally Distributed Random Data

Principal components eigenvalues are often used to determine the number of common factors. This is the default in most statistical software packages, and it is the primary practice in the literature. It is also the method used by many factor analysis experts, including Cattell, who often examined principal components eigenvalues in his scree plots to determine the number of common factors. But others believe this common practice is wrong. Principal components eigenvalues are based on all of the variance in correlation matrices, including both the variance that is shared among variables and the variances that are unique to the variables. In contrast, principal axis eigenvalues are based solely on the shared variance among the variables. The two procedures are qualitatively different. Some therefore claim that the eigenvalues from one extraction method should not be used to determine the number of factors for the other extraction method.

The issue remains neglected and unsettled.

----- END MATRIX -----

■ TABLE 9.7

SPSS Parallel Analysis Results for the Exploratory Factor Analysis Example: 1000 Datasets with Permutations of the Raw Data

Run MATRIX procedure:

PARALLEL ANALYSIS:

PAF/Common Factor Analysis & Raw Data Permutation

Specifications for this Run:

Ncases 191
Nvars 8
Ndatasets 1000
Percent 95

To determine the number of factors to retain, compare the eigenvalues derived from the observed raw data to the eigenvalues derived from the parallel analysis computation. When the i th eigenvalue from the observed data is greater than the i th eigenvalue from the random or permuted data in the parallel analysis, then those factors are retained.

Here, we see the first two 'raw data' eigenvalues are greater than the permuted data **mean** and **percentile** eigenvalues indicating that two factors should be retained.

Root	Raw Data Eigenvalues, & Mean		& Percentile Random Data Eigenvalues	
	Raw Data	Means	Prctyle	
1.000000	{ 1.613687 }	.361951	.478716	
2.000000	.463109	.237016	.326682	
3.000000	.041523	.137709	.207489	
4.000000	.040669	.057641	.114366	
5.000000	-.085385	-.015422	.033313	
6.000000	-.088389	-.085530	-.039254	
7.000000	-.193077	-.157236	-.108595	
8.000000	-.256326	-.237790	-.184679	

Warning: Parallel analyses of adjusted correlation matrices

eg, with SMCs on the diagonal, tend to indicate more factors

than warranted (Buja, A., & Eyuboglu, N., 1992, Remarks on parallel

analysis. Multivariate Behavioral Research, 27, 509-540.).

The eigenvalues for trivial, negligible factors in the real data commonly surpass corresponding random data eigenvalues for the same roots. The eigenvalues from parallel analyses can be used to determine the real data eigenvalues that are beyond chance, but additional procedures should then be used to trim trivial factors.

Principal components eigenvalues are often used to determine the number of common factors. This is the default in most statistical software packages, and it is the primary practice

■ TABLE 9.7 (continued)

SPSS Parallel Analysis Results for the Exploratory Factor Analysis Example: 1000 Datasets with Permutations of the Raw Data

in the literature. It is also the method used by many factor analysis experts, including Cattell, who often examined principal components eigenvalues in his scree plots to determine the number of common factors. But others believe this common practice is wrong. Principal components eigenvalues are based on all of the variance in correlation matrices, including both the variance that is shared among variables and the variances that are unique to the variables. In contrast, principal axis eigenvalues are based solely on the shared variance among the variables. The two procedures are qualitatively different. Some therefore claim that the eigenvalues from one extraction method should not be used to determine the number of factors for the other extraction method.

The issue remains neglected and unsettled.

----- END MATRIX -----

9.3.2 Computing EFA With Ordinal Data Using SPSS

Next we consider an SPSS add-on, categorical principal components analysis (CATPCA), for conducting exploratory factor analysis in the case where our data is ordinal. I felt it critically important to provide this illustration in the textbook for two reasons: (1) there is an abundance of data collected and secondary data available that is ordinal in scale—specifically Likert items that measure attitude, perceptions, etc.—as well as nominal (which can also be handled with CATPCA); and (2) yet few resources are available that transparently help researchers select and use an appropriate EFA procedure and thereby avoid the pitfall of applying conventional EFA techniques to data for which it is really not appropriate. As an optimal scaling approach, nonlinear relationships between categorical variables can be modeled within CATPCA via optimal quantification in a specified dimension. *Unfortunately, CATPCA is only available as an add-on with SPSS.* Should you be renting a copy of SPSS, you likely have it. If you are accessing SPSS from an institution that purchases a finite number of SPSS licenses, you may or may not have it.

Before we conduct the analysis, let us talk about the data. The data we are using is the 2010 Survey of Doctorate Recipients (SDR, <http://www.nsf.gov/statistics/srvydoctoratework/>), available through the National Science Foundation (NSF) (2010). Thank you to NSF for making this data publicly available. This is only one of many secondary data sources available through NSF as well as other federal and nonfederal agencies,

and I encourage you to explore these extremely rich resources (particularly for multi-variate research) for your own research.

First conducted in 1973, the SDR is a “longitudinal biennial survey that provides demographic and career history information about individuals with a research doctoral degree in a science, health, or engineering (SHE) field from a U.S. academic institution. The survey follows a sample of individuals with SHE doctorates throughout their careers from the year of their degree until age 76 . . . Results are used to make decisions related to the educational and occupational achievements and career movements of the nation’s doctoral scientists and engineers” (see <http://www.nsf.gov/statistics/srvydoctoratework/>). It is important to note that the SDR is a complex sample (i.e., not a simple random sample). More specifically, it employs a stratified probability sampling design. As you will see in the dataset, the very last variable is a weight variable. When this weight is applied to the analysis, the results are adjusted for unequal selection probabilities and nonresponse (which also includes respondents who could not be located or whose eligibility was unknown) and aligned with poststratification (<http://www.nsf.gov/statistics/srvydoctoratework/>). The end results are then representative of the intended population. As stated previously, the purpose of the text is not to serve as a primer for understanding complex samples, and thus readers interested in learning more about complex survey designs are referred to resources noted earlier in the chapter.

Now, let’s review the data. We are using the SDR2010_POSTDOC.sav file. This data file has been delimited to include only individuals who completed the SDR in 2010 who were employed in a post-doctoral position during the week they responded to the survey ($n = 1080$). The size of this sample is more than sufficient to generate EFA but at the same time small enough to work with for readers who may be using a version of SPSS that limits the number of cases to 1,500. (Note: The complete SDR data file, which includes 31,362 cases, is available from the textbook’s companion website and is titled SDR2010_NSF.sav.)

You’ll notice that in both the post-doc (SDR2010_POSTDOC.sav) and full data file (SDR2010_NSF.sav) there is quite a bit of recoding that will need to be performed in order to get the data in shape for analysis. This includes defining the missing values, recoding the string variables to numeric, and where applicable, reverse coding. I’ve taken the liberty to perform this data cleaning for the variables with which we’ll be working; however, the remaining variables in the data file have been left as is so that you may practice your data cleaning skills in working with ‘real data.’

Let’s look at the data. The first variable in our dataset was the variable used to delimit the cases to only respondents who were employed in post-doctoral positions during the week of the survey. The next nine variables are those variables with which we will be analyzing for the EFA. Each row in the data set still represents one individual. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the respondents were measured. For the EFA illustration, we will work with the nine ordinal satisfaction measures.

	ACAD_POSTDOC	SAT_ADV	SAT_BEN	SAT_CHAL	SAT_IND	SAT_LOC	SAT_RESP	SAT_SAL	SAT_SEC	SAT_SOC
1	1.00	3.00	3.00	4.00	4.00	4.00	4.00	2.00	3.00	4.00
2	1.00	4.00	4.00	4.00	4.00	3.00	3.00	3.00	4.00	4.00
3	1.00	2.00	3.00	3.00	3.00	3.00	3.00	2.00	2.00	3.00
4	1.00	3.00	4.00	4.00	4.00	3.00	4.00	3.00	4.00	4.00
5	1.00	4.00	2.00	4.00	4.00	4.00	3.00	3.00	4.00	4.00
6	1.00	4.00	3.00	4.00	3.00	4.00	4.00	3.00	4.00	4.00
7	1.00	4.00	4.00	4.00	1.00	4.00	1.00	4.00	4.00	1.00
8	1.00	3.00	3.00	3.00	4.00	4.00	3.00	3.00	3.00	3.00
9	1.00	4.00	4.00	4.00	4.00	4.00	4.00	4.00	4.00	4.00

Each of the nine variables are ordinal with a four-point Likert scale such that 4 represents 'very satisfied,' 3 is 'somewhat satisfied,' 2 is 'somewhat dissatisfied,' and 1 is 'very dissatisfied.'

Reviewing the annotated SDR questionnaire available from NSF (see screenshot of questionnaire), this is a question set responding to the item, "Thinking about your principal job held during the week of October 1, please rate your satisfaction with that job's . . ." The components of the job to which they responded included (1) salary, (2) benefits, (3) job security, (4) job location, (5) opportunities for advancement, (6) intellectual challenge, (7) level of responsibility, (8) degree of independence, and (9) contribution to society.

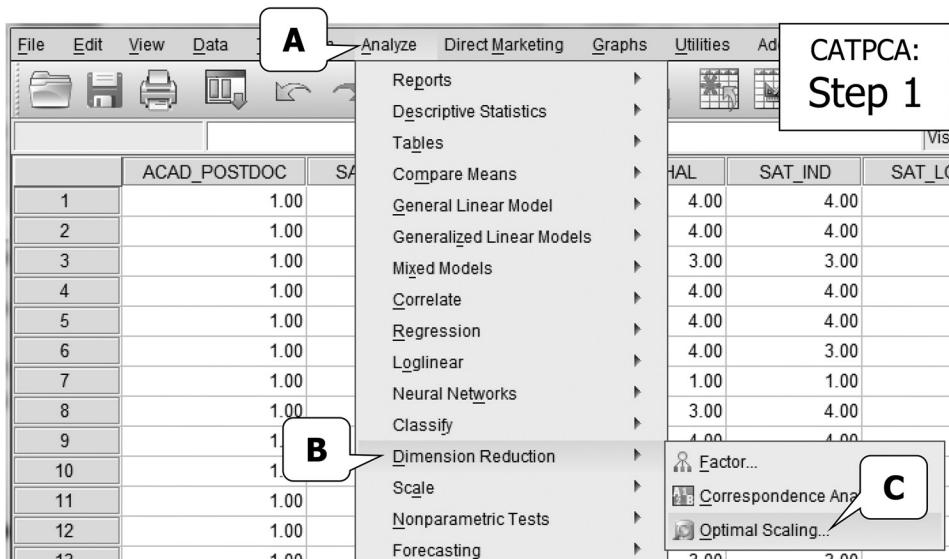
A34. Thinking about your principal job held during the week of October 1, please rate your satisfaction with that job's...

Mark one answer for each item.

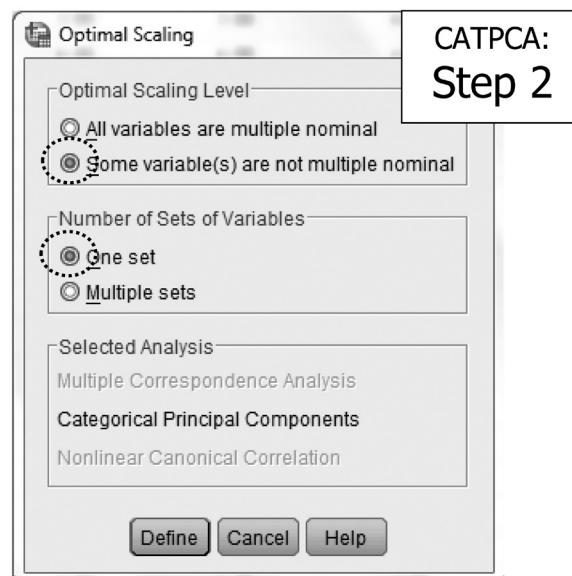
	Very satisfied	Somewhat satisfied	Somewhat dissatisfied	Very dissatisfied
1 Salary	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATSAL*
2 Benefits.....	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATBEN*
3 Job security	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATSEC*
4 Job location	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATLOC*
5 Opportunities for advancement.....	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATADV*
6 Intellectual challenge	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATCAHL*
7 Level of responsibility	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATRESP*
8 Degree of independence	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATIND*
9 Contribution to society	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/> SATSOC*

These are the column header names for the original variables in the SDR. The recoded variables which we will use include an underscore to separate SAT from the item descriptor (e.g., SAT_SAL).

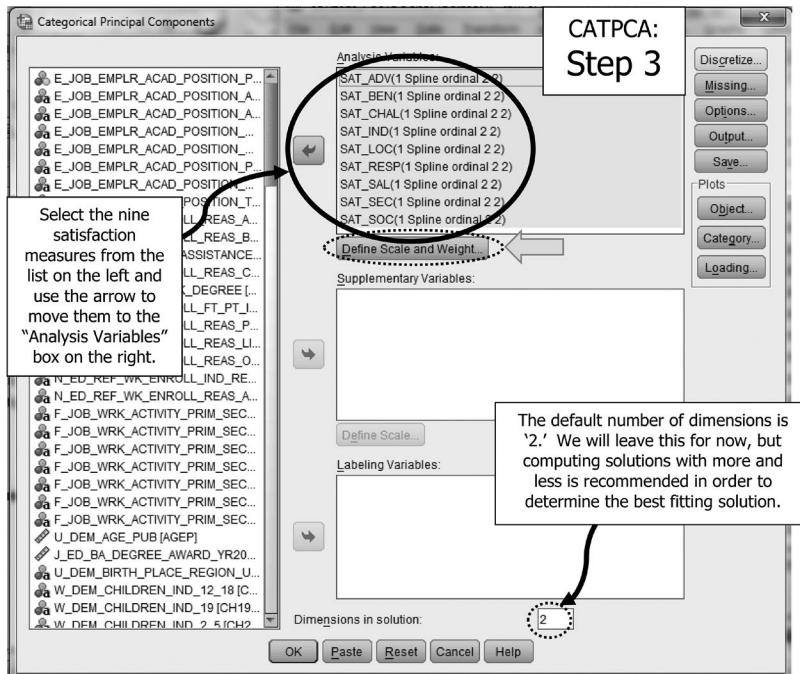
Step 1. To conduct categorical principal components analysis, go to “Analyze” in the top pull-down menu, then select “Dimension Reduction,” and then select “Optimal Scaling.” *Again, please remember that if you do not have this SPSS add-on, you will not see an option for “Optimal Scaling.”* Following the screenshot below (CATPCA: Step 1) produces the “Factor Analysis” dialog box.



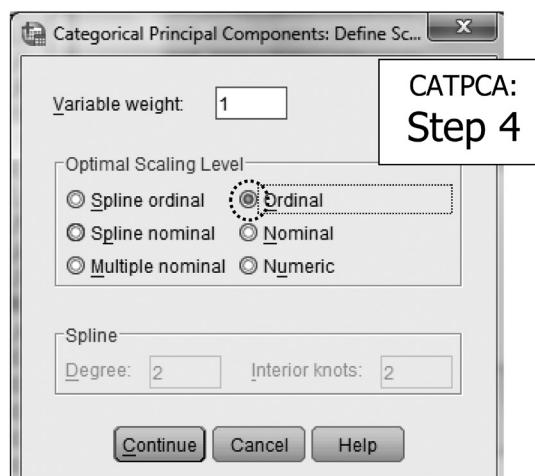
Step 2. Select the radio button for “Some variable(s) are not multiple nominal” (had all the variables been nominal, we would have selected the first option) and “one set” (see screenshot CATPCA: Step 2).



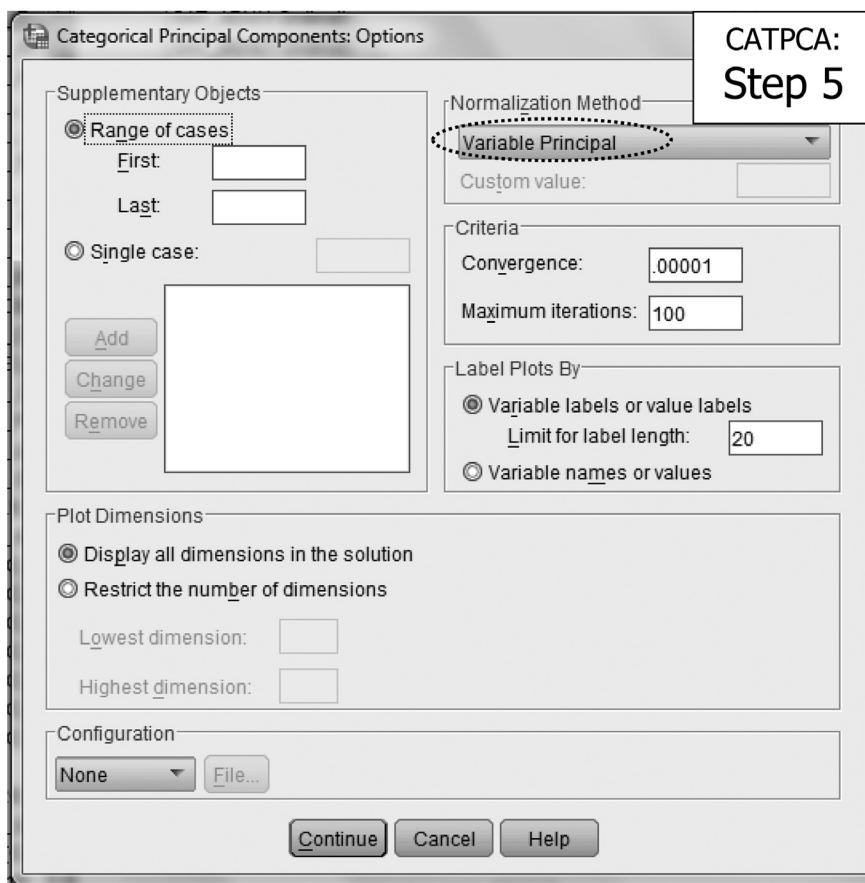
Step 3. Click the nine satisfaction measures and move into the "Analysis Variables" box by clicking the arrow button (see screenshot CATPCA: Step 3).



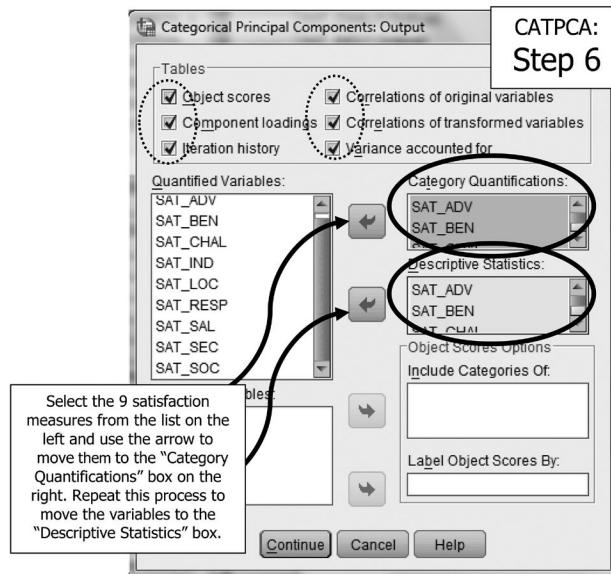
Step 4. Click in "Define scale and weight" (displayed under the Analysis Variable box) to change the optimal scaling level (see screenshot CATPCA: Step 4). The default is 'spline ordinal.' We will select the radio button for 'ordinal.' Spline ordinal and ordinal optimal scaling levels are similar in that they both preserve the order of the categories in the optimally scaled variable. Ordinal optimal scaling results in a better fitting transformation than spline ordinal but is less smooth. Click Continue to return to the main CATPCA page.



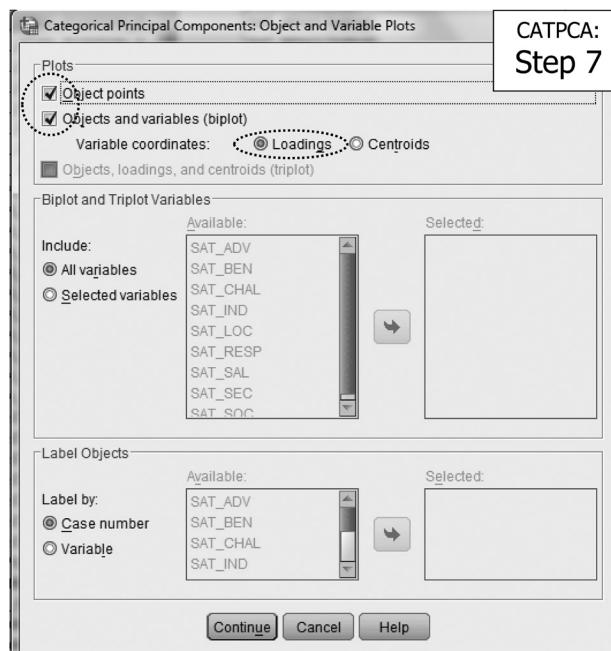
Step 5. From the CATPCA page (see screenshot CATPCA: Step 3), click on Options to bring up the Options dialog box. We will leave all default selections as is on this page (see screenshot CATPCA: Step 5). In terms of the Normalization Method, the default selection is Variable Principal. This method optimizes the relationship between variables and is an appropriate selection if the correlation between variables is your primary interest. Click Continue to return to the main CATPCA page.



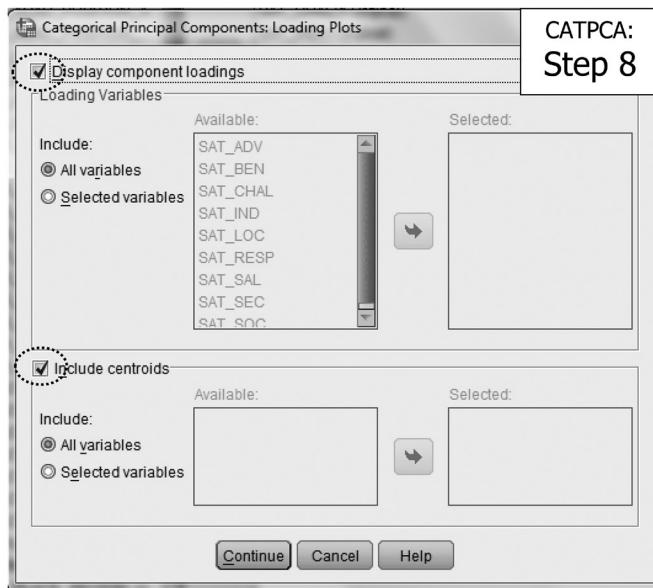
Step 6. From the CATPCA page (see screenshot CATPCA: Step 3), click on Output to bring up the Output dialog box (see screenshot CATPCA: Step 6). Object scores and Component loadings should already be selected, and we will keep those selected. Place a checkmark for the remaining tables including Iteration history, Correlations of original variables, Correlations of transformed variables, and Variance accounted for. Click Continue to return to the main CATPCA page. Move all the satisfaction variables from the Quantified Variables list to the Category Quantifications box by clicking the arrow in the middle. Repeat this process to move the variables to the Descriptive Statistics box. Click Continue to return to the main CATPCA page.



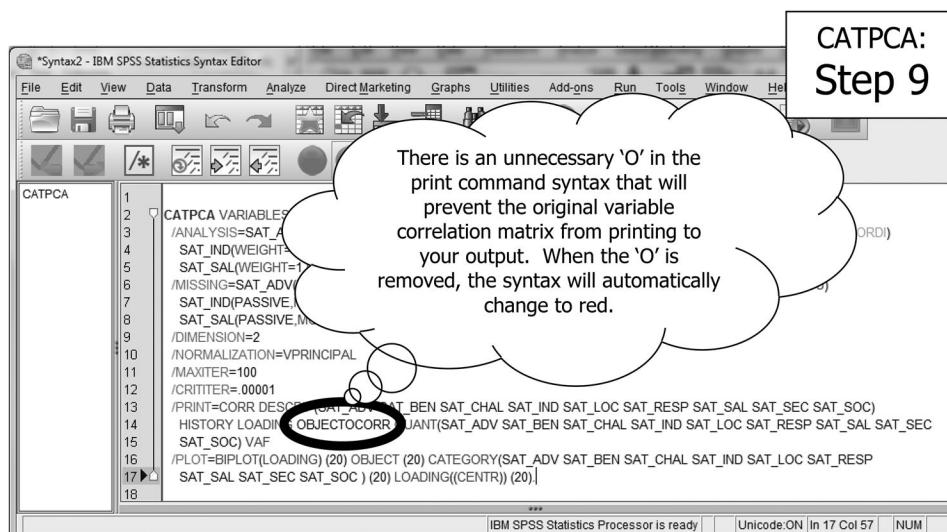
Step 7. From the CATPCA page (see screenshot CATPCA: Step 3), click on Object (listed under Plots in the right navigational menu) to bring up the Object and Variable Plots dialog box (see screenshot CATPCA: Step 7). Object points should already be selected, and we will keep that option selected. We will place a checkmark for Objects and variables (biplot) with variable coordinates Loadings. We will keep the default options selected for Biplot and Triplot Variables and Label Objects. Click Continue to return to the main CATPCA page.



Step 8. From the CATPCA page (see screenshot CATPCA: Step 3), click on Loading (listed under Plots in the right navigational menu) to bring up the Loading Plots dialog box (see screenshot CATPCA: Step 8). Display component loadings should already be selected, and we will keep that option selected. We will place a checkmark for Include Centroids. Click Continue to return to the main CATPCA page.



Step 9. From the CATPCA page (see screenshot CATPCA: Step 3), click 'paste' to open the syntax created from the commands just generated (see screenshot CATPCA:



Step 9). In some versions of SPSS, an error occurs in the print command line in specifying to print the original variable correlation matrix such that an ‘O’ is included rather than a space. If this error occurs in your syntax, remove the ‘O’ and then run the syntax to generate the output.

When the erroneous ‘O’ is removed (you must manually do this using your delete or backspace key), ‘OBJECT’ and ‘CORR’ will appear in red, indicating that the original variable correlation matrix will be printed in the output.

Interpreting the CATPCA output. Annotated results are presented in Table 9.8.

■ **TABLE 9.8**

SPSS Results for the Categorical Principal Components Analysis Example

Credit	
CATPCA	
Version 1.1	
by	
Data Theory Scaling System Group (DTSS)	
Faculty of Social and Behavioral Sciences	
Leiden University, The Netherlands	

Descriptive Statistics

Case Processing Summary

Case Processing Summary	
Valid Active Cases	1080
Active Cases with Missing Values	0
Supplementary Cases	0
Total	1080
Cases Used in Analysis	1080

F_JOB_SATISFACTION_ADVANCEMENT ^a		Frequency
Valid	Very dissatisfied	151
	Somewhat dissatisfied	310
	Somewhat satisfied ^b	414
	Very satisfied	205
	Total	1080

a. Optimal Scaling Level: Ordinal.

b. Mode.

For illustrative purposes, only one descriptive table is presented. However, the output includes descriptive stats associated with each variable.

This table presents the frequencies for each category of the variables included in the model.

TABLE 9.8 (continued)

SPSS Results for the Categorical Principal Components Analysis Example

Iteration Number	Iteration History					
	Variance Accounted For		Loss			Restriction of Centroid to Vector Coordinates
	Total	Increase	Total	Centroid Coordinates		
0 ^a	4.760364	.000005	13.239636	13.193188	.046447	
1	4.777268	.016904	13.222732	13.193188	.029544	
2	4.787877	.010609	13.212123	13.182787	.029335	
3	4.790623	.002746	13.209377	13.180107	.029270	
4	4.791586	.000962	13.208414	13.179203	.029212	
5	4.792009	.000423	13.207991	13.178830	.029161	
6	4.792228	.000219	13.207772	13.178646	.029126	
7	4.792353	.000126	13.207647	13.178540	.029107	
8	4.792430	.000077	13.207570			Eigenvalues ('variance accounted for') for each iteration are presented. 'Iteration 0' represents the solution that would have been evidenced from a <i>conventional</i> principal components analysis (i.e., not taking into consideration the measurement scale of the variables). The larger eigenvalue for the CATPCA solution (4.77, beginning with iteration 1) reflects a slightly better solution with CATPCA (which takes the ordinal measurement scale into account in the modeling) as compared to the conventional PCA.
9	4.792480	.000049	13.207520			
10	4.792512	.000033	13.207488			
11	4.792534	.000022	13.207466			
12	4.792549	.000015	13.207451			
13	4.792559	.000010	13.207441			
14 ^b	4.792566	.000007	13.207434			

a. Iteration 0 displays the statistics of the solution with all variables, except level Multiple Nominal, treated as numerical.
b. The iteration process stopped because the convergence test value was

Dimension	Model Summary			
	Cronbach's Alpha	Variance Accounted For		
		Total (Eigenvalue)	% of Variance	
1	.793	3.388	37.643	
2	.324	1.405	15.608	
Total	.890 ^a	4.793	53.251	

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

Cronbach's alpha is a measure of internal consistency, and this value is provided for each dimension (i.e., factor; labeled '1' and '2') as well as the total which represents the combination of all factors. The '% of variance' is presented for each factor and for the combination of factors (i.e., total). Both factors account for 53% of the variance in the optimally scaled items.

TABLE 9.8 (continued)

SPSS Results for the Categorical Principal Components Analysis Example

Quantifications**Table**

Category	Frequency	Quantification	Centroid Coordinates		Vector Coordinates	
			Dimension		Dimension	
			1	2	1	2
Very dissatisfied	151	-1.920	-1.232	-.387	-1.254	-.291
Somewhat dissatisfied	310	-.478	-.317	-.053	-.312	-.073
Somewhat satisfied	414	.371	.222	.143	.242	.056
Very satisfied	205	1.389	.938	.077	.907	.211

Variable Principal Normalization.

a. Optimal Scaling Level: Ordinal.

For illustrative purposes, only one descriptive table is presented. However, the output includes descriptive stats associated with each

'Centroid coordinates' reflect the average of all object scores for cases for the respective category on each factor (labeled 'dimension'). 'Vector coordinates' are the coordinates for each category when the categories are represented by a straight line between factor 1 (X axis) and factor 2 (Y axis) in a scatterplot.

The correlation matrix reflects coefficients *after* optimal scaling has been performed. These coefficients are those used in the CATPCA. If you imputed data during the CATPCA procedure, these values would reflect correlations of imputed values.

Correlations Transformed Variables

	F_JOB_SATISFACTION_ADVANCEMENT	F_JOB_SATISFACTION_BENEFITS	F_JOB_SATISFACTION_CHALLENGE	F_JOB_SATISFACTION_INDEPENDENCE	F_JOB_SATISFACTION_LOCATION	F_JOB_SATISFACTION_RESPONSIBILITY	F_JOB_SATISFACTION_SALARY	F_JOB_SATISFACTION_SECURITY	F_JOB_SATISFACTION_SOCIETY
F_JOB_SATISFACTION_ADVANCEMENT	1.000	.163	.362	.316	.227	.392	.309	.426	.333
F_JOB_SATISFACTION_BENEFITS	.163	1.000	.077	.153	.157	.149	.428	.263	.100
F_JOB_SATISFACTION_CHALLENGE	.362	.077	1.000	.502	.220	.615	.109	.232	.495
F_JOB_SATISFACTION_INDEPENDENCE	.316	.153	.502	1.000	.207	.622	.183	.310	.434
F_JOB_SATISFACTION_LOCATION	.227	.157	.220	.207	1.000	.210	.118	.195	.177
F_JOB_SATISFACTION_RESPONSIBILITY	.392	.149	.615	.622	.210	1.000	.211	.348	.498
F_JOB_SATISFACTION_SALARY	.309	.428	.109	.183	.118	.211	1.000	.272	.093
F_JOB_SATISFACTION_SECURITY	.426	.263	.232	.310	.195	.348	.272	1.000	.241
F_JOB_SATISFACTION_SOCIETY	.333	.100	.495	.434	.177	.498	.093	.241	1.000
Dimension	1	2	3	4	5	6	7	8	9
Eigenvalue	3.388	1.405	.887	.808	.641	.600	.489	.448	.334

■ TABLE 9.8 (continued)

SPSS Results for the Categorical Principal Components Analysis Example

Objects

Object Scores

Case Number	Dimension	
	1	2
1	.906	-.667
2	1.347	.489
3	-.711	.148
4	1.150	.379
5	.971	-.213
6	1.108	.334
7	-2.188	4.974
8	.136	.651
9	1.263	1.367
10	.927	.613
....
1080	1.330	.986

For illustrative purposes, only a portion of the output is presented.

The object scores represent the coordinates associated with the respective case for each factor. Thus, there will be as many *cases* listed in the table as your sample size.

Variable Principal Normalization.

Component Loadings

Component Loadings

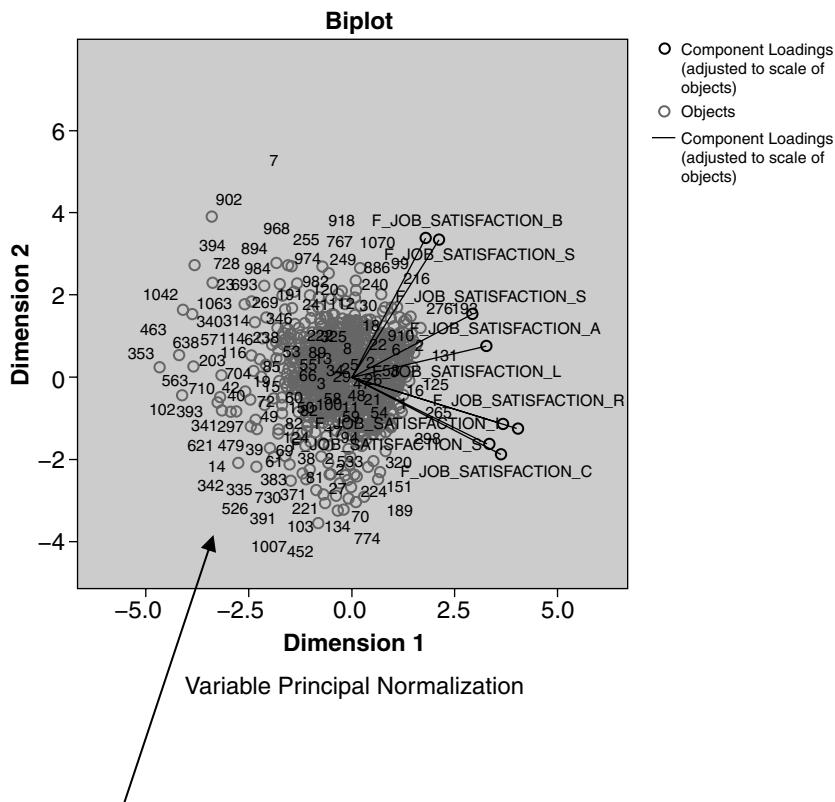
	Dimension	
	1	2
F_JOB_SATISFACTION_ADVANCEMENT	.653	.152
F_JOB_SATISFACTION_BENEFITS	.359	.684
F_JOB_SATISFACTION_CHALLENGE	.723	-.374
F_JOB_SATISFACTION_INDEPENDENCE	.735	-.225
F_JOB_SATISFACTION_LOCATION	.404	.104
F_JOB_SATISFACTION_RESPONSIBILITY	.804	-.248
F_JOB_SATISFACTION_SALARY	.422	.670
F_JOB_SATISFACTION_SECURITY	.587	.310
F_JOB_SATISFACTION_SOCIETY	.662	-.326

The component loadings scatterplot graphs the coordinates for each variable on each factor, allowing us to see how the variables relate to each other as well as the factors. Variables that lump together suggest distinguishable factors. In this case, we see that seven of the variables coalesce together on factor 1 and only two on factor 2.

Variable Principal Normalization.

■ TABLE 9.8 (continued)

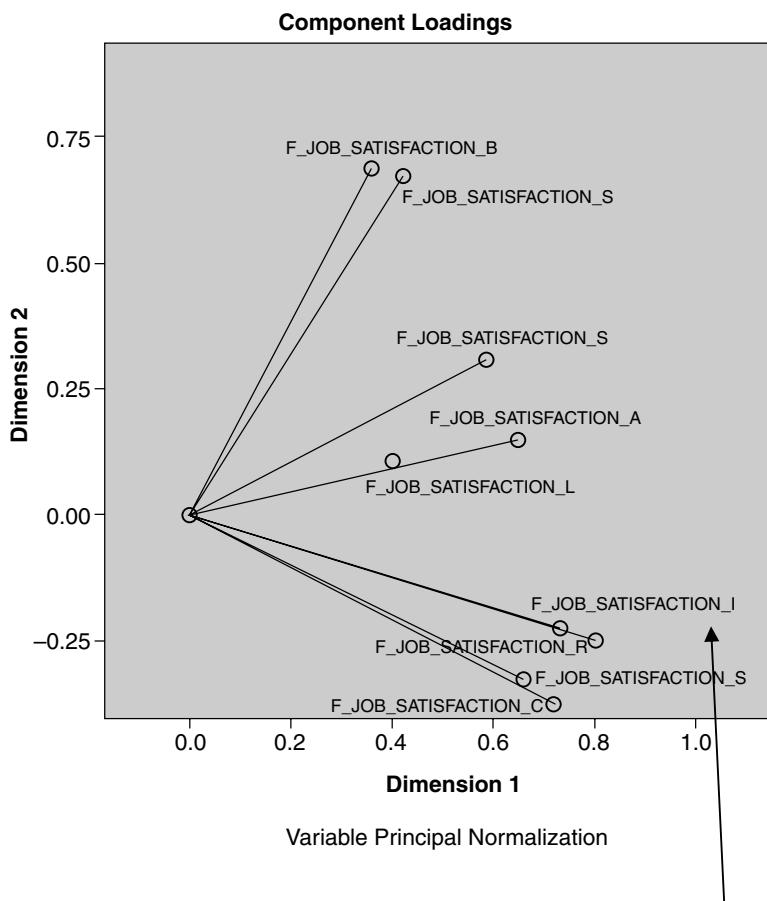
SPSS Results for the Categorical Principal Components Analysis Example



Finally, we get a scatterplot that you will see in color but is presented in grayscale here. Each variable is black and each case is green (grayscale here). Factor 1 is able to capture a bit more of the variance among the variables and cases and thus can explain the variance a bit better than factor 2 we see the variables and cases to be more tightly grouped (-4 to 4 for factor 2 as compared to -5 to 5 for factor 1). This suggests less variable variance captured.

■ TABLE 9.8 (continued)

SPSS Results for the Categorical Principal Components Analysis Example

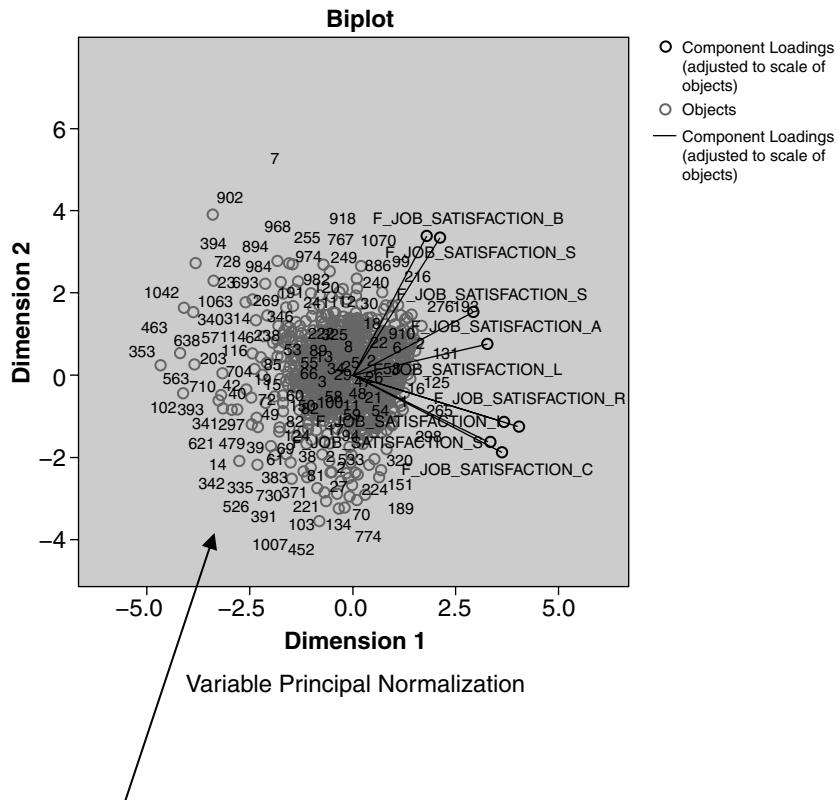


The component loadings scatterplot graphs the coordinates for each variable on each factor, allowing us to see how the variables relate to each other as well as the factors. Variables that lump together suggest distinguishable factors. In this case, we see the variables vary substantially along dimension 2 (i.e., factor 2) but tend to fall within a more narrow range of dimension 1 (i.e., factor 1) (between about .40 and .80). Here we can see how the two variables with large loadings on factor 2 are differentiating from those of factor 1. This may be where a decision is made to remove the two variables that appear to load on factor 2 and re-run. If the model improves without those items, there will be a clearer, tighter grouping of the variables on their respective factor(s).

The lines from the centroid to each variable are eigenvectors and the variable is at the eigenvalue for its vector. Thus the eigenvalue is a distance point along an eigenvector. With conventional EFA, a rotation strategy is applied to make interpretation easier. Here, we can imagine rotation such that both dimensions are rotated counterclockwise 45 degrees. In doing so, the axis of each factor (or dimension) would be going through a cloud of points (which represent the variables).

■ TABLE 9.8 (continued)

SPSS Results for the Categorical Principal Components Analysis Example



Finally, we get a scatterplot that you will see in color but is presented in grayscale here. Each variable is black and each case is green (grayscale here). Factor 1 is able to capture a bit more of the variance among the variables and cases and thus can explain the variance a bit better than factor 2 we see the variables and cases to be more tightly grouped (-4 to 4 for factor 2 as compared to -5 to 5 for factor 1). This suggests less variable variance captured.

9.4 DATA SCREENING

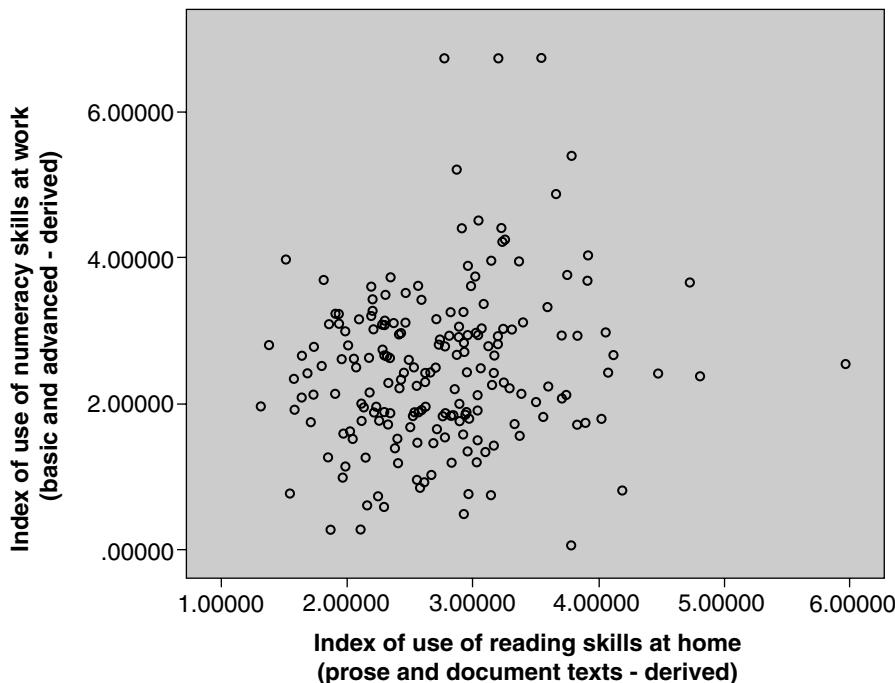
As you may recall, there were a number of assumptions associated with conventional exploratory factor analysis. These included (a) independence, (b) linearity, (c) absence of outliers (both univariate and multivariate), and (d) lack of extreme multicollinearity and singularity. Although fixed values of X were discussed in assumptions, this is not an assumption that will be tested, but is instead related to the use of the results (i.e., extrapolation and interpolation).

9.4.1 Independence

Testing for this assumption is a bit nebulous in exploratory factor analysis, as there are no independent and dependent variables that allow for this type of examination. In the absence of statistical evidence, we will rely on theoretical evidence: If the units have been randomly sampled from a population, there is evidence that the assumption of independence has been met.

9.4.2 Linearity

Linearity is an important assumption since correlation matrices underlie conventional EFA. You may recall that when you studied bivariate correlations, as well as simple and multiple regression, that scatterplots were one way that linearity could be examined. We will again use scatterplots to visually assess linearity. The challenge with EFA, as compared to other procedures where scatterplots have been applied, is the large (and often very large) number of variables, which makes review of all possible pairs of variables quite daunting and an inefficient use of your time. For example, with 10 variables (which tends to be toward the lower limit of the number of variables often applied to EFA), there are $[10(10 - 1)]/2$ or 45 different pairwise combinations of the variables, and with 20 variables there are nearly 200 combinations! One work-around for this is to generate and examine a few random scatterplots, assuming that these are representative of the entire population of scatterplots. Don't be surprised if the scatterplots do not provide picture-perfect linear relationships, and don't be ready to discard or transform variables if that is indeed the case—those consequences should be reserved only for cases where obvious curvilinearity is observed. For the PIAAC data, I ran a number of bivariate scatterplots and while not all scatterplots suggested a strong linear relationship, there does not appear to be evidence to suggest curvilinear associations. One example, graphing ‘index of use of numeracy skills at work (basic and advanced)’ with ‘index of use of reading skills at home (prose and document texts),’ is presented here:



9.4.3 Absence of Outliers

As discussed previously, factor analysis is quite robust to violations of the assumption of normality except where tests of inference are used to determine the number of factors to retain, and in this case, multivariate normality *is* an assumption. For this illustration, we are using maximum likelihood and thus will be thorough in our examination of multivariate normality.

We can examine univariate normality tests, which are less sensitive than multivariate normality tests, through skewness and kurtosis, formal tests of normality, and plots (e.g., Q-Q plots). Multivariate outliers are evidenced by statistically significant Mahalanobis distance scores ($\alpha = .001$ if you tend toward the liberal edge, which is appropriate with EFA), evaluated using a chi-square distribution with degrees of freedom equal to the number of variables. To generate Mahalanobis distance, we will generate multiple regression, applying all the variables as independent variables with the dependent variable being a binary variable coded 1 for potential outliers and 0 for all other variables within the model. The process for examining outliers is therefore to look for univariate outliers first. If any are detected, then screen for multivariate outliers. In terms of multivariate normality, a macro in SPSS (DeCarlo, 1997) (illustrated in the MANOVA chapter) can also be used to examine a number of multivariate normality indices including

(a) multivariate kurtosis (Mardia, 1970), (b) multivariate skewness and kurtosis based on Small's (1980) multivariate extension of univariate skewness and kurtosis (Looney, 1995), (c) multivariate normality omnibus test (Looney, 1995), (d) largest squared and plot of squared Mahalanobis distance, and (e) critical values for hypothesis test for a single multivariate outlier using Mahalanobis distance (Penny, 1996).

Additionally, not only are we concerned with outlying cases, but we are also concerned with outlying variables and will need to examine our data for both. These outlying variables, which can be removed from the model if/when identified, can be determined by examination of the following: (a) squared multiple correlations with all other variables and (b) weak correlations with the factors that are identified in the factor analytic model.

Reviewing univariate normality for the PIAAC data, skewness for all measures are within the range of $+/- 2.0$ and kurtosis for all measures are within $+/- 7.0$, suggesting evidence of normality.

Descriptive Statistics

	N	Skewness		Kurtosis	
	Statistic	Statistic	Std. Error	Statistic	Std. Error
Index of use of numeracy skills at home (basic and advanced—derived)	191	1.220	.176	5.663	.350
Index of use of numeracy skills at work (basic and advanced—derived)	191	.929	.176	2.516	.350
Index of use of ICT skills at home (derived)	191	1.048	.176	3.149	.350
Index of use of reading skills at home (prose and document texts—derived)	191	.837	.176	1.896	.350
Index of use of task discretion at work (derived)	191	1.431	.176	2.249	.350
Index of learning at work (derived)	191	.449	.176	-.450	.350
Index of use of planning skills at work (derived)	191	.479	.176	-1.121	.350
Index of readiness to learn (derived)	191	.878	.176	-.110	.350
Index of use of ICT skills at work (derived)	191	.457	.176	.630	.350
Index of use of influencing skills at work (derived)	191	1.018	.176	1.894	.350
Index of use of reading skills at work (prose and document texts—derived)	191	.924	.176	2.960	.350
Index of use of writing skills at work (derived)	191	1.116	.176	3.656	.350
Index of use of writing skills at home (derived)	191	.642	.176	4.601	.350
Valid N (listwise)	191				

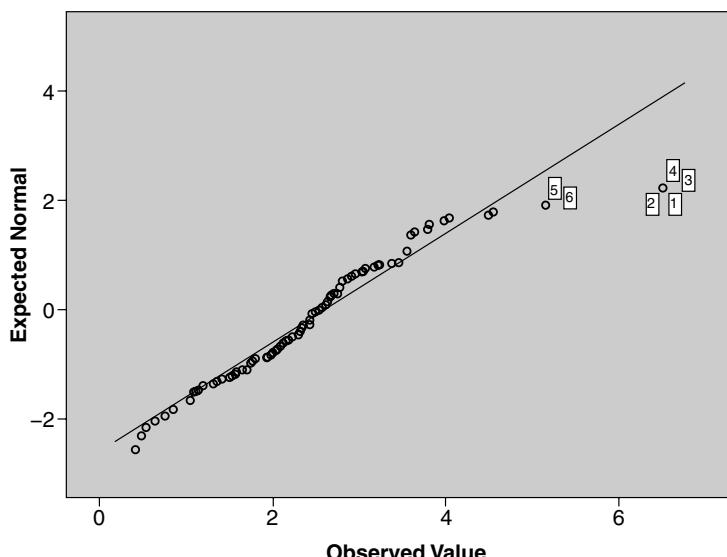
Shapiro-Wilk's formal test of normality was statistically significant for all variables suggesting evidence of nonnormality.

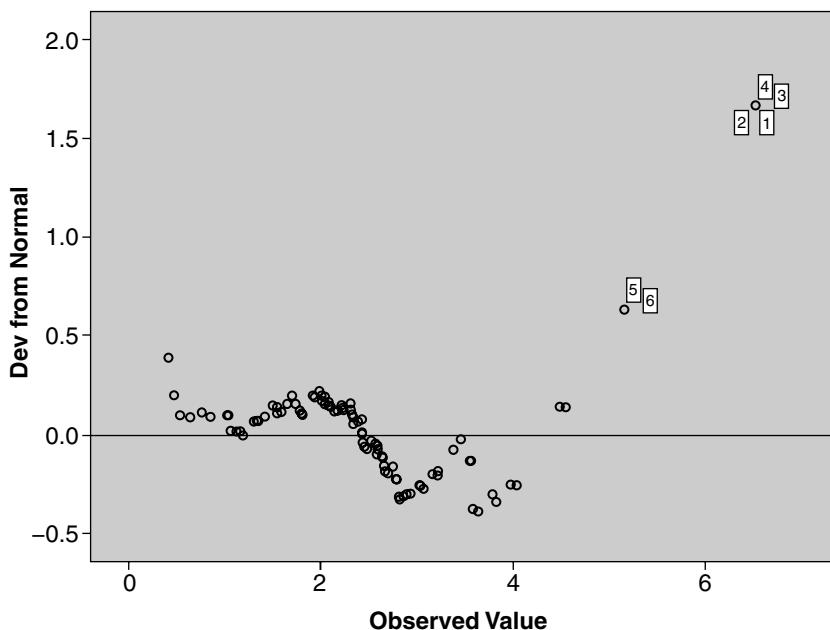
Tests of Normality

	Kolmogorov-Smirnov ^a			Shapiro-Wilk		
	Statistic	df	Sig.	Statistic	df	Sig.
Index of use of numeracy skills at home (basic and advanced—derived)	.098	191	.000	.920	191	.000
Index of use of numeracy skills at work (basic and advanced—derived)	.075	191	.011	.950	191	.000
Index of use of ICT skills at home (derived)	.068	191	.030	.951	191	.000
Index of use of reading skills at home (prose and document texts—derived)	.062	191	.075	.964	191	.000
Index of use of task discretion at work (derived)	.161	191	.000	.878	191	.000
Index of learning at work (derived)	.118	191	.000	.938	191	.000
Index of use of planning skills at work (derived)	.183	191	.000	.874	191	.000
Index of readiness to learn (derived)	.142	191	.000	.895	191	.000
Index of use of ICT skills at work (derived)	.077	191	.007	.983	191	.021
Index of use of influencing skills at work (derived)	.103	191	.000	.927	191	.000
Index of use of reading skills at work (prose and document texts—derived)	.107	191	.000	.951	191	.000
Index of use of writing skills at work (derived)	.124	191	.000	.915	191	.000
Index of use of writing skills at home (derived)	.106	191	.000	.937	191	.000

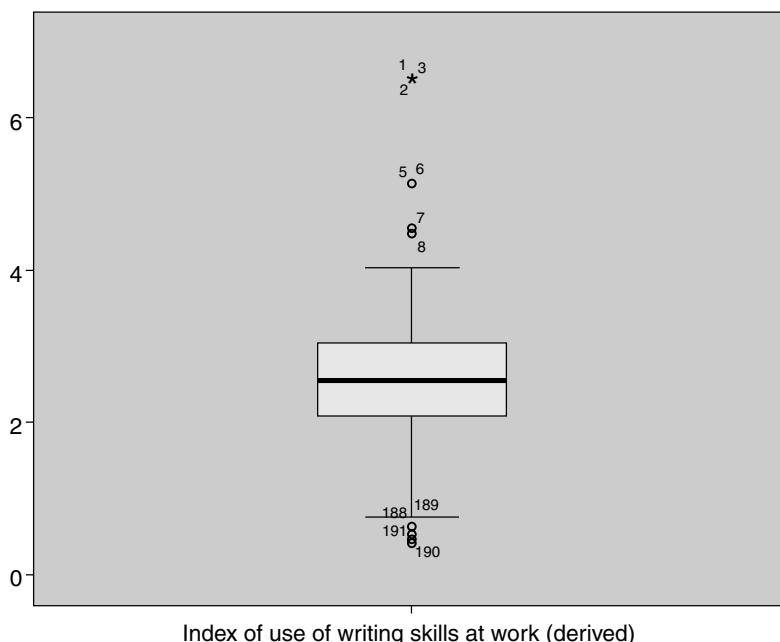
a. Lilliefors Significance Correction

The normal and detrended Q-Q plots suggest at least one potential outlying case for all 13 variables. For example, the ‘index of use of writing skills at work’ suggests that cases 1–6 may be outliers.

Normal Q-Q Plot of Index of Use of Writing Skills at Work (Derived)

Detrended Normal Q-Q Plot of Index of Use of Writing Skills at Work (Derived)


Reviewing boxplots, there are quite a few variables that have outliers suggested by the graph. For many (but not all) of the variables, these are at least some of the same cases that showed up as potential outliers in the Q-Q plots. The boxplot for the ‘index of use of writing skills at work’ suggests additional outliers that were not as evident in reviewing the Q-Q plot—not only cases 1–6 but also cases 7–8 and 188–191.



Now that we've screened for univariate outliers and have identified cases that are suggestive of outliers, we need to screen for multivariate outliers. To do so, we create a new binary variable with '1' denoting that it showed up as an outlier and '0' denoting nonoutlying cases. This has been saved in the PIAAC.EFA.sav data file and is labeled 'OUTLIER.' This binary variable will be our dependent variable in a multiple regression model with all 13 of the index variables as the independent variables. (By this point in your statistics career, it is assumed that you are familiar with creating a new variable, thus the process for doing so is not presented. Should you need a refresher on generating multiple regression, please review the earlier chapter in this text.) When generating the multiple regression model, we are not interested in the results of the analysis. Rather, we run it simply to save the Mahalanobis distance values (saved as MAH_1 in the data file). Multivariate outliers are evidenced by statistically significant Mahalanobis distance values, evaluated using a chi-square distribution with degrees of freedom equal to the number of variables. With alpha of .001, our chi-square critical value is 34.53, and our Mahalanobis distance values range from 1.84 to 56.11. Fortunately, there are only five cases with statistically significant Mahalanobis distance values.

Mahalanobis Distance

		Frequency	Percent	Valid Percent	Cumulative Percent
Valid	37.54589	1	20.0	20.0	20.0
	39.71767	1	20.0	20.0	40.0
	42.84545	1	20.0	20.0	60.0
	43.86517	1	20.0	20.0	80.0
	56.11096	1	20.0	20.0	100.0
	Total	5	100.0	100.0	

We will retain these cases for the illustration given that factor analysis is relatively robust to violations with the exception of tests of inference. (Had we filtered them out, we would see that we still end up with a two-factor solution, however only seven variables remain in the model due to the communalities greater than 1 error. For practice, you may want to try this yourself!) In this illustration, we have used maximum likelihood so we are concerned with multivariate normality. As we present our results, we will caution readers to this limitation of our data.

In terms of outlying variables, our final factor model did not suggest this was problematic (i.e., both factors had multiple items).

9.4.4 Extreme Multicollinearity and Singularity

For EFA, the simplest method to detect extreme multicollinearity and singularity is to conduct a series of multiple regression models, one regression model for each variable where that variable is the dependent variable and all remaining variables are the independent variables. If any of the resultant R_k^2 values are close to one (greater than

.9 is a good guideline to go by), then there may be an extreme collinearity problem. However, large R^2 values may also be due to small sample sizes; thus, be cautious in interpretation in cases where the number of cases is small. If the number of variables is greater than or equal to n , then perfect collinearity is a possibility. The results are not presented here for brevity; however, the largest multiple R squared values were under .50, suggesting no problems with extreme multicollinearity.

To prevent singularity, none of the variables that are being used is a composite variable for which the component variables are also included in the EFA model.

9.5 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example paragraph for the results of the exploratory factor analysis. Recall that our graduate research assistants, Addie and Oso, were assisting Dr. Wesley, a faculty member in higher education. Specifically, Dr. Wesley was interested in better understanding the underlying constructs of measures of perceived use of skills. The research question presented to Dr. Wesley from Addie and Oso included the following: What is the underlying factor structure for perceived use of skills at home and work?

Addie and Oso then assisted Dr. Wesley in conducting exploratory factor analysis, and a template for writing the research question for exploratory factor analysis is presented below.

What is the underlying factor structure for [variable set]?

It may be helpful to preface the results of the exploratory factor analysis with information on an examination of the extent to which the data were thoroughly screened.

Prior to conducting the exploratory factor analysis, the data were screened to determine the extent to which the assumptions associated with exploratory factor analysis were met. These assumptions included (a) independence, (b) linearity, (c) absence of outliers (both univariate and multivariate), and (d) lack of extreme multicollinearity and singularity. Because the data were not randomly sampled, there is a possibility that the assumption of independence has not been met. Scatterplots of each combination of variables were generated and generally suggested that the assumption of linearity was feasible, as there was no evidence of curvilinear or other nonlinear relationships. Normal and detrended Q-Q plots and boxplots suggest the presence of a few univariate outliers. These outlying points were examined as potential multivariate outliers. Mahalanobis distance values were computed using the outlying points coded as binary dependent variables and all other variables as independent variables in a multiple regression model. Five of the

cases had statistically significant Mahalanobis distance values. These items were retained, as EFA is relatively robust to violations of normality with the exception of tests of inference. However, given that maximum likelihood was the estimation method, multivariate normality was a concern. Given there is some evidence to suggest multivariate nonnormality, the model was rerun excluding the potential multivariate outliers. A two-factor solution with seven of the eight variables was achieved. Because the variable was theoretically important, it was retained in the model and the solution reflects all eight variables. Extreme multicollinearity was screened for by conducting a series of multiple regression models, one regression model for each variable where that variable is the dependent variable and all other variables are the independent variables. There were no multiple *R* squared values that were close to one; all were under .50, suggesting no problems with multicollinearity. To prevent singularity, none of the variables used are composite variables for which the component variables are also included.

Here is an example write-up of how the results for exploratory factor analysis can be presented (remember that this will be prefaced by the previous paragraph reporting the extent to which the data were thoroughly screened).

Evidence for construct validity of indices of home and work skills from the PIAAC was obtained using exploratory factor analysis.

Criteria that is often used to determine factorability of variables was applied in this analysis. These initial factorability criteria included examination of the following: (1) bivariate correlations, (2) Kaiser-Meyer-Olkin measure of sampling adequacy (overall and individual), (3) Bartlett's test of sphericity, and (4) communalities. Based on communalities above 1.0, there were five variables that were removed during this initial stage of determining factorability. The removal of these variables was done through an iterative process of removing the index with the highest communality, rerunning the EFA, and then examining the communalities. This process was repeated for each of the indices removed. The analysis presented is based on the remaining eight items.

Three of the eight items correlated at least .30 with at least one other item and an additional variable was nearly .30 (see Table 1). The overall Kaiser-Meyer-Olkin measure of sampling adequacy was .695, larger than the recommended value of .50. In addition, the measure of sampling adequacy values for the individual items were all larger than the recommended value of .50. Bartlett's test of sphericity was statistically significant [$\chi^2(28) = 193.696, p < .001$]. An additional criterion commonly used to determine factorability is that communalities should be above the recommended value of .30, providing evidence of shared variance among the items. In reviewing extracted communalities of the eight items, one-half of the variables (4 of the 8 variables) were below .30 (see Table 2). However, given the other criteria for determining factorability were met, it was determined that it was reasonable to proceed with determining the factor structure of the eight items.

Maximum likelihood estimation with promax rotation was used to extract the factors from the data. Parallel analysis was used to determine the number

■ TABLE 1Correlation Matrix for Cognitive and Work Ability Indices ($N = 191$)

Item	1	2	3	4	5	6	7
1. Index of use of numeracy skills at home (basic and advanced)	—						
2. Index of use of numeracy skills at work (basic and advanced)	.292	—					
3. Index of use of ICT skills at home (derived)	.554	.176	—				
4. Index of use of reading skills at home (prose and document texts)	.418	.166	.374	—			
5. Index of use of task discretion at work	.015	.086	.014	.049	—		
6. Index of learning at work	.001	.098	.045	.107	.037	—	
7. Index of use of planning skills at work	-.202	.013	-.111	-.048	.140	.139	—
8. Index of readiness to learn	.246	.161	.341	.296	.206	.167	.026

■ TABLE 2Factor Loadings and Communalities Based on Maximum Likelihood Analysis for Cognitive and Work Ability Indices ($N = 191$)

Item	Indices of Cognitive Skills	Indices of Work Abilities	Communality
1. Index of use of numeracy skills at home (basic and advanced)	.833	-.076	.741
2. Index of use of numeracy skills at work (basic and advanced)	.333	.124	.116
3. Index of use of ICT skills at home (derived)	.676	.145	.458
4. Index of use of reading skills at home (prose and document texts)	.535	.214	.303
5. Index of use of task discretion at work	.074	.315	.100
6. Index of learning at work	.077	.305	.094
7. Index of use of planning skills at work	-.168	.273	.121
8. Index of readiness to learn	.436	.549	.424

of factors to retain. Both 100 parallel datasets using artificial normally distributed raw data and 1,000 parallel datasets using permuted data suggested a two-factor model was appropriate (i.e., the first two raw data eigenvalues were greater than the random and permuted mean and 95th percentile eigenvalues; all other raw data eigenvalues were less in value). Although a more subjective tool for determining the number of factors, the scree plot indicated the eigenvalues leveled off after two factors, again supporting a two-factor solution. Interpretation of a two-factor solution was also plausible and was a consideration in retaining two factors. The two-factor solution represented about 30% of the variance explained when extracted. The correlation between the two extracted factors was .165.

All items contributed to a simple factor structure and had a primary factor loading above the recommended .30 with one exception—index of use of planning skills at work—which had a primary factor loading in the structure matrix of .273. One variable (index of readiness to learn) had similar factor loadings for each factor but loaded slightly stronger on factor two. All other variables had a strong primary loading with only one of the two factors in the factor structure. However, for interpretative purposes, this item was grouped with factor two. Table 2 provides the factor loading pattern matrix for the final solution. The names for the two factors are (1) Indices of Cognitive Skills and (2) Indices of Work Abilities. The results of the factor analysis lend support to internal structure validity evidence supporting the conclusion that the scores from this instrument are a valid assessment of skills and abilities, specifically Indices of Cognitive Skills and Indices of Work Abilities. Composite scores were created for the two factors by computing the mean sum of the items that loaded most strongly on each of the factors.

PROBLEMS

Conceptual Problems

1. If your research goal is to attach meaning to the identified factors, which form of factor analysis is needed?
 - a. Common factor analysis
 - b. Principal component analysis
2. What is the recommended sample size for EFA?
 - a. At least 100
 - b. At least 300
 - c. At least 500
 - d. Current research does not recommend adhering to an absolute number of cases threshold
3. Which one of the following commonly held recommendations has been shown by simulation research to often overestimate the number of factors?
 - a. Bartlett's test
 - b. Kaiser's rule
 - c. Measure of sampling adequacy
 - d. Scree plot
4. A researcher calculates KMO measure of sampling adequacy and finds a value of .60. Does this provide one form of acceptable evidence to continue the factor analysis?
 - a. Yes
 - b. No
5. Which one of the following is not used as an index to determine the initial factorability of items?
 - a. Correlations among observed items
 - b. Communalities
 - c. Measure of sampling adequacy
 - d. Scree plot

6. A researcher assumes the items they are factoring are related. Which one of the following rotation methods should be applied?
 - a. Oblique
 - b. Orthogonal
7. A researcher generates factor analysis and finds that the various indices all suggest different numbers of factors. How should the researcher determine the number of factors?
 - a. Select the fewest number of factors suggested by the results.
 - b. Select the number of factors based on where the elbow bends in the scree plot.
 - c. Apply Kaiser's rule, selecting the number of factors with eigenvalues greater than one.
 - d. Use theory to interpret results from all indices, selecting the number of factors supported statistically and defensible by theory.
8. Which one of the following is not an assumption of factor analysis?
 - a. Absence of outliers
 - b. Homogeneity of variances
 - c. Linearity
 - d. Noncollinearity
9. The measurement scale for conventional factor analysis should be at least which one of the following?
 - a. Nominal
 - b. Ordinal
 - c. Interval
 - d. Ratio
10. What factor loading is recommended for retaining a variable in a factor?
 - a. .10
 - b. .30
 - c. .60
 - d. .80

Computational Problems

1. Using the CH9_HW1_PRESCHOOL.sav dataset, conduct exploratory factor analysis following the steps in this chapter, using maximum likelihood estimation and promax rotation. Determine initial factorability using overall MSA, Bartlett's test of sphericity, and communalities. Review the pattern and structure matrix for the initial solution, and determine the variables that appear to cluster together based on the pattern matrix.
2. Using the CH9_HW2_PIAAC_NORWAY.sav dataset, conduct exploratory factor analysis following the steps in this chapter, using maximum likelihood estimation and promax rotation. Determine initial factorability using overall MSA, Bartlett's test of sphericity, and communalities. Review the pattern and structure matrix for the initial solution, and determine the variables that appear to cluster together based on the pattern matrix. (*Note: This data has been delimited to*

individuals who indicated their highest level of school was ‘above high school’ [B_Q01a_T = 3] and who were employed the year prior to completing the survey [B_Q15a = 1].)

Interpretive Problem

1. Use SPSS to conduct exploratory factor analysis with the continuous PIAAC index variables from Italy (CH9_HW_INTERPRETATIVE_ITALY.sav). The data file has been delimited to include only individuals who reported having ‘above high school’ education [B_Q01a_T = 3] and who had complete data on the index variables. Write up the results. Just for fun, compare the results using maximum likelihood estimation as compared to other estimation results. For even further fun, conduct CATPCA using the categorized index variables.

REFERENCES

- Borsboom, D. (2006). The attack of the psychometrician. *Psychometrika*, 71(3), 425–440.
- Brown, T. A. (2006). *Confirmatory factor analysis for applied research*. New York: Guilford Press.
- Bryant, F. B., & Yarnold, P. R. (1995). Principal-components analysis and exploratory confirmatory factor analysis. In L. G. Grimm & P. R. Yarnold (Eds.), *Reading and understanding multivariate statistics* (pp. 99–136). Washington, DC: American Psychological Association.
- Cattell, R. B. (1978). *The scientific use of factor analysis*. New York, NY: Plenum.
- Child, D. (2006). *The essentials of factor analysis* (3rd ed.). New York, NY: Continuum International Publishing.
- Comrey, A. L., & Lee, H. B. (1992). *A first course in factor analysis*. Hillsdale, NJ: Erlbaum.
- DeCarlo, L. T. (1997). On the meaning and use of kurtosis. *Psychological Methods*, 2, 292–307.
- Fabrigar, L. R., & Wegener, D. T. (2012). *Exploratory factor analysis: Understanding statistics*. New York, NY: Oxford University Press.
- Gaskin, C. J., & Happell, B. (2014). On exploratory factor analysis: A review of recent evidence, an assessment of current practice, and recommendations for future use. *International Journal of Nursing Studies*, 51, 511–521.
- Glorfeld, L. W. (1995). An improvement on Horn’s parallel analysis methodology for selecting the correct number of factors to retain. *Educational and Psychological Measurement*, 55, 377–393.
- Gorsuch, R. L. (1983). *Factor analysis* (2nd ed.). Hillsdale, NJ: Erlbaum.
- Guadagnoli, E., & Velicer, W. (1988). Relation of sample size to the stability of component patterns. *Psychological Bulletin*, 103(2), 265–275.
- Guilford, J. P. (1954). *Psychometric methods* (Vol. 2nd ed.). New York, NY: McGraw-Hill.
- Hahs-Vaughn, D. L. (2005). A primer for using and understanding weights with national datasets. *Journal of Experimental Education*, 73(3), 221–248.
- Hahs-Vaughn, D. L. (2006a). Analysis of data from complex samples. *International Journal of Research & Method in Education*, 29(2), 163–181.
- Hahs-Vaughn, D. L. (2006b). Weighting omissions and best practices when using large-scale data in educational research. *Association for Institutional Research Professional File*, 101, 1–9.

- Hahs-Vaughn, D. L., McWayne, C. M., Bulotskey-Shearer, R. J., Wen, X., & Faria, A. (2011a). Complex sample data recommendations and troubleshooting. *Evaluation Review*, 35(3), 304–313. doi: 10.1177/0193841X11412070
- Hahs-Vaughn, D. L., McWayne, C. M., Bulotskey-Shearer, R. J., Wen, X., & Faria, A. (2011b). Methodological considerations in using complex survey data: An applied example with the head start family and child experiences survey. *Evaluation Review*, 35(3), 269–303.
- Hendrickson, A. E., & White, P. O. (1964). Promax: A quick method for rotation to oblique simple structure. *British Journal of Statistical Psychology*, 17, 65–70.
- Henson, R. K., & Roberts, J. K. (2006). Use of exploratory factor analysis in published research: Common errors and some comment on improved practice. *Educational and Psychological Measurement*, 66, 393–416.
- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 30, 179–185.
- Hutcheson, G., & Sofroniou, N. (1999). *The multivariate social scientist: Introductory statistics using generalized linear models*. Thousand Oaks, CA: Sage.
- Kaiser, H. K. (1970). A second generation little jiffy. *Psychometrika*, 35(4), 401–415.
- Kaiser, H. K., & Rice, J. (1974). Little jiffy, mark IV. *Educational and Psychological Measurement*, 34, 111–117.
- Kish, L., & Frankel, M. R. (1973, October 17). *Inference from complex samples*. Paper presented at the annual meeting of the Royal Statistical Society.
- Kish, L., & Frankel, M. R. (1974). Inference from complex samples. *Journal of the Royal Statistical Society, Series B*, 36, 1–37.
- Kline, P. (1979). *Psychometrics and psychology*. London: Academic Press.
- Korn, E. L., & Graubard, B. I. (1995). Examples of differing weighted and unweighted estimates from a sample survey. *American Statistician*, 49, 291–305.
- Lawley, D. N., & Maxwell, A. E. (1971). *Factor analysis as a statistical method* (2nd ed.). London: Butterworths.
- Lee, E. S., Forthofer, R. N., & Lorimor, R. J. (1989). *Analyzing complex survey data*. Newbury Park, CA: Sage.
- Looney, S. W. (1995). How to use tests for univariate normality to assess multivariate normality. *American Statistician*, 49, 64–70.
- Lumley, T. (2004). Analysis of complex survey samples. *Journal of Statistical Software*, 9(8), 1–19.
- MacCallum, R. C., Widaman, K. F., Zhang, S., & Hong, S. (1999). Sample size in factor analysis. *Psychological Methods*, 4, 84–99.
- Mardia, K. V. (1970). Measures of multivariate skewness and kurtosis with applications. *Biometrika*, 57, 519–530.
- National Science Foundation. (2010). Survey of doctorate recipients 2010. Retrieved from <http://www.nsf.gov/statistics/srvydoctoratework/>
- Nunally, J. C. (1978). *Psychometric theory* (2nd ed.). New York, NJ: McGraw Hill.
- O'Connor, B. P. (2000). SPSS and SAS programs for determining the number of components using parallel analysis and Velicer's MAP test. *Behavior Research Methods, Instruments & Computers*, 32(3), 396–402.
- Penny, K. I. (1996). Appropriate critical values when testing for a single multivariate outlier by using the Mahalanobis distance. *Applied Statistics*, 45, 73–81.

- Pfeffermann, D. (1993). The role of sampling weights when modeling survey data. *International Statistical Review*, 61(2), 317–337.
- Preacher, K. J., & MacCallum, R. C. (2002). Exploratory factor analysis in behavior genetics research: Factor recovery with small sample sizes. *Behavior Genetics*, 32(2), 153–161.
- Skinner, C. J., Holt, D., & Smith, T. M. F. (Eds.). (1989). *Analysis of complex samples*. New York: Wiley.
- Small, N. J. H. (1980). Marginal skewness and kurtosis in testing multivariate normality. *Applied Statistics*, 29, 85–87.
- Steiner, D. L. (1998). Factors affecting reliability of interpretations of scree plots. *Psychological Reports*, 83, 687–694.
- Suhr, D. D. (2006). *Exploratory or confirmatory factor analysis?* Paper presented at the SAS User's Group International 31 (SUGI), San Francisco, CA.
- Thompson, B., & Daniel, L. G. (1996). Factor analytic evidence for the construct validity of scores: A historical overview and some guidelines. *Educational and Psychological Measurement*, 56(2), 197–208.
- Zwick, W. R., & Velicer, W. F. (1982). Factors influencing four rules for determining the number of components to retain. *Multivariate Behavioral Research*, 17, 253–269.
- Zwick, W. R., & Velicer, W. F. (1986). Comparison of five rules for determining the number of components to retain. *Psychological Bulletin*, 99, 432–442.

