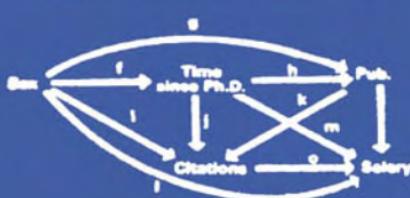


Applied Multiple Regression/ Correlation Analysis for the Behavioral Sciences

Third Edition



Jacob Cohen
Patricia Cohen
Stephen G. West
Leona S. Aiken

Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences

Third Edition

This page intentionally left blank

Applied Multiple Regression/Correlation Analysis for the Behavioral Sciences

Third Edition

Jacob Cohen

(deceased)

New York University

Patricia Cohen

New York State Psychiatric Institute

and

Columbia University College of Physicians and Surgeons

Stephen G. West

Arizona State University

Leona S. Aiken

Arizona State University



LAWRENCE ERLBAUM ASSOCIATES, PUBLISHERS
2003 Mahwah, New Jersey

London

Senior Editor: Debra Rieger

Editorial Assistant: Jason Planer

Cover Design: Kathryn Houghtaling Lacey

Textbook Production Manager: Paul Smolenski

Full-Service & Composition: UG/GGS Information Services, Inc.

Text and Cover Printer: Hamilton Printing Company

This book was typeset in 10/12 pt. Times, Italic, Bold, and Bold Italic.

The heads were typeset in Univers Italic, Medium, and Bold.

Copyright © 2003 by Lawrence Erlbaum Associates, Inc.

All rights reserved. No part of this book may be reproduced in any form, by photostat, microfilm, retrieval system, or any other means, without prior written permission of the publisher.

Lawrence Erlbaum Associates, Inc., Publishers
10 Industrial Avenue
Mahwah, New Jersey 07430

Library of Congress Cataloging-in-Publication Data

Applied multiple regression/correlation analysis for the behavioral sciences.—3rd ed./
Jacob Cohen ... [et al]

p. cm.

Rev. ed. of: Applied multiple regression/correlation analysis for the behavioral sciences
/ Jacob Cohen, Patricia Cohen. 2nd Ed. 1983.

Includes bibliographical references (p.) and indexes

ISBN 0-8058-2223-2 (hard cover: alk. paper)

1. Regression analysis. 2. Correlation (Statistics) 3. Social sciences—Statistical methods. I. Cohen, Jacob, 1923— Applied multiple regression/correlation analysis for the behavioral sciences.

HA31.3 .A67 2003

519.5'36—dc21

2002072068

Books published by Lawrence Erlbaum Associates are printed on acid-free paper, and their bindings are chosen for strength and durability.

Printed in the United States of America

10 9 8 7 6 5

Disclaimer: This eBook does not include ancillary media that was packaged with the printed version of the book.

**To Gideon, whose revisions continue
under the auspices of his own children
J.C. and P.C.**

**To Jack
S.G.W. and L.S.A.**

This page intentionally left blank

Contents

Preface

xxv

Chapter 1: Introduction

1

1.1	Multiple Regression/Correlation as a General Data-Analytic System	1
1.1.1	Overview	1
1.1.2	Testing Hypotheses Using Multiple Regression/Correlation: Some Examples	2
1.1.3	Multiple Regression/Correlation in Prediction Models	3
1.2	A Comparison of Multiple Regression/Correlation and Analysis of Variance Approaches	4
1.2.1	Historical Background	4
1.2.2	Hypothesis Testing and Effect Sizes	5
1.3	Multiple Regression/Correlation and the Complexity of Behavioral Science	6
1.3.1	Multiplicity of Influences	6
1.3.2	Correlation Among Research Factors and Partialing	6
1.3.3	Form of Information	7
1.3.4	Shape of Relationship	8
1.3.5	General and Conditional Relationships	9
1.4	Orientation of the Book	10
1.4.1	Nonmathematical	11
1.4.2	Applied	11
1.4.3	Data-Analytic	12
1.4.4	Inference Orientation and Specification Error	13
1.5	Computation, the Computer, and Numerical Results	14
1.5.1	Computation	14

viii CONTENTS

1.5.2	Numerical Results: Reporting and Rounding	14
1.5.3	Significance Tests, Confidence Intervals, and Appendix Tables	15
1.6	The Spectrum of Behavioral Science	16
1.7	Plan for the Book	16
1.7.1	Content	16
1.7.2	Structure: Numbering of Sections, Tables, and Equations	17
1.8	Summary	18
Chapter 2: Bivariate Correlation and Regression		19
2.1	Tabular and Graphic Representations of Relationships	19
2.2	The Index of Linear Correlation Between Two Variables: The Pearson Product Moment Correlation Coefficient	23
2.2.1	Standard Scores: Making Units Comparable	23
2.2.2	The Product Moment Correlation as a Function of Differences Between z Scores	26
2.3	Alternative Formulas for the Product Moment Correlation Coefficient	28
2.3.1	r as the Average Product of z Scores	28
2.3.2	Raw Score Formulas for r	29
2.3.3	Point Biserial r	29
2.3.4	Phi (ϕ) Coefficient	30
2.3.5	Rank Correlation	31
2.4	Regression Coefficients: Estimating Y From X	32
2.5	Regression Toward the Mean	36
2.6	The Standard Error of Estimate and Measures of the Strength of Association	37
2.7	Summary of Definitions and Interpretations	41
2.8	Statistical Inference With Regression and Correlation Coefficients	41
2.8.1	Assumptions Underlying Statistical Inference With B_{YX} , B_0 , \hat{Y}_i , and r_{XY}	41
2.8.2	Estimation With Confidence Intervals	42
2.8.3	Null Hypothesis Significance Tests (NHSTs)	47
2.8.4	Confidence Limits and Null Hypothesis Significance Testing	50
2.9	Precision and Power	50
2.9.1	Precision of Estimation	50
2.9.2	Power of Null Hypothesis Significance Tests	51
2.10	Factors Affecting the Size of r	53
2.10.1	The Distributions of X and Y	53
2.10.2	The Reliability of the Variables	55
2.10.3	Restriction of Range	57
2.10.4	Part-Whole Correlations	59
2.10.5	Ratio or Index Variables	60
2.10.6	Curvilinear Relationships	62
2.11	Summary	62

Chapter 3: Multiple Regression/Correlation With Two or More Independent Variables	64
3.1 Introduction: Regression and Causal Models	64
3.1.1 What Is a Cause?	64
3.1.2 Diagrammatic Representation of Causal Models	65
3.2 Regression With Two Independent Variables	66
3.3 Measures of Association With Two Independent Variables	69
3.3.1 Multiple R and R^2	69
3.3.2 Semipartial Correlation Coefficients and Increments to R^2	72
3.3.3 Partial Correlation Coefficients	74
3.4 Patterns of Association Between Y and Two Independent Variables	75
3.4.1 Direct and Indirect Effects	75
3.4.2 Partial Redundancy	76
3.4.3 Suppression in Regression Models	77
3.4.4 Spurious Effects and Entirely Indirect Effects	78
3.5 Multiple Regression/Correlation With k Independent Variables	79
3.5.1 Introduction: Components of the Prediction Equation	79
3.5.2 Partial Regression Coefficients	80
3.5.3 R , R^2 , and Shrunken R^2	82
3.5.4 sr and sr^2	84
3.5.5 pr and pr^2	85
3.5.6 Example of Interpretation of Partial Coefficients	85
3.6 Statistical Inference With k Independent Variables	86
3.6.1 Standard Errors and Confidence Intervals for B and β	86
3.6.2 Confidence Intervals for R^2	88
3.6.3 Confidence Intervals for Differences Between Independent R^2 's	88
3.6.4 Statistical Tests on Multiple and Partial Coefficients	88
3.7 Statistical Precision and Power Analysis	90
3.7.1 Introduction: Research Goals and the Null Hypothesis	90
3.7.2 The Precision and Power of R^2	91
3.7.3 Precision and Power Analysis for Partial Coefficients	93
3.8 Using Multiple Regression Equations in Prediction	95
3.8.1 Prediction of Y for a New Observation	95
3.8.2 Correlation of Individual Variables With Predicted Values	96
3.8.3 Cross-Validation and Unit Weighting	97
3.8.4 Multicollinearity	98
3.9 Summary	99

Chapter 4: Data Visualization, Exploration, and Assumption Checking: Diagnosing and Solving Regression Problems I

101

4.1	Introduction	101
4.2	Some Useful Graphical Displays of the Original Data	102
4.2.1	Univariate Displays	103
4.2.2	Bivariate Displays	110
4.2.3	Correlation and Scatterplot Matrices	115
4.3	Assumptions and Ordinary Least Squares Regression	117
4.3.1	Assumptions Underlying Multiple Linear Regression	117
4.3.2	Ordinary Least Squares Estimation	124
4.4	Detecting Violations of Assumptions	125
4.4.1	Form of the Relationship	125
4.4.2	Omitted Independent Variables	127
4.4.3	Measurement Error	129
4.4.4	Homoscedasticity of Residuals	130
4.4.5	Nonindependence of Residuals	134
4.4.6	Normality of Residuals	137
4.5	Remedies: Alternative Approaches When Problems Are Detected	141
4.5.1	Form of the Relationship	141
4.5.2	Inclusion of All Relevant Independent Variables	143
4.5.3	Measurement Error in the Independent Variables	144
4.5.4	Nonconstant Variance	145
4.5.5	Nonindependence of Residuals	147
4.6	Summary	150

Chapter 5: Data-Analytic Strategies Using Multiple Regression/Correlation

151

5.1	Research Questions Answered by Correlations and Their Squares	151
5.1.1	Net Contribution to Prediction	152
5.1.2	Indices of Differential Validity	152
5.1.3	Comparisons of Predictive Utility	152
5.1.4	Attribution of a Fraction of the XY Relationship to a Third Variable	153
5.1.5	Which of Two Variables Accounts for More of the XY Relationship?	153
5.1.6	Are the Various Squared Correlations in One Population Different From Those in Another Given the Same Variables?	154
5.2	Research Questions Answered by B Or β	154
5.2.1	Regression Coefficients as Reflections of Causal Effects	154
5.2.2	Alternative Approaches to Making B_{YX} Substantively Meaningful	154
5.2.3	Are the Effects of a Set of Independent Variables on Two Different Outcomes in a Sample Different?	157

5.2.4	What Are the Reciprocal Effects of Two Variables on One Another?	157
5.3	Hierarchical Analysis Variables in Multiple Regression/ Correlation	158
5.3.1	Causal Priority and the Removal of Confounding Variables	158
5.3.2	Research Relevance	160
5.3.3	Examination of Alternative Hierarchical Sequences of Independent Variable Sets	160
5.3.4	Stepwise Regression	161
5.4	The Analysis of Sets of Independent Variables	162
5.4.1	Types of Sets	162
5.4.2	The Simultaneous and Hierarchical Analyses of Sets	164
5.4.3	Variance Proportions for Sets and the Ballantine Again	166
5.4.4	B and β Coefficients for Variables Within Sets	169
5.5	Significance Testing for Sets	171
5.5.1	Application in Hierarchical Analysis	172
5.5.2	Application in Simultaneous Analysis	173
5.5.3	Using Computer Output to Determine Statistical Significance	174
5.5.4	An Alternative F Test: Using Model 2 Error Estimate From the Final Model	174
5.6	Power Analysis for Sets	176
5.6.1	Determining n^* for the F Test of sR_B^2 with Model 1 or Model 2 Error	177
5.6.2	Estimating the Population sR^2 Values	179
5.6.3	Setting Power for n^*	180
5.6.4	Reconciling Different n^* 's	180
5.6.5	Power as a Function of n	181
5.6.6	Tactics of Power Analysis	182
5.7	Statistical Inference Strategy in Multiple Regression/ Correlation	182
5.7.1	Controlling and Balancing Type I and Type II Errors in Inference	182
5.7.2	Less Is More	185
5.7.3	Least Is Last	186
5.7.4	Adaptation of Fisher's Protected t Test	187
5.7.5	Statistical Inference and the Stage of Scientific Investigations	190
5.8	Summary	190

Chapter 6: Quantitative Scales, Curvilinear Relationships, and Transformations

193

6.1	Introduction	193
6.1.1	What Do We Mean by Linear Regression?	193

xii CONTENTS

6.1.2	Linearity in the Variables and Linear Multiple Regression	194
6.1.3	Four Approaches to Examining Nonlinear Relationships in Multiple Regression	195
6.2	Power Polynomials	196
6.2.1	Method	196
6.2.2	An Example: Quadratic Fit	198
6.2.3	Centering Predictors in Polynomial Equations	201
6.2.4	Relationship of Test of Significance of Highest Order Coefficient and Gain in Prediction	204
6.2.5	Interpreting Polynomial Regression Results	205
6.2.6	Another Example: A Cubic Fit	207
6.2.7	Strategy and Limitations	209
6.2.8	More Complex Equations	213
6.3	Orthogonal Polynomials	214
6.3.1	The Cubic Example Revisited	216
6.3.2	Unequal n and Unequal Intervals	219
6.3.3	Applications and Discussion	220
6.4	Nonlinear Transformations	221
6.4.1	Purposes of Transformation and the Nature of Transformations	221
6.4.2	The Conceptual Basis of Transformations and Model Checking Before and After Transformation—Is It Always Ideal to Transform?	223
6.4.3	Logarithms and Exponents; Additive and Proportional Relationships	223
6.4.4	Linearizing Relationships	225
6.4.5	Linearizing Relationships Based on Strong Theoretical Models	227
6.4.6	Linearizing Relationships Based on Weak Theoretical Models	232
6.4.7	Empirically Driven Transformations in the Absence of Strong or Weak Models	233
6.4.8	Empirically Driven Transformation for Linearization: The Ladder of Re-expression and the Bulging Rule	233
6.4.9	Empirically Driven Transformation for Linearization in the Absence of Models: Box-Cox Family of Power Transformations on Y	236
6.4.10	Empirically Driven Transformation for Linearization in the Absence of Models: Box-Tidwell Family of Power Transformations on X	239
6.4.11	Linearization of Relationships With Correlations: Fisher z' Transform of r	240
6.4.12	Transformations That Linearize Relationships for Counts and Proportions	240
6.4.13	Variance Stabilizing Transformations and Alternatives for Treatment of Heteroscedasticity	244
6.4.14	Transformations to Normalize Variables	246
6.4.15	Diagnostics Following Transformation	247

6.4.16	Measuring and Comparing Model Fit	248
6.4.17	Second-Order Polynomial Numerical Example Revisited	248
6.4.18	When to Transform and the Choice of Transformation	249
6.5	Nonlinear Regression	251
6.6	Nonparametric Regression	252
6.7	Summary	253

Chapter 7: Interactions Among Continuous Variables 255

7.1	Introduction	255
7.1.1	Interactions Versus Additive Effects	256
7.1.2	Conditional First-Order Effects in Equations Containing Interactions	259
7.2	Centering Predictors and the Interpretation of Regression Coefficients in Equations Containing Interactions	261
7.2.1	Regression with Centered Predictors	261
7.2.2	Relationship Between Regression Coefficients in the Uncentered and Centered Equations	262
7.2.3	Centered Equations With No Interaction	262
7.2.4	Essential Versus Nonessential Multicollinearity	264
7.2.5	Centered Equations With Interactions	264
7.2.6	The Highest Order Interaction in the Centered Versus Uncentered Equation	266
7.2.7	Do Not Center Y	266
7.2.8	A Recommendation for Centering	266
7.3	Simple Regression Equations and Simple Slopes	267
7.3.1	Plotting Interactions	269
7.3.2	Moderator Variables	269
7.3.3	Simple Regression Equations	269
7.3.4	Overall Regression Coefficient and Simple Slope at the Mean	270
7.3.5	Simple Slopes From Uncentered Versus Centered Equations Are Identical	271
7.3.6	Linear by Linear Interactions	271
7.3.7	Interpreting Interactions in Multiple Regression and Analysis of Variance	272
7.4	Post Hoc Probing of Interactions	272
7.4.1	Standard Error of Simple Slopes	272
7.4.2	Equation Dependence of Simple Slopes and Their Standard Errors	273
7.4.3	Tests of Significance of Simple Slopes	273
7.4.4	Confidence Intervals Around Simple Slopes	274
7.4.5	A Numerical Example	275
7.4.6	The Uncentered Regression Equation Revisited	281
7.4.7	First-Order Coefficients in Equations Without and With Interactions	281
7.4.8	Interpretation and the Range of Data	282

xiv CONTENTS

7.5	Standardized Estimates for Equations Containing Interactions	282
7.6	Interactions as Partial Effects: Building Regression Equations With Interactions	284
7.7	Patterns of First-Order and Interactive Effects	285
7.7.1	Three Theoretically Meaningful Patterns of First-Order and Interaction Effects	285
7.7.2	Ordinal Versus Disordinal Interactions	286
7.8	Three-Predictor Interactions in Multiple Regression	290
7.9	Curvilinear by Linear Interactions	292
7.10	Interactions Among Sets of Variables	295
7.11	Issues in the Detection of Interactions: Reliability, Predictor Distributions, Model Specification	297
7.11.1	Variable Reliability and Power to Detect Interactions	297
7.11.2	Sampling Designs to Enhance Power to Detect Interactions—Optimal Design	298
7.11.3	Difficulty in Distinguishing Interactions Versus Curvilinear Effects	299
7.12	Summary	300

Chapter 8: Categorical or Nominal Independent Variables **302**

8.1	Introduction	302
8.1.1	Categories as a Set of Independent Variables	302
8.1.2	The Representation of Categories or Nominal Scales	302
8.2	Dummy-Variable Coding	303
8.2.1	Coding the Groups	303
8.2.2	Pearson Correlations of Dummy Variables With Y	308
8.2.3	Correlations Among Dummy-Coded Variables	311
8.2.4	Multiple Correlation of the Dummy-Variable Set With Y	311
8.2.5	Regression Coefficients for Dummy Variables	312
8.2.6	Partial and Semipartial Correlations for Dummy Variables	316
8.2.7	Dummy-Variable Multiple Regression/Correlation and One-Way Analysis of Variance	317
8.2.8	A Cautionary Note: Dummy-Variable-Like Coding Systems	319
8.2.9	Dummy-Variable Coding When Groups Are Not Mutually Exclusive	320
8.3	Unweighted Effects Coding	320
8.3.1	Introduction: Unweighted and Weighted Effects Coding	320
8.3.2	Constructing Unweighted Effects Codes	321
8.3.3	The R^2 and the r_{Yj}^2 s for Unweighted Effects Codes	324
8.3.4	Regression Coefficients and Other Partial Effects in Unweighted Code Sets	325

8.4	Weighted Effects Coding	328
8.4.1	Selection Considerations for Weighted Effects Coding	328
8.4.2	Constructing Weighted Effects	328
8.4.3	The R^2 and \bar{R}^2 for Weighted Effects Codes	330
8.4.4	Interpretation and Testing of B With Unweighted Codes	331
8.5	Contrast Coding	332
8.5.1	Considerations in the Selection of a Contrast Coding Scheme	332
8.5.2	Constructing Contrast Codes	333
8.5.3	The R^2 and \bar{R}^2	337
8.5.4	Partial Regression Coefficients	337
8.5.5	Statistical Power and the Choice of Contrast Codes	340
8.6	Nonsense Coding	341
8.7	Coding Schemes in the Context of Other Independent Variables	342
8.7.1	Combining Nominal and Continuous Independent Variables	342
8.7.2	Calculating Adjusted Means for Nominal Independent Variables	343
8.7.3	Adjusted Means for Combinations of Nominal and Quantitative Independent Variables	344
8.7.4	Adjusted Means for More Than Two Groups and Alternative Coding Methods	348
8.7.5	Multiple Regression/Correlation With Nominal Independent Variables and the Analysis of Covariance	350
8.8	Summary	351

Chapter 9: Interactions With Categorical Variables 354

9.1	Nominal Scale by Nominal Scale Interactions	354
9.1.1	The 2 by 2 Design	354
9.1.2	Regression Analyses of Multiple Sets of Nominal Variables With More Than Two Categories	361
9.2	Interactions Involving More Than Two Nominal Scales	366
9.2.1	An Example of Three Nominal Scales Coded by Alternative Methods	367
9.2.2	Interactions Among Nominal Scales in Which Not All Combinations Are Considered	372
9.2.3	What If the Categories for One or More Nominal “Scales” Are Not Mutually Exclusive?	373
9.2.4	Consideration of pr , β , and Variance Proportions for Nominal Scale Interaction Variables	374
9.2.5	Summary of Issues and Recommendations for Interactions Among Nominal Scales	374
9.3	Nominal Scale by Continuous Variable Interactions	375
9.3.1	A Reminder on Centering	375

xvi CONTENTS

9.3.2	Interactions of a Continuous Variable With Dummy-Variable Coded Groups	375
9.3.3	Interactions Using Weighted or Unweighted Effects Codes	378
9.3.4	Interactions With a Contrast-Coded Nominal Scale	379
9.3.5	Interactions Coded to Estimate Simple Slopes of Groups	380
9.3.6	Categorical Variable Interactions With Nonlinear Effects of Scaled Independent Variables	383
9.3.7	Interactions of a Scale With Two or More Categorical Variables	386
9.4	Summary	388

Chapter 10: Outliers and Multicollinearity: Diagnosing and Solving Regression Problems II

390

10.1	Introduction	390
10.2	Outliers: Introduction and Illustration	391
10.3	Detecting Outliers: Regression Diagnostics	394
10.3.1	Extremity on the Independent Variables: Leverage	394
10.3.2	Extremity on Y : Discrepancy	398
10.3.3	Influence on the Regression Estimates	402
10.3.4	Location of Outlying Points and Diagnostic Statistics	406
10.3.5	Summary and Suggestions	409
10.4	Sources of Outliers and Possible Remedial Actions	411
10.4.1	Sources of Outliers	411
10.4.2	Remedial Actions	415
10.5	Multicollinearity	419
10.5.1	Exact Collinearity	419
10.5.2	Multicollinearity: A Numerical Illustration	420
10.5.3	Measures of the Degree of Multicollinearity	422
10.6	Remedies for Multicollinearity	425
10.6.1	Model Respecification	426
10.6.2	Collection of Additional Data	427
10.6.3	Ridge Regression	427
10.6.4	Principal Components Regression	428
10.6.5	Summary of Multicollinearity Considerations	429
10.7	Summary	430

Chapter 11: Missing Data

431

11.1	Basic Issues in Handling Missing Data	431
11.1.1	Minimize Missing Data	431
11.1.2	Types of Missing Data	432
11.1.3	Traditional Approaches to Missing Data	433

11.2	Missing Data in Nominal Scales	435
11.2.1	Coding Nominal Scale X for Missing Data	435
11.2.2	Missing Data on Two Dichotomies	439
11.2.3	Estimation Using the EM Algorithm	440
11.3	Missing Data in Quantitative Scales	442
11.3.1	Available Alternatives	442
11.3.2	Imputation of Values for Missing Cases	444
11.3.3	Modeling Solutions to Missing Data in Scaled Variables	447
11.3.4	An Illustrative Comparison of Alternative Methods	447
11.3.5	Rules of Thumb	450
11.4	Summary	450

Chapter 12: Multiple Regression/Correlation and Causal Models

452

12.1	Introduction	452
12.1.1	Limits on the Current Discussion and the Relationship Between Causal Analysis and Analysis of Covariance	452
12.1.2	Theories and Multiple Regression/Correlation Models That Estimate and Test Them	454
12.1.3	Kinds of Variables in Causal Models	457
12.1.4	Regression Models as Causal Models	459
12.2	Models Without Reciprocal Causation	460
12.2.1	Direct and Indirect Effects	460
12.2.2	Path Analysis and Path Coefficients	464
12.2.3	Hierarchical Analysis and Reduced Form Equations	465
12.2.4	Partial Causal Models and the Hierarchical Analysis of Sets	466
12.2.5	Testing Model Elements	467
12.3	Models With Reciprocal Causation	467
12.4	Identification and Overidentification	468
12.4.1	Just Identified Models	468
12.4.2	Overidentification	468
12.4.3	Underidentification	469
12.5	Latent Variable Models	469
12.5.1	An Example of a Latent Variable Model	469
12.5.2	How Latent Variables Are Estimated	471
12.5.3	Fixed and Free Estimates in Latent Variable Models	472
12.5.4	Goodness-of-Fit Tests of Latent Variable Models	472
12.5.5	Latent Variable Models and the Correction for Attenuation	473
12.5.6	Characteristics of Data Sets That Make Latent Variable Analysis the Method of Choice	474
12.6	A Review of Causal Model and Statistical Assumptions	475

12.6.1	Specification Error	475
12.6.2	Identification Error	475
12.7	Comparisons of Causal Models	476
12.7.1	Nested Models	476
12.7.2	Longitudinal Data in Causal Models	476
12.8	Summary	477

Chapter 13: Alternative Regression Models: Logistic, Poisson Regression, and the Generalized Linear Model

479

13.1	Ordinary Least Squares Regression Revisited	479
13.1.1	Three Characteristics of Ordinary Least Squares Regression	480
13.1.2	The Generalized Linear Model	480
13.1.3	Relationship of Dichotomous and Count Dependent Variables Y to a Predictor	481
13.2	Dichotomous Outcomes and Logistic Regression	482
13.2.1	Extending Linear Regression: The Linear Probability Model and Discriminant Analysis	483
13.2.2	The Nonlinear Transformation From Predictor to Predicted Scores: Probit and Logistic Transformation	485
13.2.3	The Logistic Regression Equation	486
13.2.4	Numerical Example: Three Forms of the Logistic Regression Equation	487
13.2.5	Understanding the Coefficients for the Predictor in Logistic Regression	492
13.2.6	Multiple Logistic Regression	493
13.2.7	Numerical Example	494
13.2.8	Confidence Intervals on Regression Coefficients and Odds Ratios	497
13.2.9	Estimation of the Regression Model: Maximum Likelihood	498
13.2.10	Deviances: Indices of Overall Fit of the Logistic Regression Model	499
13.2.11	Multiple R^2 Analogs in Logistic Regression	502
13.2.12	Testing Significance of Overall Model Fit: The Likelihood Ratio Test and the Test of Model Deviance	504
13.2.13	χ^2 Test for the Significance of a Single Predictor in a Multiple Logistic Regression Equation	507
13.2.14	Hierarchical Logistic Regression: Likelihood Ratio χ^2 Test for the Significance of a Set of Predictors Above and Beyond Another Set	508
13.2.15	Akaike's Information Criterion and the Bayesian Information Criterion for Model Comparison	509
13.2.16	Some Treachery in Variable Scaling and Interpretation of the Odds Ratio	509

13.2.17	Regression Diagnostics in Logistic Regression	512
13.2.18	Sparseness of Data	516
13.2.19	Classification of Cases	516
13.3	Extensions of Logistic Regression to Multiple Response Categories: Polytomous Logistic Regression and Ordinal Logistic Regression	519
13.3.1	Polytomous Logistic Regression	519
13.3.2	Nested Dichotomies	520
13.3.3	Ordinal Logistic Regression	522
13.4	Models for Count Data: Poisson Regression and Alternatives	525
13.4.1	Linear Regression Applied to Count Data	525
13.4.2	Poisson Probability Distribution	526
13.4.3	Poisson Regression Analysis	528
13.4.4	Overdispersion and Alternative Models	530
13.4.5	Independence of Observations	532
13.4.6	Sources on Poisson Regression	532
13.5	Full Circle: Parallels Between Logistic and Poisson Regression, and the Generalized Linear Model	532
13.5.1	Parallels Between Poisson and Logistic Regression	532
13.5.2	The Generalized Linear Model Revisited	534
13.6	Summary	535