Chapter 10

PATH ANALYSIS, CONFIRMATORY FACTOR ANALYSIS, AND STRUCTURAL EQUATION MODELING

CHAPTER OUTLINE

10.1 What Path Analysis and Confirmatory Factor Analysis Are and How They Work	442
10.1.1 Characteristics	444
10.1.2 CFA Sample Size	460
10.1.3 CFA Effect Size	462
10.1.4 CFA Assumptions	462
10.2 Mathematical Introduction Snapshot	465
10.3 Computing Path Analysis and Confirmatory Factor Analysis Using LISREL	466
10.3.1 CFA Using LISREL	466
10.3.2 CFA Using Mplus	493
10.4 Data Screening	494
10.5 Power	494
10.5.1 Post Hoc Power for CFA Using G*Power	494
10.6 Research Question Template and Example Write-Up	496

KEY CONCEPTS

1. Exogenous variable
2. Free parameter
3. Fixed parameter
4. Goodness-of-fit
5. Latent variable
6. Measurement error

With exploratory factor analysis, we deviated from concerning our analyses with procedures that have the examination of one or more *a priori* outcomes as their goal and branched into a procedure that allows exploration of data to reduce a large number of variables into identifiable clusters of variables to better understand the structure of the data. This chapter continues that tradition a bit with the treatment of confirmatory factor analysis—a tool often used as follow-up to EFA. In this chapter, we will also touch on path analysis. Path analysis and confirmatory factor analysis (CFA) both fall under the umbrella of structural equation modeling (SEM). Structural equation modeling is a collection of regression-based statistical procedures that allow the examination of relationships of observed as well as unobserved variables. For the purposes of this text, we will focus primarily on path analysis and CFA with a brief segue into SEM. Path analysis can be considered a special case of SEM with only observed variables, and confirmatory factor analysis can be considered a special case of SEM with measurement models only.

Our objectives are that, by the end of this chapter, you will be able to (a) understand the concepts underlying path analysis and confirmatory factor analysis, (b) determine and interpret the results of path analysis and CFA, and (c) understand and evaluate how to screen data prior to conducting path analysis and CFA.

10.1 WHAT PATH ANALYSIS AND CONFIRMATORY FACTOR ANALYSIS ARE AND HOW THEY WORK

Life at the stats lab is lively today; we find that the group has been tasked with an additional analysis of the data they have been working with. Addie Venture and Oso Wyse, ones never to back down from a data exploration challenge, are ready to help.

Dr. Wesley, a faculty member from the Higher Education program who has been working with the stats lab, was extremely pleased with the EFA conducted by the team and has asked that the group continue to explore the PIAAC Survey of Adult Skills data using confirmatory factor analysis (CFA). Dr. Wesley has never conducted CFA, but a number of studies recently published in the scholarly journals in the field of higher education have applied CFA to similar data, and Dr. Wesley inquires about the possibility of applying CFA to the PIAAC data [the Survey of Adult Skills, a large data collection effort from the Organization

for Economic Cooperation and Development's Programme for the International Assessment of Adult Competencies (PIAAC)]. Addie and Oso are excited to continue to work with the PIAAC data and suggest the following research question to Dr. Wesley: *How well does the two-factor model (Cognitive Skills and Work Abilities) explain the pattern of correlations among the measured indices?* With CFA, the goal of the project is now to confirm results identified from the EFA.

Path analysis and confirmatory factor analysis (CFA) both fall under the umbrella of structural equation modeling (SEM). Structural equation modeling has many names, some of which include analysis of covariance structure, covariance structure analysis, covariance structure modeling, causal modeling, and causal analysis. The latter two synonyms for SEM are really misnomers, as ‘causal’ does not mean that causality can be assumed when conducting structural equation modeling. Rather, it is the *design of the study* that may allow for causal inference, *not* the statistical procedure that is used to analyze the data. As noted by Bollen and Pearl (2013, p. 12):

[R]esearchers do not derive causal relations from an SEM. Rather the SEM represents and relies upon the causal assumptions of the researcher. These assumptions derive from the research design, prior studies, scientific knowledge, logical arguments, temporal priorities, and other evidence that the researcher can marshal in support of them. The credibility of the SEM depends on the credibility of the causal assumptions in each application.

Indeed, SEM is not a panacea for research design and/or analytic problems, although some researchers continue to believe this. We'll not delve further into this discussion; however, those interested are encouraged to read Bollen and Pearl, as they provide a very nice historical context to this myth, among others, in SEM.

Although the origins of SEM are with the development of path analysis by biologist Sewall Wright, it wasn't until the 1980s and the advent of more user-friendly software that there was broader application of SEM (Hancock & Mueller, 2006). SEM is actually not just one statistical procedure but, as indicated previously, an umbrella for a family of regression-based procedures. It is important to note that SEM is not a stagnant family but actually a *growing* family of procedures. The majority of this chapter will focus on confirmatory factor analysis with a brief glimpse into path analysis and an overview of a number of structural models. If you find this chapter stimulating, and I hope you do, then I suggest you explore more complex procedures under the SEM umbrella such as multilevel SEM, latent class analysis, latent transition analysis, latent growth modeling (the latter three fall under another umbrella of ‘latent variable modeling’ techniques, which can be subsumed under the more general SEM umbrella), and more.

As mentioned, however, this chapter will focus primarily on confirmatory factor analysis. Before we get into the depths of CFA, a cursory knowledge of path analysis is important as it is still widely used [e.g., in about 500 SEM studies published in 16 psychological research journals between 1993 to 1997, about 25% were path analysis

models (MacCallum & Austin, 2000)]. In addition, understanding the concepts of path analysis will help build the foundation on which more complex SEM can be understood.

In the previous chapter, we learned that exploratory factor analysis (EFA) is a statistical procedure that allows us to cluster together variables and understand the factor structure. Confirmatory factor analysis (CFA) is similar to EFA in that it can be used to provide evidence of construct validity by examining the underlying constructs of variables, allowing us to work with all items simultaneously and at the same time know something about their underlying data structure. The intention of factor analysis, be it either exploratory or confirmatory, is to determine the number and nature of the latent variables (also referred to as ‘factors’) that account for the variation and covariation among observed variables (observed variables are sometimes referred to as ‘indicators’ or ‘manifest variables’). A latent variable (a new term that we’ll discuss in this chapter, also referred to as a ‘latent construct’ or ‘factor’) is an attribute that cannot be directly observed which influences more than one observed measure and accounts for the correlations among the observed variables (there are instances where a latent construct influences a single indicator, but that is not recommended). CFA provides a more parsimonious way to view the relationships among the indicators, or observed variables, because the number of latent variables, or unobserved items, is less than the number of observed and measured items.

However, CFA is different (as implied by the name of the procedure) in that it is not appropriate for exploratory situations. CFA is used when there is a theory to be tested and is often conducted *prior* to estimating structural models (i.e., models that estimate relationships between latent constructs) (and in those instances, is usually referred to as the measurement model). The global research question that is answered through SEM is this: To what extent is there consistency between the estimated population covariance matrix produced by the model and the observed covariance matrix of the sample? In application, stating your research question in this form doesn’t provide the type of specificity usually desired. More often, you want your research question stated in a way that really helps the reader understand *your* specific research. Thus, as we’ll see later, while testing the empirical covariance matrix of your sample to the estimated population covariance matrix of the model is what is accomplished in SEM, the research question we’ll offer to applied researchers will be much more descriptive and will allow readers to really understand the relationships being examined in the CFA.

10.1.1 Characteristics

10.1.1.1 Path Analysis

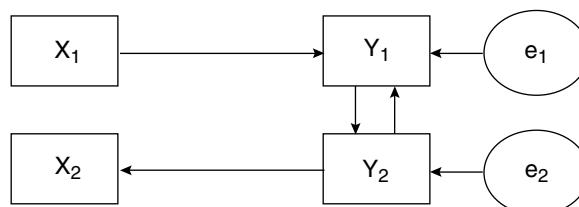
Path analysis falls under the umbrella of SEM and allows researchers to test models that include only observed variables. As such, path analysis is the most simplistic of all SEM procedures, although not all path analytic models are ‘simple,’ per se. The most basic type of path analysis model is one where an exogenous variable (X) has

a direct influence on an endogenous variable (Y) (see Figure 10.1). We'll talk more about exogenous and endogenous variables later, but endogenous variables are related to dependent variables, and exogenous variables are akin to independent variables. The path coefficient from X to Y is represented with a unidirectional arrow (with the arrow pointing from the exogenous to the endogenous variable, i.e., from the 'influenced by' to the 'influenced,' although we know from our earlier discussion that SEM alone does not infer causality) is interpreted as a regression coefficient. This unidirectional arrow indicates a *recursive* (or nonreciprocal) relationship between X and Y . The influence of X goes to Y , and Y does not also influence X . A *nonrecursive* (or reciprocal) relationship between X and Y would be indicated by one arrow from X to Y as well as another arrow from Y to X . The influence of X goes to Y , and Y also influences X . In other words, those variables are both an influencer and an influence. The error, denoted by e , is the residual term of the observed variable and represents the unexplained variation that is not modeled by the observed variable in the model. In other words, there are other possible influences on Y that are not present in the model (i.e., not explained by X in this illustration). Errors are unobserved and are thus treated as latent variables.

In path analysis and other SEM procedures, the researcher specifies relationships between variables a priori given existing theory and previous empirical research, and the effects among the variables are estimated through the path analytic model. This process is called 'model specification.' The next step is 'model identification.' In identifying the model, we ask, Can a unique set of parameters be estimated given the sample covariance matrix, S , and the theoretical model implied covariance matrix, Σ ? This entails, among other things, determining if the number of freely estimated parameters is less than or equal to the number of distinct values in the sample covariance matrix



a) *Recursive*



b) *Non-Recursive*

FIGURE 10.1

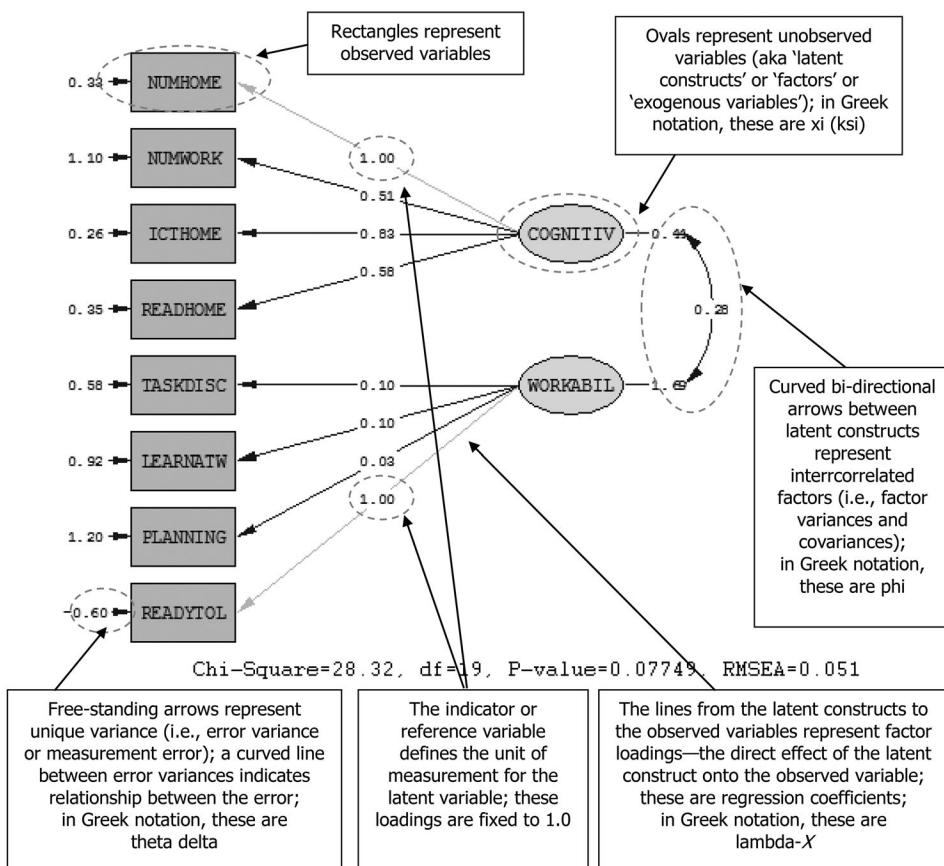
Path Analysis Example: (a) Recursive and (b) Nonrecursive

(referred to as the order condition). When the model is identified, it must then be estimated. Model estimation entails estimating the parameters of the model, and there are different methods by which models can be estimated (e.g., maximum likelihood, generalized least squares). Determining the fit of the model is the model testing process. There are a growing number of goodness-of-fit indices that can be used to test the fit of the specified model. Based on poorness-of-fit, model modification is often the next step. Modification of the model is the process of finding a better fitting model.

10.1.1.2 CFA Model Specification, Identification, Estimation, Evaluation, and Interpretation

CFA most often estimates models from a variance-covariance matrix where variances of the observed variables are on the diagonal and covariances on the off-diagonal. Unlike EFA, rotation is not required in CFA, as model specification is such that each observed variable loads on only one latent factor, creating a more parsimonious model. In other words, all cross-loadings (i.e., the relationship between observed variables and all other latent constructs) in CFA are specified to be zero, and thus CFA has fewer parameters to estimate as it attempts to reproduce the observed relationships among the observed variables. The concept of zero cross-loadings is illustrated in the path diagram in Figure 10.2, where the rectangles represent the observed variables, and the ovals represent the latent constructs. There are four observed variables ('NUMHOME' to 'READHOME') that load onto the latent construct, 'COGNITIV.' There are four observed variables ('TASK-DISC' to 'READYTOL') that load onto the latent construct, 'WORKABIL.' The first four observed variables ('NUMHOME' to 'READHOME') do not load onto 'WORKABIL,' and thus there are zero cross-loadings between these observed variables and the latent construct 'WORKABIL.' Nonzero cross-loadings would have resulted in these observed variables also having paths drawn to the latent construct 'WORKABIL.'

Whereas measurement error (i.e., unique variance, represented as arrows to the left of the observed variables in Figure 10.2) is assumed to be random in EFA, the relationship among measurement errors can be specified in CFA. Measurement error measures the imperfection of the observed variables as measures of the latent constructs. In other words, we are not assuming the manifest variables are perfect measures of the latent variables. Random measurement error assumes that the correlation between observed variables that load on the same latent construct is wholly due to the shared influence of the latent construct. Thus, when there is random measurement error, the relationship between observed variables will be zero when the latent construct is parceled out—i.e., this means that the errors are uncorrelated. CFA allows researchers to correlate this unique variance in the model when covariation due to sources other than the latent factor can be justified; for example, data collection method effects or reverse or similarly coded items, and these are referred to as 'correlated errors,' 'correlated residuals,' or 'correlated uniqueness.' In a path diagram, this would be represented by double-headed arrows between latent errors. The ability to adjust for measurement error (and not have to assume that measures are perfectly reliable) is one of the greatest advantages of CFA over other general procedures outside the SEM framework.

**FIGURE 10.2**

Initial Two-Factor Measurement Model (i.e., CFA) Example Generated Using LISREL

CFA Model Specification

Within CFA and unlike EFA, model specification is first governed by theory. Specification of a CFA model must be strongly grounded in existing theory and previous empirical research. This is because a number of decisions must be made in terms of model specification that go beyond basic selection of estimation method and other general statistical modeling decisions. These decisions include, for example, specifying which observed variables load onto which latent constructs and interrelationships between latent constructs, among others. These decisions must be informed by existing theory and previous empirical research (e.g., exploratory work such as EFA).

We learned about factor loadings and unique and common variance in our discussion of EFA. These concepts still apply in CFA. A factor loading in CFA reflects the relationship between a latent construct and an observed variable (i.e., the direct effect of the latent variable on the manifest variable) and is the slope of increase (when positive) or decrease (when negative) for each unit of increase or decrease in the latent construct.

We'll discuss this concept later in our presentation of parameter estimates. As touched on briefly already, unique variance is variance that is *not* accounted for by the latent construct and is a combination of systematic variance specific to the observed variable (e.g., reliability) and random error variance (i.e., measurement error or unreliability in the observed variable). Unique variance is assumed to be measurement error and results when the relationship between observed variables is above and beyond that shared relationship with the latent construct. The CFA counterpart to orthogonality in EFA is to fix factor covariances to zero.

Unobserved variables in SEM are latent constructs, as mentioned earlier. In the world of structural equation modeling, latent constructs can be either exogenous or endogenous. Endogenous latent variables are caused from one or more variables in the model and are related to dependent variables. In the path diagram in Figure 10.2, an endogenous latent variable would have been evidenced if there had been a path pointing from one of the other variables in the model to one of the latent constructs. Examples of CFA models that include endogenous variables are multiple indicator–multiple cause (MIMIC) models and CFA models with higher-order factors. Many times, CFA models have exogenous latent variables; latent constructs that are *not* caused by any other variable in the specified model. For purposes of this chapter, *we will explore only CFA models with exogenous latent constructs*. Exogenous latent variables are related to independent variables.

Latent constructs in CFA are generally specified to be interrelated (i.e., there is some relationship between the observed variables that load onto the latent constructs, and that relationship can be accounted for by the correlation between latent constructs given that each observed variable loads on only one latent construct) with no specification on the directionality of that relationship. Should directionality be specified, the researcher is no longer operating within just the realm of CFA.

Model complexity is the total number of parameters that can be estimated in the model. This is not to be confused with sample size. Model complexity equals the largest number of variances and unique covariances that *can be* estimated. This value can be determined by counting the number of values in the lower diagonal of the sample covariance matrix or computed as $\left(\frac{v(v+1)}{2}\right)$, where v equals the number of observed variables. The difference between the largest number of parameters to be estimated, p , and the actual number of estimated parameters, q , is the model degrees of freedom: $df_m = p - q$. Thus, an unidentified model is one where there are more parameters to be estimated than can be estimated (Kline, 2011).

CFA Model Identification

A CFA model is identified if a unique set of parameter estimates for each unknown parameter in the model (i.e., latent construct) can be derived from the known variance-covariance matrix. It is only when the measurement model is identified that the CFA

parameters can be estimated. There are two important steps that must be taken at this point. First, the number of elements in the variance-covariance matrix must be at least as great as the number of *freely* estimated parameters in the model (e.g., factor loadings, factor variances/covariances, observed variable measurement error). (A *freely* estimated parameter means that parameter is estimated in the model. A *fixed* parameter is a parameter that is not estimated—it is fixed to zero or some other value.) An *underidentified* model is one where there is more unknown than known—where the number of freely estimated parameters is greater than the number of elements in the variance-covariance matrix. An underidentified model will result in negative degrees of freedom. An *overidentified* model is one where there is more known than unknown—where the number of elements in the variance-covariance matrix is greater than the freely estimated parameters. An overidentified model will result in positive degrees of freedom. A *just-identified* model is one where the known and unknown are the same—where the number of freely estimated parameters is the same as than the number of elements in the variance-covariance matrix. In other words, in a just identified model, there is one and only one set of parameter estimates that will fit the data, and this will be a perfect fit (i.e., the input matrix will be perfectly reproduced). Just identified models have zero degrees of freedom.

The ‘rule of 3,’ which was recommended with exploratory factor analysis, applies to CFA as well: *At least three observed variables per latent construct is recommended.* This is even more important in CFA. A CFA model with only one latent construct and three observed variables will be just identified; with four or more observed variables, the model will be overidentified. Generally, a CFA model with two or more latent constructs and two observed variables per construct will be overidentified but is prone to *empirical underidentification* (parameter estimates are unattainable even though statistically the model is just or overidentified).

Second, latent constructs must be *scaled* in one of two ways: (a) fix the variance of the construct to a specified value, usually 1.0; or (b) fix the metric of the construct by scaling to an observed variable (i.e., scaling an indicator variable). The first approach standardizes the solution such that the variances are 1.0. The latter is the most common approach, and it entails selecting an observed variable through which its metric is passed to the latent construct and produces both a standardized and unstandardized solution. Regardless of scaling selected, the fit of the models will be identical.

CFA Model Estimation Procedures

Maximum likelihood (ML) is the most widely applied estimation procedure for CFA (Savalei, 2012). This may be because standard errors are produced from ML estimates, which can then be used for tests of inference of parameter estimates, and because of the goodness-of-fit indices generated through ML. ML estimation calculates parameter estimates that maximize the probability that the observed matrix would be sampled from the population. ML is an iterative process when the CFA model is overidentified, using and refining the parameter estimates to minimize the difference between

the predicted covariance matrix ($\hat{\Sigma}$) and the sample variance-covariance matrix (S). A model that converges means that any further estimates will not improve the model (i.e., the difference between the predicted covariance matrix ($\hat{\Sigma}$) and the sample variance-covariance matrix (S) cannot be improved). Estimating with ML requires meeting two conditions: (1) large sample size and (2) observed variables measured on a continuous scale. ML also requires the distribution of the observed data (i.e., those that load onto latent constructs) meet the assumption of multivariate normality. A quick note related to the use of continuous variables: There are situations where researchers may wish to standardize the observed variables included in their model. For example, standardization of observed variables is helpful when the units of measurement are difficult to interpret or when there are different units of measurement of observed variables in the model and you wish to compare coefficients. In cases where there are different units of measurement (e.g., scales measured 0–15, scales measured 100–500, and various others that are dissimilar), standardization may also help with model convergence.

If you find yourself with unfortunate receipt of ‘model failed to converge’ message in CFA, trust that you are in good company, as it happens to the best of us. There are a number of reasons why this may have occurred, some more easily solved than others. The easiest remedy is that the number of iterations simply needs to be increased. If increasing the number of iterations still does not result in model convergence, examine the scale and range of the observed variables. Observed variables with substantially different scales and/or ranges can sometimes be problematic. Standardizing the observed variables may be helpful in this situation so that they are placed on a common scale. If the model still fails to reach convergence, adequacy of the starting values for estimating the parameters may be in question. This is not something researchers would normally concern themselves with, as today’s software automatically computes starting values (however, you may find yourself in a situation where this is helpful if you do find convergence problems, particularly in more complicated models). If defining the starting values does seem like a plausible approach to reach convergence, begin by defining the error variances (Brown, 2006).

Although ML may be the most widely used estimation procedure, it is by far not the only from which to select. And, in some cases, it is not the appropriate estimation procedure to select. Recall that one of the assumptions when using ML is that the distribution of the observed variables is multivariate normal and one of the conditions is that the observed variables are continuous. When these do not occur, other estimation procedures are more appropriate. Maximum likelihood with robust standard errors and chi-square (MLM) will produce identical parameter estimates to that of ML but corrects for nonnormal data when computing the standard errors and chi-square statistic (i.e., the Satorra-Bentler scaled chi-square is produced). Unweighted least squares (ULS), weighted least squares (WLS), and robust weighted least squares (WLSMV) are other estimation methods that can be used when the observed variables are categorical and/or when the multivariate normality assumption does not hold. MLM, rather than WLS, is considered preferable based on research that suggests it outperforms with samples that are small to moderate in size (Curran, West, & Finch, 1996).

CFA Model Evaluation and Interpretation: Goodness-of-Fit Indices and Parameter Estimates

In addition to consideration of the substantive theory underlying the model (Schumacker & Lomax, 2010), evaluating model fit involves examination and interpretation of both overall model fit as well as parameter estimate evaluation. Goodness-of-fit indices for overall model fit are discussed first, and this is followed by an examination of parameter estimates.

There are a rather large number of goodness-of-fit indices that can be reviewed in the context of CFA. Various typologies for categorizing the indices exist, such as

- 1) (a) absolute fit (e.g., chi-square, SRMR, RMR), (b) model parsimony fit (e.g., RMSEA, GFI, AGFI), and (c) comparative/incremental fit (e.g., CLI, TLI, NFI) (Brown, 2006);
- 2) (a) model fit (chi-square, GFI, AGFI, RMR), (b) model comparison (TLI, NFI, NNFI, CFI), and (c) model parsimony (AGFI, PNFI, AIC);
- 3) (a) model test statistics (e.g., chi-square) and (b) approximate fit (e.g., RMSEA, GFI, CFI, SRMR) (Kline, 2011), among other typologies.

The typology for categorizing the fit indices is not important, but for sake of parsimony, we will adhere to that of Kline. What *is* important is understanding which indices are recommended to report and how to interpret them, and this is what we will concentrate on.

Model Test Statistics

Model test statistics are tests of inference related to the fit of the SEM model (in our case, CFA specifically) and are the “original fit statistics in SEM” (Kline, 2011, p. 193). The model test statistic most always reported is the model chi-square test. In terms of the theoretical model, chi-square is the only statistical test available (Schumacker & Lomax, 2010). When estimating with maximum likelihood, this is also called the likelihood ratio chi-square. The chi-square test tests the null hypothesis that the predicted covariance matrix ($\hat{\Sigma}$) and the sample variance-covariance matrix (S) are equal. When there is perfect fit, the chi-square value will be zero. In other words, there is no difference between the values in the sample covariance matrix and the model-implied covariance matrix. When statistically significant, this indicates that the predicted covariance matrix ($\hat{\Sigma}$) and the sample variance-covariance matrix (S) are not equal, suggesting poor model fit. In other words, rejecting the model test statistic is not desirable, as that means that your model does not fit the hypothesized model. While often reported, the chi-square test is a relatively weak goodness-of-fit measure [and thus sometimes referred to as a “badness-of-fit” statistic (Kline, 2011, p. 193)] as it is greatly inflated by sample size (thus statistically significant chi-square values are often an artifact of sample size alone rather than there being notable differences between S and $\hat{\Sigma}$), does not adhere to the chi-square distribution with small samples and nonnormal data (and thus results are not reliable in these situations), and has the same fault as

other hypothesis tests (i.e., absolute $\mathbf{S} = \hat{\Sigma}$ with no allowance for ‘nearly’ or ‘almost’ good fit). Despite the criticisms, the chi-square test is useful to determine best fit when comparing nested models and for calculating other goodness-of-fit indices.

Approximate Fit Indices

In comparison to model test statistics, approximate fit indices are more subjective in that they range on a continuum (e.g., to +1.0) and the higher the value (or lower, in some cases), the better the fit of the model. This should alert you to the fact that interpreting goodness-of-fit in SEM has a fair degree of subjectivity. Not only that, but some fit indices that are commonly reported in output are discouraged to use in practice. Thus, it is in your best interest to understand the range of fit indices so that educated choices in use and interpretation can be made.

‘Approximate fit indices’ is actually a fairly large umbrella under which many types of goodness-of-fit indices exist. These include absolute fit indices which test the hypothesis that the predicted variance-covariance matrix ($\hat{\Sigma}$) is equal to the sample variance-covariance matrix (\mathbf{S}) (e.g., SRMR, RMR, RMSEA), incremental fit indices which evaluate model fit to a baseline model (IFI, TLI, NFI, NNFI, CFI), and parsimony-adjusted indices which penalize model complexity (e.g., PGFI, PNFI). When comparing models (i.e., comparative measures of fit), the AIC, BIC, SABIC, and ECVI can be used.

Absolute fit indices include the SRMR, RMR, and RMSEA. The *standardized root mean square residual* (SRMR) is derived from the residual correlation matrix and reflects the mean difference between the observed and model predicted correlations. The SRMR can range from 0 (perfect fit) to +1.0 and will be lower as the number of estimated parameters increase and as the sample size increases. Different thresholds of acceptable model fit exist. Values of $< .08$ have been suggested as good model fit (Hu & Bentler, 1999), as have values $< .05$ (Schumacker & Lomax, 2010).

The *root mean square residual* (RMR) is related to the SRMR and reflects the square root of the mean squared difference between the observed and model predicted covariances. The range of the RMR is based on the scale of the observed variables, and thus interpretation of the RMR can be difficult. As correlations, rather than covariances, are usually preferred given that their interpretations are not dependent on the metric of the observed variable, the SRMR is usually preferred above the RMR. The RMR can be used, however, as a comparative fit index comparing model fit of two different models (Schumacker & Lomax, 2010).

The *root mean square error of approximation* (RMSEA) includes a “penalty function” for poor parsimony in the model, and thus gives an indication of reasonable fit rather than absolute fit (Brown, 2006, p. 83). It is for this reason that the RMSEA is perhaps the most widely used goodness-of-fit index for researchers who rely on fit indices as evidence of model fit (Marsh, Hau, & Grayson, 2005). The RMSEA relies on the noncentral chi-square distribution, allowing the degree of model misspecification to be estimated through the noncentrality parameter (NCP). An NCP of zero indicates perfect model fit, and thus the chi-square distribution holds. An NCP greater than zero

indicates less than perfect fit, relying then on the noncentral chi-square distribution. The RMSEA is generally insensitive to sample size (however the confidence intervals of the RMSEA *are* sensitive to sample size) but is sensitive to the number of parameters being estimated (given the degrees of freedom in its estimation; RMSEA confidence intervals will also be wider as the number of freely estimated parameters increases). Research also suggests that RMSEA is substantially impacted by the size of the factor loading, with the impact of model specification decreasing with larger factor loadings. More specifically, loadings of .40 or less were generally insensitive to model misspecification, and thus traditional cutoff values of .05 are too liberal when the factor loadings are weak. The RMSEA can range from 0 (perfect fit) to positive infinity (however values greater than +1.0 are uncommon). Values of .05 to .08 or smaller indicate close fit (Schumacker & Lomax, 2010). Reporting the confidence intervals as estimates of precision along with the RMSEA value is recommended (MacCallum, Browne, & Sugawara, 1996). Ideally, the lower bound of the RMSEA CI will equal zero, and the upper interval should be $< .10$ (Kline, 2011).

The test of close fit (CFit) uses the RMSEA to derive a measure of reasonable fit and is a one-tailed test, the null hypothesis being that RMSEA is .05 and the alternative being RMSEA is $> .05$. RMSEA values that are statistically significant (i.e., $p < .05$) are considered unacceptably fitting models. RMSEA in these cases is greater than .05 and the model fit is worse than close fitting. RMSEA values that are *not* statistically significant (i.e., $p > .05$) are considered acceptably fitting models. RMSEA in these cases is less than .05 and the model fit is considered close fitting. Smaller sample sizes and lower model df may be prone to decreased power with the test of close fit.

Parsimony-Adjusted Indices

Parsimony-adjusted indices penalize model complexity, and thus models that are more complex will result in poorer fit based on these indices (e.g., PGFI, PNFI).

The *parsimony goodness-of-fit* (PGFI) is based on the GFI and adjusts for degrees of freedom by penalizing models that are more complex. There are no thresholds recommended for interpreting the PGFI.

The *parsimonious normed fit index* (PNFI) is based on the NFI and adjusts for degrees of freedom by penalizing models that are more complex. There are no thresholds recommended for interpreting the PNFI.

The *goodness-of-fit* (GFI) and *adjusted goodness-of-fit* (AGFI) are other indices of parsimonious fit. For both, the closer the value is to +1.0, the better the fit. Current research discourages the use of both the GFI and AGFI given their influence on sample size (Sharma, Mukherjee, Kumar, & Dillon, 2005).

Incremental fit indices evaluate model fit to a baseline model (CFI, TLI, NNFI, NFI, IFI). Indices in this group, including the comparative fit index (CFI) (Bentler, 1990) and the Tucker-Lewis index [TLI; also known as the non-normed fit index, NNFI (Tucker &

Lewis, 1973)], compare a baseline model that generally has fixed covariances of zero for the observed variables to a more complex model identified by the researcher.

The **CFI** is based on the noncentral chi-square distribution and ranges from 0 (poor fit) to +1.0 (best fit). Values greater than .95 suggest good model fit (Hu & Bentler, 1999).

The **TLI/NNFI**, like the RMSEA, includes a penalty function for estimating parameters that do not contribute to improving the fit of the model. More specifically, the size of the TLI is dependent on the average correlation coefficient of the data. There is not a range within which values for the TLI/NNFI can fall; however, values closer to +1.0 are interpreted as better fitting models.

The Bentler-Bonett index or **normed fit index (NFI)** was the first goodness-of-fit measure proposed. NFI values above .95 suggest good model fit, between .90 and .95 marginal, and below .90 poor. The disadvantage of this index is that as the number of parameters estimated increases, the size of the NFI also increases. Thus, researchers are discouraged from interpreting the NFI for their models.

The **incremental fit index (IFI)** ranges from 0 to +1.0, indicating the worst to the best model fit, respectively.

Comparative Fit

When cross-validating or comparing models (i.e., comparative measures of fit), the Akaike information criterion (AIC), Bayesian information criterion (BIC), sample-size adjusted Bayesian information criterion (SABIC), and expected cross-validation index (ECVI) can be used. These measures allow researchers to gauge cross-validation of models. In other words, these indices can be used to evaluate model performance of competing models. These measures of fit are meaningful only when multiple models are generated and determining the better fit of the models is needed. The BIC and SABIC increase as the sample size increases; however, the SABIC does not increase at the same rate as the BIC. The AIC and BIC increase as the number of parameters estimated increases; however, the AIC penalizes at a lower rate than the BIC. AIC, BIC, and SABIC values closer to zero indicate better model fit. Thus, when comparing nested models, the model with the lowest AIC or BIC is the better fitting model. A nice feature of the AIC is that it can be used to compare nonnested models (Kaplan, 2000), as can BIC and SABIC. The ECVI is also a comparative fit index and is interpreted similarly. When there are competing models, the model with the lowest ECVI has the greatest replicability potential (Byrne, 1998).

Guidelines for Interpreting Goodness-of-Fit Indices

There are many recommendations for interpreting goodness-of-fit indices, and multiple fit indices should be reported. Reporting chi-square, RMSEA, and SRMR for all types of models and the CFI when comparing models has been recommended

(Schumacker & Lomax, 2010). Some of the most cited comes from research conducted by Hu and Bentler (1999), who suggest reporting the SRMR and a comparative fit index. A summary of indices for evaluating model fit has been provided in Box 10.1.

BOX 10.1 EVALUATING MODEL FIT USING GOODNESS-OF-FIT INDICES

Category	Indices	Interpretation	Notables
Model Test Statistics	Model chi-square	Statistical significance indicates that the predicted covariance matrix and the sample variance-covariance matrix are not equal, suggesting poor model fit	<ul style="list-style-type: none"> Greatly inflated by sample size Results are not reliable with small samples and nonnormal data No allowance for ‘nearly’ or ‘almost’ good fit
Approximate Fit Indices	SRMR	<ul style="list-style-type: none"> Ranges from 0 (perfect fit) to +1.0; lower as the number of estimated parameters increase and as the sample size increases Values of $< .08$ suggest good model fit 	
	RMR	Difficult to interpret as range is based on the scale of the observed variables	SRMR preferred over RMR as SRMR is standardized
	RMSEA (and CI)	<ul style="list-style-type: none"> Ranges from 0 (perfect fit) to positive infinity Lower interval of the RMSEA CI will ideally equal zero, and the upper interval should be $< .10$ 	<ul style="list-style-type: none"> Values greater than +1.0 are uncommon Generally insensitive to sample size (however, the confidence intervals <i>are</i> sensitive to sample size) Sensitive to the number of parameters being estimated Substantially impacted by the size of the factor loading, with the impact of model specification decreasing with larger factor loadings
	CFit	Nonstatistical significance is considered acceptable fit	Smaller sample sizes and lower model df may be prone to decreased power
Parsimony-Adjusted Indices	PGFI	No thresholds recommended	
	PNFI	No thresholds recommended	
	GFI	Closer the value is to +1.0, the better the fit	Influenced by sample size and thus the interpretation of GFI discouraged
	AGFI	Closer the value is to +1.0, the better the fit	Influenced by sample size and thus the interpretation of AGFI discouraged

BOX 10.1 (continued)

Category	Indices	Interpretation	Notables
Incremental Fit Indices	CFI	Values greater than .95 suggest good model fit	
	TLI/NNFI	Not a range within which values for the TLI/NNFI can fall, however values closer to +1.0 are interpreted as better fit	
	NFI	NFI values above .95 suggest good model fit, between .90 and .95 marginal, and below .90 poor	As the number of parameters estimated increases, the size of the NFI also increases and thus the interpretation of NFI discouraged
	IFI	Ranges from 0 to +1.0, indicating the worst to the best model fit, respectively	
Comparative Fit Indices	AIC, BIC, SABIC	Lower AIC, BIC, and SABIC values indicate better model fit	<ul style="list-style-type: none"> Meaningful only when multiple models are generated and determining the better fit of the models is needed BIC and SABIC increase as the sample size increases; however, the SABIC does not increase at the same rate as the BIC. AIC and BIC increase as the number of parameters estimated increases; however, the AIC penalizes at a lower rate than the BIC AIC can be used to compare nonnested models
	ECVI	Meaningful only when multiple models are generated and determining the better fit of the models is needed	

Parameter Estimate Evaluation

We learned about factor loadings and unique and common variance in our discussion of EFA. These concepts still apply in CFA. A factor loading reflects the relationship between a latent construct and an observed variable (i.e., the direct effect of the latent variable on the manifest variable) and is the slope of increase (when positive) or decrease (when negative) for each unit of increase or decrease in the latent construct. Thus, the factor loading is simply a regression slope for predicting the observed variable from the latent construct and thus is interpreted as a regression coefficient. Standardized factor loadings, when the observed variable loads onto only one latent construct, are interpreted as correlations and represent the relationship between the observed and latent variables. When standardized factor loadings in this situation are squared, they are interpreted as the proportion of variance of the latent variable that is explained by the observed variable. A standardized factor loading of at least .7 is recommended as, when squared, this is interpreted as the latent variable is explaining about 50% of the variance

of the observed variable—a majority of the variance of the manifest variable (Kline, 2011). Standardized factor loadings, when the observed variable loads onto more than one latent construct, are interpreted as beta weights, and squaring the loadings in this situation will not provide an interpretation of variance accounted for.

10.1.1.3 CFA Model Modification

CFA is not unlike many other procedures where, after consideration of estimates in the initial model tested, model revision is necessary. A *specification search* is the process for finding specification errors and refitting the model with the purpose of modifying the original model so that a better fitting model results. In the case of CFA, these revisions identified through the specification search are conducted for reasons such as the following: (a) to improve overall goodness-of-fit, (b) to correct an improperly specified path, and (c) to correct unspecified or improperly specified correlated errors. When the researcher has done due diligence in specifying a model grounded in and supported with theory and prior empirical research, this is less likely to occur. However, and as an example, particularly in cases where only EFA empirical research exists, the number of latent constructs in CFA may differ from what is suggested in EFA due to the lack of ability to specify correlated errors.

Potential modifications may be suggested by examining the statistical significance and squared multiple correlations of the estimated parameters. Parameters that are not statistically significant *and* of little or no theoretical concern in the model can be fixed to zero. Small squared multiple correlations for each observed variable may be an indication that they are not contributing to the model, and removing the observed variable may be warranted. Standardized residuals above about +2.0 or +2.58 (i.e., above 2 or 3 standard deviations in absolute value terms, respectively) can suggest the model is not explaining a covariance, and freeing parameters may be needed.

More objective measures for modifications exist in most SEM software. For example, LISREL reports modification indices for all nonfree parameters that include the decrease in chi-square that would result if the modification were introduced in the model. It is immensely important, however, that researchers let theoretical and previous empirical research guide any modification introduced in the model. What may be considered a small fix, such as correlating measurement errors, must be well grounded, as its introduction alters the meaning of the model. Additionally, any changes recommended by the modification indices are data driven, and thus may not replicate in another sample, further highlighting the need for cross-validation.

While there is a lack of consensus on steps for determining modifications, there are suggestions on how to more systematically conduct a specification search. These include (Schumacker & Lomax, 2010, p. 67) (Copyright (© 2010) From A Beginner's Guide to Structural Equation Modeling by Schumacker & Lomax. Reproduced by permission of Taylor and Francis Group, LLC, a division of Informa plc):

1. Let substantive theory and prior research guide your model specification.
2. When you are satisfied that Rule 1 has been met, test your implied theoretical model and move to Rule 3.

3. Conduct a specification search, first on the measurement model, and then on the structural model.
4. For each model tested, look to see if the parameters are of the expected magnitude and direction, and examine several appropriate goodness-of-fit indices.
Steps 5 through 7 can be followed in an iterative fashion. For example, you might go from Step 5 to Step 6, and successively on to Steps 7, 6, 5, and so on.
5. Examine the statistical significance of the nonfixed parameters. Look to see if any nonfixed parameters should be fixed in a subsequent model.
6. Examine the modification indices, expected parameter change statistics. Look to see if any fixed parameters should be freed in a subsequent model.
7. Examine the standardized and raw residual matrix to see if anything suspicious is occurring (larger values for a particular observed variable).
8. Once you determine a final acceptable model, cross-validate it with a new sample, or use half of the sample to find a properly specified model and use the other half to check it (cross-validate index, or CVI), or report a single sample cross-validation index (ECVI) for alternative models (Cudeck & Browne, 1983; Kroonenberg & Lewis, 1982).

10.1.1.4 Structural Equation Modeling

Structural equation modeling combines CFA (i.e., measurement models) with path models (i.e., structural models). Paths between latent variables are estimated once the measurement model yields acceptable data to model fit (Schumacker & Lomax, 2016). There have been different approaches suggested for modeling structural equation models. One was proposed by Mulaik and Millsap (2000) and included the following four steps: (1) conduct an exploratory factor analytic model, (2) conduct a confirmatory factor analytic model, (3) specify relations among the latent variables (i.e., specify the structural model), and (4) determine acceptable fit of the structural model based on model fit indices. This makes it clear how CFA builds into a structural model. There are many types of structural models including, for example, multiple group models, second order CFA models, dynamic factor models, multiple indicator–multiple cause (MIMIC) models, mixture models, multilevel SEM, and latent growth models (Schumacker & Lomax, 2016). Readers interested in furthering their knowledge may want to consult any number of excellent sources (e.g., Byrne, 1998; Marcoulides & Schumacker, 2001; Schumacker & Lomax, 2016).

Multiple Group Models

Multiple group models are appropriate for testing for parameter estimate group differences within a latent variable framework. The theoretical model is established with the data first, and then separate models are generated by group. Nonstatistically significant chi-square tests for each group indicate similarity of the groups (Schumacker & Lomax, 2016).

Second Order CFA

In this CFA, the first order latent factor (labeled as ‘sub construct’) explains a higher order (i.e., second order) latent factor (labeled as ‘main construct’). An illustration of a second order CFA is presented in Figure 10.3.

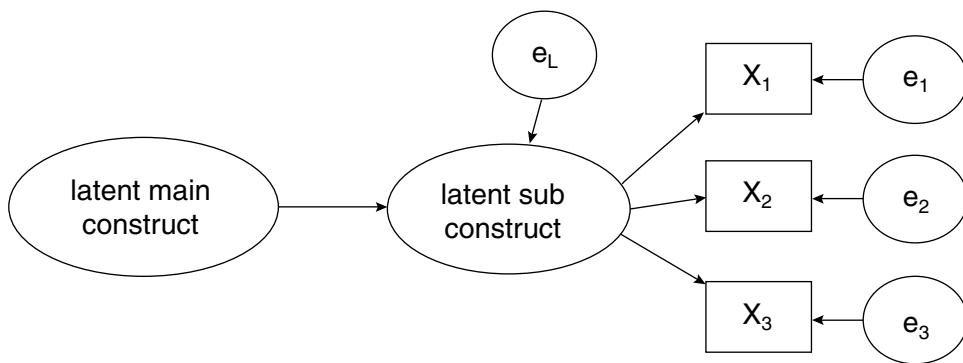


FIGURE 10.3
Second Order CFA Model

Dynamic Factor Models

These models included repeated measures—the same measurement completed by cases on two or more occasions. Dynamic factor models allow change in the latent variable to be assessed over time. Because of the repeated nature, the errors often correlate (referred to as autocorrelation), an important consideration in this modeling process (Schumacker & Lomax, 2016).

Multiple Indicator–Multiple Cause (MIMIC) Model

MIMIC models are characterized by latent variables predicted by observed variables, as illustrated in Figure 10.4. The variables that predict the latent variable are indicated on the left and are correlated, as indicated by the curved lines on the left. The measurement model is depicted by the latent variable (center) and the observed variables to the right.

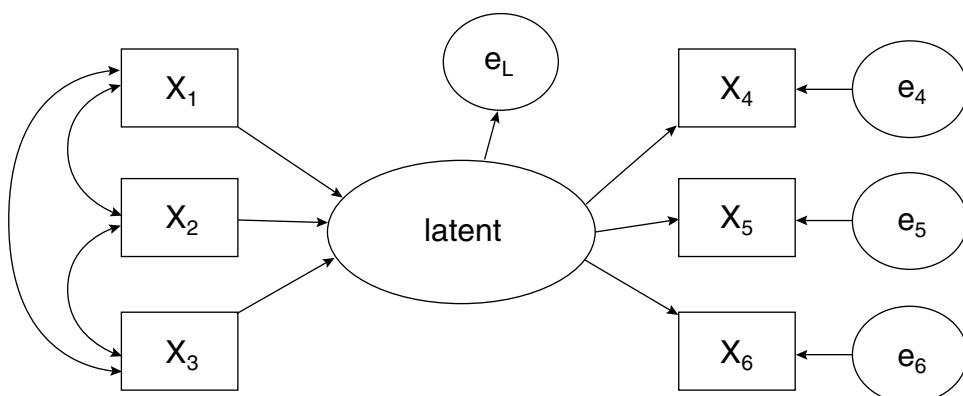


FIGURE 10.4
MIMIC Model

Mixed Variable and Latent Class Mixture Models

Though SEM was initially designed to apply continuous variables, mixed variable models allow both continuous and categorical variables to be modeled. These models can be accommodated by using matrices other than Pearson. A latent class mixture model refers to latent class analysis where latent groups are identified (Schumacker & Lomax, 2016).

Multilevel Models

Multilevel models can be estimated in a structural model framework using observed and/or latent variables, and by doing so, allow examination of variance of units within groups (i.e., between and within variation). With multilevel SEM, a null model is estimated first in order to generate the intraclass correlation coefficient (Schumacker & Lomax, 2016).

Latent Growth Models

A latent growth model (LGM) is a longitudinal design (i.e., repeated measures) that uses slopes (i.e., rate of change) and intercepts (i.e., means) to understand growth. In the case of linear trajectories, the model would reflect individual differences in the slopes and intercepts of the collection of straight lines for the sample of cases. In addition to understanding growth, the LGM can also be used to answer questions of prediction in relation to factors that relate to growth (Duncan & Duncan, 2009).

10.1.2 CFA Sample Size

Sample size requirements in CFA have tended to follow a similar path as that of EFA: These recommendations are generally based on a subject-to-parameter (referred to as subject-to-variable ratio in EFA) or absolute sample size per number of cases (N). The general recommendations that exist are very similar to guidelines for EFA, such as an absolute sample size of at least 200 (Boomsma, 1983; Boomsma & Hoogland, 2001), at least 10 cases per observed variable (Nunnally, 1967), or at least 5 cases for every parameter estimated (Bentler & Chou, 1987). More recently, and based in large part on the work of Marsh, Hau, Balla, and Grayson (1998), sample size recommendations for CFA (and structural equation modeling in general) are conditional on the data and model complexity (Boomsma & Hoogland, 2001; Gagne & Hancock, 2006). Additional research suggests an inverse relationship between the variable to construct ratio: As the number of observed variables per latent construct increases, the required sample size decreases, and this was most pronounced as the number of observed variables increased from three or four to six and less so when increasing from six to eight (Wolf, Harrington, Clark, & Miller, 2013). However, this does not operate in a vacuum, void of influence of model complexity. Contemporary research suggests that larger sample sizes, more observed variables per latent construct, and stronger factor loadings improve parameter estimation and model convergence (Gagne & Hancock,

2006). More current research has confirmed that of Marsh et al. (1998) in that there does not exist an absolute minimum sample size nor a subject to parameter ratio that when met will ensure model convergence (Gagne & Hancock, 2006). Rather, it is more effective to look at both the number of observed variables per latent construct and the magnitude of the factor loadings as a measure of “construct reliability” to determine sample size (Gagne & Hancock, 2006).

Although it is more challenging to offer sample size recommendations in this framework, Gagne and Hancock (2006) do offer tables based on construct reliability considerations (i.e., the number of observed variables per latent construct and the magnitude of the factor loadings). For example, under *homogenous* loading conditions (i.e., identical factor loadings for all observed variables), the following are a few examples of suggested sample size (Gagne & Hancock, 2006): (a) three observed variables per latent construct with factor loadings of .80, satisfactory convergence is achieved with a sample size as small as 50 and perfect convergence is achieved with a sample size of 100 (should the factor loading decrease to .40, a sample size of 400 is needed for satisfactory and 1,000 for perfect convergence); (b) five observed variables per latent construct with factor loadings of .80, satisfactory convergence is achieved with a sample size as small as 25 and perfect convergence is achieved with a sample size of 50 (should the factor loading decrease to .40, a sample size of 200 is needed for satisfactory and 400 for perfect convergence); and (c) seven observed variables per latent construct with factor loadings of .80, satisfactory convergence is achieved with a sample size as small as 25 and perfect convergence is achieved with a sample size of 50 (should the factor loading decrease to .40, a sample size of 100 is needed for satisfactory and 200 for perfect convergence). The need for large sample sizes in models with weaker factor loadings has been found by others, but models with very strong relationships also require large samples as compared to models with moderate factor loadings (Wolf et al., 2013).

Under *heterogeneous* loading conditions (i.e., first and second loading matrices are equal), the recommendations become much more difficult to summarize succinctly and readers are referred to Gagne and Hancock (2006).

The caveat to using the tables (Gagne & Hancock, 2006) is that the guidelines set forth within them were based on simulated data which did not examine nonnormality or correlated measurement errors, among other modeling issues. Additionally, the simulated models were based on only three factors, each with correlations between latent constructs of .30, and while the number of observed variables per latent construct varied across models tested, *within each model*, the number of observed variables per latent variable was constant and every factor loading was the same. Thus, while the sample size recommendations offer a place to start for estimating required sample size,

changing any one of these model characteristics could result in a different set of sample size recommendations (although one would not expect too different), [thus] discretion should be used by applied researchers when deviating substantially

from the current model configurations; one might thus wish to err on the side of larger sample sizes.

(Gagne & Hancock, 2006, p. 80)

One of the issues noted with the work of Gagne and Hancock (2006) is that model fit (e.g., RMSEA, chi-square) was not taken into account (Jackson, Voh, & Frey, 2013). When taking model fit into account, simulation research suggests the following (these are only a few examples of those provided by Jackson et al.): (a) three observed variables per latent construct with factor loadings of .80 and either three or six latent constructs will achieve only sufficient Swain corrected maximum likelihood chi-square ($ML_S\chi^2$) model fit with a sample size as small as 50 (a factor loading of .40 will achieve both ML and chi-square model fit with a sample of 400 or more); (b) five observed variables per latent construct with factor loadings of .80 and three latent constructs will achieve only sufficient $ML_S\chi^2$ model fit with a sample size as small as 50 (a factor loading of .40 will achieve both ML and $ML_S\chi^2$ model fit with a sample of 200 or more); and (c) seven observed variables per latent construct with factor loadings of .80 and three latent constructs will achieve only sufficient $ML_S\chi^2$ model fit with a sample size of 100 or greater (a factor loading of .40 will only achieve chi-square model fit with a sample of 100 or more). Jackson et al. (2013) do offer a range of conditions and required sample size to both meet convergence and satisfactory model fit, and thus may be helpful as a starting point to estimate requisite sample size. The caveat again to using the tables (Jackson et al., 2013) is that the guidelines set forth within them were based on simulated data, which did not examine all possible modeling issues.

10.1.3 CFA Effect Size

Squared multiple correlations are computed in CFA and represent the proportion of variance of the observed variable that is accounted for by the latent construct. Conceptually, this is equivalent to a communality estimate in EFA. Applying Cohen's (1988) guidelines for interpreting the values, squared multiple correlations of .01 are small effects, .09 are moderate, and .25 are large.

10.1.4 CFA Assumptions

CFA shares assumptions that are common to EFA and regression. These include (a) independence, (b) linearity, (c) multivariate normality and absence of outliers, and (d) lack of multicollinearity and singularity. A data condition for conventional SEM is continuous data (i.e., interval or ratio).

10.1.4.1 Independence

The first assumption is concerned with **independence** of the observations. Violations of this assumption can detrimentally impact standard error values and thus any resulting hypothesis tests. Testing for this assumption in CFA is tricky, as there are no

independent and dependent variables that allow for this type of examination. In the absence of statistical evidence, we will rely on theoretical evidence: If the units have been randomly sampled from a population, there is evidence that the assumption of independence has been met.

10.1.4.2 Linearity

As you recall, CFA and other procedures that fall under the SEM umbrella are regression-based procedures, and thus CFA uses relationships among the variables as the basis for determining factors with conventional CFA doing so via covariances. Therefore, it is assumed there is a linear relationship among the variables. Bivariate scatterplots can be examined to determine the extent to which this assumption is held. Violation of the linearity assumption can be detected through residual plots. The residuals should be located within a band of ± 2 standard errors, indicating no systematic pattern of points. Polynomial transformations may help in situations where nonlinear relationships are evident.

10.1.4.3 Multivariate Normality and Absence of Outliers (Both Bivariate and Multivariate)

Screening for both univariate and multivariate normality is important. The following can be used to detect univariate normality violations: frequency distributions, normal probability (Q-Q) plots, and skewness statistics. The simplest procedure involves checking for symmetry in a histogram, frequency distribution, boxplot, or skewness and kurtosis statistics. Although **nonzero kurtosis** (i.e., a distribution that is either flat, platykurtic, or has a sharp peak, leptokurtic) will have minimal effect on the parameter estimates, **nonzero skewness** (i.e., a distribution that is not symmetric with either a positive or negative skew) will have much more impact on these estimates. Thus, finding asymmetrical distributions is a must. One recommendation is to be concerned if the skewness value is larger than an absolute value of 2.0 in magnitude.

Another useful graphical technique is the normal probability plot (or Q-Q plot). With normally distributed residuals, the points on the normal probability plot will fall along a straight diagonal line, whereas nonnormal data will not. There is a difficulty with this plot because there is no criterion with which to judge deviation from linearity. It is recommended that skewness and/or the normal probability plot be considered at a minimum when determining normality evidence.

Transformations can also be used to normalize nonnormal data. However, again there is the problem of dealing with transformed variables measured along some other scale than that of the original variables.

Outliers in CFA operate in an unfavorable fashion, just as they do in other procedures. One or more outlying cases (either univariate or multivariate) can have undue and unwanted influence on the model. Outliers in CFA can be screened in a similar fashion

as we did with EFA. In addition to the ways we've screened for outliers in previous procedures (e.g., boxplots), they can also be screened by reviewing standard scores of the variables as we did when screening for outliers in EFA. Standardized scores with absolute values of 3.29 or greater (which equates to values more than 3–1/4 standard deviation units from the mean; about .05% of cases are above and below this point in a standardized normal distribution) should be flagged as outliers.

In CFA, as it was with EFA, it is also possible to have outlying variables—variables that are unrelated to others in the factor model. These outlying variables can be determined by reviewing the squared multiple correlations. Outlying variables that are identified can be disregarded. Maximum likelihood estimation is relatively robust to moderate departures of multivariate normality.

10.1.4.4 Lack of Extreme Multicollinearity and Singularity

As you recall, multicollinearity is a very strong linear relationship between two or more of the predictors. Given we are again working with matrix algebra and inversion, extreme multicollinearity is quite problematic. Singularity is a special case of multicollinearity; it is perfect multicollinearity and occurs when two or more items/variables perfectly predict and are therefore perfectly redundant. This can occur in CFA (just as it did in EFA and multiple regression), for example, when a composite variable as well as its component variables are used as predictors in the same factor analytic model (e.g., including GRETOT, GRE-Quantitative, and GRE-Verbal as predictors).

How do we detect violations of this assumption? The simplest procedure is to conduct a series of special regression analyses, one for each observed variables, where that variable, X , is predicted by all of the remaining X 's. If any of the resultant squared multiple correlation values are close to one (greater than .9 is a good rule of thumb), then there may be a collinearity problem.

Another statistical method for detecting collinearity is to compute a variance inflation factor (VIF) for each observed variable (using the same series of special regression analyses). The VIF is defined as the inflation that occurs for each regression coefficient above the ideal situation of uncorrelated predictors. Many suggest that the largest VIF should be less than 10 in order to satisfy this assumption (Myers, 1990; Stevens, 2009; Wetherill, 1986). Tolerance values, computed as $1 - R^2$, represent the proportion of total standardized variance unexplained by the other variables and should be greater than .10.

There are several possible methods for dealing with a collinearity problem. One solution is to simply remove one of the redundant variables from the model. Combining the variables to create a composite may also be an option.

10.1.4.5 Concluding Thoughts on Assumptions

As mentioned in previous chapters, there is no rule stating that research that violates assumptions must be scrapped. However, researchers who face violations of

assumptions must handle these situations on a case-by-case basis, considering both the goal of the analyses and the extent to which the assumptions were violated and the resulting effect of violation. It is also important that researchers present the evidence found, along with justification for decisions that were made. A summary of the assumptions and the effects of their violation for CFA is presented in Table 10.1.

■ TABLE 10.1

Assumptions and Violation of Assumptions: Confirmatory Factor Analysis

Assumption	Effect of Assumption Violation
Independence	Influences standard errors of the model and thus hypothesis tests
Linearity	Reduces interpretability of the CFA solution
Multivariate normality and absence of outlying cases and variables	Minimal effect when violated with exceptions including (a) when hypothesis testing is conducted as part of the CFA, (b) when maximum likelihood is used to estimate the model, and (c) with small sample sizes
Lack of extreme multi-collinearity	Reduces ability to separate effects of variables

10.2 MATHEMATICAL INTRODUCTION SNAPSHOT

The overarching goal of CFA is to estimate the parameters of the model that produce a predicted population covariance matrix ($\hat{\Sigma}$) that reproduces the sample variance-covariance matrix (\mathbf{S}). A very simplistic two-variable sample variance-covariance matrix is presented as follows:

$$\mathbf{S} = \begin{bmatrix} \text{cov}(X, X) & \text{cov}(Y, X) \\ \text{cov}(X, Y) & \text{cov}(Y, Y) \end{bmatrix}$$

This simplifies to the sample covariance matrix:

$$\mathbf{S} = \begin{bmatrix} \text{var } X \\ \text{cov}(X, Y) & \text{var } Y \end{bmatrix}$$

The model-based covariance matrix is as follows:

$$\Sigma(\theta) = \begin{bmatrix} b^2 \text{var}(X) + \text{var}(e) & b \text{var}(X) \\ b \text{var}(X) & \text{var}(X) \end{bmatrix}$$

If the model parameters perfectly produce the population covariance matrix, (Σ), then $\Sigma(\theta) = \Sigma$. Since the population parameters are not known, sample data are used to estimate the unknown parameters and predict the population variance covariance matrix.

Now in less technical language. The parameters in SEM (i.e., the regression coefficients and variances and covariances of the exogenous variables) are calculated from

the observed data and thus represent an estimate of the population. These estimated parameters are combined via matrix algebra to generate an estimated population covariance matrix. The difference between the estimated population covariance matrix and the sample covariance matrix is compared. The smaller the difference between the two (and hopefully this difference is also not statistically significant) the better, as this suggests consistency between the sample and population matrices. Zero in the residual matrix indicates that there is perfect fit between the sample covariance matrix and the predicted population covariance matrix. The goal is to minimize the difference (i.e., the residual matrix) between the sample covariance matrix and the model-implied covariance matrix.

The matrix equation for the latent dependent variable in a measurement model (i.e., CFA) is as follows:

$$Y = \Lambda_y \eta + \varepsilon$$

Where Y = latent dependent variable

Λ_y = capital lambda sub y , the matrix of the factor loadings for Y , i.e., the relationships between the observed variables and latent variables for Y

η = eta, the latent dependent variable which is a vector ($m \times 1$) of latent dependent variables

ε = lowercase epsilon, vector of measurement errors for Y

As this section is designed to simply provide a snapshot into the underlying mathematics of SEM, you may be left feeling it is a bit cursory. Readers interested in more detailed accounts of the mathematics may wish to review the Schumacker and Lomax (2010) chapter on the matrix approach to SEM.

10.3 COMPUTING PATH ANALYSIS AND CONFIRMATORY FACTOR ANALYSIS USING LISREL

We are branching out to new software programs to work with confirmatory factor analysis. We will be using LISREL and Mplus, software ideally suited for latent variable modeling and procedures that align (e.g., path analysis and CFA).

10.3.1 CFA Using LISREL

For this illustration, we'll use the free student edition, LISREL version 9.1. Before we conduct the analysis, let us talk about the data. The data we are using is the 2013 Survey of Adult Skills (<http://www.oecd.org/site/piaac/surveyofadultskills.htm>), available through the Organisation for Economic Co-operation and Development (OECD). This

data was introduced in the chapter on exploratory factor analysis and the specifics will not be reiterated here. Readers are referred back to that chapter for details.

We are using the **CH10_PIAAC_CFA.sav** file. This is data from the U.S., and the data file has been delimited to include only individuals who were between the ages of 25 and 29 [AGEG5LFS = 3], who were employed or participated in education or training during the 12 months prior to completing the survey [NEET = 0], and who reported having ‘above high school’ education [B_Q01a_T = 3] ($n = 288$). The size of this sample is sufficient to generate CFA, and it creates at least an intuitively homogenous sample that would be anticipated to respond similarly on the items. (Note: The complete PIAAC Survey of Adult Skills data file, which includes 5,010 cases, is available from the textbook’s companion website and is titled **PIAAC_SurveyOfAdultSkills.sav**.)

Before we run the data, it’s always important to examine frequency distributions of the variables that will be used in the model to assess missing data, potential data entry problems, and similar. With this data, we have some missing data (it has already been coded by the survey collectors as 9996), and thus I’ve taken the liberty to perform listwise deletion on the missing items (resulting in $n = 191$); however, the remaining variables in the data file have been left as is so that you may practice your data cleaning skills in working with ‘real data.’

Let’s look at the data. For the CFA illustration, we’ll be working with the first 8 of 13 indices (variables 1–8 in your SPSS file), each of which is measured on a continuous scale. We are selecting only the first eight, as those were the variables that were identified in the EFA model; the other five variables were screened out due to communality issues. (Note that we are using the same data for illustrative purposes as what was used for EFA; however, in practice a cross-validation sample—i.e., a second sample of data—is recommended.) All variables are retained in this data file, however, in the event you want to test other models in LISREL. Because the student version of LISREL supports only up to 16 observed variables, the dataset has been condensed to the following:

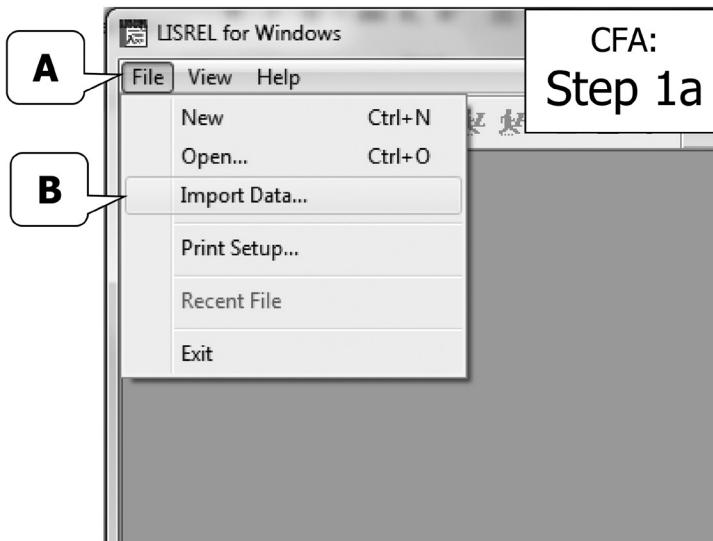
1. Index of use of numeracy skills at home (basic and advanced—derived)
2. Index of use of numeracy skills at work (basic and advanced—derived)
3. Index of use of ICT skills at home (derived)
4. Index of use of reading skills at home (prose and document texts—derived)
5. Index of use of task discretion at work (derived)
6. Index of learning at work (derived)
7. Index of use of planning skills at work (derived)
8. Index of readiness to learn (derived)
9. Index of use of ICT skills at work (derived)
10. Index of use of influencing skills at work (derived)
11. Index of use of reading skills at work (prose and document texts—derived)
12. Index of use of writing skills at work (derived)
13. Index of use of writing skills at home (derived)

The first 13 variables are the indices for CFA (*only the first eight of which will be used in this illustration*). This is followed by two variables in the dataset that represent the country and participant ID variables. I've left those in the data file just in case you are interested in merging variables from the full dataset with this smaller, delimited file. Each row in the data set still represents one individual. As seen in the screenshot below, the SPSS data is in the form of multiple columns that represent the variables on which the respondents were measured. For the CFA illustration, we will work with the first eight continuous index measures.

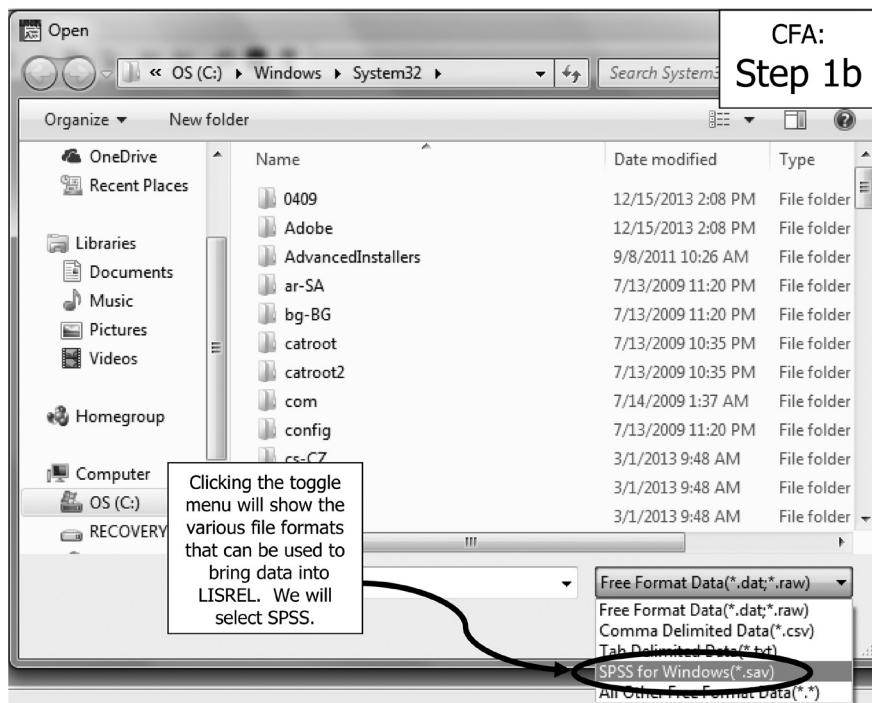
I	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNNETWORK	PLANNING	READYTOLEARN	ICTWORK	INFLUENCE	READWORK
1	3.81842	6.72630	3.09433	3.54036	2.23419	4.34696	3.82347	5.00418	4.05051	5.78929	6.21393
2	3.61768	2.13456	3.57861	3.38760	1.70196	1.80888	1.17814	3.21841	4.70515	2.84993	3.61270
3	2.01930	1.99907	2.89841	2.88970	3.16957	2.42758	2.66925	5.00418	2.85857	2.13151	3.56752
4	2.20013	2.41922	2.85982	4.0719	2.95203	1.36681	3.82347	2.61279	2.13376	5.78929	2.71635
5	2.99929	2.53842	3.68743	5.94446	1.92449	2.50955	2.56925	5.00418	6.28512	2.79571	2.66331
6	2.66278	3.61267	3.90046	2.9354	1.30146	4.34696	2.22069	2.29040	2.38056	2.71789	4.30735
7	2.79634	2.08671	2.25562	1.39575	1.63605	3.08547	2.22069	1.22396	1.71389	2.93267	2.86236
8	3.06415	2.96927	3.22884	4.5796	2.98049	2.42758	2.22069	2.35854	4.28065	2.63603	3.25087
9	3.24173	3.68996	2.94005	3.0577	2.07432	4.34696	1.56394	3.07297	3.11809	3.05326	4.53684
10	2.44410	1.73260	1.61112	2.68728	2.15631	4.34696	1.43576	2.64003	2.38056	1.76144	2.86718

We will conduct CFA using the first 8 index measures.

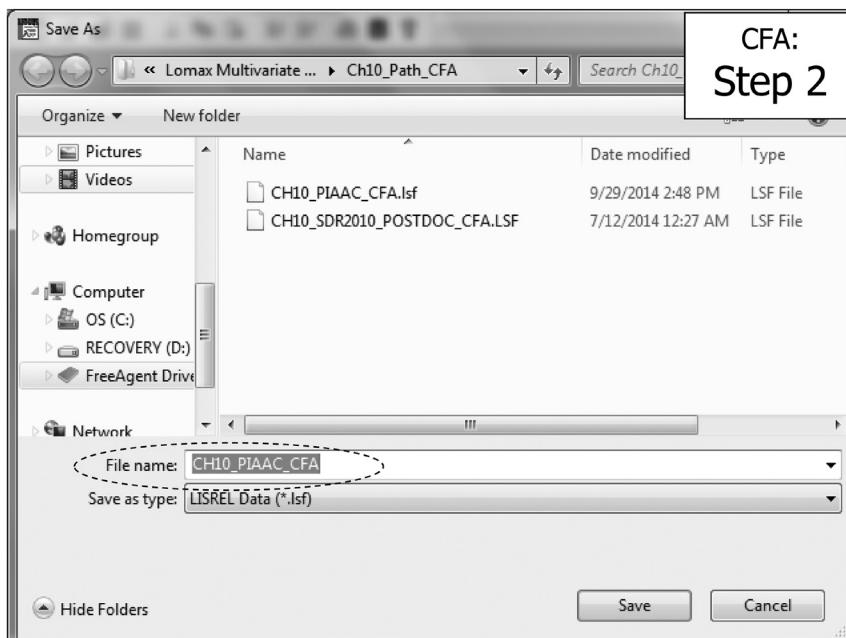
Step 1a. To conduct confirmatory factor analysis using LISREL, we have to first bring in our data. To do that, go to “File” in the top pull-down menu, then select “Import Data.” Following the screenshot below (see screenshot Step 1a) will bring up the dialog box that allows you to search and find your data file (see screenshot Step 1b).



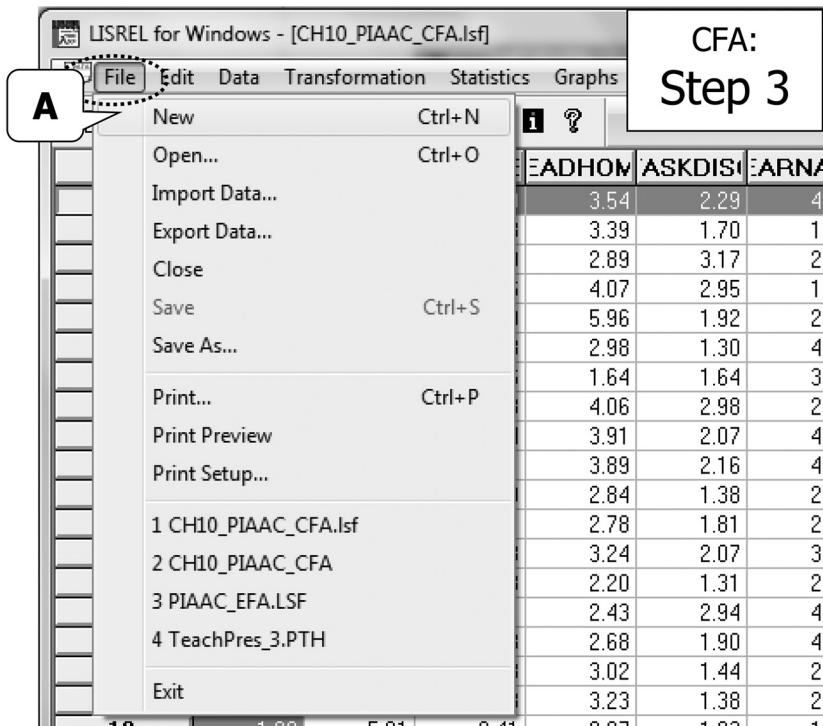
Step 1b. The dialog box that allows you to search and find your data file will pop up (see screenshot Step 1b). A number of different file formats can successfully be imported into LISREL. For this illustration, we use the toggle menu to select SPSS format.



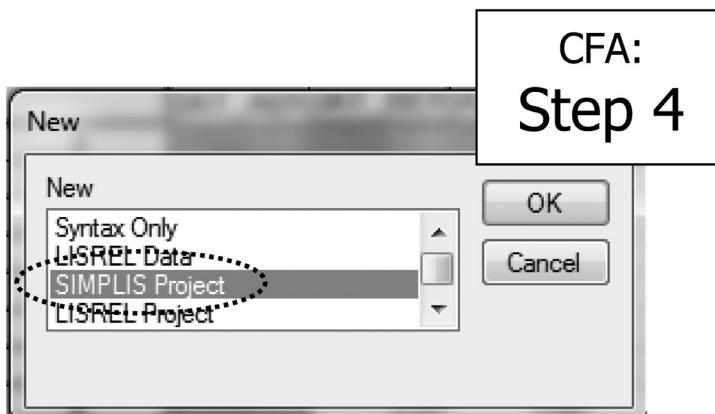
Step 2. Once we browse and select the filename, a “Save As” dialog box will appear. At this point, we will define our data file name and then click “Save.” The format it will be saved in will be recognized as data within LISREL (specifically, .lsf) (see screenshot Step 2).