Chapter 14: Random Coefficient Regression and Multilevel Models

536

14.1	Clustering Within Data Sets	536
14.1.1	Clustering, Alpha Inflation, and the Intraclass Correlation	537
14.1.2	Estimating the Intraclass Correlation	538
14.2	Analysis of Clustered Data With Ordinary Least Squares Approaches	539
14.2.1	Numerical Example, Analysis of Clustered Data With Ordinary Least Squares Regression	541
14.3	The Random Coefficient Regression Model	543
14.4	Random Coefficient Regression Model and Multilevel Data Structure	544
14.4.1	Ordinary Least Squares (Fixed Effects) Regression Revisited	544
14.4.2	Fixed and Random Variables	544
14.4.3	Clustering and Hierarchically Structured Data	545
14.4.4	Structure of the Random Coefficient Regression Model	545
14.4.5	Level 1 Equations	546
14.4.6	Level 2 Equations	547
14.4.7	Mixed Model Equation for Random Coefficient Regression	548
14.4.8	Variance Components—New Parameters in the Multilevel Model	548
14.4.9	Variance Components and Random Coefficient Versus Ordinary Least Squares (Fixed Effects) Regression	549

XX CONTENTS

14.4.10	Parameters of the Random Coefficient Regression Model: Fixed and Random Effects	550
14.5	Numerical Example: Analysis of Clustered Data With Random Coefficient Regression	550
14.5.1	Unconditional Cell Means Model and the Intraclass Correlation	551
14.5.2	Testing the Fixed and Random Parts of the Random Coefficient Regression Model	552
14.6	Clustering as a Meaningful Aspect of the Data	553
14.7	Multilevel Modeling With a Predictor at Level 2	553
14.7.1	Level 1 Equations	553
14.7.2	Revised Level 2 Equations	554
14.7.3	Mixed Model Equation With Level 1 Predictor and Level 2 Predictor of Intercept and Slope and the Cross-Level Interaction	554
14.8	An Experimental Design as a Multilevel Data Structure: Combining Experimental Manipulation With Individual Differences	555
14.9	Numerical Example: Multilevel Analysis	556
14.10	Estimation of the Multilevel Model Parameters: Fixed Effects, Variance Components, and Level 1 Equations	560
14.10.1	Fixed Effects and Variance Components	560
14.10.2	An Equation for Each Group: Empirical Bayes Estimates of Level 1 Coefficients	560
14.11	Statistical Tests in Multilevel Models	563
14.11.1	Fixed Effects	563
14.11.2	Variance Components	563
14.12	Some Model Specification Issues	564
14.12.1	The Same Variable at Two Levels	564
14.12.2	Centering in Multilevel Models	564
14.13	Statistical Power of Multilevel Models	565
14.14	Choosing Between the Fixed Effects Model and the Random Coefficient Model	565
14.15	Sources on Multilevel Modeling	566
14.16	Multilevel Models Applied to Repeated Measures Data	566
14.17	Summary	567

Chapter 15: Longitudinal Regression Methods

568

15.1	Introduction	568
15.1.1	Chapter Goals	568
15.1.2	Purposes of Gathering Data on Multiple Occasions	569
15.2	Analyses of Two-Time-Point Data	569
15.2.1	Change or Regressed Change?	570
15.2.2	Alternative Regression Models for Effects Over a Single Unit of Time	571
15.2.3	Three- or Four-Time-Point Data	573
15.3	Repeated Measure Analysis of Variance	573

15.3.1	Multiple Error Terms in Repeated Measure Analysis of Variance	574
15.3.2	Trend Analysis in Analysis of Variance	575
15.3.3	Repeated Measure Analysis of Variance in Which Time Is Not the Issue	576
15.4	Multilevel Regression of Individual Changes Over Time	578
15.4.1	Patterns of Individual Change Over Time	578
15.4.2	Adding Other Fixed Predictors to the Model	582
15.4.3	Individual Differences in Variation Around Individual Slopes	583
15.4.4	Alternative Developmental Models and Error Structures	584
15.4.5	Alternative Link Functions for Predicting Y From Time	586
15.4.6	Unbalanced Data: Variable Timing and Missing Data	587
15.5	Latent Growth Models: Structural Equation Model Representation of Multilevel Data	588
15.5.1	Estimation of Changes in True Scores	589
15.5.2	Representation of Latent Growth Models in Structural Equation Model Diagrams	589
15.5.3	Comparison of Multilevel Regression and Structural Equation Model Analysis of Change	594
15.6	Time Varying Independent Variables	595
15.7	Survival Analysis	596
15.7.1	Regression Analysis of Time Until Outcome and the Problem of Censoring	596
15.7.2	Extension to Time-Varying Independent Variables	599
15.7.3	Extension to Multiple Episode Data	599
15.7.4	Extension to a Categorical Outcome: Event-History Analysis	600
15.8	Time Series Analysis	600
15.8.1	Units of Observation in Time Series Analyses	601
15.8.2	Time Series Analyses Applications	601
15.8.3	Time Effects in Time Series	602
15.8.4	Extension of Time Series Analyses to Multiple Units or Subjects	602
15.9	Dynamic System Analysis	602
15.10	Statistical Inference and Power Analysis in Longitudinal Analyses	604
15.11	Summary	605

Chapter 16: Multiple Dependent Variables: Set Correlation 608

16.1	Introduction to Ordinary Least Squares Treatment of Multiple Dependent Variables	608
16.1.1	Set Correlation Analysis	608

16.1.2	Canonical Analysis	609
16.1.3	Elements of Set Correlation	610
16.2	Measures of Multivariate Association	610
16.2.1	$R_{Y,X}^2$, the Proportion of Generalized Variance	610
16.2.2	$T_{Y,X}^2$ and $P_{Y,X}^2$, Proportions of Additive Variance	611
16.3	Partialing in Set Correlation	613
16.3.1	Frequent Reasons for Partialing Variable Sets From the Basic Sets	613
16.3.2	The Five Types of Association Between Basic Y and X Sets	614
16.4	Tests of Statistical Significance and Statistical Power	615
16.4.1	Testing the Null Hypothesis	615
16.4.2	Estimators of the Population $R_{Y,X}^2$, $T_{Y,X}^2$, and $P_{Y,X}^2$	616
16.4.3	Guarding Against Type I Error Inflation	617
16.5	Statistical Power Analysis in Set Correlation	617
16.6	Comparison of Set Correlation With Multiple Analysis of Variance	619
16.7	New Analytic Possibilities With Set Correlation	620
16.8	Illustrative Examples	621
16.8.1	A Simple Whole Association	621
16.8.2	A Multivariate Analysis of Partial Variance	622
16.8.3	A Hierarchical Analysis of a Quantitative Set and Its Unique Components	623
16.8.4	Bipartial Association Among Three Sets	625
16.9	Summary	627

APPENDICES

Appendix 1:	The Mathematical Basis for Multiple Regression/Correlation and Identification of the Inverse Matrix Elements	631
A1.1	Alternative Matrix Methods	634
A1.2	Determinants	634
Appendix 2:	Determination of the Inverse Matrix and Applications Thereof	636
A2.1	Hand Calculation of the Multiple Regression/Correlation Problem	636
A2.2	Testing the Difference Between Partial β s and B s From the Same Sample	640
A2.3	Testing the Difference Between β s for Different Dependent Variables From a Single Sample	642

Appendix Tables	643
Table A <i>t</i> Values for $\alpha = .01, .05$ (Two Tailed)	643
Table B <i>z'</i> Transformation of <i>r</i>	644
Table C Normal Distribution	645
Table D <i>F</i> Values for $\alpha = .01, .05$	646
Table E <i>L</i> Values for $\alpha = .01, .05$	650
Table F Power of Significance Test of <i>r</i> at $\alpha = .01, .05$ (Two Tailed)	652
Table G <i>n*</i> to Detect <i>r</i> by <i>t</i> Test at $\alpha = .01, .05$ (Two Tailed)	654
References	655
Glossary	671
Statistical Symbols and Abbreviations	683
Author Index	687
Subject Index	691

This page intentionally left blank

Preface

Origins and Background

This book had its origin over 30 years ago, when it became apparent to Jack Cohen that there were relationships between multiple regression and correlation (MRC) on the one hand and the analysis of variance (ANOVA) on the other which were undreamed of (or at least did not appear) in the standard textbooks with which he was familiar. On the contrary, the texts of the era treated MRC and ANOVA as wholly distinct systems of data analysis intended for types of research that differed fundamentally in design, goals, and types of variables. Some research, both statistical and bibliographic, confirmed the relationships noted and revealed yet others. These relationships served to enrich both systems in many ways, but it also became clear that multiple regression/correlation was potentially a very general system for analyzing data in the behavioral sciences, one that could incorporate the analysis of variance and covariance as special cases. An article outlining these possibilities was published in the quantitative methods section of *Psychological Bulletin* (J. Cohen, 1968),¹ and it has gone on to become one of the most cited articles in the *Bulletin's* history (Sternberg, 1992).² The volume and sources of early reprint requests and requests to reprint the article suggested that a responsive chord had been struck among behavioral scientists in diverse areas. It was also obvious that a book-length treatment was needed for adequacy of both systematic coverage and expository detail.

In 1969 Jack and Pat were married and began a happy collaboration, one of whose chief products is this book. (Another has been saluted on the dedication page of each edition.) During the preparation of the first edition of the book, the ideas of the 1968 paper were expanded, further systematized, tried out on data, and hardened in the crucible of our teaching, research, and consulting. We find this system, which has now attained broad usage in the behavioral sciences, to be surprisingly easy to teach and learn. The first edition of this book was published in 1975 and, following further development and revision of the ideas, the second edition was published in 1983.

Despite the continuing popularity of the second edition of this text, by the early 1990s Jack and Pat were very aware of the need to update and extend its coverage of new methods, options,

¹Cohen, J. (1968). Multiple regression as a general data-analytic system. *Psychological Bulletin*, 70, 426–443.

²Sternberg, R. J. (1992). Psychological Bulletin's top 10 "hit parade". *Psychological Bulletin*, 112, 387–388.

and graphics in the regression field. The methods that Jack had done so much to promote were becoming so familiar to scientists in the field that they no longer needed an extensive elaboration of their virtues. New improved methods of regression diagnostics and graphics, multilevel analyses, logistic and other nominal dependent variable methods of regression, and treatment of longitudinal data made updating seem critical. New generations of computer software for multiple regression analysis had been developed, including “point and click” computer programs that now reside on the desktop computer of every researcher and student. These new programs have combined the wonderful virtues of flexibility and ease of use that have made multiple regression analysis even more accessible. However, this increased accessibility has also increased the risk that multiple regression may be used in a mindless manner by those unfamiliar with its basic concepts. This development made the need for continued clear coverage of the basic concepts of multiple regression even more apparent.

Because Jack and Pat were aware of the magnitude of the revisions that should be made, and because they wanted statistical experts who were working in fields central to psychological research, they invited Drs. Leona Aiken and Stephen West, another “multivariate couple,” to collaborate on the revision. Jack and Pat particularly admired Aiken and West’s book *Multiple Regression: Testing and Interpreting Interactions* and its extensive use of graphical presentations. Not surprisingly, when we all started to work together, we found that the revisions that Jack and Pat had originally envisioned were not sufficient to cover all the changes that we collectively thought to be important.

Jack’s death in 1998 has made this revision much more difficult for the remaining three of us and, in some ways, more important. The four of us had planned the changes together, divided the tasks, and were well started, but there was still a lot of work to be done. We wanted to decrease the emphasis on significance tests and increase the recommendations for confidence intervals and effect size measures, which Jack was so active in promoting. Some of his last writing could be incorporated, but in other cases we needed to work these ideas in without his help.

The Audience for the Book

To describe the primary audience for whom this book is intended requires two dimensions. Substantively, this book is addressed to behavioral and social scientists. These terms have no sharply defined reference, but we intend them in the most inclusive sense to include the academic sciences of psychology, sociology, economics, branches of biology, political science, anthropology, and social epidemiology, as well as to those in business, health sciences, communication and other applied fields of behavioral research.

The other dimension of our intended audience, amount of background in statistics and research, covers an equally broad span. We are very pleased that the book serves both as a textbook for students and as a handbook for researchers. One particular feature of this book will be appreciated by many in both groups of readers: Its orientation is nonmathematical, applied, and data-analytic. This orientation is discussed in the introductory chapter, and will not be belabored here. Our experience has been that with few exceptions, both students and substantive researchers in the behavioral and social sciences often approach statistics with considerable wariness, and profit most from a verbal-conceptual exposition, rich in redundancy and concrete examples. This we have sought to supply.

As a textbook, whether used in a course at the graduate or advanced undergraduate level, it is assumed that students have already had a semester’s introductory statistics course. Roughly the first half of the book can be used as the foundation of a solid one semester regression course. We anticipate that two full semesters would be necessary to cover the entire book. As a

manual, the researcher's specific interests may dictate the sequence of reference to the chapters that follow. As much as possible, we attempted to write the book in a manner that minimizes the need to refer to previous chapters.

The Third Edition: Old and New

The text of this revision remains very much in the spirit of the previous editions: that is, in the words of one of us, "multiple regression in the service of the ego" (J. Cohen, 1964).³ We are delighted that the behavioral and social sciences have increasingly moved away from the use of statistical tests as a kind of mathematical blessing on an investigation, and toward a data-analytic point of view that focuses on answering the questions that motivated the investigation. While we have tried to keep the overall conceptual tone of previous editions of the book, we have felt it necessary to make many modifications to reflect the current and developing practices in the field. These have included an increased emphasis on graphical presentations throughout. We have somewhat down-played the comparison of MRC to ANOVA, feeling that the battle between these methods is as close to being resolved as it will ever be, and that MRC is now clearly recognized as the more general method (as well as sometimes being supplanted by other methods that better address failure of its assumptions). In addition, we recognize that although ANOVA may still hold sway in some experimental fields, many of our readers in other behavioral sciences will not be very familiar with these statistical models. We also believe that the behavioral sciences have developed to the point where it is appropriate to begin emphasizing the reporting of regression results in meaningful units rather than relying heavily on correlational statistics.

Because of the widely positive response to our presentation of the basic regression ideas, Chapters 1 through 3 are only moderately modified. Chapter 1 begins with an outline of the general system of data analysis made possible by multiple regression/correlation methods. Beginning students may benefit from rereading Chapter 1 after they gain greater familiarity with the specifics of this system through their reading of later chapters. Chapter 2 begins "from scratch" with bivariate correlation and regression, and reviews elementary statistical concepts and terminology. Chapter 2 is not really intended to be a thorough, basic exposition, but rather its purpose is to refresh the reader's memory, and to affix the basic meanings of regression and correlation so firmly in the reader's mind that later exposition will be readily understood. Chapter 3 extends this basic understanding to models with multiple independent variables. Chapter 4, new to this edition, considers the assumptions of multiple regression and outlines alternative procedures that may be taken when they are not met. We have organized old and new issues of data-analytic strategy into an independent Chapter 5. The chapters on curvilinear relationships and transformations (Chapter 6) and interactions between continuous variables (Chapter 7) have been substantially rewritten to reflect a variety of new developments in these areas. Chapter 8 on nominal (group) variables has been moderately rewritten and Chapter 9 on group by continuous variable interactions has been extensively rewritten to reflect new developments. Chapter 10, new to this edition, considers the potential problems of outliers and multicollinearity and their remedies. Chapter 11 on missing data now incorporates the full armamentarium of new methods that have been developed to cope with this common problem. Chapter 12 on causal analysis is updated. Chapter 13, new to this edition, covers a variety of techniques including logistic regression and Poisson regression that are useful

³Cohen, J. (1964). Lecture given at the New York Academy of Medicine. Published as J. Cohen (1968). Prognostic factors in functional psychosis: A study in multivariate methodology. *Transactions of the New York Academy of Sciences*, 30(6), 833–840.

in addressing data sets in which the dependent variables are binary, ordered categories, or counts. Chapter 14, also new to this edition, provides an introduction to multilevel models for clustered data. Another new Chapter (15) introduces the reader to the many methods of answering research questions using longitudinal data. Finally, Jack's Chapter 16 on set correlation analysis of multiple dependent variables has been revised based on his presentation in the Keren and Lewis (1993) volume.⁴

We have many to thank for help in this revision, including Dorothy Castille who suggested and created the first entry list for the technical glossary, Henian Chen who programmed the multilevel longitudinal examples, Tom Crawford who matched the SEM models of longitudinal data with the multilevel models and assisted on other path models, Larkin McReynold who volunteered with the SPSS programming of many examples. Steven C. Pitts developed many of the key numerical examples in Chapters 2, 3, 7, and 10 as well as graphical displays used in Chapters 6, 7, 10, 13, and 14, and Kathy Gordon did these jobs for other examples in Chapters 2, 3, 5, 9, 11, 12, and 15. Jennifer L. Krull developed the multilevel example in Chapter 14. Oi-Man Kwok and Jonathan Butner wrote the computer syntax for examples in a number of chapters. Technical production help was also provided by Jonathan Butner, Kathy Gordon, Steven C. Pitts, and Justin McKenzie. We are thankful for the substantial guidance on each new or rewritten chapter that we were able to get from experts, many of whom were drawn from our friends and colleagues in the Society for Multivariate Experimental Psychology. Thanks are due to the following colleagues who reviewed some of the chapters for us, including (in alphabetical order) Jeremy Biesanz, Barbara Byrne, Dianne Chambliss, Patrick Curran, Mindy Erchull, John W. Graham, Fumiaki Hamagami, Siek-Toon Khoo, Bruce Levin, William Mason, John J. McArdle, Roger Millsap, John Nesselroade, Jason Newsom, Abigail Panter, Mark Reiser, David Rindskopf, Patrick Shrout, and Aaron Taylor. We give special thanks to William F. Chaplin, Daniel W. King, Lynda A. King, Harry Reis, Steven Reise, and Leland Wilkinson, who provided insightful and informative reviews of the entire volume. Of course, the errors that may remain are entirely our own.

JACOB COHEN
PATRICIA COHEN
STEPHEN G. WEST
LEONA S. AIKEN

⁴Keren, G., and Lewis, C. (Eds.). (1993). *A handbook for data analysis in the behavioral sciences: Statistical issues*. Hillsdale, NJ: Erlbaum.

**Applied Multiple
Regression/Correlation
Analysis for the
Behavioral Sciences**

Third Edition

This page intentionally left blank

1

Introduction

1.1 MULTIPLE REGRESSION/CORRELATION AS A GENERAL DATA-ANALYTIC SYSTEM

1.1.1 Overview

Multiple regression/correlation analysis (MRC) is a highly general and therefore very flexible data analytic system. Basic MRC may be used whenever a quantitative variable, the dependent variable (Y), is to be studied as a function of, or in relationship to, any factors of interest, the independent variables (IVs).¹ The broad sweep of this statement is quite intentional.

1. The form of the relationship is not constrained: it may be simple or complex, for example, straight line or curvilinear, general or conditional, or combinations of these possibilities.
2. The nature of the research factors expressed as independent variables is also not constrained. They may be quantitative or qualitative, main effects or interactions in the analysis of variance (ANOVA) sense, or covariates in the analysis of covariance (ANCOVA) sense. They may be correlated with each other or uncorrelated as in balanced factorial designs in ANOVA commonly found in laboratory experiments. They may be naturally occurring (“organismic” variables) like sex, diagnosis, IQ, extroversion, or years of education, or they may be planned experimental manipulations (treatment conditions). In short, virtually any information whose bearing on the dependent variable is of interest may be expressed as research factors.
3. The nature of the dependent variable is also not constrained. Although MRC was originally developed for scaled dependent variables, extensions of the basic model now permit appropriate analysis of the full range of dependent variables including those that are of the form of categories (e.g., ill vs. not ill) or ordered categories.
4. Like all statistical analyses, the basic MRC model makes assumptions about the nature of the data that are being analyzed and is most confidently conducted with “well-behaved” data that meet the underlying assumptions of the basic model. Statistical and graphical methods now part of many statistical packages make it easy for the researcher to determine whether

¹In this book we typically employ Y to indicate a dependent variable and IV to represent an independent variable to indicate their role in the statistical analysis without any *necessary* implication of the existence or direction of causal relationship between them.

2 1. INTRODUCTION

estimates generated by the basic MRC model are likely to be misleading and to take appropriate actions. Extensions of the basic MRC model include appropriate techniques for handling “badly behaved” or missing data and other data problems encountered by researchers.

The MRC system presented in this book has other properties that make it a powerful analytic tool. It yields measures of the magnitude of the total effect of a factor on the dependent variable as well as of its partial (unique, net) relationship, that is, its relationship over and above that of other research factors. It also comes fully equipped with the necessary apparatus for statistical hypothesis testing, estimation, construction of confidence intervals, and power analysis. Graphical techniques allow clear depictions of the data and of the analytic results. Last, but certainly not least, MRC is a major tool in the methods of causal (path, structural equation) analysis. Thus, MRC is a versatile, all-purpose system of analyzing the data over a wide range of sciences and technologies.

1.1.2 Testing Hypotheses Using Multiple Regression/Correlation: Some Examples

Multiple regression analysis is broadly applicable to hypotheses generated by researchers in the behavioral sciences, health sciences, education, and business. These hypotheses may come from formal theory, previous research, or simply scientific hunches. Consider the following hypotheses chosen from a variety of research areas:

1. In health sciences, Rahe, Mahan, and Arthur (1970) hypothesized that the amount of major life stress experienced by an individual is positively related to the number of days of illness that person will experience during the following 6 months.
2. In sociology, England, Farkas, Kilbourne, and Dou (1988) predicted that the size of the positive relationship between the number of years of job experience and workers' salaries would depend on the percentage of female workers in the occupation. Occupations with a higher percentage of female workers were expected to have smaller increases in workers' salaries than occupations with a smaller percentage of female workers.
3. In educational policy, there is strong interest in comparing the achievement of students who attend public vs. private schools (Coleman, Hoffer, & Kilgore, 1982; Lee & Bryk, 1989). In comparing these two “treatments” it is important to control statistically for a number of background characteristics of the students such as prior academic achievement, IQ, race, and family income.
4. In experimental psychology, Yerkes and Dodson (1908) proposed a classic “law” that performance has an inverted U-shaped relationship to physiological arousal. The point at which maximum performance occurs is determined by the difficulty of the task.
5. In health sciences, Aiken, West, Woodward, and Reno (1994) developed a predictive model of women's compliance versus noncompliance (a binary outcome) with recommendations for screening mammography. They were interested in the ability of a set of health beliefs (perceived severity of breast cancer, perceived susceptibility to breast cancer, perceived benefits of mammography, perceived barriers to mammography) to predict compliance over and above several other sets of variables: demographics, family medical history, medical input, and prior knowledge.

Each of these hypotheses proposes some form of relationship between one or more factors of interest (independent variables) and an outcome (dependent) variable. There are usually other variables whose effects also need to be considered, for reasons we will be discussing in

this text. This book strongly emphasizes the critical role of theory in planning MRC analyses. The researcher's task is to develop a statistical model that will accurately estimate the relationships among the variables. Then the power of MRC analysis can be brought to bear to test the hypotheses and provide estimations of the size of the effects. However, this task cannot be carried out well if the actual data are not evaluated with regard to the assumptions of the statistical model.

1.1.3 Multiple Regression/Correlation in Prediction Models

Other applications of MRC exist as well. MRC can be used in practical prediction problems where the goal is to forecast an outcome based on data that were collected earlier. For example, a college admissions committee might be interested in predicting college GPA based on high school grades, college entrance examination (SAT or ACT) scores, and ratings of students by high school teachers. In the absence of prior research or theory, MRC can be used in a purely exploratory fashion to identify a collection of variables that strongly predict an outcome variable. For example, coding of the court records for a large city could identify a number of characteristics of felony court cases (e.g., crime characteristics, defendant demographics, drug involvement, crime location, nature of legal representation) that might predict the length of sentence. MRC can be used to identify a minimum set of variables that yield the best prediction of the criterion for the data that have been collected (A. J. Miller, 1990). Of course, because this method will inevitably capitalize on chance relationships in the original data set, replication in a new sample will be critical. Although we will address purely predictive applications of MRC in this book, our focus will be on the MRC techniques that are most useful in the testing of scientific hypotheses.

In this chapter, we initially consider several issues that are associated with the application of MRC in the behavioral sciences. Some disciplines within the behavioral sciences (e.g., experimental psychology) have had a misperception that MRC is only suitable for nonexperimental research. We consider how this misperception arose historically, note that MRC yields identical statistical tests to those provided by ANOVA yet additionally provides several useful measures of the size of the effect. We also note some of the persisting differences in data-analytic philosophy that are associated with researchers using MRC rather than ANOVA. We then consider how the MRC model nicely matches the complexity and variety of relationships commonly observed in the behavioral sciences. Several independent variables may be expected to influence the dependent variable, the independent variables themselves may be related, the independent variables may take different forms (e.g., rating scales or categorical judgments), and the form of the relationship between the independent and dependent variables may also be complex. Each of these complexities is nicely addressed by the MRC model. Finally, we consider the meaning of causality in the behavioral sciences and the meanings of control. Included in this section is a discussion of how MRC and related techniques can help rule out at least some explanations of the observed relationships. We encourage readers to consider these issues at the beginning of their study of the MRC approach and then to reconsider them at the end.

We then describe the orientation and contents of the book. It is oriented toward practical data analysis problems and so is generally nonmathematical and applied. We strongly encourage readers to work through the solved problems, to take full advantage of the programs for three major computer packages and data sets included with the book, and, most important, to learn MRC by applying these techniques to their own data. Finally, we provide a brief overview of the content of the book, outlining the central questions that are the focus of each chapter.

1.2 A COMPARISON OF MULTIPLE REGRESSION/CORRELATION AND ANALYSIS OF VARIANCE APPROACHES

MRC, ANOVA, and ANCOVA are each special cases of the *general linear model* in mathematical statistics.² The description of MRC in this book includes extensions of conventional MRC analysis to the point where it is essentially equivalent to the general linear model. It thus follows that any data analyzable by ANOVA/ANCOVA may be analyzed by MRC, whereas the reverse is not the case. For example, research designs that study how a scaled characteristic of participants (e.g., IQ) and an experimental manipulation (e.g., structured vs. unstructured tasks) jointly influence the subjects' responses (e.g., task performance) cannot readily be fit into the ANOVA framework. Even experiments with factorial designs with unequal cell sample sizes present complexities for ANOVA approaches because of the nonindependence of the factors, and standard computer programs now use a regression approach to estimate effects in such cases. The latter chapters of the book will extend the basic MRC model still further to include alternative statistical methods of estimating relationships.

1.2.1 Historical Background

Historically, MRC arose in the biological and behavioral sciences around 1900 in the study of the natural covariation of observed characteristics of samples of subjects, including Galton's studies of the relationship between the heights of fathers and sons and Pearson's and Yule's work on educational issues (Yule, 1911). Somewhat later, ANOVA/ANCOVA grew out of the analysis of agricultural data produced by the controlled variation of treatment conditions in manipulative experiments. It is noteworthy that Fisher's initial statistical work in this area emphasized the multiple regression framework because of its generality (see Tatsuoka, 1993). However, multiple regression was often computationally intractable in the precomputer era: computations that take milliseconds by computer required weeks or even months to do by hand. This led Fisher to develop the computationally simpler, equal (or proportional) sample size ANOVA/ANCOVA model, which is particularly applicable to planned experiments. Thus multiple regression and ANOVA/ANCOVA approaches developed in parallel and, from the perspective of the substantive researchers who used them, largely independently. Indeed, in certain disciplines such as psychology and education, the association of MRC with nonexperimental, observational, and survey research led some scientists to perceive MRC to be less scientifically respectable than ANOVA/ANCOVA, which was associated with experiments.

Close examination suggests that this guilt (or virtue) by association is unwarranted—the result of the confusion of data-analytic method with the logical considerations that govern the inference of causality. Experiments in which different treatments are applied to randomly assigned groups of subjects and there is no loss (attrition) of subjects permit unambiguous inference of causality; the observation of associations among variables in a group of randomly selected subjects does not. Thus, interpretation of a finding of superior early school achievement of children who participate in Head Start programs compared to nonparticipating children depends on the design of the investigation (Shadish, Cook, & Campbell, 2002; West, Biesanz, & Pitts, 2000). For the investigator who randomly assigns children to Head Start versus Control programs, attribution of the effect to program content is straightforward. For the investigator who simply observes whether children whose parents select Head Start programs have higher school achievement than those who do not, causal inference becomes less certain. Many other possible differences (e.g., child IQ; parent education) may exist between

²For the technically minded, our primary focus will be on the "fixed" version of these models, representing the most common usage of the general linear model in the behavioral sciences.

the two groups of children that could potentially account for any findings. But each of the investigative teams may analyze their data using either ANOVA (or equivalently a t test of the mean difference in school achievement) or MRC (a simple one-predictor regression analysis of school achievement as a function of Head Start attendance with its identical t test). The logical status of causal inference is a function of how the data were produced, not how they were analyzed (see further discussion in several chapters, especially in Chapter 12).