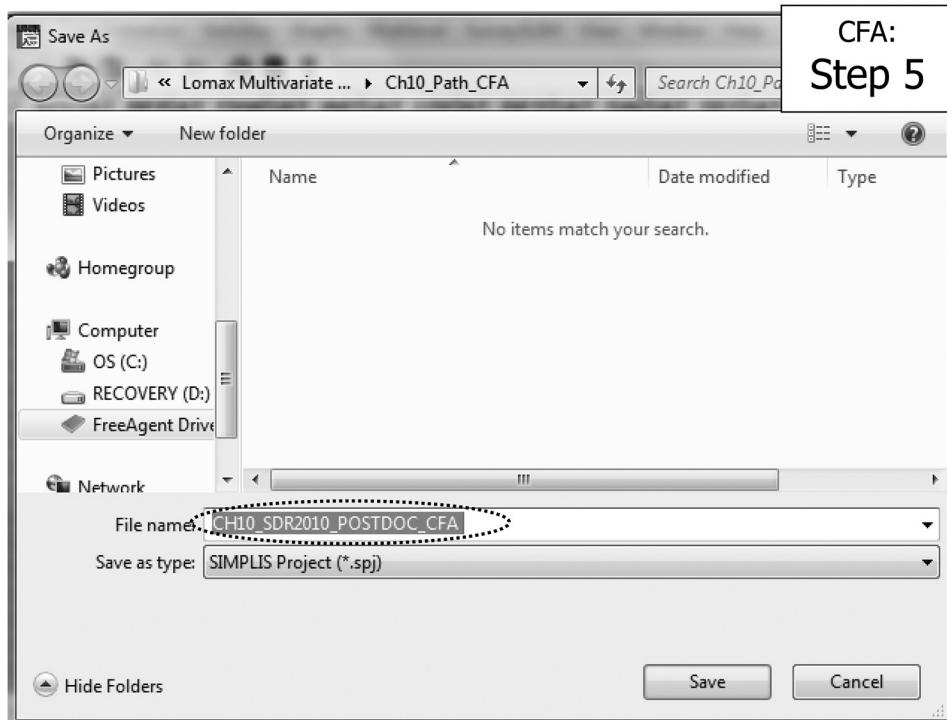
Step 3. From the LISREL drop down menu (see screenshot Step 3), click on "File" then "New" to bring up various options for generating statistics in LISREL.



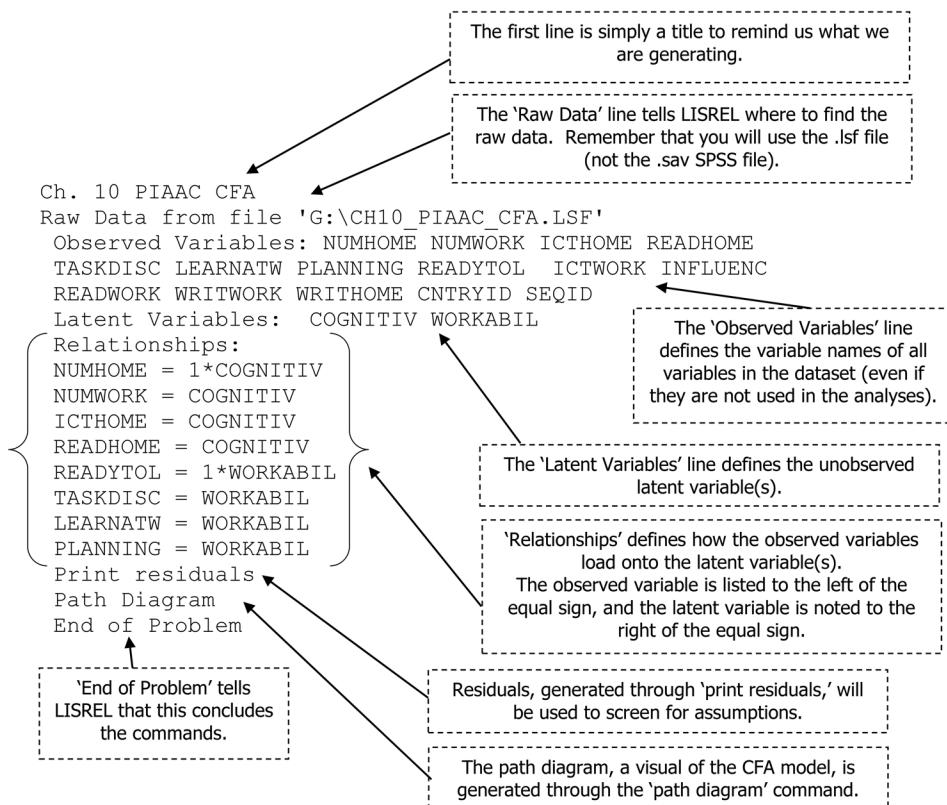
Step 4. From the New drop down menu (see screenshot Step 4), click on "SIMPLIS Project" then click OK.



Step 5. A dialog box will appear that will allow you to save the name of your new SIMPLIS project (see screenshot Step 5). After defining the file name, click OK to produce a blank SIMPLIS project file.



Step 6. A blank SIMPLIS field will appear, within which you will tell LISREL how to generate the CFA. The screenshot illustrates how the code will appear in LISREL—as simple text. The annotated code will help you understand how to write the code for your own data. LISREL can be a bit finicky. Make sure that the names of the observed variables appear exactly as they do in your data file (note that LISREL truncates to a maximum of 8 characters, even if the variable name was longer in SPSS; I recommend creating your SPSS file so that the variable names—as in the column headers—are no longer than 8 characters so that when you bring your file into LISREL, it will be ready to work with). The SIMPLIS command language for the original model tested, based on a two-factor solution evidenced in EFA, is presented here. It is important to note in the LISREL syntax that no factor correlation or factor variances are specified. This is because the program makes assumptions about these parameters being free. Not all software programs do this. When using LISREL, therefore, these parameters do not need to be specified, while they may need to be specified if using a different software program.



The SIMPLIS command language for the *revised* model tested, based on a one-factor solution with correlated measurement errors, follows.

```

Ch. 10 PIAAC CFA
Raw Data from file 'G:\CH10_PIAAC_CFA.LSF'
Observed Variables: NUMHOME NUMWORK ICTHOME READHOME
TASKDISC LEARNATW PLANNING READYTOL ICTWORK INFLUENC
READWORK WRITWORK WRITHOME CNTRYID SEQID
Latent Variables: SKILABIL
Relationships:
NUMHOME = 1*SKILABIL
NUMWORK = SKILABIL
ICTHOME = SKILABIL
READHOME = SKILABIL
READYTOL = SKILABIL
TASKDISC = SKILABIL
  
```

```

LEARNATW = SKILABIL
PLANNING = SKILABIL
{ Let error covariance between READYTOL and NUMHOME correlate
  Let error covariance between READYTOL and TASKDISC correlate
  Let error covariance between READYTOL and LEARNATW correlate
  Print residuals
  Path Diagram
  End of Problem
}

```

The additional three lines of code correlate the measurement errors between these observed variables

For illustrative purposes, we have specified correlated errors in the revised model as the fit of the model improved (RMSEA from .06 to .02, for example). However, it is immensely important that researchers have theoretical and/or previous empirical support to respecify the model in any way, including correlating errors. If you are working with full version of LISREL, you'll notice that modification indices are provided in the output. Modification indices are not provided through the student version of LISREL.

10.3.1.1 Interpreting the Output

Annotated results are presented in Tables 10.2 (initial two-factor model tested), 10.3 (respecified two-factor model with no correlated errors), and 10.4 (respecified two-factor model with correlated errors). Corresponding path diagrams are respectively presented in Figures 10.2 (presented earlier in the chapter), 10.5, and 10.6. Note in the final two-factor model with correlated errors, we have in the sample variance-covariance matrix \mathbf{S} the following number of distinct values given that there are eight observed variables:

$$\frac{p(p+1)}{2} = \frac{8(8+1)}{2} = 36$$

A count of the free parameters (i.e., those parameters that are being estimated in the model) is as follows:

- 7 factor loadings (recall one factor loading was fixed to 1.0)
- 7 measurement error variances
- 3 measurement error covariance correlations

Thus, there are 17 freely estimated parameters (this corresponds to the degrees of freedom reported in the output under goodness-of-fit statistics). Therefore, in terms of model identification, this model is *overidentified*, since there are more parameters in the sample variance-covariance matrix than parameters to be estimated (i.e., the free parameters).

■ TABLE 10.2

LISREL Results for the Confirmatory Factor Analysis Example

L I S R E L 9.10 (STUDENT)

BY

Karl G. Jöreskog & Dag Sörbom

This program is published exclusively by
Scientific Software International, Inc.
<http://www.ssicentral.com>

Copyright by Scientific Software International, Inc., 1981-2012
Use of this program is subject to the terms specified in the
Universal Copyright Convention.

The following lines were read from file G:\Lomax Multivariate
Textbook\Ch10_Path_CFA\CH10_PIAAC_CFA_Orig.spj:

Ch. 10 PIAAC CFA
Raw Data from file 'G:\Lomax Multivariate
Textbook\Ch10_Path_CFA\CH10_PIAAC_CFA.LSF'
Observed Variables: NUMHOME NUMWORK ICTHOME READHOME
TASKDISC LEARNATW PLANNING READYTOL ICTWORK INFLUENC
READWORK WRITWORK WRITHOME CNTRYID SEQID
Latent Variables: COGNITIV WORKABIL
Relationships:
NUMHOME = 1*COGNITIV
NUMWORK = COGNITIV
ICTHOME = COGNITIV
READHOME = COGNITIV
READYTOL = 1*WORKABIL
TASKDISC = WORKABIL
LEARNATW = WORKABIL
PLANNING = WORKABIL
Print residuals
Path Diagram
End of Problem

The SIMPLIS code is
reprinted in the
output.

Sample Size = 191

Ch. 10 PIAAC CFA

Covariance Matrix

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
NUMWORK	0.276	1.218				

■ TABLE 10.2 (continued)

LISREL Results for the Confirmatory Factor Analysis Example

ICTHOME	0.368	0.136	0.564			
READHOME	0.258	0.117	0.203	0.499		
TASKDISC	0.014	0.078	0.004	0.026	0.592	0.934
LEARNATW	0.011	0.100	0.035	0.081	0.019	0.135
PLANNING	-0.186	0.019	-0.094	-0.034	0.112	0.170
READYTOL	0.219	0.193	0.261	0.208	0.171	

Covariance Matrix

	PLANNING	READYTOL
PLANNING	1.198	
READYTOL	0.036	1.090

Total Variance = 6.862 Generalized Variance = 0.0675

Largest Eigenvalue = 1.845 Smallest Eigenvalue = 0.267

Condition Number = 2.628

Ch. 10 PIAAC CFA

Number of Iterations = 16

LISREL Estimates (Maximum Likelihood)

Measurement Equations

NUMHOME = 1.000*COGNITIV, Errorvar.= 0.327 , R² = 0.573Standerr (0.0629)
Z-values 5.204
P-values 0.000NUMWORK = 0.513*COGNITIV, Errorvar.= 1.102 , R² = 0.0950Standerr (0.117)
Z-values 9.450
P-values 0.000ICTHOME = 0.826*COGNITIV, Errorvar.= 0.265 , R² = 0.531Standerr (0.0455)
Z-values 5.810
P-values 0.000READHOME = 0.583*COGNITIV, Errorvar.= 0.349 , R² = 0.300Standerr (0.0415)
Z-values 8.411
P-values 0.000TASKDISC = 0.101*WORKABIL, Errorvar.= 0.575 , R² = 0.0290Standerr (0.0612)
Z-values 9.401
P-values 0.000

■ TABLE 10.2 (continued)

LISREL Results for the Confirmatory Factor Analysis Example

LEARNATW = 0.0971*WORKABIL, Errorvar.= 0.918 , R² = 0.0171
 Standerr (0.111) (0.0950)
 Z-values 0.872 9.664
 P-values 0.383 0.000

PLANNING = 0.0250*WORKABIL, Errorvar.= 1.197 , R² = 0.000886
 Standerr (0.0551) (0.122)
 Z-values 0.455 9.775
 P-values 0.649 0.000

READYTOL = 1.000*WORKABIL, Errorvar.= -0.603 , R² = 1.553
 Standerr (1.721)
 Z-values -0.350
 P-values 0.726

W_A_R_N_I_N_G : Error variance is negative.

Covariance Matrix of Independent Va

	COGNITIV	WORKABIL
COGNITIV	0.440 (0.088) 4.970	
WORKABIL	0.282 (0.064) 4.396	1.693 (1.718) 0.986

The warning of a negative error variance for 'readiness to learn' suggests we have a problem. Bandaid solutions, such as fixing the error variance to zero, are not recommended as they do not address the root of the problem which may be insufficient sample size, non-normality, multicollinearity, or a misspecified model. In this example, we will respecify the model to a one-factor solution.

Log-likelihood

Estimated Model

Number of free parameters(t)	17
-2ln(L)	1041.564
AIC (Akaike, 1974)*	1075.564
BIC (Schwarz, 1978)*	1130.853

1085.245
1202.327

*LISREL uses AIC= 2t - 2ln(L) and BIC = tln(N) - 2ln(L)

Goodness of Fit Statistics

Degrees of Freedom for (C1)-(C2)	19
Maximum Likelihood Ratio Chi-Square (C1)	28.319 (P = 0.0775)
Browne's (1984) ADF Chi-Square (C2_NT)	27.791 (P = 0.0875)
Estimated Non-centrality Parameter (NCP)	9.319
90 Percent Confidence Interval for NCP	(0.0 ; 27.703)
Minimum Fit Function Value	0.148
Population Discrepancy Function Value (F0)	0.0488
90 Percent Confidence Interval for F0	(0.0 ; 0.145)
Root Mean Square Error of Approximation (RMSEA)	0.0507
90 Percent Confidence Interval for RMSEA	(0.0 ; 0.0874)
P-Value for Test of Close Fit (RMSEA < 0.05)	0.450

■ TABLE 10.2 (continued)

LISREL Results for the Confirmatory Factor Analysis Example

Expected Cross-Validation Index (ECVI)	0.326
90 Percent Confidence Interval for ECVI	(0.277 ; 0.423)
ECVI for Saturated Model	0.377
ECVI for Independence Model	1.323
Chi-Square for Independence Model (28 df)	236.775
Normed Fit Index (NFI)	0.880
Non-Normed Fit Index (NNFI)	0.934
Parsimony Normed Fit Index (PNFI)	0.597
Comparative Fit Index (CFI)	0.955
Incremental Fit Index (IFI)	0.957
Relative Fit Index (RFI)	0.823
Critical N (CN)	243.819
Root Mean Square Residual (RMR)	0.0550
Standardized RMR	0.0605
Goodness of Fit Index (GFI)	0.965
Adjusted Goodness of Fit Index (AGFI)	0.933
Parsimony Goodness of Fit Index (PGFI)	0.509

Ch. 10 PIAAC CFA

Fitted Covariance Matrix

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
NUMWORK	0.226	1.218				
ICTHOME	0.363	0.186	0.564			
READHOME	0.256	0.131	0.212	0.499		
TASKDISC	0.028	0.015	0.023	0.017	0.592	
LEARNATW	0.027	0.014	0.023	0.016	0.017	0.934
PLANNING	0.007	0.004	0.006	0.004	0.004	0.004
READYTOL	0.282	0.145	0.233	0.165	0.171	0.164

Fitted Covariance Matrix

	PLANNING	READYTOL
PLANNING	1.198	
READYTOL	0.042	1.090

Fitted Residuals

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.000					
NUMWORK	0.051	0.000				
ICTHOME	0.005	-0.050	0.000			
READHOME	0.002	-0.015	-0.009	0.000		
TASKDISC	-0.015	0.064	-0.019	0.009	0.000	
LEARNATW	-0.017	0.086	0.012	0.065	0.003	0.000
PLANNING	-0.193	0.016	-0.100	-0.038	0.107	0.131
READYTOL	-0.063	0.048	0.028	0.044	0.000	0.006

■ TABLE 10.2 (continued)

LISREL Results for the Confirmatory Factor Analysis Example

Fitted Residuals

	PLANNING	READYTOL
PLANNING	0.000	
READYTOL	-0.006	0.000

Summary Statistics for Fitted Residuals

Smallest Fitted Residual = -0.193
 Median Fitted Residual = 0.000
 Largest Fitted Residual = 0.131

Stemleaf Plot

- 1|9
 - 1|0
 - 0|65
 - 0|422111100000000000
 0|1111234
 0|55679
 1|13

Standardized Residuals

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.000					
NUMWORK	1.646	0.000				
ICTHOME	0.745	-1.709	0.000			
READHOME	0.132	-0.357	-0.716	0.000		
TASKDISC	-0.404	1.077	-0.599	0.275	0.000	
LEARNATW	-0.325	1.146	0.264	1.426	0.057	0.000
PLANNING	-2.860	0.180	-1.724	-0.689	1.783	1.725
READYTOL	-2.437	0.700	1.127	1.217	-0.013	0.562

Standardized Residuals

	PLANNING	READYTOL
PLANNING	0.000	
READYTOL	-0.435	0.000

Summary Statistics for Standardized Residuals

Smallest Standardized Residual = -2.860
 Median Standardized Residual = 0.000
 Largest Standardized Residual = 1.783

Stemleaf Plot

- 2|9
 - 2|4
 - 1|77
 - 1|
 - 0|776

■ TABLE 10.2 (continued)

LISREL Results for the Confirmatory Factor Analysis Example

```
- 0|4443000000000  
0|11233  
0|677  
1|11124  
1|678  
Largest Negative Standardized Residuals  
Residual for PLANNING and NUMHOME -2.860
```

Time used 0.047 seconds

■ TABLE 10.3

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

L I S R E L 9.10 (STUDENT)

BY

Karl G. Jöreskog & Dag Sörbom

This program is published exclusively by
Scientific Software International, Inc.
<http://www.ssicentral.com>

Copyright by Scientific Software International, Inc., 1981-2012
Use of this program is subject to the terms specified in the
Universal Copyright Convention.

The following lines were read from file G:\Lomax Multivariate
Textbook\Ch10_Path_CFA\CH10_PIAAC_CFA.spj:

```
Ch. 10 PIAAC CFA  
Raw Data from file 'G:\Lomax Multivariate  
Textbook\Ch10_Path_CFA\CH10_PIAAC_CFA.LSF'  
Observed Variables: NUMHOME NUMWORK ICTHOME READHOME  
TASKDISC LEARNATW PLANNING READYTOL ICTWORK INFLUENC  
READWORK WRITWORK WRITHOME CNTRYID SEQID  
Latent Variables: SKILABIL  
Relationships:  
NUMHOME = 1*SKILABIL  
NUMWORK = SKILABIL  
ICTHOME = SKILABIL  
READHOME = SKILABIL  
READYTOL = SKILABIL  
TASKDISC = SKILABIL  
LEARNATW = SKILABIL  
PLANNING = SKILABIL  
Print residuals  
Path Diagram  
End of Problem
```

The SIMPLIS code
is reprinted in the
output. We have
respecified the
model to be a one-
factor solution and
have renamed the
factor or latent
construct to be
"Skills and Abilities."

Sample Size = 191

■ TABLE 10.3 (continued)

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

Ch. 10 PIAAC CFA

Covariance Matrix

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
NUMWORK	0.276	1.218				
ICTHOME	0.368	0.136	0.564			
READHOME	0.258	0.117	0.203	0.499		
TASKDISC	0.014	0.078	0.004	0.026	0.592	
LEARNATW	0.011	0.100	0.035	0.081	0.019	0.934
PLANNING	-0.186	0.019	-0.094	-0.034	0.112	0.135
READYTOL	0.219	0.193	0.261	0.208	0.171	0.170

Covariance Matrix

	PLANNING	READYTOL
PLANNING	1.198	
READYTOL	0.036	1.090

Total Variance = 6.862 Generalized Variance = 0.0675

Largest Eigenvalue = 1.845 Smallest Eigenvalue = 0.0675

Condition Number = 2.628

ML estimates for the model are presented in this section.

Ch. 10 PIAAC CFA

Number of Iterations = 7

LISREL Estimates (Maximum Likelihood)

Measurement Equations

NUMHOME = 1.000 * SKILABIL, Errorvar. = 0.320 , R² = 0.582Standerr (0.0622)
Z-values 5.149
P-values 0.000NUMWORK = 0.514 * SKILABIL, Errorvar. = 1.100 , R² = 0.0969Standerr (0.116)
Z-values 9.447
P-values 0.000ICTHOME = 0.812 * SKILABIL, Errorvar. = 0.270 , R² = 0.521Standerr (0.0449)
Z-values 6.021
P-values 0.000READHOME = 0.579 * SKILABIL, Errorvar. = 0.349 , R² = 0.300Standerr (0.0414)
Z-values 8.423
P-values

Covariance matrix of the observed variables is provided

Each observed variable is expressed as a linear function of the latent construct. NUMHOME was the indicator variable reflected in a parameter estimate of 1.0 with no standard error reported. The unstandardized parameter estimate is akin to a regression coefficient. For NUMWORK, the magnitude of the unstandardized parameter estimate (.514) indicates the change in the latent construct for a one unit change in the observed variable. Changing NUMWORK by one unit results in about ½ unit decrease in SKILABIL.

The error variance is the measurement error.

Evaluate the individual parameter estimates to determine individual parameter fit to the model. Observed variables with statistically significant parameter estimates indicate they are related to the underlying latent construct. Non-statistically significant manifest variables may be considered for removal in the model modification process, however theory should guide this choice.

■ TABLE 10.3 (continued)

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

P-values	0.000	0.000	
TASKDISC	= 0.0713*SKILABIL, Errorvar.= 0.590	, R ² = 0.00383	
Standerr	(0.0953)	(0.0604)	
Z-values	0.748	9.761	
P-values	0.455	0.000	
LEARNATW	= 0.126*SKILABIL, Errorvar.= 0.927	, R ² = 0.00756	
Standerr	(0.120)	(0.0951)	
Z-values	1.049	9.750	
P-values	0.294	0.000	
PLANNING	= - 0.265*SKILABIL, Errorvar.= 1.167	, R ² = 0.0261	
Standerr	(0.136)	(0.120)	
Z-values	-1.944	9.692	
P-values	0.052	0.000	
READYTOL	= 0.634*SKILABIL, Errorvar.= 0.910	R ² = 0.165	
Standerr	(0.134)	(0.0993)	
Z-values	4.747	9.167	
P-values	0.000	0.000	
Variances of Independent Variables			
SKILABIL			

0.447			
(0.089)			
5.045			
As seen on the path diagram, the variance of the latent construct			
Log-likelihood Values			
	Estimated Model	Saturated Model	
-----	-----	-----	
Number of free parameters(t)	16	36	
-2ln(L)	1050.546	1013.245	
{ AIC (Akaike, 1974)*	1082.546	1085.245	
BIC (Schwarz, 1978)*	1134.582	1202.327	

*LISREL uses AIC= 2t - 2ln(L) and BIC = tln(N) - 2ln(L)

R squared is the squared multiple correlation of the observed variable with the latent construct. Adhering to Cohen's (1988) guidelines for effect, .10 = small, .30 = moderate, and .5 = large for correlations thus for squared correlations, .01 = small, .09 = moderate, and .25 = large.

As seen on the path diagram, the variance of the latent construct

Comparative measures of fit (AIC, BIC)

■ TABLE 10.3 (continued)

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

Goodness of Fit Statistics		<p>Reject null hypothesis for the model chi-square (labeled "ML ratio chi-square"), which indicates poor fit. RMSEA and its CI suggest reasonable fit; while the lower bound is > 0, the upper bound is = .10 which is just within what is considered acceptable.</p> <p>The CFI is < .95 which suggests poor fit. The SRMR is < .08 which suggests good fit.</p>	
Degrees of Freedom for (C1)–(C2)	20		
Maximum Likelihood Ratio Chi-Square (C1)	37.300 ($P = 0.0105$)		
Browne's (1984) ADF Chi-Square (C2_NT)	40.261 ($P = 0.0046$)		
Estimated Non-centrality Parameter (NCP)	17.300		
90 Percent Confidence Interval for NCP	(3.885 ; 38.518)		
Minimum Fit Function Value	0.195		
Population Discrepancy Function Value (FO)	0.0906		
90 Percent Confidence Interval for FO	(0.0203 ; 0.202)		
Root Mean Square Error of Approximation (RMSEA)	0.0673		
90 Percent Confidence Interval for RMSEA	(0.0319 ; 0.100)		
P-Value for Test of Close Fit (RMSEA < 0.05)	0.183		
{ Expected Cross-Validation Index (ECVI) 90 Percent Confidence Interval for ECVI ECVI for Saturated Model ECVI for Independence Model		0.363 (0.293 ; 0.474)	Comparative measures of fit (ECVI)
Chi-Square for Independence Model (28 df)	236.775		
{ Normed Fit Index (NFI) Non-Normed Fit Index (NNFI) Parsimony Normed Fit Index (PNFI) Comparative Fit Index (CFI) Incremental Fit Index (IFI) Relative Fit Index (RFI)		0.842 0.883 0.601 0.917 0.920 0.778	Incremental fit indices (NFI, NNFI, CFI, IFI) and parsimony-adjusted index (PNFI)
Critical N (CN)	192.354		
{ Root Mean Square Residual (RMR) Standardized RMR Goodness of Fit Index (GFI) Adjusted Goodness of Fit Index (AGFI) Parsimony Goodness of Fit Index (PGFI)		0.0603 0.0652 0.950 0.910 0.528	Absolute fit indices (RMR, SRMR) Parsimony-adjusted index (PGFI)

TABLE 10.3 (continued)

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

Ch. 10 PIAAC CFA

Fitted Covariance Matrix						
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
NUMWORK	0.230	1.218				
ICTHOME	0.363	0.186	0.564			
READHOME	0.259	0.133	0.210	0.499		
TASKDISC	0.032	0.016	0.026	0.018	0.592	
LEARNATW	0.056	0.029	0.046	0.033	0.004	0.934
PLANNING	-0.118	-0.061	-0.096	-0.068	-0.008	-0.015
READYTOL	0.283	0.146	0.230	0.164	0.020	0.036

Fitted Covariance Matrix						
	PLANNING	READYTOL				
PLANNING	1.198					
READYTOL	-0.075	1.090				

Fitted Residuals						
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.000					
NUMWORK	0.047	0.000				
ICTHOME	0.006	-0.050	0.000			
READHOME	-0.001	-0.016	-0.007	0.000		
TASKDISC	-0.018	0.062	-0.022	0.008	0.000	
LEARNATW	-0.046	0.071	-0.011	0.048	0.015	0.000
PLANNING	-0.068	0.080	0.002	0.035	0.120	0.150
READYTOL	-0.064	0.047	0.031	0.044	0.150	0.134

Fitted Residuals						
	PLANNING	READYTOL				
PLANNING	0.000					
READYTOL	0.111	0.000				

Summary Statistics for Fitted Residuals						
Smallest Fitted Residual =	-0.068					
Median Fitted Residual =	0.001					
Largest Fitted Residual =	0.150					

Stemleaf Plot

```

- 6|84
- 4|06
- 2|2
- 0|861710000000
 0|2685
 2|15
 4|4778

```

The fitted covariance matrix is the model-implied covariance matrix compared to the sample covariance matrix. The values are thus the estimated variances and covariances based on the model.

The fitted residuals represent the difference between the sample covariances and the fitted covariances. For example, NUMWORK fitted residual of .047 is simply .276 - .230 (with difference based on rounding). Negative fitted residuals indicate *overfitting* where the actual covariance has been overestimated. Positive fitted residuals indicate *underfitting* where the actual covariance has been underestimated. Fitted residuals may be difficult to interpret as they are dependent on the metric of the variable.

■ TABLE 10.3 (continued)

LISREL Results for the Confirmatory Factor Analysis Example (No Correlated Errors)

6 21																																																																
8 0																																																																
10 1																																																																
12 04																																																																
14 00																																																																
Standardized Residuals	{	}																																																														
<table border="0"> <thead> <tr><th></th><th>NUMHOME</th><th>NUMWORK</th><th>ICTHOME</th><th>READHOME</th><th>TASKDISC</th><th>LEARNATW</th></tr> </thead> <tbody> <tr><td>NUMHOME</td><td>0.000</td><td></td><td></td><td></td><td></td><td></td></tr> <tr><td>NUMWORK</td><td>1.545</td><td>0.000</td><td></td><td></td><td></td><td></td></tr> <tr><td>ICTHOME</td><td>0.799</td><td>-1.666</td><td>0.000</td><td></td><td></td><td></td></tr> <tr><td>READHOME</td><td>-0.062</td><td>-0.394</td><td>-0.542</td><td>0.000</td><td></td><td></td></tr> <tr><td>TASKDISC</td><td>-0.756</td><td>1.079</td><td>-0.911</td><td>0.248</td><td>0.000</td><td></td></tr> <tr><td>LEARNATW</td><td>-1.526</td><td>0.996</td><td>-0.371</td><td>1.265</td><td>0.289</td><td>0.000</td></tr> <tr><td>PLANNING</td><td>-2.052</td><td>0.997</td><td>0.051</td><td>0.808</td><td>2.009</td><td>2.004</td></tr> <tr><td>READYTOL</td><td>-2.510</td><td>0.688</td><td>1.208</td><td>1.226</td><td>2.925</td><td>2.088</td></tr> </tbody> </table>		NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW	NUMHOME	0.000						NUMWORK	1.545	0.000					ICTHOME	0.799	-1.666	0.000				READHOME	-0.062	-0.394	-0.542	0.000			TASKDISC	-0.756	1.079	-0.911	0.248	0.000		LEARNATW	-1.526	0.996	-0.371	1.265	0.289	0.000	PLANNING	-2.052	0.997	0.051	0.808	2.009	2.004	READYTOL	-2.510	0.688	1.208	1.226	2.925	2.088	<p>Standardized residuals are computed as the fitted residual divided by their standard error and are interpreted as a conventional standardized value (i.e., mean = 0, $SD = 1$). Thus, conventions of 'large' standardized values can be applied (e.g., $\geq \pm 2.0$ or ± 3.0).</p> <p>The standardized residuals are particularly helpful for diagnosing potential modifications to the model.</p>
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW																																																										
NUMHOME	0.000																																																															
NUMWORK	1.545	0.000																																																														
ICTHOME	0.799	-1.666	0.000																																																													
READHOME	-0.062	-0.394	-0.542	0.000																																																												
TASKDISC	-0.756	1.079	-0.911	0.248	0.000																																																											
LEARNATW	-1.526	0.996	-0.371	1.265	0.289	0.000																																																										
PLANNING	-2.052	0.997	0.051	0.808	2.009	2.004																																																										
READYTOL	-2.510	0.688	1.208	1.226	2.925	2.088																																																										
Summary Statistics for Standardized Residuals	{	}																																																														
<table border="0"> <tbody> <tr><td>Smallest Standardized Residual =</td><td>-2.51</td></tr> <tr><td>Median Standardized Residual =</td><td>0.02</td></tr> <tr><td>Largest Standardized Residual =</td><td>2.92</td></tr> </tbody> </table>	Smallest Standardized Residual =	-2.51	Median Standardized Residual =	0.02	Largest Standardized Residual =	2.92	<p>The 'largest positive standardized residuals' suggest that an error covariance between READYTOL and TASKDISC may improve model fit. It is up to the researcher to decide if this is theoretically appropriate.</p> <p>Specific modification indices suggested by LISREL. These may or may not be theoretically plausible.</p>																																																									
Smallest Standardized Residual =	-2.51																																																															
Median Standardized Residual =	0.02																																																															
Largest Standardized Residual =	2.92																																																															
Stemleaf Plot																																																																
<table border="0"> <tbody> <tr><td>- 2 51</td></tr> <tr><td>- 1 75</td></tr> <tr><td>- 0 98544100000000</td></tr> <tr><td>0 123788</td></tr> <tr><td>1 00122355</td></tr> <tr><td>2 0019</td></tr> </tbody> </table>	- 2 51	- 1 75	- 0 98544100000000	0 123788	1 00122355	2 0019																																																										
- 2 51																																																																
- 1 75																																																																
- 0 98544100000000																																																																
0 123788																																																																
1 00122355																																																																
2 0019																																																																
Largest Positive Standardized Residuals	{	}																																																														
<table border="0"> <tbody> <tr><td>Residual for READYTOL and TASKDISC</td><td>2.925</td></tr> </tbody> </table>	Residual for READYTOL and TASKDISC	2.925																																																														
Residual for READYTOL and TASKDISC	2.925																																																															
The Modification Indices Suggest to Add an Error Covariance Between and Decrease in Chi-Square New Estimate	{	}																																																														
<table border="0"> <tbody> <tr><td>READYTOL</td><td>TASKDISC</td><td>8.5</td><td>0.16</td></tr> </tbody> </table>	READYTOL	TASKDISC	8.5	0.16																																																												
READYTOL	TASKDISC	8.5	0.16																																																													
Time used 0.062 seconds																																																																

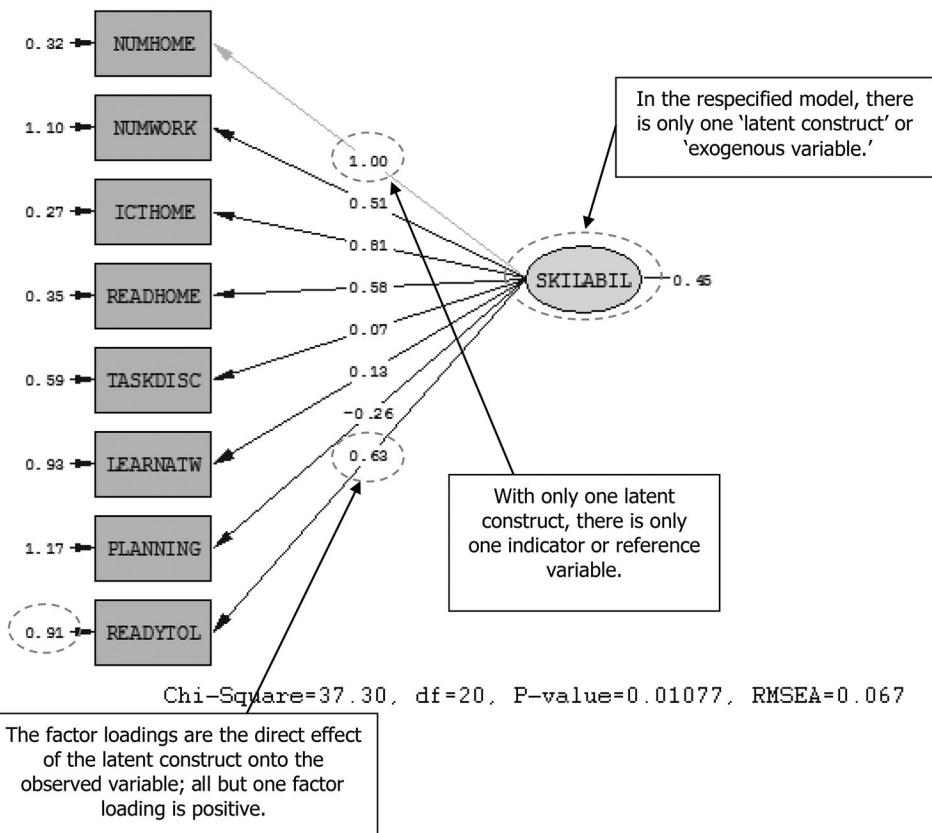


FIGURE 10.5

Respecified One-Factor Measurement Model (i.e., CFA) Example With No Correlated Errors

■ TABLE 10.4

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

L I S R E L 9.10 (STUDENT)

BY

Karl G. Jöreskog & Dag Sörbom

This program is published exclusively by
 Scientific Software International, Inc.
<http://www.ssicentral.com>

Copyright by Scientific Software International, Inc., 1981-2012
 Use of this program is subject to the terms specified in the
 Universal Copyright Convention.

The following lines were read from file G:\CH10_PIAAC_CFA.spj:

Ch. 10_PIAAC_CFA
 Raw Data from file 'G:\CH10_PIAAC_CFA.LSP'
 Observed Variables: NUMHOME NUMWORK ICTHOME READHOME
 TASKDISC LEARNATW PLANNING READYTOL ICTWORK INFLUENC
 READWORK WRITWORK WRITHOME CNTRYID SEQID

Latent Variables: SKILABIL

Relationships:

```

NUMHOME = 1*SKILABIL
NUMWORK = SKILABIL
ICTHOME = SKILABIL
READHOME = SKILABIL
READYTOL = SKILABIL
TASKDISC = SKILABIL
LEARNATW = SKILABIL
PLANNING = SKILABIL
  
```

Let error covariance between READYTOL and NUMHOME correlate
 Let error covariance between READYTOL and TASKDISC correlate
 Let error covariance between READYTOL and LEARNATW correlate
 Print residuals
 Path Diagram
 End of Problem

Sample Size = 191

Sample size

We are using raw data, however other input (e.g., correlation or covariance matrices) can be used

Observed variable names appear as they do in the column headers; LISREL truncates to 8 characters

Relationship commands define the observed to latent relations.
 Manifest variables are left of the '=' and latent variables are to the right of the '='

The last 3 lines tell LISREL to print residuals (we'll use for diagnostic purposes), print the path diagram, and know where the program stops.

Defines the correlated errors

The SIMPLIS code is reprinted in the output. We have respecified the model to be a one-factor solution with correlated errors and have renamed the factor or latent construct to be "Skills and Abilities." Since the latent variable is not observed, the researcher can name it what they choose as long as it is within the parameters of the program (e.g., 8 characters).

Ch. 10 PIAAC CFA

Covariance Matrix

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
NUMWORK	0.276	1.218				
ICTHOME	0.368	0.136	0.564			
READHOME	0.258	0.117	0.203	0.499		
TASKDISC	0.014	0.078	0.004	0.026	0.592	
LEARNATW	0.011	0.100	0.035	0.081	0.019	0.934
PLANNING	-0.186	0.019	-0.094	-0.034	0.112	0.135
READYTOL	0.219	0.193	0.261	0.208	0.171	0.170

Covariance Matrix

	PLANNING	READYTOL
PLANNING	1.198	
READYTOL	0.036	1.090

Variances and covariances of the observed variables. The saturated model equals the number of variances and covariances of the observed variables (in this case, 36).

Total Variance = 6.862 Generalized Variance = 0.0675

Largest Eigenvalue = 1.845 Smallest Eigenvalue = 0.267

Condition Number = 2.628

■ TABLE 10.4 (continued)

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

<p>Ch. 10 PIAAC CFA</p> <p>Number of Iterations = 6</p> <p>LISREL Estimates (Maximum Likelihood)</p> <p>Measurement Equations</p>	<p>ML estimates for the model are presented in this section.</p>
<p>NUMHOME = 1.000*SKILABIL, Errorvar. = 0.233 Standerr (0.0759) Z-values 3.066 P-values 0.002</p> <p>NUMWORK = 0.472*SKILABIL, Errorvar. = 1.092 Standerr (0.125) Z-values 3.762 P-values 0.000</p> <p>ICTHOME = 0.694*SKILABIL, Errorvar. = 0.307, R² = 0.456 Standerr (0.106) Z-values 6.532 P-values 0.000</p> <p>READHOME = 0.500*SKILABIL, Errorvar. = 0.365, R² = 0.268 Standerr (0.0876) Z-values 5.708 P-values 0.000</p> <p>TASKDISC = 0.0315*SKILABIL, Errorvar. = 0.592, R² = 0.0008 Standerr (0.0860) Z-values 0.366 P-values 0.714</p>	<p>Each observed variable is expressed as a linear function of the latent construct. NUMHOME was the indicator variable, and this is reflected in a parameter estimate of 1.0 with no standard error reported.</p> <p>R² = 0.697</p> <p>The unstandardized parameter estimate is akin to a regression coefficient. For NUMWORK, the magnitude of the unstandardized parameter estimate (.472) indicates the change in the latent construct for a one unit change in the observed variable. Changing NUMWORK by one unit results in about ½ unit increase in SKILABIL.</p> <p>R² = 0.0978</p> <p>Other than NUMHOME (indicator variable), all other loadings are freely estimated, and the p values indicate that 5 of the observed variables have statistically significant direct effects of the latent construct on the observed variable.</p> <p>Evaluate the individual parameter estimates to determine individual parameter fit to the model. Observed variables with statistically significant parameter estimates indicate they are related to the underlying latent construct. Non-statistically significant manifest variables may be considered for removal in the model modification process, however theory should guide this choice.</p>
<p>LEARNATW = 0.0683*SKILABIL, Errorvar. = 0.932, R² = 0.00267 Standerr (0.108) Z-values 0.631 P-values 0.528</p> <p>PLANNING = -0.259*SKILABIL, Errorvar. = 1.162, R² = 0.029 Standerr (0.120) Z-values -2.148 P-values 0.032</p> <p>READYTOL = 0.695*SKILABIL, Errorvar. = 0.825, R² = 0.238 Standerr (0.139) Z-values 5.001 P-values 0.000</p>	<p>The error variance is the measurement error. All are statistically significant, which is not of particular interest.</p> <p>The z-values are computed as the factor loading (i.e., regression coefficient) divided by the standard error. In this case, .0683/.108.</p> <p>R squared is the squared multiple correlation of the observed variable with the latent construct. Adhering to Cohen's (1988) guidelines for effect, .10 = small, .30 = moderate, and .5 = large for correlations thus for squared correlations, .01 = small, .09 = moderate, and .25 = large. The closer the R squared value is to +1.0, the better the observed variable is as an indicator of the latent construct.</p> <p>In this example, there are 4 large effects (NUMHOME, ICTHOME, READHOME, READYTOL), 1 moderate effect (NUMWORK), 1 small effect (PLANNING), and 2 negligible effects near zero (TASKDISC, LEARNATW).</p> <p>Negative coefficient indicates the latent construct has a negative effect on the observed variable.</p>

■ TABLE 10.4 (continued)

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

Error Covariance for READYTOL and NUMHOME = -0.148
 (0.0629)
 -2.350

Error Covariance for READYTOL and TASKDISC = 0.154
 (0.0561)
 2.754

Error Covariance for READYTOL and LEARNATW = 0.123
 (0.0688)
 1.794

Variances of Independent Variables

SKILABIL
 0.534
 (0.104)
 5.158

As seen on the path diagram, the variance of the latent construct

Saturated model =
 number of parameters to be estimated = number of variances and covariances of the observed variables

$$\left(\frac{k(k+1)}{2} \right) = \\ \left(\frac{8(8+1)}{2} \right) = 36$$

Log-likelihood Values

Estimated Model
 Number of free parameters(t)
 $\begin{cases} -2\ln(L) \\ AIC \text{ (Akaike, 1974)*} \\ BIC \text{ (Schwarz, 1978)*} \end{cases}$
 19

1031.756
 1069.756
 1131.549

Saturated Model

36
 1013.245
 1085.245
 1202.327

Comparative measures of fit (AIC, BIC)

*LISREL uses $AIC = 2t - 2\ln(L)$ and $BIC = t\ln(N) - 2\ln(L)$

'Free' parameters are the parameters actually estimated in the model. In this example, there are 7 coefficients, 8 measurement errors, 3 correlated errors, and 1 disturbance ($7 + 8 + 3 + 1 = 19$).

■ TABLE 10.4 (continued)

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

<p>Goodness of Fit Statistics</p> <p>Degrees of Freedom for (C1)-(C2) Maximum Likelihood Ratio Chi-Square (C1) Browne's (1984) ADF Chi-Square (C2_NT) Estimated Non-centrality Parameter (NCP) 90 Percent Confidence Interval for NCP</p> <p>Minimum Fit Function Value Population Discrepancy Function Value (F0) 90 Percent Confidence Interval for F0 Root Mean Square Error of Approximation (RMSEA) 90 Percent Confidence Interval for RMSEA P-Value for Test of Close Fit (RMSEA < 0.05)</p>	
	<p>Reject null hypothesis for the model chi-square (labeled 'ML ratio chi-square'), which indicates good fit. RMSEA and its CI suggest good fit; the lower bound is = 0, the upper bound is $\leq .10$. The CFI is = .99 which suggests good fit. The SRMR is = .05 which suggests good fit.</p>
	<p>17 18.511 ($P = 0.3573$) 19.502 ($P = 0.3005$) 1.511 (0.0 ; 16.216)</p> <p>0.0969 0.00791 (0.0 ; 0.0849) 0.0216 (0.0 ; 0.0707) 0.783</p>
	<p>We have a small RMSEA value (.02). The lower bound for the RMSEA confidence interval is 0, which is recommended, and the upper bound is less than .10, also recommended. The test of close fit is <i>not</i> statistically significant, indicating that RMSEA is statistically less than .05 and evidence of close fit.</p>
<p>Expected Cross-Validation Index (ECVI) 90 Percent Confidence Interval for ECVI ECVI for Saturated Model ECVI for Independence Model</p> <p>Chi-Square for Independence Model ([28 df])</p> <p>Normed Fit Index (NFI) Non-Normed Fit Index (NNFI) Parsimony Normed Fit Index (PNFI) Comparative Fit Index (CFI) Incremental Fit Index (IFI) Relative Fit Index (RFI)</p> <p>Critical N (CN)</p> <p>Root Mean Square Residual (RMR) Standardized RMR Goodness of Fit Index (GFI) Adjusted Goodness of Fit Index (AGFI) Parsimony Goodness of Fit Index (PGFI)</p> <p>Null model df where all observed variables are uncorrelated</p>	<p>0.296 (0.288 ; 0.373) 0.377 1.323</p> <p>236.775</p> <p>With 28 df the chi-square critical value is 41.372, thus the χ^2 for independence model is statistically significant. This is desirable as it tells us the variables are related (the null hypothesis being that they are unrelated).</p> <p>0.921 0.988 0.559 0.993 0.993 0.871</p> <p>343.939</p> <p>0.9482 0.0502, 0.975 0.947 0.460</p> <p>$\frac{(k(k-1)}{2} = \frac{8(7)}{2} = 28$</p>
	<p>Comparative measures of fit (ECVI)</p> <p>Incremental fit indices (NFI, NNFI, CFI, IFI) and parsimony-adjusted index (PNFI)</p> <p>Absolute fit indices (RMR, SRMR)</p> <p>Parsimony-adjusted index (PGFI)</p>

■ TABLE 10.4 (continued)

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

Ch. 10 PIAAC CFA

Fitted Covariance Matrix						
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.767					
	0.252	1.218				
	0.371	0.175	0.564			
	0.267	0.126	0.186	0.499		
	0.017	0.008	0.012	0.008	0.592	
	0.036	0.017	0.025	0.018	0.001	0.934
	-0.138	-0.065	-0.096	-0.069	-0.004	-0.009
READYTOL	0.224	0.175	0.258	0.186	0.166	0.149

Fitted Covariance Matrix						
	PLANNING	READYTOL				
PLANNING	1.198					
READYTOL	-0.096	1.083				

Fitted Residuals						
	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.000					
	0.024	0.000				
	-0.003	-0.039	0.000			
	-0.009	-0.009	0.018	0.000		
	-0.003	0.070	-0.007	0.018	0.000	
	-0.026	0.083	0.009	0.063	0.018	0.000
	-0.048	0.084	0.002	0.035	0.116	0.145
READYTOL	-0.004	0.018	0.003	0.022	0.004	0.021

Fitted Residuals						
	PLANNING	READYTOL				
PLANNING	0.000					
READYTOL	0.132	0.007				

Summary Statistics for Fitted Residuals						
Smallest Fitted Residual =	-0.048					
Median Fitted Residual =	0.004					
Largest Fitted Residual =	0.145					

Stemleaf Plot

```

- 4|8
- 2|96
- 0|997433000000
  0|234798888
  2|1245
  4|
  6|30

```

The fitted covariance matrix is the model-implied covariance matrix compared to the sample covariance matrix. The values are thus the estimated variances and covariances based on the model.

The fitted residuals represent the difference between the sample covariances and the fitted covariances. For example, NUMWORK fitted residual of .024 is simply .276 - .256. Negative fitted residuals indicate *overfitting* where the actual covariance has been overestimated in the model. Positive fitted residuals indicate *underfitting* where the actual covariance has been underestimated in the model.

Fitted residuals may be difficult to interpret as they are dependent on the metric of the variable.

■ TABLE 10.4 (continued)

LISREL Results for the Confirmatory Factor Analysis Example with Correlated Errors

8		34
10		6
12		2
14		5

	NUMHOME	NUMWORK	ICTHOME	READHOME	TASKDISC	LEARNATW
NUMHOME	0.000					
NUMWORK	0.948	0.000				
ICTHOME	-0.318	-1.121	0.000			
READHOME	-0.771	-0.218	1.321	0.000		
TASKDISC	-0.184	1.222	-0.272	0.552	0.000	
LEARNATW	-1.245	1.153	0.273	1.567	0.340	0.000
PLANNING	-1.746	1.049	0.042	0.783	1.941	1.929
READYTOL	-0.603	0.274	0.145	0.683	0.325	1.065

	Standardized Residuals	
	PLANNING	READYTOL
PLANNING	0.000	
READYTOL	1.944	0.605

Standardized residuals are: computed as the fitted residual divided by their standard error; represent the number of SD units the observed residual differs from zero (i.e., differs from the perfect model fit); and are interpreted as a conventional standardized values (i.e., mean = 0, SD = 1). Thus, conventions of 'large' standardized values can be applied (e.g., $\geq \pm 2.0$ or ± 3.0).

Summary Statistics for Standardized Residuals

{ Smallest Standardized Residual = -1.746 }	{ Median Standardized Residual = 0.209 }
{ Largest Standardized Residual = 1.944 }	

The standardized residuals may be helpful to review to determine modifications to the model. It is up to the researcher to decide what modifications are theoretically appropriate.

Stemleaf Plot

```

- 1|7
- 1|21
- 0|86
- 0|332200000000
  0|13333
  0|66789
  1|01223
  1|6999
  
```

If you are using the full version of LISREL (*i.e.*, not the student version), additional modification suggestions will be provided that look something like this:

$\left\{ \begin{array}{l} \text{Largest Positive Standardized Residuals} \\ \text{Residual for READYTOL and TASKDISC } 2.925 \end{array} \right\}$	<p>The 'largest positive standardized residuals' suggest that an error covariance between READYTOL and TASKDISC may improve model fit. It is up to the researcher to decide if this is theoretically appropriate.</p>
$\left\{ \begin{array}{lll} \text{The Modification Indices Suggest to Add an Error Covariance} & & \\ \text{Between } \text{READYTOL} \text{ and } \text{TASKDISC} & \text{Decrease in Chi-Square} & \text{New Estimate} \\ \text{READYTOL} & \text{TASKDISC} & 8.5 & 0.16 \end{array} \right\}$	<p>Specific modification indices suggested by LISREL. These may or may not be theoretically plausible.</p>

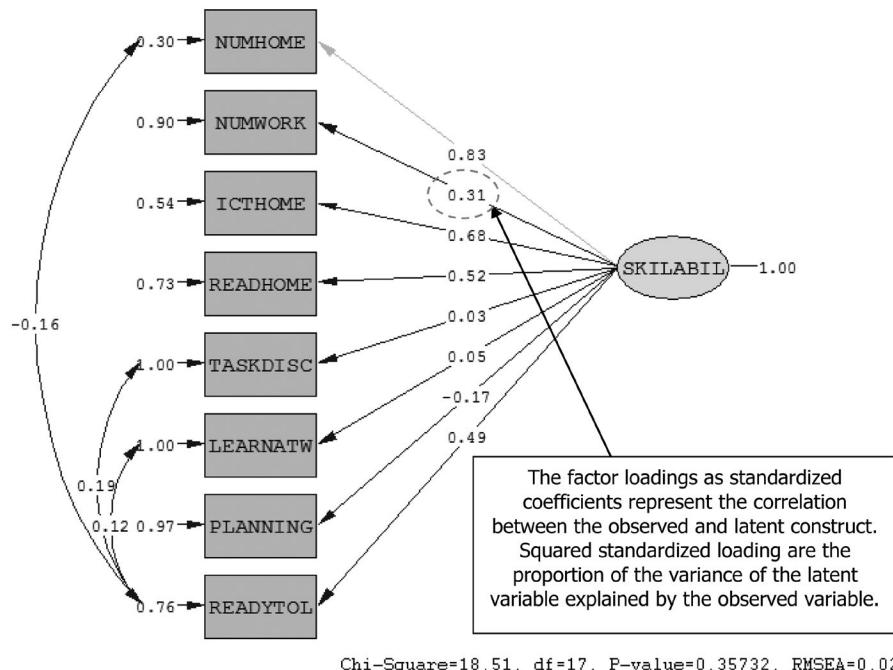
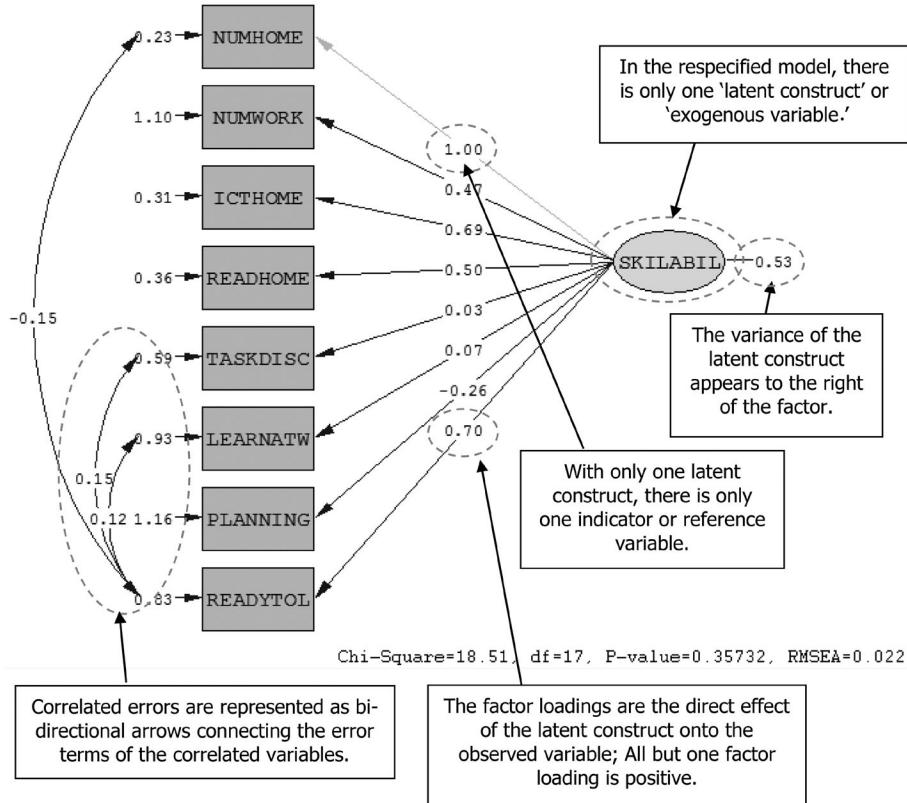


FIGURE 10.6

Respecified Measurement Model (i.e., CFA) Path Diagram Example With Correlated Errors: (a) Unstandardized and (b) Standardized Factor Loadings

10.3.2 CFA Using Mplus

For this illustration, Mplus version 7.4 was used. The code for computing the CFA is as follows. Note that for brevity, only the variables used in the CFA are included in this syntax (and the respective data file). Should additional variables be in the data file, the respective names would be listed in the ‘variable’ command.

```

TITLE:      PIAAC CFA
DATA:       FILE IS C:\Users\dhahs\Desktop\Ch10_PIAAC_CFA.
            dat;
VARIABLE:   NAMES ARE NUMHOME NUMWORK ICTHOME READHOME
            TASKDISC LEARNWK PLANNING RDYLEARN;
MODEL:      COGNIT BY NUMHOME NUMWORK ICTHOME READ-HOME;
            WORKABIL BY TASKDISC LEARNWK PLANNING RDY-LEARN;
OUTPUT:     standardized sampstat;

```

The path diagram generated from Mplus is provided in Figure 10.7.

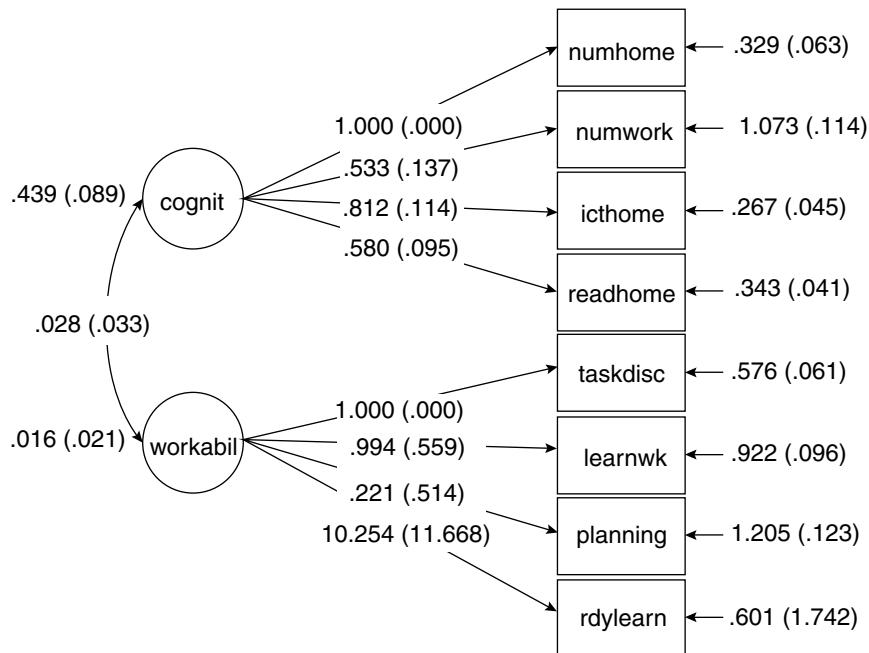


FIGURE 10.7

Initial Two-Factor Measurement Model Example Generated Using Mplus

10.4 DATA SCREENING

CFA shares assumptions that are common to EFA and regression. These include (a) independence, (b) linearity, (c) multivariate normality and absence of outliers, and (d) lack of multicollinearity and singularity. A data condition for conventional SEM is continuous data (i.e., interval or ratio). As these assumptions were screened in the EFA chapter, they will not be reiterated here. Readers are referred back to chapter 9 for data assumption screening details.

10.5 POWER

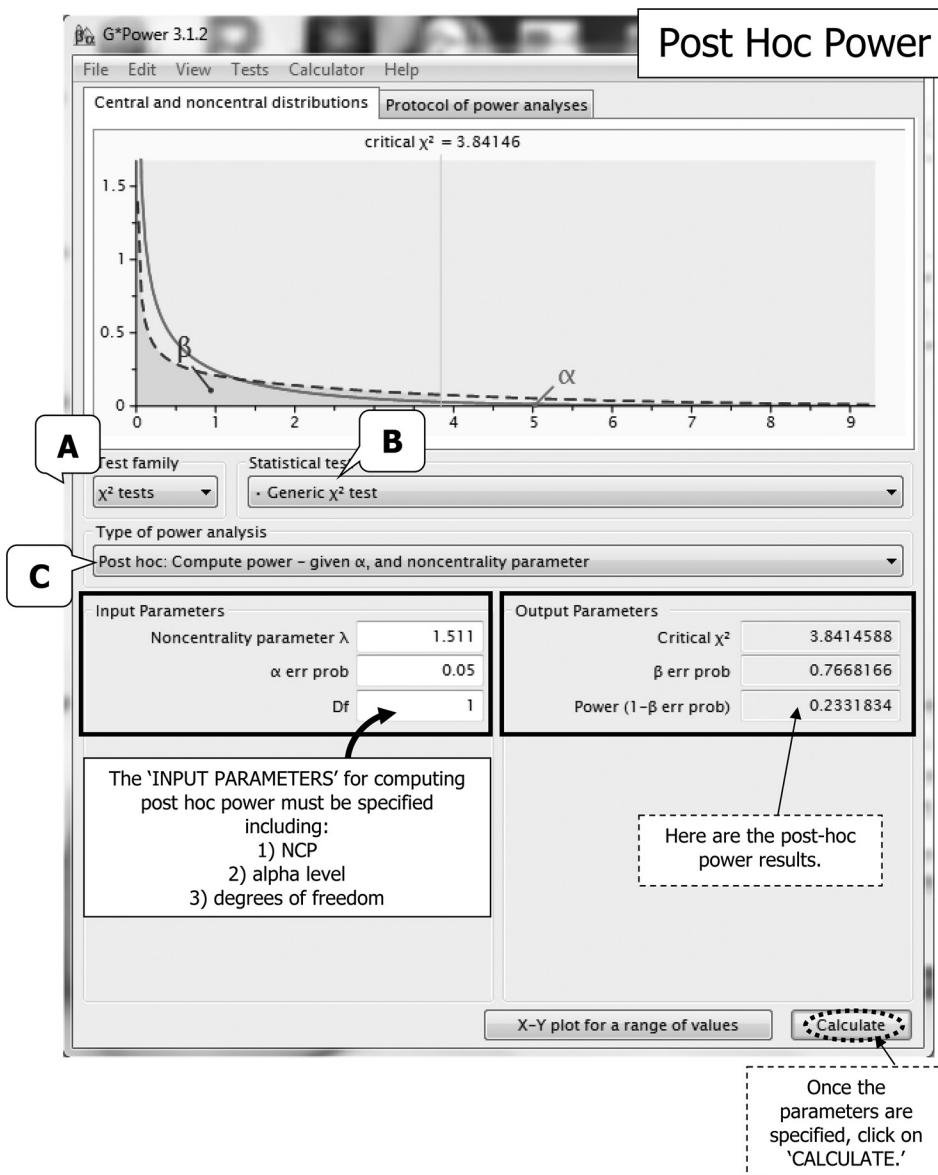
Determining power in models that fall under the SEM umbrella, like CFA, can be complicated (Schumacker & Lomax, 2010) as the models that are usually defined are just that—complicated. Syntax for generating power in standard statistical software is available (e.g., SAS, SPSS) (Schumacker & Lomax, 2010). Saris and Satorra (1993) provide a method for estimating power that can be easily calculated using G*Power and the noncentrality parameter (NCP) estimated from testing a theoretical model. The NCP is calculated as normal theory weighted least squares $\chi^2 - df_{model}$.

10.5.1 Post Hoc Power for CFA Using G*Power

The first thing that must be done when using G*Power for computing post hoc power is to select the correct test family. In our case, we conducted a CFA, and we are using the NCP to compute power. Thus, we need to change the 'Test Family' in the drop down menu to ' χ^2 tests' and then change the 'Statistical test' to 'Generic χ^2 test.' The 'Type of power analysis' desired then needs to be selected. To compute post hoc power, we need to select 'Post hoc: Compute achieved power—given α , and noncentrality parameter.'

The 'INPUT PARAMETERS' must then be specified (see Post Hoc Power screenshot). The noncentrality parameter was estimated when we generated the CFA. In the final model, the NCP was 1.511. In our example, the alpha level we used was .05, and the degrees of freedom for testing the hypothesis of model fit is 1 (which corresponds to a chi-square tabled critical value of 3.841 at an alpha of .5; the hypothesis being that the theoretical model fits the data). Once the parameters are specified, click on 'CALCULATE' to find the power statistics.

The 'OUTPUT PARAMETERS' provide the relevant statistics given the input just specified (see Post Hoc Power screenshot). In this example, we were interested in determining post hoc power for the CFA given the estimated NCP, alpha level of .05, and chi-square degrees of freedom. Based on those criteria, the post hoc power for the main effect of attractiveness was .23. In other words, the probability of rejecting the null hypothesis when it is really false was about 23%, which would be considered poor



power (sufficient power is often .80 or above). In this case, being underpowered may make it look as though our sample data does not differ from the hypothesized model, when it really may.

Although conducting power analysis a priori is recommended so that you avoid a situation where, post hoc, you find that the sample size was not sufficient to reach the desired level of power (given the observed parameters), G*Power does not currently offer the capability of computing a priori power for CFA through the generic χ^2 test family.

10.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example write-up for the results of the confirmatory factor analysis. Because of the complexity of conducting CFA and the numerous decisions that must be made in the CFA process, a number of guidelines have been published with recommendations on reporting CFA results so that both the methods and the results are thoroughly and transparently reported. The example write-up presented here includes elements as outlined in Brown (2006) and reprinted with permission in Table 10.5.

■ TABLE 10.5

Information to Report in a CFA Study

Model Specification

- Conceptual/empirical justification for the hypothesized model
- Complete description of the parameter specification of the model
 - List the indicators for each factor
 - Indicate how the metric of the factors was defined (e.g., specify which observed variables were used as marker indicators)
 - Describe all freely estimated, fixed, and constrained parameters (e.g., factor loadings and cross-loadings, random and correlated indicator errors, factor correlations, intercepts, and factor means)
- Demonstrate that the model is identified (e.g., positive model df , scaling of latent variables, absence of empirical underidentification)

Input Data

- Description of sample characteristics, sample size, and sampling method
- Description of the type of data used (e.g., nominal, interval; scale range of indicators)
- Tests of estimator assumptions (e.g., multivariate normality of input indicators)
- Extent and nature of missing data and the method of missing data management (e.g., direct ML, multiple imputation)
- Provide sample correlation matrix and indicator SDs (and means, if applicable), or make such data available upon request

Model Estimation

- Indicate the software and version used (e.g., LISREL 8.72)
- Indicate the type of data/matrices analyzed (e.g., variance-covariance, tetrachoric correlations/asymptotic covariances)
- Indicate the estimator used (e.g., ML, weighted least squares; as justified by properties of the input data)

Model Evaluation

- Overall goodness-of-fit
 - Report model χ^2 along with its df and p values
 - Report multiple fit indices (e.g., SRMR, RMSEA, CFI) and indicate cutoffs used (e.g., RMSEA < .06); provide confidence intervals, if applicable (e.g., RMSEA)
- Localized areas of ill fit
 - Report strategies used to assess for focal strains in the solution (e.g., modification indices/Lagrange multipliers, standardized residuals, Wald tests, EPC values)
 - Report absence of areas of ill fit (e.g., largest modification index) or indicate the areas of strain in the model (e.g., modification index, EPC value)

- If model is respecified, provide a compelling substantive rationale for the added or removed parameters and clearly document (improvement in) fit of the modified models
- Parameter estimates
 - Provide all parameter estimates (e.g., factor loadings, error variances, factor variances) including any nonsignificant estimates
 - Consider the clinical as well as the statistical significance of the parameter estimates (e.g., are all indicators meaningfully related to the factors?)
 - Ideally, include the standard errors or confidence intervals of the parameter estimates
- If necessary (e.g., suitability of N could be questioned), report steps taken to verify the power and precision of the model estimates (e.g., Monte Carlo evaluation using the model estimates as population values)

Substantive Conclusions

- Discuss CFA results in regard to their substantive implications, directions for future research, and so on.
- Interpret the findings in context of study limitations (e.g., range and properties of the indicators and sample) and other important considerations (e.g., equivalent CFA models)

Recall that Addie and Oso were assisting Dr. Wesley in exploring the factor structure of the PIAAC Survey of Adult Skills data using confirmatory factor analysis. The research question presented to Dr. Wesley was the following: How well does the two-factor model (Cognitive Skills and Work Abilities) explain the pattern of correlations among the measured indices?

Addie and Oso then assisted Dr. Wesley in conducting confirmatory factor analysis, and a template for writing the research question for CFA is presented below.

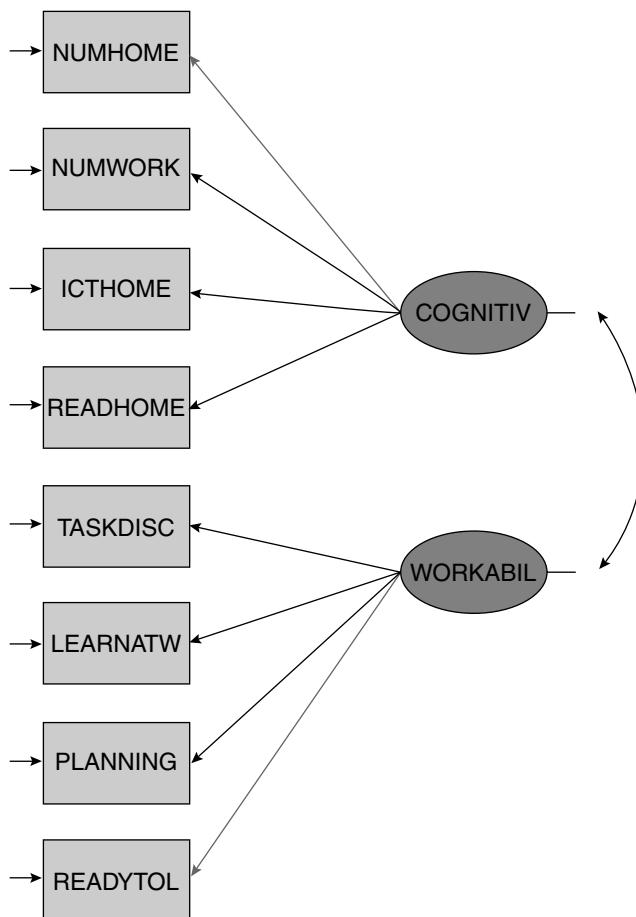
How well does the [number of latent constructs to be tested]-factor model [name(s) of latent constructs] explain the pattern of correlations among the measured variables?

It may be helpful to preface the results of the confirmatory factor analysis with information on an examination of the extent to which the data were thoroughly screened.

Prior to conducting the confirmatory factor analysis (CFA), the data were screened to determine the extent to which the assumptions associated with CFA were met. These assumptions included (a) independence, (b) linearity, (c) multivariate normality, and (d) lack of multicollinearity and singularity. [See chapter 9 for example write-up of data screening.]

Here is an example summary of results for confirmatory factor analysis (remember that this will be prefaced by a section reporting the extent to which the data were thoroughly screened).

A confirmatory factor analysis was performed to assess the factor structure of the PIAAC cognitive and work abilities indexes. A path diagram of the initial two-factor model tested is presented in Figure A.

**FIGURE A**

Initial Two-Factor Model

The following observed variables were related to the latent construct 'cognitive skills': numeracy skills at home, numeracy skills at work, ICT skills at home, reading at home. The following observed variables were related to the latent construct 'work abilities': ready to learn, task discretion, learning at work, planning. 'Numeracy skills at home' was used as the indicator variable to scale the latent construct 'cognitive skills,' and 'ready to learn' was used as the indicator variable to scale the latent construct 'work abilities.' All other manifest variables were freely estimated and loaded on only one latent variable. The latent constructs were allowed to correlate. All manifest variables were continuous. The sample correlation matrix is presented in Table A (the covariance matrix of the sample data was the input for LISREL). Means and standard deviations are provided in Table B. LISREL 9.10 student version with maximum likelihood estimation was used to compute the models.

Model fit for the *two-factor* model was less than desirable, as a negative error variance was generated. Thus, a *one-factor* model was tested that remedied the negative error variance; however, model fit was poor. The model χ^2 was 37.300 ($p = .011$), indicating poor fit between the model implied and

■ TABLE ACorrelation Matrix for Cognitive and Work Ability Indices ($N = 191$)

Item	1	2	3	4	5	6	7
1. Index of use of numeracy skills at home (basic and advanced)	—						
2. Index of use of numeracy skills at work (basic and advanced)	.292	—					
3. Index of use of ICT skills at home (derived)	.554	.176	—				
4. Index of use of reading skills at home (prose and document texts)	.418	.166	.374	—			
5. Index of use of task discretion at work	.015	.086	.014	.049	—		
6. Index of learning at work	.001	.098	.045	.107	.037	—	
7. Index of use of planning skills at work	-.202	.013	-.111	-.048	.140	.139	—
8. Index of readiness to learn	.246	.161	.341	.296	.206	.167	.026

■ TABLE BMeans (Standard Deviations) and Minimum and Maximum Values for Cognitive and Work Ability Indices ($N = 191$)

Item	Mean (<i>SD</i>)	Range
1. Index of use of numeracy skills at home (basic and advanced)	2.52 (.88)	.31, 6.91
2. Index of use of numeracy skills at work (basic and advanced)	2.48 (1.10)	.06, 6.73
3. Index of use of ICT skills at home (derived)	2.63 (.75)	1.04, 6.48
4. Index of use of reading skills at home (prose and document texts)	2.74 (.70)	1.31, 5.96
5. Index of use of task discretion at work	1.89 (.77)	.48, 4.42
6. Index of learning at work	2.57 (.97)	.36, 4.35
7. Index of use of planning skills at work	2.14 (1.10)	.30, 3.82
8. Index of readiness to learn	2.83 (1.04)	1.18, 5.00

sample covariance matrices. Goodness-of-fit indices were mixed with some suggesting good model fit and others an ill-fitting model: (a) RMSEA = .067, RMSEA 90% CI = (.032, .100), which does not meet the lower threshold of zero nor the upper threshold of less than .10; (b) CFit test $p = .183$, not statistically significant indicating that RMSEA is statistically significantly greater than the threshold of .05; and (c) SRMR = .065, less than the recommended .08 suggesting good model fit; CFI = .917, less than the recommended .95 also suggesting good (but not great) model fit. Standardized residuals suggested freeing the covariance between 'ready to learn' and other manifest variables would improve model fit. Thus, error terms between the following pairs of variables were allowed to correlate: (a) 'ready to learn' and 'numeracy skills at home,' (b) 'ready to learn' and 'task discretion,' and (c) 'ready to learn' and 'learn at work.' Comparative model fit improved from the two-factor model with no correlated errors, AIC = 1082.546, BIC = 1134.582, ECVI = .363, ECVI 90% CI = (.293, .474); to the modified two-factor model with correlated errors, AIC = 1069.756, BIC = 1131.549, ECVI = .296, ECVI 90% CI = (.288, .373).

For the one-factor model with correlated errors, the model χ^2 was 18.511 ($p = .357$), indicating good fit between the model implied and sample

covariance matrices. Goodness-of-fit indices suggested good model fit: (a) RMSEA = .022, RMSEA 90% CI = (.0, .071), which meets the lower threshold of zero and the upper threshold of less than .10; (b) CFit test $p = .783$, not statistically significant indicating that RMSEA is statistically significantly less than .05; and (c) SRMR = .05, lower than the recommended .08 indicating good fit; CFI = .99, greater than the recommended .95, also indicating good fit.

Standardized residuals were reviewed to determine localized areas of poor fit. The largest standardized residual was less than 2.0, and thus further modifications were not considered. Parameter estimates are reported in Table C. Five of the observed variables were statistically significantly related to the latent construct. 'Task discretion' and 'learning at work' were neither statistically (based on p values) nor clinically meaningfully (based on R squared) related to the factor. Power was calculated using G*Power version 3.1.2 and was computed to be .23.

■ TABLE C

Parameter Estimates

	Coefficient	SE	<i>p</i>	Error Variance	R^2
Numeracy skills at home	1.000			.233	.697
Numeracy skills at work	.472	.125	.000	1.099	.098
ICT skills at home	.694	.106	.000	.307	.456
Reading at home	.500	.088	.000	.365	.268
Task discretion	.315	.086	.714	.592	.001
Learning at work	.068	.108	.528	.932	.003
Planning	-.259	.120	.032	1.162	.029
Ready to learn	.695	.139	.000	.825	.238
Factor variance = .534 (SE = .104)					
Error covariance for ready to learn and numeracy skills at home: -.148 (SE = .063)					
Error covariance for ready to learn and task discretion: .154 (SE = .056)					
Error covariance for ready to learn and learning at work: .123 (SE = .069)					

The results suggest a one-factor model with seven manifest variables. [Substantive implications and directions for future research presented here.]