1.2.2 Hypothesis Testing and Effect Sizes

Any relationship we observe, whether between independent variables (treatments) and an outcome in an experiment or between independent variables and a “dependent” variable in an observational study, can be characterized in terms of the strength of the relationship or its effect size (ES). We can ask how much of the total variation in the dependent variable is produced by or associated with the independent variables we are studying. One of the most attractive features of MRC is its automatic provision of regression coefficients, proportion of variance, and correlational measures of various kinds, all of which are kinds of ES measures. We venture the assertion that, despite the preoccupation of the behavioral sciences, the health sciences, education, and business with quantitative methods, the level of consciousness in many areas about strength of observed relationships is at a surprisingly low level. This is because concern about the statistical significance of effects has tended to pre-empt attention to their magnitude (Harlow, Mulaik, & Steiger, 1997). Statistical significance only provides information about whether the relationship exists at all, often a question of trivial scientific interest, as has been pointed out in several commentaries (e.g., J. Cohen, 1994; Meehl, 1967). The level of statistical significance reflects the sample size, incidental features of the design, the sampling of cases, and the nature of the measurement of the dependent variable; it provides only a very pale reflection of the effect size. Yet many research reports, at least implicitly, confuse the issues of effect size and level of statistical significance, using the latter as if it meant the former (Gigerenzer, 1993).

Part of the reason for this unfortunate tendency is that traditional ANOVA/ANCOVA yields readily interpretable F and t ratios for significance testing and differences between cell means for interpretation of the direction of the effect, but no standardized index of effect size. When the dependent measure is in commonly understood units, such as yield of cotton per acre in agricultural research or dollars of income in economic research, the difference in means provides an informative measure. In the social sciences mean differences may also be informative, providing that some method of establishing meaningful measurement units has been accomplished. However, such unit establishment is often not the case, a problem discussed further in Section 5.2. In such a case standardized measures of effect size provided by the MRC analysis often permit more straightforward interpretation. Indeed, researchers in the ANOVA/ANCOVA tradition have become aware of standardized measures of effect size because of the rise of meta-analytic approaches that provide quantitative summaries of entire research literatures (e.g., Rosenthal, 1991). Some journal editors have also begun to encourage or even require inclusion of standardized effect size measures in articles published in their journals.

In addition to effect size measures in original (raw) and standardized units, the MRC system routinely provides several measures of the proportion of variance accounted for (the squares of simple, multiple, partial, and semipartial correlation coefficients). These measures of effect size are unit free and are easily understood and communicated. Each of the measures comes with its significance test value for the null hypothesis (F or t) so that no confusion between the two issues of whether and how much need arise.

1.3 MULTIPLE REGRESSION/CORRELATION AND THE COMPLEXITY OF BEHAVIORAL SCIENCE

The greatest virtue of the MRC system is its capacity to represent, with high fidelity, the types and the complexity of relationships that characterize the behavioral sciences. The word *complexity* is itself used here in a complex sense to cover several issues.

1.3.1 Multiplicity of Influences

The behavioral sciences inherited from older branches of empirical inquiry the simple experimental paradigm: Vary a single presumed causal factor (C) and observe its effects on the dependent variable (Y) while holding constant other potential factors. Thus, $Y = f(C)$; that is, to some degree, variation in Y is a function of controlled variation in C . This model has been, and continues to be, an effective tool of inquiry in the physical sciences, engineering, and in some areas of the behavioral sciences. A number of areas within the physical sciences and engineering typically deal with a few distinct causal factors, each measured in a clear-cut way, and each in principle independent of others.

However, as one moves to the broad spectrum of the basic and applied behavioral sciences ranging from physiological psychology to cultural anthropology to evaluation of educational programs, the number of potential causal factors increases, their representation in measures becomes increasingly uncertain, and weak theories abound and compete. Consider the following set of dependent variables from selected areas of the behavioral sciences, health sciences, education, and business: number of presidential vetoes (political science), extent of women's participation in the labor force (sociology), distance from home to work (geography), reaction time (experimental psychology), migration rate (demography), depression (clinical psychology), kinship system (anthropology), new business startups (economics), compliance with medical regime (health sciences), school achievement (education), and personnel turnover (business). A few moment's reflection about the context in which each of these is embedded suggests the multiplicity of both the potential causal factors and the forms of their relationships to the dependent variables. Given several research factors, C , D , E , etc., to be studied, one might use the single-factor paradigm repeatedly in a program of research: $Y = f(C)$, then $Y = f(D)$, then $Y = f(E)$, etc. But MRC permits the far more efficient simultaneous examination of the influences of multiple factors; that is, $Y = f(C, D, E, \text{etc.})$. Moreover, techniques such as structural equation analysis use interlocking regression equations to estimate formal models of causal processes derived from complex substantive theories.

1.3.2 Correlation Among Research Factors and Partialing

A far more important type of complexity than the sheer multiplicity of research factors lies in the effect of relationships among them. The simplest condition is that in which the factors C , D , E , . . . are statistically unrelated (orthogonal) to each other, as is the case in experiments in which the subject's level on each factor is under the experimenter's control and equal (or proportional) numbers of subjects are represented at each combination of factors. The overall importance of each factor in the experiment can be unambiguously determined because its independence of the other factors assures that its effects on Y cannot overlap with the effects of the others. Consider an experiment in which the apparent age (30 vs. 40) and sex (male, female) of a communicator are manipulated and their separate and joint effects on attitude change of male subjects is observed. The orthogonality of the factors is assured by having equal numbers of subjects in each of the four cells defined by the possible combinations of

gender and age of the communicator (30-year-old male, 30-year-old female, 40-year-old male, 40-year-old female). No part of the difference in overall Y means for the two communicator ages can be attributed to their gender, nor can any part of the difference in the overall Y means for the two sexes be attributed to their ages.

Complexity arises when one departs from equal or proportional numbers of subjects in different conditions, because the independent variables are no longer independent. If in an experiment, the majority of the 40-year-olds were male and the majority of the 30-year-olds were female, then any difference between male and female communicators in the overall Y means would be confounded with (correlated with) communicator age. The age and sex effects would no longer be additive. Many issues in the behavioral sciences are simply inaccessible to true experiments and can only be addressed by the systematic observation of phenomena as they occur in their natural context. In nature, factors that influence Y are generally correlated with one another. Thus, if attitudes toward abortion (Y) are studied in a sample of survey respondents as a function of political party (C), religious background (D), and socioeconomic status (E), it is likely that C , D , and E will be correlated with each other. Relationships with Y , taken singly, will not accurately represent their separate influences, because of correlations among the factors (see Section 3.4). This is the familiar phenomenon of redundancy among correlated independent variables with regard to what they explain. The Y relationship with each of the independent variables overlaps to some degree with their relationships with other variables in the statistical model. This, in turn, requires a concept of the unique (“partialed”) relationship of each variable with Y , in the context of the other variables in the model. This picture is often sharply different from that provided by looking at each factor singly. For example, it might be argued that the apparent influence of political party on attitudes toward abortion is entirely attributable to the relationship of party affiliation to religious preference or socioeconomic status. Such a pattern of results suggests that the apparent influence of political party on attitudes when appraised by itself may be “spurious”; that is, within subgroups that are homogeneous with regard to religious background and socioeconomic status, there is no difference on the average between members of one party and members of the other. Detailed attention to the relationships among potentially causal independent variables and how these bear on Y is the hallmark of causal analysis, and may be accomplished by MRC.

MRC’s capability for assessing unique or partial relationships is perhaps its most important feature. Even a small number of research factors define many alternative causal systems. Some of these causal systems will be implausible because of considerations of prior research findings, logic, or research design (e.g., in a longitudinal design variables that occur later in time may be ruled out as potential causes of earlier variables). However, selection among the remaining causal systems is greatly facilitated by the ability, using MRC, of assessing the unique effect of a research factor, statistically controlling for (partialing) the effects of any desired set of other factors. Correlation does not prove causation; however, the absence of correlation implies the absence of the existence of a causal relationship. Thus, the skillful use of MRC can invalidate causal alternatives, assist researchers in choosing between competing theories, and help disentangle multiple influences through its partialing feature.

1.3.3 Form of Information

Variables employed in MRC may represent several different levels of measurement, of which it is often useful to distinguish the following (S. S. Stevens, 1951, 1958):

1. **Ratio scales.** These are equal interval scales with a true zero point, a point at which there is none of whatever the scale is measuring. Only such scales make statements such as “John weighs twice as much as Jim” or “Mary earns two-thirds as much as Jane” sensible.

Some examples of ratio scale measures include inches, pounds, seconds, size of group, dollars, distance from hospital, years in prison, and literacy rate.

2. Interval scales. These scales have equal intervals but are measured from an arbitrary point. For example, the Fahrenheit temperature scale uses the temperature at which a certain concentration of salt water freezes to represent 0. Values on the scale of less than 0 can and do occur. Many psychological and sociological indices are at this level, for example, scores on tests of intelligence, special abilities, achievement, personality, temperament, vocational interest, and social attitude. Such scales may not have a meaningful zero value at all.

3. Ordinal scales. Only the relative positions within a specific collection are signified by the values of ordinal scales. These scales do not have either equal intervals or a true zero point. Examples of ordinal scales include simple rankings of subjects in a sample as well as re-expressions of such rankings into percentiles, deciles, and quartiles.

4. Nominal scales. Nominal (categorical) scales involve simple classification of subjects into categories. The categories of nominal scales represent distinguishable qualities without a natural order or other quantitative properties. Examples include ethnic group, experimental treatment condition, place of birth, religion, marital status, psychiatric diagnosis, type of family structure, political party, public versus private sector, and gender. The set of categories are usually mutually exclusive and exhaustive. Thus, nominal scales are sets of groups that differ on some qualitative attribute.

This classification scheme is not exhaustive of quantitative scales, and others have been proposed. For example, psychological test scores are unlikely to measure with exactly equal intervals and it may be argued that they fall between interval and ordinal scales. Also, some rating scales frequently used in psychological research are not covered by Stevens' conception of levels of measurement. For example, scales like "0 = never, 1 = seldom, 2 = sometimes, 3 = often, and 4 = always" have a defined zero point, but intervals of dubious equality, although for most purposes they are treated as if they are approximately equal.

Basic MRC analysis can potentially consider information at any single level or any mixture of these levels of measurement. Ratio- and interval-level independent variables can be directly included in MRC models. Nominal variables can be expressed as coded variables (e.g., male = 0; female = 1), as will be discussed in Chapters 2, 8, and 9. Ordinal IVs may be treated as if they were interval variables in MRC models, and the results of the analyses may often be satisfactory. However, such an employment of these variables requires special caution, as is discussed further in Chapter 4. On the dependent variable side, Y may be measured at any of the levels of measurement, but the basic MRC model will usually work best if the data are interval or ratio. Some types of dependent variables may lead to violations of basic assumptions of the MRC model. In such cases, generalizations of the basic MRC model (the generalized linear model) can lead to improvements in the accuracy of the results over the basic MRC model (discussed in Chapter 13). This capacity of MRC and its generalizations to use information in almost any form, and to mix forms as necessary, is an important part of its adaptive flexibility.

1.3.4 Shape of Relationship

Consider the relationship $Y = f(C)$, where Y is a measure of poor health such as number of days of illness per year. For some factors the relationship may be well described by a straight line on the usual graph, for example, if C is daily cigarette consumption. Or, adequate description may require that the line be curved; for example, if C is age in years, the very young and the elderly are more often sick than young and middle-aged adults. Or, the shape

may not be definable, as when C is a nominal variable like sex, ethnic background, or religion. When multiple research factors are being studied simultaneously, each may relate to Y (and each other) in any of these ways. Thus, when we write $Y = f(C, D, E, \dots)$, f (as a function of) potentially covers a variety of complex functions that are readily brought under the sway of MRC.

How so? Many readers will know that MRC is often (and properly) referred to as *linear* MRC and may well be under the impression that correlation and regression are restricted to the study of straight-line relationships. This mistaken impression is abetted by the common usage of *linear* to mean “rectilinear” (straight line) and *nonlinear* to mean “curvilinear” (curved line). What is meant by *linear* in the MRC framework is any relationship of the form

$$(1.1.1) \quad Y = a + bU + cV + dW + eX + \dots$$

where the lowercase letters are constants (either positive or negative) and the capital letters are variables. Y is said to be “linear in the variables U , V , etc.” because it may be estimated by taking certain amounts (b , c , etc.) of each variable, and the constant a , and simply adding them together. In the fixed regression model framework in which we operate, there is no constraint on the nature of the IVs.³ To illustrate this, consider substituting other variables for specific variables in the equation. For example, we could replace U and V in Eq. (1.1.1) with U and V^2 , resulting in $Y = a + bU + cV^2$. Or, we could replace W with the logarithm of Z , resulting in $Y = a + d \log(Z)$. Or, we could replace X with a code variable representing sex (S , which takes values 0 = male and 1 = female), $Y = a + eS$. As our substitutions illustrate, the variables may be chosen to define relationships of *any* shape, rectilinear or curvilinear, or of no shape at all for unordered nominal independent variables, as well as all the complex combinations of these which multiple factors can produce.

Multiple regression equations are, indeed, linear; they are exactly of the form of Eq. (1.1.1). Yet they can be used to describe such a complex relationship as the length of psychiatric hospital stay as a function of ratings of patient symptoms on admission, diagnosis, age, sex, and average length of prior hospitalizations. This complex relationship is patently not rectilinear (straight line), yet it is readily described by a linear multiple regression equation.

To be sure, most relationships studied in the behavioral sciences are not of this order of complexity. But, the critical point is the capacity of MRC to represent any degree or type of shape—complexity is yet another of the important features which make it truly a *general* data-analytic system.

1.3.5 General and Conditional Relationships

Some relationships between Y and some factor C remain the same in regard to both degree and form despite variation in other factors D , E , F . In the MRC context, we will call such relationships *general* or *unconditional*: Readers familiar with ANOVA will know them as *main effects*. For example, suppose Y is a measure of visual acuity and C is age. In our example, both the form and degree of the relationship between visual acuity and age may remain the same under varying conditions of education level (D), ethnic group (E), and sex (F). The relationship between Y and C can then be said to be general insofar as the other specific factors are concerned. Note that this generality holds regardless of the form and degree of relationship between Y (visual acuity) and D , E , and F , between C (age) and D , E , and F , or among

³As we will note in Section 3.3, the “fixed” model we use throughout much of this book implies that we have generated or preselected the values of the IVs to which we wish to generalize.

D , E , and F . The $Y-C$ relationship can thus be considered unconditional with regard to, or independent of, D , E , and F .

Now consider the same research factors, but with Y as a measure of attitudes toward abortion. The form and/or degree of relationship of age to Y is now almost certain to vary as a function of one or more of the other factors: it may be stronger or shaped differently at lower educational levels than higher (D), and/or in one ethnic group or another (E), and/or for men compared to women (F). The relationship of Y to C is now said to be conditional on D and/or E and/or F . In ANOVA contexts, such relationships are called *interactions*. For example, if the $C-Y$ relationship is not constant over different values of D , there is said to be a $C \times D$ (age by educational level) interaction. Greater complexity is also possible: The $C-Y$ relationship may be constant over levels of D taken by themselves, and over levels of E taken by themselves, yet may be conditional on the *combination* of D and E levels. Such a circumstance would define a “three-way” interaction, represented as $C \times D \times E$. Interactions of even higher order, and thus even more complex forms of conditionality, are theoretically possible, although rarely reliably found because of the very large sample size typically required to detect them.

Some behavioral science disciplines have found it useful to discriminate two types of conditional relationships.⁴ *Moderation* indicates that the strength of the relationship between C and Y is reduced as the value of D increases. For example, researchers interested in the relationship between stress and illness report that social support moderates (weakens or buffers) this relationship. In contrast, *augmentation* or *synergy* means that the strength of the relationship between C and Y is increased as the value of D increases. Thus, moderation and augmentation describe particular forms of conditional relationships.

One facet of the complexity of the behavioral sciences is the frequency with which such conditional relationships are encountered. Relationships among variables often change with changes in experimental conditions (treatments, instructions, even experimental assistants), age, sex, social class, ethnicity, diagnosis, religion, personality traits, geographic area, etc. As essential as is the scientific task of estimating relationships between independent and dependent variables, it is also necessary to identify the conditions under which these estimates hold or change.

In summary, the generality of the MRC system of data analysis appropriately complements the complexity of the behavioral sciences, which complexity includes multiplicity and correlation among potential causal influences, a variety of forms in which information is couched, and variations in the shape and conditionality of relationships. Multiple regression/correlation also provides a full yield of measures of effect size with which to quantify various aspects of the strength of relationships (proportions of variance and correlation and regression coefficients). Finally, these measures are subject to statistical hypothesis testing, estimation, construction of confidence intervals, and power-analytic procedures.

1.4 ORIENTATION OF THE BOOK

This book was written to serve as a textbook and manual in the application of the MRC system for data analysis by students and practitioners in diverse areas of inquiry in the behavioral sciences, health sciences, education, and business. As its authors, we had to make many

⁴Elsewhere moderation may be used to describe both forms of conditional relationship. Whether a relationship may be considered to be moderated or augmented in the sense used here is entirely dependent on the (often arbitrary) direction of scoring of the IVs involved.

decisions about its level, breadth, emphasis, tone, and style of exposition. Readers may find it useful, at the outset, to have our orientation and the basis for these decisions set forth.

1.4.1 Nonmathematical

Our presentation of MRC is generally as conceptually oriented and nonmathematical as we could make it. Of course, **MRC is itself a product of mathematical statistics, based on matrix algebra, calculus, and probability theory.** There is little question that such a background makes possible a level of insight otherwise difficult to achieve. However, it is also our experience that some mathematically sophisticated scientists may lack the conceptual frame that links the mathematical procedures to the substantive scientific task in a particular case. When new mathematical procedures are introduced, we attempt to convey an intuitive conceptual rather than a rigorous mathematical understanding of the procedure. We have included a glossary at the end of the book in which the technical terms employed repeatedly in the book are given a brief conceptual definition. We hope that this aid will enable readers who have forgotten the meaning of a term introduced earlier to refresh their memories. Of course, most of these same terms also appear in the index with notation on the many times they may have been used. A separate table at the end of the book reviews the abbreviations used for the statistical terms in the book.

We thus abjure mathematical proofs, as well as unnecessary offhand references to mathematical concepts and methods not likely to be understood by the bulk of our audience. In their place, we heavily emphasize detailed and deliberately redundant verbal exposition of concrete examples. Our experience in teaching and consulting convinces us that our audience is richly endowed in the verbal, logical, intuitive kind of intelligence that makes it possible to understand how the MRC system works, and thus use it effectively (Dorothy Parker said, “**Flattery will get you anywhere.**”) This kind of understanding is eminently satisfactory (as well as satisfying), because it makes possible effective use of the system. We note that to drive a car, one does not need to be a physicist, nor an automotive engineer, nor even a highly skilled auto mechanic, although some of the latter’s skills are useful when one is stuck on the highway. That is the level we aim for.

We seek to make up for the absence of mathematical proofs by providing demonstrations instead. For example, the regression coefficient for a dichotomous or binary (e.g., male–female) independent variable that is scored 0–1 equals the difference between the two groups’ Y means. Instead of offering the six or seven lines of algebra that would constitute a mathematical proof, we demonstrate that it holds using a small set of data. True, this proves nothing, because the result may be accidental, but curious readers can check it out using their own or our data (and we urge that such checks be made throughout). Whether it is checked or not, we believe that most of our audience will profit more from the demonstration than the proof. If the absence of formal proof bothers some readers from Missouri (the “show me” state), all we can do is pledge our good faith.

1.4.2 Applied

The first word in this book’s title is *applied*. Our heavy stress on illustrations serves not only the function of clarifying and demonstrating the abstract principles being taught, but also that of exemplifying the kinds of applications possible. We attend to statistical theory only insofar as sound application makes it necessary. The emphasis is on “how to do it.” This opens us to the charge of writing a “cookbook,” a charge we deny because we do not neglect the whys and

wherefores. If the charge is nevertheless pressed, we can only add the observation that in the kitchen, cookbooks are likely to be more useful than textbooks in organic chemistry.

1.4.3 Data-Analytic

Mathematical statisticians proceed from exactly specified premises such as independent random sampling, normality of distributions, and homogeneity of variance. Through the exercise of ingenuity and appropriate mathematical theory, they arrive at exact and necessary consequences (e.g., the F distribution, statistical power functions). They are, of course, fully aware that no set of real data will exactly conform to the formal premises from which they start, but this is not properly their responsibility. As all mathematicians do, they work with abstractions to produce formal models whose “truth” lies in their self-consistency. Borrowing their language, we might say that inequalities are symmetrical: Just as behavioral scientists are not mathematicians, mathematicians are not behavioral scientists.

The behavioral scientist relies very heavily on the fruits of the labors of theoretical statisticians. Taken together with contributions from substantive theory and previous empirical research, statistical models provide guides for teasing out meaning from data, setting limits on inference, and imposing discipline on speculation (Abelson, 1995). Unfortunately, in the textbooks addressed to behavioral scientists, statistical methods have often been presented more as harsh straightjackets or Procrustean beds than as benign reference frameworks. Typically, a method is presented with some emphasis on its formal assumptions. Readers are advised that the failure of a set of data to meet these assumptions renders the method invalid. Alternative analytic strategies may not be offered. Presumably, the offending data are to be thrown away.

Now, this is, of course, a perfectly ridiculous idea from the point of view of working scientists. Their task is to contrive situations that yield information about substantive scientific issues—*they must and will analyze their data*. In doing so, they will bring to bear, in addition to the tools of statistical analysis and graphical display of the data, their knowledge of theory, past experience with similar data, hunches, and good sense, both common and uncommon (Krantz, 1999). They attempt to apply the statistical model that best matches their data; however, they would rather risk analyzing their data using a less than perfect model than not at all. For them, data analysis is not an end in itself, but the next-to-last step in a scientific process that culminates in providing information about the phenomenon. This is by no means to say that they need not be painstaking in their efforts to generate and perform analyses of the data. They need to develop statistical models to test their preferred scientific hypothesis, to rule out as many competing explanations for the results as they can, and to detect new relationships that may be present in the data. But, at the end they must translate these efforts into substantive information.

Most happily, the distinction between data analysis and statistical analysis has been made and given both rationale and respectability by one of our foremost mathematical statisticians, John Tukey. In his seminal *The Future of Data Analysis* (1962), Tukey describes data analysis as the special province of scientists with substantial interest in methodology. Data analysts employ statistical analysis as the most important tool in their craft, but they employ it together with other tools, and in a spirit quite different from that which has come to be associated with it from its origins in mathematical statistics. Data analysis accepts “inadequate” data, and is thus prepared to settle for “indications” rather than conclusions. It risks a greater frequency of errors in the interest of a greater frequency of occasions when the right answer is “suggested.” It compensates for cutting some statistical corners by using scientific as well as mathematical judgment, and by relying upon self-consistency and repetition of results. Data

analysis operates like a detective searching for clues that implicate or exonerate likely suspects (plausible hypotheses) rather than seeking to prove out a balance. In describing data analysis, Tukey has provided insight and rationale into the way good scientists have always related to data.

The spirit of this book is strongly data-analytic, in exactly this sense. We offer a variety of statistical models and graphical tools that are appropriate for common research questions in the behavioral sciences. We offer straightforward methods of examining whether the assumptions of the basic fixed-model MRC are met, and provide introductions to alternative analytic approaches that may be more appropriate when they are not. At the same time, we are aware that some data sets will fail to satisfy the assumptions of any standard statistical model, and that even when identified there may be little that the data analyst can do to bring the data “into line.” We recognize the limits on inference in such cases but are disposed to treat the limits as broad rather than narrow. We justify this by mustering whatever technical evidence there is in the statistical literature (especially evidence of the “robustness” of statistical tests), and by drawing upon our own and others’ practical experience, even upon our intuition, all in the interest of getting on with the task of making data yield their meaning. If we risk error, we are more than compensated by having a system of data analysis that is general, sensitive, and fully capable of reflecting the complexity of the behavioral sciences and thus of meeting the needs of behavioral scientists. And we will reiterate the injunction that no conclusions from a given set of data can be considered definitive: Replication is essential to scientific progress.

1.4.4 Inference Orientation and Specification Error

As noted earlier, perhaps the single most important reason for the broad adoption of MRC as a data-analytic tool is the possibility that it provides for taking into account—“controlling statistically or partialing”—variables that may get in the way of inferences about the influence of other variables on our dependent variable Y . These operations allow us to do statistically what we often cannot do in real life—separate the influences of variables that often, or even usually, occur together. This is often critically important in circumstances in which it is impossible or unethical to actually control one or more of these related variables. However, the centrality of this operation makes it critically important that users of these techniques have a basic, sound understanding of what partialing influences does and does not entail.

In emphasizing the extraction of meaning from data we will typically focus primarily on potential problems of “specification error” in the estimates produced in our analyses. Specification errors are errors of inference that we make because of the way we analyze our data. They include the assumption that the relationship between the dependent variable Y and each of the independent variables (IVs) is linear (constant over the range of the independent variables) when it is not, and that the relationships of some IVs to Y do not vary as a function of other IVs, when they do. When we attempt to make causal inferences on the basis of the relationships expressed in our MRC analyses, we may also make other kinds of specification errors, including assuming that Y is dependent on the IVs when some of the IVs are dependent on Y , or that the relationship between Y and certain IVs is causal when these relationships reflect the influence of common causes or confounders. Or assuming that the estimated relationship reflects the relationship between Y and the theoretically implicated (“true”) IV when it only reflects the relationship between Y and an imperfectly measured representative of the theoretically implicated IV. More technically, specification errors may include the conclusion that some relationship we seem to have uncovered in our sample data generalizes to the population, when our statistical analyses are biased by distributional or nonindependence problems in the data.

1.5 COMPUTATION, THE COMPUTER, AND NUMERICAL RESULTS

1.5.1 Computation

Like all mathematical procedures, MRC makes computational demands. The amount of computation increases with the size of the problem. Indeed, Darlington and Boyce (1982) estimate that computation time increases roughly with k^5 , where k is the number of IVs. Early in the book, in our exposition of bivariate correlation and regression and MRC with two independent variables, we give the necessary details with small worked examples for calculation by hand calculator. This is done because the intimate association with the arithmetic details makes plain to the reader the nature of the process: *exactly* what is being done, with what purpose, and with what result. With one to three independent variables, where the computation is easy, not only can one see the fundamentals, but a basis is laid down for generalization to many variables.

With most real problems, MRC requires the use of a computer. An important reason for the rapid increase in the use of MRC during the past three decades is the computer revolution. Widely available computers conduct analyses in milliseconds that would have taken months or even years in Fisher's time. Statistical software has become increasingly user friendly, with versions that allow either simple programming or "point and click" analysis. Graphical routines that permit insightful displays of the data and the results of statistical analyses have become increasingly available. These advances have had the beneficial effect of making the use of MRC analysis far faster and easier than in the past.

We have deliberately placed the extensive calculational details of the early chapters outside the body of the text to keep them from distracting attention from our central emphasis: understanding how the MRC system works. We strongly encourage readers to work through the details of the many worked illustrations using both a hand calculator and a statistical package. These can help provide a basic understanding of the MRC system and the statistical package.

But readers should then apply the methods of each chapter to *data of their own* or data with which they are otherwise familiar. The highest order of understanding is achieved from the powerful synergism of the application of unfamiliar methods to familiar data.