PROBLEMS

Conceptual Problems

- Which one of the following is specified as part of the measurement model in confirmatory factor analysis?
 - Direct effects between latent variables
 - Indirect effects between latent constructs
 - Relationships between endogenous and exogenous variables
 - Relationships between observed variables and latent constructs

2. Latent constructs are represented in path diagrams as which one of the following?
 - a. Curved bidirectional arrows
 - b. Ovals
 - c. Rectangles
 - d. Unidirectional arrows
3. Observed variables are represented in path diagrams as which one of the following?
 - a. Curved bidirectional arrows
 - b. Ovals
 - c. Rectangles
 - d. Unidirectional arrows
4. Correlations or covariances are depicted in path diagrams as which one of the following?
 - a. Curved bidirectional arrows
 - b. Ovals
 - c. Rectangles
 - d. Unidirectional arrows
5. Direct effects of latent constructs on observed variables are depicted in path diagrams as which one of the following?
 - a. Curved bidirectional arrows
 - b. Ovals
 - c. Rectangles
 - d. Unidirectional arrows
6. Specifying a reference or marker indicator does which one of the following?
 - a. Defines a covariance/correlation between variables
 - b. Fixes the metric of the latent construct to the same as the indicator
 - c. Identifies areas of localized poor fit
 - d. Standardizes the model
7. Which one of the following does NOT represent a goodness-of-fit index that might be reported with confirmatory factor analysis?
 - a. RMSEA
 - b. SRMR
 - c. WLS
 - d. χ^2
8. Which one of the following is a goodness-of-fit index for model parsimony?
 - a. RMSEA
 - b. SRMR
 - c. WLS
 - d. χ^2
9. Failing to specify the correlated errors means which one of the following?
 - a. All measurement error is random
 - b. All the covariations among the observed variables loading on a given latent construct are due to the observed variables
 - c. Some shared variance is external to the model
 - d. None of the above. Correlated errors cannot be specified in CFA.

10. Which one of the following represents the interpretation of a completely standardized factor loading?
 - a. Proportion of variance of the observed variable that is explained by the latent construct
 - b. Square root of the factor variance
 - c. Standardized regression coefficient
 - d. Unstandardized regression coefficient

Computational Problems

1. Using the CH10_HW1_2_PIAAC_NORWAY.sav dataset, conduct confirmatory factor analysis following the steps in this chapter, using maximum likelihood estimation and testing a *one-factor* model. Set ‘ICT skills at home’ as the indicator variable. Review the standardized factor loadings. Interpret goodness-of-fit using overall chi-square, RMSEA, SRMR, and CFI. Suggest modification indices to improve model fit, if needed. (*Note: This data has been delimited to individuals who indicated their highest level of school was ‘above high school’ [B_Q01a_T = 3] and who were employed the year prior to completing the survey [B_Q15a = 1]. All cases with missing data on one or more index were removed from this dataset.*)
2. Using the CH10_HW1_2_PIAAC_NORWAY.sav dataset, conduct confirmatory factor analysis following the steps in this chapter, using maximum likelihood estimation and testing a *two-factor* model. Influenced by the first latent construct are the following observed variables: ICTHOME, WRITHOME, NUMHOME, READHOME, READYTOL. Influenced by the second latent construct are the following observed variables: READWORK, WRITWORK, LEARNATW, PLANNING, ICTWORK. Set ‘ICTHOME’ as the indicator variable for factor 1, and ‘READWORK’ as the indicator variable for factor 2. Review the standardized factor loadings. Interpret goodness-of-fit using overall chi-square, RMSEA, SRMR, and CFI. Suggest modification indices to improve model fit, if needed. (*Note: This data has been delimited to individuals who indicated their highest level of school was ‘above high school’ [B_Q01a_T = 3] and who were employed the year prior to completing the survey [B_Q15a = 1]. All cases with missing data on one or more index were removed from this dataset.*)

Interpretive Problem

1. Use SPSS to conduct confirmatory factor analysis with the continuous PIAAC index variables from Italy (CH9_HW_INTERPRET_ITALY.sav). The data file has been delimited to include only individuals who reported having ‘above high school’ education [B_Q01a_T = 3] and who had complete data on the index variables. (*Note: If you’re using the student version of LISREL, you will need to remove the variables not included in the model, as there are more variables in the data file than the student version of LISREL will allow.*) Write up the results.

REFERENCES

- Bentler, P.M. (1990). Comparative fit indices in structural models. *Psychological Bulletin, 107*, 238–246.
- Bentler, P.M., & Chou, C. H. (1987). Practical issues in structural modeling. *Sociological Methods and Research, 16*, 78–117.
- Bollen, K. A., & Pearl, J. (2013). Eight myths about causality and structural equation modeling. In S. L. Morgan (Ed.), *Handbook of causal analysis for social research* (pp. 301–328). New York, NY: Springer.
- Boomsma, A. (1983). *On the robustness of LISREL against small sample size and nonnormality*. Amsterdam: Sociometric Research Foundation.
- Boomsma, A., & Hoogland, J. J. (2001). The robustness of LISREL modeling revisited. In R. Cudeck, S. du Toit, & D. Sorbom (Eds.), *Structural equation modeling: Present and future*. Chicago, IL: Scientific Software International.
- Brown, T. A. (2006). *Confirmatory factor analysis for applied research*. New York: Guilford Press.
- Byrne, B. M. (1998). *Structural equation modeling with LISREL, PRELIS, and SIMPLIS: Basic concepts, applications, and programming*. Mahwah, NJ: L. Erlbaum Associates.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, NJ: Lawrence Erlbaum.
- Cudeck, R., & Browne, M. W. (1983). Cross-validation of covariance structures. *Multivariate Behavioral Research, 18*, 147–167.
- Curran, P. J., West, S. G., & Finch, J. F. (1996). The robustness of test statistics to nonnormality and specification error in confirmatory factor analysis. *Psychological Methods, 1*, 16–29.
- Duncan, T. E., & Duncan, S. C. (2009). The ABC's of LGM: An introductory guide to latent variable growth curve modeling. *Social and Personality Psychology Compass, 3*(6), 979–991.
- Gagne, P., & Hancock, G. R. (2006). Measurement model quality, sample size, and solution propriety in confirmatory factor models. *Multivariate Behavioral Research, 41*(1), 65–83.
- Hancock, G. R., & Mueller, R. O. (2006). Introduction. In G. R. Hancock & R. O. Mueller (Eds.), *Structural equation modeling: A second course* (pp. 1–13). Greenwich, CT: Information Age Publishing.
- Hu, L., & Bentler, P.M. (1999). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural Equation Modeling, 6*(1), 1–55.
- Jackson, D. L., Voh, J., & Frey, M. (2013). A note on sample size and solution propriety for confirmatory factor analytic models. *Structural Equation Modeling, 20*, 86–97.
- Kaplan, D. (2000). *Structural equation modeling: Foundations and extensions*. Thousand Oaks, CA: Sage Publications.
- Kline, R. B. (2011). *Principles and practice of structural equation modeling* (3rd ed.). New York, NY: Guilford.
- Kroonenberg, P.M., & Lewis, C. (1982). Methodological issues in the search of a factor model: Exploration through confirmation. *Journal of Educational Statistics, 7*, 69–89.
- MacCallum, R., & Austin, J. T. (2000). Applications of structural equation modeling in psychological research. *Annual Review of Psychology, 51*, 201–226.
- MacCallum, R. C., Browne, M. W., & Sugawara, H. M. (1996). Power analysis and determination of sample size for covariance structure modeling. *Psychological Methods, 1*(2), 130–149.

- Marcoulides, G. A., & Schumacker, R. E. (2001). *New developments and techniques in structural equation modeling*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Marsh, H. W., Hau, K. T., Balla, J. R., & Grayson, D. (1998). Is more ever too much? The number of indicators per factor in confirmatory factor analysis. *Multivariate Behavioral Research*, 33, 181–220.
- Marsh, H. W., Hau, K. T., & Grayson, D. (2005). Goodness of fit in structural equation models. In A. Maydeu-Olivares & J. J. McArdle (Eds.), *Contemporary psychometrics* (pp. 275–340). Mahwah, NJ: Lawrence Erlbaum.
- Mulaik, S. A., & Millsap, R. E. (2000). Doing the four-step right. *Structural Equation Modeling*, 7, 36–73.
- Myers, R. H. (1990). *Classical and modern regression with applications* (2nd ed.). Boston, MA: Duxbury.
- Nunnally, J. C. (1967). *Psychometric theory*. New York, NY: McGraw-Hill.
- Saris, W. E., & Satorra, A. (1993). Power evaluation in structural equation modeling. In K. A. Bollen & J. S. Long (Eds.), *Testing structural equation models* (pp. 181–204). Newbury Park, CA: Sage.
- Savalei, V. (2012). The relationship between root mean square error of approximation and model misspecification in confirmatory factor analysis models. *Educational and Psychological Measurement*, 72(6), 910–932.
- Schumacker, R. E., & Lomax, R. G. (2010). *A beginner's guide to structural equation modeling* (3rd ed.). New York: Routledge.
- Schumacker, R. E., & Lomax, R. G. (2016). *A beginner's guide to structural equation modeling*. New York, NY: Routledge Taylor & Francis.
- Sharma, S., Mukherjee, S., Kumar, A., & Dillon, W. R. (2005). A simulation study to investigate the use of cutoff values for assessing model fit in covariance structure models. *Journal of Business Research*, 58, 935–943.
- Stevens, J. P. (2009). *Applied multivariate statistics for the social sciences* (5th ed.). New York: Psychology Press.
- Tucker, L. R., & Lewis, C. (1973). A reliability coefficient for maximum likelihood factor analysis. *Psychometrika*, 38, 1–10.
- Wetherill, G. B. (1986). *Regression analysis with applications*. London: Chapman & Hall.
- Wolf, E. J., Harrington, K. M., Clark, S. L., & Miller, M. W. (2013). Sample size requirements for structural equation models: An evaluation of power, bias, and solution propriety. *Educational and Psychological Measurement*, 73(6), 913–934.

Chapter 11

MULTILEVEL LINEAR MODELING

CHAPTER OUTLINE

11.1 What Multilevel Linear Modeling Is and How It Works	507
11.1.1 Characteristics	508
11.1.2 Sample Size	521
11.1.3 Power	522
11.1.4 Effect Size	523
11.1.5 Assumptions	525
11.1.6 Conditions	529
11.2 Mathematical Introduction Snapshot	530
11.3 Computing Multilevel Modeling Using HLM	531
11.3.1 Random Intercepts Model	543
11.3.2 Random Intercepts Random Slopes Model	545
11.3.3 Random Intercepts Model With Level 2 Predictor	548
11.3.4 Random Intercepts Model With Nonrandomly Varying Slopes Effects	550
11.4 Data Screening	553
11.4.1 Level 1 Residuals	553
11.4.2 Level 2 Residuals	553
11.4.3 Graphing	556
11.4.4 Model Fit	558
11.5 Power Using Optimal Design	559
11.6 Research Question Template and Example Write-Up	561

KEY CONCEPTS

1. Between-groups
2. Centering
3. Clustering
4. Estimation method
5. Fixed effects
6. Intraclass correlation coefficient
7. Model fit
8. Nonrandomly varying effects
9. Random effects
10. Within-groups

Multilevel modeling, which you have likely been introduced to through hierarchical ANOVA, provides the opportunity to examine clustered or nested data. The study of individuals within organizations, such as students within classrooms and schools, employees within companies, patients within hospitals, and residents within neighborhoods, are all common examples of clustered data. In many situations, researchers are interested not only in modeling factors at the individual level but also at the group level (e.g., classroom, organization) because context—at least in education and the social and health sciences, as well as many other disciplines—is important. These higher-level factors are contextual effects that may (and usually do) contribute to or influence the individual outcomes being studied. All of us have been involved in many different types of group settings. It is probably easy for you to imagine, and you have likely experienced this for yourself, how the context of the group (e.g., the personalities of the other individuals in your group, the resources available to you, the physical environment) can greatly influence your experience in that setting. For example, the school level percentage of students on free and reduced lunch may help explain the outcomes of students within schools and the median housing price may help explain outcomes of residents within neighborhoods. Even if the interest is not in contextual effects, the hierarchical or clustered nature of the data must still be accounted for in the modeling, as nested units are more likely to produce similar outcomes as compared to units that are randomly sampled. Ignoring the clustering (assuming independence of errors) most often produces estimates with incorrect and decreased standard errors, because it is usually the case that clustered data violates the assumption of independent and identically distributed errors. Think back to the very simple one-sample t test, probably the most elementary inferential procedure you know. The denominator of the t test statistic formula is the standard error. Holding all else constant, as the standard error decreases, the test statistic increases. As the test statistic increases, the likelihood of rejecting the null hypothesis increases. This same concept applies with other, more advanced statistical procedures. The point is that ignoring clustering and hierarchical nature of data will usually result in the increased probability of a Type I error—rejecting the null hypothesis when it is really true. This translates to interpretations suggesting evidence of real effects/differences/relationships, when in fact the results are simply an artifact

of failing to account for the clustered design of the data. This is a good time to bring up the concept of *ecological fallacy*, interpreting relationships that are observed in groups to hold also for individual units (e.g., Freedman, 1999), and *atomistic fallacy*, interpreting relationships that are observed in individual units to hold also for groups (Hox, 2002). These fallacies are “a problem of *inference*, not of *measurement*” (Luke, 2004, p. 6, emphasis original), and can be combated with multilevel modeling.

The purpose of this chapter is to introduce multilevel modeling through a linear perspective focusing on a two-level model, but the hope is that it will whet your appetite to learn other applications of multilevel modeling, as the options available for conducting various types of multilevel models (e.g., linear and nonlinear regression, structural equation) have grown exponentially in a very short period of time. There are a multitude of excellent resources available for learning more about multilevel modeling, a few of which include the Centre for Multilevel Modeling (<http://www.bristol.ac.uk/cmm/>; multilevel software MLwiN is produced through the Center) and UCLA’s Academic Technology Services (<http://www.ats.ucla.edu/stat/>), through which there are links to many of the major statistical software programs (*R*, SAS, SPSS, Stata) and multilevel capabilities of each, respectively. Our objectives are that, by the end of this chapter, you will be able to (a) understand the basic concepts of multilevel linear modeling, (b) determine and interpret the results of multilevel linear modeling, and (c) understand and evaluate the assumptions of multilevel linear modeling.

11.1 WHAT MULTILEVEL LINEAR MODELING IS AND HOW IT WORKS

It may be nearing the end of the semester, but the graduate research assistants in the stats lab have their hands full. Today, we find them working together on a most exciting project.

Challie Lenge, Ott Lier, Oso Wyse, and Addie Venture are nearing the end of their graduate studies, and thus their graduate research positions as well. The excitement in the lab, however, continues. The lab’s faculty supervisor and mentor has asked the group to consult with Dr. Maeve, a psychology faculty member who is examining data on the Big Five personality factors (Goldberg, 1992). Dr. Maeve is using public secondary data from individuals in 19 North American countries who completed online personality tests anonymously. The research questions that Dr. Maeve is interested in examining include (1) *Is there variation in conscientiousness between countries?* (2) *Can conscientiousness be predicted by intellect/imagination?* and (3) *To what extent does the developing status of the country moderate the relationship between intellect and conscientiousness?* After speaking with Dr. Maeve and previewing the data, a two-level multilevel model is recommended by the graduate student team.

Multilevel linear models—also known as mixed models, mixed-effects models (Pinheiro & Bates, 2000), hierarchical models (Raudenbush & Bryk, 2002), random

coefficients models (Longford, 1993), variance components models, and intercepts and slopes as outcomes—differ from ordinary least squares regression because they contain not only *fixed effects* but also *random effects*. The random effects allow modeling of the variation between clusters or groups. Conceptually, multilevel linear models are quite similar to OLS regression. They have as their goal the prediction of some continuous dependent variable based on a function of one or more independent variables. The twist with multilevel models as compared to OLS regression is that this prediction model operates on two or more levels. As mentioned previously, this chapter considers a two-level linear model.

11.1.1 Characteristics

11.1.1.1 Fixed, Random, and Nonrandomly Varying Effects

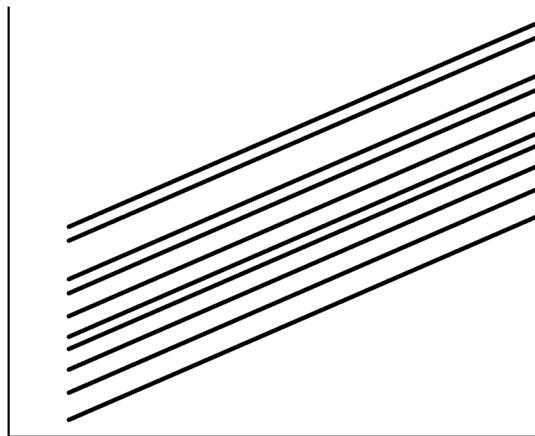
It is important to understand the distinction between fixed, random, and nonrandomly varying effects. *Fixed effects* are parameters that do not vary across clusters or groups. In other words, these parameters consider only one value for the entire sample, and thus the effects are assumed constant across groups. Traditionally, single-level analyses focus on fixed effects, but multilevel models also include fixed effects. In MLM, intercepts and slopes can be fixed, although fixing both the intercept and slope does not make sense in MLM, as a fixed intercept and fixed slope results in a straight line (i.e., there becomes no variability in the intercept or slope and thus every case has the same parameters, resulting in a single regression line which is estimable by OLS regression). In terms of notation, fixed effects for the slope are denoted by the lack of a residual (or random) component in the level 2 equations:

$$\text{Level 1: } Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + r_{ij}$$

$$\begin{aligned}\text{Level 2: } \beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10}\end{aligned}$$

This level 2 equation, with random intercepts (denoted by u_{0j}) and fixed slopes, graphically appears in Figure 11.1. The intercepts, the average outcomes between groups, vary between groups, and this is suggested by lines that would cross the Y axis at different values; however, the slopes are fixed. The fixed slopes are suggested by parallel lines—the rate of change in Y is the same across all groups.

In comparison, *random effects* are parameters that *do* vary across clusters. In a multilevel model, both the intercept and slope (i.e., β_{0j} and β_{1j}) can also be random effects (also known as random coefficients) and can be allowed to vary or differ across the clusters or groups (i.e., level 2 or the higher order groups). In other words, effects that are random vary across clusters. The idea behind modeling the random effects is that units or cases within a particular group or cluster tend to be more similar (i.e., homogenous) to other units within that cluster than to units or

**FIGURE 11.1**

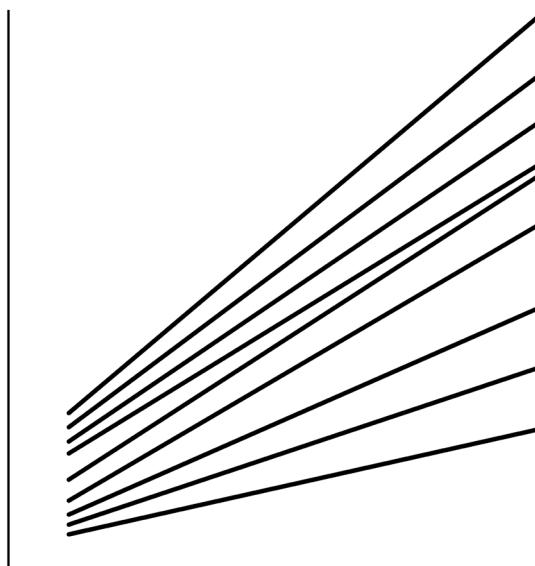
Random Intercepts, Fixed Slopes

cases in different groups. This creates a dependency between units that statistically needs to be teased out in order for the resulting parameter estimates to be interpreted without bias. Even mild violations of independence can result in substantial Type I error problems, usually increased probabilities of making a Type I error (i.e., rejecting the null hypothesis when it is really true) (Kish, 1965). In other words, the assumption of independence (i.e., the assumption of independent and identically distributed—*iid*—variables, which says that units can be considered as randomly selected from a probability distribution) is violated when this homogeneity is not modeled. This dependency can be estimated with the intraclass correlation coefficient (ICC), which we will consider next. In terms of notation, random intercepts and random slopes are reflected by the inclusion of the residual terms, respectively (u_{0j} and u_{1j}) in the level 2 equations:

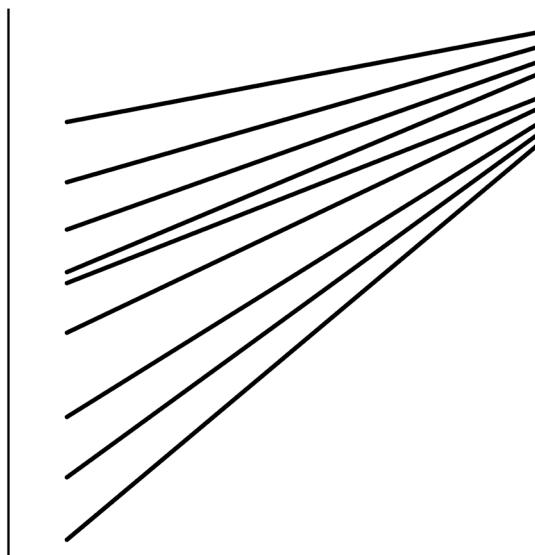
$$\text{Level 1: } Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + r_{ij}$$

$$\begin{aligned}\text{Level 2: } \beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10} + u_{1j}\end{aligned}$$

When there are both random intercepts and random slopes, a relationship may or may not exist between the intercepts and slopes. A positive relationship between intercepts and slopes will result in a graph that follows a pattern similar to Figure 11.2. In comparison, a negative relationship between intercepts and slopes will result in a graph that follows a pattern similar to Figure 11.3. When no relationship exists between the intercepts and slopes, a graph of the individual groups will display the regression lines without an evident pattern—simply random lines—similar to Figure 11.4.

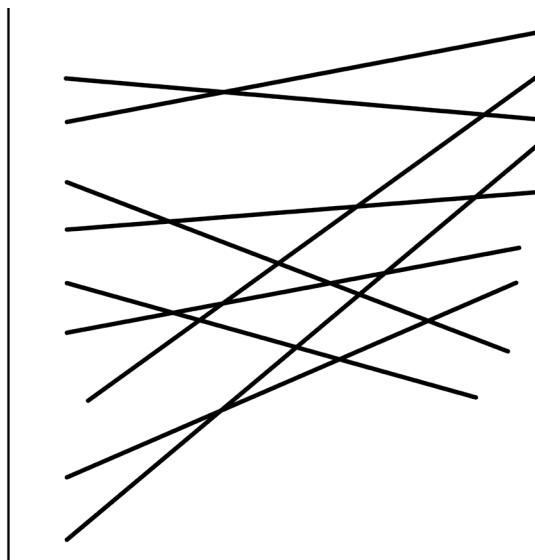
**FIGURE 11.2**

Random Intercepts, Random Slopes; Positive Relationship Between Intercepts and Slopes

**FIGURE 11.3**

Random Intercepts, Random Slopes; Negative Relationship Between Intercepts and Slopes

Nonrandomly varying effects indicate that a coefficient varies between groups, but not randomly. This variation can be modeled by predictors in the model. In these equations, the inclusion of a level 2 predictor, W_j , models these parameters as varying, but they vary as a function of W_j , not randomly.

**FIGURE 11.4**

Random Intercepts, Random Slopes; No Relationship Between Intercepts and Slopes

$$\text{Level 1: } Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + r_{ij}$$

$$\begin{aligned}\text{Level 2: } \beta_{0j} &= \gamma_{00} + \gamma_{01}W_j \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}W_j\end{aligned}$$

An equation that includes both a nonrandomly varying and a randomly varying component is considered a random effect, even though it also includes a nonrandomly varying component. *Let's consider the level 2 equations that follow.* The intercept, β_{0j} , is randomly varying, and this is reflected by u_{0j} . The slope, β_{1j} , is also random. Even though this slope has a nonrandomly varying component, specifically $\gamma_{11}W_j$, the random slope effect of u_{1j} makes this a random effect. The slope, β_{2j} , is a nonrandomly varying component, and this is reflected by the inclusion of a level 2 predictor, $\gamma_{21}W_j$, and no random effect. This slope, therefore, varies between groups as a function of covariate W_j , but this variation is not random—it can be explained by W_j . The slope, β_{3j} , is a fixed effect, and this is reflected by the absence of a residual in the equation. We interpret the fixed slope to mean that these particular slopes are the same across all schools (i.e., these slopes do not vary between schools).

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}W_j + u_{1j} \\ \beta_{2j} &= \gamma_{20} + \gamma_{21}W_j \\ \beta_{3j} &= \gamma_{30}\end{aligned}$$

11.1.1.2 Intraclass Correlation Coefficient

In the case of nested data, the correlation between variables within the same group or cluster will be more strongly correlated than correlations between variables for cases in different groups. Multilevel modeling is the ideal solution to approaching situations where dependency is an issue, as it allows the simultaneous estimation of effects at both the unit or case level (e.g., person level, level 1) and the group or cluster level (level 2). The ICC reflects the proportion of variance in the dependent variable that is between groups (e.g., level 2 units). The ICC will be greater than zero only in situations where the intercept variance (τ_{00}) is greater than one. In this respect, a test of the ICC is equivalent to a test of the intercept variance (Snijders & Bosker, 2012). Thinking about a basic level 2 model, the ICC is also interpreted as the variance over the level 2 units. The ICC is applicable to random intercept models only.

The intraclass correlation coefficient, the proportion of variance in the outcome that is between groups, can be computed for a two-level model as follows:

$$\rho = \frac{\tau_{00}}{(\sigma^2 + \tau_{00})}$$

Where σ^2 = level 1 variance or that total variance within group or cluster (i.e., level 2 unit) that can be explained by the level 1 model; this parameter captures the within-group variability;

τ_{00} = level 2 variance component (i.e., that total variation explained at level 2—the cluster or group level); this parameter captures the between-group variability.

If ICC reflects the proportion of variance in the outcome that is *between* groups, then it stands to reason that the complement of ICC (i.e., $1 - \text{ICC}$) reflects the proportion of variance *within* groups (i.e., the proportion of variation in the outcome that is at level 1).

A question you may be asking is, ‘At what ICC value is multilevel modeling justified?’ There are different opinions on this, and this is ultimately a decision that the researcher must make and defend. ICC values of .05 to .15 have been applied as rules of thumb for common ICC values in educational research (Hedges & Hedberg, 2007). In planning cluster-randomized experiments in educational settings, larger ICC values (.15–.25) occur in national samples. In educational settings that specifically examine low-SES schools, a range of .05–.15 is more applicable to ICCs of unconditional models that examine low-achieving schools and ICCs of covariate-adjusted models (Hedges & Hedberg, 2007).

An ICC of zero means there is no variation in the outcome between groups. In this case, it is justified *not* to use a multilevel model. However, a model with a positive ICC value will be better served to be estimated in a multilevel framework (as compared to OLS regression), as the standard errors of the coefficients in OLS will be inaccurate (Snijders & Bosker, 2012).

11.1.1.3 Centering

In simple terms, centering is a rescaling of the independent variables and is done to increase interpretability of parameters. The concept of centering is not usually on the forefront of conversation for most statistical procedures but is a critical component to consider in multilevel modeling. This is because the choice of centering will impact the estimated value and interpretation of parameters. We will discuss centering in the context of the illustrative model of this chapter, a basic two-level model with individuals or cases at level 1 and clusters or groups at level 2. Other types of model (e.g., growth models) entail different centering considerations, and readers are referred in those cases to other sources (e.g., Singer & Willett, 2003). Additionally, we will constrain our conversation to mean centering, which is by far the most common option, although there are other types of centering (e.g., median).

In multilevel modeling, level 1 variables can be uncentered (i.e., using the raw metric) or centered at the grand (i.e., across all groups) or group (i.e., within cluster) mean. Going back to consideration of a basic two-level model, at level 2, there are essentially only two options: uncentered or grand mean centered. Group mean centering is not an option at level 2 in a two-level model given that all level 1 units share the same value on the level 2 predictor.

Uncentered

Selecting not to center means that the raw metric of the variable will be applied in the model—there is no rescaling of the variable. This only makes sense if a value of zero is a plausible score on the response scale of the explanatory variable. An uncentered independent variable will produce an intercept that represents the outcome when the independent variable is zero. In the case of binary variables where responses have been coded 0 and 1, where zero reflects the absence of the characteristic, selecting not to center may make sense as zero really does reflect a plausible value. However, other centering options—even for binary variables—may improve the interpretations of your parameters.

Grand Mean Centering

In grand mean centering, the average for the covariate (irrespective of cluster) is subtracted from the value of the level 1 predictor: $(X_{ij} - \bar{X}_\cdot)$. The intercept, under grand mean centering, is interpreted as an average adjusted mean, the predicted value for a case that is at the grand mean. When all level 1 predictors are grand mean centered, we can interpret the level 1 intercept as adjusted for the level 1 predictors (Heck & Thomas, 2000), and the intercept and slope variances can be interpreted as the expected variance when all predictors equal zero (i.e., the expected variance for the individual who is at the average) (Hox, 2010). Grand mean centering allows us to say, in other words, that if every group were the same in terms of the level 1

characteristics that have been specified in the model (i.e., the level 1 predictors), the expected value for the average outcome score of any group or cluster would be β_{0j} (Heck & Thomas, 2000). Grand mean centering “equalizes” the level 1 characteristics across groups, which in turn greatly impacts the variation of the cluster means (Heck & Thomas, 2000, p. 81). The challenge with grand mean centering is that the regression slope becomes difficult to interpret, as it includes a relationship component between level 1 and level 2.

Group Mean Centering

In group mean centering, the average for the covariate for the cluster to which the case belongs (i.e., cluster average) is subtracted from the value of the level 1 predictor: $(X_{ij} - \bar{X}_{.j})$. The scores are therefore relative to units *within the same group*, rather than all units as in the grand mean. The intercept, under group mean centering, is interpreted as an average unadjusted mean, the predicted value for a case that is at the group or cluster average. The regression slope in group mean centering is a “pure estimate” of the relationship between the covariate and the outcome at level 1 (Enders & Tofiqhi, 2007, p. 127).

Centering Recommendations

There is not a necessarily correct versus incorrect choice in centering, but the selection needs to be based on the research question that is attempting to be addressed (Enders & Tofiqhi, 2007). Recommendations have been offered based on the goal of the research question, but a research study with multiple questions may therefore likely require different centering, dependent on the focus of the research question.

- 1) Level 1 predictor has primary substantive value. If this is the goal, then **group mean centering** is recommended as a “pure” estimate of the pooled within-group regression coefficient produced, given that all the between-cluster variation is removed from the predictor (Enders & Tofiqhi, 2007, p. 128).
- 2) Level 2 predictor has primary substantive value. **Grand mean centering** is recommended when this is the research emphasis, as it quantifies the influence of the level 2 predictor in the regression coefficient, controlling for the level 1 predictor. Since better estimation at level 2 is the goal, any distortion that may occur with the level 1 estimates is of secondary concern (Enders & Tofiqhi, 2007).
- 3) Influence of predictor at level 1 and level 2. The question this addresses is one of differential effect: whether the influence of the predictor on the outcome at level 1 is the same as the influence of the aggregate group mean predictor at level 2. In this situation, an aggregate group mean score of the level 1 predictor is computed and included in the model. Thus, the predictor is included both at level 1 (e.g., individual level scores) and level 2 (aggregated group mean scores). Centering for examination of differential effects can be accomplished with **either group or grand mean centering**.

- 4) Interaction effects. The question this addresses is one of moderating effects or cross-level interaction effects: Does the level 2 predictor moderate the relationship between a level 1 predictor and the outcome? In other words, cross-level interactions occur when a higher-level model (e.g., level 2) moderates the strength of the relationship at the lower level (e.g., level 1). Including a level 2 predictor to the slope formula in the model will allow examination of cross-level interactions. When interaction effects (aka moderating effects or cross-level interactions) are the primary concern, **group mean centering** is recommended, as it offers a “pure” estimate of the moderation of the level 2 predictor on the level 1 relationship between the level 1 predictor and the outcome (Enders & Tofghi, 2007).

11.1.1.4 Model Estimation

We approach estimating multilevel models (not to be confused with model estimation methods, which we address later) as a multistep process with each step informing the last. The first model is a null model, from which the intraclass correlation is computed. This is followed by testing the level 1 model, and then finally the level 2 contextual model. Higher-order models (e.g., three- or four-level models) can be built similarly. Note, however, that different approaches to modeling can be taken, such as including *all* variables in the model at once and removing variables in a backward fashion (similar to stepwise regression). This, again, is the researcher’s call, but remember that how the model is built will need to be defended.

Null Model: The One-Way Random Effects ANOVA

The null, intercept-only, or unconditional model is of no substantive interest but does provide extremely helpful initial information, including:

1. The variation in the outcome is within and between level 2 groups (i.e., provides information about the outcome variability at each level)—this is used to estimate the extent to which there is sufficient within- and between-variation such that conducting multilevel modeling is needed (estimated via the intraclass correlation coefficient);
2. The reliability of each level 2 cluster’s sample mean as an estimate of its true population mean; and
3. A point estimate for the grand mean, γ_{00} .

In this unconditional model, the level 1 equation includes only the intercept term (β_{0j}); independent variables are not included in the model. In this way, the model is simply a one-way random effects ANOVA model. The outcome (Y_{ij}) is the sum of the respective group mean and a residual or individual error effect (r_{ij}) that reflects the deviation of unit i from group j . The level 1 equation is noted as follows, and this equation exists uniquely for each group in our data:

$$Y_{ij} = \beta_{0j} + r_{ij}$$

Thus, the level 1 outcome equates to the cluster's mean plus the cluster's error, and at level 2, the intercept therefore equals the mean of the cluster means. The level 2 equation is noted as follows:

$$\beta_{0j} = \gamma_{00} + u_{0j}$$

At level 2, the outcome (β_{0j}) is the sum of the average outcome across all level 2 groups (γ_{00}) plus an error effect (i.e., random intercept effect) or deviation of the group's mean from the grand mean (u_{0j}).

Each of the groups in our model has an OLS regression, where the coefficients represent the effects of the level 1 predictors on the level 1 outcome within each level 2 group. In aggregate, these regression models by group become a system of estimates whose overall mean and variance are averaged across all the level 2 groups. The outcomes in the level 2 model, therefore, are the intercept and slope(s) at level 1. This is why multilevel models are sometimes referred to as ‘intercepts and slopes as outcomes.’

The level 1 and level 2 models can be combined into one equation by substitution. This is what is referred to as a mixed model.

$$Y_{ij} = \gamma_{00} + u_{0j} + r_{ij}$$

In terms of variance components, the variance within clusters (σ^2) is the level 1 residual (r_{ij}), and the variance between clusters (τ_{00}) is the level 2 residual (u_{0j}).

Random Intercepts Model

In the model building process, a random intercepts model is considered next. In this model, the intercepts vary randomly across groups, but the slopes do not vary randomly across groups. In other words, the slope (β_{1j}) is now included (via inclusion of a level 1 predictor) but is fixed—it has the same effect for all groups at level 1. Building a theoretically sound model is critical, and should that not be done, “these models can quickly become oversaturated and obscure information necessary for testing hypotheses related to research questions . . . parsimony should be a key principle in model development” (Heck and Thomas, 2000, p. 79). This step is where the model begins to be built in relation to inclusion of predictors. In other words, one or more level 1 predictors are included in the model.

Using notation, a random intercept model with group mean centering of the predictor (X) has a level 1 model noted as follows (the term, ($X_{ij} - \bar{X}_{..}$), would have been included had grand mean centering been applied):

$$Y_{ij} = \beta_{0j} + \beta_{1j} (X_{ij} - \bar{X}_{..}) + r_{ij}$$

And a level 2 model, where the intercept and slope from the level 1 model become the outcomes, denoted as follows:

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10}\end{aligned}$$

Note that the slope outcome does not include a random component given this is simply a random intercept model.

The types of questions that can be answered with a random intercept model include:

1. What is the average intercept and slope for the level 2 groups/clusters and are these values statistically significantly different from zero?
2. How much do the regression equations (intercepts and slopes) vary from group to group?
3. How much is within-school variability (σ^2) reduced after including the level 1 predictor(s)?
4. Is there remaining statistically significant between-group variability (τ_{00})?
5. How much is the reduction in variation between groups (i.e., ICC)?

Random Coefficients Model: Random Intercepts and Random Slopes Model

In the model building process, a random intercepts and random slopes model is considered next. In this model, the intercepts vary randomly across groups, and the slopes also vary randomly across groups. In other words, the slope (β_{1j}) is now included (via inclusion of a level 1 predictor) and is allowed to be random—it is being tested to see if there are differences in slopes for the level 2 groups.

Using notation, a random intercept model with group mean centering of the predictor has a level 1 model noted as follows (the term, $(X_{ij} - \bar{X}_{..})$, would have been included had grand mean centering been applied):

$$Y_{ij} = \beta_{0j} + \beta_{1j} (X_{ij} - \bar{X}_{..}) + r_{ij}$$

And a level 2 model, where the intercept and slope from the level 1 model become the outcomes, denoted as follows, where the slope now includes a random effect:

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10} + u_{1j}\end{aligned}$$

With this model, in addition to the questions that could be addressed with the random intercepts model, we are now also able to determine the relationship between the intercepts and slopes.

Intercepts and Slopes as Outcomes Model

An intercepts and slopes as outcome model is considered next. In this model, the intercepts vary randomly across groups, and the slopes also vary randomly across groups. In addition, we add level 2 predictors to assist in explaining the variation in the intercepts and slopes

Using notation, a random intercept model with group mean centering of the predictor has a level 1 model noted as follows (the term, $(X_{ij} - \bar{X}_{..})$, would have been included had grand mean centering been applied):

$$Y_{ij} = \beta_{0j} + \beta_{1j}(X_{ij} - \bar{X}_{..}) + r_{ij}$$

And a level 2 model, where the intercept and slope from the level 1 model become the outcomes, denoted as follows, where the slope now includes a random effect, and both the intercepts and slopes equations include a group-level predictor, W , to assist in explaining the random effects at the group level.

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + \gamma_{01}(W) + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}(W) + u_{1j}\end{aligned}$$

Additional Models

As models increased in complexity by adding both level 1 and level 2 predictors, there are nearly infinite possibilities in model building. For example, a researcher may choose to model nonrandomly varying slopes as noted here:

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + \gamma_{01}(W) + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}(W)\end{aligned}$$

In other instances, a researcher estimating multiple slope coefficients may choose, based on theory and on model testing, to estimate some slopes as nonrandomly varying, some as random, and others as fixed as seen in this example, respectively. In addition, again with theory guiding these decisions, you may choose to have different variables explaining different parameters (e.g., W_1 and W_2).

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + \gamma_{01}(W_1) + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}(W_1) + \gamma_{12}(W_2) \\ \beta_{2j} &= \gamma_{20} + \gamma_{21}(W_1) + u_{2j} \\ \beta_{3j} &= \gamma_{30} + u_{3j}\end{aligned}$$

11.1.1.5 Estimation Methods

The estimation methods used in multilevel modeling assume that the fixed effects and variance components must be simultaneously estimated, given the complexity of

multilevel modeling. In OLS estimation, in comparison, the variance components can be estimated independent of estimating fixed effects (Swaminathan & Rogers, 2008). There are also various estimation methods applied at different levels and for different components in MLM. Generalized least squares (GLS) is the estimation for the level 2 coefficients (gamma). With GLS, each group's data is weighted by the inverse of its variance-covariance matrix. Empirical Bayes or Bayesian estimators (EB) are used to estimate randomly varying level 1 coefficients (i.e., β_{qj}). With EB, the estimates are weighted according to the reliability of the coefficients.

Maximum likelihood (ML) methods, including full maximum likelihood (FML) and restricted maximum likelihood (RML), are the most common multilevel estimation methods for the variance components, and each has different computational algorithms that can be used (though both are computationally iterative). In large samples, these estimation methods will produce quite similar results, but may produce very different results in small samples (Raudenbush & Bryk, 2002). Without getting into matrix algebra, the goal of full maximum likelihood is to do just that—maximize the likelihood. For multilevel modeling, this means that an integral matrix is first examined and then it is maximized. Thus ‘maximizing the likelihood’ is the process of maximizing the integral matrix with respect to the vector of estimated parameters (i.e., all variances and covariances and all fixed effects) so that inferences about those estimated parameters can be made (Raudenbush & Bryk, 2002). In other words, FML produces parameter estimates that maximize the probability of finding the observed sample data using vectors and matrix algebra. If you are using HLM to estimate your MLM, FML is the default estimation method for estimating the variance components. In comparison, restricted maximum likelihood (RML) estimates are computed from data that is transformed, so that the likelihood function is not influenced by parameters that are of no interest (e.g., within-group variance).

Neither FML nor RML is unequivocally better than the other. The decision on which to use to estimate your model may best be based on the group sample size and how you plan to use hypothesis testing as a component of your model building. When the number of groups is large, both FML and RML will produce similar estimates. When the number of groups is small, FML may produce biased variance estimates as compared to RML—more specifically, variance estimates that are too small. If you plan to use hypothesis testing as a component of your model building, FML may provide an advantage if you are looking to compare model effects in terms of both fixed and random effects. RML allows model comparisons only in terms of random effects (i.e., the models being compared must have the same fixed effects). If you find RML is the most appropriate estimation method *and* you want to use hypothesis testing, then simply use the RML results for the coefficients and variance estimates and rerun the models using FML to pull only the hypothesis testing results.

11.1.1.6 Model Fit

There are a few model fit indices that can be used to compare models based on changes in fixed and/or random effects. These are considered multiparameter tests, as they

allow nonnested models (i.e., models that are not a subset of the other) to be compared. The model fit tests include (a) chi-square deviance test, (b) Akaike information criteria (AIC), (c) Bayesian information criteria (BIC), and (d) sample size adjusted Bayesian information criteria (SBIC).

Deviance Test

The difference deviances likelihood ratio test is used for comparing nested models (e.g., a restricted model to a more parameterized model) and is considered a lack of fit or badness of fit between the model and the data. This is because the larger the deviance, the poorer the model fit. The deviance test follows a chi-square distribution, and statistically significant results indicate evidence that the more parameterized model has better model fit. This test can be computed using HLM by inputting the deviance value and number of parameters from a restricted model (estimated using full maximum likelihood) when generating a more parameterized model (also estimated using full maximum likelihood).

AIC

The AIC, which can be computed for nested or nonnested models, is computed as follows:

$$AIC = D + 2p$$

Where D is the deviance statistic computed using full maximum likelihood and p represents the number of parameters in the model. By including the number of parameters in the model, the AIC compensates for the parameters estimated and encourages parsimony. Models with low AIC are preferred. However, there are no criteria for determining how much the AIC has to decrease in order for one model to have better fit as compared to another.

BIC

The BIC, which can be computed for nested or nonnested models, is computed as follows:

$$BIC = D + \ln(n)(2p)$$

Where D is the deviance statistic computed using full maximum likelihood, p represents the number of parameters in the model, and $\ln(n)$ represents the natural log of the sample size. There is a lack of consensus on which sample size (e.g., within group or number of clusters) to use; however, using the larger sample size (i.e., the number of level 1 units) requires a more parsimonious model and thus may be preferred. Models with low BIC are preferred. Raftery (1995) developed guidelines for interpreting changes in BIC between a less-parameterized and more-parameterized model (i.e., $BIC_{Model1} - BIC_{Model2}$). These include:

- < 2 = weak evidence to favor model 2 over model 1
- 2–6 = positive evidence to favor model 2 over model 1

- 6–10 = strong evidence to favor model 2 over model 1
- > 10 = very strong evidence to favor model 2 over model 1

SBIC

The SBIC, which can be computed for nested or nonnested models, is computed as follows:

$$SBIC = D + \left[\ln\left(\frac{n+2}{24}\right) \right] (2p)$$

Where D is the deviance statistic computed using full maximum likelihood, p represents the number of parameters in the model, and \ln represents the natural log—this time it is the natural log of the ratio of the sample size plus 2 to 24. Again, there is a lack of consensus on which sample size (e.g., within group or number of clusters) to use; however, using the larger sample size (i.e., the number of level 1 units) requires a more parsimonious model and thus may be preferred. Models with low SBIC are preferred.

11.1.2 Sample Size

Sample size is often a difficult question to answer with single-level analyses, and the question of sufficient sample size becomes even more complex to answer with multilevel models. In general, in MLM, the sample size at the highest level is primarily of most concern, because the sample size at that level is always smaller than at the lowest level (Maas & Hox, 2005). In a two-level model, this would be the sample size at the group or cluster level. There are a few recommendations that exist for minimum group sample size, including more than 10 groups (Snijders & Bosker, 1999), assuming restricted maximum likelihood is the estimation method, and a minimum of 30 groups (Kreft & de Leeuw, 1998). Sample size at the lowest level is less a concern, and groups with even just one observation in them should be retained. While those groups will not contribute to the within-group variances, they will contribute to the between-group variance and overall average.

In addition to considering the sample sizes at each level, the proportion of variation in the outcome between groups (i.e., ICC) as well as the estimation method (FML or RML), are also considerations. Simulation research has been conducted that has conditioned on estimation methods (FML, RML), number of groups (30, 50, 100), size of groups (5, 30, 50), and ICC (.1, .2, .3) (Maas & Hox, 2005). In all conditions, regression coefficients and variance components are unbiased. However, the standard errors of the variances at level 2 are underestimated when the number of groups is less than 100; however, the bias is, “in practice, probably acceptable” (Maas & Hox, 2005, p. 91). Conditions were also tested with only 10 groups based on work by Snijders and Bosker (1999). While the regression coefficients and level 1 variance components were unbiased, the level 2 variance components were overestimated, and the standard errors

were unacceptably underestimated, suggesting that 10 groups at level 2 is insufficient for estimating MLM (Maas & Hox, 2005). Optimal Design is a free program available online that is designed to estimate power and sample size in group randomized designs and can be used a priori or post hoc (see <https://sites.google.com/site/optimaldesign-software/home>). Even if you are not in a situation where randomization of groups will be or has been done, Optimal Design may provide the best available information for estimating sample size.

11.1.3 Power

Power, in basic terms, is the probability of correctly rejecting the null hypothesis when it is false—detecting a difference if it is really there. In single-level models, we are accustomed to power being a function of one element (e.g., sample size) when holding the other elements (e.g., effect size, alpha) constant. This is not the case with MLM. As with sample size recommendations, estimation of power in MLM is quite complicated (Kreft & de Leeuw, 1998) and depends on a number of elements that are working together simultaneously (e.g., ICC, level 1 and level 2 sample sizes and covariates, estimation method). On top of this, power differs for parameters depending on fixed versus random effects, variance components, and cross-level interactions; and for a given level of power, the requisite sample sizes differ depending on the parameter of interest (Scherbaum & Ferreter, 2009). As stated by Spybrook (2014, p. 352), “the fact that different levels drive the power makes it nearly impossible to power for both.”

Simulation research suggests that sample sizes at level 1 of at least 20–40 and at level 2 of at least 30 are required to achieve power of level 1 predictors of .80 or greater (Bell et al., 2010). Generally, “increasing the sample size at the highest level (i.e., sampling more groups) will do more to increase power than increasing the number of individuals in the groups” (Scherbaum & Ferreter, 2009, p. 352), and research suggests this is the case at varying ICC values (e.g., Browne & Draper, 2000). As indicated previously, simulation research suggests that level 2 (i.e., group) sample sizes less than 30 result in standard errors for fixed effects that are underestimated (Maas & Hox, 2005). In research conducted by Mass and Hox (2005), as stated previously, even though the standard errors of the level 2 variances are underestimated when the number of groups is less than 100 (which leads to overestimated power), this is, “in practice, probably acceptable” (Maas & Hox, 2005, p. 91).

Power to detect cross-level interactions is primarily impacted by the magnitude of the cross-level effect, the variability of the level 1 slope coefficients, and the average level 1 and level 2 sample sizes (Mathieu, Aguinis, Culpepper, & Chen, 2012), with the power to detect cross-level interaction effects driven by the within-group sample size (Spybrook, 2014). Thus, in contrast to the general recommendation that the sample size of higher-level groups is more important than the sample size of lower levels, when cross-level interactions are of utmost importance, more emphasis is needed on gaining larger sample sizes within groups. “There is about a 3:2 premium on the average size

of the lower level samples, as compared to the upper level sample size” (Mathieu et al., 2012, p. 961). It is important to know that the power for a cross-level interaction (i.e., moderating effect) will always be less than the power for a main effect. Given this, in order to have sufficient power to detect moderation, a study interested in cross-level interaction may be overpowered in relation to detecting main effects (Spybrook, 2014).

Given the multitude of considerations, what is the recommendation on power? Power should be conducted on the parameter that is of most interest. For example, if fixed effects are of greatest interest (which is usually the case), then estimating power of these effects should be conducted first.

11.1.3.1 Power for Cluster Randomized Trials

The discussion on power thus far has been generic in the sense of the design of the study. However, you may find yourself in the fortunate situation of a cluster randomized design, where randomization occurs and treatment is applied at the group level rather than at the individual level; however, the outcome of interest is at the individual (not group) level. In this situation, power is a function of group size (i.e., the number of units per group), total number of groups, within-group variation, between-group variation, and true mean difference between groups (Spybrook, 2008), with power for cluster randomized MLM designs being dominated by the number of groups (Snijders & Bosker, 1993).

11.1.3.2 Software for Computing Power

There are a few software programs available for computing MLM power. *ML Power Tool* (R Development Core Team, 2011) is an *R* executable program for a priori estimation of cross-level interaction power. Optimal Design (Raudenbush, 2011) (http://sitemaker.umich.edu/group-based/optimal_design_software) is a freely available tool for computing power for person-level (single-level, multisite, or repeated measures), group-level randomization with either person-level or group-level outcomes, and meta-analysis.

11.1.4 Effect Size

Effect size in MLM can be computed for overall model components as well as for individual coefficients. A summary of the effects, the components of the model for which they are helpful in defining the effect (e.g., overall, within- or between-groups), the equations, and the interpretations are offered in Box 11.1.

11.1.4.1 Effect Size: Overall Model

The *proportion reduction in variance explained within-groups* (i.e., at level 1) answers the question: How much is the within-group variability reduced after adding predictors? This effect is estimated as follows, a pseudo-variance explained effect:

$$R^2_{\text{Level 1}} = \frac{(\sigma_{\text{Null Model}}^2 - \sigma_{\text{Full Model}}^2)}{\sigma_{\text{Null Model}}^2}$$

Sigma squared of the null model is applied as the model to which the conditional or contextual model is compared, as it represents the total within-group variation that can be explained by any level 1 model.

The *proportion reduction in variance explained between-groups* (i.e., at level 2) answers the question: How much is the between-group variability reduced after including one or more predictors? This effect is estimated as follows, a pseudo-variance explained effect:

$$R^2_{\text{Level 2}} = \frac{(\tau_{00 \text{ (null)}} - \tau_{00 \text{ (full model)}})}{\tau_{00 \text{ (null)}}}$$

One small word of caution in the application of the proportion reduction effect size values. Negative R squared values are computationally possible and can be produced when a variable that has almost no variation at one of the levels is included in the models.

11.1.4.2 Effect Size: Within-Group

Effect sizes for the individual beta coefficients within-groups (i.e., for use with level 1 predictors) can also be computed. Measured in standard deviation units (specifically the within-group standard deviation), small effects are generally interpreted as less than .20; moderate effects, .20 to .50; and large effects, more than .50 (Lee, Loeb, & Lubeck, 1998).

$$ES_{\text{Within-Group}} = \frac{\text{Beta Coefficient}}{\sqrt{\sigma^2}}$$

11.1.4.3 Effect Size: Between-Groups for Intercepts

Effect sizes for the between-group gamma coefficients for the intercepts (i.e., for use with level 2 predictors used to predict the level 2 intercept) can also be computed. Measured in standard deviation units (specifically the between-group standard deviation), small effects are generally interpreted as less than .20; moderate effects, .20 to .50; and large effects, more than .50 (Lee et al., 1998).

$$ES_{\text{Between-Group Gamma}} = \frac{\text{Gamma Coefficient}}{SD_{\text{Between-Groups}}} = \frac{\gamma}{\sqrt{\tau_{00}}}$$

11.1.4.4 Effect Size: Between-Groups for Slopes

Effect sizes for the between-group gamma coefficients for the slopes (i.e., for use with level 2 predictors used to predict the level 2 slopes) can also be computed. Measured in standard deviation units (specifically the between-group standard deviation of the slope), small effects are generally interpreted as less than .20; moderate effects, .20 to .50; and large effects, more than .50 (Lee et al., 1998).

$$ES_{\text{Between-Group Slopes}} = \frac{\text{Gamma Coefficient}}{SD_{\text{slope}}} = \frac{\gamma}{(SE_{\text{gamma}})(\sqrt{N_{\text{group}}})}$$

BOX 11.1 MULTILEVEL MODELING EFFECT SIZE MEASURES

Effect Size	Equation	Interpretation
Overall Model	$R^2_{\text{Level 1}} = \frac{(\sigma_{\text{Null Model}}^2 - \sigma_{\text{Full Model}}^2)}{\sigma_{\text{Null Model}}^2}$	This is a pseudo-variance measure, the <i>proportion reduction in variance explained within-groups</i> (i.e., at level 1). This effect indicates the amount of within-group variability that is reduced after including level 1 predictors.
	$R^2_{\text{Level 2}} = \frac{(\tau_{00(\text{null})} - \tau_{00(\text{full model})})}{\tau_{00(\text{null})}}$	This is a pseudo-variance measure, the <i>proportion reduction in variance explained between-groups</i> (i.e., at level 2). This effect indicates the amount of between-group variability that is reduced after including level 1 and/or level 2 predictors.
Within-Group	$ES_{\text{Within-Group}} = \frac{\text{Beta Coefficient}}{\sqrt{\sigma^2}}$	This is the effect size for the individual beta coefficients within-groups (i.e., for use with level 1 predictors). This effect is measured in standard deviation units (specifically the within-group standard deviation). <ul style="list-style-type: none"> • Small effects < .20 • Moderate effects > .20 and < .50 • Large effects > .50
Between-Group	$ES_{\text{Between-Group Gamma}} = \frac{\text{Gamma Coefficient}}{SD_{\text{Between-Groups}}} = \frac{\gamma}{\sqrt{\tau_{00}}}$	This is the effect size for the between-group gamma coefficients for the intercepts (i.e., for use with level 2 predictors used to predict the level 2 intercept). This effect is measured in standard deviation units (specifically the between-group standard deviation). <ul style="list-style-type: none"> • Small effects < .20 • Moderate effects > .20 and < .50 • Large effects > .50

11.1.5 Assumptions

Multilevel linear modeling assumptions include (1) linearity (i.e., linear relationship between variables), (2) residuals (i.e., random effects) at level 1 (i.e., r_{ij}) are normally

distributed and have equal variances (i.e., σ^2) for every level 1 unit within every level 2 group, (3) residuals at level 2 (i.e., u_{0j}, u_{1j}) are multivariate normal and homoscedastic (i.e., residual variance is constant—it does not relate to the predictor variables), and (4) predictors at a respective level are uncorrelated with random effects at the other level (e.g., level 2 predictor, W , is uncorrelated with the level 1 random effect, r_{ij} ; and level 1 predictor, X , is uncorrelated with the level 2 random effect, u_{0j}).

11.1.5.1 Linearity

The first assumption of linearity obviously does not hold in nonlinear multilevel models (e.g., hierarchical generalized linear models). Linearity can be examined in scatterplots of level 1 predicted values to level 1 residuals. Linearity is assumed when there is a random display of points. Procedures to deal with nonlinearity include transformations (of one or more of the predictor and/or of outcome variable) and other regression models (e.g., hierarchical generalized linear models).

Using the level 2 residuals, scatterplots of empirical Bayes residuals to level 2 predictors can assist in assessing linearity. Linearity can be assumed when there is a random distribution of points around zero (Raudenbush & Bryk, 2002).

11.1.5.2 Normality

Normality is an assumption that is important to consider in MLM when the outcome is continuous and a linear relationship is expected. Examining normality is done within the context of examination of residuals. Computation of residuals in multilevel modeling is more complex than OLS however, as each level of a multilevel model generates its own residuals (Van der Leeden, Busing, & Meijer, unknown). More specifically in MLM, the level 1 residuals, r_{ij} , are normally distributed and the level 2 residuals (i.e., $(u_{0j}, u_{1j}, \dots, u_{pj})$, the level 2 random coefficients) are also normally distributed.

We've already discussed how normality of residuals occurs at each level of the multilevel model. Given this, examining normality at each level is important. For level 1 residuals, the tools with which we have examined normality in previous procedures (e.g., skewness and kurtosis, normal probability plots, formal tests of normality, box-plots, and histograms) can be applied.

Using the level 2 residuals, plots of empirical Bayes residuals to level 2 predictors can also assist in assessing normality, and more specifically, potential outliers (Raudenbush & Bryk, 2002). Probability plots of Mahalanobis distance can also be used to gauge the extent of normality. The closer the points adhere to the diagonal line, the more evidence there is to suggest normality at level 2 is reasonable. Confidence in the specification of the fixed effects is important prior to consideration of the normal probability plots. This is because normal probability plots can result not only in cases of nonnormal residuals but also in cases where the fixed effects are misspecified (Eberly & Thackeray, 2005).

In terms of normality, even with relatively severe residual nonnormality at both level 1 and level 2, *fixed* parameter estimates are relatively robust, but downward biased *random* parameter estimates are produced (Shieh, Fouladi, & Pullman, 1999). In other words, the regression coefficients and the standard errors associated with the regression coefficients have little or no bias in cases of nonnormality; however, the variance components and the standard errors of the variance components may be substantially biased (Maas & Hox, 2004; Shieh et al., 1999; Van der Leeden et al., unknown). Other research suggests that violations of normality at level 2 lead to biased fixed effect estimates; however, this improves with more groups (i.e., more level 2 units) but greater sample sizes of the groups doesn't matter (Maas & Hox, 2004). Nonnormality still has a greater effect on the random components, however (Maas & Hox, 2004). Should robust standard errors differ markedly from the asymptotic standard errors, nonnormality may be evident (Raudenbush & Bryk, 2002).

11.1.5.3 Homogeneity of Variance

In multilevel modeling, homogeneity of variance occurs when the level 1 random effects have equal variance (i.e., σ^2) across the level 2 groups. Violation of this assumption in terms of *random* variation of the level 1 variance over level 2 groups has little impact on the estimated coefficients at level 2 in terms of bias—the results are unbiased and the standard errors remain accurate. On the other hand, if there is *systematic* (i.e., nonrandom) variation of the level 1 variance over level 2 groups, bias can result and more problematically, heterogeneity can suggest model misspecification (Raudenbush & Bryk, 2002). Omission of critical level 1 predictors and erroneous treatment of fixed, random, and nonrandomly varying effects (e.g., a fixed effect should actually be specified as a random effect) can create heterogeneity (Raudenbush & Bryk, 2002). Thus, the link between model specification and assumptions is very clear. While violating this assumption is less problematic in the MLM framework relative to biased parameter estimates and inaccurate standard errors, violation should be a light bulb for closer examination of the specifications of the model.

11.1.5.4 Homoscedasticity

Regarding level 2, the level 2 residuals (or random coefficients) are also assumed to have a constant covariance matrix—also known as *homoscedasticity*. This is the assumption that the within-group and between-group intercept and slope residuals have constant variances. Failing to recognize nonconstant variances, heteroscedasticity can result in erroneous decisions on hypothesis tests that involved the variables responsible for the nonconstant variance (Snijders & Bosker, 2012). Homoscedasticity can be examined by plotting standardized level 2 residuals to level 2 predictors. Should evidence suggest that the assumption of homoscedasticity is violated, including additional theoretically driven predictors at either level 1 or level 2, it may account for differential residual variance at level 1. Another option that may be most applicable in situations where the dependent variable is nonnormally distributed is to apply a nonlinear transformation, such as the national log or square root, to the dependent

variable. In cases of heteroscedasticity where the outcome has a limited number of categories, multilevel ordered logit may be applied or the dependent variable may be dichotomized with the binary outcome being estimated with a multilevel logistic model (Snijders & Bosker, 1999).

11.1.5.5 Uncorrelated Predictors and Random Effects

To speak of uncorrelated predictors and random effects is really to address issues of model specification (and misspecification). Some researchers indicate that model specification is an assumption, not a condition, as misspecification of a model can have serious consequences on the results and ensuing interpretations. Model misspecification at level 1 has consequences for the level 2 estimates, including biased intercept and slope parameter estimates of level 2 predictors (Raudenbush & Bryk, 2002). The implication of the assumption of the level 2 random coefficient having a population mean of zero, conditioned on the predictors, is that the random intercepts and random slopes are uncorrelated with the covariates. In situations where this assumption does not hold (i.e., there is a nonzero correlation between the random intercepts and slopes with the predictors), the inclusion of predictors can remedy the issue. Nonzero correlations between the level 2 residuals and level 1 predictors, X , can be addressed by including an interaction term of the level 1 predictor and the level 1 predictor's group mean, $X(\bar{X})$. Nonzero correlations between the level 2 residuals and level 2 predictors, W , can be addressed by including an interaction term of the level 1 and level 2 predictor, XW .

Model misspecification at level 2, specifically omission of a level 2 predictor that is related to a level 1 predictor, will produce a biased coefficient for the level 1 predictor and result in correlation between the level 2 random effects and the level 1 predictor. Two approaches can deal with this model misspecification at level 2. First, include at level 2 the group mean of the level 1 predictor. This removes confounding between the level 2 omitted variable and the level 1 predictor. Second, group mean centers the level 1 predictor, which will eliminate the covariance between the group mean of the predictor and the random effect (Raudenbush & Bryk, 2002).

Model specification at level 2 means that the random level 2 effects, i.e., u_{gj} , in each level 2 equation are unrelated to the predictors in the respective level 2 equation. When model misspecification occurs (e.g., in the presence of a confounding variable, a significant predictor that is related to one or more level 2 predictors in the model), the level 2 coefficient estimates will be biased (Raudenbush & Bryk, 2002). A relatively easy, but not necessarily theoretically plausible nor iteratively achievable, solution to this issue is the inclusion of the same set of predictors for all level 2 equations. When this occurs, the results are unbiased, as the level 2 coefficients in the slopes are independent of the level 2 intercept (Raudenbush & Bryk, 2002). Another way to examine potential model specification at level 2 is to fix the level 2 slopes (which eliminates the covariance between the slope and intercept) and compare those results to a level 2 random effects slopes model. If the gamma coefficients are markedly different, this suggests model misspecification (Raudenbush & Bryk, 2002).

11.1.6 Conditions

The obvious condition that is applicable to multilevel modeling is nested data structure with enough groups at the highest level to support a multilevel analysis (see the discussion on sample size for further considerations regarding group sample size). Illustrated in this chapter is the most basic situation of units nested within groups, which can be answered through a two-level linear model. Higher levels of nesting, however, can be examined, including three- and four-level models. Additionally, while units within groups are common, it may be the nesting is measurement occasions within person (or other unit of analysis) or the unit is simultaneously nested (e.g., an employee is nested both within the organization they work for as well as the neighborhood in which they live).

In previous procedures, a common condition of the test relates to the measurement scale of the variables applied. Multilevel modeling, in the broadest sense, is different in that all measurement scales can be accommodated. However, given a specific type of multilevel model (e.g., multilevel linear model), a condition of the test *does* become that only certain scales of measurement are appropriate for the outcome. A two-level linear model was illustrated in this chapter. Given that, and as we learned with multiple regression, a condition of the test is that the scale of measurement of the dependent variable is at the interval or ratio level. The independent variables, just like in OLS, can be interval or ratio but multilevel modeling can also accommodate independent variables which are categorical—nominal or ordinal in scale—as long as they are dummy coded with one of the categories removed to serve as the reference category. However, there is one caveat to the measurement scale of the independent variable. Keep in mind, however, that multilevel modeling is extremely flexible; because they fall under the general linear framework, all types of analyses can be performed within the broad umbrella of ‘multilevel models.’ Outcomes that are nonlinear (e.g., binary or multinomial) can easily be accommodated through multilevel counterparts to logistic and multinomial logistic modeling (i.e., hierarchical generalized linear modeling). Multiple outcomes that are examined simultaneously can also be accommodated (i.e., hierarchical multivariate linear modeling) as can latent variable models (i.e., the multilevel counterpart of structural equation modeling). These are just a few. The point is, if you can examine it as a single-level model, there is likely a multilevel cousin that can be applied when the data have a nested structure, so don’t let the scale of measurement of your variables or the distributional shape deter you from considering a multilevel model when there is a hierarchical structure to your data.

In summary, the conditions of a multilevel linear model that is a two-level linear model are as follows: (a) most importantly, the data have a nested structure where one unit is nested within a group or cluster, (b) one or more independent variables of any measurement scale, (c) independent variables that are categorical are dummy coded with one of the categories held out as the reference category, and (d) there is one dependent variable that is measured at least at the interval level.

11.2 MATHEMATICAL INTRODUCTION SNAPSHOT

The mathematical structure of a two-level multilevel linear model with one predictor at each level and random effects for both the intercepts and slopes is as follows:

$$\text{Level 1: } Y_{ij} = \beta_{0j} + \beta_{1j}X_{ij} + r_{ij}$$

$$\begin{aligned}\text{Level 2: } \beta_{0j} &= \gamma_{00} + \gamma_{01}W_j + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}W_j + u_{1j}\end{aligned}$$

The level 1 equation should be quite familiar as (with the exception of the subscripts) it is an OLS equation. In multilevel modeling, the level 1 equation models the relationship *within* the lower-level units. Let's review each component of the level 1 equation:

Y_{ij} = dependent variable for each i th level 1 unit nested within the j th level 2 unit (i.e., the groups or clusters)

β_{0j} = intercept for the j th level 2 unit

β_{1j} = regression coefficient for predictor X_{ij} for the j th level 2 unit

X_{ij} = value on the level 1 independent variable (uncentered)

r_{ij} = random error associated with the i th level 1 unit nested within the j th level 2 unit

Should the independent variable be group mean centered, the notation for the independent variable would be more specifically reflected as $(X_{ij} - X_{..})$, and with grand mean centering, the notation for the independent variable would be more specifically reflected as $(X_{ij} - X_{..})$.

The level 2 equation models the relationship *between* the lower-level units (i.e., across multiple groups or clusters). There are two equations at level 2 in this particular illustration, one to reflect each of the level 1 coefficients, specifically β_{0j} and β_{1j} , and these level 1 coefficients are now the outcomes in the level 2 model. Should additional parameters be included in the level 1 model (e.g., additional independent variables), that will be reflected by additional level 2 equations. Let's review each component of the level 2 equations. Keep in mind that interpretations of the overall average intercept, in particular, are dependent on the centering of the level 1 predictors. Reviewing the discussion of centering in relation to interpretation of the overall intercept, γ_{00} , is recommended.

β_{0j} = intercept for the j th level 2 unit (random effect)

γ_{00} = overall average intercept, with the specific interpretation dependent on the centering of the level 1 predictors (fixed effect)

γ_{01} = regression coefficient for W_j relative to the level 1 intercept, β_{0j} (fixed effect)

W_j = value on the level 2 independent variable

u_{0j} = random effect of the j th level 2 unit adjusted for W_j on the level 2 intercept, β_{0j}

β_{1j} = slope for the j th level 2 unit (random effect)

γ_{10} = overall average slope

γ_{11} = regression coefficient for W_j relative to the level 1 slope, β_{1j} (fixed effect)

W_j = value on the level 2 independent variable

u_{1j} = random effect of the j th level 2 unit adjusted for W_j on the level 2 slope, β_{1j}

The dependence of errors is evident in the level 2 equations. The random error terms in level 2, specifically u_{0j} , model the dependency among the level 1 units that are clustered within each level 2 unit. These are residuals at the group or cluster level and are group effects that are unexplained by level 2 predictors (i.e., they model the unexplained variability at level 2) (Snijders & Bosker, 2012).

The level 1 and level 2 equations can be combined into one model (i.e., a mixed model) that includes both fixed and random effects. More specifically, this mixed model incorporates both level 1 and level 2 predictors, a cross-level or interaction term (i.e., $W_j X_{ij}$ and implied by the inclusion of X and W), and a composite error term ($u_{1j} + u_{0j} + r_{ij}$) as follows:

$$Y_{ij} = \gamma_{00} + \gamma_{10}X_{ij} + \gamma_{01}W_j + \gamma_{11}W_j + u_{0j} + u_{1j} + r_{ij}$$

11.3 COMPUTING MULTILEVEL MODELING USING HLM

Next, we consider HLM software for multilevel linear modeling. Before we conduct the analysis, let us review the data. For this example, data was drawn from a personality testing site (<http://personality-testing.info/>). Users complete online personality tests anonymously, and the raw data is collected and available for public use (http://personality-testing.info/_rawdata/). For this illustration, we are using data collected on the Big Five Personality Test (Goldberg, 1992). At the time the data were downloaded, the online files were updated through 5/18/2014. I have taken the liberty to clean up the data (e.g., recoding dummy variables, creating composites for the five personality types, recoding ages of nine participants who indicated their age was over 100 to ‘100,’ deleting one case with a missing country code, deleting two cases with missing age, deleting one case with missing Big Five scores). Additionally, for the multilevel illustration, we have delimited the sample to include only individuals who indicated they were from a country in North America (individuals $n = 9,868$; countries $n = 19$). To this data, I have added a variable to denote if the country is a developing country based on the International Statistical Institute’s (ISI) standards (see <http://www.isi-web.org/component/content/article/5-root/root/81-developing>). The countries denoted as developing were defined by the ISI and hold that classification from 1/1/2015 through 12/31/2015. ‘Developing country’ was defined as having gross national income per capita per year of \$11,905 or less. Because the student version of HLM limits the number of level 1 units to 8,000, a random sample of U.S. cases were drawn (there were initially 8,772 U.S. participants), and participants from the remaining countries were used in their entirety (individuals $n = 7,994$; countries $n = 19$). The filename is **Big5_N_7994.sav**.

The inclusion of an identification variable in your data is critical. The ID variable matches the lower-level units to the higher-level groups, and thus the ID defines the nesting structure. In a two-level model, only one ID is needed, and it represents the groups within which the lower-level units are nested. In a three-level model, two IDs are needed. One ID will match the lowest-level units to the level 2 units, and the second ID will match the level 2 units within the level 3 units. Data must be sorted in ascending order by the ID. When two or more IDs are required, sort by the highest-level group followed by the next highest, and so forth.

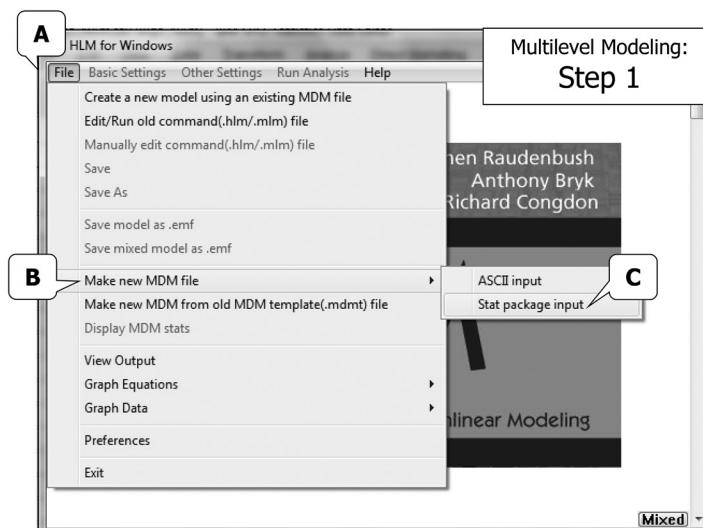
The Level 2, or country-level, variable includes whether or not the country is a developing country (coded as 1) per ISI standards

Level 1, or person-level, variables include:

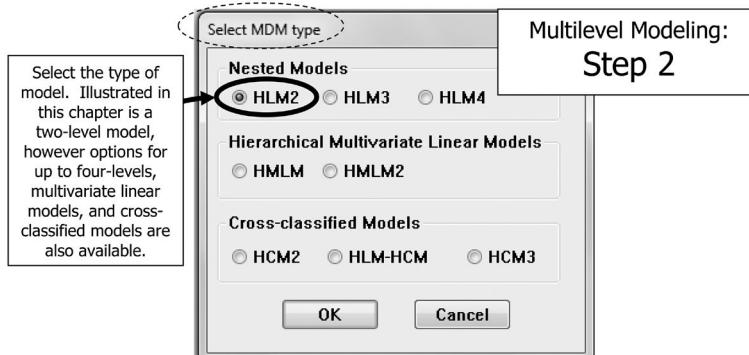
1. English is native language
2. Write with right hand
3. Male
4. Caucasian
5. Age
6. Big 5 subscale scores
 - a. Extraversion
 - b. Emotional stability
 - c. Agreeableness
 - d. Conscientiousness
 - e. Intellect/Imagination

	country_recode	N_Am	Developing	English_native	RightHand	Male	Caucasian	age	Extraversion	EmotionalStability	Agreeableness	Conscientiousness	Intellect_Imagination
1	5.00	1.00	0	1.00	1.00	.00	00	24	2.80	1.70	4.30	3.60	3.60
2	15.00	1.00	0	1.00	1.00	1.00	00	16	1.50	1.90	3.90	2.70	3.40
3	15.00	1.00	0	1.00	1.00	1.00	00	36	3.90	4.30	4.60	5.00	4.60
4	25.00	1.00	0	1.00	1.00	.00	00	20	2.40	3.00	3.80	4.60	4.00
5	25.00	1.00	0	1.00	1.00	.00	00	39	3.90	3.70	3.60	3.60	4.00
6	28.00	1.00	1.00	1.00	1.00	.00	00	20	3.20	4.10	3.10	3.80	3.80
7	28.00	1.00	1.00	1.00	1.00	.00	00	22	2.30	2.70	4.50	3.00	3.60
8	28.00	1.00	1.00	1.00	1.00	.00	00	22	3.30	3.40	4.60	4.50	4.00
9	28.00	1.00	1.00	1.00	1.00	1.00	00	23	3.70	3.10	3.10	3.60	3.50
10	28.00	1.00	1.00	1.00	1.00	1.00	00	25	2.90	3.40	3.90	4.40	3.30

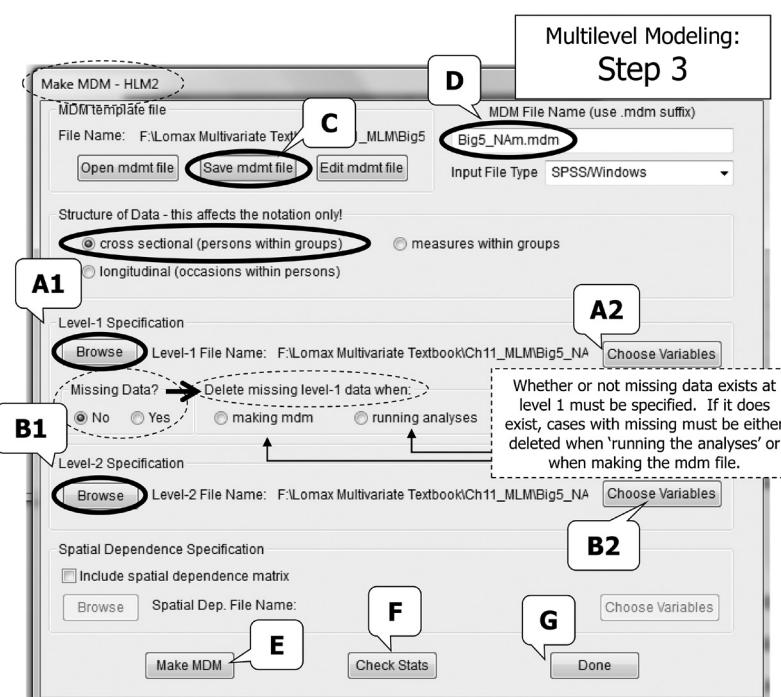
Step 1: We will use the student version of HLM to estimate the multilevel model. We first need to pull in our SPSS file. To do this, go to “File” in the top pull-down menu, then select “Make new MDM file,” and then select “Stat package input.” Following the screenshot (Step 1) produces the “Select MDM type” dialog box.