Finally, we caution readers about an unintended by-product of the ease of use of current statistical packages: Users can now easily produce misleading results. Some simple commonsense checks can often help avoid errors. Careful initial examination of simple statistics (means; correlations; number of cases) and graphical displays can often provide a good sense of the data, providing a baseline against which the results of more complicated analyses can be compared. We encourage readers using new software to try out the analysis first on a previously analyzed data set, and we include such data sets for the worked examples in the book, for which analyses have been carried out on the large SAS, SPSS, and SYSTAT statistical programs. Achieving a basic understanding of the MRC system and the statistical packages as well as careful checking of one's results is an important prerequisite to publication. There is no guarantee that the peer review process in journals will detect incorrect analyses.

1.5.2 Numerical Results: Reporting and Rounding

Statistical packages print out numerical results to several decimal places. For comparison purposes, we follow the general practice in this book of reporting computed correlation and regression coefficients rounded to two places (or significant digits) and squared coefficients rounded to three. When working with a hand calculator, the reader should be aware that small rounding errors will occur. Checks that agree within a few points in the third decimal may thus be taken as correct.

Following Ehrenberg (1977), we encourage readers to be conservative in the number of significant digits that are reported in their research articles. Despite the many digits of accuracy that characterize modern statistical programs, this level of accuracy only applies to the sample data. Estimates of population parameters are far less accurate because of sampling error. For the sample correlation (r) to provide an estimate of the population correlation (ρ) that is accurate to two decimal places would require as many as 34,000 cases (J. Cohen, 1990).

1.5.3 Significance Tests, Confidence Intervals, and Appendix Tables

Most behavioral scientists employ a hybrid of classical Fisherian and Neyman-Pearson null hypothesis testing (see Gigerenzer, 1993; Harlow, Mulaik, & Steiger, 1997), in which the probability of the sample result given that the null hypothesis is true, p , is compared to a prespecified significance criterion, α . If $p <$ (is less than) α , the null hypothesis is rejected and the sample result is deemed statistically significant at the α level of significance. The null hypothesis as typically specified is that the value of the parameter corresponding to the sample result is 0; other values can be specified based on prior research.

A more informative way of testing hypotheses in many applications is through the use of confidence intervals. Here an interval is developed around the sample result that would theoretically include the population value $(1 - \alpha)\%$ of the time in repeated samples. Used in conjunction with MRC procedures, the center of the confidence interval provides an estimate of the strength of the relationship and the width of the confidence interval provides information about the accuracy of that estimate. The lower and upper limits of the confidence interval show explicitly just how small and how large the effect size in the population (be it a regression coefficient, multiple R^2 , or partial r) might be. Incidentally, if the population value specified by the null hypothesis is not contained in the confidence interval, the null hypothesis is rejected.

The probability of the sample result given that the null hypothesis is true, p , is based on either the t or F distribution in basic MRC. Nearly all statistical packages now routinely compute exact values of p for each significance test. We also provide tables of F and t for $\alpha = .05$ and $\alpha = .01$. These values are useful for the construction of confidence intervals and for simple problems which can be solved with a hand calculator. The $\alpha = .05$ criterion is widely used as a standard in the behavioral sciences. The $\alpha = .01$ criterion is sometimes used by researchers as a matter of taste or tradition in their research area. We support this tradition when there are large costs of falsely rejecting the null hypothesis; however, all too frequently researchers adopt the $\alpha = .01$ level because they erroneously believe that this decision will necessarily make their findings stronger and more meaningful. The $\alpha = .01$ level is often used as a partial control on the incidence of spuriously significant results when a large number of hypothesis tests are being conducted. The choice of α also depends importantly on considerations of statistical power (the probability of rejecting the null hypothesis), which is discussed in several places, particularly in Section 4.5. We present tables for statistical power analysis in the Appendix; several programs are commercially available for conducting statistical power analyses on personal computers (e.g., Borenstein, Cohen, & Rothstein, 2001).

The statistical tables in the Appendix were largely abridged from Owen (1962) and from J. Cohen (1988). The entry values were selected so as to be optimally useful over a wide range of MRC applications. In rare cases in which the needed values are not provided, linear interpolation is sufficiently accurate for almost all purposes. Should more extensive tables be required, Owen (1962) and Pearson and Hartley (1970) are recommended. Some statistical packages will also compute exact p values for any specified df for common statistical distributions such as t , F , and χ^2 .

1.6 THE SPECTRUM OF BEHAVIORAL SCIENCE

When we address behavioral scientists, we are faced with an exceedingly heterogeneous audience. They range in level from student to experienced investigator and possess from modest to fairly advanced knowledge of statistical methods. With this in mind, we assume a minimum background for the basic exposition of the MRC system. When we must make assumptions about background that may not hold for some of our readers, we try hard to keep everyone on board. In some cases we use boxes in the text to present more technical information, which provides a greater understanding of the material. The boxes can be skipped on first reading without loss of continuity.

But it is with regard to substantive interests and investigative methods and materials that our audience is of truly mind boggling diversity. Behavioral science itself covers areas of “social”, “human”, and even “life” sciences—everything from the physiology of behavior to cultural anthropology, in both their “basic science” and “applied science” aspects. Add in health sciences, education, and business, and the substantive range becomes immense. Were it not for the fact that the methodology of science is inherently more general than its substance, a book of this kind would not be possible. This permits us to address substantive researchers whose primary interests lie in a bewildering variety of fields.

We have sought to accommodate to this diversity, even to capitalize upon it. Our illustrative examples are drawn from different areas, assuring the comfort of familiarity for most of our readers at least some of the time. Their content is presented at a level that makes them intellectually accessible to nonspecialists. We try to use the nontechnical discussion of the examples in a way that may promote some methodological cross-fertilization between fields of inquiry. Our hope is that this discussion may introduce better approaches to fields where data have been analyzed using traditional rather than more optimal procedures.

1.7 PLAN FOR THE BOOK

1.7.1 Content

Following this introductory chapter, we continue by introducing the origins and meanings of the coefficients that represent the relationship between two variables (Chapter 2). Chapter 3 extends these concepts and measures first to two independent variables and then to any larger number of independent variables. Chapter 4 expands on the graphical depiction of data, and particularly on the identification of data problems, and methods designed to improve the fit of the data to the assumptions of the statistical model. Chapter 5 describes the strategies that a researcher may use in applying MRC analyses to complex substantive questions, including selecting the appropriate statistical coefficients and significance tests. It continues by describing two widely useful techniques, hierarchical (sequential) analyses of data and the analysis of independent variables grouped into structural or functional sets.

Chapters 6 and 7 describe and illustrate the methods of identifying nonlinear and conditional relationships between independent variables and Y , beginning with methods for representing curvilinearity in linear equations. This chapter is followed by detailed presentations of the treatment and graphic display of interactions between scaled variables in their relationship with Y . Chapter 8 continues with the consideration of sets of independent variables representing mutually exclusive categories or groups. Relationships between scaled measures and Y may vary between sample subgroups; techniques for assessing and describing these interactions are reviewed in Chapter 9.

Chapter 10 presents the problem of multicollinearity among predictors and methods of controlling its extent. Chapter 11 details the full range of methods for coping with missing data in MRC, and the considerations appropriate for choosing among them.

Chapter 12 expands on the discussion of MRC applications to causal hypotheses that is found in earlier chapters and introduces the reader to some of the more complex methods of estimating such models and issues relevant to their employment.

Chapter 13 describes uses of the generalized linear model to analyze dependent variables that are dichotomous, ordered categories, or counts of rare phenomena. Chapter 14 introduces the reader to the multilevel analysis of data clusters arising from nonindependent sampling or treatment of participants.

Chapter 15 provides an introduction to a whole range of methods of analyzing data characterized by multiple observations of units over time. Beginning with simple repeated measure ANOVA and two time-point MRC, the chapter presents an overview of how the substantive questions and the structure of the data combine to suggest a choice among available sophisticated data analytic procedures.

The final chapter presents a multivariate method called set correlation that generalizes MRC to include sets (or partialled sets) of dependent variables and in so doing, generalizes multivariate methods and yields novel data-analytic forms.

For a more detailed synopsis of the book's contents, the reader is referred to the summaries at the ends of the chapters. The data for almost all examples in the book are also provided on the accompanying CD-ROM, along with the command codes for each of the major statistical packages that will yield the tabular and other findings presented in the chapters.

A note on notation. We have tried to keep the notation simultaneously consistent with the previous editions of this book and with accepted practice, insofar as possible. In general, we employ Greek letters for population estimates, but this convention falls down in two places. First, β is used conventionally both for the standardized regression coefficient and for the power: We have followed these conventions. Second, the maximum likelihood estimations methods discussed in Chapters 13 and 14 use a range of symbols, including Greek letters, designed to be distinct from those in use in OLS. We also use P and Q ($= P - 1.0$) to indicate proportions of samples, to distinguish this symbol from p = probability.

We have attempted to help the reader keep the major concepts in mind in two ways. We have included a glossary of technical terms at the end of the book, so that readers of later chapters may refresh their recall of terms introduced earlier in the book. We have also included a listing of the abbreviations of statistical terms, tests, and functions. In addition there are two technical appendices, as well as the appendix Tables.

One more difference between this edition and previous editions may be noted. In the introductory Chapter 2 we originally introduced equations using the sample standard deviation, with n in the denominator. This forced us into repeated explanations when later statistics required a shift to the sample-based population estimate with $n - 1$ in the denominator. The advantage was simplicity in the early equations. The serious disadvantage is that every statistical program determines sd with $n - 1$ in the denominator, and so students trying to check sds , z scores and other statistics against their computer output will be confused. In this edition we employ the population estimate sd consistently and adjust early equations as necessary.

1.7.2 Structure: Numbering of Sections, Tables, and Equations

Each chapter is divided into major sections, identified by the chapter and section numbers, for example, Section 5.4.3 ("Variance Proportions for Sets and the Ballantine Again") is the third

subsection of Section 5.4. Further subdivisions are not numbered, but titled with an italicized heading.

Tables, figures, and equations within the body of the text are numbered consecutively within major sections. Thus, for example, Figure 5.4.1 is the first figure in Section 5.4, and Eq. (2.6.5) is the fifth equation in Section 2.6. We follow the usual convention of giving equation numbers in parentheses. A similar plan is followed in the two appendices. The reference statistical tables make up a separate appendix and are designated as Appendix Tables A through G.

On the accompanying data disk each chapter has a folder; within that folder each example for which we provide data and syntax/command files in SAS, SPSS, and SYSTAT has a folder.

1.8 SUMMARY

This introductory chapter begins with an overview of MRC as a data-analytic system, emphasizing its generality and superordinate relationship to the analysis of variance/covariance (Section 1.1). MRC is shown to be peculiarly appropriate for the behavioral sciences in its capacity to accommodate the various types of complexity that characterize them: the multiplicity and correlation among causal influences, the varieties of form of information and shape of relationship, and the frequent incidence of conditional (interactive) relationships. The special relevance of MRC to the formal analysis of causal models is described (Section 1.2).

The book's exposition of MRC is nonmathematical, and stresses informed application to scientific and technological problems in the behavioral sciences. Its orientation is "data analytic" rather than statistical analytic, an important distinction that is discussed. Concrete illustrative examples are heavily relied upon (Section 1.3).

The popularity of MRC in the analysis of nonexperimental data for which manipulation of variables is impossible or unethical hinges on the possibility of statistical control or partialing. The centrality of this procedure, and the various kinds of errors of inferences that can be made when the equations include specification error are discussed (Section 1.4).

The means of coping with the computational demands of MRC are briefly described and largely left to the computer, with details relegated to appendices so as not to distract the reader's attention from the conceptual issues (Section 1.5). We acknowledge the heterogeneity of background and substantive interests of our intended audience, and discuss how we try to accommodate to it and even exploit it to pedagogical advantage (Section 1.6).

The chapter ends with a brief outline of the book and the scheme by which sections, tables, figures, and equations are numbered.

2

Bivariate Correlation and Regression

One of the most general meanings of the concept of a relationship between a pair of variables is that knowledge with regard to one of the variables carries information about the other. Information about the height of a child in elementary school has implications for the probable age of the child, and information about the occupation of an adult can lead to more accurate guesses about her income level than could be made in the absence of that information.

2.1 TABULAR AND GRAPHIC REPRESENTATIONS OF RELATIONSHIPS

Whenever data have been gathered on two quantitative variables for a set of subjects or other units, the relationship between the variables may be displayed graphically by means of a scatterplot.

For example, suppose we have scores on a vocabulary test and a digit-symbol substitution task for 15 children (see Table 2.1.1). If these data are plotted by representing each child as a point on a graph with vocabulary scores on the horizontal axis and the number of digit symbols on the vertical axis, we would obtain the scatterplot seen in Fig. 2.1.1. The circled dot, for example, represents Child 1, who obtained a score of 5 on the vocabulary test and completed 12 digit-symbol substitutions.

When we inspect this plot, it becomes apparent that the children with higher vocabulary scores tended to complete more digit symbols (d-s) and those low on vocabulary (v) scores were usually low on d-s as well. This can be seen by looking at the average of the d-s scores, M_{d_s} , corresponding to each v score given at the top of the figure. The child receiving the lowest v score, 5, received a d-s score of 12; the children with the next lowest v score, 6, obtained an average d-s score of 14.67, and so onto the highest v scorers, who obtained an average of 19.5 on the d-s test. A parallel tendency for vocabulary scores to increase is observed for increases in d-s scores. The form of this relationship is said to be positive, because high values on one variable tend to go with high values on the other variable and low with low values. It may also be called linear because the tendency for a unit increase in one variable to be accompanied by a constant increase in the other variable is (fairly) constant throughout the scales. That is, if we

TABLE 2.1.1
Illustrative Set of Data on Vocabulary
and Digit-Symbol Tests

Child (no.)	Vocabulary	Digit-symbol
1	5	12
2	8	15
3	7	14
4	9	18
5	10	19
6	8	18
7	6	14
8	6	17
9	10	20
10	9	17
11	7	15
12	7	16
13	9	16
14	6	13
15	8	16

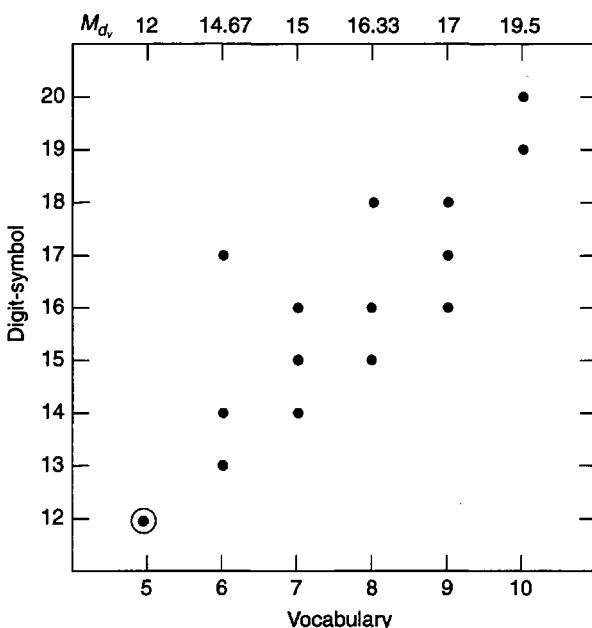


FIGURE 2.1.1 A strong, positive linear relationship.

were to draw the straight line that best fits the average of the d-s values at each v score (from the lower left-hand corner to the upper right-hand corner) we would be describing the trend or shape of the relationship quite well.

Figure 2.1.2 displays a similar scatterplot for age and the number of seconds needed to complete the digit-symbol task. In this case, low scores on age tended to go with high test time in seconds and low test times were more common in older children. This relationship may be

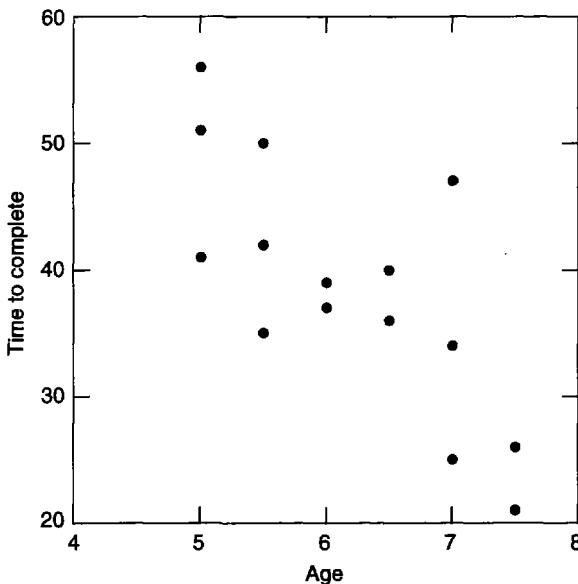


FIGURE 2.1.2 A negative linear relationship.

said to be negative and linear. It should also be clear at this point that whether a relationship between two variables is positive or negative is a direct consequence of the direction in which the two variables have been scored. If, for example, the vocabulary scores from the first example were taken from a 12-item test, and instead of scoring the number correct a count was made of the number wrong, the relationship with d-s scores would be negative. Because such scoring decisions in many cases may be essentially arbitrary, it should be kept in mind that any positive relationship becomes negative when either (but not both) of the variables is reversed, and vice versa. Thus, for example, a negative relationship between age of oldest child and income for a group of 30-year-old mothers implies a positive relationship between age of first becoming a mother and income.¹

Figure 2.1.3 gives the plot of a measure of motivational level and score on a difficult d-s task. It is apparent that the way motivation was associated with performance score depends on whether the motivational level was at the lower end of its scale or near the upper end. Thus, the relationship between these variables is curvilinear. Finally, Fig. 2.1.4 presents a scatterplot for age and number of substitution errors. This plot demonstrates a general tendency for higher scores on age to go with fewer errors, indicating that there is, in part, a negative linear relationship. However, it also shows that the decrease in errors that goes with a unit increase in age was greater at the lower end of the age scale than it was at the upper end, a finding that indicates that although a straight line provides some kind of fit, clearly it is not optimal.

Thus, scatterplots allow visual inspection of the form of the relationship between two variables. These relationships may be well described by a straight line, indicating a rectilinear (negative or positive) relationship, or they may be better described by a line with one or more curves. Because approximately linear relationships are very common in all sorts of data, we will concentrate on these in the current discussion, and will present methods of analyzing nonlinear relationships in Chapter 6.

¹Here we follow the convention of naming a variable for the upper end of the scale. Thus, a variable called *income* means that high numbers indicate high income, whereas a variable called *poverty* would mean that high numbers indicate much poverty and therefore low income.

22 2. BIVARIATE CORRELATION AND REGRESSION

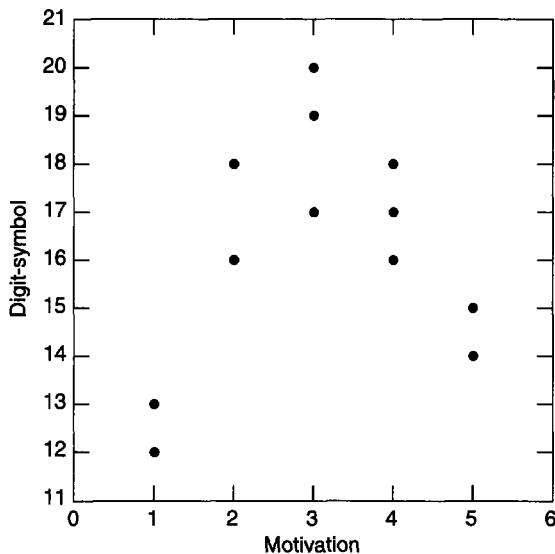


FIGURE 2.1.3 A positive curvilinear relationship.

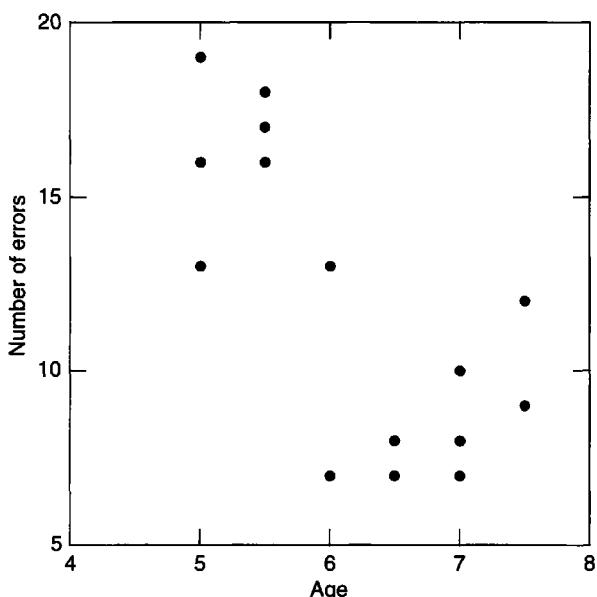


FIGURE 2.1.4 A negative curvilinear relationship.

Now suppose that Fig. 2.1.1 is compared with Fig. 2.1.5. In both cases the relationship between the variables is linear and positive; however, it would appear that vocabulary provided better information with regard to d-s completion than did chronological age. That is, the degree of the relationship with performance seems to be greater for vocabulary than for age because one could make more accurate estimates of d-s scores using information about vocabulary than using age. To compare these two relationships to determine which is greater, we need an index of the degree or strength of the relationship between two variables that will be comparable from one pair of variables to another. Looking at the relationship between v and d-s scores,

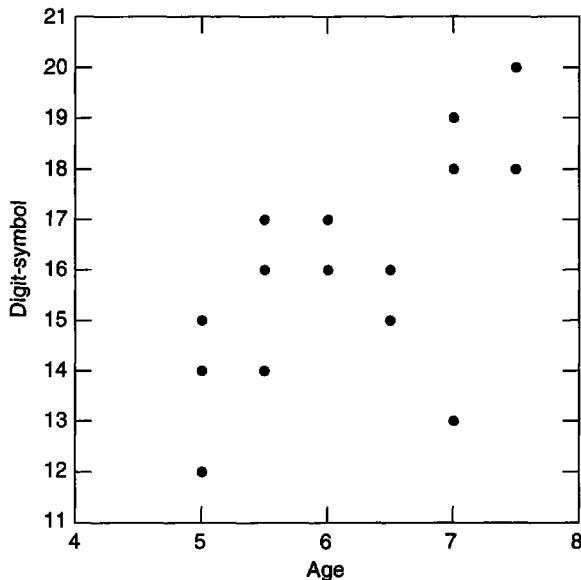


FIGURE 2.1.5 A weak, positive linear relationship.

other questions come to mind: Should this be considered a strong or weak association? On the whole, how great an increase in digit-symbol score is found for a given increase in vocabulary score in this group? If $d-s$ is estimated from v in such a way as to minimize the differences between our estimations and the actual $d-s$ scores, how much error will, nevertheless, be made? If this is a random sample of subjects from a larger population, how much confidence can we have that v and $d-s$ are linearly related in the entire population? These and other questions are answered by correlation and regression methods. In the use and interpretation of these methods the two variables are generally treated as interval scales; that is, constant differences between scale points on each variable are assumed to represent equal "amounts" of the construct being measured. Although for many or even most scales in the behavioral sciences this assumption is not literally true, empirical work (Baker, Hardyck, & Petrinovich, 1966) indicates that small to moderate inequalities in interval size produce little if any distortion in the validity of conclusions based on the analysis. This issue is discussed further in Chapter 6.

2.2 THE INDEX OF LINEAR CORRELATION BETWEEN TWO VARIABLES: THE PEARSON PRODUCT MOMENT CORRELATION COEFFICIENT

2.2.1 Standard Scores: Making Units Comparable

One of the first problems to be solved by an index of the degree of association between two variables is that of measurement unit. Because the two variables are typically expressed in different units, we need some means of converting the scores to comparable measurement units. It can be readily perceived that any index that would change with an arbitrary change in measurement unit—from inches to centimeters or age in months to age in weeks, for example—could hardly be useful as a general description of the strength of the relationship between height and age, one that could be compared with other such indices.

TABLE 2.2.1
**Income and Major Household Appliances in Original Units,
 Deviation Units, and z Units**

Household	Income	Appliances	$I - M_I = i$	$A - M_A = a$	i^2	a^2	Rank I	Rank A
1	24,000	3	-3,500	-1.75	12,250,000	3.0625	1	1
2	29,000	7	+1,500	+2.25	2,250,000	5.0625	3	4
3	27,000	4	-500	-.75	250,000	.5625	2	2
4	30,000	5	+2,500	+.25	6,250,000	.0625	4	3
Sum (Σ)	110,000	19	0	0	21,000,000	8.75		
Mean	27,500	4.75						
$sd_I^2 = \Sigma i^2 / (n - 1) = 7,000,000$								
$sd_I = \sqrt{\Sigma i^2 / (n - 1)} = 2,645.75$								
	$i / sd_I = z_I$	$a / sd_A = z_A$	z_I^2	z_A^2				
1	-1.323	-1.025	1.750	1.050				
2	+0.567	+1.317	0.321	1.736				
3	-0.189	-0.439	0.036	0.193				
4	+0.945	+0.146	0.893	0.021				
Σ	0	0	3.00	3.00				

To illustrate this problem, suppose information has been gathered on the annual income and the number of major household appliances of four households (Table 2.2.1).² In the effort to measure the degree of relationship between income (I) and the number of appliances (A), we will need to cope with the differences in the nature and size of the units in which the two variables are measured. Although Households 1 and 3 are both below the mean on both variables and Households 2 and 4 are above the mean on both (see i and a , scores expressed as deviations from their means, with the means symbolized as M_I and M_A , respectively), we are still at a loss to assess the correspondence between a difference of \$3500 from the mean income and a difference of 1.5 appliances from the mean number of appliances. We may attempt to resolve the difference in units by ranking the households on the two variables—1, 3, 2, 4 and 1, 4, 2, 3, respectively—and noting that there seems to be some correspondence between the two ranks. In so doing we have, however, made the differences between Households 1 and 3 (\$3000) equal to the difference between Households 2 and 4 (\$1000); two ranks in each case.

To make the scores comparable, we clearly need some way of taking the different variability of the two original sets of scores into account. Because the standard deviation (sd) is an index of variability of scores, we may measure the discrepancy of each score from its mean (x) relative to the variability of all the scores by dividing by the sd :

$$(2.2.1) \quad sd_x = \sqrt{\frac{\sum x^2}{n - 1}},$$

²In this example, as in all examples that follow, the number of cases (n) is kept very small in order to facilitate the reader's following of the computations. In almost any serious research, the n must, of course, be very much larger (Section 2.9).

where Σx^2 means "the sum of the squared deviations from the mean."³ The scores thus created are in standard deviation units and are called *standard* or *z* scores:

$$(2.2.2) \quad z_X = \frac{X - M_X}{sd_X} = \frac{x}{sd_X}.$$

In Table 2.2.1 the *z* score for income for Household 1 is -1.323 , which indicates that its value (\$24,000) falls about $1\frac{1}{3}$ income standard deviations (\$2646) *below* the income mean (\$27,500). Although income statistics are expressed in dollar units, the *z* score is a pure number; that is, it is unit-free. Similarly, Household 1 has a *z* score for number of appliances of -1.025 , which indicates that its number of appliances (3) is about 1 standard deviation (1.71) below the mean number of appliances (4.75). Note again that -1.025 is not expressed in number of appliances, but is also a pure number. Instead of having to compare \$24,000 and 3 appliances for Household 1, we can now make a meaningful comparison of -1.323 (z_I) and -1.025 (z_A), and note incidentally the similarity of the two values for Household 1. This gives us a way of systematically approaching the question of whether a household is as relatively wealthy as it is relatively "appliance."

It should be noted that the rank of the *z* scores is the same as that of the original scores and that scores that were above or below the mean on the original variable retain this characteristic in their *z* scores. In addition, we note that the difference between the incomes of Households 2 and 3 ($I_2 - I_3 = \$2000$) is twice as large, and of opposite direction to the difference between Households 2 and 4 ($I_2 - I_4 = -\$1000$). When we look at the *z* scores for these same households, we find that $z_{I2} - z_{I3} = .567 - (-.189) = .756$ is twice as large and of opposite direction to the difference $z_{I2} - z_{I4} = .567 - .945 = -.378$ (i.e., $.756 / -.378 = -2$). Such proportionality of differences or distances between scores,

$$(2.2.3) \quad \frac{X_i - X_j}{X_m - X_n} = \frac{z_{X_i} - z_{X_j}}{z_{X_m} - z_{X_n}}$$

is the essential element in what is meant by retaining the original relationship between the scores. This can be seen more concretely in Fig. 2.2.1, in which we have plotted the pairs of scores. Whether we plot *z* scores or raw scores, the points in the scatterplot have the same relationship to each other.

The *z* transformation of scores is one example of a linear transformation. A linear transformation is one in which every score is changed by multiplying or dividing by a constant or adding or subtracting a constant or both. Changes from inches to centimeters, dollars to francs, and Fahrenheit to Celsius degrees are examples of linear transformations. Such transformations will, of course, change the means and *sds* of the variables upon which they are performed. However, because the *sd* will change by exactly the same factor as the original scores (that is, by the constant by which scores have been multiplied or divided) and because *z* scores are created by subtracting scores from their mean, all linear transformations of scores will yield the same set of *z* scores. (If the multiplier is negative, the signs of the *z* scores will simply be reversed.)

Because the properties of *z* scores form the foundation necessary for understanding correlation coefficients, they will be briefly reviewed:

³As noted earlier, this edition employs the population estimate of *sd* with $n - 1$ in the denominator throughout to conform with computer program output, in contrast to earlier editions, which employed the sample *sd* with n in the denominator in earlier equations in the book and moved to the population estimate when inferences to the population involving standard errors were considered, and thereafter.

Also note that the summation sign, Σ , is used to indicate summation over all n cases here and elsewhere, unless otherwise specified.

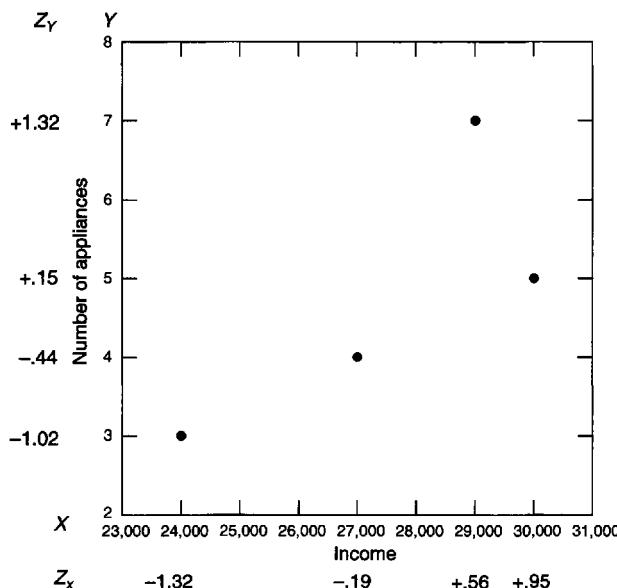


FIGURE 2.2.1 Household income and number of appliances.

1. The sum of a set of z scores (Σz) (and therefore also the mean) equals 0.
2. The variance (sd^2) of the set of z scores equals 1, as does the standard deviation (sd).
3. Neither the shape of the distribution of X nor its absolute correlation with any other variable is affected by transforming it to z (or any other linear transformation).

2.2.2 The Product Moment Correlation as a Function of Differences Between z Scores

We may now define a perfect (positive) relationship between two variables (X and Y) as existing when all z_X and z_Y pairs of scores consist of two exactly equal values. Furthermore, the degree of relationship will be a function of the departure from this "perfect" state, that is, a function of the differences between pairs of z_X and z_Y scores. Because the average difference between paired z_X and z_Y and is necessarily zero (because $M_{z_Y} = M_{z_X} = 0$), the relationship may be indexed by finding the average⁴ of the squared discrepancies between z scores, $\Sigma(z_X - z_Y)^2/n$.

For example, suppose that an investigator of academic life obtained the (fictitious) data shown in Table 2.2.2. The subjects were 15 randomly selected members of a large university department, and the data include the time in years that had elapsed since the faculty member's Ph.D. was awarded and the number of publications in professional journals.

Several things should be noted in this table. Deviation scores ($x = X - M_X$ and $y = Y - M_Y$) sum to zero. So do z_X and z_Y . The standard deviations, sd_{z_X} and sd_{z_Y} , are both 1, M_{z_X} and M_{z_Y} are both 0 (all of which are mathematical necessities), and these equalities reflect the equal footing on which we have placed the two variables.

We find that the squared differences (Σs^2) between z scores sums to 9.614, which when divided by the number of paired observations equals .641. How large is this relationship? We have stated that if the two variables were perfectly (positively) related, all z score differences

⁴Because we have employed the sample-based estimate of the population sd , with a divisor of $n - 1$, when z scores have been based on this sd this equation should also use $n - 1$.

TABLE 2.2.2
***z* Scores, *z* Score Differences, and *z* Score Products on Data Example**

Case	X Time since Ph.D.	Y No. of publications	$\frac{X_i - M_X}{sd_X} = z_{X_i}$	$\frac{Y_i - M_Y}{sd_Y} = z_{Y_i}$	$z_X - z_Y$	$z_X z_Y$
1	3	18	-1.020	-.140	-.880	.142
2	6	3	-.364	-1.225	.861	.446
3	3	2	-1.020	-1.297	.278	1.322
4	8	17	.073	-.212	.285	-.015
5	9	11	.291	-.646	.938	-.188
6	6	6	-.364	-1.008	.644	.367
7	16	38	1.821	1.307	.514	2.380
8	10	48	.510	2.030	-1.520	1.035
9	2	9	-1.238	-.791	-.447	1.035
10	5	22	-.583	.150	-.732	-.087
11	5	30	-.583	.728	-1.311	-.424
12	6	21	-.364	.077	-.441	-.028
13	7	10	-.146	-.719	.573	.105
14	11	27	.728	.511	.217	.372
15	18	37	2.257	1.235	1.023	2.787
Σ	115	299	0	0	0	
Σ squared	1235	8635	14	14	9.614	
M	7.67	19.93			.641	.613
sd^2	19.55	178.33	1	1		
sd	4.42	13.35	1	1		

would equal zero and necessarily their sum and mean would also be zero. A perfect negative relationship, on the other hand, may be defined as one in which the *z* scores in each pair are equal in absolute value but opposite in sign. Under the latter circumstances, it is demonstrable that the average of the squared discrepancies times $n/(n - 1)$ always equals 4. It can also be proved that under circumstances in which the pairs of *z* scores are on the average equally likely to be consistent with a negative relationship as with a positive relationship, the average squared difference times $n/(n - 1)$ will always equal 2, which is midway between 0 and 4. Under these circumstances, we may say that there is no linear relationship between *X* and *Y*.⁵

Although it is clear that this index, ranging from 0 (for a perfect positive linear relationship) through 2 (for no linear relationship) to 4 (for a perfect negative one), does reflect the relationship between the variables in an intuitively meaningful way, it is useful to transform the scale linearly to make its interpretation even more clear. Let us reorient the index so that it runs from -1 for a perfect negative relationship to +1 for a perfect positive relationship. If we divide the sum of the squared discrepancies by $2(n - 1)$ and subtract the result from 1, we have

$$(2.2.4) \quad r = 1 - \left(\frac{\sum (z_X - z_Y)^2}{2(n - 1)} \right),$$

⁵Note that this equation is slightly different from that in earlier editions. The $n/(n - 1)$ term is necessary because the *sd* used here is the sample estimate of the population *sd* rather than the sample *sd* which uses *n* in the denominator.

which for the data of Table 2.2.2 gives

$$r = r = 1 - \left(\frac{9.614}{28} \right) = .657.$$

r is the product moment correlation coefficient, invented by Karl Pearson in 1895.⁶ This coefficient is the standard measure of the linear relationship between two variables and has the following properties:

1. It is a pure number and independent of the units of measurement.
2. Its value varies between zero, when the variables have no linear relationship, and +1.00 or -1.00, when each variable is perfectly estimated by the other. The absolute value thus gives the degree of relationship.
3. Its sign indicates the direction of the relationship. A positive sign indicates a tendency for high values of one variable to occur with high values of the other, and low values to occur with low. A negative sign indicates a tendency for high values of one variable to be associated with low values of the other. Reversing the direction of measurement of one of the variables will produce a coefficient of the same absolute value but of opposite sign. Coefficients of equal value but opposite sign (e.g., +.50 and -.50) thus indicate equally strong linear relationships, but in opposite directions.

2.3 ALTERNATIVE FORMULAS FOR THE PRODUCT MOMENT CORRELATION COEFFICIENT

The formula given in Eq. (2.2.4) for the product moment correlation coefficient as a function of squared differences between paired z scores is only one of a number of mathematically equivalent formulas. Some of the other versions provide additional insight into the nature of r ; others facilitate computation. Still other formulas apply to particular kinds of variables, such as variables for which only two values are possible, or variables that consist of rankings.

2.3.1 r as the Average Product of z Scores

It follows from algebraic manipulation of Eq. (2.2.4) that

$$(2.3.1) \quad r_{XY} = \frac{\sum z_X z_Y}{n - 1}$$

The product moment correlation is therefore seen to be the mean of the products of the paired z scores.⁷ In the case of a perfect positive correlation, because $z_X = z_Y$,

$$r_{XY} = \frac{\sum z_X z_Y}{n - 1} = \frac{\sum z^2}{n - 1} = 1.$$

For the data presented in Table 2.2.1, these products have been computed and $r_{XY} = 9.193/14 = .657$, necessarily as before.

⁶The term *product moment* refers to the fact that the correlation is a function of the product of the *first moments*, of X and Y , respectively. See the next sections.

⁷If we used z s based on the *sample sd* which divides by n , this *average* would also divide by n .

2.3.2 Raw Score Formulas for r

Because z scores can be readily reconverted to the original units, a formula for the correlation coefficient can be written in raw score terms. There are many mathematically equivalent versions of this formula, of which the following is a convenient one for computation by computer or calculator:

$$(2.3.2) \quad r_{XY} = \frac{n \sum XY - \sum X \sum Y}{\sqrt{[n \sum X^2 - (\sum X)^2][n \sum Y^2 - (\sum Y)^2]}}.$$

When the numerator and denominator are divided by n^2 , Eq. (2.3.2) becomes an expression of r in terms of the means of each variable, of each squared variable, and of the XY product:

$$(2.3.3) \quad r_{XY} = \frac{M_{XY} - M_X M_Y}{\sqrt{(M_X^2 - M_{X^2})(M_Y^2 - M_{Y^2})}}.$$

It is useful for hand computation to recognize that the denominator is the product of the variables' standard deviations, thus an alternative equivalent is

$$(2.3.4) \quad r_{XY} = \frac{\sum xy/(n - 1)}{sd_X sd_Y}$$

This numerator, based on the product of the *deviation* scores is called the *covariance* and is an index of the tendency for the two variables to *covary* or go together that is expressed in deviations measured in the original units in which X and Y are measured (e.g., income in *dollars* and *number* of appliances). Thus, we can see that r is an expression of the covariance between standardized variables, because if we replace the *deviation* scores with *standardized* scores, Eq. (2.3.4) reduces to Eq. (2.3.1).

It should be noted that r inherently is *not* a function of the number of observations and that the $n - 1$ in the various formulas serves only to cancel it out of other terms where it is hidden (for example, in the sd). By multiplying Eq. (2.3.4) by $(n - 1)/(n - 1)$ it can be completely canceled out to produce a formula for r that does not contain any vestige of n :

$$(2.3.5) \quad r_{XY} = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

2.3.3 Point Biserial r

When one of the variables to be correlated is a dichotomy (it can take on only two values), the computation of r simplifies. There are many dichotomous variables in the behavioral sciences, such as yes or no responses, left- or right-handedness, and the presence or absence of a trait or attribute. For example, although the variable "gender of subject" does not seem to be a quantitative variable, it may be looked upon as the presence or absence of the characteristics of being female (or of being male). As such, we may decide, arbitrarily, to score all females as 1 and all males as 0. Under these circumstances, the sd of the gender variable is determined by the proportion of the total n in each of the two groups; $sd = \sqrt{PQ}$, where P is the proportion in one group and $Q = 1 - P$, the proportion in the other group.⁸ Because r indicates a relationship between two standardized variables, it does not matter whether we choose 0 and 1 as the two values or any other pair of different values, because any pair will yield the same absolute z scores.

⁸Note that here the sd is the sample sd (divided by n) rather than the sample-based estimate of the population σ . As noted earlier, because the ns in the equation for r cancel, this difference is immaterial here.

TABLE 2.3.1
Correlation Between a Dichotomous and a Scaled Variable

Subject no.	Stimulus condition (X)	Task score (Y)	X _A	X _B	Z _Y	Z _A	Z _B	Z _Y Z _A	Z _Y Z _B
1	NONE	67	0	50	-0.41	-.802	.802	0.329	-0.329
2	NONE	72	0	50	1.63	-.802	.802	-1.307	1.307
3	NONE	70	0	50	0.81	-.802	.802	-0.650	0.650
4	NONE	69	0	50	0.41	-.802	.802	-0.329	0.329
5	STIM	66	1	20	-0.81	1.069	-1.069	-0.866	0.866
6	STIM	64	1	20	-1.63	1.069	-1.069	-1.742	1.742
7	STIM	68	1	20	0	1.069	-1.069	0	0
Sum		476	3	260	0	0	0	-4.565	4.565
Mean		68	.429	37.14	0	0	0		
sd in sample		2.45	.495	14.9		M _Y NONE = 69.5 M _Y STIM = 66.0			

For example, Table 2.3.1 presents data on the effects of an interfering stimulus on task performance for a group of seven experimental subjects. As can be seen, the absolute value of the correlation remains the same whether we choose (X_A) 0 and 1 as the values to represent the absence or presence of an interfering stimulus or choose (X_B) 50 and 20 as the values to represent the same dichotomy. The sign of *r*, however, depends on whether the group with the higher mean on the other (*Y*) variable, in this case the no-stimulus group, has been assigned the higher or lower of the two values. The reader is invited to try other values and observe the constancy of *r*.

Because the *z* scores of a dichotomy are a function of the proportion of the total in each of the two groups, the product moment correlation formula simplifies to

$$(2.3.6) \quad r_{pb} = \frac{(M_{Y_1} - M_{Y_0})\sqrt{PQ}}{sd_Y},$$

where *M_{Y₁}* and *M_{Y₀}* are the *Y* means of the two groups of the dichotomy and the *sd_Y* is the sample value, which is divided by *n* rather than *n* - 1. The simplified formula is called the point biserial *r* to take note of the fact that it involves one variable (*X*) whose values are all at one of two points and one continuous variable (*Y*). In the present example,

$$(2.3.7) \quad r_{pb} = \frac{(66.0 - 69.5)\sqrt{(.429)(.571)}}{2.45} = -.707.$$

The point biserial formula for the product moment *r* displays an interesting and useful property. When the two groups of the dichotomy are of equal size, *p* = *q* = .5, so $\sqrt{PQ} = .5$. The *r_{pb}* then equals half the difference between the means of the *z* scores for *Y*, and so $2r_{pb}$ equals the difference between the means of the standardized variable.

2.3.4 Phi (ϕ) Coefficient

When both *X* and *Y* are dichotomous, the computation of the product moment correlation is even further simplified. The data may be represented by a fourfold table and the correlation computed directly from the frequencies and marginals. For example, suppose a study investigated the

TABLE 2.3.2
Fourfold Frequencies for Candidate Preference
and Homeowning Status

		Candidate U	Candidate V	Total
Homeowners	A	B		
	19	54		$73 = A + B$
Nonhomeowners	C	D		
	60	52		$112 = C + D$
Total	$79 = A + C$	$106 = B + D$		$185 = n$

difference in preference of homeowners and nonhomeowners for the two candidates in a local election, and the data are as presented in Table 2.3.2. The formula for r here simplifies to the difference between the product of the diagonals of a fourfold table of frequencies divided by the square root of the product of the four marginal sums:

$$(2.3.8) \quad r_{\phi} = \frac{BC - AD}{\sqrt{(A+B)(C+D)(A+C)(B+D)}} \\ = \frac{(54)(60) - (19)(52)}{\sqrt{(73)(112)(79)(106)}} = -.272$$

Once again it may be noted that this is a computing alternative to the z score formula, and therefore it does not matter what two values are assigned to the dichotomy because the standard scores, and hence the absolute value of r_{ϕ} will remain the same. It also follows that unless the division of the group is the same for the two dichotomies ($P_Y = P_X$ or Q_X), their z scores cannot have the same values and r_{ϕ} cannot equal 1 or -1. A further discussion of this limit is found in Section 2.10.1.

2.3.5 Rank Correlation

Yet another simplification in the product moment correlation formula occurs when the data being correlated consist of two sets of ranks. Such data indicate only the ordinal position of the subjects on each variable; that is, they are at the ordinal level of measurement. This version of r is called the Spearman rank correlation (r_S). Because the sd of a complete set of ranks is a function only of the number of objects being ranked (assuming no ties), some algebraic manipulation yields

$$(2.3.9) \quad r_S = 1 - \frac{6 \sum d^2}{n(n^2 - 1)},$$

where d is the difference in the ranks of the pair for an object or individual. In Table 2.3.3 a set of 5 ranks is presented with their deviations and differences. Using one of the general formulas (2.3.4) for r ,

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} \\ = \frac{-3}{\sqrt{10} \sqrt{10}} = -.300$$



TABLE 2.3.3
Correlation Between Two Sets of Ranks

I.D.	X	Y	x	x^2	y	y^2	xy	d	d^2
1	4	2	1	1	-1	1	-1	2	4
2	2	1	-1	1	-2	4	2	1	1
3	3	4	0	0	1	1	0	-1	1
4	5	3	2	4	0	0	0	2	4
5	1	5	-2	4	2	4	-4	-4	16
Sum	15	15	0	10	0	10	-3	0	26

TABLE 2.3.4
Product Moment Correlation Coefficients
for Special Kinds of Data

Data type	Coefficient
A scaled variable and a dichotomous variable	Point biserial $r (r_{pb})$
Two dichotomous variables	ϕ or r_ϕ
Two ranked variables	Spearman rank order $r (r_S)$

The rank order formula (2.3.9) with far less computation yields

$$\begin{aligned}
 r_S &= 1 - \frac{6(26)}{5(24)} \\
 &= 1 - \frac{156}{120} = -.300,
 \end{aligned}$$

which agrees with the result from Eq. (2.3.4).

We wish to stress the fact that the formulas for r_{pb} , r_ϕ , and r_S are simply computational equivalents of the previously given general formulas for r that result from the mathematical simplicity of dichotomous or rank data (Table 2.3.4). They are of use when computation is done by hand or calculator. They are of no significance when computers are used, because whatever formula for r the computer uses will work when variables are scored 0–1 (or any other two values) or are ranks without ties. It is obviously not worth the trouble to use special programs to produce these special-case versions of r when a formula such as Eq. (2.3.2) will produce them.

2.4 REGRESSION COEFFICIENTS: ESTIMATING Y FROM X

Thus far we have treated the two variables as if they were of equal status. It is, however, often the case that variables are treated asymmetrically, one being thought of as a dependent variable or criterion and the other as the independent variable or predictor. These labels reflect the reasons why the relationship between two variables may be under investigation. There are two reasons for such investigation; one scientific and one technological. The primary or scientific question looks upon one variable as potentially causally dependent on the other, that is, as in part an effect of or influenced by the other. The second or technological question has for its goal forecasting, as for example, when high school grades are used to predict college

grades with no implication that the latter are actually caused by the former. In either case the measure of this effect will, in general, be expressed as the number of units of change in the Y variable per unit change in the X variable.

To return to our academic example of 15 faculty members presented in Table 2.2.2, we wish to obtain an estimate of Y , for which we use the notation \hat{Y} , which summarizes the average amount of change in the number of publications for each year since Ph.D. To find this number, we will need some preliminaries. Obviously, if the relationship between publications and years were perfect and positive, we could provide the number of publications corresponding to any given number of years since Ph.D. simply by adjusting for differences in scale of the two variables. Because, when $r_{XY} = 1$, for any individual j , the estimated \hat{z}_Y simply equals z_X , then

$$\frac{\hat{Y}_j - M_Y}{sd_Y} = \frac{X_j - M_X}{sd_X},$$

and solving for j 's estimated value of Y ,

$$\hat{Y}_j = \frac{sd_Y(X_j - M_X)}{sd_X} + M_Y,$$

and because M_X , M_Y , and sd_Y are known, it remains only to specify X_j and then \hat{Y}_j may be computed.

When, however, the relationship is not perfect, we may nevertheless wish to show the estimated \hat{Y} that we would obtain by using the best possible "average" conversion or prediction rule from X in the sense that the computed values will be as close to the actual Y values as is possible with a linear conversion formula. Larger absolute differences between the actual and estimated scores ($Y_j - \hat{Y}_j$) are indicative of larger errors. The average error $\Sigma(Y - \hat{Y})/N$ will equal zero whenever the overestimation of some scores is balanced by an equal underestimation of other scores. That there be no consistent over- or underestimation is a desirable property, but it may be accomplished by an infinite number of conversion rules. We therefore define as close as possible to correspond to the least squares criterion so common in statistical work—we shall choose a conversion rule such that not only are the errors balanced (they sum to zero), but also the sum of the squared discrepancies between the actual Y and estimated \hat{Y} will be minimized, that is, will be as small as the data permit.

It can be proven that the linear conversion rule which is optimal for converting z_X to an estimate of \hat{z}_Y is

$$(2.4.1) \quad \hat{z}_Y = r_{XY} z_X.$$

To convert to raw scores, we substitute for $\hat{z}_Y = (\hat{Y} - M_Y)/sd_Y$ and for $\hat{z}_X = (X - M_X)/sd_X$.

Solving for \hat{Y} gives

$$(2.4.2) \quad \hat{Y} = r_{XY} sd_Y \frac{(X - M_X)}{sd_X} + M_Y.$$

It is useful to simplify and separate the elements of this formula in the following way. Let

$$(2.4.3) \quad B_{YX} = r_{XY} \frac{sd_Y}{sd_X},$$

and

$$(2.4.4) \quad B_0 = M_Y - B_{YX} M_X,$$

34 2. BIVARIATE CORRELATION AND REGRESSION

from which we may write the regression equation for estimating Y from X as

$$(2.4.5) \quad \hat{Y} = B_{YX}X + B_0.$$

Alternatively, we may write this equation in terms of the original Y variable by including an “error” term e , representing the difference between the predicted and observed score for each observation:

$$(2.4.6) \quad Y = B_{YX}X + B_0 + e$$

These equations describe the regression of Y on X . B_{YX} is the regression coefficient for estimating Y from X and represents the rate of change in Y units per unit change in X , the constant by which you multiply each X observation to estimate Y . B_0 is called the regression constant or Y intercept and serves to make appropriate adjustments for differences in size between X and Y units. When the line representing the best linear estimation equation (the Y on X regression equation) is drawn on the scatterplot of the data in the original X and Y units, B_{YX} indicates the slope of the line and B_0 represents the point at which the regression line crosses the Y axis, which is the estimated \hat{Y} when $X = 0$. (Note that B_0 is sometimes represented as A or A_{YX} in publications or computer output.)

For some purposes it is convenient to *center* variables by subtracting the mean value from each score.⁹ Following such subtraction the mean value will equal 0. It can be seen by Eq. (2.4.4) that when both the dependent and independent variables have been centered so that both means = 0, the $B_0 = 0$. This manipulation also demonstrates that the predicted score on Y for observations at the mean of X must equal the mean of Y . When only the IV is centered, the B_0 will necessarily equal M_Y . For problems in which X does not have a meaningful zero point, centering X may simplify interpretation of the results (Wainer, 2000). The slope B_{YX} is unaffected by centering.

The slope of a regression line is the measure of its steepness, the ratio of how much Y rises (or, when negative, falls) to any given amount of increase along the horizontal X axis. Because the “rise over the run” is a constant for a straight line, our interpretation of it as the number of units of change in Y per unit change in X meets this definition.

Now we can deal with our example of 15 faculty members with a mean of 7.67 and a *sd* of 4.58 years since Ph.D. (Time) and a mean of 19.93 and a *sd* of 13.82 publications (Table 2.2.2). The correlation between time and publications was found to be .657, so

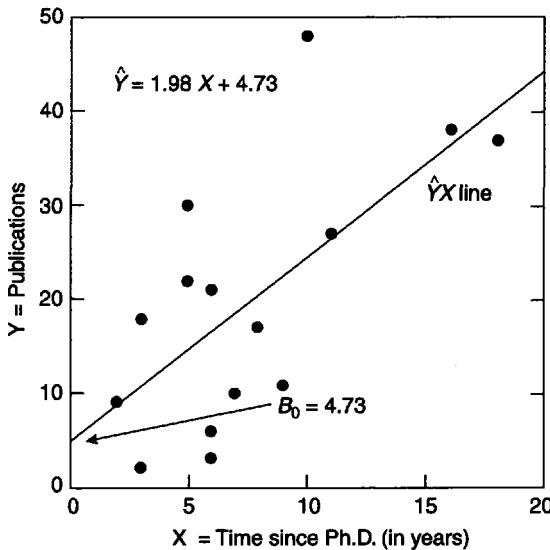
$$B_{YX} = .657(13.82/4.58) = 1.98,$$

$$B_0 = 19.93 - 1.98(7.67) = 4.73.$$

The regression coefficient, B_{YX} , indicates that for each unit of increase in Time (X), we estimate a change of +1.98 units (publications) in Y (i.e., about two publications per year), and that using this rule we will minimize our errors (in the least squares sense). The B_0 term gives us a point for starting this estimation—the point for a zero value of X , which is, of course, out of the range for the present set of scores. The equation $\hat{Y}_X = B_{YX}X + B_0$ may be used to determine the predicted value of Y for each value of X , and graphed as the $\hat{Y}X$ line in a scatterplot, as illustrated for these data in Fig. 2.4.1.

We could, of course, estimate X from Y by interchanging X and Y in Eqs. (2.4.3) and (2.2.2). However, the logic of regression analysis dictates that the variables are not of equal status, and estimating an independent or predictor variable from the dependent or criterion variable

⁹As will be seen in Chapters 6, 7, and 9, centering on X can greatly simplify interpretations of equations when relationships are curvilinear or interactive.



X	2	3	4	5	6	7	8	9	10
\hat{Y}	8.70	10.68	—	14.64	16.63	18.61	20.59	22.58	24.56
z_X	-1.24	-1.02	—	-0.58	-0.36	-0.15	0.07	0.29	0.51
Mz_Y	-0.84	-0.69	—	-0.40	-0.25	-0.10	0.05	0.20	0.35
X	11	12	13	14	15	16	17	18	
\hat{Y}	26.54	—	—	—	—	36.46	—	40.42	
z_X	0.73	—	—	—	—	1.88	—	2.34	
Mz_Y	0.50	—	—	—	—	1.24	—	1.53	

FIGURE 2.4.1 Regression of publications on time since Ph.D.

makes no sense. Suffice it to say that were we to do so, the line estimating X from Y (the X on Y regression) would not be the same as the line estimating Y from X (the Y on X regression). Neither its slope nor its intercept would be the same.

The meaning of the regression coefficient may be seen quite well in the case in which the independent variable is a dichotomy.¹⁰ If we return to the example from Table 2.3.1 where the point biserial $r = -.707$ and calculate

$$B_{YX} = -.707 \left(\frac{2.45}{.495} \right) = -3.5,$$

we note that this is exactly the difference between the two group means on Y , 66 – 69.5. Calculating the intercept, we get

$$B_0 = 68 - (-3.5)(.428) = 69.5,$$

which is equal to the mean of the group coded 0 (the no-stimulus condition). This must be the case because the best (least squares) estimate of Y for each group is its own mean, and the

¹⁰Chapter 8 is devoted to the topic of categorical IVs, for which we provide only a brief introduction here.

regression equation for the members of the group represented by the 0 point of the dichotomy is solved as

$$\hat{Y} = B_{YX}(0) + B_0 = B_0 = M_Y.$$

2.5 REGRESSION TOWARD THE MEAN

A certain amount of confusion exists in the literature regarding the phenomenon of regression toward the mean. It is sometimes implied that this is an artifact attributable to regression as an analytic procedure. On the contrary, it is a mathematical necessity that whenever two variables correlate less than perfectly, cases that are at one extreme on one of the variables will, on the average, be less extreme on the other. There are many examples in the literature where investigators mistakenly claim that some procedure results in a beneficial result when only the regression effect is operating (Campbell & Kenny, 1999). Consider a research project in which a neuroticism questionnaire is administered to an entering class and the students with the poorest scores are given psychotherapy, retested, and found to improve greatly. The “artifact” is the investigator’s claim of efficacy for the treatment when, unless the scores remained exactly the same so that the correlation between pretest and posttest was 1.0, they were certain to have scores closer to the mean than previously.