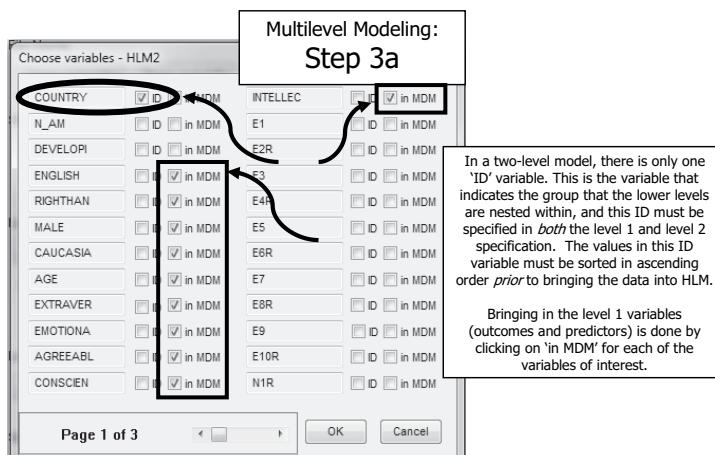
Step 2. A two-level linear model is illustrated in this chapter, so we click the radio button for 'HLM2' (see screenshot Step 2). Click on "OK" to bring up the Make MDM dialog box.



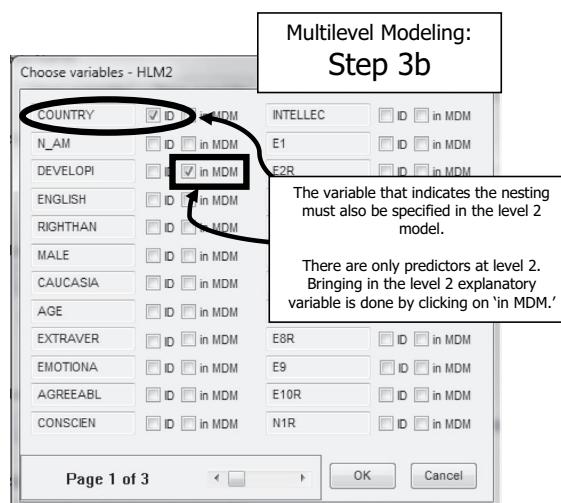
Step 3. From the Make MDM dialog box (see screenshot Step 3), there are a number of steps that must be completed and decisions that must be made. Specifying the variables at level 1 (labeled A1 and A2) and level 2 (labeled B1 and B2) is the first step. Clicking "Browse" will let you search for the data file that includes the level 1 variables. This may also be the data file that includes the level 2 variables; however, should the level 2 variables be contained in a different data file, HLM can accommodate that (in fact, previous versions of HLM actually required they be two separate data files—this is no longer the case).



Step 3a. Clicking on Browse for the level 1 specification (see screenshot Step 3a) will bring up the dialog box that allows selection of the level 1 variables (nesting variable, outcomes, and predictors). Only one 'ID' variable can be selected. This ID is the variable that specifies the nesting structure of the data (i.e., the higher order groups in which the lower level units are clustered). The outcomes and independent variables at level 1 are specified by placing a checkmark in the box for 'in MDM.' This step tells HLM that these variables, and only these variables, are to be brought into the next step. Click on "OK" to return to the Make MDM dialog box.



Step 3b. Clicking on Browse for the level 2 specification (see screenshot Step 3b) will bring up the dialog box that allows selection of the level 2 variables (nesting variable and predictors). The ID variable must again be specified in this level. The independent variables at level 2 are specified by placing a checkmark in the box for 'in MDM.' This step tells HLM that these variables, and only these variables, are to be brought into the next step. Click on "OK" to return to the Make MDM dialog box.



Step 3 (again!). From the Make MDM dialog box (see screenshot Step 3), after specifying the variables for level 1 and level 2, there are still a number of steps that must be taken. For the level 1 specification, the researcher must click the radio button that corresponds to the extent to which there is missing data in their data file. HLM does not work well with missing data, and essentially cases that have missing data are not included in the analyses (thus, it is in your best interest to deal with missing data *prior* to bringing your data into HLM). These cases are excluded at the point of either making the MDM file or running the analyses. If there are a minimal number of cases with missing data, and you anticipate *not* using the variables that have missing data, it is advantageous to select to delete the missing level 1 data when running the analyses. The default selection for 'Structure of Data' is cross-sectional (persons within groups). Should a different type of model be examined (e.g., longitudinal), then the selection for generating models with appropriate notation can be made by clicking the respective radio button.

Step 3 (again!). The next step is 'Save mdmt file' (see screenshot Step 3 item C). This command saves all the information that has just been input into HLM, and can later be reopened using the 'Open mdmt file' button in the top left of screenshot Step 3.

Step 3 (again!). Next, save the MDM file by type in the file name in the top right text box (see screenshot Step 3 items D and E). Note that '.mdm' must be included in the file name. The default input file type is SPSS. Should a different file type be used (e.g., SAS or Stata), the input file type can be changed by clicking the toggle bar. Once the name is specified, the MDM file can be saved by click on 'Make MDM.' This will allow you to browse to select where you want to save the MDM file. Once the MDM file is created, a box will appear briefly that indicates the level 1 and level 2 variables included, as well as basic descriptives on each. The most important information at this step is double-checking the level 1 sample size and number of groups at level 2. Should those numbers not correspond to the data that you had prior to bringing into HLM, this is the point where you need to go back to your original data and see why the data was not brought in properly (e.g., failing to sort by level 2 units, failing to deal with missing data at level 1). Another item that is noticeable at this point is the truncation of HLM of the variable names to eight characters. Should you have variable names that have the same first eight characters prior to bringing them into HLM (e.g., Question1, Question2, Question3), HLM will redefine the variable names with generic labels. Given this, it is in your best interest to adhere to eight character naming conventions prior to bringing your data into HLM.

LEVEL-1 DESCRIPTIVE STATISTICS					
VARIABLE NAME	N	MEAN	SD	MINIMUM	MAXIMUM
ENGLISH	7994	0.88	0.32	0.00	1.00
RIGHTHAN	7994	0.88	0.33	0.00	1.00
MALE	7994	0.34	0.47	0.00	1.00
CAUCASIA	7994	0.67	0.47	0.00	1.00
AGE	7994	26.88	12.84	13.00	100.00
EXTROVERT	7994	3.03	0.56	1.00	4.00
EMOTIONAL	7994	2.96	0.88		
AGREEABL	7994	3.98	0.73		
CONSCIEN	7994	3.41	0.74		
INTELLEC	7994	3.95	0.62		

LEVEL-2 DESCRIPTIVE STATISTICS					
VARIABLE NAME	N	MEAN	SD	MINIMUM	MAXIMUM
DEVELOPI	19	0.68	0.48		

7994 level-1 records have been processed
19 level-2 records have been processed

Once the MDM file is created, a box will appear briefly that indicates the level-1 and level-2 variables included as well as basic descriptives on each.

Step 3 (again!). It is advisable to save the file that includes the sample and group sizes and basic descriptive statistics in the event you later need to reference these values. To do so, click on 'Check Stats' to generate the same information (see screenshot Step 3 items F and G)—this time in a format that can be saved as a text file. Clicking on "Done" (see screenshot Step 3) will take you to the dialog box that allows the model of interest to be defined, based on the previous settings specified.

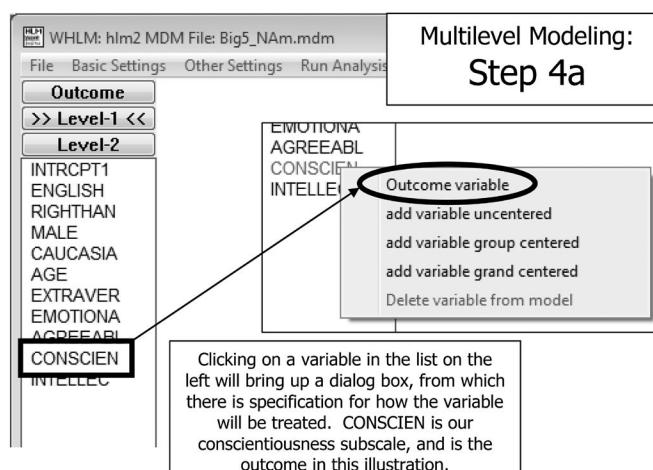
LEVEL-1 DESCRIPTIVE STATISTICS					
VARIABLE NAME	N	MEAN	SD	MINIMUM	MAXIMUM
ENGLISH	7994	0.88	0.32	0.00	1.00
RIGHTHAN	7994	0.88	0.33	0.00	1.00
MALE	7994	0.34	0.47	0.00	1.00
CAUCASIA	7994	0.67	0.47	0.00	1.00
AGE	7994	26.88	12.84	13.00	100.00
EXTRAVER	7994	3.03	0.96	1.00	5.00
EMOTIONA	7994	2.96	0.88	1.00	5.00
AGREEABL	7994	3.90	0.73	1.00	5.00
CONSCIEN	7994	3.41	0.74	1.00	5.00
INTELLEC	7994	3.95	0.62	1.40	5.00

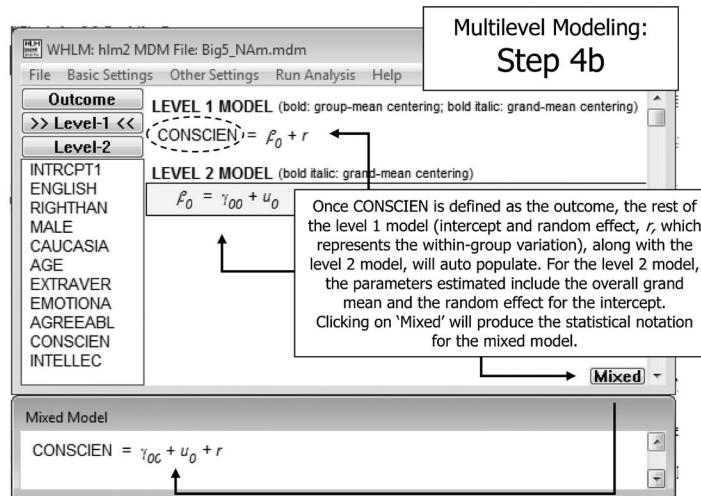
LEVEL-2 DESCRIPTIVE STATISTICS					
VARIABLE NAME	N	MEAN	SD	MINIMUM	MAXIMUM
DEVELOPI	19	0.68	0.48	0.00	5.00

MDM template: F:\Lomax Multivariate Textbook\ch11
MDM file name: Big5_NAm.mdm
Date: Mar 24, 2015
Time: 14:57:19

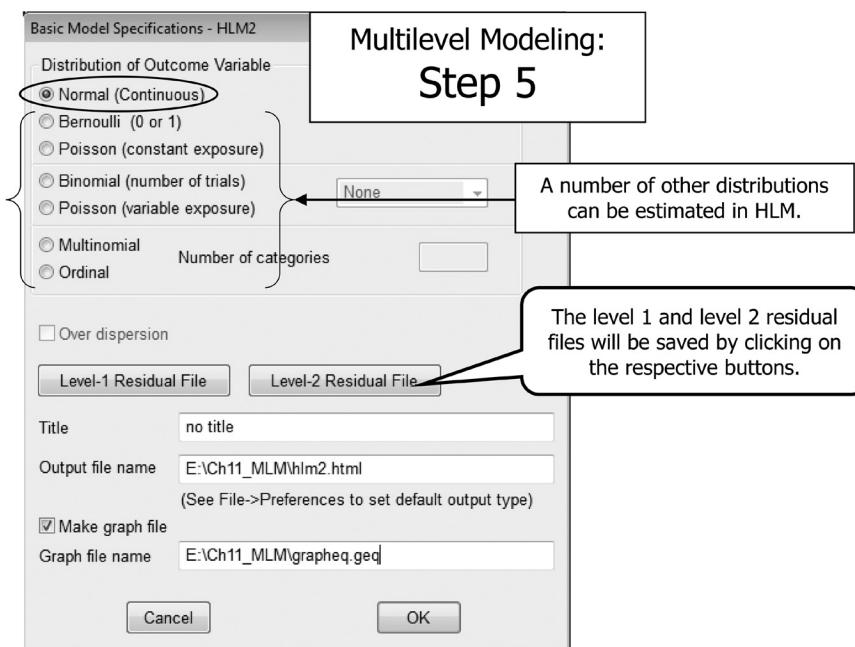
Information that appears when you 'check stats.'

Steps 4a and 4b. The dialog box at this point is split between level 1 and level 2, and this is where you specify the model specifics (see screenshots Steps 4a and 4b). The first thing that must be done is to specify the outcome. For this illustration, we will specify the Big Five subscale of Conscientiousness [CONSCIEN] as the outcome. This is done by clicking on the variable we wish to define as the outcome. Clicking on any of the variables will bring up the same options: (a) outcome variable, (b) add variable uncentered, (c) add variable group centered, (d) add variable grand centered, and (e) delete variable from model. The outcome variable must be defined first, thus we select 'Outcome variable' (see screenshot Step 4a). Once the outcome is defined, the model is automatically populated as seen in screenshot Step 4b.

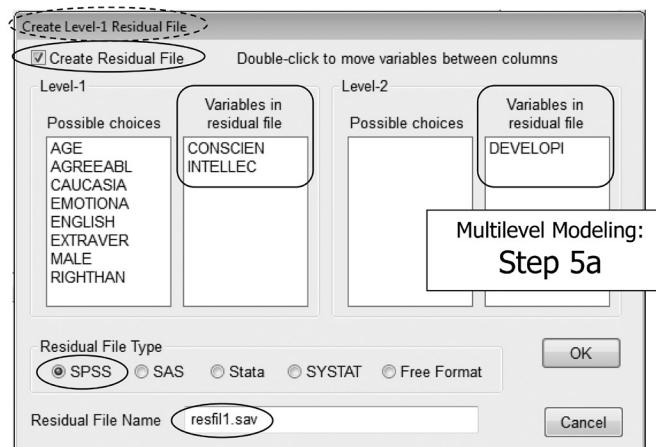




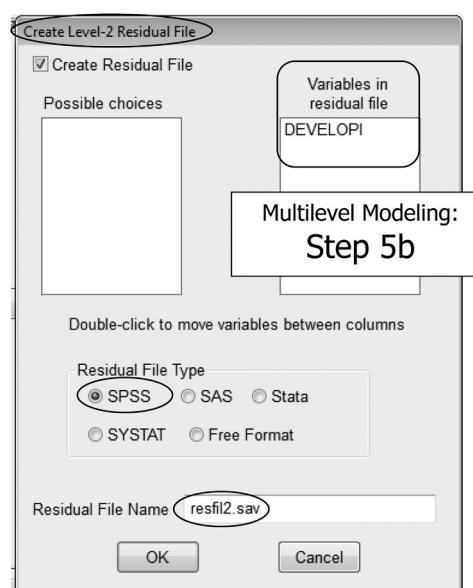
Step 5. Prior to running the model, it is important to check (and change as needed), model settings. These can be accessed from the top menu seen in Step 4b. Clicking on “Basic Settings” will bring up the model specification screen (see screenshot Step 5). From here, the distribution of the outcome variable can be defined. The default is ‘Normal (Continuous),’ and that is applicable for this illustration. However, many different nonlinear outcomes can be examined in the multilevel framework. From this screen, we will also request that the residuals be saved so they can later be examined for assumption evidence. This is done by clicking on the respective residual file buttons (specifically Level-1 Residual File and Level-2 Residual File).



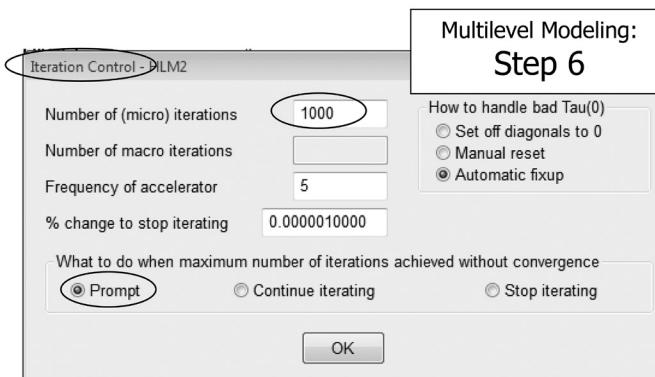
Step 5a. Clicking on Level-1 Residual File will bring up the Create Level-1 Residual File dialog box (see screenshot Step 5a). We want to include in our residual file the variables that have been included in levels 1 and 2 of our model, and this is done by listing the respective variables in the boxes labeled 'Variables in residual file.' This includes both the outcome and level 1 and level 2 predictors. The default file type is SPSS, and a default file name is provided. For purposes of illustration, we will leave those as is. Click "OK" to return to the basic setting dialog box.



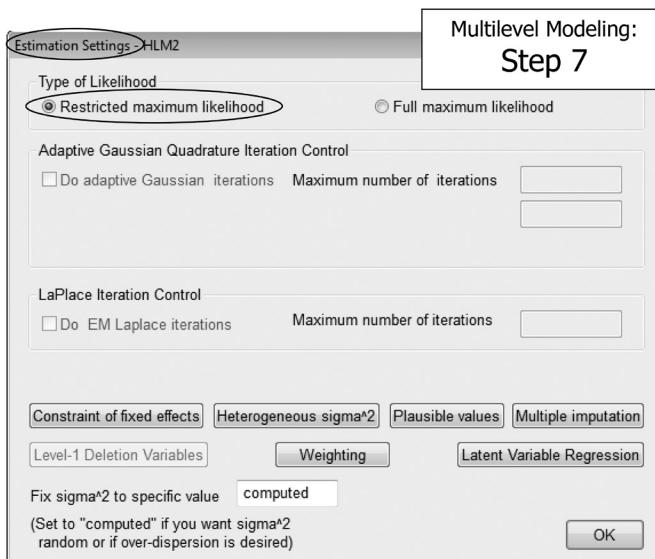
Step 5b. Clicking on Level-2 Residual File will bring up the Create Level-2 Residual File dialog box (see screenshot Step 5b). We want to include in our residual file the variables that have been included in level 2 of our model, and this is done by listing the respective variables in the boxes labeled 'Variables in residual file.' The default file type is SPSS, and a default file name is provided. For purposes of illustration, we will leave those as is. Click "OK" to return to the basic setting dialog box.



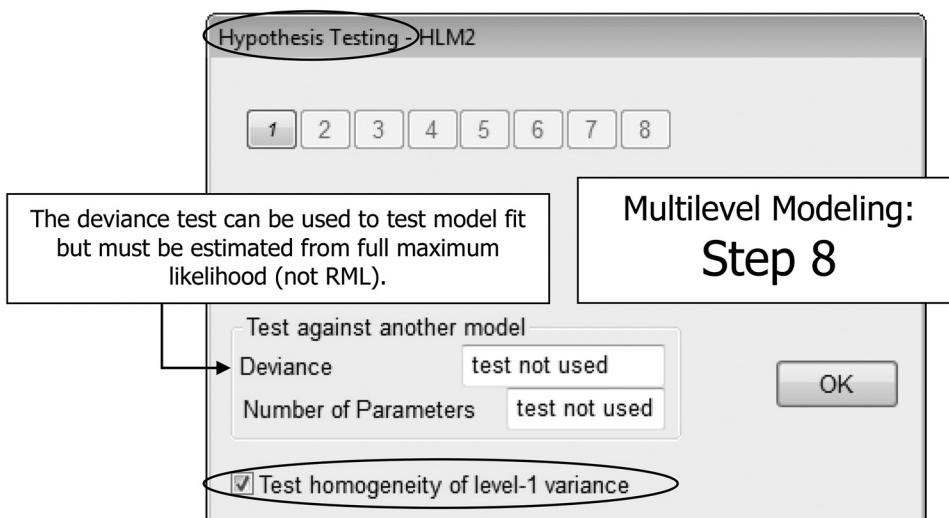
Step 6. Prior to running the model, you also want to change, as needed, the number of iterations. This can be accessed from the top menu seen in Step 4b. Clicking on “Other Settings” will bring up the iteration specification screen (see screenshot Step 6). From here, the distribution of the outcome variable can be defined. The default is ‘100,’ and we will change this to 1000. The default action when the maximum number of iterations is achieved is to ‘Prompt,’ and we will leave this setting as is. Click “OK” to return to the main dialog box (see screenshot Step 4).



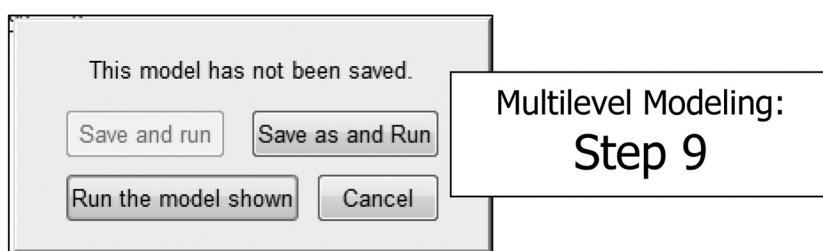
Step 7. Prior to running the model, you also want to change, as needed, the estimation method. This can be accessed from the top menu seen in Step 4b. Clicking on “Other Settings” will bring up the estimation settings screen (see screenshot Step 7). The default is ‘full maximum likelihood,’ and we will change this to ‘Restricted maximum likelihood’ by clicking the radio button. Recall the RML produces more accurate standard errors when there are fewer level 2 groups (as is the case). From this screen, you can also select weights when complex survey data has been examined, deal with multiple imputed and plausible value data, and more. Clicking “OK” will return to the main dialog box (see screenshot Step 4).



Step 8. Clicking on “Other Settings” from the top menu seen in Step 4a will also bring up the hypothesis testing screen (see screenshot Step 8). There are no hypothesis tests that are selected by default, and we will place a check to select ‘Test homogeneity of level 1 variance.’ The test of deviance is also conducted within this screen. Remember that the deviance test is a model fit test that compares one less parameterized model to a more parameterized model. Should the deviance test be desired, the model will need to be generated using full maximum likelihood, as restricted maximum likelihood allows comparison only of random effects. Clicking “OK” will return to the main dialog box (see screenshot Step 4a).



Step 9. From the main dialog page (see screenshot Step 4a), we are now ready to run our model. To do so, click on “Run Analysis.” The model can be run without saving (i.e., ‘Run the model shown’) or simultaneously saved and run (‘Save as and run’) (see screenshot Step 9). Once the model is run, the output will automatically open as an html file.



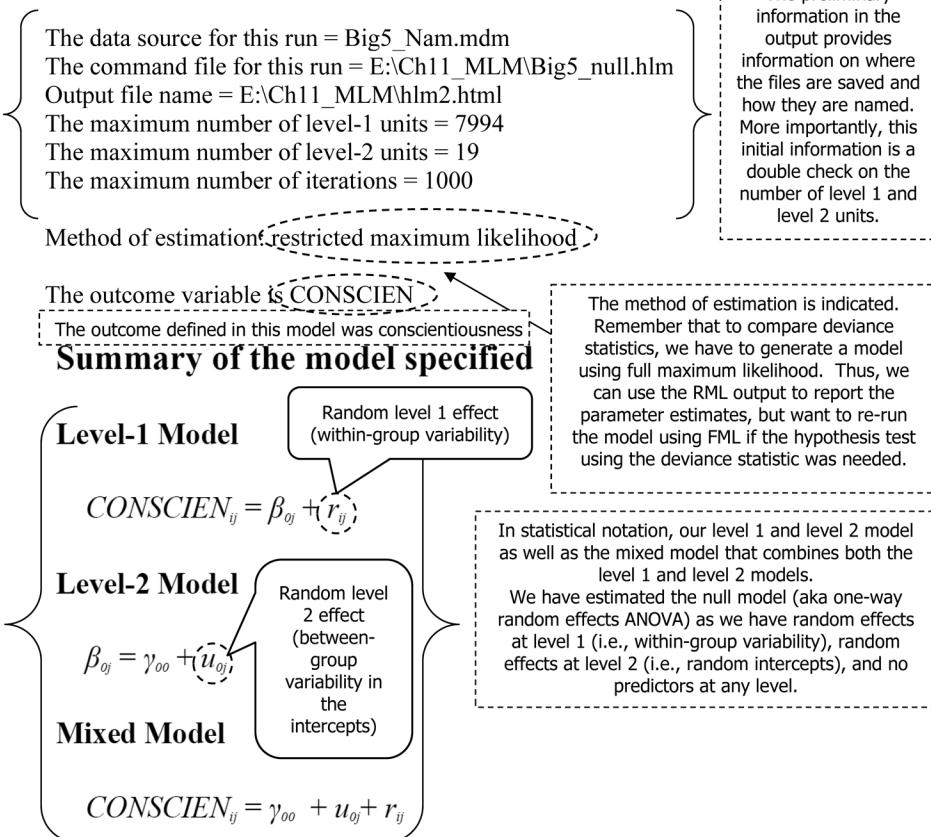
Interpretations of the one-way random effects ANOVA are provided in Table 11.1.

TABLE 11.1

HLM Results for the One-Way Random Effects Anova Model (i.e., the Null Model)

Specifications for this HLM2 run

Problem Title: no title



Final Results - Iteration 11

Iterations stopped due to small change in likelihood function

$\sigma^2 = 0.54732$	Within-group variability. This is the variance in the outcome that is within-groups or clusters.
τ INTRCPT1, β_0 0.01604	Between-group variability in the intercepts. This is the variance in the level 2 group means around the grand mean (gamma 00).
Random level-1 coefficient INTRCPT1, β_0	Reliability estimate 0.261

The value of the log-likelihood function at iteration

This reliability is a measure of the overall reliability of the OLS estimates for each of the intercepts.

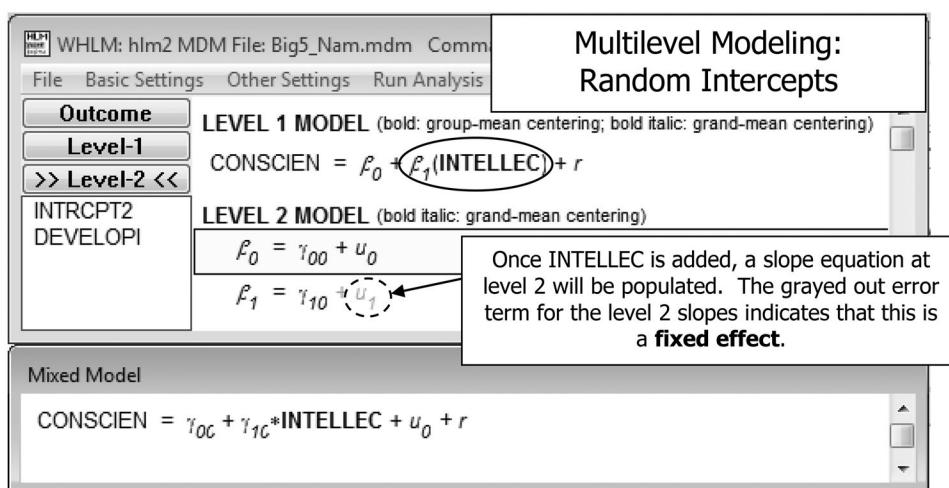
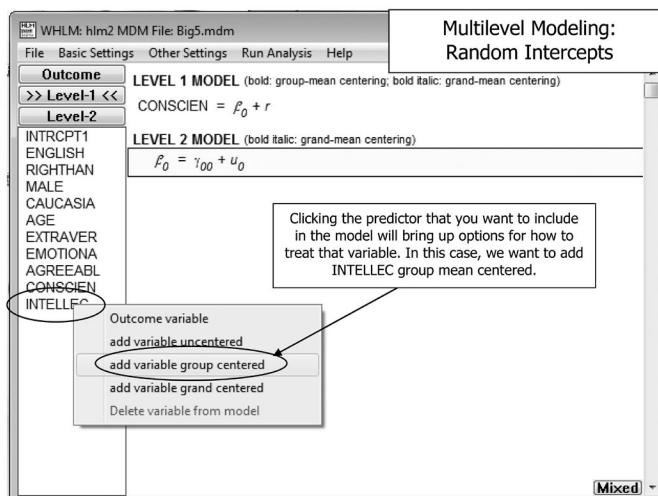
TABLE 11.1 (continued)

HLM Results for the One-Way Random Effects Anova Model (i.e., the Null Model)

Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0					
INTRCPT2, γ_{00}	3.390355	0.056886	59.600	18	<0.001
The average intercept across the level 2 groups, i.e., the grand mean CONSCIEN score, is about 3.40 and statistically significantly different than zero.					
Final estimation of fixed effects (with robust standard errors)					
Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0					
INTRCPT2, γ_{00}	3.390355	0.053117	63.828	18	<0.001
A large number of level 2 groups are needed to produce accurate robust standard errors. Should robust and asymptotic standard errors differ substantially, that may be a signal that MLM assumptions have been violated.					
The robust standard errors are appropriate for datasets having a moderate to large number of level 2 units. These data do not meet this criterion.					
Final estimation of variance components					
Random Effect	Standard Deviation	Variance Component	d.f.	χ^2	p-value
INTRCPT1, u_0	0.12667	0.01604	18	45.14910	<0.001
level-1, r	0.73981	0.54732			
Statistics for current covariance components model					
Deviance = 17881.452243 Number of estimated parameters = 2			This is the hypothesis test for the between-group variability. It answers the question: is the estimated value of tau (i.e., the between-group variability) significantly different from zero? This chi-squared test with $j - 1$ degrees of freedom (where j = the number of level 2 groups) is statistically significant. This tells us that there is statistically significant variation between states in conscientiousness.		
Do these variance values sound familiar? They should! This is sigma squared and tau that were provided earlier in the output.					
From these variance estimates, we can compute the ICC (the proportion of variation in the outcome that is between groups). In this example, about 3% of the variation in CONSCIEN is between states:					
$\text{ICC} = \frac{\tau_{00}}{\sigma^2 + \tau_{00}} = \frac{.016}{.547 + .016} = .025$					

11.3.1 Random Intercepts Model

The model just estimated was the one-way random effects ANOVA. It is primarily used to estimate the proportion of variance in the outcome that is between groups. Based on the ICC in the model we've estimated, it is appropriate to proceed with multilevel model. Adding to and changing the model from this point forward is quite simple. For the random intercepts model, we will start by including one level 1 predictor, intellect/imagination [INTELLEC], that is group mean centered. At level 2, we see that an equation for the slopes has automatically populated, and the random slopes effect term, u_1 , is grayed out. This means that the slopes are fixed. For this example, we leave the slopes as fixed effects. This will estimate the model with random intercepts only. Next, from the main dialog page (see screenshot Step 4a), we are now ready to run our model. To do so, click on "Run Analysis." The model can be run without saving (i.e., 'Run the model shown') or simultaneously saved and run ('Save as and run') (see screenshot Step 9). Once the model is run, the output will automatically open as an hml file.



Interpretations of the random intercepts model are provided in Table 11.2.

TABLE 11.2

HLM Results for the Random Intercepts (Select Results)

TABLE 11.2
HLM Results for the Random Intercepts (Select Results)

Summary of the model specified					
Level-1 Model	Random level 1 effect (within-group variability)				
$CONSCIEN_y = \beta_{0j} + \beta_{1j} * (INTELLEC_y) + r_y$	In statistical notation, our level 1 and level 2 model as well as the mixed model that combines both the level 1 and level 2 models.				
Level-2 Model	We have estimated a random intercept model as we have random effects at level 1 (i.e., within-group variability), random effects at level 2 (i.e., random intercepts), and one level 1 predictor with fixed slopes at level 2.				
$\beta_{0j} = \gamma_{00} + u_{0j}$	Random level 2 effect (between-group variability in the intercepts)				
$\beta_{1j} = \gamma_{10}$					
INTELLEC has been centered around the group mean.					
Mixed Model	The level 1 predictor is centered around the group mean.				
$CONSCIEN_y = \gamma_{00}$ + $\gamma_{10} * INTELLEC_y + u_{0j} + r_y$	Within-group variability. This is the variance in the outcome that is within-groups or clusters.				
$\sigma^2 = 0.54583$					
τ	Between-group variability in the intercepts. This is the variance in the level 2 group means around the grand mean (gamma 00).				
INTRCPT1, β_0	0.01623				
Random level-1 coefficient	Reliability estimate				
INTRCPT1, β_0	0.263				
The value of the log-likelihood function at iteration 11 = -8.933706E+003					
Final estimation of fixed effects:					
Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0	INTRCPT2, γ_{00}	3.390350	0.057058	59.419	18 <0.001
For INTELLEC slope, β_1	INTRCPT2, γ_{10}	0.063555	0.013338	4.765	7974 <0.001
For every one point increase in INTELLEC, CONSCIEN increases by about .06 points. This is a statistically significant increase.					

TABLE 11.2 (continued)

HLM Results for the Random Intercepts (Select Results)

Random Effect	Standard Deviation	Variance Component	d.f.	χ^2	p-value
INTRCPT1, u_0	0.12741	0.01623	18	45.27468	<0.001
level-1, r	0.73881	0.54583			

Do these variance values sound familiar? They should! This is sigma squared and tau that were provided earlier in the output.

From these variance estimates, we can compute the ICC (the proportion of variation in the outcome that is between groups). In this example, about 2% of the variation in CONSCIEN is between states:

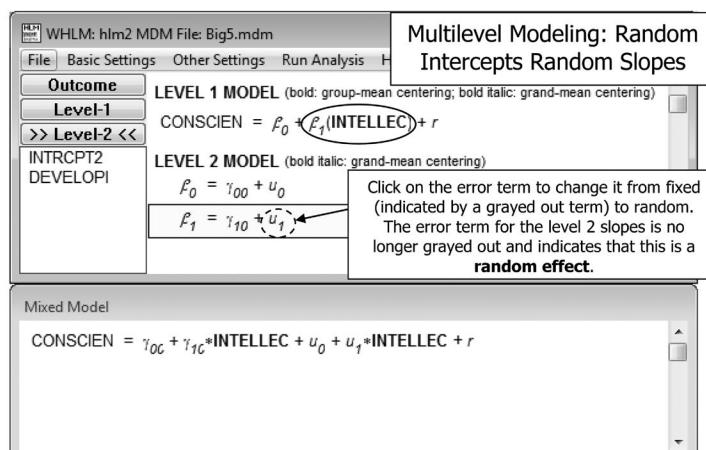
$$\text{ICC} = \frac{\tau_{00}}{\sigma^2 + \tau_{00}} = \frac{.016}{.546 + .016} = .024$$

This is the hypothesis test for the between-group variability. It answers the question: is the estimated value of tau (i.e., the between-group variability) significantly different from zero?

This chi-square test with $j - 1$ degrees of freedom (where j = the number of level 2 groups) is statistically significant. This tells us that there is statistically significant variation between states in conscientiousness.

11.3.2 Random Intercepts Random Slopes Model

For the random intercepts random slopes model, we will simply click the random slopes effect so that the random effect is no longer grayed out. This will estimate the model with both random intercepts and random slopes. Next, from the main dialog page (see screenshot Step 4a), we are now ready to run our model. To do so, click on "Run Analysis." The model can be run without saving (i.e., 'Run the model shown') or simultaneously saved and run ('Save as and run') (see screenshot Step 9). Once the model is run, the output will automatically open as an hml file.



Interpretations of the random intercepts random slopes model are provided in Table 11.3.

TABLE 11.3

HLM Results for the Random Intercepts and Random Slopes Model (Select Output)

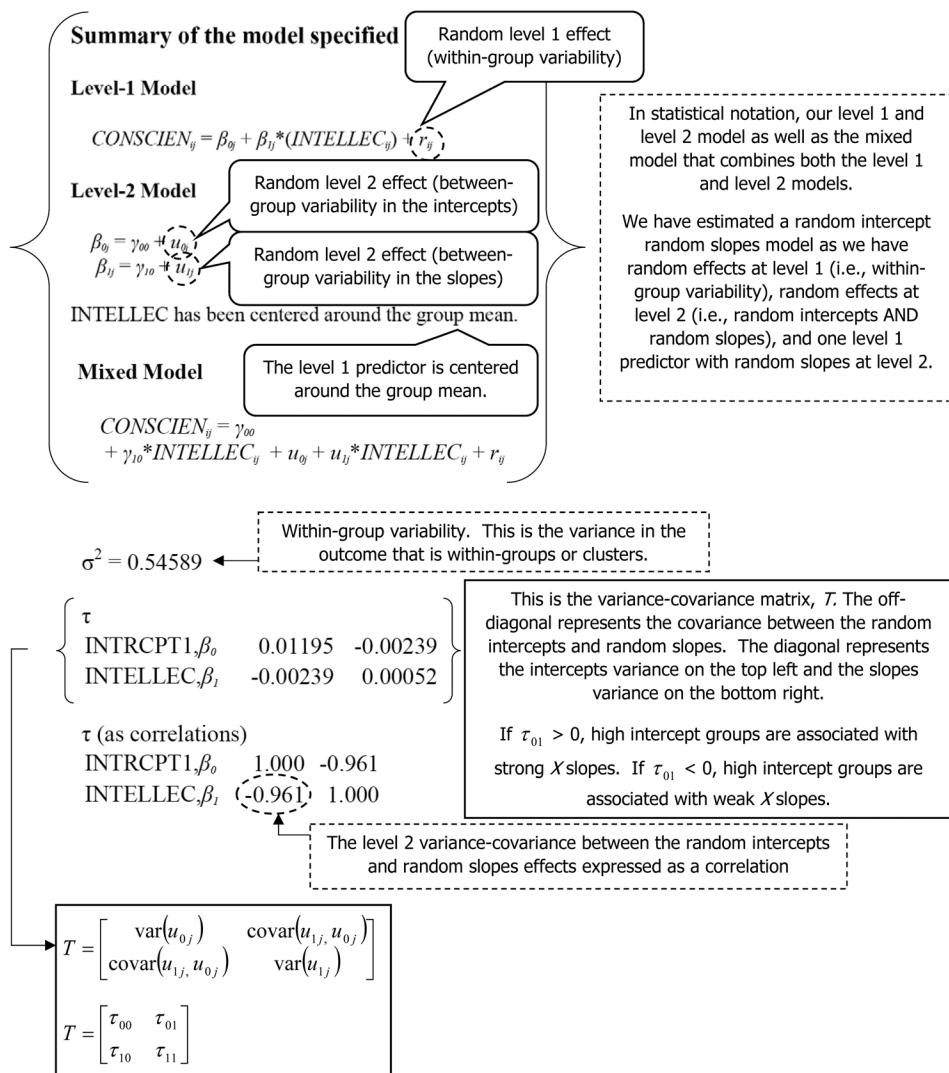


TABLE 11.3 (continued)

HLM Results for the Random Intercepts and Random Slopes Model (Select Output)

Random level-1 coefficient	Reliability estimate
INTRCPT1, β_0	0.272
INTELLEC, β_1	0.064

Slopes with reliability less than .10 may signal the need to fix slopes (rather than leave random). These slopes may not be helpful in distinguishing groups based on their relationship between X and the outcome.

Note: The reliability estimates reported above are based on only 16 of 19 units that had sufficient data for computation. Fixed effects and variance components are based on all the data.

The average intercept across the level 2 groups, i.e., the grand mean CONSCIEN score, is about 3.40 and statistically significantly different than zero. Since the level 1 predictor is group mean centered, we interpret this as the average CONSCIEN score for a person who is at the average INTELLEC for their group. The confidence interval can be computed as

$$\gamma_{00} \pm 1.96/\sqrt{\tau_{00}} = 3.39 \pm 1.96\sqrt{.01} = (3.19, 3.59)$$

Final estimation of fixed effects:

Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0					
INTRCPT2, γ_{00}	(3.391489)	0.051835	65.428	18	(<0.001)
For INTELLEC slope, β_1					
INTRCPT2, γ_{10}	(0.067949)	0.017674	3.845	18	(0.001)

For every one point increase in INTELLEC, CONSCIEN increases by about .07 points. This is a statistically significant increase.

Between group variability in the intercepts (u_{0j}) and slopes (u_{1j}).

The non-statistically significant random effect for the slopes tells us that between groups, the relationship between INTELLEC and CONSCIEN (i.e., the change in CONSCIEN as a function of INTELLEC) does not vary significantly. This means that we can fix the slope effect rather than run them as randomly varying.

Final estimation of variance components

Random Effect	Standard Deviation	Variance Component	d.f.	χ^2	p-value
INTRCPT1, u_0	0.10932	0.01195	15	39.60007	(<0.001)
INTELLEC slope, u_1	0.02275	0.00052	15	21.03868	(0.135)
level-1, r	0.73884	0.54589			

Note: The chi-square statistics reported above are based on only 16 of 19 units that had sufficient data for computation. Fixed effects and variance components are based on all the data.

This message may appear for univariate chi square tests, reliabilities, and least-squares estimates. The main results (fixed and variance-covariance components) use all data. If you receive this message, it means that the least squares estimator did not exist for the groups that were dropped (in this case, 3 countries were not included). This may have been caused by failing to achieve invertibility or positive determinant and is often caused by a small number of cases in the group.

TABLE 11.3 (continued)

HLM Results for the Random Intercepts and Random Slopes Model (Select Output)

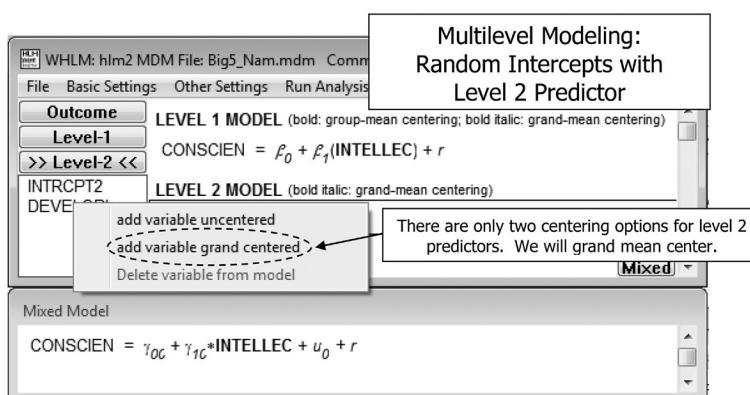
Test of homogeneity of level-1 variance

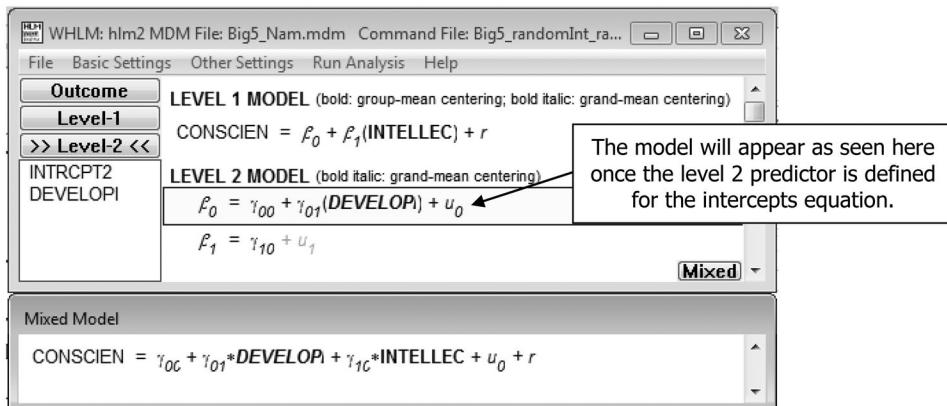
χ^2 statistic = 7.66133
 degrees of freedom = 11
 p -value = >.500

The test of level 1 variance tells us that the level 1 variances in this model have equal variances (or similar enough that they are not statistically significantly different from each other). When there is heterogeneity, this may be a sign that one or more important level 1 covariates have been omitted.

11.3.3 Random Intercepts Model With Level 2 Predictor

Adding to and changing the level 2 model from this point forward is also quite simple. The random intercepts random slopes results suggest that fixing the slopes are appropriate (recall the random slopes effect was not statistically significant). Thus, we'll fix the slopes (i.e., remove the random effects term for the slopes equation), and we will add to the model by including one level 2 predictor, DEVELOP, that is grand mean centered. At level 2, we can add level 2 variables in any one or more level 2 equations that we choose, and we do this by clicking on the name of the level 2 variable of interest in the left list of variables. In this illustration, we will include DEVELOP as a predictor of the intercepts only since it was determined in the previous model that fixing the slopes is appropriate. In other words, there is not variation between countries in the slopes that requires level 2 predictors to be included to explain that between-group slope variation (should there be theoretical justification, however, we could defend including group-level predictors, and we will illustrate this later). Including the level 2 predictor to explain the intercepts will allow us to determine if the country-level predictor, whether or not the country is a developing country [DEVELOP], can assist in explaining the variation in intercepts at level 2. Once the variable name is clicked, the options for centering appear. For this illustration, we will grand mean center. Next, from the main dialog page (see screenshot Step 4a), we are now ready to run our model. To do so, click on "Run Analysis." The model can be run without saving (i.e., 'Run the model shown') or simultaneously saved and run ('Save as and run') (see screenshot Step 9). Once the model is run, the output will automatically open as an hml file.

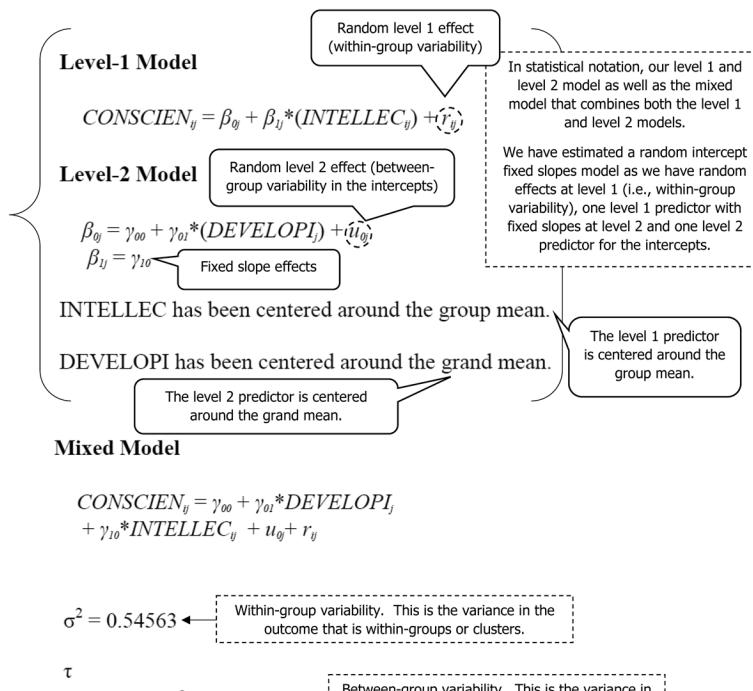




Interpretations of the random intercepts model that includes a group level covariate are provided in Table 11.4.

TABLE 11.4

HLM Results for the Random Intercepts and Fixed Slopes Model with Level 2 Predictor (Select Output)



Random level-1 coefficient	Reliability estimate
INTRCPT1, β_{0y}	0.352

TABLE 11.4 (continued)

HLM Results for the Random Intercepts and Fixed Slopes Model with Level 2 Predictor (Select Output)

Final estimation of fixed effects					
Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0					
INTRCPT2, γ_{00}	3.377134	0.073616	45.875	17	<0.001
DEVELOPI, γ_{01}	-0.090372	0.145196	-0.622	17	0.542
For INTELLEC slope, β_1					
INTRCPT2, γ_{10}	0.063555	0.013335	4.766	7974	<0.001

For every one point increase in INTELLEC, CONSCIEN increases by about .06 points. This is a statistically significant increase.

Non significant level 2 predictor suggests that whether or not the person is in a developing country doesn't relate to CONSCIEN. Unless there is theoretical justification to leave this predictor in the model, it could be removed.

Final estimation of variance components

Random Effect	Standard Deviation	Variance Component	d.f.	χ^2	p-value
INTRCPT1, u_0	0.18555	0.03443	17	47.77295	<0.001
level-1, r	0.73867	0.54563			

Test of homogeneity of level-1 variance

$$\chi^2 \text{ statistic} = 7.42112$$

degrees of freedom = 16

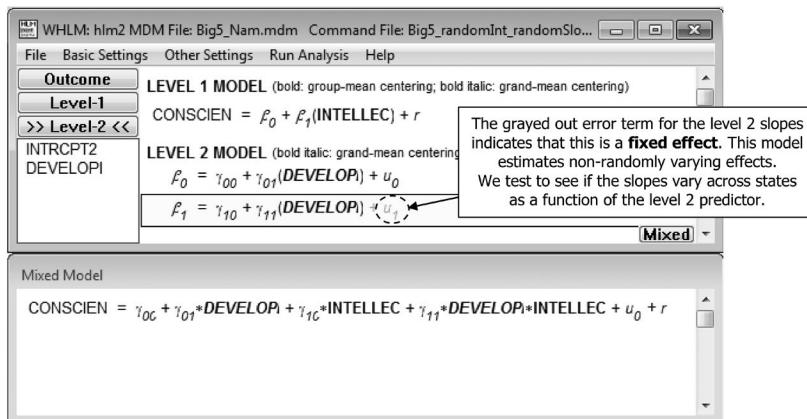
p-value = >.500

There is still statistically significant variation in the intercepts. This suggests that additional level 1 predictors are needed to continue to explain this variation.

11.3.4 Random Intercepts Model With Nonrandomly Varying Slopes Effects

In this model, we will include the developing status of the country [DEVELOP] as a level 2 predictor and estimate nonrandomly varying slopes effects. We do this for illustration only, as the previous model (Table 11.4) suggested that whether or not the country was a developing country was not a statistically significant level 2 predictor, and thus excluding it from the model is justified. Again, for the sake of illustration, we will continue. Additionally, it may be that there is sufficient theory to support adding a level 2 predictor when there are not random effects, and thus we will illustrate this process. We add the level 2 predictor by clicking on the name of the level 2 variable of interest in the left list of variables. In this illustration, we will include the developing status of the country as a predictor of the slopes. Including the level 2 predictor to explain the intercepts

will allow us to determine if the country-level predictor, developing status, can assist in explaining the slopes. In the case of nonrandomly varying, we are saying that there is variation between countries, but that the variation is not random—rather, it is explained by whether or not the country is a developing country. Once the variable name is clicked, the options for centering appear. For this illustration, we will grand mean center. Next, from the main dialog page (see screenshot Step 4a), we are now ready to run our model. To do so, click on “Run Analysis.” The model can be run without saving (i.e., ‘Run the model shown’) or simultaneously saved and run (‘Save as and run’) (see screenshot Step 9). Once the model is run, the output will automatically open as an hml file.



Interpretations of the random intercepts model that includes slope effects that are non-randomly varying are provided in Table 11.5.

TABLE 11.5

HLM Results for the Random Intercepts and Fixed Slopes Model with Level 2 Predictor for Intercepts and Slopes (Select Output)

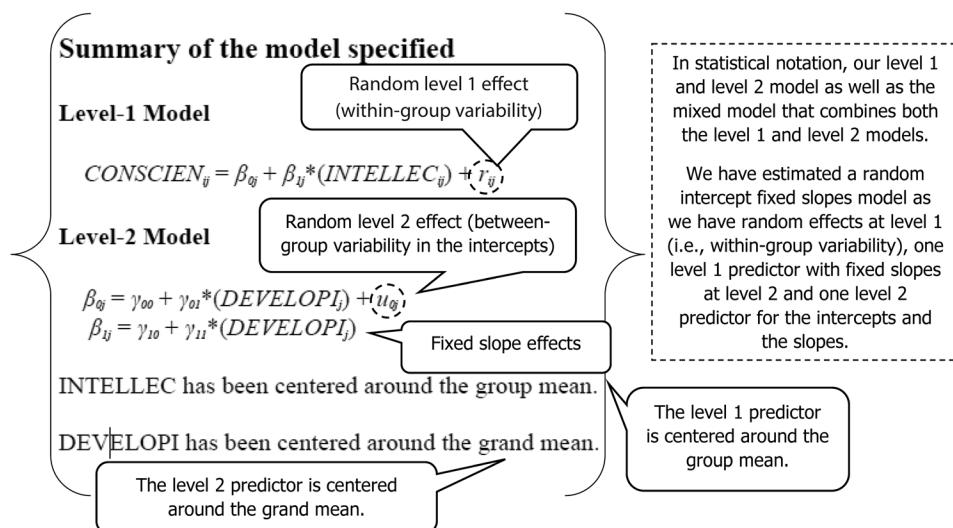


TABLE 11.5 (continued)

HLM Results for the Random Intercepts and Fixed Slopes Model with Level 2 Predictor for Intercepts and Slopes (Select Output)

Mixed Model

$$CONSCIEN_{ij} = \gamma_{00} + \gamma_{01}*DEVELOPI_j \\ + \gamma_{10}*INTELLEC_{ij} + \gamma_{11}*DEVELOPI_j*INTELLEC_{ij} \\ + u_{0j} + r_{ij}$$

$\sigma^2 = 0.54562$	Within-group variability. This is the variance in the outcome that is within-groups or clusters.
τ INTRCPT1, β_0 0.03443	Between-group variability. This is the variance in the outcome that is between-groups or clusters.

Random level-1 coefficient	Reliability estimate
INTRCPT1, β_0	0.352

The average intercept across the level 2 groups, i.e., the grand mean CONSCIEN score, is about 3.38 and statistically significantly different than zero. This is the average CONSCIEN score for a person who is at the average INTELLEC for their group (since the level 1 predictor is group mean centered, we interpret) who is not in a developing country (since the level 2 predictor was uncentered). The confidence interval can be computed as

$$\gamma_{00} \pm 1.96\sqrt{\tau_{00}} = 3.38 \pm 1.96\sqrt{0.03} = (3.04, 3.72)$$

Final estimation of fixed effects:

Fixed Effect	Coefficient	Standard error	t-ratio	Approx. d.f.	p-value
For INTRCPT1, β_0					
INTRCPT2, γ_{00}	3.377133	0.073618	45.874	17	<0.001
DEVELOPI, γ_{01}	-0.090379	0.145199	-0.622	17	0.542
For INTELLEC slope, β_1					
INTRCPT2, γ_{10}	0.012082	0.072342	-0.167	7973	0.867
DEVELOPI, γ_{11}	-0.113187	0.106401	-1.064	7973	0.287

After controlling the slopes for whether or not the country is developing, INTELLEC is no longer statistically significant.

Non significant level 2 predictor suggests that whether or not the person is in a developing country doesn't relate to CONSCIEN. Unless there is theoretical justification to leave this predictor in the model, it could be removed.

Final estimation of variance components

Random Effect	Standard Deviation	Variance Component	d.f.	χ^2	p-value
INTRCPT1, u_0	0.18555	0.03443	17	47.77668	<0.001
level-1, r	0.73866	0.54562			

There is still statistically significant variation in the intercepts. This suggests that additional level 1 predictors are needed to continue to explain this variation.

Test of homogeneity of level-1 variance

χ^2 statistic = 7.21676
 degrees of freedom = 16
 p-value = >.500

11.4 DATA SCREENING

As stated previously, multilevel linear modeling assumptions include (1) linearity (i.e., linear relationship between variables), (2) normally distributed level residuals with equal variances, (3) multivariately normally distributed residuals at level 2 that are homoscedastic (i.e., residual variance is constant—it does not relate to the predictor variables), and (4) predictors at a respective level are uncorrelated with random effects at the other level. We will test the assumptions using the level 1 and level 2 residuals in addition to other variables that were saved when requesting residuals.

11.4.1 Level 1 Residuals

We have already demonstrated how we can test for equal variances of level 1 residuals using HLM with the test of homogeneity of variances. Using the level 1 residual file, we can test for normality (see screenshot Step 5 for saving residuals). For this example, we will use the residual file from the random intercepts model with the level 2 predictor (Table 11.4). More specifically, we will use the second variable, L1RESID, and generate basic evidence such as skewness and kurtosis, Q-Q plots, formal test of normality (e.g., Shapiro-Wilk), histogram, and boxplots. Doing this, we find skewness and kurtosis that are within the range of normal (−.149 and −.426, respectively). However, the Kolmogorov-Smirnov formal test of normality suggests a nonnormal distribution ($K-S = .020$, $df = 7994$, $p < .001$). The Q-Q plots and histogram suggest some nonnormality in the tails of the distribution (left tail in particular). The boxplot suggests potential outliers on the lower end of the scale. We can assess linearity by reviewing scatterplots of level 1 residuals and fitted values (i.e., predicted values). A random display of points suggests linearity.

The variables in the level 1 residual file include:

1. level 2 identification number
2. level 1 residuals (i.e., discrepancies between the observed and fitted values)
3. fitted values (FV) for each level 1 unit (i.e., values predicted on the basis of the model)
4. square root of sigma squared (i.e., the standard deviation)
5. value of CONSCIEN in the data file
6. value of INTELLEC in the data file
7. value of DEVELOPI in the data file

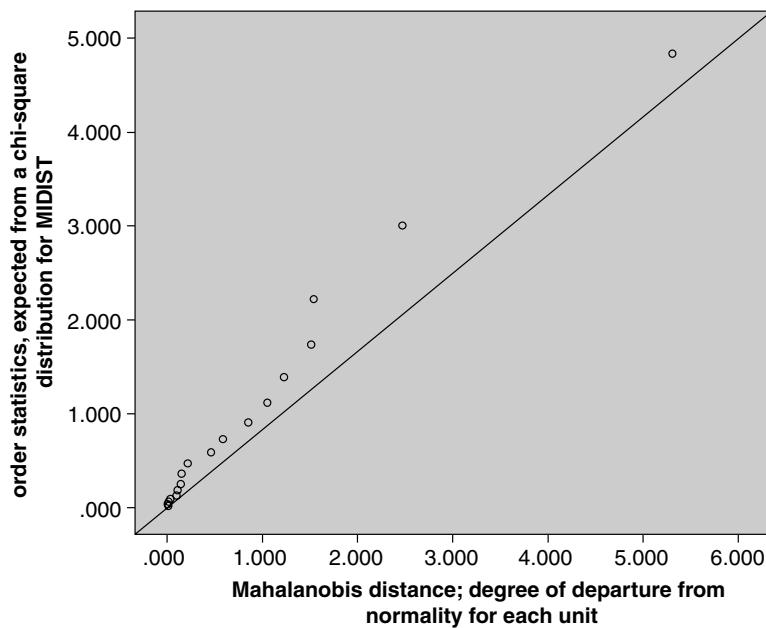
11.4.2 Level 2 Residuals

There is one row of level 2 residuals for each group in the analyses. Thus, this illustration has 19 rows of data in the level 2 residuals, one representing each country. The explanation of each variable is provided in the screenshot that follows.

	Name	Type	Width	Decimals	Label
1	L2ID	Numeric	8	0	level 2 ID
2	NJ	Numeric	8	0	within group sample size (i.e., sample size for a particular school)
3	CHIPCT	Numeric	12	3	order statistics, expected from a chi-square distribution for MDIST
4	MDIST	Numeric	12	3	Mahalanobis distance; degree of departure from normality for each ...
5	LNTOTVAR	Numeric	12	3	natural log of the total SD within each group (i.e., school)
6	OLSRSVAR	Numeric	12	3	natural log of the OLS residual SD within each group
7	MDRSVAR	Numeric	12	3	natural log of the final MODEL residual SD within each unit
8	EBINTRCPT1	Numeric	12	3	Empirical Bayes residuals of the intercept
9	OLINTRCPT1	Numeric	12	3	OLS residuals, intercept
10	FVINTRCPT1	Numeric	12	3	Fitted values of the intercept for this model, each unit
11	FVINTELLEC	Numeric	12	3	Fitted values of the slope for this model, each unit
12	ECINTRCPT1	Numeric	12	3	Empirical Bayes coefficients for the intercepts
13	ECINTELLEC	Numeric	12	3	Empirical Bayes coefficients for the slopes
14	PV0_0	Numeric	12	3	posterior (ending values) for variance for the estimate of the intercepts
15	PVC0_0	Numeric	12	3	posterior covariance
16	DEVELOPI	Numeric	12	3	variable included in the level 2 model

Examination of linearity at level 2 can be done by plotting the empirical Bayes residuals of the slope to the predictor variable included in the level 2 model [DEVLOPI]. A random distribution of points around zero suggest linearity is reasonable. In this illustration, the model we've generated does not provide for the computation of empirical Bayes residuals for the slopes, thus examination of linearity at level 2 is not required.

Normality of the level 2 residuals can be examined via a scatterplot of Mahalanobis distance (MDIST) and the expected values of the order statistics (CHIPCT). Points that generally adhere to the diagonal line show evidence of multivariate normally distributed data. Groups that have large Mahalanobis distance values, i.e., substantially distanced from the reference line, are also potential outliers. In this example, there are several groups that are suggestive of outliers, and the points are not adhering as closely to the diagonal as multivariate normal would suggest. However, these plots should only be considered accurate when the level 1 sample sizes within groups are at least moderately large (Raudenbush, Bryk, Cheong, Congdon, & du Toit, 2011)



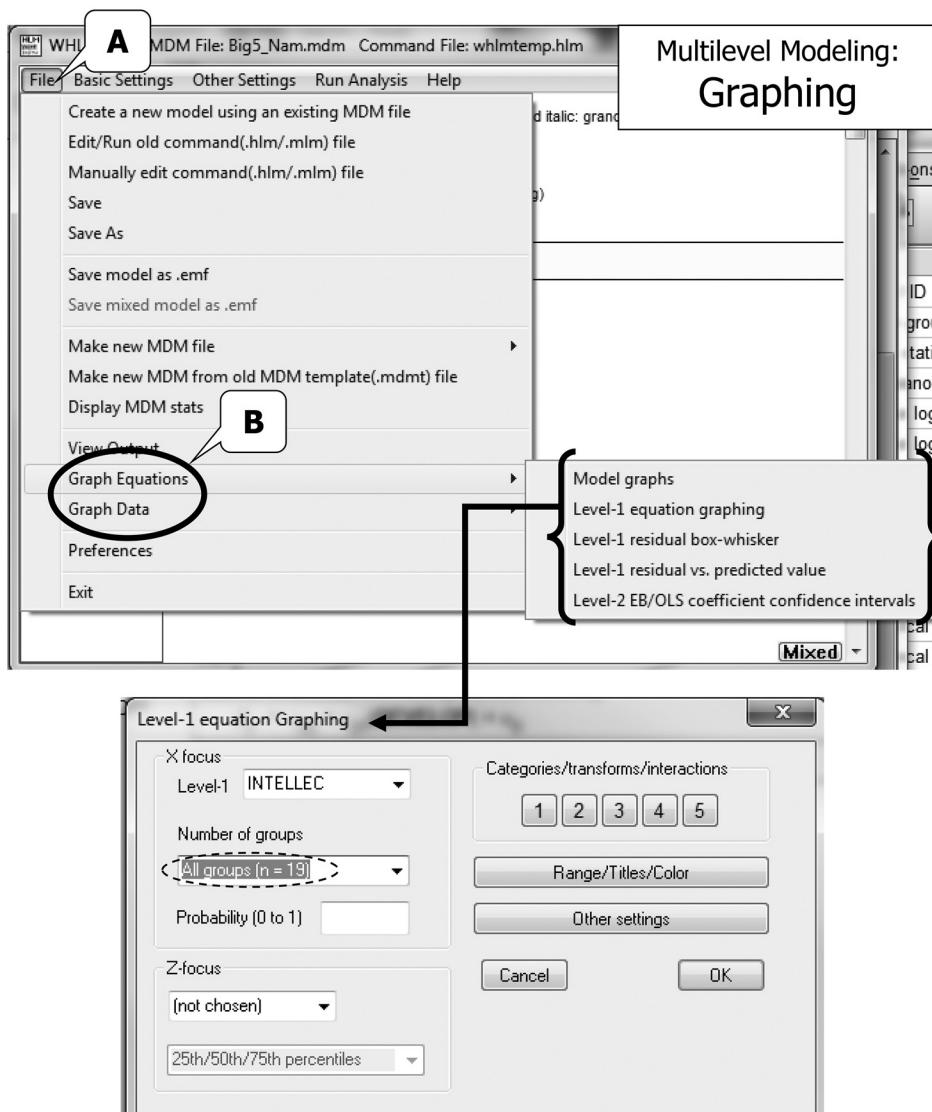
Heteroscedasticity can lead to model misspecification, and evidence of this can be examined by inclusion of nonlinear and interaction effects.

Random slopes were not included in this model. Had they been, the residual file would have been expanded to include additional residuals as we see in the following, and additional data screening can be performed.

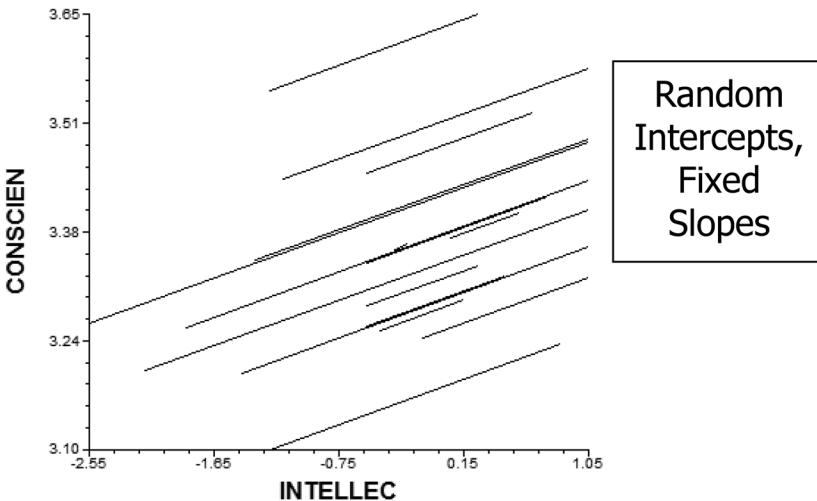
	Name	Type	Width	Decimals	Label
1	L2ID	Numeric	8	0	level 2 ID
2	NJ	Numeric	8	0	within group sample size (i.e., sample size for a particular school)
3	CHIPCT	Numeric	12	3	order statistics, expected from a chi-square distribution for MDIST
4	MDIST	Numeric	12	3	Mahalanobis distance; degree of departure from normality for each unit
5	LNTOTVAR	Numeric	12	3	natural log of the total SD within each group (i.e., school)
6	OLSRSVAR	Numeric	12	3	natural log of the OLS residual SD within each group
7	MDRSVAR	Numeric	12	3	natural log of the final MODEL residual SD within each unit
8	EBINTRCPT1	Numeric	12	3	Empirical Bayes residuals of the intercept
9	EBINTELLEC	Numeric	12	3	Empirical Bayes residuals of the slope for INTELLEC (for this particular model)
10	OLINTRCPT1	Numeric	12	3	OLS residuals, intercept
11	OLINTELLEC	Numeric	12	3	OLS residuals, slope for INTELLEC
12	FVINTRCPT1	Numeric	12	3	Fitted values of the intercept for this model, each unit
13	FVINTELLEC	Numeric	12	3	Fitted values of the slope for this model, each unit
14	ECINTRCPT1	Numeric	12	3	Empirical Bayes coefficients for the intercepts
15	ECINTELLEC	Numeric	12	3	Empirical Bayes coefficients for the slopes
16	PV0_0	Numeric	12	3	posterior (ending values) for variance for the estimate of the intercepts
17	PV1_0	Numeric	12	3	posterior covariance for estimates of the intercepts and slopes
18	PV1_1	Numeric	12	3	posterior variances for the slopes
19	PVC0_0	Numeric	12	3	posterior covariance
20	PVC1_0	Numeric	12	3	posterior covariance
21	PVC1_1	Numeric	12	3	posterior covariance
22	DEVELOPI	Numeric	12	3	variable included in the level 2 model
23					

11.4.3 Graphing

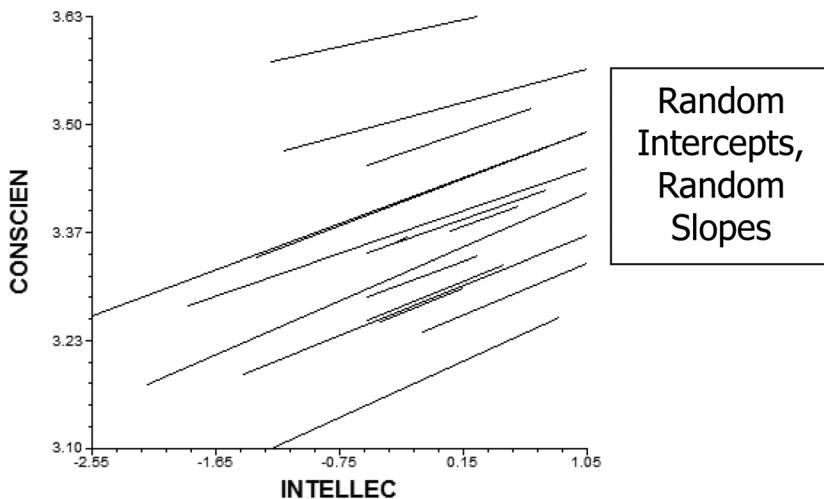
There are a number of excellent graphing options in HLM, and we encourage you to explore these to assist in visually describing and understanding your data. For brevity, only a few will be illustrated; however, all the graphing options can be accessed from either Graph Equations or Graph Data options listed from the File menu. For example, to view the graphs of the countries in our data file, based on random intercepts and fixed slopes, we go to "Graph Equations" then "Level 1 equation graphing." In this illustration, we have a manageable number of level 2 groups, so creating the graph with 'all groups' is reasonable. You may find, however, that looking at a random percentage of groups provides better interpretation. Also keep in mind that the graph that is generated will be based on the *last model run*. Thus, should you change any fixed or random effects, you will want to rerun the graph after doing so.



Given that we had fixed slopes, we see that all countries have common slopes, but the random effects for the intercepts are apparent in the CONSCIEN scores by group.

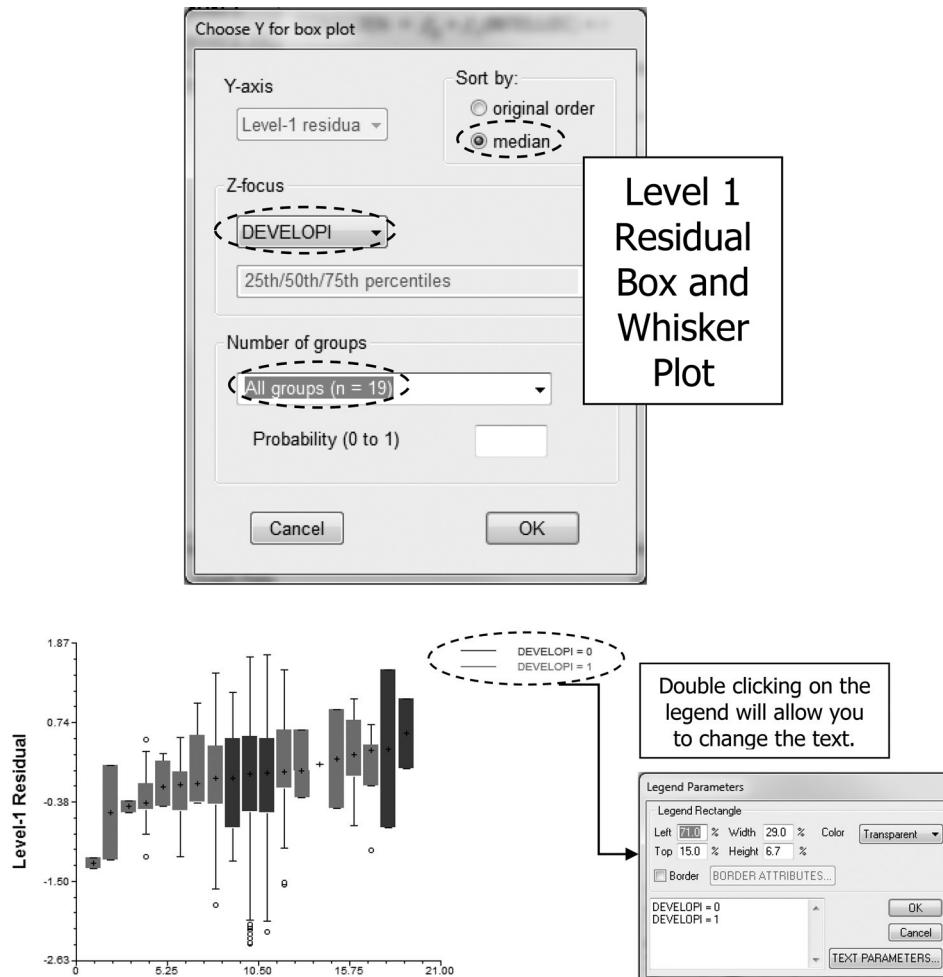


Plotting the same model but with *random slopes* produces the following graph. It is easy to see why the slopes were not statistically significant. The groups have very similar slopes as evidenced by nearly parallel lines.



Another useful graph for examining assumptions is the Level 1 residual box and whisker. In this, we graph the level 1 residuals by group and include our level 2 predictor for the Z focus. This delineates the two groups in our analyses. Sorting by median

helps to display patterns (e.g., developing countries, displayed as 1, tend to have lower residuals, i.e., lower observed to fitted values).