Although the number of cases in a small data set may be too small to show this phenomenon reliably at each data point, examination of the z_X and z_Y values in Fig. 2.4.1 will illustrate the point. The median of time since Ph.D. for the 15 professors is 6 years. If we take the 7 cases above the median, we find that their mean z score is $.82$, whereas the mean z score for the 5 professors below the median is $-.92$. Now, the mean z score for *number of publications* for the older professors is only $.52$ and the mean z score for publications for the younger professors is $-.28$. The cases high and low in years since Ph.D. (X) are distinctly less so on publications (Y); that is, they have “regressed” toward the mean. The degree of regression toward the mean in any given case will vary with the way we define *high* and *low*. That is, if we defined high time since Ph.D. as more than 12 years, we would expect an even greater difference between their mean z on time and the mean z on publications. The same principle will hold in the other direction: Those who are extreme on number of publications will be less extreme on years since Ph.D. As can be seen from these or any other bivariate data that are not perfectly linearly related, this is in no sense an artifact, but a necessary corollary of less than perfect correlation.

A further implication of this regression phenomenon is evident when one examines the consequences of selecting extreme cases for study. In the preceding paragraph, we found that those whose Ph.D.s were no more than 5 years old had a mean z score for years since Ph.D. of $-.92$, but a mean z score for number of publication of $-.28$. An investigator might well be tempted to attribute the fact that these new Ph.D.s are so much closer to the mean on number of publications than they are on years since Ph.D. to their motivation to catch up in the well-documented academic rat race. However, recognition that a less than perfect correlation is a necessary and sufficient condition to produce the observed regression toward the mean makes it clear that any specific substantive interpretation is not justified. (There is a delicious irony here: the lower the correlation, the greater the degree of regression toward the mean, and the more to “interpret,” spuriously, of course.)

Because regression toward the mean *always* occurs in the presence of an imperfect linear relationship, it is also observed when the variables consist of the same measure taken at two points in time. In this circumstance, unless the correlation is perfect, the extreme cases at Time 1 will be less extreme at Time 2. If the means and sds are stable, this inevitably means that low scores improve and high scores deteriorate. Thus, on the average over time, overweight people lose weight, low IQ children become brighter, and rich people become poorer. To ask why these

examples of regression to the mean occur is equivalent to asking why correlations between time points for weight, IQ, and income are not equal to +1.00. Of course, measurement error is one reason why a variable will show a lower correlation with itself over time, or with any other variables. However, regression to the mean is not solely dependent on measurement error, but on any mechanism whatsoever that makes the correlation less than perfect. Campbell and Kenny (1999) devote an entire volume to the many ways in which regression to the mean can lead to complexities in understanding change.

The necessity for regression toward the mean is not readily accessible to intuition but does respond to a simple demonstration. Expressed in standard scores, the regression equation is simply $\hat{z}_Y = r_{XY}z_X$ (Eq. 2.4.1). Because an r of +1 or -1 never occurs in practice, \hat{z}_Y will necessarily be absolutely smaller than z_X , because r is less than 1. Concretely, when $r = .40$, whatever the value of z_X , \hat{z}_Y must be .4 as large (see a comparable set of values below Fig. 2.4.1). Although for a single individual the actual value of z_Y may be larger or smaller than z_X , the expected or average value of the z_Y s that occur with z_X , that is, the value of \hat{z}_Y , will be .4 of the z_X value (i.e., it is “regressed toward the mean”). The equation holds not only for the expected value of z_Y for a single individual’s z_X , but also for the expected value of the mean z_Y for the mean z_X of a group of individuals. Of course, this holds true even when Y is the same variable measured at a later time than X . Unless the correlation over time is perfect, indicating no change, or the population mean and sd increase, *on the average*, the fat grow thinner, the dull brighter, the rich poorer, and vice versa.

2.6 THE STANDARD ERROR OF ESTIMATE AND MEASURES OF THE STRENGTH OF ASSOCIATION

In applying the regression equation $\hat{Y} = B_{YX}X + B_0$, we have of course only approximately matched the original Y values. How close is the correspondence between the information provided about Y by X (i.e., \hat{Y}), and the actual Y values? Or, to put it differently, to what extent is Y associated with X as opposed to being independent of X ? How much do the values of Y , as they vary, coincide with their paired X values, as they vary: equivalently, how big is e in Eq. (2.4.6)?

As we have noted, variability is indexed in statistical work by the sd or its square, the variance. Because variances are additive, whereas standard deviations are not, it will be more convenient to work with sd^2_Y . What we wish to do is to partition the variance of Y into a portion associated with X , which will be equal to the variance of the estimated scores, $sd^2_{\hat{Y}}$, and a remainder not associated with X , $sd^2_{Y-\hat{Y}}$, the variance of the discrepancies between the actual and the estimated Y scores (e). (Those readers familiar with ANOVA procedures may find themselves in a familiar framework here.) $sd^2_{\hat{Y}}$ and $sd^2_{Y-\hat{Y}}$ will sum to sd^2_Y , provided that \hat{Y} and $Y - \hat{Y}$ are uncorrelated. Intuitively it seems appropriate that they should be uncorrelated because \hat{Y} is computed from X by the optimal (OLS¹¹) rule. Because $\hat{Y} = B_{YX}X + (\text{a constant})$, it is just a linear transformation of X and thus necessarily correlates perfectly with X . Nonzero correlation between \hat{Y} and $Y - \hat{Y}$ would indicate correlation between X (which completely determines \hat{Y}) and $Y - \hat{Y}$, and would indicate that our original rule was not optimal. A simple algebraic proof confirms this intuition; therefore:

$$(2.6.1) \quad sd^2_Y = sd^2_{\hat{Y}} + sd^2_{Y-\hat{Y}} = sd^2_{\hat{Y}} + sd^2_e,$$

¹¹We introduce the term ordinary least squares (OLS) here, to represent the model that we have described, in which simple weights of predictor variable(s) are used to estimate Y values that collectively minimize the squared discrepancies of the predicted from the observed Y s, so that any other weights would result in larger average discrepancy.

and we have partitioned the variance of Y into a portion determined by X and a residual portion not linearly related to X . If no linear correlation exists between X and Y , the optimal rule has us ignore X because $B_{YX} = 0$, and minimize our errors of estimation by using M_Y as the best guess for every case. Thus we would be choosing that point about which the squared errors are a minimum and $sd_{Y-\hat{Y}}^2 = sd_Y^2$. More generally we may see that because (by Eq. 2.4.1) $\hat{z}_Y = r_{XY}z_X$,

$$sd_{z_Y}^2 = \frac{\sum (r_{XY}z_X)^2}{n-1} = r_{XY}^2 \frac{\sum z_X^2}{n-1} = r_{XY}^2,$$

and because $sd_{z_Y}^2 = 1$, and

$$(2.6.2) \quad sd_{z_Y}^2 = r_{XY}^2 + sd_{z_Y - \hat{z}_Y}^2,$$

then r_{XY}^2 is the proportion of the variance of Y linearly associated with X , and $1 - r_{XY}^2$ is the proportion of the variance of Y *not* linearly associated with X .

It is often helpful to visualize a relationship by representing each variable as a circle.¹² The area enclosed by the circle represents its variance, and because we have standardized each variable to a variance of 1, we will make the two circles of equal size (see Fig. 2.6.1). The degree of linear relationship between the two variables may be represented by the degree of overlap between the circles (the shaded area). Its proportion of either circle's area equals r^2 , and $1 - r^2$ equals the area of the nonoverlapping part of either circle. Again, it is useful to note the equality of the variance of the variables once they are standardized: the size of the overlapping and nonoverlapping areas, r^2 , and $1 - r^2$, respectively, must be the same for each. If one wishes to think in terms of the variance of the original X and Y , one may define the circles as representing 100% of the variance and the overlap as representing the proportion of each variable's variance associated with the other variable. We can also see that it does not matter in this form of expression whether the correlation is positive or negative because r^2 must be positive.

We will obtain the variance of the residual (nonpredicted) portion when we return to the original units by multiplying by sd_Y^2 to obtain

$$(2.6.3) \quad sd_{Y-\hat{Y}}^2 = sd_Y^2(1 - r^2).$$

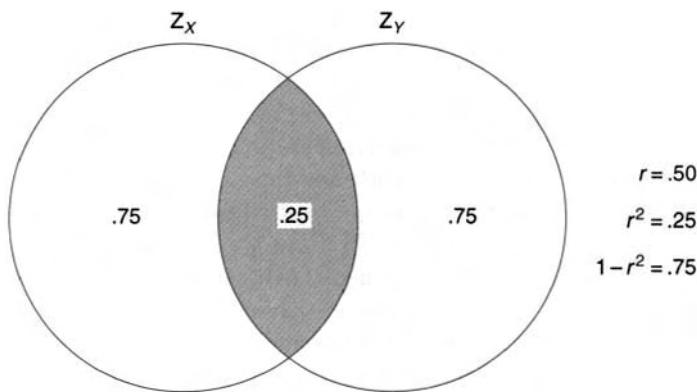


FIGURE 2.6.1 Overlap in variance of correlated variables.

¹²Such figures are called Venn diagrams in mathematical statistics. Here we call them “ballantines,” a name taken from a logo for a now-defunct beer company, because we use them illustratively only, and do not wish to imply the mathematical precision that should accompany a Venn diagram.

The standard deviation of the residuals e , that is, of that portion of Y not associated with X is therefore given by

$$(2.6.4) \quad sd_{Y-\hat{Y}} = sd_Y \sqrt{1 - r^2}.$$

For example, when $r = .50$, the proportion of shared variance $= r^2 = .25$, and .75 of sd_Y^2 is not linearly related to X . If the portion of Y linearly associated with X is removed by subtracting $B_{YX}X + B_0$ ($= \hat{Y}$) from Y , the sd of the residual is reduced compared to the original sd_Y to $sd_{Y-\hat{Y}} = sd_Y \sqrt{.75} = .866 sd_Y$.

We see that, in this case, although $r = .50$, only 25% of the variance in Y is associated with X , and when the part of Y which is linearly associated with X is removed, the standard deviation of what remains is .866 as large as the original SD_Y .

To make the foregoing more concrete, let us return to our academic example. The regression coefficient B_{YX} was found to be 1.98, the intercept B_0 was 4.73, and r_{XY} was .657. Table 2.6.1 gives the Y , X , and z_Y values and estimated \hat{Y} and \hat{z}_Y from the regression equations (2.4.5) and (2.4.1), which for these values are:

$$\begin{aligned}\hat{Y} &= 1.98 X_0 + 4.73 \quad \text{and} \\ \hat{z}_Y &= .657 z_X.\end{aligned}$$

The $Y - \hat{Y}$ values are the residuals for Y estimated from X or the errors of estimate in the sample. Because \hat{Y} is a linear transformation of X , $r_{Y\hat{Y}}$ must equal r_{XY} ($= .657$). The correlations between $Y - \hat{Y}$ and \hat{Y} must, as we have seen, equal zero. Parallel entries are given for the standardized \hat{z}_Y values where the same relationships hold.

Turning our attention to the variances of the variables, we see that

$$(2.6.5) \quad \frac{sd_{\hat{Y}}^2}{sd_Y^2} = \frac{sd_{\hat{z}_Y}^2}{1} = r^2 = .657^2 = .4312.$$

The ratio $sd_{Y-\hat{Y}}/sd_Y = \sqrt{1 - r^2} = .754$, which is called the coefficient of alienation, is the part of sd_Y that remains when that part of Y associated with X has been removed. It can also be thought of as the coefficient of noncorrelation, because r is the coefficient of correlation. The standard deviation of the residual scores is given by Eq. (2.6.4) as $sd_{Y-\hat{Y}} = sd_Y \sqrt{1 - r^2} = 13.35(.754) = 10.07$, as shown in Table 2.6.1. For the bivariate case, the population *variance error of estimate* or *residual variance* has $df = n - 2$ and is given by

$$(2.6.6) \quad SE_{Y-\hat{Y}}^2 = \frac{\sum (Y - \hat{Y})^2}{n - 2} = \frac{(1 - r_{XY}^2) \sum (Y - M_Y)^2}{n - 2}.$$

For the two summations, Table 2.6.1 gives in its $\Sigma \sqrt{x^2}$ row, 1521.51 for the $Y - \hat{Y}$ column and 2674.93 for the Y column. Substituting, we get

$$SE_{Y-\hat{Y}}^2 = \frac{1521.51}{15 - 2} = \frac{(1 - .657^2)2674.93}{15 - 2},$$

and both equations give 117.04. When we take square roots, we obtain the *standard error of estimate*:

$$(2.6.7) \quad SE_{Y-\hat{Y}} = \sqrt{\frac{\sum (Y - \hat{Y})^2}{n - 2}} = \sqrt{\frac{(1 - r_{XY}^2) \sum (Y - M_Y)^2}{n - 2}},$$

which equals 10.82. Here, too, $df = n - 2$.



TABLE 2.6.1
Estimated and Residual Scores for Academic Example

X Time since Ph.D.	Y No. of publications	\hat{Y}	$Y - \hat{Y}$	\hat{z}_Y	$z_Y - \hat{z}_Y$	\hat{Y}_W	$Y - \hat{Y}_W$	\hat{Y}_V	$Y - \hat{Y}_V = e$
3	18	10.68	7.32	-.67	.53	10.60	7.40	11.07	6.93
6	3	16.63	-13.63	-.24	-.99	16.60	-13.60	16.77	-13.77
3	2	10.68	-8.68	-.67	-.63	10.60	-8.60	11.07	-9.07
8	17	20.59	-3.59	.05	-.26	20.60	-3.60	20.57	-3.57
9	11	22.58	-11.58	.19	-.84	22.60	-11.60	22.47	-11.47
6	6	16.63	-10.63	-.24	-.77	16.60	-10.60	16.77	-10.77
16	38	36.46	1.54	1.20	.11	36.60	1.40	35.77	2.23
10	48	24.56	23.44	.33	1.70	24.60	23.40	24.37	23.63
2	9	8.70	0.30	-.81	.02	8.60	.40	9.17	-.17
5	22	14.65	7.36	-.38	.53	14.60	7.40	14.87	7.13
5	30	14.65	15.36	-.38	1.11	14.60	15.40	14.87	15.13
6	21	16.63	4.37	-.24	.32	16.60	4.40	16.77	4.23
7	10	18.61	-8.61	-.10	-.62	18.60	8.60	18.67	8.67
11	27	26.54	0.46	.48	.03	26.60	.40	26.27	73
18	37	40.42	-3.42	1.48	-.25	40.60	3.60	39.57	2.57
M	7.67	19.93	19.93	0	0	19.93	0	19.93	0
sd	4.577	13.82	8.77	10.07	.657	.754		10.072	8.40
sd^2	19.56	178.3	76.98	101.42	.431	.569		101.44	70.60
$\Sigma x_i $			120.29				116.40		120.07
$\Sigma\sqrt{x_i^2}$	2674.93		1521.51						

$$r_{Xz_X} = r_{Yz_Y} = r_{XY} = r_{z_X z_Y} = r_{\hat{Y}X} = 1.$$

$$r_{XY} = r_{z_X z_Y} = r_{Y\hat{Y}} = .657$$

$$r_{Y(Y-\hat{Y})}^2 = .5689; r_{(Y-\hat{Y})\hat{Y}} = r_{(Y-\hat{Y})X} = 0.$$

Finally, \hat{Y}_W and \hat{Y}_V in Table 2.6.1 have been computed to demonstrate what happens when any other regression coefficient or weight is used. The values $B_{WX} = 2.0$ and $B_{VX} = 1.9$ were chosen to contrast with $B_{YX} = 1.98$ (the regression constants have been adjusted to keep the estimated values centered on Y). The resulting sd^2 for the sample residuals was larger in each case, 101.44 and 101.57, respectively as compared to 101.42 for the least squares estimate. The reader is invited to try any other value to determine that the squared residuals will in fact always be larger than with 1.98, the computed value of B_{YX} .

Examination of the residuals will reveal another interesting phenomenon. If one determines the *absolute* values of the residuals from the true regression estimates and from the \hat{Y}_W , it can be seen that their sum is smaller for both $Y - \hat{Y}_W$ (116.40) and $Y - \hat{Y}_V$ (120.07) than it is for the true regression residuals (120.29). Whenever residuals are not exactly symmetrically distributed about the regression line there exists an absolute residual minimizing weight different from B_{YX} . To reiterate, B_{YX} is the weight that minimizes the squared residuals, not their absolute value. This is a useful reminder that ordinary least squares (OLS), although very useful, is only one way of defining discrepancies from estimation, or error.¹³

¹³Chapter 4 will introduce alternative methods, which are further presented in later chapters.

2.7 SUMMARY OF DEFINITIONS AND INTERPRETATIONS

The product moment r_{XY} is the rate of linear increase in z_Y per unit increase or decrease in z_X (and vice versa) that best fits the data in the sense of minimizing the sum of the squared differences between the estimated and observed scores.

r^2 is the proportion of variance in Y associated with X (and vice versa).

B_{YX} is the regression coefficient of Y on X . Using the original raw units, it is the rate of linear change in Y per unit change in X , again best fitting in the least squares sense.

B_0 is the regression intercept that serves to adjust for differences in means, giving the predicted value of the dependent variable when the independent variable's value is zero.

The coefficient of alienation, $\sqrt{1 - r^2}$, is the proportion of sd_Y remaining when that part of Y associated with X has been subtracted; that is, $sd_{Y-\hat{Y}}/sd_Y$.

The standard error of estimate, $SE_{Y-\hat{Y}}$, is the estimated population standard deviation (σ) of the residuals or errors of estimating Y from X .

2.8 STATISTICAL INFERENCE WITH REGRESSION AND CORRELATION COEFFICIENTS

In most circumstances in which regression and correlation coefficients are determined, the intention of the investigator is to provide valid inferences from the sample data at hand to some larger universe of potential data—from the statistics obtained for a sample to the parameters of the population from which it is drawn. Because random samples from a population cannot be expected to yield sample values that exactly equal the population values, statistical methods have been developed to determine the confidence with which such inferences can be drawn. There are two major methods of statistical inference, estimation using confidence intervals and null hypothesis significance testing. In Section 2.8.1, we consider the formal model assumptions involved. In Section 2.8.2, we describe confidence intervals for B_{YX} , B_0 , r_{XY} , for differences between independent sample values of these statistics. In Section 2.8.3, we present the null hypothesis tests for simple regression and correlation statistics. Section 2.8.4 critiques null hypothesis testing and contrasts it with the approach of confidence limits.

2.8.1 Assumptions Underlying Statistical Inference with B_{YX} , B_0 , \hat{Y}_i , and r_{XY}

It is clear that no assumptions are necessary for the computation of correlation, regression, and other associated coefficients or their interpretation when they are used to describe the available sample data. However, the most useful applications occur when they are statistics calculated on a sample from some population in which we are interested. As in most circumstances in which statistics are used inferentially, the addition of certain assumptions about the characteristics of the population substantially increases the useful inferences that can be drawn. Fortunately, these statistics are *robust*; that is, moderate departure from these assumptions will usually result in little error of inference.

Probably the most generally useful set of assumptions are those that form what has been called the *fixed linear regression model*. This model assumes that the two variables have been distinguished as an independent variable X and a dependent variable Y . Values of X are treated as “fixed” in the analysis of variance sense, that is, as selected by the investigator rather than

sampled from some population of X values.¹⁴ Values of Y are assumed to be randomly sampled for each of the selected values of X . The residuals (“errors”) from the mean value of Y for each value of X are assumed to be normally distributed in the population, with equal variances across the full range of X values. It should be noted that no assumptions about the shape of the distribution of X and the total distribution of Y per se are necessary, and that, of course, the assumptions are made about the population and not about the sample. This model, extended to multiple regression, is used throughout the book.

2.8.2 Estimation With Confidence Intervals

A *sampling* distribution is a distribution of the values of a sample *statistic* that would occur in repeated random sampling of a given size, n , drawn from what is conceived as an infinite population. Statistical theory makes possible the estimation of the shape and variability of such sampling distributions. We estimate the population value (*parameter*) of the sample statistic we obtained by placing it within a *confidence interval* (*CI*) to provide an estimate of the margin of error (*me*), based on these distributions.

Confidence Interval for B_{YX}

We have seen that B_{YX} is a regression coefficient that gives the slope of the straight line that estimates Y from X . We will see that, depending on the context, it can take on many meanings in data analysis in MRC, including the size of a difference between two means (Section 2.4), the degree of curvature of a regression line (Chapter 6), or the effect of a datum being missing (Chapter 11).

Continuing our academic example, we found in Section 2.4 that for this sample the least squares estimate of $B_{YX} = 1.98$, indicating that for each additional year since Ph.D. we estimate an increase of 1.98 publications, that is, an increase of about two publications. If we were to draw many random samples of that size from the population, we would get *many* values of B_{YX} in the vicinity of +1.98. These values constitute the *sampling distribution* of B_{YX} and would be approximately normally distributed. The size of the vicinity is indicated by the standard deviation of this distribution, which is the *standard error* (SE) of B_{YX} :

$$(2.8.1) \quad SE_{B_{YX}} = \frac{sd_Y}{sd_X} \sqrt{\frac{1 - r_{YX}^2}{n - 2}}$$

Substituting,

$$SE_{B_{YX}} = \frac{13.82}{4.58} \sqrt{\frac{1 - .657^2}{15 - 2}} = .632.$$

Because this is a very small sample, we will need to use the t distribution to determine the multiplier of this *SE* that will yield estimates of the width of this interval. Like the normal distribution, the t distribution is a symmetrical distribution but with a relatively higher peak in the middle and higher tails. The t model is a family of distributions, each for a different number of *degrees of freedom* (df). As the df increase from 1 toward infinity, the t distribution becomes progressively less peaked and approaches the shape of the normal distribution. Looking in

¹⁴In the “multilevel” models discussed in Chapters 14 and 15 this assumption is not made for all independent variables.

Appendix Table A, we find that the necessary t at the two-tailed 5% level for 13 df is 2.16. Multiplying .632 by 2.16 gives 1.36, the 95% *margin of error (me)*. Then, the 95% *confidence limits (CLs)* are given as $1.98 \pm 1.36 = +.62$ as its lower limit and $+3.34$ as its upper limit. If 1.98 is so much smaller than the population value of B_{YX} that only 2.5% of the possible sample B_{YX} values are smaller still, then the population value is 1.36 publications *above* 1.98, that is, 3.34 (see Fig. 2.8.1), and if 1.98 is so much larger that only 2.5% of the possible sample B_{YX} values are larger still, then the population value is 1.36 publications *below* 1.98, that is, .62 (see Fig. 2.8.2). Thus, the 95% CI is $+.62$ to $+3.34$. This CI indicates our 95% certainty that the population value falls between $.62$ and $+3.34$. Note for future reference the fact that the CI for B_{XY} in this sample does *not* include 0 (see Section 2.8.3).

Although the *single most likely* value for the change in number of publications per year since Ph.D. is the sample value 1.98, or about 2 publications per year, we are 95% confident that the true change falls between .62 and 3.34 publications per year since Ph.D. This may be too large an interval to be of much use, as we should have expected when we examined so

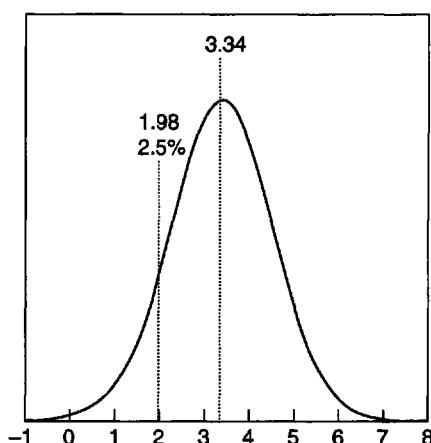


FIGURE 2.8.1 Expected distribution of B s from samples of 15 subjects when the population $B = 3.34$.

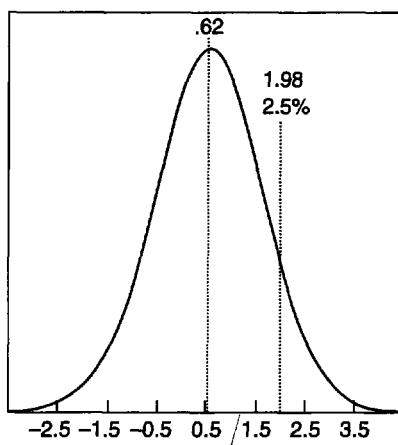


FIGURE 2.8.2 Expected distribution of B s from samples of 15 subjects when the population $B = 0.62$.

small a sample. Were we to have found the same sample value of 1.98 on a sample as large as 62, the standard error of B_{YX} would go down to .294 (Eq. 2.8.1). When $n = 62$, the df for $SE_{B_{YX}}$ is $n - 2 = 60$, so t for the 95% CI = 2.00 (Appendix Table A). The me (margin of error) is now $2.00(.294) = .588$, less than half as large as before, so the 95% CI is now $1.98 \pm (.588) = 1.40$ to 2.56, from about 1.5 to 2.5 publications per year since Ph.D., distinctly narrower and thus more useful.

Although 95% CIs are the most frequently used, other degrees of confidence, greater or smaller, may be preferred. A multiplier of 2.6 will give an approximate 99% CI, and 1.3 an approximate 80% interval for all but the smallest samples. Since standard errors are always reported in computer output, and should always be reported in research reports, one can easily approximate a CI that includes 68% (about $\frac{2}{3}$) of the cases in the sampling distribution by taking the me for the sample B_{YX} value to equal its SE , so the approximate 68% CI is $B_{YX} \pm SE_{B_{YX}}$. The odds are then approximately 2 to 1 that the population B_{YX} value falls between those limits.

Confidence Interval for B_0

B_0 is the regression coefficient that gives the Y *intercept*, the value of \hat{Y} when the $\hat{Y}X$ regression line that estimates Y from X is at $X = 0$. Although in many behavioral science applications this coefficient is ignored, because the means of the variables are essentially on an arbitrary scale, there are applications in which it is of interest. When zero on the X scale has a useful meaning, and is within the range of the observations, it tells us what the expected value of Y is for $X = 0$. In Section 2.4, we found using Eq. (2.4.4) that for our running example the intercept $B_0 = M_Y - B_{YX}M_X = 19.93 - 1.98(7.67) = 4.73$, indicating a predicted value of 4.73 publications when years since Ph.D. equals 0, that is, the individual has just obtained a Ph.D. Of course, such a predicted value is not to be trusted under the circumstances in which it falls outside the observed data, as it does here.

The standard error of B_0 is given by

$$(2.8.2) \quad SE_{B_0} = SE_{Y-\hat{Y}} \sqrt{\frac{1}{n} + \frac{M_X^2}{(n-1)sd_X^2}}.$$

We found from Eq. (2.6.7) that for this example, $SE_{Y-\hat{Y}} = 10.82$. Substituting from Table 2.6.1 for $n = 15$, $M_X = 7.67$, and $sd^2 = 4.58^2 = 20.95$.

$$SE_{B_0} = 10.82 \sqrt{\frac{1}{15} + \frac{7.67^2}{(14)(20.95)}} = 5.59.$$

We generate CIs for B_0 as before, using the t distribution for $n - 2 = 13$ df. For the 95% CI, Appendix Table A gives $t = 2.16$, so the $me = 2.16(5.59) = 12.07$ and the 95% CI = $4.73 \pm 12.07 = -7.34$ to 16.80. The table gives for 13 df, $t = 1.35$ for the 80% CI, so $me = 1.35(5.59) = 7.55$, and the 80% CI = $4.73 \pm 7.55 = -2.82$ to 12.28. These large CIs, with their negative lower limits, mean that with such a small sample we cannot even confidently say whether, on the average, faculty members had published before they got their degrees!