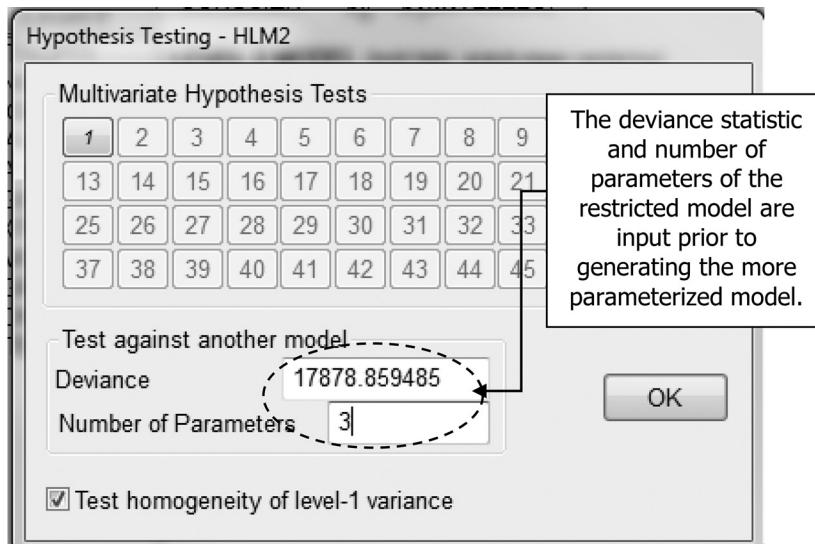
11.4.4 Model Fit

As indicated previously, there are a number of model fit indices that can be applied as multiparameter tests to gauge changes in model fit when adding and/or removing predictors and/or fixed and/or random effects. Using the Big Five data, we will look at the difference in deviances likelihood ratio test.

11.4.4.1 Difference in Deviances Likelihood Ratio Test

For this example, we will compare the Big Five null model, which estimated conscientiousness as the outcome, to a more parameterized model that includes intellect/imagination as a level 1 predictor with random intercepts and fixed slopes (see Table 11.2 for notation; however, Table 11.2 was generated using RML). *To compute the test, we must estimate both models using FML.* That output is not provided here, but the deviance value of the reduced model that is input here was based on FML estimation. The FML deviance and number of parameters estimated from the null model (or the

less parameterized model) are input in Hypothesis Testing, which is accessible from Other Settings in the top toolbar from the main HLM dialog page (see screenshot Step 4a). Then, the more parameterized model is generated.



This will generate the hypothesis test results similar to what is seen here. In this illustration, the more parameterized model that includes intellect/imagination as a predictor at level 1 is statistically significantly a better-fitting model than the restricted model (which was the null model).

Deviance = 17856.198242

Number of estimated parameters = 4

Model comparison test

χ^2 statistic = 22.66124

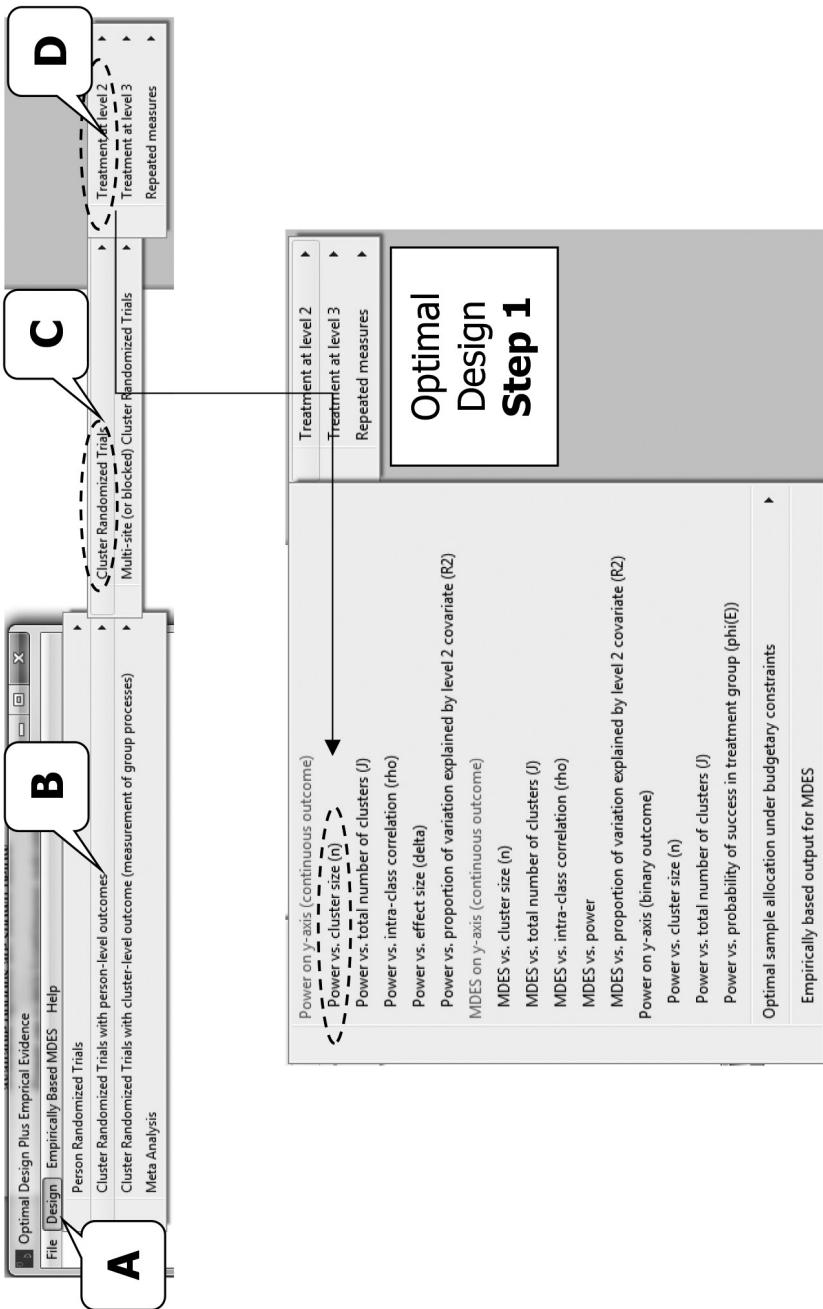
Degrees of freedom = 1

p-value = <0.001

11.5 POWER USING OPTIMAL DESIGN

Optimal Design (http://sitemaker.umich.edu/group-based/optimal_design_software) is a software program for computing power for individual- and group-level randomized designs and is freely accessible online. In the context of multilevel analyses, we are more interested in this software for its capacity in determining power for cluster or group randomization. We will illustrate the software for a simple two-level group randomized design, where individuals are nested within group and the randomization occurs at the group or cluster level (which is often applicable in education and the social sciences).

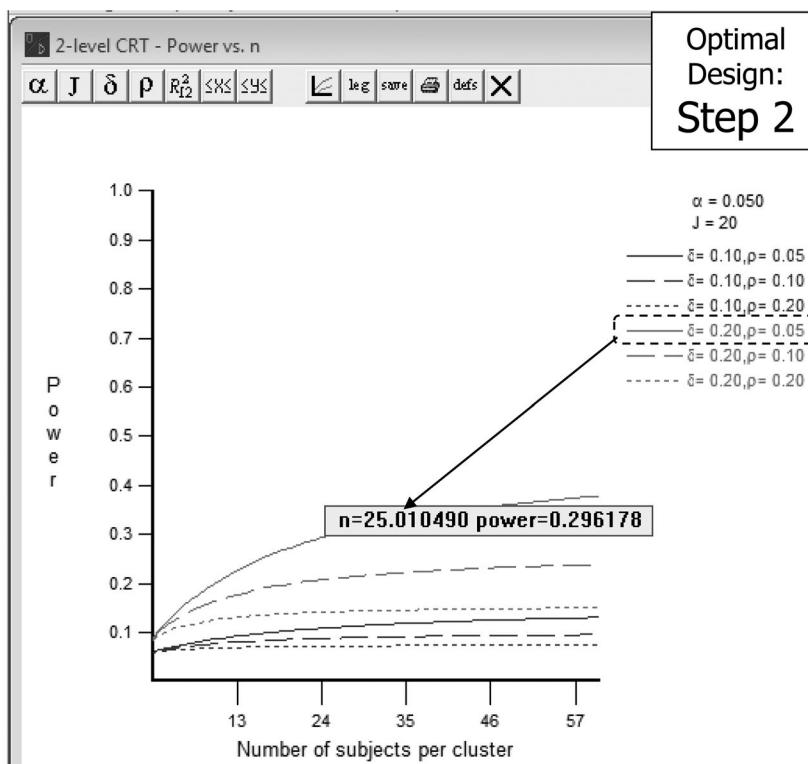
We first click on "Design," then "cluster randomized trials with person-level outcomes," then "cluster randomized trials," then "treatment at level 2." This will bring up a number of options for computing power (see Step 1). In this illustration, we will select 'power vs. cluster size (n)'.



The icons in the top left graphing dialog box allow (see Step 2) changes to be made in the settings. In this example, let's say that we want to estimate the number of units within group based on the number of groups (J). We will define the following parameters:

- Alpha = .05
- $J = 20$
- Effect size delta. Optimal Design allows up to three effect size estimates to be included. In this example, we will use .10 and .20.
- Rho. Optimal Design allows up to three values of rho to be included. In this example, we will use rho values of .05, .10, and .20.

Applying these settings, the following power curves are created. Clicking at any point on any one curve will provide the estimated power and required sample size given the parameters in that curve. For example, we see that when delta is .20 and rho is .05, the estimated sample size per group is 25 with power of approximately 30% (insufficient power).



11.6 RESEARCH QUESTION TEMPLATE AND EXAMPLE WRITE-UP

Finally, here is an example write-up for the results of the multilevel linear model. Because of the complexity of conducting multilevel modeling and the numerous

decisions that must be made in the process, guidelines have been published with recommendations on reporting multilevel results, so that both the methods and the results are thoroughly and transparently reported. The example write-up presented here includes, to the extent possible, elements as outlined in Ferron et al. (2008), that should be reported with MLM results including reporting: (a) study components, (b) model specification components, (c) estimation and inference components, (d) data components, and (e) results. *Study relevant components* that are important to report include the purpose of the study and research questions; relevant literature; sampling technique and, if applicable, sampling weights; sample size at each level of the model; how the lower-level cases are distributed across the higher-level groups; a priori power analysis; research design; variables included in the model; and reliability evidence for the variables (Ferron et al., 2008). *Model specification components* that should be reported include the number of models estimated, fixed effects and covariance structure for each model tested, the process that was used to determine the fixed effects and covariance structure for each model tested, centering, and how model fit was determined (Ferron et al., 2008). *Estimation and inference components* to report include the name of software and version used to generate the models, estimation method and any problems that occurred while identifying the models and how the problems were addressed, and methods used to make statements of inference (Ferron et al., 2008). *Data components* to report include data clustering structure (e.g., homeowner nested within neighborhood), distribution of the variables and coefficients at level 1, correlation matrix of variables included in the model, and data screening efforts (e.g., if/ how missing data were dealt, if/how missing data impacted the results, how data were screened for outliers and how they were addressed if found, extent to which assumptions were met, and consideration of how violations of assumptions may impact the results) (Ferron et al., 2008). *Results* should be reported for each research question. Additionally, the ICC, parameter estimates, precision of estimates (e.g., standard error, confidence intervals), and the extent to which limitations may impact how the results are interpreted should all be reported (Ferron et al., 2008).

Getting back to our scenario . . . Recall that the stats lab grad assistants, Challie Lenge, Ott Lier, Oso Wyse, and Addie Venture, were assisting Dr. Maeve in exploring Big Five personality data completed by individuals from multiple countries using multilevel modeling. The research questions presented to Dr. Maeve included:

- 1) Is there variation in conscientiousness between countries?
- 2) Can conscientiousness be predicted by intellect/imagination?
- 3) To what extent does the developing status of the country moderate the relationship between intellect and conscientiousness?

The team then assisted Dr. Maeve in conducting multilevel linear modeling, and a template for writing the research questions for a two-level MLM follows.

Is there variation in [OUTCOME] between [GROUP]?
Can [OUTCOME] be predicted by [LEVEL 1 PREDICTORS]?

To what extent do [LEVEL 2 PREDICTORS] moderate the relationship between [LEVEL 1 PREDICTORS] and [OUTCOME]?

It may be helpful to preface the results of the multilevel model with information on an examination of the extent to which the data were thoroughly screened.

Prior to conducting the multilevel linear model, the data were screened to determine the extent to which the assumptions associated with MLM were met. These assumptions included (1) linearity, (2) residuals (i.e., random effects) at level 1 are normally distributed and have equal variances, (3) residuals at level 2 are multivariate normal and homoscedastic, and (4) predictors at a respective level are uncorrelated with random effects at the other level.

Linearity at level 1 was reviewed by plotting the level 1 residuals to fitted values. A random display of points suggested this assumption was met. Linearity at level 2 was assessed by examining empirical Bayes residuals to level 2 predictors. A random display of points above and below zero suggests that homoscedasticity is reasonable.

In terms of normality at level 1, skewness and kurtosis were within the range of normal (−.149 and −.426, respectively). However, the Kolmogorov–Smirnov formal test of normality suggests a nonnormal distribution ($K-S = .020$, $df = 7994$, $p < .001$). The Q-Q plots and histogram suggest some nonnormality in the tails of the distribution (left tail in particular). The boxplot suggests potential outliers on the lower end of the scale.

Multivariate normality was assessed by a scatterplot of Mahalanobis distance (MDIST) and the expected values of the order statistics (CHIPCT). Points that generally adhere to the diagonal line show evidence of multivariate normally distributed data. There are several groups in the data that are suggestive of outliers, and the points are not adhering as closely to the diagonal as multivariate normal would suggest. However, these plots should only be considered accurate when the level 1 sample sizes within groups are at least moderately normal—as is the case with this data.

The hypothesis test for homogeneity of variances at level 1 suggested that equal variances between countries is plausible. Homoscedasticity at level 2 was examined by plotting empirical Bayes residuals to the level 2 predictor. A random display of points above and below zero suggests that homoscedasticity is reasonable.

Here is an example summary of results for the multilevel results (remember that this will be prefaced by a section reporting the extent to which the data were thoroughly screened).

A two-level multilevel linear model using full maximum likelihood in HLM version 7.01 was estimated. This nonexperimental design examined individuals (i.e., level 1, $N = 7994$) within North American countries (i.e., level 2, $N = 19$). There were four models estimated. The unconditional model, a one-way random effects ANOVA, was estimated first. Based on this model, the intraclass correlation coefficient (ICC) was .028, suggesting that about 3% of the variation in conscientiousness is between countries, and about 97% is within countries. Thus, it

was deemed reasonable to proceed with the multilevel linear model. There was statistically significant variation between countries ($u_0 = .016, p < .001$).

Second, a random intercepts fixed slopes model was estimated. In this model, intellect/imagination, a continuous variable, was entered in the model as a level 1 predictor that was group mean centered. The overall average conscientiousness for a person who was at their country's average intellect/imagination was 3.39 (CI 3.11, 3.67), and this was statistically significantly different from zero ($SE = .06, p < .001$). On average and across countries, intellect/imagination is positive and statistically significantly related to conscientiousness. The average effect (i.e., slope) across countries for intellect/imagination is represented as an increase of .06 points for every one point increase in conscientiousness. Adding intellect/imagination as a predictor of conscientiousness reduces the within-country variability by less than 0.5%. Variance between countries has remained the same at about 3% when controlling for intellect/imagination. Statistically significant variation in the country means still exists ($u_0 = .016, p < .001$).

Third, a random intercepts random slopes model was estimated. Intellect/imagination was entered in the model as a level 1 predictor that was group mean centered. The overall average conscientiousness for a person who was at their country's average intellect/imagination was 3.39 (CI 3.11, 3.67), and this was statistically significantly different from zero ($SE = .05, p < .001$). On average and across countries, intellect/imagination is positive and statistically significantly related to conscientiousness. The average effect (i.e., slope) across countries for intellect/imagination is represented as an increase of .07 points for every one point increase in conscientiousness. Variance between countries has remained the same at about 3% when controlling for intellect/imagination. Statistically significant variation in the country means still exists ($u_0 = .016, p < .001$). This suggests that differences between the countries in conscientiousness still exists. This between-country variability may be explained by incorporating level 2 predictors. The variance among intellect/imagination-conscientiousness slopes is not statistically significant ($SE = .0005, p = .14$), suggesting that there are similar intellect/imagination-conscientiousness relationships across countries. Thus, the slopes will be fixed in the next model. The correlation between the intercepts and slopes is $- .96$ and indicates that countries with high conscientiousness are associated with intellect/imagination-conscientiousness slopes. The proportion reduction of within-country variation was less than 1%. The proportion reduction in variation between-countries was about 25%.

Fourth, a random intercepts fixed slopes model was estimated as follows:

$$\begin{aligned} \text{Level 1: } & \text{CONSCIENTIOUSNESS} = \beta_{0j} \\ & + \beta_{1j}(\text{INTELLECT / IMAGINATION}) + r_{ij} \end{aligned}$$

$$\begin{aligned} \text{Level 2: } & \beta_{0j} = \gamma_{00} + \gamma_{01}(\text{DEVELOPING COUNTRY}) + u_{0j} \\ & \beta_{1j} = \gamma_{10} \end{aligned}$$

Intellect/imagination was entered in the model as a level 1 predictor that was group mean centered. Developing status of the country was grand mean centered as a level 2 predictor for the intercepts. The results for the model are provided in Table X. The overall average conscientiousness for a person

who was at their country's average intellect/imagination and who was not living in a developing country was 3.38 (plausible value range, 3.02, 3.74), and this was statistically significantly different from zero ($SE = .07, p < .001$). On average and across countries, intellect/imagination is positive and statistically significantly related to conscientiousness. The average effect (i.e., slope) across countries for intellect/imagination is represented as an increase of .06 points for every one point increase in conscientiousness ($SE = .01, p < .001, ES = .08$, plausible value range, -.30, .42). On average and across countries, developing status of the country where the person lives does not statistically relate to their conscientiousness (coefficient = $-.09, SE = .15, p = .54$, plausible value range = .45, .27). Variance between countries has increased to 6% when controlling for intellect/imagination and developing status of the country. Statistically significant variation in the country means still exists ($u_0 = .03, p < .001$). This suggests that differences between the countries in conscientiousness still exists. The proportion reduction of within-country variation, relative to model 3, was substantially less than 1% (.0005). The variation between countries, relative to model 3, has increased (from .01 to .03).

TABLE 11

Contextual Model

Fixed Effects	Coefficient (<i>SE</i>)	<i>t</i>	<i>p</i>	<i>ES</i> ^a	Reliability
Mean country conscientiousness (β_0)	3.38 (.07)	45.88	.001	—	.35
Intercept (γ_{00})					
Developing status (γ_{01})	−.09 (.15)	−.62	.54	—	—
Model for intellect/imagination slope (β_1)	.06 (.01)	4.77	<.001	.08	—
Intercept (γ_{10})					

Random Effects (Variance Components)	Variance	<i>df</i>	chi-square	<i>p</i>
Variance between country intercepts (τ_{00}) (u_0)	.03	17	47.77	<.001
Variance within countries (σ^2) (r)	.55			

a. Effect size calculated as $coefficient / \sqrt{\sigma^2}$

PROBLEMS

Conceptual Problems

1. Which of the following effects is evidenced in the situation where all groups have the same rate of change in the outcome?
 - a. Fixed intercepts
 - b. Fixed slopes
 - c. Nonrandomly varying intercepts
 - d. Nonrandomly varying slopes
 - e. Random intercepts
 - f. Random slopes

2. Which of the following equations reflects nonrandomly varying coefficients?
 - a. $\beta_{0j} = \gamma_{00} + u_{0j}$
 - b. $\beta_{1j} = \gamma_{10}$
 - c. $\beta_{2j} = \gamma_{20} + \gamma_{21}(\%Female)_{2j} + u_{2j}$
 - d. $\beta_{3j} = \gamma_{30} + \gamma_{31}(Rural)_{3j}$
3. Within-group variability is reflected by which one of the following?
 - a. β_{0j}
 - b. γ_{ij}
 - c. r_{ij}
 - d. u_{ij}
4. Moderation effects can be tested by which one of the following?
 - a. A comparison of the null model to the conditional model
 - b. Including a level 1 predictor with random intercepts
 - c. Including a level 1 and level 2 predictor with random slopes
 - d. Including within-group effects and random intercept effects
5. A model is estimated that has one level 1 predictor, random intercepts, and fixed slopes. Which of the following parameters are estimated in this model under full maximum likelihood estimation?
 - a. between-group variability, within-group variability, slope coefficient, intercept coefficient
 - b. within-group variation, variability in the mean outcome between groups
 - c. intercept coefficient, slope coefficient, within-group random effect, between-group random effect for the slopes
 - d. intercept coefficient, between-group random effect for the intercepts, between-group random effect for the slopes
6. A researcher is examining the influence of coffee consumption (measured as the number of cups of coffee drunk per day) on the amount of time spent on social media (measured as the number of minutes per day) for homeowners within neighborhoods (nesting based on zip code). A report based on this study reports an alpha of .05 and a p value of .02 for the slope variance for coffee consumption, (u_i). Which one of the following is a correct interpretation of this?
 - a. Neighborhoods with more time spent on social media have statistically significantly more coffee consumption.
 - b. The average amount of time spent on social media is about 2%.
 - c. There is statistically significant variability across neighborhoods in the relationship between coffee consumption and amount of time spent on social media.
 - d. Within neighborhoods, there is statistically significant variability in the amount of time spent on social media.
7. Which of the following situations would allow the correlation between the intercept and slope to be examined?
 - a. A one-way random effects ANOVA model
 - b. A random intercept and fixed slope model
 - c. A random intercept and random slope model
 - d. An unconditional model

8. A researcher finds an intraclass correlation coefficient of .05. Which one of the following reflects a correct interpretation of this?
 - a. The observed probability level is .05
 - b. The probability of rejecting the null hypothesis when it is really false is 5%
 - c. The proportion of variation in the outcome that is between groups is 5%
 - d. The variability in the outcome that is within groups is about 5%
9. A researcher examines level 1 residuals and finds the following: (a) skewness = 6.20, (b) kurtosis = 12.50, (c) Shapiro-Wilk $p = .02$, and (d) a boxplot with no outliers. Which one of the following interpretations best fits these results?
 - a. Evidence supports linearity of the data
 - b. The data are homoscedastic
 - c. There is no evidence to suggest good model fit has been achieved
 - d. The evidence generally suggests nonnormality at level 1
10. A researcher computes a null model then a random intercept fixed slope model. The researcher finds a BIC value of .03. Which one of the following is suggested?
 - a. The more parameterized model is statistically significantly a better-fitting model than the null model
 - b. The conditional model has positive evidence to support it over the null model
 - c. There is weak evidence to support the random intercept fixed slope model relative to the null model
 - d. There is strong evidence to support the null model relative to the more parameterized model

Computational Problems

1. Using the PIAAC_merged_2percentrandom_nomiss.sav dataset, conduct multi-level linear modeling following the steps in this chapter, using restricted maximum likelihood estimation and testing a two-level random intercepts model with participants nested within country [CNTRYID]. ‘Problem solving scale score’ is the outcome. ‘Age’ [AGE] and ‘female’ [FEMALE] are the level 1 predictors. Age is group mean centered and female is uncentered.

The specified model is noted as:

$$\text{Level 1: } \text{PROBLEMS} = \beta_{0j} + \beta_{1j}(\text{AGE}) + \beta_{2j}(\text{FEMALE}) + r_{ij}$$

$$\beta_{0j} = \gamma_{00} + u_{0j}$$

$$\text{Level 2: } \beta_{1j} = \gamma_{10}$$

$$\beta_{2j} = \gamma_{20}$$

Compute the ICC for the null model. Interpret the fixed and random effects of the random intercepts model. Test for equal variances at level 1. Interpret goodness-of-fit using the deviance test.

2. Using the nurses.sav dataset (available from <http://www.joophox.net/mlbook2/MLbook.htm>), conduct multilevel linear modeling following the steps in this chapter, using restricted maximum likelihood estimation and testing a two-level random intercepts random slopes model with nurses nested within hospital [HOSPITAL]. ‘Nurse stress score’ [STRESS] is the outcome. ‘Experimental vs. control condition’ [EXPCON] is the level 1 predictor, uncentered. ‘Hospital size’ [HOSPSIZE] is the level 2 predictor, grand mean centered, and is estimated to predict both the intercepts and the slopes. The specified model is noted as:

$$\text{Level 1: } \text{STRESS} = \beta_{0j} + \beta_{1j}(\text{EXPCON}) + r_{ij}$$

$$\begin{aligned}\text{Level 2: } \beta_{0j} &= \gamma_{00} + \gamma_{01}(\text{HOSPSIZE}) + u_{0j} \\ \beta_{1j} &= \gamma_{10} + \gamma_{11}(\text{HOSPSIZE}) + u_{1j}\end{aligned}$$

Compute the ICC for the null model. Interpret the fixed and random effects of the random intercepts model. Test for equal variances at level 1. Interpret goodness-of-fit using the deviance test.

Interpretive Problem

1. Use HLM to conduct multilevel linear modeling with the additional variables (Big5_N_7994.sav). Write up the results.

REFERENCES

- Bell, B. A., Morgan, G. B., Schoeneberger, J. A., Loudermilk, B. L., Kromrey, J. D., & Ferron, J. M. (2010). *Dancing the sample size limbo with mixed models: How low can you go?* Paper presented at the SAS Global Forum 2010, Seattle, WA.
- Browne, W. J., & Draper, D. (2000). Implementation and performance issues in the Bayesian and likelihood fitting of multilevel models. *Computational Statistics*, 15(3), 391–420.
- Eberly, L. E., & Thackeray, L. M. (2005). On Lange and Ryan’s plotting techniques for diagnosing non-normality of random effects. *Statistics and Probability Letters*, 75, 77–85.
- Enders, C. K., & Tofghi, D. (2007). Centering predictor variables in cross-sectional multilevel models: A new look at an old issue. *Psychological Methods*, 12(2), 121–138.
- Ferron, J. M., Hogarty, K. Y., Dedrick, R., Hess, M., Niles, J., & Kromrey, J. D. (2008). Reporting results from multilevel analysis. In A. A. O’Connell & D. B. McCoach (Eds.), *Multilevel modeling of educational data* (pp. 391–426). Charlotte, NC: Information Age Publishing.
- Freedman, D. A. (1999). Ecological inference and the ecological fallacy. In N. J. Smelser & P. B. Baltes (Eds.), *International encyclopedia of the social and behavioral sciences* (Vol. 6, pp. 4027–4030). New York, NY: Elsevier.
- Goldberg, L. R. (1992). The development of markers for the big-five factor structure. *Psychological Assessment*, 4(1), 26–42.

- Heck, R. H., & Thomas, S. L. (2000). *An introduction to multilevel modeling techniques*. Mahwah, NJ: Lawrence Erlbaum.
- Hedges, L. V., & Hedberg, E. C. (2007). Intraclass correlation values for planning group-randomized trials in education. *Educational Evaluation and Policy Analysis*, 29(1), 60–87.
- Hox, J. J. (2002). *Multilevel analysis*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Hox, J. J. (2010). *Multilevel analysis: Techniques and applications*. New York, NY: Routledge.
- Kish, L. (1965). *Survey sampling*. New York: Wiley.
- Kreft, I., & de Leeuw, J. (1998). *Introducing multilevel modeling*. Thousand Oaks, CA: Sage.
- Lee, V. E., Loeb, S., & Lubeck, S. (1998). Contextual effects of prekindergarten classrooms for disadvantaged children on cognitive development: The case of chapter 1. *Child Development*, 69(2), 479–494.
- Longford, N. T. (1993). *Random coefficients models*. New York, NY: Oxford University Press.
- Luke, D. A. (2004). *Multilevel modeling*. Thousand Oaks, CA: Sage.
- Maas, C. J. M., & Hox, J. J. (2004). Robustness issues in multilevel regression analysis. *Statistica Neerlandica*, 58(2), 127–137.
- Maas, C. J. M., & Hox, J. J. (2005). Sufficient sample size for multilevel modeling. *Methodology: European Journal of Research Methods for the Behavioral and Social Sciences*, 1(3), 86–92.
- Mathieu, J. E., Aguinis, H., Culpepper, S. A., & Chen, G. (2012). Understanding and estimating the power to detect cross-level interaction effects in multilevel modeling. *Journal of Applied Psychology*, 97(5), 951–966.
- Pinheiro, J. C., & Bates, D. M. (2000). *Mixed-effects models in S and S-Plus*. New York, NY: Springer.
- Raftery, A. E. (1995). Bayesian model selection in social research. *Sociological Methodology*, 25, 111–163.
- Raudenbush, S. W. (2011). Optimal design software for multilevel and longitudinal research (version 3.01). Retrieved from www.wtgrantfoundation.org
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical linear models: Applications and data analysis methods* (2nd ed.). Thousand Oaks, CA: Sage.
- Raudenbush, S. W., Bryk, A., Cheong, Y. F., Congdon, R., & du Toit, M. (2011). *HLM 7: Hierarchical linear and nonlinear modeling*. Lincolnwood, IL: Scientific Software International, Inc.
- R Development Core Team. (2011). *R: A language and environment for statistical computing*. Vienna, Austria: R Foundation for Statistical Computing.
- Scherbaum, C. A., & Ferreter, J. M. (2009). Estimating statistical power and required sample sizes for organizational research using multilevel modeling. *Organizational Research Methods*, 12(2), 347–367.
- Shieh, Y., Fouladi, R. T., & Pullman, T. (1999). *The effect of error term nonnormality on multilevel model parameter estimates and standard errors*. Paper presented at the American Educational Research Association annual meeting, Montreal, Canada.
- Singer, J. D., & Willett, J. B. (2003). *Applied longitudinal data analysis: Modeling change and event occurrence*. New York: Oxford University Press.
- Snijders, T. A. B., & Bosker, R. J. (1993). Standard errors and sample sizes for two-level research. *Journal of Educational Statistics*, 18(3), 237–259.

- Snijders, T. A. B., & Bosker, R. J. (1999). *Multilevel analysis: An introduction to basic and advanced multilevel modeling*. Thousand Oaks, CA: Sage.
- Snijders, T. A. B., & Bosker, R. J. (2012). *Multilevel analysis: An introduction to basic and advanced multilevel modeling* (2nd ed.). Thousand Oaks, CA: SAGE.
- Spybrook, J. (2008). Power, sample size, and design. In A. A. O'Connell & D. B. McCoach (Eds.), *Multilevel modeling of educational data*. Charlotte, NC: Information Age Publishing.
- Spybrook, J. (2014). Detecting intervention effects across context: An examination of the precision of cluster randomized trials. *Journal of Experiential Education*, 82(3), 334–357.
- Swaminathan, H., & Rogers, H. J. (2008). Estimation procedures for hierarchical linear models. In A. A. O'Connell & D. B. McCoach (Eds.), *Multilevel modeling of educational data* (pp. 469–520). Charlotte, NC: Information Age Publishing.
- Van der Leeden, R., Busing, F.M.T.A., & Meijer, E. (unknown). Bootstrap methods for two-level models. Leiden University. The Netherlands. Retrieved from <http://www.soziologie.uni-halle.de/langer/buecher/mehrebenen/literatur/busing1997.pdf>

Chapter 12

PROPENSITY SCORE ANALYSIS

CHAPTER OUTLINE

12.1 What Propensity Score Analysis Is and How It Works	572
12.1.1 Characteristics	574
12.1.2 Sample Size	581
12.1.3 Assumptions	581
12.2 Mathematical Introduction Snapshot	582
12.3 Computing Propensity Score Analysis Using R	582
12.3.1 PSA Using MatchIt in R	583
12.4 Example Write-Up	594

KEY CONCEPTS

1. Balance
2. Common support
3. Counterfactual framework of causality
4. Ignorable treatment assignment assumption
5. Matching algorithms
6. Propensity score
7. Sensitivity
8. Stable unit treatment value assumption (SUTVA)

In the previous chapters, with the exception of factor analytic types of procedures, we have generally considered procedures that produce ‘end results.’ In other words,

models are estimated and tested and inferences are made from those results. In comparison, propensity score analysis (PSA) is a precursor to estimating outcomes, the results from which can be applied to any number of procedures for examining causal outcomes. For cooking enthusiasts, this could be considered akin to marinating. The marinade is something applied prior to the cooking procedure, a tool used to help produce a better outcome. Likewise with PSA. The theme of this chapter is how to conduct PSA. The procedures from many of the previous chapters can then be applied to answer the substantive efficacy questions that PSA is often created for. In this chapter, we will, however, focus on getting the data prepared—through propensity score analysis—for that substantive analysis. In the propensity score process, there are a number of analytic decisions that must be made when conducting PSA, and several new foundational concepts that come with this. Like other procedures covered in the text, propensity score matching is relatively complex. This text is meant to serve as an introduction. Readers that are interested in additional coverage may want to review other excellent sources (Guo & Fraser, 2010, 2015). My objectives are that, by the end of this chapter, you will be able to (a) understand the basic concepts of PSA, (b) estimate propensity scores, (c) understand various conditioning strategies and be able to match units on their propensity scores, and (d) check model adequacy.

12.1 WHAT PROPENSITY SCORE ANALYSIS IS AND HOW IT WORKS

It is the last hoorah, so to speak, for the graduate research assistants in the stats lab. Today, we find them working together on a most exciting project.

Challie Lenge, Ott Lier, Oso Wyse, and Addie Venture are nearing the end of their graduate studies, and thus their research positions as well. Their faculty supervisor and mentor has asked the group to consult with the local school district's Chief Research Operations Officer, Dr. Mayfield, who is interested in examining academic performance of students who attend public versus private schools using public data available from the U.S. Department of Education. Dr. Mayfield's research question is as follows: *Is there a mean difference in reading performance for students who attend public versus private school?* Dr. Mayfield is interested in exploring the national data prior to and as a way to inform work they are doing locally in their own district. In this effort, Dr. Mayfield has sought advice from the research lab on how to best approach this question. Challie and Ott, the lead researchers on the project, shared with Dr. Mayfield that because there was not random assignment to groups—assignment to attend either a public or private school—selection bias was likely. As such, propensity score analysis for matching the students before the substantive efficacy analyses was suggested by the research assistants.

Propensity score analysis is defined as the conditional probability of treatment assignment given observed covariates, and Rosenbaum and Rubin (1983) are attributed as the first to propose it (although its mathematical roots have a much earlier history). The propensity score is one value—a scalar that summarizes information from observed

covariates and reflects the probability of being ‘treated.’ Propensity scores are also balancing scores. This means that distribution of the observed covariates that were used to compute the propensity score are the same in the treated and untreated groups at each value of the propensity score. Because of this, subjects with similar propensity scores can be grouped, similar to a randomized experiment. The value in a randomized experiment is the ability to infer causation.

This is a great segue to introduce the counterfactual framework, more specifically referred to as the *Neyman-Rubin counterfactual framework of causality*, although other researchers beyond Neyman and Rubin have contributed to the framework (e.g., Holland, 1986; Little & Rubin, 2000; Neyman, 1923, 1990; Rubin, 1974, 1978, 2005). A counterfactual is a potential outcome—an outcome that would have occurred absent the cause (Shadish, Cook, & Campbell, 2002). In the case of an intervention, the counterfactual for the treatment units would be the potential outcome under control conditions, and the counterfactual for the control units being the outcome under treatment conditions. The counterfactual is not observable. It is not possible to know what would have been the outcome if the treatment units would have received the control and likewise for the control. The Neyman-Rubin counterfactual framework of causality posits that the counterfactual can indeed be estimated, with the framework expressed as follows:

$$Y_i = W_i Y_{1i} + (1 - W_i) Y_{0i}$$

Where Y_i is the outcome, Y_{1i} is the outcome under the treated state, and Y_{0i} is the outcome under control conditions. W_i is a binary variable with 1 indicating treatment receipt and 0 indicating no treatment. W_i serves to turn the treatment on or off. To infer causality between the treatment ($W_i = 1$) and the outcome (Y_i), the outcome under control conditions (i.e., Y_{0i} when $W_i = 0$) must be examined in relation to the outcome under treated conditions (i.e., Y_{1i} when $W_i = 1$). The counterfactual is estimated as the difference in the average outcome for the treated and untreated units, also known as the *standard estimator for the average treatment effect* (ATE), expressed as follows:

$$\tau = E(Y_1 | W = 1) - E(Y_0 | W = 0)$$

Where the treatment effect (τ) is the difference between $E(Y_1 | W = 1)$, which is the average observable outcome for treated units, and $E(Y_0 | W = 0)$, which is the average observable outcome for the untreated units. In this equation, $E(Y_0 | W = 0)$ is used to estimate the counterfactual [i.e., $E(Y_0 | W = 1)$], and that is of most interest. ATE is analogous to *intent-to-treat*, the effectiveness of an intervention when implemented under actual conditions (Shadish et al., 2002).

The counterfactual argument just posed is overly simple, and does not take into account consideration of potential covariates. When this is done, the *ignorable treatment assignment assumption* is required (Rosenbaum & Rubin, 1983), which means that the assignment to treatment or nontreatment is independent of the nontreated (i.e., Y_0) and treated (i.e., Y_1) outcomes when condition on covariates (i.e., X). In other words,

assignment to groups is independent of the potential outcomes when the observed covariates are held constant (Guo & Fraser, 2010). The *ignorable treatment assignment assumption* is usually quite easily met in randomized controlled trials through the balancing created by the randomization, but it is usually violated in quasi-experiments and observational studies due to the lack of random assignment to groups (i.e., the creation of groups in nonrandomized studies is confounded with the outcomes). The *nonignorable treatment assignment* thus suggests that there are other factors that determine W , the receipt or nonreceipt of, and these factors should be taken into account in the statistical model to avoid biased results (Guo & Fraser, 2010).

The outcomes in the counterfactual framework are sufficiently represented only when the *stable unit treatment value assumption (SUTVA)* holds (Rubin, 1980, 1986). SUTVA is an a priori assumption that the outcome for treatment unit i remains the same regardless of the treatment assignment mechanism for i and regardless of treatments received by other units.

When faced with nonignorable treatment assignment, one of the ways that the counterfactual can be estimated is through propensity score matching (Rosenbaum & Rubin, 1983), which allows nonrandomized studies to approximate randomization (Rubin, 2008). It is important to note that propensity score analysis does not “correct” for nonignorable treatment assignment, but rather PSA controls for overt (not hidden) selection bias (Guo & Fraser, 2010, p. 38).

12.1.1 Characteristics

12.1.1.1 Analytic Decisions in Propensity Score Analysis

Randomized controlled trials (RCTs) have the unique characteristic of random assignment to groups. RCTs are considered the most rigorous research design because they ensure subjects are balanced on all covariates (both observed and unobserved) that may potentially impact the outcome. In many cases, random assignment to groups of interest is not feasible or ethical. Observational studies are those empirical examinations of causal effects when random assignment to allocate units to groups (e.g., treatment or control) is not performed. Consider, for example, the effect of mental disorder on recidivism or the effect of preschool attendance on academic outcomes. In neither example would random assignment to group (i.e., diagnosis of mental disorder or no mental disorder diagnosis; preschool attendance or nonattendance) be possible. Propensity score analysis is a rigorous quasi-experimental design that allows matching of subjects in observational studies that lack random assignment and produces unbiased treatment estimates and that therefore can be a useful tool prior to estimating causal effects (Rubin, 1974).

Rosenbaum and Rubin (1983) produced one of the seminal works on propensity score methods. As they defined it, the propensity score is the probability of assignment to treatment conditioned on observed covariates measured on or before baseline. Propensity

scores are often termed ‘balancing scores,’ as the distribution of the observed covariates (which are used to estimate the propensity score) are similar between the units within the groups (i.e., the treated and untreated units) conditional on the propensity score. In other words, units (e.g., individuals) that have similar propensity scores will also have similar distributions of observed covariates. In experiments where random assignment is performed, the propensity score is known and defined by the research design. In observational studies, the propensity score can be estimated using covariates measured at or before the treatment, Z_i .

Propensity score analysis is a multistep procedure, the ultimate goal of which is to produce matched treated and untreated units from which intervention efficacy on some outcome can be determined. The first step in propensity score analysis is to estimate the propensity score (i.e., the predicted probability of treatment). This entails selecting the observed covariates measured at or before treatment, as well as the propensity score estimation method. After the propensity score is estimated, the second step is to determine model adequacy. This is termed ‘balance’ and answers the question of whether there are comparable treatment groups such that the distributions of covariates are similar. If there is balance, the third step is to actually condition (i.e., match) on the propensity scores, creating matched sets of treated and untreated units that have similar propensity scores. The fourth and final step, the postmatching analysis, is to estimate treatment efficacy by examining the outcomes of the matched units. In this chapter, we will examine the first three steps in more detail.

Estimating the Propensity Score

The first step in propensity score analysis is to **estimate the propensity score** (i.e., the predicted probability of treatment). Any model that has a binary causal variable, such as attendance at public versus private school or participation in a voluntary training program, which can be related to a set of pretreatment covariates, can be used to estimate propensity scores. There are two primary decisions that must be made at this step. One is the selection of how the propensity score will actually be estimated, and the other is the selection of observed covariates measured at or before treatment that will be used in the model that estimates the propensity score.

Logistic regression is the most commonly used estimation method for propensity scores in the social sciences (Thoemmes & Kim, 2011) and will be used as the estimation method for the illustrations within this chapter. Logistic regression is a non-linear model. The application of a link function, such as a logit function, allows the model to be expressed as a generalized linear model. Transformation of the treatment assignment (W) through a link function creates a linear function (Guo & Fraser, 2010). With the rise in application of propensity score analysis, however, additional methods are being examined (Lee, Lessler, & Stuart, 2010; McCaffrey, Ridgeway, & Morral, 2004; Setoguchi, Schneeweiss, Brookhart, Glynn, & Cook, 2008). Although logistic regression has been shown to generally perform well under varying conditions, other estimation methods (random forests, bagged and boosted classification and regression

trees (CART)) may reduce additional bias in some conditions (e.g., when interactions are not included, nonlinearity) (Lee et al., 2010). It is important to remember that the goal of whatever estimation is selected is to obtain the propensity score, not the parameter estimates produced as a result of estimating the propensity score.

The next decision in estimating the propensity score is **covariate selection**. Unlike most other statistical procedures, estimating the propensity score is *not* one that requires parsimony, although thoughtful consideration for covariates is still required. “The literature on propensity score matching is nearly unanimous in its emphasis on the importance of including in models carefully chosen and appropriate conditioning variables in the correct functional form” (Guo & Fraser, 2010, p. 139). Researchers (Brookhart et al., 2006) have found that optimal PS estimation models include all variables that are known or thought to be related to the outcome (e.g., measures of the outcome taken prior to treatment assignment such as pretest or baseline measures), regardless of if, or to what extent, they may be related to treatment assignment. However, including a variable that is related to the treatment assignment but unrelated to the outcome does not decrease bias but increases variance of the estimated exposure effect and thereby introduces noise that may cause misclassification of units to treatment condition or mismatched units (Brookhart et al., 2006). Excluding variables that are measured after treatment assignment is therefore best practice. Thus while the number of covariates included in the PS estimation is not capped, care does need to be taken in the variables included so that the assumption of strongly ignorable treatment assignment is met. If the strongly ignorable treatment assignment holds, assignment to group (treated or untreated) is independent of the outcomes given the covariates—the difference in the mean outcome between the treated and untreated groups at a given propensity score are unbiased estimates of treatment effect at that propensity score, and therefore matching on the propensity score produces unbiased average treatment effect estimates—and there is a nonzero probability of receiving each treatment given the covariates (Rosenbaum & Rubin, 1983). Research has also demonstrated that inclusion of a small number of demographics yields poor results (Shadish, Clark, & Steiner, 2008), and an upper limit for the number of covariates to include does not seem to exist. Examples exist of studies that have estimated propensity scores with upwards of 200 (Gasper, DeLuca, & Estacion, 2012) and over 200 covariates (Hong & Raudenbush, 2005).

Checking Model Adequacy

Model adequacy includes both an examination of balance and common support (Thoemmes & Kim, 2011). **Balance checks** are done to determine if balance of means and variances on the observed covariates has been achieved. In other words, balanced propensity scores have distributions of covariates that are comparable for the treated and untreated groups. Common tools to measure model adequacy include computing the standardized mean difference (i.e., Cohen’s d effect size, which should be near zero after matching) and variance ratio (which should be near 1 after matching). There should also be consideration of the inclusion of high-order polynomials (e.g., squared

term) and/or cross-product interaction terms (Rosenbaum & Rubin, 1984, 1985). The inclusion of these terms is usually done after matching on the propensity scores (i.e., the PS is reestimated with the new terms included) and only if there is insufficient overlap between groups after estimating the PS (Rubin, 1997). Should there be small effects between treated and untreated groups, evidence of balance exists.

Graphs can also be used to check balance. Boxplots of covariates by group (treated and untreated) as well as histograms and Q-Q plots are visual (but somewhat subjective) tools for examining balance.

The region of **common support** is the area of overlap between propensity score distributions. Ideally, there should be much overlap, which allows estimates of causal effects to apply to all or nearly all of the sample (Thoemmes & Kim, 2011). Cases that are outside the region of common support should be considered for exclusion (Imai, King, & Stuart, 2008).

Conditioning on the Propensity Score

Conditioning or matching on the propensity score refers to using the estimated propensity score to create sets of treated and untreated subjects or units. This is the step where the treated and untreated units are actually matched based on the estimated propensity score. It is usually assumed that there are more untreated than treated subjects (Gu & Rosenbaum, 1993). There are a number of different ways that conditioning or matching on the propensity score can be accomplished. Rosenbaum and Rubin (1983) provided three general methods: matching, stratification, and covariance adjustment. Weighting on the propensity score is an additional way that PS can be applied to balance groups (Luellen, Shadish, & Clark, 2005). Should the researcher decide to apply the propensity score as a sampling weight, this step of matching on the propensity score is skipped (as there is no matching entailed), and the researcher proceeds directly to efficacy analyses. Covariance adjustment is not recommended given assumption constraints, specifically that a linear relationship exists between the propensity score and the outcome and that there is no relationship between the propensity score and the treatment (Rosenbaum & Rubin, 1983), and so will not be explained further. Thus, the conditioning methods that we will focus include matching and stratification.

Matching

Conditioning on the propensity score using a matching technique means that treated and untreated units are matched based on the same/similar propensity scores. Researchers who use matching as the general conditioning strategy must make additional decisions based on distance, algorithm, and structure (Gu & Rosenbaum, 1993):

- 1) algorithm—whether or not the match minimizes the average distance on the PS (optimal matching versus greedy matching);

- 2) distance—selection criteria for how close in value the propensity score is between the treated and untreated subjects (approximate, termed ‘nearest neighbor,’ or exact); in situations where cases are matched based on approximate propensity scores, a caliper can be defined that allows the match to be done within a specified range (‘nearest neighbor within caliper’) or one or more of the untreated subjects can be randomly selected (‘nearest neighbor’);
- 3) structure—selection criteria for untreated subjects (i.e., the number of treated to control units to match such as 1:1 or 1:many).

An additional use of propensity scores is through inverse probability of treatment weighting (IPTW).

Algorithms

Matching algorithms create sets of treated (n) and untreated (m) subjects based on propensity score distance. The broad umbrellas for algorithms include greedy matching, optimal matching, and fine balance. Greedy and optimal matching both match pairs or sets of pairs using one single value—the propensity score. Fine balance, in comparison, does not.

Greedy Matching. In greedy matching, treated subjects are placed in random order and the first in order is then paired to the untreated subject whose propensity score is closest in value. In this process, there is no regard to whether the untreated subject would be a closer match to a subsequent treated subject (i.e., there is no concern to minimizing the distance between the propensity scores of the matched set), and thus the term ‘greedy.’ There are criticisms to greedy matching algorithms. For example, the treated and untreated units require a good amount of overlapping common support region in order to be ensured that matching will be found (the region of common support is the range of propensity scores within which both the treated and untreated units have). Even when applying a caliper, matches may be difficult to find. This results in cases being excluded from the analysis as sufficient matches could not be identified. There are a number of matching algorithms that fall under the umbrella of greedy matching. These include Mahalanobis metric matching and nearest neighbor matching.

Mahalanobis Metric Matching. The original Mahalanobis metric matching method was conceived prior to propensity score matching (Cochran & Rubin, 1973). With this method, units are matched based on distance between treated and untreated units, with distance defined as Mahalanobis distance. Matches may become difficult to find using this method if there are many covariates, as the average Mahalanobis distance increases as the number of covariates increase. A modification of this method includes the propensity score as an additional covariate (Guo & Fraser, 2010).

Nearest Neighbor. Nearest neighbor is one of the most frequently used matching algorithms (Thoemmes & Kim, 2011). With 1:1 nearest neighbor matching, each treated unit is matched with one and only one untreated subject based on minimized absolute

difference of propensity scores between all pairs or propensity scores. Restrictions are not placed on how close or how far the ‘neighbors’ are, only that the estimated propensity score is closest for the matched pair. With 1:many nearest neighbor matching, once each treated subject has been matched with one untreated subject, the process repeats with the first treated subject assigned to the nearest of the remaining $m - n$ untreated units. *Nearest neighbor within caliper* imposes a restriction on the distance between the propensity scores of the treated and untreated units. The recommended caliper is .25 or less standard deviation of the estimated propensity score (Rosenbaum & Rubin, 1985).

Selecting to match untreated units *without* replacement means that once an untreated subject has been matched to a treated unit, it will no longer be considered as a potential match for any other treated units. Matching *with* replacement allows untreated units to be considered for more than one matched set. Because of the duplication of untreated subjects in the matching process, adjustments must be made to estimating the variance (Hill & Reiter, 2006).

The best of all worlds is obtained through *nearest available Mahalanobis metric matching within calipers defined by the propensity score* matching method, which combines Mahalanobis metric matching and nearest neighbor. This method has been noted as producing the best covariate balance (Rosenbaum & Rubin, 1985).

Optimal Matching. Optimal matching methods were created as a result of criticisms of greedy matching methods (Guo & Fraser, 2010). Optimal matching forms sets so that the total distance in propensity scores between pairs is minimized but also optimized, as there is consideration throughout the matching decision process such that each decision informs the next, and later decisions go back and inform earlier decisions (Guo & Fraser, 2010). Restrictions can be imposed in optimal matching (e.g., 1-to-1 matching, matching within ratio of treated to untreated, variable matching specifying minimum and maximum untreated to treated units, or full matching); however, the same restrictions imposed in optimal and greedy algorithms will produce different matched sets (Guo & Fraser, 2010). The following recommendations in terms of optimal matching structure have been offered: (1) efficiency is gained by having two untreated for each treated as compared to a one-to-one match, (2) efficiency is not substantially increased with a large number of untreated units, and (3) efficiency is not harmed when there is variation in the number of matched units among all strata (Haviland, Nagin, & Rosenbaum, 2007). In the event that there are multiple untreated to treated units, testing different restrictions to determine the best bias reduction and efficiency is quite common (Guo & Fraser, 2010).

Fine Balance. Unlike greedy and optimal, which use the propensity score to match, fine balance exactly balances on a nominal variable.

Distance

Matching algorithms use the distance between propensity scores to assign treated and untreated units. The traditional greedy algorithm described previously is approximate,

termed ‘nearest neighbor,’ or exact; in situations where cases are matched based on approximate propensity scores, a caliper can be defined that allows the match to be done within a specified range (‘nearest neighbor within caliper’) or one or more of the untreated subjects can be randomly selected (‘nearest neighbor’). A caliper of .25 standard deviations of the propensity score is suggested to remove at least 90% of the bias (Rosenbaum & Rubin, 1985).

Structure

When researchers match 1:1 (also called pair matching), each treated subject is matched to one untreated subject, and unmatched units are discarded. With 1:many matching, each treated subject can be matched to more than one untreated subject. Full matching, a variant of 1:many matching, entails creating matched sets that have either one treated and one or more untreated *or* one untreated and one or more treated units.

Stratification

While matching may require the most decisions, the additional conditioning methods also require decisions that should be reported. Stratification involves ranking cases based on their estimated propensity scores and then stratifying into mutually exclusive and equal-sized subgroups. Treated and untreated units have approximately equal values on the propensity score within each strata or subgroup, which then means that the observed covariates on which they were matched are distributed similarly between treated and untreated groups. This means, then, that within strata, the distribution of measured covariates are about the same between treated and untreated units. Once stratification is performed, treatment effects are estimated within each stratum, and the stratum estimates can then be pooled to generate an estimated overall treatment effect (Rosenbaum & Rubin, 1984). These estimates are often weighted by the proportion of cases within the respective stratum. ATE is estimated when stratum weights of the respective stratum, computed as the inverse of the number of strata (i.e., $1/(\text{number of strata})$), are used when pooling the treatment effects across stratum. ATT is estimated when stratum weights that are equal to the percentage of treated units within stratum are applied. Rosenbaum and Rubin (1984) is the classic study cited on propensity score stratification, which builds from Cochran (1968), providing evidence that 90% of the bias based on measured confounders is removed when five strata are created, assuming linearity between propensity and outcome. In the case of stratification, the number of strata and how the strata were formed should be reported (Thoemmes & Kim, 2011).

Inverse Probability of Treatment Weighting (IPTW)

Inverse probability of treatment weighting, proposed by Rosenbaum (1987), applies the propensity score as a weight. The weight is equal to the inverse of the probability of receiving the treatment (or no treatment) assigned. In doing this, the sample is such that the distribution of observed baseline covariates that were applied in creating the

propensity score is independent of treatment assignment. Cases with very low probability or receiving the assigned treatment (or no treatment) may produce weights that are inaccurate (Austin, 2011)

12.1.2 Sample Size

Propensity score analysis is a large sample size procedure. However, very little is actually written about requisite sample size for respective conditioning strategies and/or matching algorithms, and thus operational definitions of ‘large sample size’ in relation to PSA are nonexistent. What researchers must keep in mind is that PSA is quite carnivorous in respect to data—in terms of requiring both a large sample size to support the number of variables that can be applied in the process of estimating the propensity score and also the sample size to support matching algorithms that may throw out a lot of data. Considering the algorithm for estimating the propensity score, recall from our coverage of logistic regression that simulation research suggests that logistic regression is best used with large samples. Samples of size 100 or greater are needed to accurately conduct tests of significance for logistic regression coefficients (Long, 1997). Additionally, as we learned earlier, many cases may be tossed when using greedy matching, as sufficient matches may be difficult to find. Thus, while the application of logistic regression in the PSA process is not as a test of inference, considering 100 as a minimum sample size for PSA may be an appropriate starting point. However, research to support minimal sample size in PSA is needed.

12.1.3 Assumptions

Propensity score analysis does not have assumptions per se, but rather the assumptions that govern PSA are those that are applicable to the procedures applied during the PSA process. For example, the assumptions associated with the method selected to estimate the propensity score (e.g., logistic regression) should be reviewed. Likewise, the assumptions for the procedure used in the final step, to actually estimate treatment efficacy after PS matching, should also be examined. As stated previously, logistic regression is the most common method for estimating propensity scores. Readers who apply this method of estimation are referred to the assumptions that are detailed in the applicable chapter.

12.1.3.1 Conditions

Conditional independence assumes that self-selection into the treatment (or nontreatment) is based on (and influenced by) observable characteristics and the outcome, after conditioning on the observables, does not depend on the treatment. Omitting key variables that may influence the self-selection process may bias the results. As stated earlier, covariates that are considered for use in estimating the propensity score should be measured before treatment or be time invariant (or at least anticipated to be very stable over time).

12.2 MATHEMATICAL INTRODUCTION SNAPSHOT

Rosenbaum and Rubin (1983) defined the propensity score as the probability of assignment to treatment conditioned on pretreatment observed covariates. Standard notation for this statement is as follows:

$$e_i = \Pr(Z_i = 1 | X_i)$$

Where e_i is the estimated propensity score,

Z_i indicates treatment received (i.e., treatment or control) for individual i , and

X_i is the observed covariates from which the propensity score is estimated.

The Neyman-Rubin counterfactual framework of causality posits that the counterfactual can indeed be estimated, with the framework expressed and explained earlier. The counterfactual is estimated as the difference in the average outcome for the treated and untreated units, also known as the *standard estimator for the average treatment effect* (ATE). ATE is analogous to *intent to treat*, the effectiveness of an intervention when implemented under actual conditions (Shadish et al., 2002).

The ATE and related intention to treat (ITT) concept were introduced earlier. As you may recall, ITT is of interest when the researcher wants to determine if the treatment is effective for all units who were *assigned* to receive the treatment. Thus, the *intent* to treat implies that all assigned to treatment may not have actually received it (e.g., non-adherence). Depending on the situation, a researcher may be interested in the average treatment effect for the treated (ATT) or the treatment on the treated (i.e., as *received* in comparison to ‘as assigned’). ATT and ATE are not equivalent. In ATT, the interest is on those who actually received the treatment (which may not be all those that were *assigned* to receive treatment). ATT for the *treated* can be estimated as:

$$E(Y_1 - Y_0) | X, W = 1$$

And ATT for the *untreated* can be estimated as:

$$E(Y_1 - Y_0) | X, W = 0$$

12.3 COMPUTING PROPENSITY SCORE ANALYSIS USING R

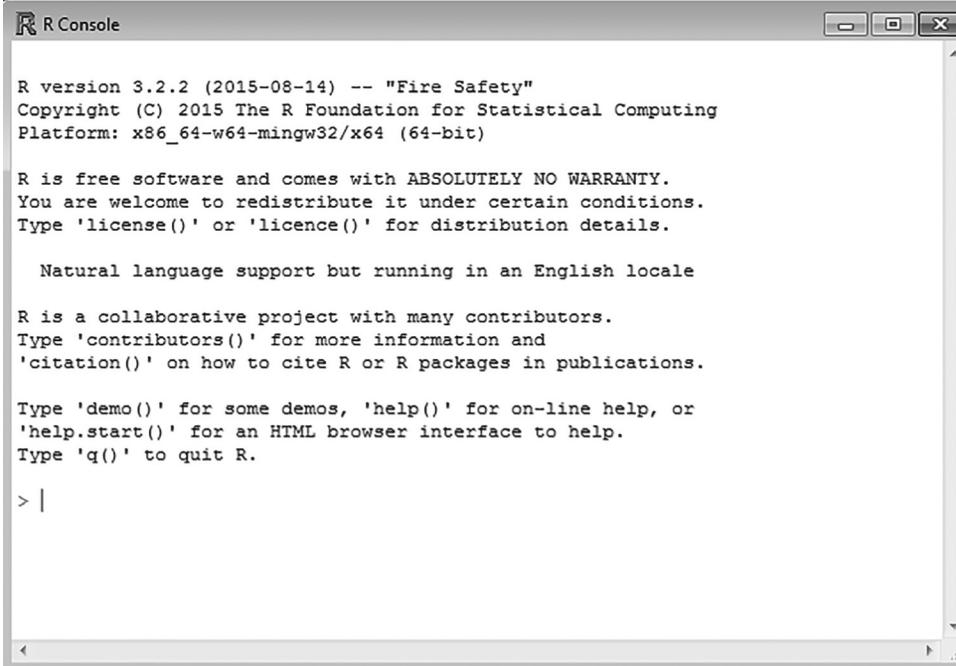
Next we consider **R** software for propensity score matching. Before we conduct the analysis, let us review the data. The data we are using is the Early Childhood Longitudinal Study Kindergarten Class of 1998–1999 (ECLS-K) (<http://nces.ed.gov/ecls/kindergarten.asp>), available through the U.S. Department of Education National Center for Education Statistics Institute of Education Sciences (NCES IES). The ECLS-K is a longitudinal study, following a nationally representative sample of public and private school children from kindergarten through eighth grade. Further details on the

ECLS-K are provided in the MANOVA chapter, and readers that need a refresher are referred back to the respective chapter. The data for this illustration were presented in Hahs-Vaughn (2015) and the data can be accessed from http://people.duke.edu/~wp40/data_code.html.

For this illustration, the dependent variable (r) is reading performance, measured by the reading IRT scale score [C7R4RSCL] collected during the spring of the student's eighth-grade year. The causal variable of interest is attendance in public (coded as '1') versus private school (coded as '0'). I have delimited this particular dataset to include only children who did not change school type (e.g., moving from public to private or vice versa). Thus, the causal variable represents public or private school attendance from kindergarten through eighth grade. The covariates that we will use to estimate the propensity score were measured during the first wave of the ECLS-K data collection (i.e., fall of the child's kindergarten year). While these covariates were measured after treatment had commenced (i.e., in the fall semester of the year in which the child was enrolled in public or private school), the covariates selected were such that it could be assumed that the values reported by the parents are those that would have been measured prior to treatment and were unaffected by the child being in kindergarten at the time the parent reported the information (e.g., number of nonparental hours in pre-K; has the child ever attended Head Start). The covariates included represent a mix of variables that have been identified through previous research to be related to reading performance, as well as variables that were hypothesized to be related to the outcome, such as demographic types of variables (e.g., age at kindergarten entry, socioeconomic level, gender, race, disability status, biological mother's age at first birth, non-English, Census region), cognitive-related variables (e.g., reading performance, mathematics performance, general knowledge), motor skills (fine and gross), and parent and teacher reports on the child (e.g., self-control, approaches to learning, interpersonal, externalizing problem behaviors, internalizing problem behaviors). We will work with a total of 60 covariates (40 continuous, 20 categorical) to estimate the propensity score.

12.3.1 PSA Using MatchIt in R

Next, we will use **R** for propensity score analysis, specifically MatchIt (Ho, Imai, & Stuart, 2007). As we learned, there are different conditioning strategies that can be used with propensity scores, one of which is optimal matching, and that will be illustrated here. **R** is a freely accessible downloadable open source statistical software program that can be used with Windows, Mac OS X, and Linus operating systems. **R** is the core program that runs behind various packages, one of which is MatchIt, which we will use to conduct propensity score matching. Given this, **R** must be installed first. Instructions for downloading and installing **R** can be found at <https://www.r-project.org/>. We are using version 3.2.2. Once **R** is installed, MatchIt must be installed and loaded. If you are using optimal matching (as we are in this illustration), you also need to install optmatch. The detailed steps for using MatchIt will begin at this point (i.e., we assume that **R** is installed).



R version 3.2.2 (2015-08-14) -- "Fire Safety"
 Copyright (C) 2015 The R Foundation for Statistical Computing
 Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
 You are welcome to redistribute it under certain conditions.
 Type 'license()' or 'licence()' for distribution details.

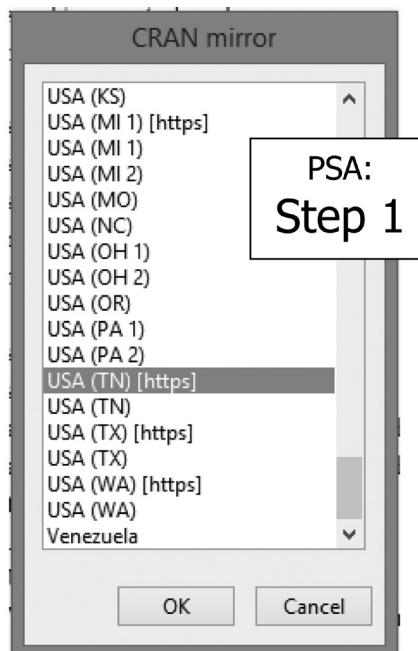
Natural language support but running in an English locale

R is a collaborative project with many contributors.
 Type 'contributors()' for more information and
 'citation()' on how to cite R or R packages in publications.

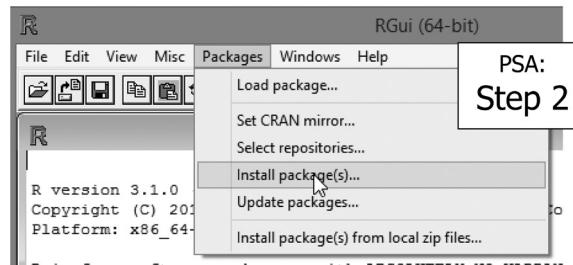
Type 'demo()' for some demos, 'help()' for on-line help, or
 'help.start()' for an HTML browser interface to help.
 Type 'q()' to quit R.

> |

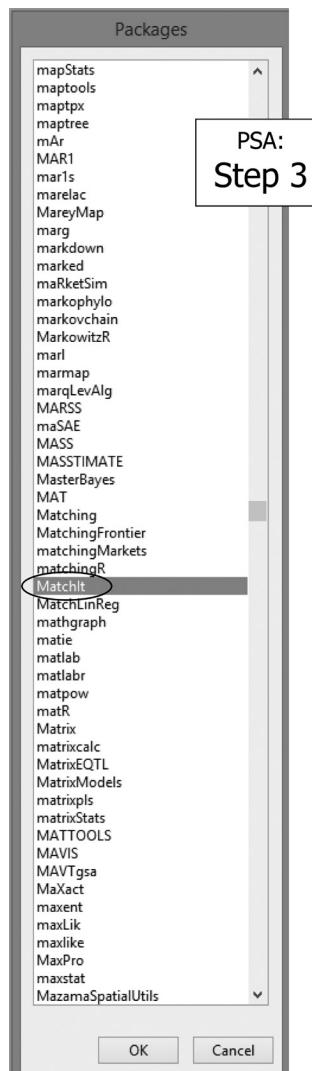
Step 1. From the CRAN Mirror dialog box, select a site from the list. Per R's recommendation (see <https://cran.r-project.org/mirrors.html>), select a CRAN mirror that is in a location close to you. For this example, I've selected the CRAN mirror located at the National Institute for Computational Sciences in Oak Ridge, TN (see screenshot Step 1). Click OK.



Step 2. To install MatchIt, open **R**. From the top menu in **R**, select **Packages** then **Install Package(s)** (see screenshot Step 2).

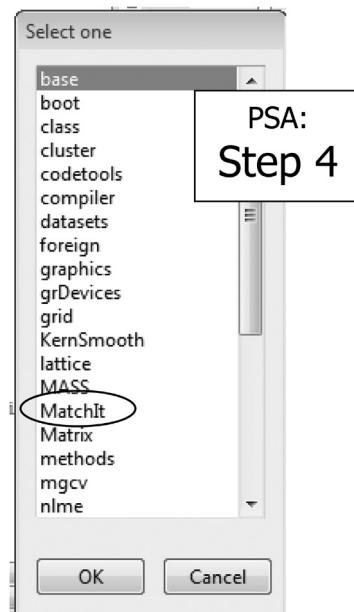


Step 3. The list that appears includes all the packages that can be loaded and run using **R**. Scroll down and click on **MatchIt** to *install* this package (see screenshot Step 3).



This same process must be done to install optmatch. For brevity, those screenshots are not presented here.

Step 4. From the Packages dropdown menu, select Load Package. From the available packages listed in the Select One window, click on MatchIt and then OK. This will *load* the MatchIt package. Please note that installation of MatchIt is done only once (i.e., the download); however, every time you want to use MatchIt, you will need to load it (see screenshot Step 4). This same process must be done to load optmatch. For brevity, those screenshots are not presented here.



Once MatchIt and optmatch load, you will receive a message in R that is similar to this:

```
– Please select a CRAN mirror for use in this session–  
trying URL 'https://mirrors.nics.utk.edu/cran/bin/windows/  
contrib/3.2/MatchIt_2.4-21.zip'  
Content type 'application/zip' length 78218 bytes (76 KB)  
downloaded 76 KB
```

```
package 'MatchIt' successfully unpacked and MD5 sums checked
```

```
The downloaded binary packages are in  
C:\Users\dhahs\AppData\Local\Temp\RtmpS20g7K\  
downloaded_packages  
> utils:::menuInstallPkgs()  
trying URL 'https://mirrors.nics.utk.edu/cran/bin/win-  
dows/contrib/3.2/optmatch_0.9-5.zip'  
Content type 'application/zip' length 973119 bytes (950 KB)  
downloaded 950 KB
```

package 'optmatch' successfully unpacked and MD5 sums checked

The downloaded binary packages are in
 C:\Users\dhahs\AppData\Local\Temp\RtmpS2Og7K\downloaded_packages

Step 5. This next step will bring your data into **R**. There are very specific data needs at this point. First, the data that are brought into MatchIt must include a binary causal variable (i.e., the grouping variable) and one or more matching variables. The causal variable is the variable on which you want the units to be matched. In this illustration, the causal variable is public or private school attendance. Second, missing values must be addressed prior to bringing the data into **R**. For convenience, we have saved the data as a .csv file and will bring in the data to **R** using the csv file. Other file formats can be brought into **R**; however, csv files have the simplest process.

The first line of the following code will read in the data and rename your data **mydata**. In this step, it is important to map the exact location of where your data is stored so that you are directing **R** to the correct folder location. The next two lines of syntax will print the variable names and the first specified number of cases in the data (10 in this example). This is an optional step but a helpful double check to make sure the data are being brought into **R** correctly.

```
mydata <- read.table("C:/Users/dhahs/Desktop/Ch12_ECLSK_N1909_private.csv", header=TRUE,
sep=",", row.names="CHILDDID")
attach(mydata)
mydata[1:10,]
```

The first 10 cases by variable are provided in R (see screenshot Step 5).

PSA:
Step 5

c7r4rscl	PRIVATE	pinumpla	plageent	wknumprk	p1hrsprk	wksesl	
1165022C	205.28	1	1	66	1	30	1.84
0598001C	198.61	1	1	62	1	35	1.94
1040003C	197.01	1	1	61	2	50	-0.34
1165018C	193.90	1	1	64	3	42	2.33
0887012C	191.47	1	2	62	6	26	2.69
0338017C	193.59	1	1	69	1	18	0.78
1165021C	207.69	1	1	63	1	24	2.12
0338019C	207.70	1	1	70	0	0	1.37
0473019C	200.35	0	1	68	2	50	1.26
0390013C	165.79	1	1	65	2	100	0.58
p1_impor	v24_a	p1_often	ciheight	ciweight	cir4rscl	cir4mscl	
1165022C	4.166667	3.666667	3.142857	46.50	46.50	37.45	46.90
0598001C	3.166667	4.000000	3.428571	46.00	45.50	50.01	32.73
1040003C	3.500000	4.000000	3.285714	44.38	53.50	39.02	29.07
1165018C	3.333333	3.333333	3.000000	45.25	50.00	36.94	30.37
0887012C	3.166667	4.000000	2.000000	44.50	68.75	41.61	24.97
0338017C	4.500000	3.333333	3.000000	48.50	49.50	47.59	46.02
1165021C	2.833333	3.333333	2.428571	47.75	55.00	47.90	47.80
0338019C	3.833333	3.333333	3.428571	44.00	37.50	31.09	26.09
0473019C	5.000000	3.000000	2.571429	45.88	44.00	34.62	24.49
0390013C	3.833333	3.666667	2.571429	46.50	62.50	37.04	30.40
c1rgscal	cifmotor	cigmotor	p1oldmom	w1momscr	w1dadscr	p1himage	p1hdage
1165022C	25.691	9	8	31	52.54	59.00	37
0598001C	28.411	6	8	37	59.00	72.10	42

Step 6. This next step is to conduct the matching. The code is presented first, followed by annotation to understand its purpose. The default estimation is logistic regression. Logit can be specified, as can other estimation methods. In this illustration, we have not specified an estimation method, thus logistic will be used.

```
m.out <- matchit(PRIVATE ~ plnumpla + plageent +
wknumprk + plhrsprk + wksesl + wkincome + p1_impor +
v24_a + p1_often + clheight + clweight + clr4rscl +
clr4mscl + clrgscal + clfmotor + clgmotor + ploldmom +
w1momscr + wldadscr + plhimage + plhdage + plnumnow +
plchlbbo + plchlaud + pllearn + plcontro + plsocial +
plimpuls + tlrarsli + tlrarsma + tlrarsge + t1learn +
t1contro + t1interp + t1extern + t1intern + plhmafb +
plagefrs + plnumsib + plless18 + r_sex0 + MOMED -
LESSBACH + DAGED_LESSBACH + r_race40 + r_race10 +
r_race20 + r_race30 + RACE_HAW_IND_MORE + r_disab0 +
r_pover0 + r_urban0 + r_urban2 + r_enrol1 + r_homla0
+r_readp0 + r_readp8 + r_mathp0+ r_work0 + r_momem2
+ PARENT + r_prekc0 + r_kprep1 + r_foods0 + r_money0 +
r_wic0 + r_afdc0 + r_fdsta0 + r_regiol + r_regio2 + r_
regio3, data = mydata, method = "optimal", ratio = 1)
summary(m.out, standardize=T)
plot(summary(m.out, standardize=T), interactive=F)
plot(m.out, type = "jitter", interactive=F)
plot(m.out, type = "hist")
plot(m.out)
m.data1 <- match.data(m.out)
write.csv(m.data1, file = "c:/Users/dhahs/Desktop/ECLSK_
matched.csv")
```

The line of code starting with `m.out` and ending with `data = mydata, method = "optimal", ratio = 1` tells **R** the causal variable (in this case, `PRIVATE`) and the covariates (`plnumpla` to `r_regio3`). The code `data = mydata` tells **R** what data to use. The last set of code in this section, specifically `method = "optimal", ratio = 1`, tells **R** to use optimal matching with a ratio of 1:1 (the ratio could be set to a different value, such as 2).

The additional code includes the following:

`summary(m.out)` This syntax produces numerical summaries to check balance. Adding `standardize=T` to this command will produce absolute standardized mean differences.

`plot(summary(m.out, standardize=T), interactive=F)` will plot absolute standardized mean differences for all covariates, before and after matching.

`plot(m.out, type = "jitter", interactive=F)` The jitter plot shows the overall distribution of propensity scores for the treated and untreated groups.

`plot(m.out, type = "hist")` This syntax produces histograms of the distance measure.

`plot(m.out)` This syntax produces plots of the distance measure and Q-Q plots for each covariate. Evidence of empirical distribution balance of treated and control groups after matching is evidenced by points that fall on the diagonal line in the Q-Q plots. The jitter plot shows the overall distribution of propensity scores for the treated and untreated groups.

The last two lines of code extract the data to a new csv file that includes only the matched units. These lines include:

```
m.data1 <- match.data(m.out)
write.csv(m.data1, file = "c:/Users/dhahs/Desktop/ECLSK_matched.csv")
```

Interpretations. Interpretations of the one-way **R** output are provided next.

Once the output is run,
pressing 'next' will advance
the graphs that appear in R.

```
R Gui (64-bit)
R Console

r_kprep1      42.3246  0.0000 41.9355  0.0000
r_foods0      94.1233  0.0000 94.1176  0.0000
r_money0      76.8281  0.0000 77.7778  0.0000
r_wic0        84.1730  0.0000 84.2105  0.0000
r_afdc0       98.3843  0.0000 100.0000 100.0000
r_fdsta0      84.5018  0.0000 86.2069  0.0000
r_regio1      80.3154  0.0000 80.3922  0.0000
r_regio2      73.1120  0.0000 74.3590  0.0000
r_regio3     -76.6093  0.0000 -75.0000  0.0000

Sample sizes:
          Control Treated
All           1573    328
Matched       656     328
Unmatched     917     0
Discarded      0      0

> plot(m.out, type = "jitter")
[1] "To identify the units, use first mouse button; to stop, use second."
integer(0)
> plot(m.out, type = "hist")
> plot(m.out)
Waiting to confirm page change...
Waiting to confirm page change...
```

Excerpts of the output are presented. Summaries of balance for all the data, for the matched data, and for the improvement (expressed in standardized mean differences) are provided.