Confidence Interval for an Estimated \hat{Y}_i Value

When we employ the regression equation

$$(2.4.5) \quad \hat{Y} = B_{YX}X + B_0$$

to estimate a particular \hat{Y}_i from a particular value of X_i , what we find is the Y coordinate of the point on the $\hat{Y}X$ regression line for that value of X . In the sample data, the Y values are scattered

above and below the regression line and their distances from the line are the *residuals* or *errors*. The *standard error of estimate* (Eq. 2.6.7) estimates their variability in the population. In our running example estimating number of publications from number of years since Ph.D., we found $SE_{Y-\hat{Y}}$ to equal 10.82. Let's write the regression equation to estimate \hat{Y}_i , the number of publications estimated for a specific faculty member with 9 years since Ph.D. The equation for these values was found as $\hat{Y}_i = 1.98X + 4.73$. Substituting $X_i = 9$, we find $\hat{Y}_i = 22.58$.

It is useful to realize that, whatever sampling error was made by using the sample B_{YX} ($= 1.98$) instead of the (unavailable) population regression coefficient, it will have more serious consequences for X values that are more distant from the X mean than for those near it. For the sake of simplicity, let us assume that both X and Y are z scores with means of 0 and standard deviations of 1. Suppose that $B_{YX} = .20$ for our sample, whereas the actual population value is .25. For new cases that come to our attention with $X_i = .1$, we will estimate \hat{Y}_i at .02 when the actual mean value of Y for all $X_i = .1$ is .025, a relatively small error of .005. On the other hand, new values of $X_i = 1.0$ will yield estimated \hat{Y}_i values of .20 when the actual mean value of Y for all $X_i = 1$ is .25, the error (.05) being 10 times as large.

When a newly observed X_i is to be used to estimate \hat{Y}_i we may determine the standard error and thus confidence limits for this \hat{Y}_i . The standard error of \hat{Y}_i is given by

$$(2.8.3) \quad SE_{\hat{Y}_i} = SE_{Y-Y_i} \sqrt{\frac{1}{n} + \frac{(X_i - M_X)^2}{(n-1)sd_X^2}},$$

where SE_{Y-Y_i} (Eq. 2.6.7) is the standard error of estimate and is based on $n - 2df$. We found from the regression equation that for $X_i = 9$ years since Ph.D., we estimate $\hat{Y}_i = 22.58$ publications. We find its standard error by substituting in Eq. (2.8.3),

$$SE_{\hat{Y}_i} = 10.82 \sqrt{\frac{1}{15} + \frac{(9 - 7.67)^2}{(14)(20.95)}} = 2.92$$

For the 95% CI, Appendix Table A gives $t = 2.16$ for 13 df, so the *me* is $2.16(2.92) = 6.30$ and the 95% CI = $22.58 \pm 6.30 = 16.3$ to 28.9 publications (rounding). For the 80% CI, the table gives $t = 1.35$ for 13 df, so the *me* = $1.35(2.92) = 3.94$ and the CI is $22.58 \pm 3.94 = 18.6$ to 26.5 publications (rounding). These CIs are uselessly large because of the large SE_{Y_0} , due mostly in turn to the smallness of the sample.

Confidence Interval for r_{XY}

The approach we used in generating CIs for B_{YX} and B_0 will not work for r_{XY} because the sampling distribution for r_{XY} is not symmetrical except when ρ_{YX} (the population r_{XY}) equals 0. That is, the lower and upper limits for a CI for r_{XY} do not fall at equal distances from the obtained sample value. The reason for this is that, unlike $SE_{B_{YX}}$, the SE_r varies with ρ_{YX} , which is, of course, unknown. To solve this problem, R. A. Fisher developed the z prime (z') transformation of r :

$$(2.8.4) \quad z' = \frac{1}{2}[\ln(1+r) - \ln(1-r)],$$

where \ln is the natural (base e) logarithm.

The sampling distribution of z' depends only on the sample size and is nearly normal even for relatively small values of n . The standard error of a sample z' is given by

$$(2.8.5) \quad SE_{z'} = \frac{1}{\sqrt{n-3}}$$

Appendix Table B gives the r to z' transformation directly, with no need for computation.

To find the *CI* for a sample r , transform the r to z' and, using the $SE_{z'}$ and the appropriate multiplier for the size of the *CI* desired, find the *me* and then the lower and upper limits of the *CI* for z' . Then transform them back to r . For our academic example, we found the r between years since Ph.D. and number of publications to be .657. In Appendix Table B we find the z' transformation to be approximately $z' = .79$. With $n = 15$, we find from (2.8.4) that

$$SE_{z'} = \frac{1}{\sqrt{15 - 3}} = .289.$$

Then, using the multiplier 1.96 from the *normal distribution* for the 95% limits (Appendix Table C), we find $1.96(.289) = .57$ as the *me* for z' , so $.79 \pm .57$ gives the 95% limits for z' as .22 and 1.36. But what we want are the 95% limits for r , so using Appendix Table B we transform these z' values back to r and obtain $r = .22$ (from .22) and .88 (from 1.36). Thus, we can expect with 95% confidence that the population r is included in the approximate *CI* .22 to .88. Note that these limits are not symmetrical about the sample r of .657.

The 95% *CI* for r in this example, .22 to .88, is very wide, as are all the *CIs* for this small sample of $n = 15$.¹⁵ The odds of inclusion here are 95 : 5 (that is, 19 to 1). For narrower and thus less definitive limits, the 80% *CI* gives 80 : 20 (4 to 1) odds of inclusion. To find it, we proceed as before, using the normal curve multiplier for an 80% *CI* of 1.28 (Appendix Table C). We first find the confidence limits for z' by subtracting and adding the $me = 1.28 (.29) = .38$ to the sample z' of .79, obtaining .41 and 1.17. From Appendix Table B we convert z' to r to find the approximate 80% *CI* for r to be .39 (from .41) to .82 (from 1.17). This is yet another object lesson in precision (or, rather, its lack) with small samples. For most purposes, limits as wide as this would not be of much use.

Confidence Interval for the Difference Between Regression Coefficients: $B_{XY_V} - B_{XY_W}$

Given the many uses to which regression coefficients are put, the size of the difference between a pair of B_{YX} sample values coming from different groups is often a matter of research interest. The *SE* of the difference between two independent B_{YX} values is a function of their standard errors, whose formula we repeat here for convenience:

$$(2.8.1) \quad SE_{B_{YX}} = \frac{sd_Y}{sd_X} \sqrt{\frac{1 - r_{YX}^2}{n - 2}}.$$

Assume that the sample in Section 2.4 in which we found the regression coefficient describing the relationship between time since Ph.D. and number of publications, 1.98, was drawn from University V and numbered 62 cases. Substituting the sample values found in Section 2.4 in Eq. (2.8.1), we find its standard error to be .294. Now assume that in a random sample of 143 cases from University W, we find $sd_{Y_w} = 13.64$, $sd_{X_w} = 3.45$, and $r_w = .430$. Substituting these values in Eq. (2.4.3), we find $B_w = 1.70$, and in Eq. (2.8.1) we find $SE_{B_w} = .301$. Now, the difference between B_v and B_w is $1.98 - 1.70 = .28$. The standard error of the difference between the two coefficients is

$$(2.8.6) \quad SE_{B_v - B_w} = \sqrt{(SE_{B_v})^2 + (SE_{B_w})^2}$$

Substituting, we find

$$SE_{B_v - B_w} = \sqrt{(.294)^2 + (.301)^2} = .42$$

¹⁵Indeed, it would be foolish to place any serious faith in the adequacy of the estimate based on such a small sample, which is employed here only for illustrative purposes.

Using the multiplier 2 (a reasonable approximation of 1.96) for the 95% *CI*, we find the *me* for the difference between the *B* values, $2(.42) = .84$, and obtain the approximate 95% *CI* for $B_V - B_W$ as $.28 \pm .84 = -.56$ to $+1.12$. This means that the confidence limits go from University V's slope being .56 (about $\frac{1}{2}$ of a publication) *smaller* per year since Ph.D. to being 1.12 (about 1) publication larger. Take particular note of the fact that the 95% *CI* includes 0. Thus, we cannot conclude that there is *any* difference between the universities in the number of publications change per year since Ph.D. at this level of confidence.

Equation (2.8.6) gives the standard error of the difference between regression coefficients coming from different populations as the square root of the sum of their squared standard errors. This property is not unique to regression coefficients but holds for *any* statistic—means, standard deviations, and, as we see in the next section, correlation coefficients as well.

Confidence Interval for $r_{XY_V} - r_{XY_W}$

We cannot approach setting confidence limits for differences between *rs* using the *z'* transformation because of the nonlinear relationship between them—equal distances along the *r* scale do not yield equal distances along the *z'* scale (which can be seen in Appendix Table B).

Recent work by Olkin and Finn (1995) has provided relatively simple means for setting confidence intervals for various functions of correlation coefficients. For *large* samples, the difference between r_{YX} in two independent samples, V and W, is normally distributed and is given approximately by

$$(2.8.7) \quad SE_{r_V - r_W} = \sqrt{\frac{1 - r_V^2}{n_V} + \frac{1 - r_W^2}{n_W}}.$$

Returning to the example in which we compared the regression coefficients for our running problem, we can estimate confidence intervals for the difference between the correlations of .657 for University V ($n_V = 62$) and .430 for University W ($n_W = 143$). Substituting in Eq. (2.8.7),

$$SE_{r_V - r_W} = \sqrt{\frac{1 - .657^2}{62} + \frac{1 - .430^2}{143}} = .122$$

The difference between the *rs* is $.657 - .430 = .277$. Assuming normality, the 95% *CI* uses 1.96 as the multiplier, so the 95% *me* is $1.96(.122) = .239$. Then the approximate 95% *CI* is $.277 \pm .239 = +.04$ to $.52$. We interpret this to mean that we can be 95% confident that the ρ_{YX} of time since Ph.D with number of publications for University V is .04 to .52 *larger* than that for University W. Note here that the confidence interval of the difference between the *rs* of the two universities does not include 0, but the *CI* of the difference between their regression coefficients does. This demonstrates that correlation and regression coefficients are different measures of the degree of linear relationship between two variables. Later, we will argue that regression coefficients are often more stable across populations, in contrast to *rs* that reflect population differences in variability of *X*. In the preceding example, we saw $sd_X = 4.58$ in the original University V and $sd_X = 3.45$ in the comparison University W. The smaller *r* in University W is apparently attributable to their faculty's constricted range of years since Ph.D.

2.8.3 Null Hypothesis Significance Tests (NHSTs)

In its most general meaning, a null hypothesis (H_0) is a hypothesis that a population effect size (ES) or other parameter has some value specified by the investigator. The term “null” arises from R. A. Fisher’s statistical strategy of formulating a proposition that the research data may

be able to *nullify* or reject. By far, the most popular null hypothesis that is tested is the one that posits that a population effect size, such as a correlation coefficient or a difference between means, is *zero*, and the adjective “*null*” takes on the additional meaning of no relationship or no effect. We prefer to use the term “*nil*” hypothesis to characterize such propositions for reasons that will become clear later (J. Cohen, 1994).

The Nil Hypothesis Test for B_{YX}

In our running example of the 15 faculty members, we found that the regression coefficient for the number of publications on number of years since Ph.D. was 1.98 ($= B_{YX}$), which means that, on the average in this sample, each additional year since Ph.D. was associated with about two publications. The standard error of the coefficient ($SE_{B_{YX}}$) from Eq. (2.8.1) was .632. Let’s perform a *t* test of the nil hypothesis that in the population, each additional year since Ph.D. is associated on the average with *no* additional publications, that is, that there is no linear relationship between years since Ph.D. and publications. We will perform this test at the $p < .05$ ($= \alpha$) significance level. The general form of the *t* test is

$$(2.8.8) \quad t = \frac{\text{sample value} - \text{null-hypothetical value}}{\text{standard error}}$$

which, for regression coefficients, is

$$(2.8.9) \quad t = \frac{B_{YX} - H_0}{SE_{B_{YX}}}.$$

Substituting,

$$t = \frac{1.98 - 0}{.632} = 3.14,$$

which, for $df = n - 2 = 13$ readily meets the $\alpha = .05$ significance criterion of $t = 2.16$ (Appendix Table A). We accordingly reject H_0 and conclude that there is a greater than zero relationship between years since Ph.D. and number of publications in the population. Note, however, that neither the size nor the statistical significance of the *t* value provides information about the *magnitude* of the relationship. Recall, however, that when we first encountered the $SE_{B_{YX}}$ at the beginning of Section 2.8.2, we found the 95% *CI* for B_{YX} to be +.62 to +3.34, which *does* provide a magnitude estimate. Moreover, note that the 95% *CI* *does not include 0*. After we have determined a *CI* for B_{YX} , a *t* test of the nil hypothesis for B_{YX} is unnecessary—once we have a *CI* that does not include 0, we know that the nil hypothesis can be rejected at that significance level (here, $\alpha = .05$). However, if the only relevant information about a population difference is whether it has some specified value, or whether it exists at all, and there are circumstances when that is the case, then *CIs* are unnecessary and a null hypothesis test is in order.

For example, assume that we wish to test the proposition as a non-nil null hypothesis that the population regression coefficient is 2.5 publications per year since Ph.D.: H_0 : population $B_{YX} = 2.5$. We can proceed as before with Eq. (2.8.9) to find $t = (1.98 - 2.5)/.632 = .82$, which is not significant at $\alpha = .05$, and we can conclude that our results are consistent with the possibility that the population value is 2.5. But since the 95% *CI* (+.62 to +3.34) contains the null-hypothetical value of 2.5, we can draw the same conclusion. However, by obtaining the 95% *CI* we have the *range* of B_{YX} values for which the H_0 cannot be rejected at $\alpha = .05$. Not only 2.5 or 0, but *any value* in that range cannot be rejected as a H_0 . Therefore, one may think of a *CI* as a range of values within which the H_0 *cannot* be rejected and outside of which H_0 *can* be rejected on the basis of this estimate. The *CI* yields more information than the NHST.

The Null Hypothesis Test for B_0

In the previous section, we found the Y intercept for our running example $B_0 = 4.73$ and, using its standard error (Eq. 2.8.2), found $SE_{B_0} = 5.59$. We can perform a t test for 13 df of the H_0 that the population intercept equals 0 in the usual fashion. Using Eq. (2.8.7) for B_0 and substituting in Eq. (2.8.7), we find

$$t = \frac{4.73 - 0}{5.59} = .85,$$

which fails to meet conventional significance criteria. (In Section 2.8.2 we found 95% and 80% CIs, both of which included 0.)

The Null Hypothesis Test for r_{XY}

When ρ_{XY} (the population r_{XY}) = 0, the use of the Fisher z' transformation is unnecessary. The t test of the *nil* hypothesis for r_{XY} , $H_0: \rho_{XY} = 0$, is

$$(2.8.10) \quad t = \frac{r_{XY}\sqrt{n-2}}{\sqrt{1-r_{XY}^2}} \quad \text{with } df = n-2.$$

Returning to our running example, the r_{XY} between years since Ph.D. and publications for the sample of 15 faculty members was .657. Substituting,

$$t = \frac{.657\sqrt{15-2}}{\sqrt{1-.657^2}} = 3.14.$$

The $\alpha = .05$ significance criterion for t with 13 df is 2.16, readily exceeded by 3.14. We conclude that $\rho_{XY} \neq 0$. (The 95% CI was found via the Fisher z' transformation in the previous section to be .22 to .88.)

The Null Hypothesis Test for the Difference Between Two Correlations with Y : $r_{XY_V} - r_{XY_W}$

In Section 2.8.2 we presented a method for setting approximate confidence intervals for differences between independent rs suitable for large samples. For an approximate *nil* hypothesis test, suitable for samples of any size, we again resort to the Fisher z' transformation. The relevant data for the two universities are

University	N	r_{XY}	z'_{XY}
V	62	.657	.79
W	143	.430	.46

To test the H_0 that the difference between the population correlations: $\rho_V - \rho_W = 0$, we test the equivalent $H_0: z'_V - z'_W = 0$ by computing the normal curve deviate

$$(2.8.11) \quad z = \frac{z'_V - z'_W}{\sqrt{1/(n_V - 3) + 1/(n_W - 3)}}.$$

Substituting,

$$z = \frac{.79 - .46}{\sqrt{1/(62 - 3) + 1/(143 - 3)}} = 2.13,$$

which exceeds 1.96, the two-tailed $\alpha = .05$ significance criterion for the normal distribution (see Appendix Table C), and we can conclude that University V's ρ_{XY} is probably larger than University W's. The reason that we can test for z 's and conclude about ρ s is that there is a one-to-one correspondence between z' and ρ so that when the z 's are not equal, the ρ s are necessarily also not equal. (The 95% CI for the difference between the r s was previously found to be +.04 to +.52.)

2.8.4 Confidence Limits and Null Hypothesis Significance Testing

For more than half a century, NHST has dominated statistical inference in its application in the social, biological, and medical sciences, and for just as long, it has been subject to severe criticism by methodologists including Berkson (1946), Yates (1951), Rozeboom (1960), Meehl (1967), Lykken (1968), and Tukey (1969), among others. More recently, many methodologists, including J. Cohen (1990, 1994) and a committee of the American Psychological Association (Wilkinson of the APA Task Force on Statistical Inference, 1999), among others, have inveighed against the excessive use and abuse of NHST.

We have seen repeatedly that when confidence intervals on statistics or effect sizes are available, they include the information provided by null hypothesis tests. However, there may be a useful role for NHST in cases where the direction of systematic differences is of much more interest than their magnitude and the information provided by confidence intervals may simply be distracting (Harlow, Mulaik, & Steiger, 1997). In addition, as we will see in subsequent chapters, significance tests are useful guides to the decision as to whether certain variables are or are not needed for the explanation of Y . Abelson (1995) notes the usefulness of NHST in making categorical claims that add to the background substantive scientific lore in a field under study.

2.9 PRECISION AND POWER

For research results to be useful, they must be accurate or, at least, their degree of accuracy must be determinable. In the preceding material, we have seen how to estimate regression parameters and test null hypothesis after the sample data have been collected. However, we can plan to determine the degree of precision of the estimation of parameters or of the probability of null hypothesis rejection that we shall be able to achieve.

2.9.1 Precision of Estimation

The *point estimate* of a population parameter such as a population B or ρ is the value of the statistic (B , r) in the sample. The margin of error in estimation is the product of the standard error and its multiplier for the degree of inclusion (95%, 80%) of the confidence interval. The standard error is a function of the sample size, n . We show how to estimate n^* , the sample size necessary to achieve the desired degree of precision of the statistics covered in Section 2.8.2.

We begin by drawing a trial sample of the data for whose statistics we wish to determine CIs. The sample of $n = 15$ cases we have been working with is much too small to use as a trial sample, so let's assume that it had 50 rather than 15 cases so that we can use the same statistics as before: $M_X = 7.67$, $sd_X = 4.58$, $sd_Y = 13.82$, $r_{XY} = .657$, $B_{YX} = 1.98$, $B_0 = 4.73$, and $SE_{Y-\hat{Y}} = 10.82$.

We use the approximate multipliers (t , z) of the standard errors to determine the inclusion of the confidence limits: 99%, 2.6; 95%, 2; 80%, 1.3; and 68%, 1. The standard errors for the regression/correlation statistics of our $n = 50$ sample are as follows:

Estimated B_{YX}

$$\text{Eq. (2.8.1)} \quad SE_{B_{YX}} = \frac{13.82}{4.58} \sqrt{\frac{1 - .657^2}{50 - 2}} = .329$$

Estimated intercept

$$\text{Eq. (2.8.2)} \quad SE_{B_0} = 10.82 \sqrt{\frac{1}{50} + \frac{7.67^2}{(50 - 1)(20.95)}} = .301$$

Estimated value of \hat{Y} for a case where $X = 9$

$$\text{Eq. (2.8.3)} \quad SE_{\hat{Y}_i} = 10.82 \sqrt{\frac{1}{50} + \frac{(9 - 7.67)^2}{(50 - 1)(20.95)}} = 1.59.$$

Estimated r_{YX}

$$\text{Eq. (2.8.5)} \quad SE_{z'} = \frac{1}{\sqrt{50 - 3}} = .146$$

Estimated difference between B in two populations

$$\text{Eq. (2.8.6)} \quad SE_{B_v - B_w} = \sqrt{.329^2 + .329^2} = \sqrt{.2165} = .465.$$

Estimated difference between r 's in two large samples from different populations

$$\text{Eq. (2.8.7)} \quad SE_{r_v - r_w} = \sqrt{\frac{1 - .657^2}{50} + \frac{1 - .430^2}{50}} = \sqrt{.01136 + .01630} = .166.$$

The SE is inversely proportional to \sqrt{n} to a sufficient approximation when n is not small. Quadrupling n cuts SE approximately in half. To make a standard error x times as large as that for $n = 50$, compute $n* = n/x^2$, where $n*$ is the necessary sample size to attain x times the SE . For example, we found $SE_{B_{YX}} = .329$ for our sample of $n = 50$ cases. To make it half (.5) as large, we would need $n* = 50/.5^2 = 200$.

To change a standard error from SE to $SE*$, find $n* = n(SE/SE*)^2$. For example, to change the $SE_{B_{YX}}$ from .329 (for $n = 50$) to $SE* = .20$, we would need $n* = 50 (.329/.20)^2 = 135$ cases.

For differences between B s and r s, use their statistics from the trials to determine the desired changes in the SE s for the two samples and compute the anticipated SE of the difference (Eqs. 2.8.6 and 2.8.7). Adjust the ns as necessary.

2.9.2 Power of Null Hypothesis Significance Tests

In Section 2.8.3, we presented methods of appraising sample data in regard to α , the risk of mistakenly rejecting the null hypothesis when it is true, that is, drawing a spuriously positive conclusion (Type I error). We now turn our attention to methods of determining β ,¹⁶ the probability of failing to reject the null hypothesis when it is false (Type II error), and ways in which it can be controlled in research planning.

¹⁶We have been using β to represent the standardized regression coefficient. It is used here with a different meaning for consistency with the literature.

52 2. BIVARIATE CORRELATION AND REGRESSION

Any given test of a null hypothesis is a complex relationship among the following four parameters:

1. The power of the test, the probability of rejecting H_0 , defined as $1 - \beta$.
2. The region of rejection of H_0 as determined by the α level and whether the test is one-tailed or two-tailed. As α increases, for example from .01 to .05, power increases.
3. The sample size n . As n increases, power increases.
4. The magnitude of the effect in the population, or the degree of departure from H_0 . The larger this is, the greater the power.

These four parameters are so related that when any three of them are fixed, the fourth is completely determined. Thus, when an investigator decides for a given research plan the significance criterion α and n , the power of the test is determined. However, the investigator does not know what this power is without also knowing the magnitude of the effect size (ES) in the population, the estimation of which is the whole purpose of the study. The methods presented here focus on the standardized effect size, r in the present case.

There are three general strategies for estimating the size of the standardized population effect a researcher is trying to detect as "statistically significant":

1. To the extent that studies have been carried out by the current investigator or others which are closely similar to the present investigation, the ES s found in these studies reflect the magnitude that can be expected. Thus, if a review of the relevant literature reveals rs ranging from .32 to .43, the population ES in the current study may be expected to be somewhere in the vicinity of these values. Investigators who wish to be conservative may determine the power to detect a population ρ of .25 or .30.

2. In some research areas an investigator may posit some minimum population effect size that would have either practical or theoretical significance. An investigator may determine that unless $\rho = .05$, the importance of the relationship is insufficient to warrant a change in the policy or operations of the relevant institution. Another investigator may decide that a population correlation of .10 would have a material import for the adequacy of the theory within which the experiment has been designed, and thus would wish to plan the experiment so as to detect such an ES . Or a magnitude of B_{YX} that would be substantively important may be determined and other parameters estimated from other sources to translate B_{YX} into ρ .

3. A third strategy in deciding what ES values to use in determining the power of a study is to use certain suggested conventional definitions of *small*, *medium*, and *large* effects as population $\rho = .10$, $.30$, and $.50$, respectively (J. Cohen, 1988). These conventional ES s, derived from the average values in published studies in the social sciences, may be used either by choosing one of these values (for example, the conventional medium ES of $.30$) or by determining power for all three populations. If the latter strategy is chosen, the investigator would then revise the research plan according to an estimation of the relevance of the various ES s to the substantive problem. This option should be looked upon as the default option only if the earlier noted strategies are not feasible.

The point of doing a power analysis of a given research plan is that when the power turns out to be insufficient the investigator may decide to revise these plans, or even drop the investigation entirely if such revision is impossible. Obviously, because little or nothing can be done after the investigation is completed, determination of statistical power is of primary value as a preinvestigation procedure. If power is found to be insufficient, the research plan may be revised in ways that will increase it, primarily by increasing n , or increasing the number of levels or variability of the independent variable, or possibly by increasing α . A more complete general discussion of the concepts and strategy of power analysis may be found in J. Cohen (1965, 1988). It is particularly useful to use a computerized program for calculating the statistical

power of a proposed research plan, because such a program will provide a graphic depiction of the effect of each of the parameters (ES , n , α) on the resulting power to reject a false null hypothesis.

2.10 FACTORS AFFECTING THE SIZE OF r

2.10.1 The Distributions of X and Y

Because $r = 1.00$ only when each $z_X = z_Y$, it can only occur when the shapes of the frequency distributions for X and Y are exactly the same (or exactly opposite for $r = -1.00$). The greater the departure from distribution similarity, the more severe will the restriction be on the maximum possible r . In addition, as such distribution discrepancy increases, departure from homoscedasticity—equal error for different predicted values—must also necessarily increase. The decrease in the maximum possible value of (positive) r is especially noticeable under circumstances in which the two variables are skewed in opposite directions. One such common circumstance occurs when the two variables being correlated are each dichotomies: With very discrepant proportions, it is not possible to obtain a large positive correlation.

For example, suppose that a group of subjects has been classified into “risk takers” and “safe players” on the basis of behavior in an experiment, resulting in 90 risk takers and 10 safe players. A correlation is computed between this dichotomous variable and self classification as “conservative” versus “liberal” in a political sense, with 60 of the 100 subjects identifying themselves as conservative (Table 2.10.1). Even if all political liberals were also risk takers in the experimental situation, the correlation will be only (by Eq. 2.3.6):

$$r_{\phi} = \frac{400 - 0}{\sqrt{90 \cdot 10 \cdot 40 \cdot 60}} = .272.$$

It is useful to divide the issue of the distribution of variables into two components, those due to differences in the distribution of the underlying constructs and those due to the scales on which we have happened to measure our variables. Constraints on correlations associated with differences in distribution inherent in the constructs are not artifacts, but have real interpretive meaning. For example, gender and height for American adults are not perfectly correlated, but we need have no concern about an artificial upper limit on r attributable to this distribution difference. If gender completely determined height, there would only be two heights, one for men and one for women, and r would be 1.00.

TABLE 2.10.1
Bivariate Distribution of Experimental and Self-Reported
Conservative Tendency

		Experimental		Total:
Self-report		Risk takers	Safe players	Total:
		Liberal	Conservative	
	Liberal	40	0	40
	Conservative	50	10	60
Total:		90	10	100

Similarly the observed correlation between smoking and lung cancer is about .10 (estimated from figures provided by Doll & Peto, 1981). There is no artifact of distribution here; even though the risk of cancer is about 11 times as high for smokers, the vast majority of both smokers and nonsmokers alike will not contract lung cancer, and the relationship is low because of the nonassociation in these many cases.

Whenever the concept underlying the measure is logically continuous or quantitative¹⁷—as in the preceding example of risk taking and liberal versus conservative—it is highly desirable to measure the variables on a many-valued scale. One effect of this will be to increase the opportunity for reliable and valid discrimination of individual differences (see Section 2.10.2). To the extent that the measures are similarly distributed, the risk of underestimating the relationship between the conceptual variables will be reduced (see Chapter 4). However, the constraints on r due to unreliability are likely to be much more serious than those due to distribution differences on multivalued scales.

The Biserial r

When the only available measure of some construct X is a dichotomy, d_X , an investigator may wish to know what the correlation would be between the underlying construct and some other quantitative variable, Y . For example, X may be ability to learn algebra, which we measure by d_X , pass–fail. If one can assume that the “underlying” continuous variable X is normally distributed, and that the relationship with Y is linear, an estimate of the correlation between X and Y can be made, even though only d_X and Y are available. This correlation is estimated as

$$(2.10.1) \quad r_b = \frac{(M_{Y_p} - M_{Y_q})PQ}{h(sd_Y)} = r_{pb} \frac{\sqrt{PQ}}{h},$$

where M_{Y_p} and M_{Y_q} are the Y means for the two points of the dichotomy, P and Q ($= 1 - P$) are the proportions of the sample at these two points, and h is the ordinate (height) of the standard unit normal curve at the point at which its area is divided into P and Q portions (see Appendix Table C).

For example, we will return to the data presented in Table 2.3.1, where r_{pb} was found to be $-.707$. We now take the dichotomy to represent not the presence or absence of an experimentally determined stimulus but rather gross (1) versus minor (0) naturally occurring interfering stimuli as described by the subjects. This dichotomy is assumed to represent a continuous, normally distributed variable. The biserial r between stimulus and task score will be

$$r_b = \frac{(66 - 69.5)(.428)(.572)}{.392(2.45)} = -.893$$

where $.392$ is the height of the ordinate at the $.428, .572$ break, found by linear interpolation in Appendix Table C and $r_{pb} = -.707$.

The biserial r of $-.893$ may be taken to be an estimate of the product moment correlation that would have been obtained had X been a normally distributed continuous measure. It will always be larger than the corresponding point biserial r and, in fact, may even nonsensically exceed 1.0 when the Y variable is not normally distributed. When there is no overlap between the Y scores of the two groups, the r_b will be at least 1.0. It will be approximately 25% larger than the corresponding r_{pb} when the break on X is $.50 - .50$. The ratio of r_b/r_{pb} will increase

¹⁷ *Continuous* implies a variable on which infinitely small distinctions can be made; *quantitative* or *scaled* is more closely aligned to real measurement practice in the behavioral sciences, implying an ordered variable of many, or at least several, possible values. Theoretical constructs may be taken as continuous, but their measures will be quantitative in this sense.

as the break on X is more extreme; for example with a break of .90 – .10, r_b will be about two-thirds larger than r_{pb} .

Confidence limits are best established on r_{pb} or, equivalently, on the difference between the Y means corresponding to the two points of d_X .

Tetrachoric r

As we have seen, when the relationship between two dichotomies is investigated, the restriction on the maximum value of r_ϕ when their breaks are very different can be very severe. Once again, we can make an estimate of what the linear correlation would be if the two variables were continuous and normally distributed. Such an estimate is called the tetrachoric correlation. Because the formula for the tetrachoric correlation involves an infinite series and even a good approximation is a laborious operation, tetrachoric r s are obtained by means of computer programs. Tetrachoric r will be larger than the corresponding phi coefficient and the issues governing their interpretation and use are the same as for r_b and r_{pb} .

Caution should be exercised in the use of biserial and tetrachoric correlations, particularly in multivariate analyses. Remember that they are not observed correlations in the data, but rather hypothetical ones depending on the normality of the distributions underlying the dichotomies. Nor will standard errors for the estimated coefficients be the same as those for the product moment coefficients presented here.

2.10.2 The Reliability of the Variables

In most research in the behavioral sciences, the concepts that are of ultimate interest and that form the theoretical foundation for the study are only indirectly and imperfectly measured in practice. Thus, typically, interpretations of the correlations between variables as measured should be carefully distinguished from the relationship between the constructs or conceptual variables found in the theory.

The reliability of a variable (r_{XX}) may be defined as the correlation between the variable as measured and another equivalent measure of the same variable. In standard psychometric theory, the square root of the reliability coefficient $\sqrt{r_{XX}}$ may be interpreted as the correlation between the variable as measured by the instrument or test at hand and the “true” (error-free) score. Because true scores are not themselves observable, a series of techniques has been developed to estimate the correlation between the obtained scores and these (hypothetical) true scores. These techniques may be based on correlations among items, between items and the total score, between other subdivisions of the measuring instrument, or between alternative forms. They yield a reliability coefficient that is an estimate (based on a sample) of the population reliability coefficient.¹⁸ This coefficient may be interpreted as an index of how well the test or measurement procedure measures whatever it is that it measures. This issue should be distinguished from the question of the test’s *validity*, that is, the question of whether *what it measures is what the investigator intends that it measure*.

The discrepancy between an obtained reliability coefficient and a perfect reliability of 1.00 is an index of the relative amount of measurement error. Each observed score may be thought of as composed of some true value plus a certain amount of error:

$$(2.10.2) \quad X = X_t + X_e.$$

¹⁸Because this is a whole field of study in its own right, no effort will be made here to describe any of its techniques, or even the theory behind the techniques, in any detail. Excellent sources of such information include McDonald (1999) and Nunnally & Bernstein (1993).

56 2. BIVARIATE CORRELATION AND REGRESSION

These error components are assumed to have a mean of zero and to correlate zero with the true scores and with true or error scores on other measures. Measurement errors may come from a variety of sources, such as errors in sampling the domain of content, errors in recording or coding, errors introduced by grouping or an insufficiently fine system of measurement, errors associated with uncontrolled aspects of the conditions under which the test was given, errors due to short- or long-term fluctuation in individuals' true scores, errors due to the (idiosyncratic) influence of other variables on the individuals' responses, etc.

For the entire set of scores, the reliability coefficient equals the proportion of the observed score variable that is true score variance

$$(2.10.3) \quad r_{XX} = \frac{sd_{X_t}^2}{sd_X^2}$$

Because, as we have stated, error scores are assumed not to correlate with anything, r_{XX} may also be interpreted as that proportion of the measure's variance that is available to correlate with other measures. Therefore, the correlation between the observed scores (X and Y) for any two variables will be numerically smaller than the correlation between their respective unobservable true scores (X_t and Y_t). Specifically,

$$(2.10.4) \quad r_{XY} = r_{X_t Y_t} \sqrt{r_{XX} r_{YY}}.$$

Researchers sometimes wish to estimate the correlations between two theoretical constructs from the correlations obtained between the imperfect observed measures of these constructs. To do so, one corrects for attenuation (unreliability) by dividing r_{XY} by the square root of the product of the reliabilities (the maximum possible correlation between the imperfect measures). From Eq. (2.10.4),

$$(2.10.5) \quad r_{X_t Y_t} = \frac{r_{XY}}{\sqrt{r_{XX} r_{YY}}}.$$

Thus, if two variables, each with a reliability of .80, were found to correlate .44,

$$r_{X_t Y_t} = \frac{.44}{\sqrt{(.80)(.80)}} = .55.$$

Although correlations are subject to attenuation due to unreliability in either or both variables, bivariate regression coefficients are not affected by unreliability in Y . This can be seen from the following, where we consider unreliability only in Y . The regression coefficient expressed as the relationship between the perfectly reliable variables [by Eq. (2.4.3)] is

$$(2.10.6) \quad B_{Y_t X_t} = r_{X_t Y_t} \left(\frac{sd_{Y_t}}{sd_{X_t}} \right)$$

By Eq. (2.10.5), when $r_{XX} = 1.0$, $r_{XY} = r_{X_t Y_t} \sqrt{r_{YY}}$. By Eq. (2.10.3),

$$r_{YY} = \frac{sd_{Y_t}^2}{sd_{Y_t}^2 + sd_{Y_e}^2} \quad \text{and} \quad sd_Y = \sqrt{sd_{Y_t}^2 + sd_{Y_e}^2}$$

so

$$r_{XY_t} = \frac{r_{XY}}{\sqrt{sd_{Y_t}^2 / (sd_{Y_t}^2 + sd_{Y_e}^2)}} \quad \text{and} \quad r_{XY} = r_{XY_t} \frac{sd_{Y_t}}{\sqrt{sd_{Y_t}^2 + sd_{Y_e}^2}}.$$

Therefore, using Eq. (2.4.3) where $B_{YX} = r_{XY}(sd_Y/sd_X)$, substituting:

$$B_{YX} = r_{XY} \left(\frac{sd_{Y_t}}{\sqrt{sd_{Y_t}^2 + sd_{Y_e}^2}} \right) \left(\frac{\sqrt{sd_{Y_t}^2 + sd_{Y_e}^2}}{sd_X} \right)$$

and canceling

$$= r_{XY} \left(\frac{sd_{Y_t}}{sd_X} \right) = B_{Y,X}$$

As is generally true for coefficients based on a series of estimates, caution must be used in interpreting attenuation-corrected coefficients, because each of the coefficients used in the equation is subject to sampling error (as well as model assumption failure). Indeed, it is even possible to obtain attenuation-corrected correlations larger than 1.0 when the reliabilities come from different populations than r_{XY} , are underestimated, or when the assumption of uncorrelated error is false. Obviously, because the disattenuated r is hypothetical rather than based on real data, its confidence limits are likely to be very large.¹⁹

To reiterate, unreliability in variables as classically defined is a sufficient reason for low correlations; it *cannot* cause correlations to be spuriously high. Spuriously high correlations may, of course, be found when sources of *bias* are shared by variables, as can happen when observations are not “blind,” when subtle selection factors are operating to determine which cases can and cannot appear in the sample studied, and for yet other reasons.

2.10.3 Restriction of Range

A problem related to the question of reliability occurs under conditions when the range of one or both variables is restricted by the sampling procedure. For example, suppose that in the data presented in Table 2.2.2 and analyzed in Table 2.6.1 we had restricted ourselves to the study of faculty members who were less extreme with regard to years since Ph.D., occupying the restricted range of 5 to 11 years rather than the full range of 3 to 18 years. If the relationship is well described by a straight line and homoscedastic, we shall find that the variance of the Y scores about the regression line, $sd_{Y-\hat{Y}}^2$, remains about the same. Because when $r \neq 0$, sd_Y^2 will be decreased as an incidental result of the reduction of sd_X^2 , and because $sd_Y^2 = sd_{\hat{Y}}^2 + sd_{Y-\hat{Y}}^2$, the proportion of sd_Y^2 associated with X , namely, $sd_{\hat{Y}}^2$, will necessarily be smaller, and therefore, $r^2 (= sd_{\hat{Y}}^2/sd_Y^2)$ and r will be smaller. In the current example, r decreases from .657 to .388, and r^2 , the proportion of variance, from .432 to .151. (See Table 2.10.2.) When the relationship is completely linear, the regression coefficient, B_{YX} , will remain constant because the decrease in r will be perfectly offset by the increase in the ratio sd_Y/sd_X . It is 2.456 here, compared to 1.983 before. (It increased slightly in this example, but could just as readily have decreased slightly.) The fact that regression coefficients tend to remain constant over changes in the variability of X (providing the relationship is fully linear and the sample size sufficiently large to produce reasonable estimates) is an important property of regression coefficients. It is shown later how this makes them more useful as measures of relationship than correlation coefficients in some analytic contexts (Chapter 5).



CH02EX06

¹⁹Current practice is most likely to test “disattenuated” coefficients via latent variable models (described in Section 12.5.4), although the definition and estimation is somewhat different from the reasoning presented here.

TABLE 2.10.2
**Correlation and Regression of Number of Publications
on a Restricted Range of Time Since Ph.D.**

Publications	Time since Ph.D.	
<i>Y</i>	<i>X</i>	
3	6	
17	8	
11	9	
6	6	$r_{XY} = .388 (.657)^a$
48	10	
22	5	$r_{XY}^2 = .150 (.431)$
30	5	
21	6	$sd_{Y-\bar{Y}} = 11.10 (10.42)$
10	7	
27	11	$B_{YX} = 2.456 (1.983)$
<i>M</i>	19.50	7.30
<i>sd</i>	12.04	1.31
<i>sd</i> ²	144.94	1.71

^aParenthetic values are those for the original (i.e., unrestricted) sample.

Suppose that an estimate of the correlation that would be obtained from the full range is desired, when the available data have a curtailed or restricted range for *X*. If we know the *sd_Y* of the unrestricted *X* distribution as well as the *sd_{X_c}* for the curtailed sample and the correlation between *Y* and *X* in the curtailed sample ($r_{X_c Y}$), we may estimate r_{XY} by

$$(2.10.7) \quad \tilde{r}_{YX} = \frac{r_{YX_c} (sd_X / sd_{X_c})}{\sqrt{1 + r_{YX_c}^2 ((sd_X^2 / sd_{X_c}^2) - 1)}}$$

For example, $r = .25$ is obtained on a sample for which $sd_{X_c} = 5$ whereas the sd_X of the population in which the investigator is interested is estimated to be 12. Situations like this occur, for example, when some selection procedure such as an aptitude test has been used to select personnel and those selected are later assessed on a criterion measure. If the finding on the restricted (employed) sample is projected to the whole group originally tested, \tilde{r}_{XY} would be estimated to be

$$\tilde{r}_{XY} = \frac{.25(12/5)}{\sqrt{1 + .25^2[(12/5)^2 - 1]}} = \frac{.60}{\sqrt{1.2975}} = .53$$

It should be emphasized that .53 is an estimate and assumes that the relationship is linear and homoscedastic, which might not be the case. There are no appropriate confidence limits on this estimate.

It is quite possible that restriction of range in either *X* or *Y*, or both, may occur as an incidental by-product of the sampling procedure. Therefore, it is important in any study to report the *sds* of the variables used. Because under conditions of homoscedasticity and linearity regression coefficients are not affected by range restriction, comparisons of different samples using the same variables should usually be done on the regression coefficients rather than on the correlation coefficients when *sds* differ. Investigators should be aware, however, that the questions answered by these comparisons are not the same. Comparisons of correlations

answer the question, Does X account for as much of the variance in Y in group E and in Group F ? Comparisons of regression coefficients answer the question, Does a change in X make the same amount of score difference in Y in group E as it does in group F ?

Although the previous discussion has been cast in terms of restriction in range, an investigator may be interested in the reverse—the sample in hand has a range of X values that is large relative to the population of interest. This could happen, for example, if the sampling procedure was such as to include disproportionately more high- and low- X cases and fewer middle values. Equation (2.10.7) can be employed to estimate the correlation in the population of interest (whose range in X is less) by reinterpreting the subscript C in the equation to mean changed (including increased) rather than curtailed. Thus, r_{YX_C} and sd_{X_C} are the “too large” values in the sample, sd_Y is the (smaller) sd of the population of interest, and the estimated r in that population will be smaller. Note that the ratio sd_Y/sd_{X_C} , which before was greater than one, is now smaller than one. Because the correlation (ES) will be higher in a sample with a larger sd , sampling in order to produce a larger sd , as in studies in which the number of “cases” is larger than in a random sample of the general population, is a major strategy for increasing the statistical power of a study.

2.10.4 Part-Whole Correlations

Occasionally we will find that a correlation has been computed between some variable J and another variable W , which is the sum of scores on a set of variables including J . Under these circumstances a positive correlation can be expected between J and W due to the fact that W includes J , even when there is no correlation between J and $W - J$. For example, if k test items of equal sd and zero r with each other are added together, each of the items will correlate exactly $1/\sqrt{k}$ with the total score. For the two-item case, therefore, each item would correlate .707 with their sum, W , when neither correlates with the other. On the same assumptions of zero correlation between the variables but with unequal sds , the variables are effectively weighted by their differing sds and the correlation of J with W will be equal to $sd_J/\sqrt{\sum sds_i^2}$, where sds are summed over the items. Obviously, under these circumstances $r_{J(W-J)} = 0$. In the more common case where the variables or items are correlated, the correlation of J with $W - J$ may be obtained by

$$(2.10.8) \quad r_{J(W-J)} = \frac{r_{JW}sd_W - sd_J}{\sqrt{sd_W^2 + sd_J^2 - 2r_{JW}sd_Wsd_J}}$$

This is not an estimate and may be tested via the usual t test for the significance of r .

Given these often substantial spurious correlations between elements and totals including the elements, it behooves the investigator to determine $r_{J(W-J)}$, or at the very least determine the expected value when the elements are uncorrelated before interpreting r_{JW} . Such a circumstance often occurs when the interest is in the correlation of a single item with a composite that includes that item, as is carried out in psychometric analysis.

Change Scores

It is not necessary that the parts be literally added in order to produce such spurious correlation. If a subscore is subtracted, a spurious negative component in the correlation will also be produced. One common use of such difference scores in the social sciences is the use of post-minus pretreatment (change) scores. If such change scores are correlated with the pre- and posttreatment scores from which they have been obtained, we will typically find that subjects initially low on X will have larger gains than those initially high on X , and that those with the

highest final scores will have made greater gains than those with lower final scores. Again, if $sd_{pre} = sd_{post}$ and $r_{pre\ post} = 0$, the $r_{pre\ change} = -.707$ and $r_{post\ change} = +.707$. Although in general, we would expect the correlation between pre- and posttreatment scores to be some positive value, it will be limited by their respective reliabilities (Section 2.10.2) as well as by individual differences in true change.

If the post- minus pretreatment variable has been created in order to control for differences in pretreatment scores, the resulting negative correlations between pretreatment and change scores may be taken as a failure to remove all influence of pretreatment scores from posttreatment scores. This reflects the regression to the mean phenomenon discussed in Section 2.5 and the consequent interpretive risks. The optimal methods of handling this and related problems are the subject of a whole literature (Collins & Horn, 1993) and cannot be readily summarized. However, the appropriate analysis, as always, depends on the underlying causal model. (See Chapters 5, 12, and 15 for further discussion of this problem.)

2.10.5 Ratio or Index Variables

Ratio (index or rate) scores are those constructed by dividing one variable by another. When a ratio score is correlated with another variable or with another ratio score, the resulting correlation depends as much on the denominator of the score as it does on the numerator. Because it is usually the investigator's intent to "take the denominator into account" it may not be immediately obvious that the correlations obtained between ratio scores may be spurious—that is, may be a consequence of mathematical necessities that have no valid interpretive use. Ratio correlations depend, in part, upon the correlations between all numerator and denominator terms, so that $r_{(Y/Z)X}$ is a function of r_{YZ} and r_{XZ} as well as of r_{YX} , and $r_{(Y/Z)(X/W)}$ depends on r_{YW} and r_{XZ} as well as on the other four correlations. These correlations also involve the coefficients of variation

$$(2.10.9) \quad v_X = \frac{sd_X}{M_X}$$

of each of the variables. Although the following formula is only a fair approximation of the correlation between ratio scores (requiring normal distributions and homoscedasticity and dropping all terms involving powers of v greater than v^2), it serves to demonstrate the dependence of correlations between ratios on all vs and on rs between all variable pairs:

$$(2.10.10) \quad r(Y/Z)(X/W) = \frac{r_{YX}v_Yv_X - r_{YW}v_Yv_W - r_{XZ}v_Xv_Z - r_{ZW}v_Zv_W}{\sqrt{v_Y^2 + v_Z^2 - 2r_{YZ}v_Yv_Z}\sqrt{v_X^2 + v_W^2 - 2r_{XW}v_Xv_W}}$$

When the two ratios being correlated have a common denominator, the possibility of spurious correlations becomes apparent. Under these circumstances, the approximate formula for the correlation simplifies, because $Z = W$. If all coefficients of variation are equal when all three variables are uncorrelated we will find $r_{(Y/Z)(X/Z)} \approx .50$.

Because the coefficient of variation depends on the value of the mean, it is clear that whenever this value is arbitrary, as it is for many psychological scores, the calculated r is also arbitrary. Thus, ratios should not be correlated unless each variable is measured on a ratio scale, a scale for which a zero value means literally none of the variable (see Chapters 5 and 12). Measures with ratio scale properties are most commonly found in the social sciences in the form of counts or frequencies.

At this point it may be useful to distinguish between rates and other ratio variables. Rates may be defined as variables constructed by dividing the number of instances of some phenomenon by the total number of opportunities for the phenomenon to occur; thus, they are literally

proportions. Rates or proportions are frequently used in ecological or epidemiological studies where the units of analysis are aggregates of people or areas such as counties or census tracts. In such studies, the numerator represents the incidence or prevalence of some phenomenon and the denominator represents the population at risk. For example, a delinquency rate may be calculated by dividing the number of delinquent boys ages 14–16 in a county by the total number of boys ages 14–16 in the county. This variable may be correlated across the counties in a region with the proportion of families whose incomes are below the poverty level, another rate. Because, in general, the denominators of these two rates will reflect the populations of the counties, which may vary greatly, they can be expected to be substantially correlated. In other cases the denominators may actually be the same—as, for example, in an investigation of the relationship between delinquency rates and school dropout rates for a given age-gender group. The investigator will typically find that these rates have characteristics that minimize the problem of spurious correlation. In most real data, the coefficients of variation of the numerators will be substantially larger than the coefficients of variation of the denominators, and thus the correlation between rates will be determined substantially by the correlation between the numerators. Even in such data, however, the resulting proportions may not be optimal for the purpose of linear correlation. Chapter 6 discusses some nonlinear transformations of proportions, which may be more appropriate for analysis than the raw proportions or rates themselves.

Experimentally produced rates may be more subject to problems of spurious correlation, especially when there are logically alternative denominators. The investigator should determine that the correlation between the numerator and denominator is very high (and positive), because in general the absence of such a correlation suggests a faulty logic in the study. In the absence of a large correlation, the coefficients of variation of the numerator should be substantially larger than that of the denominator if the problem of spurious correlation is to be minimized.

Other Ratio Scores

When the numerator does not represent some subclass of the denominator class, the risks involved in using ratios are even more serious, because the likelihood of small or zero correlations between numerators and denominators and relatively similar values of v is greater. If the variables do not have true zeros and equal intervals, correlations involving ratios should probably be avoided altogether, and an alternative method for removing the influence of Z from X or Y should be chosen, as presented in Chapters 3 and 12.

The difficulties that may be encountered in correlations involving rates and ratios may be illustrated by the following example. An investigator wishes to determine the relationship between visual scanning and errors on a digit-symbol (d-s) task. All subjects are given 4 minutes to work on the task. Because subjects who complete more d-s substitutions have a greater opportunity to make errors, the experimenter decides, reasonably enough, to determine the error rate by dividing the number of errors by the number of d-s substitutions completed. Table 2.10.3 displays the data for 10 subjects. Contrary to expectation, subjects who completed more d-s tasks did not tend to produce more errors ($r_{ZX} = -.105$), nor did they scan notably more than did low scorers ($r_{ZY} = .023$). Nevertheless, when the two ratio scores are computed, they show a substantial positive correlation (.427) in spite of the fact that the numerators showed slight negative correlation (−.149), nor is there any tendency for scanning and errors to be correlated for any given level of d-s task completion. Thus, because $r_{ZZ} = 1$, the $r_{(X/Z)(Y/Z)}$ may here be seen to be an example of spurious correlation.²⁰

²⁰An alternative method of taking into account the number completed in considering the relationship between errors and number of scans might be to partial Z (see subsequent chapters).

TABLE 2.10.3
An Example of Spurious Correlation Between Ratios

Subject	No. completed d-s tasks (Z)	No. errors (X)	No. scans (Y)	Error rate (X/Z)	Scan rate (Y/Z)
1	25	5	24	.20	.96
2	29	3	30	.10	1.03
3	30	3	27	.10	.90
4	32	4	30	.12	.94
5	37	3	18	.08	.49
6	41	2	33	.05	.80
7	41	3	27	.07	.66
8	42	5	21	.12	.50
9	43	3	24	.07	.56
10	43	5	33	.12	.77

$r_{ZX} = -.105$, $r_{ZY} = .106$, $r_{XY} = -.149$
 $r_{(X/Z)(Y/Z)} = .427$

2.10.6 Curvilinear Relationships

When the relationship between the two variables is only moderately well fitted by a straight line, the correlation coefficient that indicates the degree of linear relationship will underestimate the predictability from one variable to the other. Frequently the relationship, although curvilinear, is monotonic; that is, increases in Z are accompanied by increases (or decreases) in Y , although not at a constant rate. Under these circumstances, some (nonlinear) monotonic transformation of X or Y or both may straighten out the regression line and provide a better indication of the size of the relationship between the two variables (an absolutely larger r). Because there are several alternative ways of detecting and handling curvilinear relationships, the reader is referred to Chapters 4 and 6 for a detailed treatment of the issues.

2.11 SUMMARY

A linear relationship exists between two quantitative variables when there is an overall tendency for increases in the value of one variable to be accompanied by increases in the other variable (a positive relationship), or for increases in the first to be accompanied by decreases in the second (a negative relationship); (Section 2.1). Efforts to index the degree of linear relationship between two variables must cope with the problem of the different units in which variables are measured. Standard (z) scores are a conversion of scores into distances from their own means, in standard deviation units, and they render different scores comparable. The Pearson product moment correlation coefficient, r , is a measure of the degree of relationship between two variables, X and Y , based on the discrepancies of the subjects' paired z scores, $z_X - z_Y$. r varies between -1 and $+1$, which represent perfect negative and perfect positive linear relationships, respectively. When $r = 0$, there is no linear correlation between the variables (Section 2.2).

r can be written as a function of z score products, a function of variances and covariance, or in terms of the original units. Special simplified formulas are available for r when one variable is a dichotomy (point biserial r), when both variables are dichotomies (r_ϕ), or when the data are two sets of complete ranks (Spearman rank order correlation); (Section 2.3).

The regression coefficient, B_{YX} , gives the optimal rule for a linear estimate of Y from X , and is the change in Y units per unit change in X , that is, the slope of the regression line. The intercept, B_0 , gives the predicted value of Y for a zero value of X . B_{YX} and B_0 are optimal in the sense that they provide the smallest squared discrepancies between Y and estimated \hat{Y} . r is the regression coefficient for the standardized variables. When X is centered, $B_0 = M_Y$ (Section 2.4). Unless $r = 1$, it is a mathematical necessity that the average score for a variable being estimated (e.g., \hat{Y}) will be relatively closer to M_Y than the value from which it is being estimated (e.g., X) will be to its mean (M_X) when both are measured in sd units (Section 2.5).

When Y is estimated from X the sd of the difference between observed scores and the estimated scores (the sample standard error of estimate) can be computed from r and sd_Y . The coefficient of alienation represents the error as a proportion of the original sd_Y . r^2 equals the proportion of the variance (sd^2) of each of the variables that is shared with or can be estimated from the other (Sections 2.6 and 2.7).

The two major methods of statistical inference are estimation and null hypothesis testing. The formal model assumptions are presented (Section 2.8.1), confidence intervals are given for B_{YX} , B_{Y0} , r_{XY} , for differences between independent sample values of these statistics, and for the estimated \hat{Y}_i (Section 2.8.2). Given α , confidence intervals provide the range of values within which the corresponding population values can be expected to fall. In Section 2.8.3, we present the null hypothesis tests for simple regression and correlation statistics. Section 2.8.4 critiques null hypothesis testing and contrasts it with the use of confidence intervals.

The degree of accuracy (precision) in the estimation of parameters is reflected in the statistic's confidence interval. The probability of null hypothesis rejection (statistical power) can be assessed before the research sample is collected (Section 2.9). Methods of finding the sample size to produce a margin of error for a given degree of inclusion in the confidence interval (95%, 80%) are presented (Section 2.9.1) and methods are given for determining the sample size needed for the desired statistical power, that is, the probability of rejecting the null hypothesis (Section 2.9.2).

A number of characteristics of the X and Y variables will affect the size of the correlation between them. Among these are differences in the distribution of the X and Y variables (Section 2.10.1), unreliability in one or both variables (Section 2.10.2), and restriction of the range of one or both variables (Section 2.10.3). When one variable is included as a part of the other variable, the correlation between them will reflect this overlap (Section 2.10.4). Scores obtained by dividing one variable by another will produce spurious correlation with other variables under some conditions (Section 2.10.5). The r between two variables will be an underestimate of the magnitude of their relationship when a curved rather than a straight line best fits the bivariate distribution (Section 2.10.6). Under such circumstances, transformation of one or both variables or multiple representation of one variable will provide a better picture of the relationship between the variables.