Summary of balance for all data:

	Means Treated	Means Control	SD	Control Std.	Mean Diff.
distance	0.4114	0.1227	0.1518	1.1630	
plnumppla	1.6555	1.9870	1.1348	-0.3935	
plageent	66.6921	65.8270	4.0417	0.1971	
wknumprk	1.4364	1.2257	0.9408	0.2202	
plhrsprk	26.8810	23.6888	21.0988	0.1586	
wksesl	0.6933	0.0812	0.6895	0.9116	
wkincome	91109.7561	51606.8709	42291.3846	0.4850	
p1_impor	3.9538	3.9943	0.4818	-0.0750	
v24_a	3.2490	3.1224	0.6103	0.2273	
p1_often	2.7859	2.7424	0.4984	0.0887	
clheight	45.1475	44.5655	2.1131	0.2690	
clweight	47.1648	46.0105	8.7012	0.1296	

The mean difference (or standardized mean difference) of 'all data' relative to 'matched data' is the important information to focus on in the table. Standardized mean differences between public and private school students on the covariates *prior* to matching show dissimilarities between groups prior to matching.

Summary of balance for matched data:

	Means Treated	Means Control	SD	Control Std.	Mean Diff.
distance	0.4114	0.2504	0.1617	0.6486	
plnumppla	1.6555	1.6978	0.8602	-0.0502	
plageent	66.6921	66.3855	4.0899	0.0699	
wknumprk	1.4364	1.3777	0.9686	0.0614	
plhrsprk	26.8810	25.7463	21.2978	0.0564	
wksesl	0.6933	0.4774	0.6825	0.3215	
wkincome	91109.7561	72406.7073	48134.7262	0.2296	
p1_impor	3.9538	3.9855	0.5237	-0.0587	
v24_a	3.2490	3.2025	0.5680	0.0835	
p1_often	2.7859	2.7522	0.4561	0.0688	
clheight	45.1475	44.8981	2.1573	0.1153	
clweight	47.1648	46.7012	8.7011	0.0521	

Standardized mean differences between public and private school students on the covariates *after* matching show greatly reduced differences between groups.

The standardized percent balance improvement results.

Percent Balance Improvement:

	Std. Mean Diff.	eCDF Med	eCDF Mean	eCDF Max
distance	44.2293	47.2411	46.9100	40.0642
plnumppla	87.2471	89.1006	77.3799	69.5531
plageent	64.5617	69.5142	59.5135	45.9603
wknumprk	72.1299	93.5979	72.0122	43.9710
plhrsprk	64.4542	74.8905	44.9127	29.5192
wksesl	64.7330	56.9041	60.7645	60.2304
wkincome	52.6540	52.9687	64.5494	68.4554
p1_impor	21.6890	59.9813	31.5701	19.7667
v24_a	63.292	67.6806	70.2180	57.3098
p1_often	22.465	0.9238	-9.0716	-27.4062
clheight	57.143	64.4900	58.5699	43.6285
clweight	59.839	7.5069	52.7782	41.9143

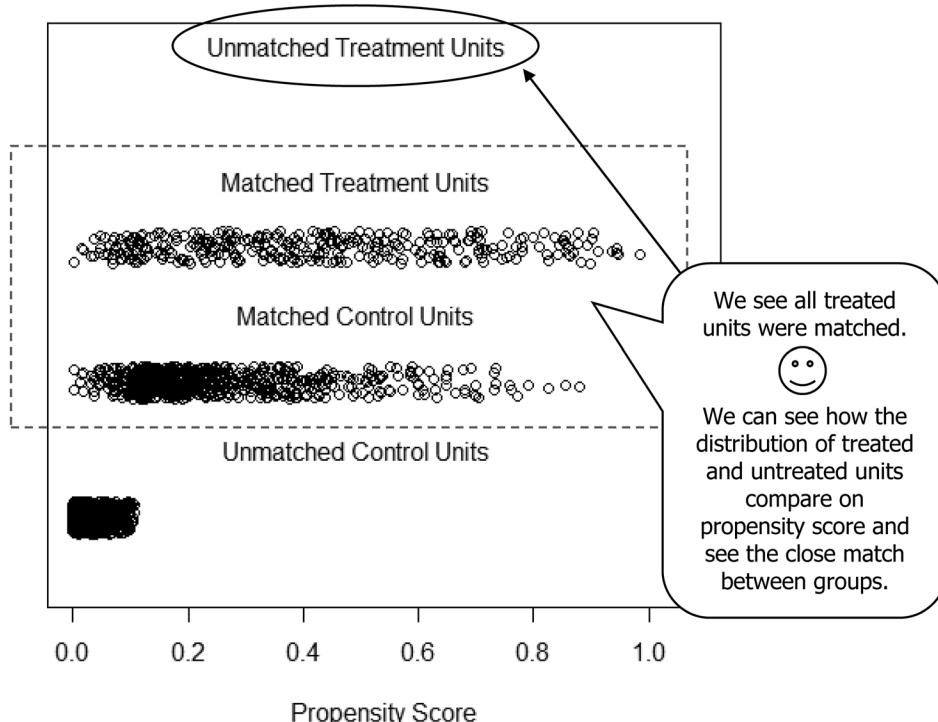
The larger the percentage balance improvement, the better the matching performed. In other words, closer to zero reflects poorer matching and closer to 100 reflects better matching.

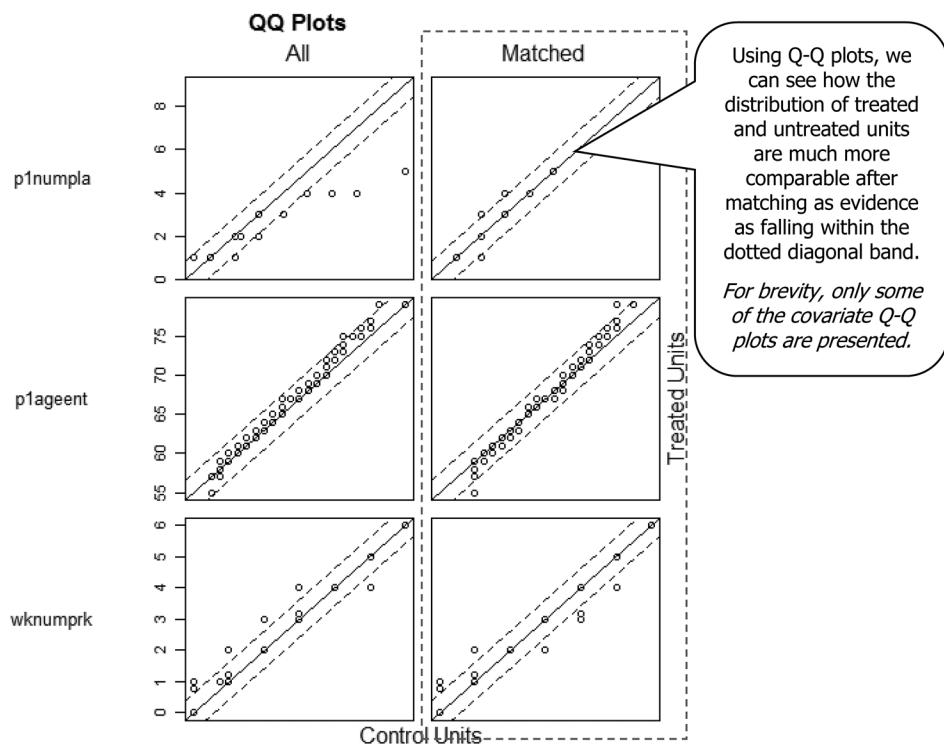
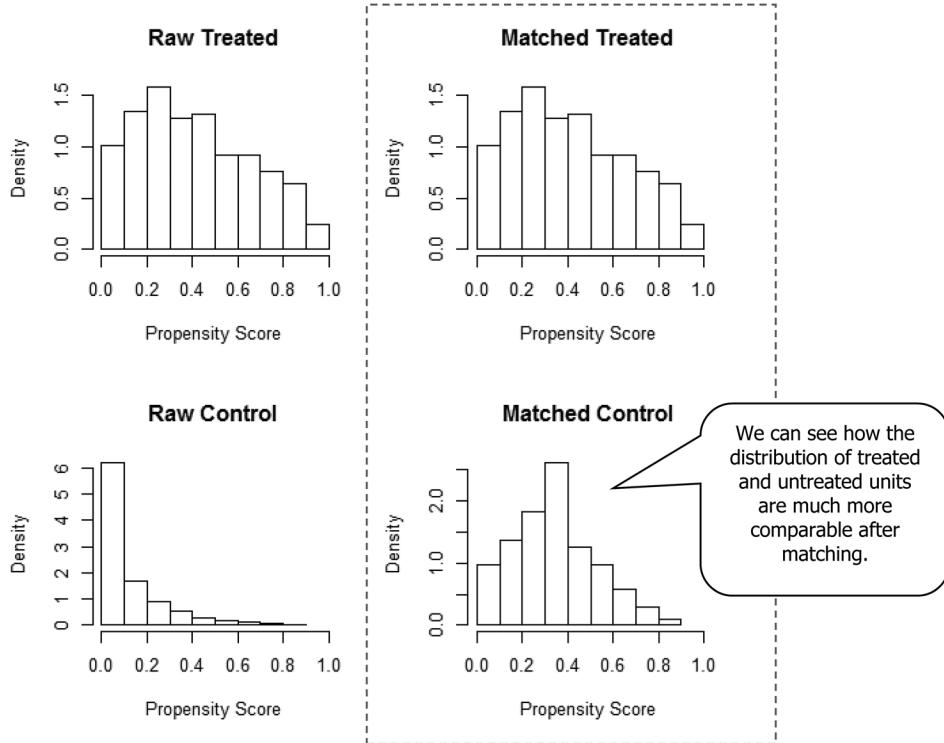
The unstandardized percent balance improvement results.

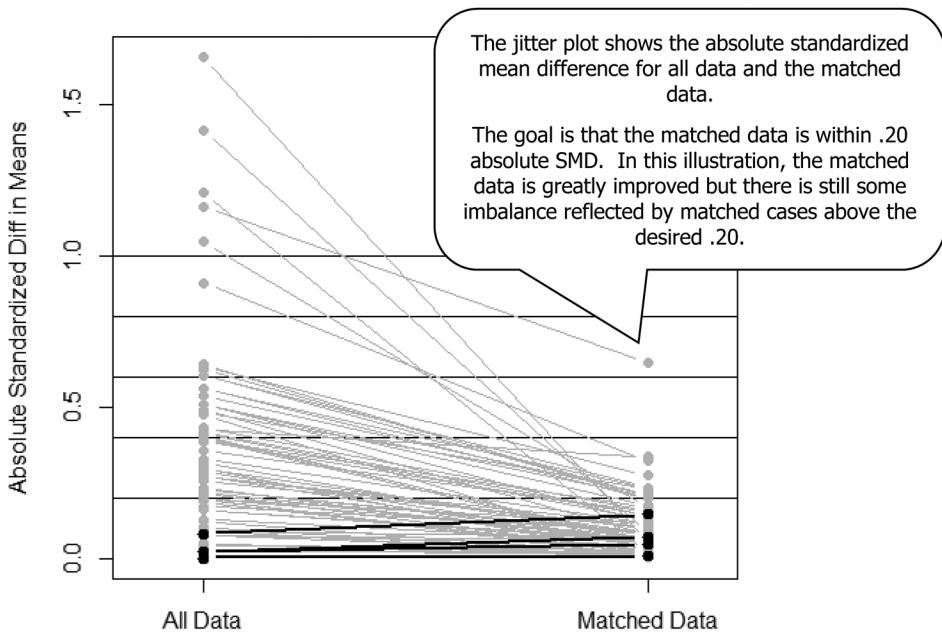
	Percent Balance Improvement:	Mean Diff.	eQQ Med	eQQ Mean	eQQ Max
distance	44.2293	41.5498	43.9107	32.6859	
plnumpla	87.2471	0.0000	82.6342	75.0000	
plageent	64.5617	100.0000	57.0447	0.0000	
wknumprk	72.1299	0.0000	66.3062	0.0000	
pihrsprk	64.4542	75.0000	42.9659	-20.0000	
wksesl	64.7330	65.1515	64.2848	61.0390	
wkincome	52.6540	65.5172	51.6835	0.0000	
p1_impor	21.6890	0.0000	29.4118	49.9999	
v24_a	63.2926	0.0000	62.4876	0.0000	
p1_often	22.4653	0.0000	-1.8637	49.9999	
clheight	57.1435	50.0000	55.8709	47.8873	
clweight	59.8368	50.0000	49.3198	-4.1667	

Median, mean, and maximum quartile differences between public and private students. Smaller Q-Q values indicate better matching.

Distribution of Propensity Scores



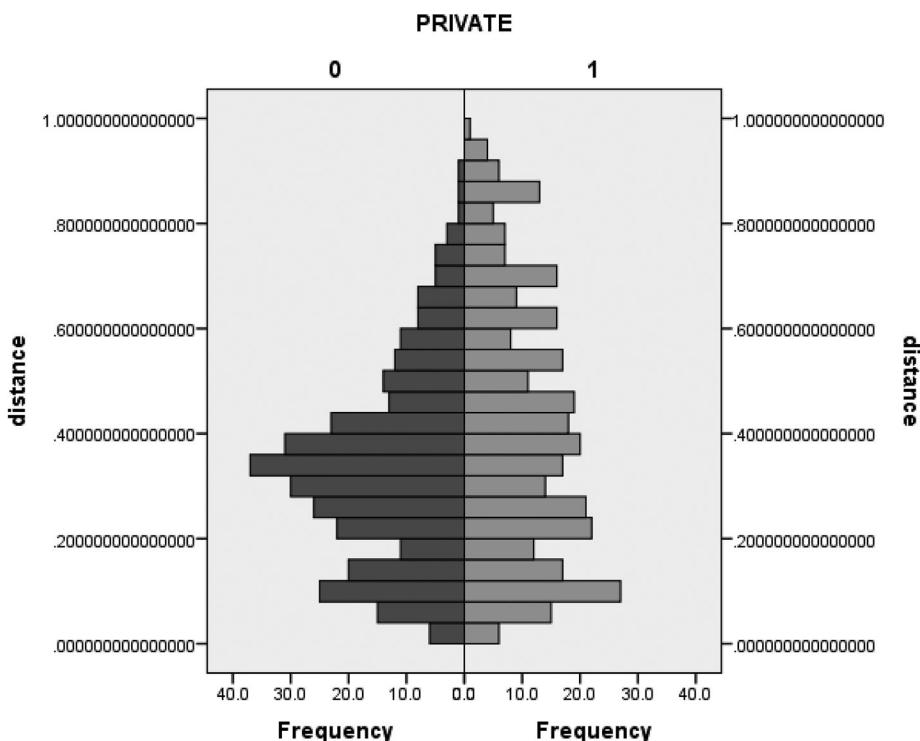




The matched data was saved to a .csv file that can be pulled into SPSS for analyzing the primary question of interest using only the cases that were matched in the propensity score process.

A	B	C	D	E	F
1	c7r4rscl	PRIVATE	p1numpla	p1ageent	wknumprj
2	1165022C	205.28	1	1	66
3	0598001C	198.61	1	1	62
4	1040003C	197.01	1	1	61
5	1165018C	193.9	1	1	64
6	0887012C	191.47	1	2	62
7	0338017C	193.59	1	1	69
8	1165021C	207.69	1	1	63
9	0338019C	207.7	1	1	70
10	0473019C	200.35	0	1	68

A histogram of the propensity scores (labeled 'distance' in the .csv file) of the matched units, created in SPSS, is another way to check the quality of matching. Distributions that are reflective of each other suggest matching has been successful.



The 'distance' created from MatchIt, which was saved in the .csv file, is the propensity score. Thus this graph reflects overlap of public and private school students after matching.

12.4 EXAMPLE WRITE-UP

In comparison to other procedures we've reviewed, propensity score analysis is a pre-processor, so to speak (akin to a marinade, as you may remember). As such, there is no research question needed for the propensity score analysis (i.e., the matching process)—and thus no research question template provided here. However, a thorough treatment in your narrative of conducting the propensity score analysis is essential.

The propensity of children to attend public versus private school was estimated using data from the ECLS-K. Previous empirical and theoretical research, as well as correlations, were used to determine the covariates for matching. Those included were predictive of reading performance and occurred prior to the outcome. [Explain the process and description of the variables included in the matching process.] Logistic regression was used to estimate the predicted probability of attending a public versus private school. Matching on the propensity scores was conducted using optimal matching in MatchIt using **R** software. The students were matched 1:1 without replacement within a caliper

of .20. This yielded a 100% matched set of children attending public versus private school ($n = 328$ public and $n = 328$ private).

Means and percentage bias were compared before and after matching. Of the 75 covariates, absolute standardized mean differences close to zero are preferable as that indicates small differences between the treated and untreated units. Thresholds of less than .10 suggest balance. In this study, absolute standardized mean differences were above .10 for about 36% of the covariates suggesting some imbalance in the model. The Q-Q plots, comparing the distribution of baseline covariates between students that attended public versus private school, also suggest some imbalance remains.

[Include jitter plot, histogram of propensity score distribution of treated to untreated, and other plots as deemed relevant to make the case that matching was successful.]

In aggregate, matching on the propensity score resulted in a matched sample such that many of the baseline covariates are very similar between students who attended public school and students who attended private school; however, some evidence of imbalance of covariates remains, suggesting that additional exploration of the matching may lead to better balanced units.

PROBLEMS

Conceptual Problems

1. Which of the following covariates should be included in estimating the propensity score?
 - a. Covariates related to the outcome
 - b. Covariates related to treatment assignment
 - c. Covariates related to both treatment assignment and outcome
 - d. Covariates unrelated to both treatment assignment and outcome
2. Parsimony is a consideration when conducting propensity score analysis.
 - a. True
 - b. False
3. The extent of overlap in the common support region is a consideration for selecting the matching algorithm.
 - a. True
 - b. False
4. Matching that entails each treatment be matched with one and only one untreated unit is which one of the following?
 - a. 1:1
 - b. 1:2
 - c. 1:many
 - d. Stratified
5. Which algorithm uses the distance between the propensity scores to assign treated to untreated?
 - a. Covariance adjustment
 - b. Matching
 - c. Stratification
 - d. Weighting

6. Which one of the following is a correct interpretation of PSA?
 - a. The conditional probability of treatment assignment given the observed covariates
 - b. The estimated probability of detecting bias in cases of nonrandom assignment
 - c. The observed probability of random assignment given covert bias
 - d. The predicted probability of the outcome given the independent variables
7. The inclusion of higher order polynomials is usually done at what stage?
 - a. After matching on the propensity score
 - b. At the point of determining treatment efficacy
 - c. When considering which theoretically important variables to include
 - d. When estimating the propensity score
8. Which one of the following is analogous to intent to treat?
 - a. Average treatment effect (ATE)
 - b. Average treatment effect for the treated (ATT)
 - c. Matching without replacement
 - d. Propensity score analysis (PSA)
9. Which one of the following is NOT a method to condition on the propensity score?
 - a. Covariance adjustment
 - b. Matching
 - c. Stratification
 - d. Weighting
10. Stratification with five strata removes 90% of the bias.
 - a. True
 - b. False

Computational Problems

1. Using the ECLS-K dataset used to illustrate the concepts in this chapter, conduct propensity score analysis using ‘private’ as the treatment group and covariates including continuous SES [WKSESL], biological mother’s age at birth [P1OLD-MOM], mother’s education less than a bachelor’s degree [MOMED_LESS-BACH], and disability status [R_DISAB0]. Use 1:1 optimal matching. Interpret the standardized mean difference and Q-Q plots in relation to balance achieved.
2. Using the ECLS-K dataset used to illustrate the concepts in this chapter, conduct propensity score analysis using ‘small town and rural’ [R_URBAN0] as the treatment group and covariates including continuous SES [WKSESL], biological mother’s age at birth [P1OLDMOM], mother’s education less than a bachelor’s degree [MOMED_LESSBACH], and southern geographical region [R_REGIO2]. Use 1:1 optimal matching. Interpret the standardized mean difference and Q-Q plots in relation to balance achieved.

Interpretative Problem

1. Use a dataset from a previous chapter in the textbook. Define a ‘treatment,’ select at least five covariates measured prior to the treatment, and conduct propensity score matching. Determine the extent to which balance was achieved.

REFERENCES

- Austin, P. C. (2011). An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate Behavioral Research*, 46(3), 399–424.
- Brookhart, M. A., Schneeweiss, S., Rothman, R. J., Glynn, R. J., Avorn, J., & Sturmer, R. (2006). Variable selection for propensity score models. *American Journal of Epidemiology*, 163(12), 1149–1156.
- Cochran, W. G. (1968). The effectiveness of adjustment by subclassification in removing bias in observational studies. *Biometrics*, 24(2), 295–313.
- Cochran, W. G., & Rubin, D. B. (1973). Controlling bias in observation studies: A review. *Sankya, Series A*, 35, 417–446.
- Gasper, J., DeLuca, S., & Estacion, A. (2012). Switching schools: Revisiting the relationship between school mobility and high school dropout. *American Educational Research Journal*, 49(3), 487–519.
- Gu, X. S., & Rosenbaum, P. R. (1993). Comparison of multivariate matching methods: Structures, distances, and algorithms. *Journal of Computational and Graphical Statistics*, 2(4), 405–420.
- Guo, S., & Fraser, M. W. (2010). *Propensity score analysis: Statistical methods and applications*. Thousand Oaks, CA: Sage Publications, Inc.
- Guo, S., & Fraser, M. W. (2015). *Propensity score analysis: Statistical methods and applications* (2nd ed.). Thousand Oaks, CA: Sage.
- Hahs-Vaughn, D. L. (2015). Propensity score analysis with complex survey samples. In W. Pan & H. Bai (Eds.), *Propensity score analysis: Fundamentals, developments, and extensions* (pp. 236–264). New York, NY: Guilford.
- Haviland, A., Nagin, D. S., & Rosenbaum, P. R. (2007). Combining propensity score matching and group-based trajectory analysis in an observational study. *Psychological Methods*, 12(3), 247–267.
- Hill, J., & Reiter, J. P. (2006). Interval estimation for treatment effects using propensity score matching. *Statistics in Medicine*, 25, 2230–2256.
- Ho, D. E., Imai, K. G., & Stuart, E. A. (2007). Matching as nonparametric preprocessing for reducing model dependence in parametric causal inference. *Political Analysis*, 15, 199–236.
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association*, 81(396), 945–960.
- Hong, G., & Raudenbush, S. W. (2005). Effects of kindergarten retention policy on children's cognitive growth in reading and mathematics. *Educational Evaluation and Policy Analysis*, 27(3), 205–224.
- Imai, K., King, G., & Stuart, E. A. (2008). Misunderstandings among experimentalists and observationalists about causal inference. *Journal of the Royal Statistical Society A*, 177, 481–502.
- Lee, B. K., Lessler, J., & Stuart, E. A. (2010). Improving propensity score weighting using machine learning. *Statistics in Medicine*, 29, 337–346.
- Little, R. J., & Rubin, D. B. (2000). Causal effects in clinical and epidemiological studies via potential outcomes: Concepts and analytical approaches. *Annual Review of Public Health*, 21, 121–145.
- Luellen, J. K., Shadish, W. R., & Clark, M. H. (2005). Propensity scores: An introduction and experimental test. *Evaluation Review*, 29(6), 530–558.
- McCaffrey, D. F., Ridgeway, G., & Morral, A. R. (2004). Propensity score estimation with boosted regression for evaluating causal effects in observational studies. *Psychological Methods*, 9(4), 403–425.

- Neyman, J. (1923). On the application of probability theory to agricultural experiments. Essay on principles. Section 9. *Annals of Agricultural Sciences*, 1-51.
- Neyman, J. (1990). On the application of probability theory to agricultural experiments. Essay on principles. Section 9. *Statistical Science*, 5(4), 465–480.
- Rosenbaum, P. R. (1987). Model-based direct adjustment. *Journal of the American Statistical Association*, 82(398), 387–394.
- Rosenbaum, P. R., & Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1), 41–55.
- Rosenbaum, P. R., & Rubin, D. B. (1984). Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American Statistical Association*, 79(387), 516–524.
- Rosenbaum, P. R., & Rubin, D. B. (1985). Constructing a control group using multivariate matched sampling methods that incorporate the propensity score. *The American Statistician*, 39(1), 33–38.
- Rubin, D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66, 688–701.
- Rubin, D. B. (1978). Bayesian inference for causal effects: The role of randomization. *Annals of Statistics*, 6, 34–58.
- Rubin, D. B. (1980). Discussion of ‘randomization analysis of experimental data in the Fisher randomization test’ by Basu. *Journal of the American Statistical Association*, 75, 591–593.
- Rubin, D. B. (1986). Which ifs have causal answers? *Journal of the American Statistical Association*, 81, 961–962.
- Rubin, D. B. (1997). Estimating causal effects from large data sets using propensity scores. *Annals of Internal Medicine*, 127, 757–763.
- Rubin, D. B. (2005). Causal inference using potential outcomes. *Journal of the American Statistical Association*, 100(469), 322–331.
- Rubin, D. B. (2008). For objective causal inference, design trumps analysis. *Annals of Applied Statistics*, 2, 808–840.
- Setoguchi, S., Schneeweiss, S., Brookhart, M. A., Glynn, R. J., & Cook, E. F. (2008). Evaluating the uses of data mining techniques in propensity score estimation: A simulation study. *Pharmacoepidemiology and Drug Safety*, 17, 546–555.
- Shadish, W. R., Clark, M. H., & Steiner, P.M. (2008). Can nonrandomized experiments yield accurate answers? A randomized experiment comparing random and nonrandom assignments. *Journal of the American Statistical Association*, 103(484), 1334–13565.
- Shadish, W. R., Cook, T. D., & Campbell, D. T. (2002). *Experimental and quasi-experimental designs for generalized causal inference*. Boston: Houston Mifflin.
- Thoemmes, F. J., & Kim, E. S. (2011). A systematic review of propensity score methods in the social sciences. *Multivariate Behavioral Research*, 46(1), 90–118.

Appendix A

AN INTRODUCTION TO MATRIX ALGEBRA

CHAPTER OUTLINE

A.1 Matrices	600
A.2 Calculations With Matrices	600
A.2.1 Matrix Addition and Subtraction	601
A.2.2 Matrix Multiplication and (Almost) Division	601
A.3 Types of Matrices	602
A.3.1 Vector	602
A.3.2 Square Matrix	602
A.3.3 Diagonal Matrix	603
A.3.4 Symmetric Matrix	603
A.3.5 Identity Matrix	603
A.3.6 Singular Matrix	604
A.4 Matrices and Multivariate Statistics	605

KEY CONCEPTS

1. Matrix
2. Dimension
3. Scalar
4. Singularity

Before embarking on a journey into multivariate statistics, it is fitting to have a general understanding of matrix algebra. Readers wishing to have a comprehensive and technical understanding of matrix algebra will likely find this chapter cursory at best, and they are referred to more focused readings on matrix algebra (e.g., Namboodiri, 1984).

Our objectives are that, by the end of this chapter, you will be able to (a) understand the purpose of matrix algebra, (b) understand basic types of matrices, and (c) understand basic matrix calculation.

A.1 MATRICES

Although many of the statistical procedures that you have learned to this point may be a bit complex and unwieldy to calculate by hand, most were still doable with basic linear algebra. Multivariate statistics, on the other hand, encompass much more complex problems and require more elaborate linear algebra, that is, matrix algebra. Matrix algebra provides a compact way to notate multiple linear equations at one time.

Before understanding matrix algebra, however, it's important to understand the basics of mathematical matrices. Visually, a matrix is a square or rectangular array of numbers that are displayed in the rows and columns. The **dimension** or order of the matrix indicates the number of rows (listed first) and number of columns (listed second). For example, the dimension of the following matrix is 2 (rows, i) x 3 (columns, j).

$$\begin{bmatrix} 2 & 22 & 11 \\ 4 & 4 & 9 \end{bmatrix}$$

The position of an entry in the matrix, A , can be referred to by its row and column as seen here. Entries are referred to as a specific value in the i th row and j th column where the i and j refer to the order of the element.

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}$$

Each entry or element in the matrix is a **scalar**, which is simply a real number. Transposing a column simply means flipping the rows and columns of the matrix. In other words, the first *row* of matrix A becomes the first *column* of the transposed matrix, A' .

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} = \begin{bmatrix} 2 & 22 & 11 \\ 4 & 4 & 9 \end{bmatrix} \Rightarrow A' = \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \\ a_{13} & a_{23} \end{bmatrix} = \begin{bmatrix} 2 & 4 \\ 22 & 4 \\ 11 & 9 \end{bmatrix}$$

A.2 CALCULATIONS WITH MATRICES

Mathematical calculations with which we are familiar, such as adding, subtracting, and multiplication, can also be performed with matrices.

A.2.1 Matrix Addition and Subtraction

Let's start with the most simple of operations: adding and subtracting. Matrices that have the same dimensions (i.e., same number of rows and columns) can be summed by adding together the elements in like positions. Both the commutative (i.e., order of matrix multiplication must be preserved) (i.e., $A + B = B + A$) and associative [i.e., $(A + B) + C = A + (B + C)$] rules hold for matrix addition. Subtraction works similarly (though not illustrated here). An example of matrix addition is presented here:

$$A = \begin{bmatrix} 2 & 22 & 11 \\ 4 & 4 & 9 \end{bmatrix}, B = \begin{bmatrix} 3 & 1 & 8 \\ 7 & 5 & 1 \end{bmatrix}$$

$$A + B = \begin{bmatrix} 2+3 & 22+1 & 11+8 \\ 4+7 & 4+5 & 9+1 \end{bmatrix} = \begin{bmatrix} 5 & 23 & 19 \\ 11 & 9 & 10 \end{bmatrix}$$

A.2.2 Matrix Multiplication and (Almost) Division

Matrices can also be multiplied. Scalar multiplication is multiplication of an entire matrix by a **scalar** (i.e., a constant value). With scalar multiplication, each element in the matrix is multiplied by the constant (i.e., a scalar), α , and the commutative law applies to scalar multiplication as well.

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix} \Rightarrow \alpha A = \begin{bmatrix} \alpha a_{11} & \alpha a_{12} & \alpha a_{13} \\ \alpha a_{21} & \alpha a_{22} & \alpha a_{23} \end{bmatrix}$$

$$2 \times A = \begin{bmatrix} 2 \times 2 & 2 \times 22 & 2 \times 11 \\ 2 \times 4 & 2 \times 4 & 2 \times 9 \end{bmatrix} = \begin{bmatrix} 4 & 44 & 22 \\ 8 & 8 & 18 \end{bmatrix}$$

Multiplication of matrices is possible in one special situation: when the number of columns of the first matrix equals the number of rows in the second matrix, termed conformable. In this situation, the first *row* of the first matrix is multiplied with the elements in the first *column* of the second matrix. Each successive row of the first matrix is multiplied with the first column of the second matrix. This process is repeated for each column in the second matrix. The resulting matrix has the number of rows of the first matrix and the number of columns as the second matrix. Matrix multiplication, unlike scalar multiplication, is not commutative but *is* associative [i.e., $(AB)C = A(BC)$].

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \end{bmatrix}, C = \begin{bmatrix} a_{11} & a_{21} \\ a_{12} & a_{22} \\ a_{13} & a_{23} \end{bmatrix}$$

$$A \times C = \begin{bmatrix} (a_{11} \times a_{11}) + (a_{12} \times a_{12}) + (a_{13} \times a_{13}) & (a_{11} \times a_{21}) + (a_{12} \times a_{22}) + (a_{13} \times a_{23}) \\ (a_{21} \times a_{11}) + (a_{22} \times a_{12}) + (a_{23} \times a_{13}) & (a_{21} \times a_{21}) + (a_{22} \times a_{22}) + (a_{23} \times a_{23}) \end{bmatrix}$$

$$\begin{aligned}
 A \times C &= \begin{bmatrix} 2 & 22 & 11 \\ 4 & 4 & 9 \end{bmatrix} \times \begin{bmatrix} 1 & 0 \\ 5 & 6 \\ 3 & 8 \end{bmatrix} \\
 &= \begin{bmatrix} (2 \times 1) + (22 \times 5) + (11 \times 3) & (2 \times 0) + (22 \times 6) + (11 \times 8) \\ (4 \times 1) + (4 \times 5) + (4 \times 3) & (4 \times 0) + (4 \times 6) + (4 \times 8) \end{bmatrix} \\
 &= \begin{bmatrix} (2 + 110 + 33) & (0 + 132 + 88) \\ (4 + 20 + 12) & (0 + 24 + 32) \end{bmatrix} = \begin{bmatrix} 145 & 220 \\ 36 & 56 \end{bmatrix}
 \end{aligned}$$

Division, as we normally think of it, is not possible with matrices. The closest proxy to division with matrices, accomplishing the same thing as arithmetic division, is multiplication by an inverse (which will be discussed in square matrices).

A.3 TYPES OF MATRICES

There are a number of different types of mathematical matrices. Some of the most common include (a) vector, (b) square, (c) diagonal, (d) symmetric, (e) identity, and (f) singular.

A.3.1 Vector

A vector is a special type of matrix in that it is represented by a single column (i.e., $n \times 1$), referred to as a column vector, or single row (i.e., $1 \times n$), referred to as a row vector. An example of a column vector is presented here:

$$A_{\text{vector}} = \begin{bmatrix} 22 \\ 4 \\ 26 \\ 12 \end{bmatrix}$$

A.3.2 Square Matrix

A matrix that has the same number of rows (i) as columns (j) is referred to as a square matrix. An example of a 3×3 square matrix is presented here:

$$A_{\text{square}} = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}$$

Determinants, denoted as $|A|$, are scalars defined by and for square matrices and are used to calculate inverted matrices. A determinant that equals zero is a singular matrix or has ‘singularity,’ meaning there is a lack of independence in the rows and columns

(i.e., one or more of the rows or columns is a linear transformation of another row or column). Singularity is a problem as the inverted matrix cannot be calculated, and this then translates to no unique solution. Singularity may be caused by modeling error. Both determinants and inverted matrices can be used with complex statistical procedures.

Some square matrices can be inverted, and these are denoted as A^{-1} . In arithmetic, a number raised to a negative indicates division by that number (e.g., $10^{-1} = \frac{1}{10}$). Thus, as stated previously, inversion is the matrix proxy for division.

A.3.3 Diagonal Matrix

A diagonal matrix results when only the diagonal of the matrix has values other than zero. For example:

$$A_{\text{diagonal}} = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 11 & 0 & 0 \\ 0 & 0 & 22 & 0 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

A.3.4 Symmetric Matrix

A matrix that is the same as its transpose is referred to as a symmetric matrix. You are likely already familiar with symmetric matrices through previous work with correlation matrices. A correlation coefficient matrix is a symmetric matrix. Variance-covariance matrices are also symmetric matrices, as are identity matrices.

In a symmetric matrix, $a_{ij} = a_{ji}$. By default, because these matrices are identical, the matrix must also be a square matrix (i.e., the same number of rows and columns). Because the off-diagonal values are all zero, every square diagonal matrix is also a symmetric matrix. If the values above the diagonal are reflective of the values below the diagonal, then a symmetric matrix will result. For example:

$$A_{\text{symmetric}} = \begin{bmatrix} \# & a & b & c \\ a & \# & d & e \\ b & d & \# & f \\ c & e & f & \# \end{bmatrix} = \begin{bmatrix} 7 & 2 & 10 & 6 \\ 2 & 9 & 4 & 3 \\ 10 & 4 & 1 & 5 \\ 6 & 3 & 5 & 8 \end{bmatrix} \Rightarrow A' = \begin{bmatrix} 7 & 2 & 10 & 6 \\ 2 & 9 & 4 & 3 \\ 10 & 4 & 1 & 5 \\ 6 & 3 & 5 & 8 \end{bmatrix}$$

A.3.5 Identity Matrix

An identity matrix is a special type of diagonal matrix where ones are on the diagonal. The remainder of the matrix are zero values, as is the case with other diagonal matrices. An example of an identity matrix is seen here:

$$I = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

A.3.6 Singular Matrix

A square matrix that has no inverse is referred to as a singular matrix. In advanced statistics, singularity errors can occur. These are quite problematic, as they indicate that a solution cannot be obtained as the matrix cannot be inverted. Singular matrices can occur in the following situations: (a) when all the values in one row or one column are zero, (b) when two rows or two columns are exactly the same, (c) when two rows or two columns are proportional, and/or (d) when one row can be written as a linear function of other rows or when one column can be written as a linear function of other columns. What does this mean in more practical terms as you move forward with your understanding of multivariate statistics? See Box A.1 for thoughts on this.

BOX A.1 SINGULARITY PRACTICAL CONSIDERATIONS

<i>What is singularity in conceptual terms?</i>	Singularity is the most extreme case of multicollinearity (i.e., extremely strong bivariate or multivariate correlations).
<i>Why is singularity a problem?</i>	Beyond the fact that singularity indicates redundancy in the independent variables, mathematically singularity means that matrix inversion is impossible.

The first two situations of singularity are relatively easy to understand. In the matrix here,

$$A_{\text{singular}} = \begin{bmatrix} 0 & 2 & 2 \\ 0 & 4 & 4 \\ 0 & 6 & 6 \end{bmatrix}$$

we see the first column is represented by all zeros, the singularity condition defined by (a). We also see that column 2 (2, 4, 6) and column 3 (2, 4, 6) are exactly the same, the singularity condition defined by (b).

A matrix that has proportionality singularity, (c), is illustrated here:

$$A_{\text{singular}} = \begin{bmatrix} 1 & 7 & 5 \\ 3 & 8 & 15 \\ 5 & 6 & 25 \end{bmatrix}$$

In this matrix, the third column can be calculated by multiplying the first row by 5.

An illustration of (d), linear function singularity, is provided in this matrix:

$$A_{\text{singular}} = \begin{bmatrix} 18 & 4 & 6 \\ 23 & 5 & 8 \\ 16 & 3 & 7 \end{bmatrix}$$

We see that the first column is actually a linear function of columns 2 and 3. More specifically, column 1 can be calculated as the sum of the value in row 3 plus the value in row 2 multiplied by 3 (e.g., $18 = 6 + (3)(4)$).

A.4 MATRICES AND MULTIVARIATE STATISTICS

The value of matrices with multivariate statistics is that multiple linear equations, as are needed with multivariate statistics, can be represented with matrix equations. For example, the following two linear equations

$$\begin{aligned} Y_1 &= a_{11}x_1 + a_{12}x_2 \\ Y_2 &= a_{21}x_1 + a_{22}x_2 \end{aligned}$$

can efficiently be arrayed in matrix form:

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}, X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}, Y = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix}$$

Where A = matrix of the values of the independent variables

X = column vector for coefficients for the variables

Y = column vector of the values of the dependent variable

PROBLEMS

Conceptual Problems

1. Which of the following best describes a matrix?
 - a. A concise way to compute multiple linear equations at one time
 - b. A square or rectangular array of numbers displayed in rows and columns
 - c. A vector with at least two rows
 - d. An array of numbers represented with ones on the diagonal

2. When is the one condition under which matrices be multiplied?
 - a. Only when the number of columns and rows is exactly the same in both matrices
 - b. When the number of columns of the first matrix equals the number of rows in the second matrix
 - c. When the number of rows of the first matrix equals the number of columns in the second matrix
 - d. When at least one of the matrices is a vector
3. Multiplication of a matrix with a constant requires which of the following two elements?
 - a. Identity matrix and diagonal matrix
 - b. Matrix and scalar
 - c. Two symmetrical matrices
 - d. Vector and singularity matrix
4. What type of matrix results when the matrix equals its transpose?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector
5. What type of matrix always has the same number of rows as columns?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector
6. What is the matrix referred to as that has all zeros in one column?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector
7. What type of matrix has only one column?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector

8. What type of matrix has zeros in all places except the diagonal?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector
9. What type of matrix has reflectivity above and below the diagonal?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector
10. What type of matrix will result when two rows are exactly the same?
 - a. Diagonal
 - b. Idempotent
 - c. Identity
 - d. Singular
 - e. Square
 - f. Symmetric
 - g. Vector

Computational Problems

1. You are given the following matrices. Add the matrices.

$$A = \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \end{bmatrix}, B = \begin{bmatrix} 9 & 8 & 7 \\ 6 & 5 & 4 \\ 3 & 2 & 1 \end{bmatrix}$$

2. You are given the following matrices. Multiply the matrices.

$$A = \begin{bmatrix} 2 & 4 \\ 6 & 8 \end{bmatrix}, B = \begin{bmatrix} 1 & 3 \\ 5 & 7 \end{bmatrix}$$

REFERENCE

Namboodiri, K. (1984). *Matrix algebra: An introduction*. Beverly Hills, CA: SAGE.



Taylor & Francis

Taylor & Francis Group

<http://taylorandfrancis.com>

Appendix B

ANSWERS TO ODD-NUMBERED CONCEPTUAL & COMPUTATIONAL QUESTIONS

CHAPTER 1

Answers to Conceptual Problems

1. d—structural equation modeling allows the examination of relationships or prediction with multiple dependent variables
3. c—multiple linear regression provides for examining a continuous outcome
5. c—multilevel linear modeling is used to analyze relationships of units nested within groups (e.g., children within preschool)
7. d—confirmatory factor analysis allows for grouping of constructs when there is strong theoretical evidence to support the relationships between variables
9. d—propensity score matching allows units to be matched; the matched groups can then be used for later inferential analyses

CHAPTER 2

Answers to Conceptual Problems

1. c—see definition
3. c— d is an effect size index, a measure of practical significance
5. c—equal sample sizes are not required
7. c—the intercept is 37000, which represents average salary when cumulative GPA is zero
9. d>null hypothesis does not consider SS values

CHAPTER 3

Answers to Conceptual Problems

1. false—the assumption of independence is the assumption that is most closely aligned with the sampling design

3. d—kurtosis statistic that is within an absolute value of 7 is one form of evidence of normality
5. c—the same plots used to screen for independence can be used to screen for homogeneity of variance; a scatterplot of residuals to predicted values; random display of points indicates evidence of the assumption being met
7. c—homoscedasticity applies to multiple linear regression and assumes the variation in scores for one continuous variable are approximately equal to the variation in scores for another continuous variable
9. d—VIF value of 30 suggests multicollinearity

CHAPTER 4

Answers to Conceptual Problems

1. b—partial correlations correlate two variables while holding constant a third
3. c—perfect prediction when the standard error = 0
5. false—adding an additional predictor can result in the same R^2
7. no— R^2 is higher when the predictors are uncorrelated
9. c—given there is theoretical support, the best method of selection is hierarchical regression
11. no—the purpose of the adjustment is to take the number of predictors into account; thus R_{adj}^2 may actually be smaller for most predictors

Answers to Computational Problems

1. intercept = 28.0952, $b_1 = .0381$, $b_2 = .8333$, $SS_{res} = 21.4294$, $SS_{reg} = 1,128.5706$, $F = 105.3292$ (reject at .01), $s^2_{res} = 5.3574$, $s(b_1) = .0058$, $s(b_2) = .1545$, $t_1 = 6.5343$ (reject at .01), $t_2 = 5.3923$ (reject at .01)
3. in order, the t values are 0.8 (not significant), 0.77 (not significant), -8.33 (significant)
5. $r_{1(2,3)} = -.0140$
7. $r_{12,3} = -.8412$, $r_{1(2,3)} = -.5047$
9. intercept = -1.2360, $b_1 = .6737$, $b_2 = .6184$, $SS_{res} = 58.3275$, $SS_{reg} = 106.6725$, $F = 15.5453$ (reject at .05), $s^2_{res} = 3.4310$, $s(b_1) = .1611$, $s(b_2) = .2030$, $t_1 = 4.1819$ (reject at .05), $t_2 = 3.0463$ (reject at .05)

CHAPTER 5

Answers to Conceptual Problems

1. c—The measurement scale of the dependent variable
3. a—Employment status (employed; unemployed, not looking for work; unemployed, looking for work) as there are more than two groups or categories
5. a—True
7. a—The log odds become larger as the odds increase from 1 to 100
9. d—Wald test (assesses significance of individual predictors)

Answers to Computational Problems

- $-2LL = 7.558$; $b_{HSGPA} = -.366$; $b_{athlete} = 22.327$; $b_{constant} = .219$; $se(b_{HSGPA}) = 1.309$; $se(b_{athlete}) = 20006.861$; odds ratio_{HSGPA} = .693; odds ratio_{athlete} < .001; Wald_{HSGPA} = .078; Wald_{athlete} = .000

CHAPTER 6

Answers to Conceptual Problems: One-Way and *k*-Way MANOVA Models

- a—with three dependent variables and one categorical independent variable, MANOVA is appropriate
- b—discriminant analysis is the procedure of choice to determine group differences after a statistically significant omnibus MANOVA is found
- b—only factor A is statistically significant ($p = .03$) at alpha of .05
- d—Wilks's lambda is a common MANOVA omnibus test
- d—per Cohen's (1988) interpretations, partial eta squared of .20 is generally considered a large effect

Answers to Conceptual Problems: Repeated Measures MANOVA

- a—with three dependent variables measured at two points in time and one categorical independent variable, repeated measures MANOVA is appropriate
- d—two outcomes measured at two points in time and one grouping (i.e., between-subjects factor) make this a viable scenario for repeated measures MANOVA
- d—spread-versus-level plots provide a helpful as a visual inspection of homogeneity of variance-covariance
- a—a nonstatistically significant Box's test provides evidence that the assumption of homogeneity of variance-covariance has been met
- d—Wilks's lambda is a common MANOVA omnibus test

Answers to Computational Problems: One-Way and *k*-Way MANOVA Models

- Use SPSS to conduct MANOVA (ECLSK_MANOVA_N1393.sav). The dependent variables are externalizing problems [C5SDQEXR] and internalizing problems [C5SDQINR]. The independent variables are male [MALE] and two-parent family structure [TWOPARENT_GR3]. Report the results for testing the assumption of homogeneity of variance-covariance matrices as well as the omnibus multivariate tests for the main effects and interaction.

Based on Box's M test, the assumption of homogeneity of variance-covariance matrices has not been met [$M = 30.01$, $F(9, 4251218.11) = 3.32$, $p < .001$].

The overall Wilks's lambda was statistically significant for the main effect for males indicating that the combined dependent variables differed, on average, between males and females, $\Lambda = .957$, $F(2, 1388) = 30.847$, $p < .001$, partial $\eta^2 = .043$.

The overall Wilks's lambda was statistically significant for the main effect for two-parent homes, indicating that the combined dependent variables differed, on average, between children from two-parents as compared to other family structures, $\Lambda = .963, F(2, 1388) = 26.775, p < .001$, partial $\eta^2 = .037$.

The overall Wilks's lambda was not statistically significant for the interaction effect for males by two-parent homes, indicating that the combined dependent variables did not differ, on average, based on the interaction of males and two-parent homes, $\Lambda = .998, F(2, 1388) = 1.168, p = .311$, partial $\eta^2 = .002$.

Answers to Computational Problems: Repeated Measures MANOVA

1. Use SPSS to conduct repeated measures MANOVA with only within-factors (ECLSK_REPEATEDMANOVA_N1344.sav). There are two repeated dependent variables, *reading* and *general knowledge*. More specifically, reading IRT scale score in fall kindergarten [C1R4RSCL] and reading IRT scale score in spring kindergarten [C2R4RSCL] and general knowledge IRT scale score in fall kindergarten [C1RGSCAL] and general knowledge IRT scale score in spring kindergarten [C2RGSCAL]. Report the results for the omnibus multivariate tests for the main effect of time as well as effect size.

The overall Wilks's lambda was statistically significant for the main effect of time indicating that the combined dependent variables differed, on average, over time, $\Lambda = .221, F(2, 1339) = 2363.267, p < .001$, partial $\eta_p^2 = .779$. Partial eta squared indicates a large effect, based on Cohen's interpretations.

CHAPTER 7

Answers to Conceptual Problems

1. a—The assumptions of logistic regression are more relaxed relative to those with discriminant analysis
3. b—Since continuous variables are not appropriate to use for dependent variables in discriminant analysis, income (measured in whole numbers) would not be an appropriate outcome
5. d—Wilks's lambda is appropriate for measuring discrimination in discriminant analysis, but not appropriate as a measure of acceptable classification
7. c—The canonical correlation represents the relationship between groups and the discriminant function, and when squared, provides a measure of proportion of variance of the discriminant function explained by groups of the dependent variable
9. a—Kappa statistics range from -1.0 to $+1.0$, with values closer to an absolute value of 1.0 suggesting stronger prediction. A recommendation, which follows Cohen's interpretations for correlations (1988), is that values of $.50$ or higher suggest strong prediction.

Answers to Computational Problems

- Wilks's lambda = .874, $p = .623$; Box's $M = .43$, $p = .929$; canonical correlation = .356; squared canonical correlation = .127; cross-validation classification rate = 40%

CHAPTER 8

Answers to Conceptual Problems

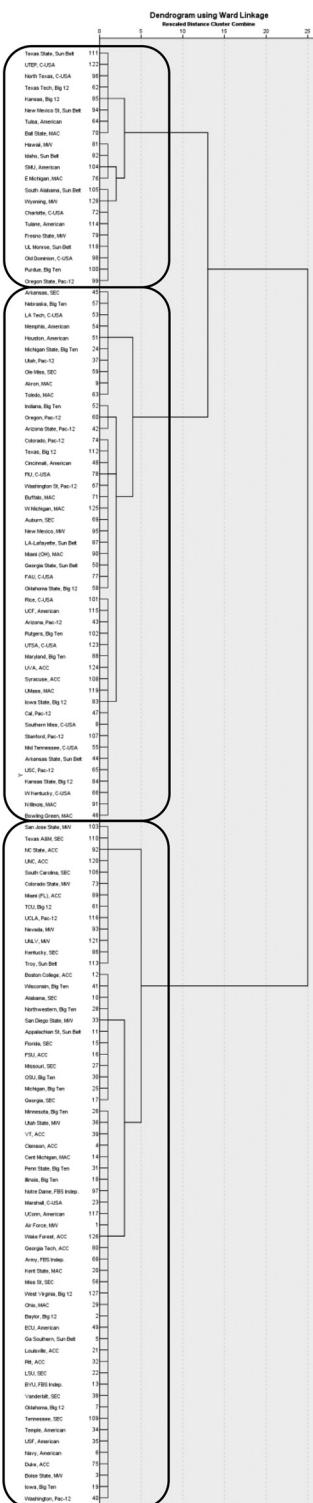
- b—cluster analysis is similar to EFA in that it is a grouping procedure
- c—single linkage is another name for nearest neighbor
- d—nonhierarchical methods begin with the researcher defining the number of clusters
- a—the first step in nonhierarchical clustering procedures is to set the cluster seed, the step where the researcher specifies the number of clusters to create
- b—clustering algorithms are the methods through which determination is made on how clusters are formed

Answers to Computational Problems

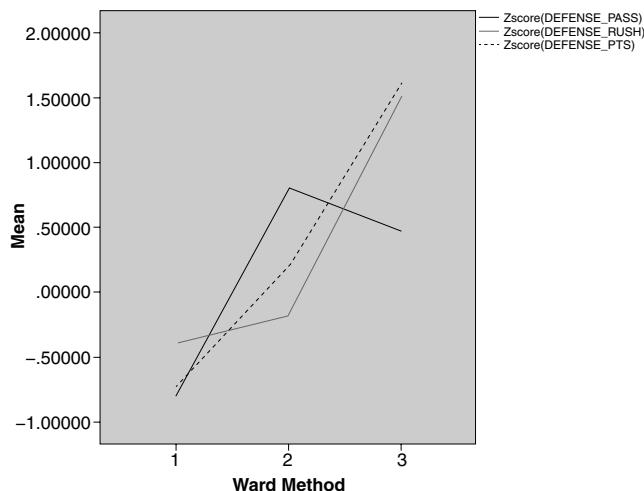
- Using the FBS_2015.sav dataset, conduct cluster analysis with the following three variables: (1) defensive passing yards, (2) defensive rushing yards, and (3) defensive points. Use Ward's method and squared Euclidean distance. Standardize the variables using z scores. Determine the best cluster solution using the agglomeration schedule and the dendrogram and categorize the clusters based on means of the clustering variables.

Based on the agglomeration schedule and dendrogram, a three- to five-cluster solution seems reasonable. These results are based on the three-cluster solution.

110	61	92	31.438	90	72	114
111	62	64	33.097	104	55	122
112	72	105	35.037	103	51	120
113	10	11	37.496	88	106	123
114	61	103	40.482	110	86	125
115	9	45	43.587	92	87	124
116	1	3	47.373	95	105	118
117	8	43	52.373	107	101	119
118	1	4	58.199	116	102	123
119	8	48	64.891	117	100	121
120	72	76	72.433	112	109	122
121	8	42	84.026	119	108	124
122	62	72	96.121	111	120	126
123	1	10	113.677	118	113	125
124	8	9	135.948	121	115	126
125	1	61	161.730	123	114	127
126	8	62	233.223	124	122	127
127	1	8	381.000	125	126	0



Examining a line graph of the standardized variable, the three cluster solution suggests: a) cluster one is low mean defense pass, rush, and points; b) cluster two is low mean defense rush yards, moderate mean defense points, and moderate mean defense pass; and c) cluster three is high mean defense rush and pass and moderate mean defense points.



CHAPTER 9

Answers to Conceptual Problems

1. a—Common factor analysis allows you to interpret meaning of the factors; principal components analysis is only a data reduction technique
3. b—Although easy to apply, Kaiser's rule has been shown to poorly estimate the number of factors
5. d—Scree plots are helpful in determining the number of factors to retain but not for determining initial factorability
7. d—the statistics can assist in making decisions in factor analysis, but these decisions must also be supported by theory
9. c—interval scales are needed for conventional factor analysis, although many researchers bend the rules to apply factor analysis to Likert scales, which are ordinal

Answers to Computational Problems

1. Using the CH9_HW1_PRESCHOOL.sav dataset, conduct exploratory factor analysis following the steps in this chapter, using maximum likelihood estimation and promax rotation. Determine initial factorability using overall MSA, Bartlett's test of sphericity, and communalities. Review the pattern and structure matrix for the initial solution, and determine the variables that appear to cluster together based on the pattern matrix.

The overall Kaiser-Meyer-Olkin measure of sampling adequacy was .811, larger than the recommended value of .50. Bartlett's test of sphericity was statistically significant [$\chi^2 (55) = 1641.885, p < .001$], suggesting factorability was appropriate. Communalities

should be above the recommended value of .30, providing evidence of shared variance among the items. In reviewing extracted communalities, four of the eleven variables were below .30, however two of these were within rounding of .30 (.290 and .295).

Both the pattern and structure matrix suggest the following two clusters of variables:

■ FACTOR 1

- Teacher rating of behavior problems: Hyperactive
- Teacher rating of behavior problems: Aggressive
- Social skills score
- Teacher rating of behavior problems: Withdrawn

■ FACTOR 2

- WJ applied problems
- PPVT ability
- WJ letter-word
- Teacher rating of behavior
- Naming colors score
- Social awareness
- Drawing score

KMO and Bartlett's Test

Kaiser-Meyer-Olkin Measure of Sampling Adequacy.	.811
Bartlett's Test of Sphericity	1641.885
df	55
Sig.	.000

Pattern Matrix^a

	Factor	
	1	2
Teacher rating of behavior problems: Hyperactive	.802	.014
Teacher rating of behavior problems: Aggressive	.792	.046
Social skills score	-.708	.060
Teacher rating of behavior problems: Withdrawn	.538	-.001
Woodcock Johnson applied problems standard score	-.005	.714
PPVT W ability score	.049	.609
Woodcock Johnson letter-word standard score	.071	.595
Teacher rating of behavior	-.082	.569
Naming colors score	-.047	.525
Social awareness	.048	.404
Drawing score	-.117	.305

Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

a. Rotation converged in three iterations.

Structure Matrix

	Factor	
	1	2
Teacher rating of behavior problems: Hyperactive	.797	-.263
Teacher rating of behavior problems: Aggressive	.776	-.228
Social skills score	-.729	.306
Teacher rating of behavior problems: Withdrawn	.538	-.188
Woodcock Johnson applied problems standard score	-.253	.716
Teacher rating of behavior	-.279	.597
PPVT W ability score	-.162	.592
Woodcock Johnson letter-word standard score	-.135	.571
Naming colors score	-.229	.541
Social awareness	-.092	.387
Drawing score	-.223	.346

Extraction Method: Maximum Likelihood.

Rotation Method: Promax with Kaiser Normalization.

a. Rotation converged in three iterations.

Initial Estimates

KMO and Bartlett's Test

Kaiser-Meyer-Olkin Measure of Sampling Adequacy.	.775
Bartlett's Test of Sphericity	Approx. Chi-Square 2490.999
df	78
Sig.	.000

CHAPTER 10

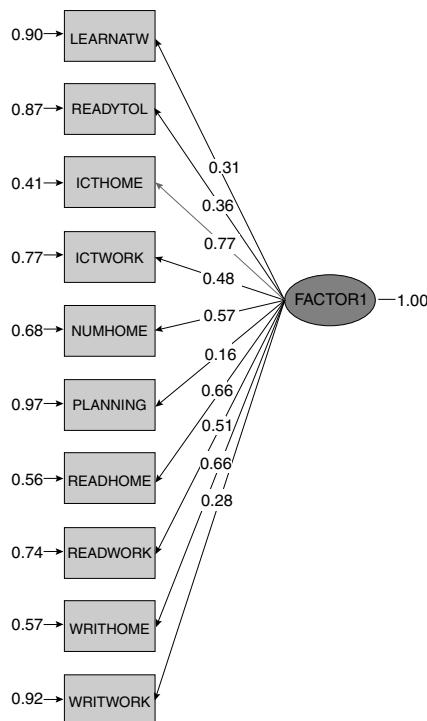
Answers to Conceptual Problems

1. d—the relationships between observed variables and latent constructs are defined in the measurement model; the structural model defines various relationships between latent constructs
3. c—rectangles are the geometric shape used in path diagrams to define observed variables
5. d—unidirectional arrows are depicted in path diagrams to specify direct effects of latent constructs on observed variables
7. c—WLS is an estimation method, *not* a goodness-of-fit index
9. a—all measurement error will be calculated as random when the correlated errors are not specified

Answers to Computational Problems

- Using the CH10_HW1_2_PIAAC_NORWAY.sav dataset, conduct confirmatory factor analysis following the steps in this chapter, using maximum likelihood estimation and testing a one-factor model. Set 'ICT skills at home' as the indicator variable. Review the standardized factor loadings. Interpret goodness-of-fit using overall chi-square, RMSEA, SRMR, and CFI. Suggest modification indices to improve model fit, if needed. (*Note: This data has been delimited to individuals who indicated their highest level of school was 'above high school' [B_Q01a_T = 3] and who were employed the year prior to completing the survey [B_Q15a = 1]. All cases with missing data on one or more index were removed from this dataset.*)

The standardized path diagram appears here. The standardized loadings represent the correlation between each observed variable and the corresponding factor. Recall that ICTHOME was the indicator variable (and is represented as the highest standardized loading in the path diagram and the grayed out unidirectional arrow).



Chi-square=324.64, df=35, p-value=0.00000, RMSEA=0.101

Select LISREL output for the one-factor model follows:

The following lines were read from file E:\CH10_HW1_2_PIAAC_NORWAY.spj:

Ch. 10 HW2 PIAAC NORWAY
Raw Data from file 'E:\CH10_HW1_2_PIAAC_NORWAY.psf'
Observed Variables: LEARNATW READYTOL ICTHOME ICTWORK
INFLUENC NUMHOME
NUMWORK PLANNING READHOME READWORK TASKDISC WRITHOME
WRITWORK
Latent Variables: FACTOR1
Relationships:
ICTHOME = 1*FACTOR1
WRITHOME = FACTOR1
NUMHOME = FACTOR1
READHOME = FACTOR1
READYTOL = FACTOR1
READWORK = FACTOR1
WRITWORK = FACTOR1
LEARNATW = FACTOR1
PLANNING = FACTOR1
ICTWORK = FACTOR1
Print residuals
Path Diagram
End of Problem

Sample Size = 820

LISREL Estimates (Maximum Likelihood)
Measurement Equations

LEARNATW = 0.45*FACTOR1, Errorvar.= 0.50, R² = 0.097

(0.055)	(0.025)
8.08	19.75

READYTOL = 0.54*FACTOR1, Errorvar.= 0.53, R² = 0.13

(0.058)	(0.027)
9.35	19.56

ICTHOME = 1.00*FACTOR1, Errorvar.= 0.18, R² = 0.59

(0.014)
13.50

ICTWORK = 0.73*FACTOR1, Errorvar.= 0.47, R² = 0.23
(0.058) (0.025)
12.57 18.86

NUMHOME = 0.80*FACTOR1, Errorvar.= 0.37, R² = 0.32
(0.054) (0.020)
14.71 18.09

PLANNING = 0.29*FACTOR1, Errorvar.= 0.84, R² = 0.027
(0.069) (0.042)
4.25 20.11

READHOME = 0.72*FACTOR1, Errorvar.= 0.18, R² = 0.44
(0.042) (0.011)
17.13 16.62

READWORK = 0.49*FACTOR1, Errorvar.= 0.19, R² = 0.26
(0.037) (0.010)
13.27 18.64

WRITHOME = 0.96*FACTOR1, Errorvar.= 0.33, R² = 0.43
(0.056) (0.020)
17.01 16.72

WRITWORK = 0.35*FACTOR1, Errorvar.= 0.37, R² = 0.080
(0.047) (0.019)
7.36 19.84

Variances of Independent Variables

Interpret goodness-of-fit using overall chi-square, RMSEA, SRMR, and CFI.

RMSEA = .10 (CI .09, .11)
 CFit $p < .001$
 SRMR = .06
 CFI = .90

RMSEA CI is larger than desired (ideally, lower bound = and upper bound $< .10$). CFit is statistically significant, suggesting unacceptable fit. SRMR suggests good model fit ($< .08$). CFI suggests poor fit ($< .95$).

FACTOR1

0.27
 (0.02)
 11.79

Goodness-of-Fit Statistics

Degrees of Freedom = 35

Minimum Fit Function Chi-Square = 288.25 ($P = 0.0$)
 Normal Theory Weighted Least Squares Chi-Square = 324.64
 ($P = 0.0$)
 Estimated Non-centrality Parameter (NCP) = 289.64
 90 Percent Confidence Interval for NCP = (235.67 ; 351.08)

Minimum Fit Function Value = 0.35
 Population Discrepancy Function Value (F0) = 0.35
 90 Percent Confidence Interval for F0 = (0.29 ; 0.43)
 Root Mean Square Error of Approximation (RMSEA) = 0.10
 90 Percent Confidence Interval for RMSEA = (0.091 ; 0.11)
 P-Value for Test of Close Fit (RMSEA < 0.05) = 0.00

Expected Cross-Validation Index (ECVI) = 0.45
 90 Percent Confidence Interval for ECVI = (0.38 ; 0.52)
 ECVI for Saturated Model = 0.13
 ECVI for Independence Model = 3.21
 Chi-Square for Independence Model with 45 Degrees of Freedom = 2607.08

Independence AIC = 2627.08
 Model AIC = 364.64
 Saturated AIC = 110.00
 Independence CAIC = 2684.17
 Model CAIC = 478.83
 Saturated CAIC = 424.01

Normed Fit Index (NFI) = 0.89
 Non-Normed Fit Index (NNFI) = 0.87
 Parsimony Normed Fit Index (PNFI) = 0.69
 Comparative Fit Index (CFI) = 0.90
 Incremental Fit Index (IFI) = 0.90
 Relative Fit Index (RFI) = 0.86

Critical N (CN) = 163.92

Root Mean Square Residual (RMR) = 0.032
 Standardized RMR = 0.064
 Goodness-of-Fit Index (GFI) = 0.93
 Adjusted Goodness-of-Fit Index (AGFI) = 0.88
 Parsimony Goodness-of-Fit Index (PGFI) = 0.59

Standardized Residuals

	LEARNATW	READYTOL	ICTHOME	ICTWORK	NUMHOME	PLANNING
LEARNATW	--					
READYTOL	2.50	--				
ICTHOME	-3.34	-0.06	--			
ICTWORK	0.08	-1.22	3.04	--		
NUMHOME	-4.14	0.66	4.36	0.02	--	
PLANNING	4.02	0.93	-5.16	1.79	-0.86	--
READHOME	0.50	0.26	-1.87	-3.74	-0.68	1.09
READWORK	5.95	0.92	-6.17	5.63	-1.34	3.73
WRITHOME	-0.93	-1.84	5.01	-4.67	-0.94	-1.47
WRITWORK	4.19	0.13	-3.39	2.25	-2.46	4.21

Standardized Residuals

	READHOME	READWORK	WRITHOME	WRITWORK
READHOME	--			
READWORK	2.72	--		
WRITHOME	3.28	-3.97	--	

WRITWORK	-1.39	5.66	-0.62
			--

Summary Statistics for Standardized Residuals

Smallest Standardized Residual = -6.17

Median Standardized Residual = 0.00

Largest Standardized Residual = 5.95

Stemleaf Plot

```

-6|2
-4|2710
-2|7435
-0|985432999761000000000000
 0|113579918
 2|257037
 4|02240679

```

Largest Negative Standardized Residuals

Residual for ICTHOME and LEARNATW	-3.34
Residual for NUMHOME and LEARNATW	-4.14
Residual for PLANNING and ICTHOME	-5.16
Residual for READHOME and ICTWORK	-3.74
Residual for READWORK and ICTHOME	-6.17
Residual for WRITHOME and ICTWORK	-4.67
Residual for WRITHOME and READWORK	-3.97
Residual for WRITWORK and ICTHOME	-3.39

Largest Positive Standardized Residuals

Residual for ICTWORK and ICTHOME	3.04
Residual for NUMHOME and ICTHOME	4.36
Residual for PLANNING and LEARNATW	4.02
Residual for READWORK and LEARNATW	5.95
Residual for READWORK and ICTWORK	5.63
Residual for READWORK and PLANNING	3.73
Residual for READWORK and READHOME	2.72
Residual for WRITHOME and ICTHOME	5.01
Residual for WRITHOME and READHOME	3.28
Residual for WRITWORK and LEARNATW	4.19
Residual for WRITWORK and PLANNING	4.21
Residual for WRITWORK and READWORK	5.66

The Modification Indices Suggest to Add an Error Covariance

Between	and	Decrease in Chi-Square	New Estimate
ICTHOME	LEARNATW	11.1	-0.04
ICTWORK	ICTHOME	9.2	0.04
NUMHOME	LEARNATW	17.1	-0.07
NUMHOME	ICTHOME	19.0	0.06

PLANNING	LEARNATW	16.2	0.09
PLANNING	ICTHOME	26.7	-0.09
READHOME	ICTWORK	14.0	-0.04
READWORK	LEARNATW	35.4	0.07
READWORK	ICTHOME	38.0	-0.05
READWORK	ICTWORK	31.6	0.06
READWORK	PLANNING	13.9	0.05
WRITHOME	ICTHOME	25.1	0.07
WRITHOME	ICTWORK	21.8	-0.07
WRITHOME	READHOME	10.7	0.04
WRITHOME	READWORK	15.7	-0.04
WRITWORK	LEARNATW	17.6	0.06
WRITWORK	ICTHOME	11.5	-0.04
WRITWORK	PLANNING	17.7	0.08
WRITWORK	READWORK	32.0	0.05

Time used: 0.016 Seconds

CHAPTER 11

Answers to Conceptual Problems

1. b—fixed slopes occur when all groups or clusters have a common rate of change (i.e., slope)
3. c—notation for the within-group variation is r_{ij}
5. a—there are four parameters tested in a model that is estimated with full maximum likelihood; these include both the fixed (i.e., coefficients) and random effects (i.e., variance and covariance components)
7. c—correlations between the intercepts and slopes can only be examined in situations where there is both a random intercept and random slope
9. d—all the indices suggest normality is reasonable.

Answers to Computational Problems

1. Using the PIAAC_merged_2percentrandom_nomiss.sav dataset, conduct multi-level linear modeling following the steps in this chapter, using restricted maximum likelihood estimation and testing a two-level random intercepts model with participants nested within country [CNTRYID]. ‘Problem solving scale score’ is the outcome. ‘Age’ [AGE] and ‘female’ [FEMALE] are the level 1 predictors. Age is group mean centered and female is uncentered.

The specified model is noted as:

$$\text{Level 1: PROBLEMS} = \beta_{0j} + \beta_{1j}(\text{AGE}) + \beta_{2j}(\text{FEMALE}) + r_{ij}$$

Level 2:

$$\begin{aligned}\beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10} \\ \beta_{2j} &= \gamma_{20}\end{aligned}$$

Compute the ICC for the null model. Interpret the fixed and random effects of the random intercepts model. Test for equal variances at level 1. Interpret goodness-of-fit of the using the deviance test.

The ICC of the null model estimated with restricted maximum likelihood is:

$$ICC = \frac{\tau_{00}}{\tau_{00} + \sigma^2} \frac{20.28}{20.28 + 1517.08} = .013$$

The average problem solving scale score is 288.09 for a male who is at their country's average age ($SE = 1.89, p < .001$). For every one year increase in age, problem solving scale score statistically significantly decreases by about 1 point (coefficient = $-0.94, SE = .07, p < .001$). Problem solving scale score statistically significantly decreases for females as well, by nearly 6 points (coefficient = $-5.67, SE = 2.00, p = .005$).

The test of homogeneity of level 1 variance indicates that equal variances can be assumed (chi-square = 8.92, $df = 14, p > .500$).

The deviance test indicates that the random intercepts model is statistically significantly better than the null model (chi-square = 161.71, $df = 2, p < .001$).

CHAPTER 12

Answers to Conceptual Problems

1. a—covariates considered for inclusion in estimating the propensity score should be related to the outcome
3. a—some matching algorithms are better suited for narrow overlap than others; for example, nearest neighbor or greedy matching would likely produce nontrivial loss of matched units in situations where the common support region is quite narrow
5. b—conditioning on a propensity score with a matching technique means that the treated and untreated units are matched based on similar/same propensity scores
7. a—generally, it is after matching on the propensity score that higher-order terms are considered and then only if there is insufficient overlap between groups after estimating the PS
9. d—should a researcher apply the propensity score as a sampling weight, matching is not applicable

Answers to Computational Problems

- Using the ECLS-K dataset used to illustrate the concepts in this chapter, conduct propensity score analysis using ‘private’ as the treatment group and covariates including continuous SES [WKSESL], biological mother’s age at birth [P1OLD-MOM], mother’s education less than a bachelor’s degree [MOMED_LESS-BACH], and disability status [R_DISAB0]. Use 1:1 optimal matching. Interpret the standardized mean difference and Q-Q plots in relation to balance achieved.

The R code for the model is:

```
m.out <- matchit(PRIVATE ~ wksesl + p1oldmom + MOMED_LESSBACH + r_disab0, data = mydata, method = "optimal",
ratio = 1)

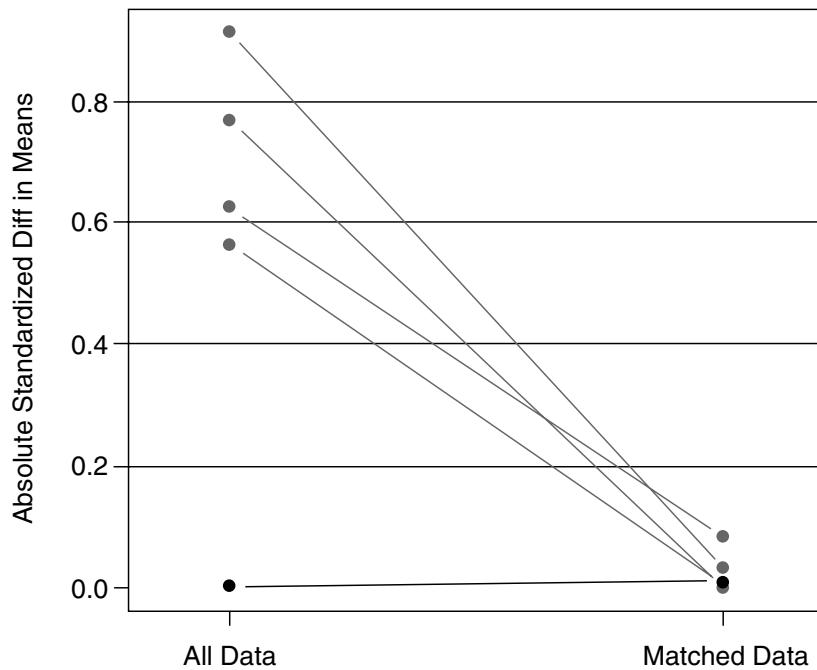
summary(m.out, standardize=T)

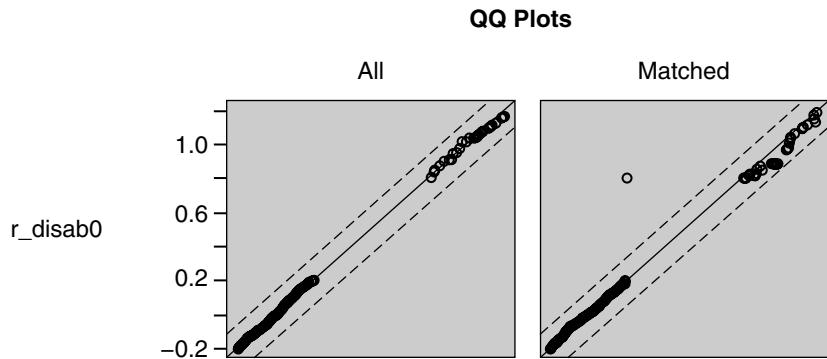
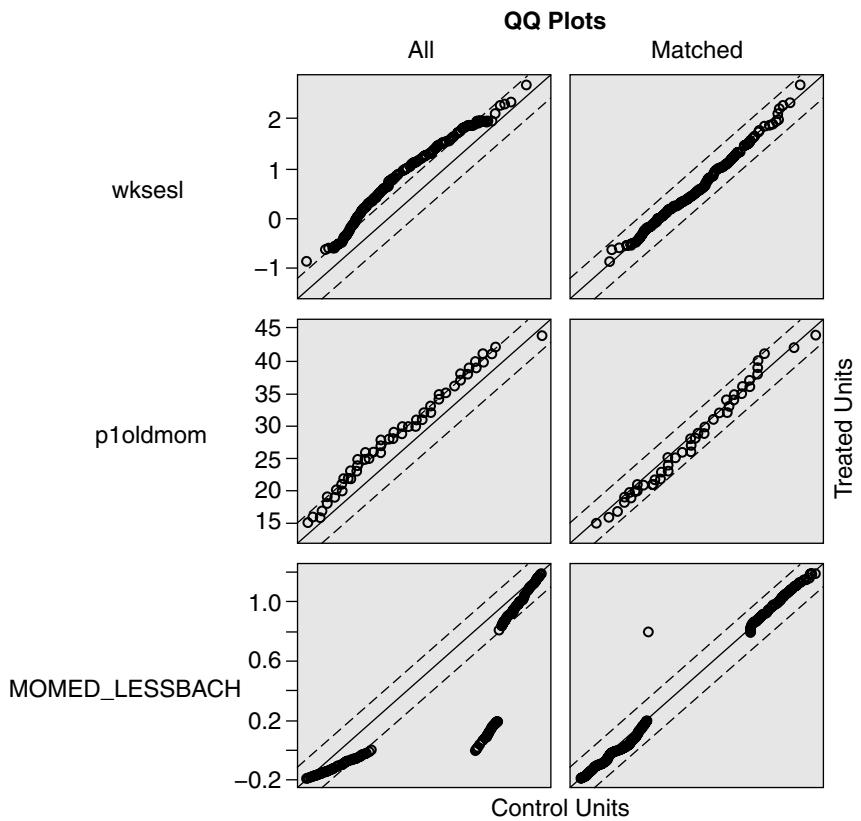
plot(summary(m.out, standardize=T), interactive=F)
plot(m.out, type = "jitter", interactive=F)
plot(m.out, type = "hist")
plot(m.out)
```

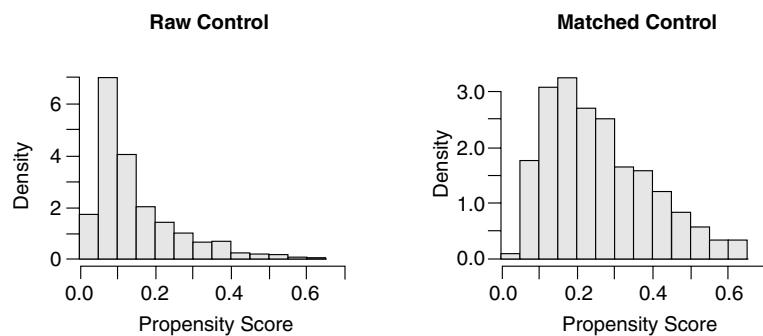
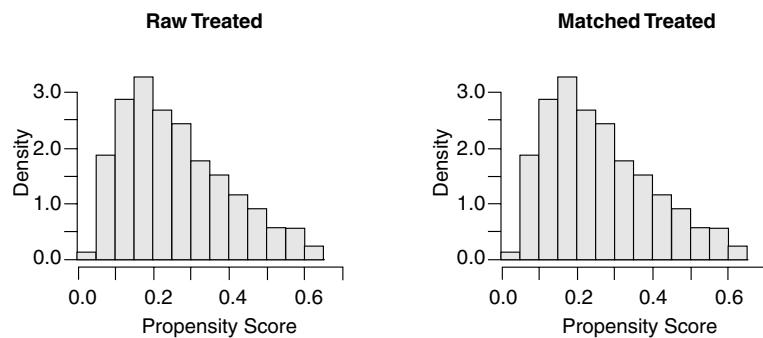
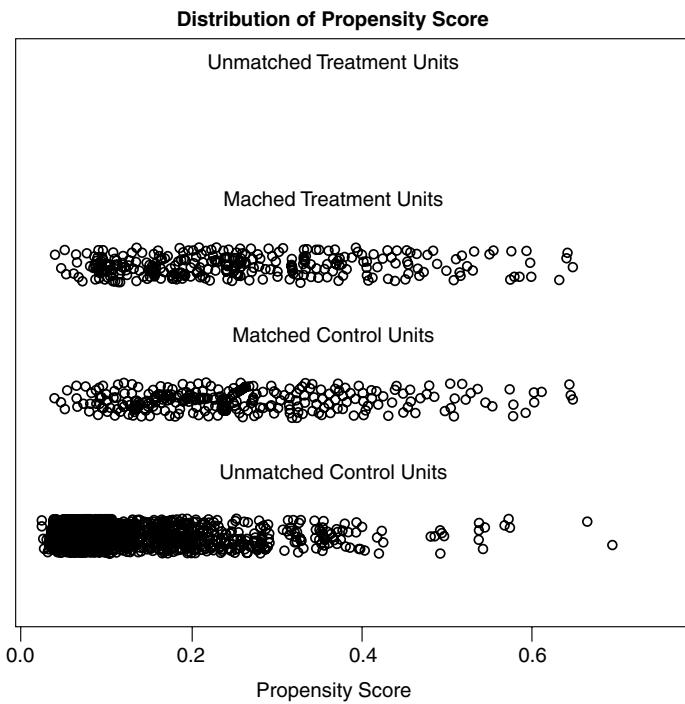
Per the output below, all covariates after matching had standardized mean differences below the recommended threshold of .10. Based on the Q-Q plots, the treated and untreated units were generally within the recommended bounds. These both suggest that there is sufficient balance.

Summary of balance for all data:						
	Means	Treated Means	Control SD	Control Std. Mean	Diff.	eCDF Med
distance	0.2593	0.1544	0.1184	0.7680	0.2783	
wksesl	0.6933	0.0812	0.6895	0.9116	0.1839	
p1oldmom	27.3228	24.0063	5.3000	0.6247	0.1688	
MOMED_LESSBACH	0.4817	0.7635	0.4251	-0.5631	0.1409	
r_disab0	0.1250	0.1252	0.3311	-0.0007	0.0001	
	eCDF	Mean	eCDF	Max		
distance	0.2512	0.4058				
wksesl	0.1918	0.4063				
p1oldmom	0.1500	0.2663				
MOMED_LESSBACH	0.1409	0.2818				
r_disab0	0.0001	0.0002				

Summary of balance for matched data:						
	Means	Treated Means	Control SD	Control Std. Mean	Diff.	eCDF Med
distance	0.2593	0.2592	0.1364	0.0011	0.0030	
wksesl	0.6933	0.6702	0.6902	0.0344	0.0122	
p1oldmom	27.3228	27.7631	4.7582	-0.0829	0.0152	
MOMED_LESSBACH	0.4817	0.4787	0.5003	0.0061	0.0015	
r_disab0	0.1250	0.1220	0.3277	0.0092	0.0015	
	eCDF	Mean	eCDF	Max		
distance	0.0025	0.0152				
wksesl	0.0158	0.0457				
p1oldmom	0.0225	0.0945				
MOMED_LESSBACH	0.0015	0.0030				
r_disab0	0.0015	0.0030				







2. Using the ECLS-K dataset used to illustrate the concepts in this chapter, conduct propensity score analysis using ‘small town and rural’ [R_URBAN0] as the treatment group and covariates including continuous SES [WKSESL], biological mother’s age at birth [P1OLDMOM], mother’s education less than a bachelor’s degree [MOMED_LESSBACH], and southern geographical region [R_REGIO2]. Use 1:1 optimal matching. Interpret the standardized mean difference and Q-Q plots in relation to balance achieved.

The R code for the model is:

```
m.out <- matchit(r_urban0 ~ wksesl + p1oldmom + MOMED_LESSBACH + r_regio2, data = mydata, method = "optimal",
ratio = 1)

summary(m.out, standardize=T)

plot(m.out, type = "jitter", interactive=F)
plot(summary(m.out, standardize=T), interactive=F)
plot(m.out, type = "hist")
plot(m.out)
```

Per the output below, only one of the four covariates after matching had a standardized mean difference below the recommended threshold of .10. Based on the Q-Q plots, histograms, and jitter plots, balance was only slightly better (if that) after matching.

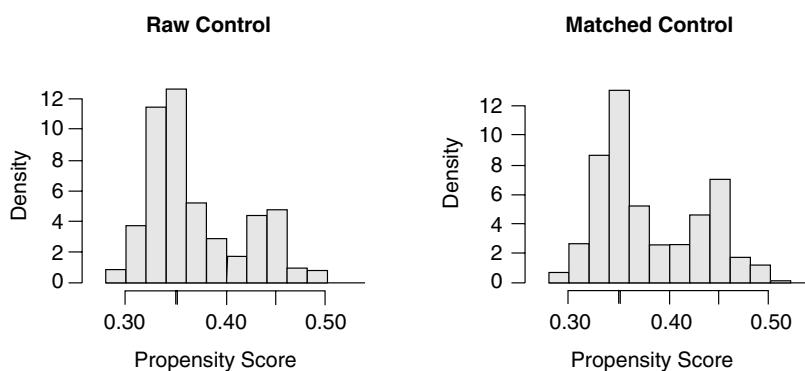
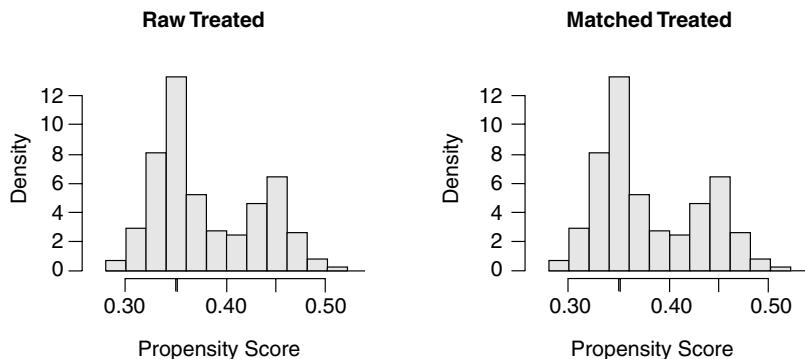
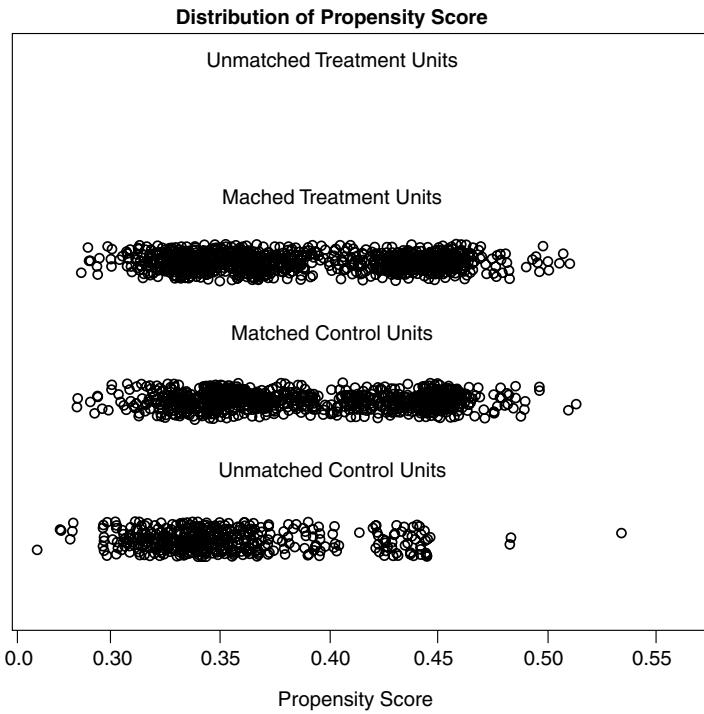
```
| Call:
matchit(formula = r_urban0 ~ wksesl + p1oldmom + MOMED_LESSBACH +
r_regio2, data = mydata, method = "optimal", ratio = 1)

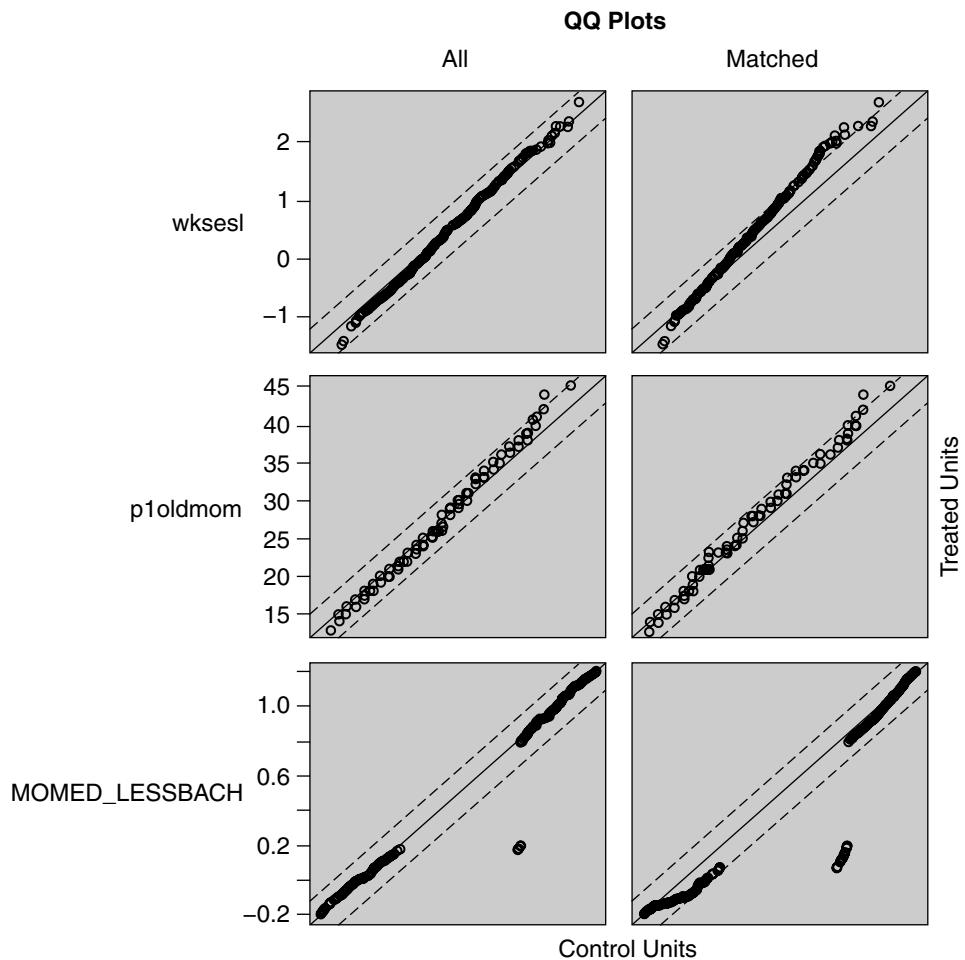
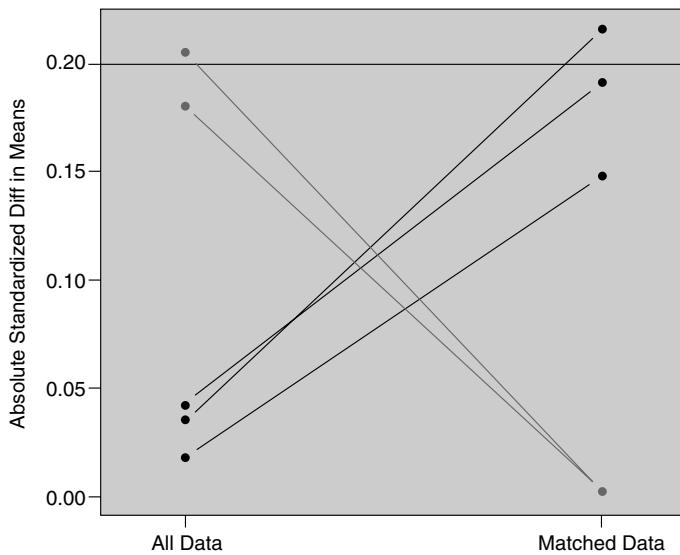
Summary of balance for all data:
      Means Treated Means Control SD Control Std. Mean Diff. eCDF Med eCDF Mean eCDF Max
distance      0.3801    0.3696   0.0482      0.2047  0.0710   0.0613   0.0947
wksesl        0.1697    0.1970   0.6929     -0.0353  0.0240   0.0303   0.0749
p1oldmom     24.5115   24.6185   5.1506     -0.0181  0.0315   0.0284   0.0605
MOMED_LESSBACH 0.7028    0.7221   0.4482     -0.0421  0.0096   0.0096   0.0193
r_regio2      0.3507    0.2645   0.4412      0.1806  0.0431   0.0431   0.0862

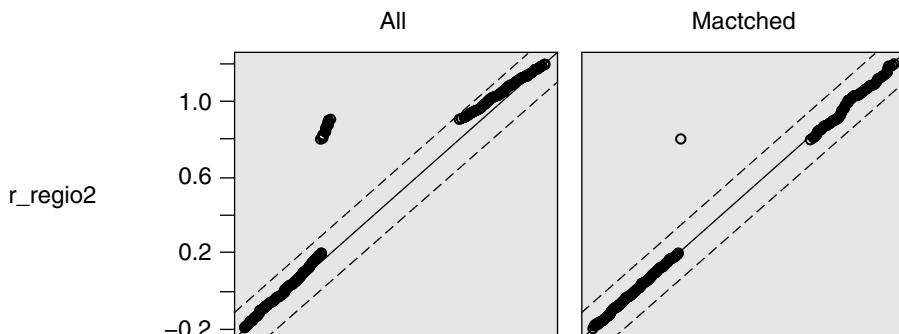
Summary of balance for matched data:
      Means Treated Means Control SD Control Std. Mean Diff. eCDF Med eCDF Mean eCDF Max
distance      0.3801    0.3799   0.0510      0.0030  0.0014   0.0025   0.0127
wksesl        0.1697    0.0025   0.6100      0.2160  0.0465   0.0500   0.1056
p1oldmom     24.5115   23.6354   5.1524      0.1482  0.0458   0.0401   0.0958
MOMED_LESSBACH 0.7028    0.7901   0.4075     -0.1909  0.0437   0.0437   0.0873
r_regio2      0.3507    0.3493   0.4771      0.0029  0.0007   0.0007   0.0014

Percent Balance Improvement:
      Std. Mean Diff. eCDF Med eCDF Mean eCDF Max
distance      98.5160  98.0154  95.8449  86.6193
wksesl       -512.4970 -94.0601 -64.7429 -40.9890
p1oldmom     -718.8974 -45.1731 -41.2768 -58.3436
MOMED_LESSBACH -353.2687 -353.2687 -353.2687 -353.2687
r_regio2      98.3665  98.3665  98.3665  98.3665

Sample sizes:
      Control Treated
All          1191    710
Matched      710     710
Unmatched    481     0
Discarded     0     0
```





QQ Plots

Summary of balance for all data:

	Means	Treated Means	Control	SD	Control Std.	Mean Diff.	eCDF Med	eCDF Mean	eCDF Max
distance	0.3801	0.3696	0.482		0.2047	0.0710	0.0613	0.0947	
wkse1	0.1697	0.1970	0.6929		-0.0353	0.0240	0.0303	0.0749	
ploldmom	24.5115	24.6185	5.1506		-0.0181	0.0315	0.0284	0.0605	
MOMED_LESSBACH	0.7028	0.7221	0.4482		-0.0421	0.0096	0.0096	0.0193	
r_regio2	0.3507	0.2645	0.4412		0.1806	0.0431	0.0431	0.0862	

Summary of balance for matched data:

	Means	Treated Means	Control	SD	Control Std.	Mean Diff.	eCDF Med	eCDF Mean	eCDF Max
distance	0.3801	0.3799	0.0510		0.0030	0.0014	0.0025	0.0127	
wkse1	0.1697	0.0025	0.6100		0.2160	0.0465	0.0500	0.1056	
ploldmom	24.5115	23.6354	5.1524		0.1482	0.0458	0.0401	0.0958	
MOMED_LESSBACH	0.7028	0.7901	0.4075		-0.1909	0.0437	0.0437	0.0873	
r_regio2	0.3507	0.3493	0.4771		0.0029	0.0007	0.0007	0.0014	

Percent Balance Improvement:

	Std. Mean Diff.	eCDF Med	eCDF Mean	eCDF Max
distance	98.5160	98.0154	95.8449	86.6193
wkse1	-512.4970	-94.0601	-64.7429	-40.9890
ploldmom	-718.8974	-45.1731	-41.2768	-58.3436
MOMED_LESSBACH	-353.2687	-353.2687	-353.2687	-353.2687
r_regio2	98.3665	98.3665	98.3665	98.3665

Sample sizes:

	Control	Treated
All	1191	710
Matched	710	710
Unmatched	481	0
Discarded	0	0

APPENDIX A**Answers to Conceptual Problems**

- b—Only when the number of columns of the first matrix equals the number of rows in the second matrix can two matrices be multiplied
- b—A constant is also referred to as a scalar; for matrix multiplication with a constant, you thus need both a matrix and a scalar
- e—Only a square matrix *always* has the same number of rows and columns
- g—A matrix with only one column or only one row is a vector
- f—A matrix that is the same above as below the diagonal is a symmetric matrix

Answers to Computational Problems

1. A 3 x 3 matrix will be generated by adding the matrices.

$$A + B = \begin{bmatrix} 10 & 12 & 14 \\ 8 & 10 & 12 \\ 6 & 8 & 10 \end{bmatrix}$$

2. A 4 x 4 matrix will be generated by multiplying the matrices.

$$A \times B = \begin{bmatrix} (2+20) & (6+28) \\ (6+30) & (24+56) \end{bmatrix} = \begin{bmatrix} 22 & 34 \\ 36 & 80 \end{bmatrix}$$

INDEX

- absolute fit indices 452
adjusted goodness-of-fit index (AGFI) 453
agglomeration schedule 343
agglomeration schedule, selection 348
agglomerative clustering 335, 338
Akaike's Information Criteria (AIC) 341, 520
Aldrich and Nelson R^2 pseudo-variance 128–9
alpha factoring 370
alternative hypothesis 10
analysis of covariance (ANCOVA) 41
analysis of variance (ANOVA) 23–5; contrasts 142; counterpart terms 171; summary table 23, 23–5; two-factor ANOVA 25–7; usage 350–1
Anderson-Rubin method 390
anti-image correlation matrix 367
approximate fit indices 452–3
a priori power, G*Power (usage): discriminant analysis 327; factorial MANOVA 255; logistic regression 162; multiple linear regression 107; omnibus test 256; repeated measures MANOVA 259
assumptions 78–82; cluster analysis 344–5; confirmatory factor analysis (CFA) 462–5; discriminant analysis 289–92; exploratory factor analysis (EFA) 378–81; ignorable treatment assignment assumption 571; k -way MANOVA models 183–6; list 83; logistic regression 134–8; MANOVA 183–7; MANOVA, violations 186; multilevel linear modeling (MLM) 525–8; one-way MANOVA models 183–6; propensity score analysis (PSA) 581; summary 82; violation 83
average linkage (between-groups linkage) 339; usage 342
average treatment effect (ATE) 580; standard estimator 573
backward elimination 69
balance 571; checks 576–7
balancing scores 575
Bartlett score, factor scores 389–90
Bartlett's test 367, 372–3, 615; examples 391, 392, 616, 617
Bayesian Information Criterion (BIC) 454
Bentler-Bonett Index (Normed Fit Index) 454
between-groups 506; linkage 339; multilevel linear modeling (MLM), effect size 524; variance, proportion reduction 523–4
BIC *See* Bayesian information criterion
binary logistic regression: execution 139; kindergarten readiness example, SPSS results 144–9
bivariate normality, confirmatory factor analysis (CFA) 463–4
box-and-whisker plot 14–16, 15
canonical correlations (discriminant analysis) 281
canonical discriminant function coefficients (discriminant analysis) 301, 308
canonical discriminant function plot, combination 285
canonical discriminant functions, summary 300–1, 307–8
CARTs *See* classification and regression trees
categorical principal components analysis (CAT-PCA) 412; Define scale and weight 416; example, SPSS results 420–3; Loading Plots 419; Normalization Method 417; output 418; output, interpretation 420, 420–3; screenshots 415–19
centered leverage values 101–2; plot 102
centering 506, 513–15; grand mean centering 513–14; group mean centering 514–15; recommendations 514–15
central tendency, measures *See* measures of central tendency
centroid method 340
centroids (discriminant analysis) 282–3; group centroids, functions 301

- CFA *See* confirmatory factor analysis
 CFI *See* comparative fit index
 CFit *See* close fit, test
 Chebychev distance 345
 CHIPCT 554–5, 563
 chi-square test 451–2, 462
 classification 274; acceptable classification (discriminant analysis) 288–9; accuracy, assessment (logistic regression) 157–60; accuracy, kappa statistic (generation) 314–15; discriminant analysis 279–80; function coefficients (discriminant analysis) 301, 309; leave-one-out classification 296; matrix 280; processing summary (discriminant analysis) 301, 308; statistics (discriminant analysis) 301–6; table, usage 145
 classification and regression trees (CARTs) 575–6
 close fit, test (CFit) 453
 cluster: algorithm, selection 339; number 342–3; single solution 350; solution, cross-validation 343–4; solution, interpretation 344
 cluster analysis 5; assumptions 344–5; characteristics 336–44; computation, SPSS (usage) 346–57; concepts 335–6, 359–60; conditions 345; data screening 358; defining/function 336–45; dendrogram 335, 342; effect size 344; example, SPSS results 354; football, statistics example 346–57; line graph 352; mathematical example 345; research question template/example write-up 358–9; sample size 344; variable selection 337; within-cluster variation 340–1
 clustering 506; agglomerative clustering 335, 338; algorithms, *K*-means clustering algorithms 335, 340–1; average linkage (between-groups linkage) 339; centroid method 340; complete linkage (furthest neighbor) 339, 351, 353; divisive clustering 335, 338; hierarchical methods 338; nonhierarchical methods 340–1; optimizing procedure 341; parallel threshold 341; procedure, selection 338–41; sequential threshold 341; single linkage (nearest neighbor) 339; two-step clustering 341; two-step clustering procedure 345; variables, selection 347, 347; Ward's method 340
 clustering methods: hierarchical clustering methods 335; partitioning clustering methods 335
 cluster randomized trials with person-level outcomes 559
 cluster similarity measure selection 338–9
 coefficients 75; canonical discriminant function coefficients 301; coefficient of multiple correlation 58, 65–6; coefficient of multiple determination 58, 65–6; discriminant analysis 320; pooled within-group regression coefficient, pure estimate 514; standardized canonical discriminant function coefficients (discriminant analysis) 300; test of significance (logistic regression) 130–2
 collinearity: diagnostics (discriminant analysis) 321; discriminant analysis 290; noncollinearity 52–4; statistics 103
 collinearity diagnostics 54; checkbox, usage 88; logistic regression example 151; multiple linear regression example 103
 common support 571, 576–7
 communality 362, 370–1; estimates 387; examples 394, 396; extraction communalities, shared/common variance 394
 comparative fit 454
 comparative fit index (CFI) 453–4; basis 454
 comparison standards (discriminant analysis) 288–9
 complete linkage (furthest neighbor) 339; examples 351, 353
 completely randomized factorial design 27
 component analysis 364
 component factor analysis 364
 confirmatory factor analysis (CFA) 5, 364, 441; assumptions 462–5; assumptions, violations 465; computation, LISREL (usage) 466–93; concepts 442, 500–2; conducting 468; continuous index measures 468; data screening 494; dataset 467–8; defining/function 442–64; differences 444; effect size 462; empirical underidentification 449; estimation 446–57; evaluation 446–57; example (LISREL results), correlated errors (usage) 486–92; example (LISREL results), correlated factor (absence) 479–85; example, LISREL results 474–7; exploratory factor analysis (EFA), comparison 444; extreme multicollinearity, absence 380; goodness-of-fit indices 451–57; heterogenous loading conditions 461; identification 446–57; independence 462–3; interpretation 446–57; just-identified model 449; linearity 463; LISREL, usage 466–81; mathematical example 465–6; Mplus, usage 493; multivariate normality 463–4; one-factor measurement model, respecification 485; one-factor solution 472–3; output, interpretation 473–92; overidentified model 449; parameter estimates 451–57, 500; post hoc power, example 495; post hoc power, G*Power (usage) 494–5; power 494–5; sample size 460–2; Save As dialog box, usage 469; second order CFA 458; second order CFA model 459; Simplis command language 472–3; SIMPLIS project file, production 471; SIMPLIS Project, selection

- 471; singularity, absence 464; SPSS file, creation 471; SPSS format, selection 468–9; statistics, generation 470; study, report information 496–7; underidentified model 449
- confirmatory factor analysis (CFA) model: estimation procedures 449–50; evaluation 451–7; identification 448–9; interpretation 451–7; model failed to converge (message) 450; modification 457–8; specification 446, 447–8
- constructs (latent constructs) 364; reliability, measurement 461
- contextual model 565
- contrast (MANOVA) 176; planned contrasts 176; post hoc contrasts 176–7
- Cook's distance 101; logistic regression 156; plot 102; saving 195, 196, 214, 214
- Cook's influence 143, 143
- correlated errors 446
- correlated residuals 446
- correlated samples *t* test 21
- correlated uniqueness 446
- correlation: coefficient of multiple correlation 58, 65–6; coefficient values 367; partial correlations 58, 59–60; Pearson product-moment correlation coefficient 28; reproduction 398; semipartial correlation (part correlation) 58, 60–61; squared multiple correlations 462; types 417
- correlation matrices 365–6; example 391, 435; inverse 392, 395
- counterfactual framework 571; Neyman-Rubin counterfactual framework 573
- covariance 27–8; variance-covariance matrices, homogeneity 42
- covariance matrices: discriminant analysis 299, 306; model-implied covariance matrix 451–2; prediction 450; representation 382; simplification 465–6
- covariance matrices, test of equality 199; discriminant analysis 299; *k*-way MANOVA models 239; one-way MANOVA models 239; repeated measures MANOVA 217, 249
- covariate selection (propensity score analysis) 576
- Cox and Snell R^2 pseudo-variance 128
- cross-level interactions, detection 522–3
- cross-products: interaction terms 577; usage 190
- cross-validation 130; cluster solution 343–4; discriminant analysis 286
- cut score (discriminant analysis) 274, 285–6
- data screening 36, 96–103; cluster analysis 358; concepts, problems 54–5; discriminant analysis 315–324; exploratory factor analysis (EFA) 427–33; *k*-way MANOVA models 227–39; logistic screening 150–60; MANOVA 227–51; multilevel linear modeling 553–9; one-way MANOVA models 227–39; repeated measures MANOVA 239–51
- decision rules 2–4; flowchart 3
- Dendrogram 335, 342; average linkage, usage 342; selection 348
- dependent samples *t* test 21
- deviances likelihood ratio test, difference 558–9
- deviance test 520; usage 540
- DFBETA: influence 143, 143; logistical regression 156–7
- DFBETA values 97; usage 102
- diagnostic plots 102
- diagonal matrix 603
- dimension 599; matrix order 600
- Dimension Reduction, selection 385, 415
- direct oblimin, oblique rotation methods 389
- discriminant analysis 5; acceptable classification 288–9; actual group, representation 302; analysis case processing summary 298, 304; assumptions 289–92; violations 292; canonical correlations 281; canonical discriminant function coefficients 301, 308; canonical discriminant functions, summary 300–1, 307–8; casewise results 296; casewise statistics 302; centroids 282–3; characteristics 277–86; classification 279–80; classification function coefficients 301, 309; classification matrix 280; classification processing summary 301, 308; classification statistics 301–6, 308–6; classification, usage 297; classify, selection 294; coefficients 320; collinearity 290; collinearity diagnostics 321; combined canonical discriminant function plot 285; combined groups 296; components, correlation 287; computation, SPSS (usage) 293–315; concepts 274–5, 331–2; conducting, steps 294–7; covariance matrices 299, 306; cross-validation 286; cut score 285–6; data screening 315–324; discriminant function 277–80; discriminant function, interpretation 280–5; discrimination 277–8; effect size 287–9; eigenvalues 280–1, 300, 307; equality of covariance matrices, test 299, 306; equality of group means, tests 298, 305; function effect size 288; group centroids, functions 301, 308; group selection 318; groups, prior probabilities 301, 309; group statistics 298, 304; highest group, representation 302; identification, scientist example data 277; independence 290–1, 316; independent variable 294; independent variables, normality 321–3; kappa statistic 289, 314–15; kappa statistic,

- interpretation 314–15; leave-one-out classification 296; linearity 290, 316–19; linearity evidence, interpretation 319; linearity (group examination), scatterplots (usage) 319; linearity, steps 317–19; logistic regression, comparison 275; MANOVA, reversal 276; mathematical example 292–3; Multivariate MANOVA: Global effects (selection) 324; multivariate normality 290–1; noncollinearity 320–2; normality, tests 322; output, interpretation 297; overall effect size 287–8; plots 283–5; Plots heading 296; pooled within-groups matrices 299, 305; post hoc power, G*Power (usage) 324–6; power 287; power, G*Power (usage) 324–7; pre self-efficacy 319; Press's *Q* 289; a priori power, G*Power (usage) 327; research question template/example write-up 328–30; sample size 286–7; separate groups 296; standardized canonical discriminant function coefficients 300, 307; standardized coefficients 278–9; standards of comparison 288; structure coefficients 282; structure matrix 300, 308; summary table 296; territorial maps 296; territorial plot 284; three-group discriminant analysis, SPSS results 305–14; two-group discriminant analysis, SPSS results 298–302; Type of power analysis 325; variance-covariance matrices, homogeneity 291, 324; Wilk's lambda 281–2, 300, 307; within-group sample size 286–7
- discrimination 274
- dispersion, measures *See* measures of dispersion
- Distance 578, 579–80
- distance matrix 335
- divisive clustering 335, 338
- doubly multivariate model (DMM) 178
- Durbin-Watson checkbox (usage) 88
- dynamic factor models 459
- ecological fallacy, concept 507
- ECVI *See* expected cross-validation index
- effect size: cluster analysis 344; confirmatory factor analysis (CFA) 462; discriminant analysis 287–9; exploratory factor analysis (EFA) 378; individual discriminant function effect size 288; *k*-way MANOVA model 181–3; measures, multilevel linear modeling (MLM) 525; multilevel linear modeling (MLM) 523–5; one-way MANOVA model 181–3; overall discriminant analysis effect size 287–8; repeated measures MANOVA 183
- eigenvalues 280–1, 362; examples 300, 307; Kaiser's Rule 372–3
- equality of covariance matrices, test 199; discriminant analysis 299, 306; *k*-way MANOVA test 239; one-way MANOVA test 239; repeated measures MANOVA 217, 249
- equality of error variances, test 201
- equality of group means, tests (discriminant analysis) 298, 305
- equality of variances, test (repeated measures MANOVA) 250
- estimated marginal means 203–4, 223–4
- Euclidean distance, equation 345
- exogenous latent constructs, usage 448
- exogenous variable 442
- expected cross-validation index (ECVI) 454
- exploratory factor analysis (EFA) 5, 342; alpha factoring 370; Anderson-Rubin method 390; anti-image matrices 391; assumptions 378–81; assumptions, violation 380, 381; Bartlett's test 367, 372–3, 391, 392; biplot variables, selection 418; case processing summary 420; characteristics 364–77; cognitive/work ability indices, correlation matrix 435; communalities 370–71, 394, 396, 435; component loadings 423; computation, SPSS (usage) 383–423; computation (continuous data), SPSS (usage) 383–412; computation (ordinal data), SPSS (usage) 412–23; concepts 362–3, 436–7; conducting 385; confirmatory factor analysis, comparison 444; correlation matrix 391; correlation matrix, inverse 392, 395; correlations transformed variables 422; data screening 427–33; defining/function 363–81; dimension reduction, selection 385; effect size 378; eigenvalues 372–3; example, SPSS results 391–8; extraction communalities, shared/common variance 394; extraction method, selection 368–9; extreme multicollinearity 432–3; factor analysis, options 390; factor analysis, rotation 389; factor correlation matrix 399; factor loadings 375–6, 435; factor matrix 397; factor model, fitting 368–71, 376–7; factor retention 371; factor retention (determination), SPSS parallel analysis (usage) 401–12; factor rotation 373–5; generalized (weighted) least squares 369; goodness-of-fit test 398; image factoring 370; independence 427; index measures 385; iteration history 421; Kaiser-Meyer-Olkin (KMO) measure 367–8, 392, 395; Kaiser's Rule 372; linearity 427–8; Mahalanobis distance 385, 432; matrices 375; maximum likelihood (ML) 369; measure of sampling adequacy (MSA) 367–8, 392; objects 423; oblique rotation 363, 374–5; ordinal variables 414; orthogonal rotation 363, 374; parallel analysis 372–3; parallel analysis, steps 401–12; pattern matrix

- 399; power 378; principal axis factoring 369–70; principal components 369; principal components analysis, contrast 364–5; reproduced correlations 398; research question template/example write-up 433–6; sample size 377–8; scree plots 371–2; singularity 432–3; specification conditions/decisions 365; sphericity, Bartlett's test 367, 392; SPSS parallel analysis results 409–12; statistical procedure 363–4, 444; structure matrix 399; total variance 396; unweighted least squares 369; variables per factor, number 373
- external validity 36–7
- extraction communalities, shared/common variance 394
- extraction method, selection 368–9
- factorability 365–8; determination 365
- factor analysis: factor scores 389; options 390; rotation 389
- factorial design 26; completely randomized factorial design 27
- factorial MANOVA: computation 191–204; dependent variables 193; global effects 251–4; interactions 254; post hoc power, G*Power (usage) 251–5; power, interactions 254; a priori power, G*Power (usage) 255
- factorial MANOVA designs 172; concepts 170
- factoring: alpha factoring 370; image factoring 370
- factor-loading matrix 381; example 382
- factor retention 363, 371; determination, SPSS parallel analysis (usage) 401–12
- factor(s) 362, 444; analytic models, factor generation 368; correlation matrix 382; correlation matrix, examples 382, 399; extraction 363, 368–9; initial factor assessment 367–8; loadings 375–6; matrix 397; model, fitting 365, 368–71, 376–7; rotation 373–5; selection 365; solution, rotation 365; structure, sample homogeneity (relationship) 366
- fixed effects 506, 508–11, 543, 551; final estimation 541–2, 547; presence 508
- fixed slopes model, HLM results 549–50, 551–2
- FML *See* full maximum likelihood
- forward selection 70
- free parameter 442
- full maximum likelihood (FML) 519, 521; default 539
- fully crossed design 26, 173
- furthest neighbor 339
- generalized (weighted) least squares 369
- General Linear Models (GLMs) 63
- GFI *See* goodness-of-fit index
- goodness-of-fit index (GFI) 451–7, 499; adjusted goodness-of-fit index (AGFI) 453; approximate fit indices 452–3; incremental fit indices 453–4; interpretation, guidelines 454–5; parsimony-adjusted indices 453–4; parsimony goodness of fit index (PGFI) 453; usage 452–97, 455–6
- grand mean centering 513–14
- greedy matching 578
- group centroids, functions (discriminant analysis) 301, 308
- group mean centering 514; recommendation 514–15
- Harrell R^2 pseudo-variance 128–9
- hierarchical ANOVA 506
- hierarchical clustering methods 335, 338
- hierarchical linear modeling (HLM): data files 533; HLM2 run, specifications 541; usage 531–52
- hierarchical models 507–8
- hierarchical regression 71–2; logistic regression usage 133
- histogram 14, 15; checkbox, usage 89
- HLM *See* hierarchical linear modeling
- homogeneity of variance 41; counterpart terms 171
- homogeneity of variance-covariance: matrices, homogeneity 42; screening 43
- homogeneity test (level-1 variance) 548, 550
- homoscedasticity 41–2, 97–8; assumption 78, 96, 108; multilevel linear modeling (MLM) 527–8; screening 42
- Hosmer and Lemeshow R^2 pseudo-variance 128
- Hosmer and Lemeshow test, contingency test 147
- Hosmer-Lemeshow goodness-of-fit test 127
- hypotheses, types 10
- hypothesis: alternative hypothesis 10; null/ statistical hypothesis 10; research hypothesis 10; scientific hypothesis 10; testing 10–11
- ICC *See* intraclass correlation coefficient
- identity matrix 603–4
- IFI *See* incremental fit index
- ignorable treatment assignment assumption 571, 573–4, 576
- image factoring 370
- incremental fit index (IFI) 453–4
- independence 36–41, 97; absence 602–3; assumption 77, 96, 108; confirmatory factor analysis (CFA) 462–3; discriminant analysis 290–1, 316; exploratory factor analysis (EFA) 427; *k*-way MANOVA models 183–4; logistic regression 155, 163–4; one-way MANOVA models 183–4; repeated measures MANOVA 239–40; screening 37–41

- independence evidence: generation 228; interpretation 40–1; interpretation (logistic regression) 155; interpretation (MANOVA) 229–30; interpretation (repeated measures MANOVA) 240
- independence of errors: assumption, logistic regression 135; evaluation (multiple linear regression) 108; logistic regression usage 136, 150
- independent samples 21; *t* test 21
- interaction effects 26, 172
- intraclass correlation coefficient (ICC) 506, 509, 512, 517; values 521–2
- inverse probability of treatment weighting (IPTW) 578, 580–1
- IPTW *See* inverse probability of treatment weighting
- jitter plot, usage 593
- Kaiser-Meyer-Olkin (KMO) measure 367–8; examples 392, 395, 616, 617
- Kaiser's Rule 372–3
- kappa statistic (discriminant analysis) 289; generation, classification accuracy usage 314–15; interpretation 314–15
- Kappa statistics 158; cell display option 159; generation 157–8; statistics option 158
- K*-means clustering algorithms 335, 340
- KMO *See* Kaiser-Meyer-Olkin
- Kolmogorov-Smirnov (K-S) test 44, 553; results 99, 231
- kurtosis: basis 429; multivariate kurtosis 45; nonzero kurtosis 44, 463–4
- k*-way MANOVA models: assumptions 183–6; assumptions, violations 186; characteristics 172–3; data screening 227–39; dependent variables, multivariate normality 184–5, 230–4; dependent variables, variance-covariance matrices (homogeneity) 185–6; file, splitting 228–9; hypotheses 173–4; independence 183–4; independence, data screening 228–30; independence evidence, interpretation 229–30; linearity 185, 238; mathematical example 188–91; multivariate normality evidence 234–5; multivariate normality evidence, interpretation 235; omnibus multivariate tests 174–5; planned/post hoc comparison procedures 176–7; power 251–5; research question template/example write-up 260–2; sample size 180; scatterplot, generation 229; univariate normality evidence 230; univariate normality evidence, interpretation 230–4; variance-covariance matrices, homogeneity 238–9; variation, partitioning 189–91
- lambda *See* Wilk's lambda: superscript 381–2
- largest root, representation 219, 219
- latent class mixture models 460
- latent constructs 364, 444, 448; exogenous latent constructs, usage 448
- latent growth model (LGM) 460
- latent variable 442
- least squares: criterion 32; estimation 32; generalized (weighted) least squares 369; unweighted least squares 369
- leave-one-out classification (discriminant analysis) 296
- level 1 residuals (multilevel linear modeling) 553
- level 1 variance, homogeneity test 548, 550
- level 2 predictor, random intercepts model 548–50; HLM results 549–50, 551–2
- level 2 residuals 553–5
- level of significance (α) 12
- Leverage values: influence 143, 143; saving 195, 196, 214, 214
- leverage values (logistic regression) 156
- LGM *See* latent growth model
- Lilliefors significance 44
- linear function singularity 605
- linearity 51; confirmatory factor analysis (CFA) 463; discriminant analysis 290, 316–19; discriminant analysis, steps 317–19; evidence, interpretation (discriminant analysis) 319; examination 554; exploratory factor analysis (EFA) 427–8; logistic regression example 152–4, 163; multilevel linear modeling (MLM) 526; multiple linear regression example 98; repeated measures MANOVA 248–9; screening 51–2
- linearity assumption: logistic regression 135–6; multiple linear regression 80, 96, 108
- Linear Structure Relations (LISREL): modification report 457; output 619–24; usage 466–93
- linkage, types 339
- LISREL *See* Linear Structure Relations
- logged odds 123; example 123
- logistic regression 4; analysis 164; analysis, assumptions 136; assumptions 134–8; casewise list 157; characteristics 121–33; classification accuracy, assessment 157–60; coefficients 151; coefficients, test of significance 130–2; collinearity diagnostics 151; computation, SPSS (usage) 139–49; concepts 118, 165–7; conditions 137–8; contrast, function 142; Cook's distance 156; cross-validation 130; data, nonseparation 137; data screening 150–60; defining/function 118–38; descriptive statistics 156; DFBETA 156–7; discriminant analysis, comparison 275; effect size 133–4; equation 121; equation variables 154; estimation

- 125–6; false negative rate 130; false positive rate 130; fixed X assumption 136; group membership, prediction 129–30; hierarchical regression 133; Hosmer-Lemeshow goodness-of-fit test 127; independence 155, 162–3; independence evidence, interpretation 155; independence of errors, assumption 136; influential points 163–4; influential points, absence 137–8; interaction term, creation 152–4; leverage values 156; linearity 152–4, 163; linearity assumption 135–6; log likelihood, change 126–7; mathematical example 138; model fit 125–6; noncollinearity 151, 163; noncollinearity assumption 135; nonzero cell counts 137; outliers 163–4; PGR_1 150; post hoc power, G*Power (usage) 160–2; power 133; power, G*Power (usage) 160–2; predictor entry, methods 132–3; predictors, group means/standard deviations 164; a priori power, G*Power (usage) 162; probability 121–2; pseudo-variance 127–9; raw data, examination 150; regression model, test of significance 126; research question template/example write-up 163–5; results 164; sample size 133; sensitivity 129; significance tests 126–30; simultaneous logistic regression 132; specificity 130; statistical project, example 118–19; stepwise logistic regression 132–3; symmetric measures 160; usage 575–6; variance inflation factor (VIF) values 163;
- logit 118, 122–5
- log likelihood, change 126–7
- log odds 122–5, 138
- Mahalanobis distances (MDIST) 101, 385, 554; exploratory factor analysis (EFA) 432; multilevel linear modeling (MLM) 554; probability plots 526; scatterplot 563; squared Mahalanobis distance, plot 379; usage 429; values, determination 380
- Mahalanobis metric matching 578
- main effects 26, 172
- manifest variables 444
- MANOVA *See* multivariate analysis of variance
- MAP *See* minimum average partial
- Matching 577–8; algorithms 571
- Matchit: installation 585; load 586–7; usage 594
- matrix (matrices) 599, 600; addition/subtraction 601; calculations 600–2; diagonal matrix 603; identity matrix 603–4; multivariate statistics 605–6; order 600; scalar matrix 599, 600; singular matrix 604–5; square matrix 602–3; symmetric matrix 603; types 602–5
- matrix algebra: calculations 600–2; concepts 599–600, 605–7; multiplication/division 601–3
- Maximum Iterations for Convergence, setting 388, 388
- maximum likelihood (ML) 369, 449–50; estimation 125–6; full maximum likelihood (FML) 519, 521; methods 519; restricted maximum likelihood (RML) 519–20, 521
- maximum likelihood method (MLM) 450
- Maximum likelihood, selection 387, 387
- MCP *See* multiple comparison procedure
- mean 19; estimated marginal means 203–4, 223–4; grand mean 203; population mean, denotation 19; sample mean, denotation 19
- measurement: error 442, 446; model path diagram, respecification 492; two-factor measurement model 447
- measure of sampling adequacy (MSA) 367–8, 392
- measures of central tendency 17–19; mean 19; median 18; mode 18; summary 19
- measures of dispersion 19–21
- median 18
- MIMIC *See* multiple indicator-multiple cause
- minimum average partial (MAP) 372–3
- ML *See* maximum likelihood
- MLM *See* maximum likelihood method; multilevel linear modeling
- MLwiN (multilevel software) 507
- MMM *See* multivariate mixed model
- mode 18
- model fit: checkbox, usage 88; evaluation, GFIs (usage) 455–6; indices 451; logistic regression 125–6; multilevel linear modeling (MLM) 506, 519–20, 558–9
- model-implied covariance matrix 451–2
- moderator variable 73
- Mplus, usage 493
- MSA *See* measure of sampling adequacy
- multilevel linear modeling (MLM) 5; Akaike's Information Criteria (AIC) 520; assumptions 525–8; Basic Settings, selection 537; Bayesian Information Criteria (BIC) 520–1; centering 506, 513–14; centering recommendations 514–15; characteristics 508–21; cluster randomized trials with person-level outcomes 559; computation, HLM (usage) 531–52; computing power, software 523; concepts 506–7, 565–7; conditions 529; contextual model 565; Create Residual File, selection 538; data screening 553–9; defining/function 507–29; Descriptive Statistics 536; deviances likelihood ratio test, difference 558–9; deviance test 520; deviance test, usage 540; effect size 523–5; effect size, between-groups 524; effect size measures

525; effect size, overall model 523–4; effect size, within-group 524; estimation methods 518–19; estimation settings 539; file, creation 535; fixed effects 508–11, 543, 551; fixed effects, final estimation 541–2, 547; framework 527; full maximum likelihood, default 539; grand mean centered options 513, 514; Graph Data options 555–6; Graph Equations options 555–6; graphing 555–7; graphing, example 556; group mean centering 513–14; HLM2 run, specifications 541; homoscedasticity 527–8; Hypothesis Testing 540, 559; intraclass correlation coefficient 512; Iteration Control 539; level 1 residual box 558; level 1 residuals 553; level-1 variance, homogeneity test 548; level 2 predictor, random intercepts model 548–50; Level-2 Residual File, selection 538, 538; level 2 residuals 553–5; linearity 526; Mahalanobis distance 554; Make MDM dialog box, usage 533, 535; mathematical example 530–1; mixed model 546; model estimation 515–18; model fit 519–20, 558–9; model misspecification 528; models 518; model specification 528; nonrandomly varying effects 508–11; normality 526–7; null model 515–16; one-way random effects ANOVA 515–16; Open mdmt file button 535; Optimal Design, steps 560; Other Settings, selection 539, 540; Outcome variable, selection 536; power 522–3; power, Optimal Design (usage) 559–61; predictor, influence 515; random effects 508–11, 528, 545; random intercepts 543; random intercepts, fixed slopes 509, 557; random intercepts, HLM results 544–5; random intercepts model 516–17, 543–5; random intercepts model, nonrandomly varying slopes effects 550–2; random intercepts, random slopes 510, 557; random intercepts random slopes model 545; randomized trials, clustering (power) 523; research question template/example write-up 561–5; restricted maximum likelihood (RML), estimation settings 539; robust standard errors 542; Run Analysis, selection 540, 543, 548, 551; running 539; sample size 521–2; sample size adjusted Bayesian Information Criteria (SBIC) 520, 521; Select MDM type dialog box 532; Stat package input, selection 532; steps 532–40; Structure of Data, default selection 535; 2-level CRT 561; uncentered options 513; uncorrelated predictors 528; Variables in residual file, label 538; variance components, final estimation 542, 545; variance, homogeneity 527; whisker plot 558 multilevel models 460 multiple comparison procedure (MCP) 25

multiple group models 458
multiple indicator-multiple cause (MIMIC) models 448, 459; flowchart 459
multiple linear regression: analysis, assumptions 77–82, 83; characteristics 59–75; computation, SPSS (usage) 87–96; concepts 58, 110–13; defining/function 58–75; effect size 76–7; post hoc power, G*Power (usage) 104–6; PRE₁ value 96; a priori power, G*Power (usage) 107; RES₁ value 96; research question, template/example write-up 107–9; SPSS step 88
multivariate analysis of variance (MANOVA) 4–5, 379; assumptions 183–7; assumptions, violation 186; characteristics 172–9; computation, SPSS (usage) 191; concepts 170, 265–7; conditions 187–8; conducting 194; contrast 176; counterpart terms 171; data screening 227–51; defining/function 170–88; effect size 181–3; mathematical example 188–91; planned contrasts 176; post hoc contrasts 176–7; power 181; research question template/example write-up 260–5; sample size 180–1; single factor designs, concepts 170; summary table 190; two-factor MANOVA, SPSS results 198–204; two-way MANOVA, partitioning variation 188
multivariate analysis of variance (MANOVA), factorial MANOVA: computation 191–204; dependent variables 193; designs, concepts 170; independent variables 193
multivariate analysis of variance (MANOVA), k-way MANOVA models: assumptions 183–6; assumptions, violations 186; characteristics 172–3; data screening 227–37; dependent variables, multivariate normality 184–5; dependent variables, variance-covariance matrices (homogeneity) 185–6; effect size 181–3; file, splitting 228–9; hypotheses 173–4; independence 183–4; independence, data screening 228–30; independence evidence, interpretation 229–30; linearity 185, 238; mathematical example 188–91; multivariate normality evidence 234–5; multivariate normality evidence, interpretation 235–7; omnibus multivariate tests 174–5; planned/post hoc comparison procedures 176–7; sample size 180; scatterplot, generation 229; univariate normality evidence 230; univariate normality evidence, interpretation 230–4; variance-covariance matrices, homogeneity 238–9; variation, partitioning 189–91
multivariate analysis of variance (MANOVA), one-way MANOVA models: assumptions 183–6; assumptions, violations 186; characteristics 172–3; data screening 227–37; dependent variables, multivariate

- normality 183–4; dependent variables, variance-covariance matrices (homogeneity) 185–6; effect size 181–3; file, splitting 228–9; hypotheses 173–4; independence 183–4; independence, data screening 228–30; independence evidence, interpretation 229–30; linearity 185, 238; mathematical example 188–91; multivariate normality evidence 234–5; multivariate normality evidence, interpretation 235–7; omnibus multivariate tests 174–5; planned/post hoc comparison procedures 176–7; sample size 180; scatterplot, generation 229; univariate normality evidence 230; univariate normality evidence, interpretation 230–4; variance-covariance matrices, homogeneity 238–9; variation, partitioning 189–91
- multivariate analysis of variance (MANOVA), repeated measures: assumptions 187; assumptions, violations 186; characteristics 177–8; computation, SPSS (usage) 207–27; data screening 239; dependent variables, multivariate normality 241–8; designs, concepts 170; effect size 183; hypothesis 178; independence 239–40; linearity 248–9; mathematical example 191; multivariate normality evidence 245; multivariate normality evidence, interpretation 246; normality, tests 242; omnibus multivariate tests 178–8; planned/post hoc comparison procedures 179; sample size 180–1; statistics 241, 242, 245; variance-covariance matrices, homogeneity 249–51
- multivariate kurtosis 45, 429
- multivariate mixed model (MMM) 178
- multivariate normality 44–5; confirmatory factor analysis (CFA) 462–4; dependent variables (MANOVA) 230–4; discriminant analysis 290–1; macro, output 49; omnibus test 45, 429; screening 47–51; SPSS macro, output 49
- multivariate normality evidence 234–5; interpretation 235–7; interpretation, repeated measures MANOVA 245–6; repeated measures MANOVA 245–6
- multivariate skewness 45, 429
- multivariate statistics 605; conceptual problems 6–8; defining 2
- multivariate test results 200
- Nagelkerke R^2 pseudo-variance
- NCP *See* noncentrality parameter
- nearest available Mahalanobis metric matching, calipers (inclusion) 579
- nearest neighbor 339, 577; propensity score analysis (PSA) 578–9; selection 349
- Neyman-Rubin counterfactual framework (causality) 573
- NFI *See* Normed Fit Index
- NNFI *See* Non-Normed Fit Index
- noncentrality parameter (NCP) 452–3
- noncollinearity 52–4; discriminant analysis 320–2; logistic regression usage 151; multiple linear regression usage 102–3; screening 53–4; VIF value 163
- noncollinearity assumption: logistic regression 135; multiple linear regression 81–2, 96, 109
- nonhierarchical methods 340–1
- nonignorable treatment assignment 574
- Non-Normed Fit Index (NNFI) 454
- nonrandomly varying effects 506, 508–11
- nonrandomly varying slopes effects (random intercepts model) 550–2
- nonrecursive path analysis example 445
- nonzero cell counts (logistic regression) 137
- nonzero cross-loadings, results 446
- nonzero kurtosis 44, 463–4; impact, absence 184
- nonzero skewness 44–5, 463–4; impact 184
- normality 43–51, 98–102; assumption 78–9, 96, 108; evidence, interpretation 98–100; multilevel linear modeling (MLM) 526–7; multivariate normality 44–5; univariate normality 43–4; violations, detection 79
- normality tests 46, 231, 430; discriminant analysis 322; K-S/S-W tests 99; repeated measures MANOVA 242
- Normalization Method 417
- Normed Fit Index (NFI) 454
- null hypothesis 10; goodness-of-fit test 398
- null hypothesis significance testing (NHST) 13
- null model, one-way random effects ANOVA 515–16
- oblique rotation 363, 374–5; methods 389
- Odds 118, 122–5; definition 122; value 122
- odds ratio (OR) 118, 133–4; Exp(B) values, equivalence 148, 148
- omnibus multivariate tests: *k*-way MANOVA models 174–5; one-way MANOVA models 174–5; repeated measures MANOVA 178–9
- one-way MANOVA models: assumptions 183–6; assumptions, violations 186; characteristics 172–3; data screening 227–39; dependent variables, multivariate normality 184–5, 230–4; dependent variables, variance-covariance matrices (homogeneity) 185–6; file, splitting 228–9; hypotheses 173–4; independence 183–4; independence, data screening 228–30; independence evidence, interpretation 229–30; linearity 185, 238; mathematical example 188–91; multivariate normality evidence 234–5; multivariate normality evidence, interpretation 235–7; omnibus multivariate tests 174–5; planned/post

- hoc comparison procedures 176–7; power 251–5; research question template/example write-up 260–2; sample size 180; scatterplot, generation 229; univariate normality evidence 230; univariate normality evidence, interpretation 230–4; variance-covariance matrices, homogeneity 238–9; variation, partitioning 189–91
- one-way random effects ANOVA 515–16; interpretation 540; model, HLM results 541–2
- Optimal Design: steps 560; usage 559–61
- optimal matching 579
- Optimal Scaling option 415
- optimizing procedure 341
- ordinary least squares (OLS) regression 118, 119; usage 516
- orthogonal rotation 363, 374
- outliers 16, 18; absence (confirmatory factor analysis) 462–4; absence (exploratory factor analysis) 428–32; absence (logistic regression) 156–7; logistic regression 163–4; methods 79; multivariate outliers, evidence 428–9
- paired samples *t* test 21
- parallel analysis (PA) 372–3; generation, SPSS syntax 403–8; output, interpretation 408–12; SPSS parallel analysis results 408–12; type, specification 402; usage 401–12
- parallel threshold 341
- parsimonious normed fit index (PNFI) 453
- parsimony-adjusted indices 453–4
- parsimony goodness of fit index (PGFI) 453
- parsimony, principle 516
- part correlation *See* semipartial correlation
- partial correlations 58, 59–60
- partial slope 62
- partitioning clustering methods 335
- path analysis 5, 441; characteristics 444–6; computation, LISREL (usage) 466–93; concepts 442, 500–2; data screening 494; defining/function 442–64; development 443; mathematical example 465–6; power 494–5; recursive/nonrecursive examples 445; research question template/example write-up 496–500
- pattern matrix 399, 616–17
- PCA *See* principal components analysis
- Pearson product-moment correlation coefficient 28
- PGFI *See* parsimony goodness of fit index
- PNFI *See* parsimonious normed fit index
- pooled within-group regression coefficient, pure estimate 514
- pooled within-groups matrices (discriminant analysis) 299, 305
- post hoc power, G*Power (usage): confirmatory factor analysis (CFA) 494–5; discriminant analysis 324–6; factorial MANOVA 251–5; logistic regression 160–2; multiple linear regression 104–6; repeated measures MANOVA 255–9
- power *See* post hoc power: comparison 200; determination 494; exploratory factor analysis (EFA) 378; *k*-way MANOVA models 251–5; logistic regression 133; MANOVA 181; multilevel linear modeling (MLM) 522–3; one-way MANOVA models 251–5; Optimal Design, usage 559–61; reduction 76; repeated measures MANOVA 255–9
- Power (1– β) 12–13
- power, G*Power (usage): discriminant analysis 324–7; logistic regression 160–2; multiple linear regression 104–7; multivariate analysis of variance (MANOVA) 251–9
- practical significance, statistical significance (contrast) 13–14
- predictors, categorical predictors 74–5
- predictors, entry: backward method 69; methods (logistic regression) 132–3; methods (multiple linear regression) 69–72
- Press's *Q* (discriminant analysis) 289
- principal axis factoring 369–70
- principal components 369
- principal components analysis (PCA) 378; exploratory factor analysis, contrast 364–5
- promax, oblique rotation methods 389
- propensity score (PS) 571; conditioning 577; distribution 591–4, 629, 631; estimation 575–6; histogram 593
- propensity score analysis (PSA) 5–6; algorithms 577, 578–9; analytic decisions 574–81; assumptions 581; average treatment effect (ATE), standard estimator 573; balance checks 576–7; balancing scores 575; causality, Neyman-Rubin counterfactual framework 573; characteristics 574–81; code, addition 588–9; common support 571, 577; computation, R (usage) 582–94; concepts 571–2, 595–6; conditions 581; covariate selection 576; CRAN mirror 584; cross-product interaction terms 577; defining/function 572–81; distance 578, 579–80; distance (creation), Matchit (usage) 594; example write-up 594–5; fine balance 579; greedy matching 578; high-order polynomials, inclusion 576–7; ignorable treatment assignment assumption 571, 573–4, 576; interpretations 589–91; inverse probability of treatment weighting (IPTW) 578, 580–1; jitter plot, usage 593; Mahalanobis metric matching 578; matching 577–8; matching, conducting 588–94; Matchit, selection 585; Matchit, usage 583–94; mathematical example 582; mean difference 590; model adequacy,

- checking 576–7; multistep procedure 575; nearest available Mahalanobis metric matching, calipers (inclusion) 579; nearest neighbor 578–9; nonignorable treatment assignment 574; optimal matching 579; output, presentation 590; packages 585; Packages dropdown menu 586; Q-Q plots, usage 592; quasi-experimental design 574; R Console 584, 587; RGui 589; sample size 581; stable unit treatment value assumption (SUTVA) 571, 574; standardized mean difference 590; stratification 580; structure 578, 580; treated/untreated units, distribution 592; treatment assignment, transformation 575–6; unmatched treatment units 591; unstandardized percent balance improvement results 591; write-up, example 594–5
- PS *See* propensity score
- PSA *See* propensity score analysis
- pseudo-variance 127–9; R^2 pseudo-variance 128–9
- quantile-quantile (Q-Q) plots 44, 46, 379, 553; examples 592, 628, 632; science knowledge pretest (discriminant analysis) 323; unstandardized residual 46; usage 79
- R^2 pseudo-variance 128–9
- R, installation 583
- R, usage 582–94; csv file, usage 587
- random coefficients model 517–18
- random effects 506, 508–11; multilevel linear modeling (MLM) 528, 545; parameters 508–9; presence 508
- random intercepts 543
- random intercepts, fixed slopes 509, 557
- random intercepts, HLM results 544–5, 546–8; Level 2 Predictor, inclusion 549–50
- random intercepts model 517–18, 543–5; level 2 predictor 548–50; nonrandomly varying slopes effects 550–2; questions 517
- random intercepts, random slopes 510, 511, 557
- randomized trials, clustering (power) 523
- random slopes model 517; HLM results 546–8
- R Console 584, 587
- recursive path analysis example 445
- regression: equations, variation 517; hierarchical regression 71–2; line, example 32; logistic regression 4; method, results 389; multiple linear regression 58; multiple regression 4; pooled within-group regression coefficient, pure estimate 514; sequential regression 69; setwise regression 71; simple linear regression 28–32; simultaneous regression 69; standardized regression model 30–1, 63–4; subsets regression 71; unstandardized regression model 61–3
- repeated measures MANOVA: assumptions 187; assumptions, violations 186; characteristics 177–8; computation, SPSS (usage) 207–24; covariance matrices, test of equality 217, 249; data screening 239–51; dependent variable, multivariate normality 241–5; dependent variable, name (inclusion) 211; designs, concepts 170; effect size 183; equality of variances, test 250; factor1, within-subject factor name (provision) 210; hypothesis 178; independence 239–40; independence evidence, interpretation 240; linearity 248–9; mathematical example 191; models, exploration 208; multivariate normality evidence 246–8; multivariate normality evidence, interpretation 246; normality, tests 242; omnibus multivariate tests 178–9; planned/post hoc comparison procedures 179; post hoc power, G*Power (usage) 255–9; post hoc power, steps 257–8; power 255–9; a priori power, G*Power (usage) 259; a priori power, within-between 259; Repeated Measures dialog box, selection 209; Repeated Measures: Post Hoc Multiple Comparisons for Observed Means dialog box, usage 213, 214; research question template, example write-up 262–5; sample size 180–1; statistics 241, 242, 245; variance-covariance matrices, homogeneity 249–51; within-between interaction 255–7; within-subject factor name 211; Within-subjects Variables (time) dialog box, usage 212, 212
- reproduced correlation matrix 398, 399
- research hypothesis 10
- research question template/example write-up: cluster analysis 358–9; discriminant analysis 328–30; exploratory factor analysis (EFA) 433–6; *k*-way MANOVA models 260–2; logistic regression 163–5; multiple linear regression 107–9; multivariate analysis of variance (MANOVA) 260–5; one-way MANOVA models 260–2; repeated measures MANOVA 262–5
- residual chi-square statistic 145
- restricted maximum likelihood (RML) 519, 521–2; estimation settings 539
- RML *See* restricted maximum likelihood
- RMR *See* root mean square residual
- root mean square error of approximation (RMSEA) 452–3, 462; results 473, 499–500
- root mean square residual (RMSR) (RMR) 378, 452
- rotation: oblique rotation 363; orthogonal rotation 363
- SABIC *See* sample-size adjusted Bayesian information criterion
- sample size: cluster analysis 344; confirmatory factor analysis (CFA) 460–2; exploratory factor analysis (EFA) 377–8; *k*-way

- MANOVA models 180; logistic regression 133; MANOVA 180–1; multilevel linear modeling (MLM) 521–2; multiple linear regression 76; one-way MANOVA models 180; propensity score analysis (PSA) 581; recommendations 460–1
- sample size adjusted Bayesian Information Criteria (SBIC) 519–20
- sampling adequacy, Kaiser-Meyer-Olkin measure 367
- Satorra-Bentler scaled chi-square, production 450
- SBIC *See* sample size adjusted Bayesian Information Criteria
- scatterplot (scattergram) 16; examples 17, 37; generation 39; generation (MANOVA) 229; simple example 40; usage (discriminant analysis) 319
- scientific hypothesis 10
- scree plots 371–2
- scree plot, selection 387
- second order CFA 458; model 459
- Select MDM dialog box, usage 532
- SEM *See* structural equation modeling
- semipartial correlation (part correlation) 58, 60–61
- sensitivity 571; logistic regression, relationship 129
- sequential regression 69; procedures, commentary 72
- sequential threshold 341
- setwise regression 71
- Shapiro-Wilk (SW) test 44, 46, 553; output 231; results 99, 231
- significance level (α) 12
- simple linear regression 28–32; model 73
- Simplis command language 471–3
- simultaneous logistic regression 132
- simultaneous regression 69
- single factor designs (MANOVA), concepts 170
- single linkage (nearest neighbor) 339
- singularity 599; absence, confirmatory factor analysis (CFA) 463–4; considerations 604; exploratory factor analysis (EFA) 432–3; linear function singularity 605; presence 602–3
- singular matrix 604–5
- skewness: multivariate skewness 45; nonzero skewness 44–5, 463–4; statistics, usage 79
- specification search 457
- specificity, logistic regression (relationship) 130
- sphericity: Bartlett's test 367, 392; Mauchly's test 199, 220
- Spread *versus* level plots checkbox, selection 197, 198, 216, 216
- square matrix 602–3
- SRMR *See* standardized root mean square residual
- stable unit treatment value assumption (SUTVA) 571, 574
- standard deviation 20
- standard error of estimate 68
- standard error of the difference between two means *See* mean
- standardized canonical discriminant function coefficients (discriminant analysis) 300, 307
- standardized regression coefficients 58
- standardized regression model 30–1, 63–4
- Standardized residuals, saving 195, 196, 214, 214
- standardized root mean square residual (SRMR) 452
- standards of comparison (discriminant analysis) 288–9
- statistical decision table 11
- statistical hypothesis 10
- Statistical Package for Social Sciences (SPSS): file, creation 471; macro 44–5, 49; parallel analysis results 408–12; Regression option, selection 87; results (binary logistic regression kindergarten readiness example) 144–9; results (cluster analysis example) 354–9; results (exploratory factor analysis) 420–3; statistics syntax editor 48, 50–1; syntax 403–8
- Statistical Package for Social Sciences (SPSS) usage: cluster analysis 346–57; discriminant analysis 293–315; exploratory factor analysis 383–423; exploratory factor analysis (EFA) 391–400; logistic regression 139–49; MANOVA 191; multiple linear regression 87–96; repeated measures MANOVA 207–27
- statistical significance, practical significance (contrast) 13–14
- stepwise logistic regression 132–3
- stepwise selection 70–1
- structural equation modeling (SEM) 5, 441, 458–60; concepts 442; coverage 443; dynamic factor models 459; framework 446; latent class mixture models 433; latent growth models (LGMs) 460; mixed variable models 460; multilevel models 460; multiple group models 458; multiple indicator-multiple cause (MIMIC) model 459; origins 443; parameters 465–6; power, determination 494; second order CFA 458; unobserved variables, latent constructs 448
- structure coefficients (discriminant analysis) 282
- structure matrix: discriminant analysis 300, 308; examples 617; exploratory factor analysis (EFA) 399
- studentized residual 38, 78
- subject-to-parameter ratio 460–1
- subject-to-variable ratio (STV) 377–8, 460–1

- sums of squares, counterpart terms 171
SUTVA *See* stable unit treatment value assumption
 symmetric matrix 603
- territorial maps (discriminant analysis) 296
 territorial plot (discriminant analysis) 283
 test of close fit (CFIt) *See* close fit
 three-group discriminant analysis, SPSS results 305–14
 thresholds, types 341
 traditional R^2 pseudo-variance 128–9
 t test: correlated samples t test 21; dependent samples t test 21; independent samples t test 21; paired samples t test 21
 Tucker-Lewis index 453
 two-factor ANOVA 25–7
 two-factor MANOVA, SPSS results 198–204
 two-factor measurement model 447; generation, Mplus (usage) 493
 two-factor model 498; LISREL output, selection 619–24; test 473
 two-group discriminant analysis, SPSS results 298–302
 two-step clustering 341; procedure 345
 two-way MANOVA, partitioning variation 188
 Type I decision error 11
 Type I error, probability (increase) 506–7
 Type II decision error 11; beta (β) 12–13
- uncentered MLM options 513
 uncorrelated predictors, multilevel linear modeling (MLM) 528
 univariate equality of error variances, tests 201
 univariate normality 43–4; presence 44; screening 45–7
 univariate normality, counterpart terms 171
 univariate normality evidence (MANOVA) 230; interpretation 230–4; statistics 230, 233
 univariate skewness, multivariate extension 429
 univariate statistics: concepts 10–14, 33–4; foundational univariate statistics 14–21
 unrotated factor solution, selection 387
 unstandardized predicted value 38, 78
 unstandardized regression coefficients 58
 unstandardized regression model 61–3
 unstandardized residual 45, 46, 47; boxplot 100; histogram 45, 99; Q-Q plot 45, 80, 100
 unweighted least squares (ULS) 369, 450
- variable(s): categorical variable 74; consistency, determination 370; dependent variable 53; equation usage 145; exclusion 145; exogenous variable 442; inclusion 145, 148; independent variables, file split 38, 39; latent variable 442; manifest variables 444; measurement scale 365–6; moderator variable 73; observation 446; reliability 370–1; variables per factor, number 373
- variance 41–3; components 516; components, final estimation 542, 545, 547; covariance 27–8; homogeneity of variance 41; homogeneity of variance, multilevel linear modeling (MLM) 527; population variance 20; proportion reduction 523–4; pseudo-variance 127–9; sample variance 21; total variance 396
- variance-covariance homogeneity, screening 43
- variance-covariance matrices, homogeneity 42; discriminant analysis 291, 324; repeated measures MANOVA 249–51
- variance-covariance matrix, elements (number) 449
- variance inflation factor (VIF) 151; confirmatory factor analysis (CFA) 465; equation (logistic regression) 135; review (multiple linear regression) 102–3; values 151, 163
- variance inflation factor (VIF), computation 52–3; multiple linear regression 82
- variance ratio criterion (VRC) 343
- vector 602
- VRC *See* variance ratio criterion
- Ward's method 340, 613, 615
- weighted (generalized) least squares (WLS) 369, 450
- whisker plot 558
- Wilk's lambda: discriminant analysis 281–2, 300, 307; representation 200, 219
- within-between: interaction, repeated measures MANOVA 255–9; a priori power, repeated measures MANOVA 260
- within-cluster variation 340–1
- within factors 208
- within-groups 506; multilevel linear modeling (MLM), effect size 523–4; variance, proportion reduction 523–4



Taylor & Francis eBooks

Helping you to choose the right eBooks for your Library

Add Routledge titles to your library's digital collection today. Taylor and Francis ebooks contains over 50,000 titles in the Humanities, Social Sciences, Behavioural Sciences, Built Environment and Law.

Choose from a range of subject packages or create your own!

Benefits for you

- » Free MARC records
- » COUNTER-compliant usage statistics
- » Flexible purchase and pricing options
- » All titles DRM-free.

REQUEST YOUR
FREE
INSTITUTIONAL
TRIAL TODAY

Free Trials Available

We offer free trials to qualifying academic, corporate and government customers.

Benefits for your user

- » Off-site, anytime access via Athens or referring URL
- » Print or copy pages or chapters
- » Full content search
- » Bookmark, highlight and annotate text
- » Access to thousands of pages of quality research at the click of a button.

eCollections – Choose from over 30 subject eCollections, including:

Archaeology	Language Learning
Architecture	Law
Asian Studies	Literature
Business & Management	Media & Communication
Classical Studies	Middle East Studies
Construction	Music
Creative & Media Arts	Philosophy
Criminology & Criminal Justice	Planning
Economics	Politics
Education	Psychology & Mental Health
Energy	Religion
Engineering	Security
English Language & Linguistics	Social Work
Environment & Sustainability	Sociology
Geography	Sport
Health Studies	Theatre & Performance
History	Tourism, Hospitality & Events

For more information, pricing enquiries or to order a free trial, please contact your local sales team:
www.tandfebooks.com/page/sales



Routledge
Taylor & Francis Group

The home of
Routledge books

www.tandfebooks.com