

**Springer Series  
in Statistics**

**Phillip Good**

# Permutation Tests

A Practical Guide to  
Resampling Methods for  
Testing Hypotheses

Second Edition



Springer

# **Springer Series in Statistics**

*Advisors:*

P. Bickel, P. Diggle, S. Fienberg, K. Krickeberg,  
I. Olkin, N. Wermuth, S. Zeger

**Springer Science+Business Media, LLC**

# **Springer Series in Statistics**

---

- Andersen/Borgan/Gill/Keiding:* Statistical Models Based on Counting Processes.
- Berger:* Statistical Decision Theory and Bayesian Analysis, 2nd edition.
- Bolfarine/Zacks:* Prediction Theory for Finite Populations.
- Borg/Groenen:* Modern Multidimensional Scaling: Theory and Applications
- Brockwell/Davis:* Time Series: Theory and Methods, 2nd edition.
- Chen/Shao/Ibrahim:* Monte Carlo Methods in Bayesian Computation.
- Efromovich:* Nonparametric Curve Estimation: Methods, Theory, and Applications.
- Fahrmeir/Tutz:* Multivariate Statistical Modelling Based on Generalized Linear Models.
- Farebrother:* Fitting Linear Relationships: A History of the Calculus of Observations 1750-1900.
- Federer:* Statistical Design and Analysis for Intercropping Experiments, Volume I: Two Crops.
- Federer:* Statistical Design and Analysis for Intercropping Experiments, Volume II: Three or More Crops.
- Fienberg/Hoaglin/Kruskal/Tanur (Eds.):* A Statistical Model: Frederick Mosteller's Contributions to Statistics, Science and Public Policy.
- Fisher/Sen:* The Collected Works of Wassily Hoeffding.
- Good:* Permutation Tests: A Practical Guide to Resampling Methods for Testing Hypotheses, 2nd edition.
- Gouriéroux:* ARCH Models and Financial Applications.
- Grandell:* Aspects of Risk Theory.
- Haberman:* Advanced Statistics, Volume I: Description of Populations.
- Hall:* The Bootstrap and Edgeworth Expansion.
- Härdle:* Smoothing Techniques: With Implementation in S.
- Hart:* Nonparametric Smoothing and Lack-of-Fit Tests.
- Hartigan:* Bayes Theory.
- Hedayat/Sloane/Stufken:* Orthogonal Arrays: Theory and Applications.
- Heyde:* Quasi-Likelihood and its Application: A General Approach to Optimal Parameter Estimation.
- Huet/Bouvier/Gruet/Jolivet:* Statistical Tools for Nonlinear Regression: A Practical Guide with S-PLUS Examples.
- Kolen/Brennan:* Test Equating: Methods and Practices.
- Kotz/Johnson (Eds.):* Breakthroughs in Statistics Volume I.
- Kotz/Johnson (Eds.):* Breakthroughs in Statistics Volume II.
- Kotz/Johnson (Eds.):* Breakthroughs in Statistics Volume III.
- Küchler/Sørensen:* Exponential Families of Stochastic Processes.
- Le Cam:* Asymptotic Methods in Statistical Decision Theory.
- Le Cam/Yang:* Asymptotics in Statistics: Some Basic Concepts.
- Longford:* Models for Uncertainty in Educational Testing.
- Miller, Jr.:* Simultaneous Statistical Inference, 2nd edition.
- Mosteller/Wallace:* Applied Bayesian and Classical Inference: The Case of the Federalist Papers.
- Parzen/Tanabe/Kitagawa:* Selected Papers of Hirotugu Akaike.
- Politis/Romano/Wolf:* Subsampling.

*(continued after index)*

Phillip Good

# Permutation Tests

A Practical Guide to Resampling Methods  
for Testing Hypotheses

Second Edition

With 14 Figures



Springer

Phillip Good  
205 W. Utica Avenue  
Huntington Beach, CA 92648  
USA  
brother\_unknown@yahoo.com

Library of Congress Cataloging-in-Publication Data  
Good, Phillip I.

Permutation tests : a practical guide to resampling methods for testing hypotheses / Phillip Good.—2nd ed.  
p. cm.—(Springer series in statistics)  
Includes bibliographical references and index.  
ISBN 978-1-4757-3237-5      ISBN 978-1-4757-3235-1 (eBook)  
DOI 10.1007/978-1-4757-3235-1  
1. Statistical hypothesis testing.   2. Resampling (Statistics)  
I. Title.   II. Series.  
QA277.G643   2000  
519.5'6—dc21                          99-16557

Printed on acid-free paper.

© 2000, 1994 Springer Science+Business Media New York  
Originally published by Springer-Verlag New York, Inc. in 2000  
Softcover reprint of the hardcover 2nd edition 2000

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher Springer Science+Business Media, LLC, except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden. The use of general descriptive names, trade names, trademarks, etc., in this publication, even if the former are not especially identified, is not to be taken as a sign that such names, as understood by the Trade Marks and Merchandise Marks Act, may accordingly be used freely by anyone.

Production managed by Timothy Taylor; manufacturing supervised by Jerome Basma.  
Typeset by TechBooks, Fairfax, VA.

9 8 7 6 5 4 3 2 1

ISBN 978-1-4757-3237-5

SPIN 10735885

# Preface to the Second Edition

In 1982, I published several issues of a *samizdat* scholarly journal called *Randomization* with the aid of an 8-bit, 1-MH personal computer with 48K of memory (upgraded to 64K later that year) and floppy disks that held 400 Kbytes. A decade later, working on the first edition of this text, I used a 16-bit, 33-MH computer with 1 Mb of memory and a 20-Mb hard disk. This preface to the second edition comes to you via a 32-bit, 300-MH computer with 64-Mb memory and a 4-Gb hard disk. And, yes, I paid a tenth of what I paid for my first computer.

This relationship between low-cost readily available computing power and the rising popularity of permutation tests is no coincidence. Simply put, it is faster today to compute an exact p-value than to look up an approximation in a table of the not-quite-appropriate statistic. As a result, more and more researchers are using Permutation Tests to analyze their data.

Of course, some of the increased usage has also come about through the increased availability of and improvements in off-the-shelf software, as can be seen in the revisions in this edition to Chapter 12 (Publishing Your Results) and Chapter 13 (Increasing Computation Efficiency).

These improvements helped persuade me it was the time to publish a first course in statistics based entirely on resampling methods (an idea first proposed by the late F.N. David). As a result, *Permutation Tests* has become two texts: one, *Resampling Methods*, designed as a first course, and this second edition aimed at upper division graduate students and practitioners who may already be familiar with the application of other statistical procedures. The popular question section at the end of each chapter now contains a number of thesis-level questions, which may or may not be solvable in their present form. While the wide applicability of permutation tests continues to be emphasized here, their limitations are also revealed. Examples include expanded sections on comparing variances (Chapter 3, Testing Hypotheses), testing interactions in balanced designs (Chapter 4, Experimental Design), and multiple regression (Chapter 7, Dependence).

Sections on Sequential Analysis (Chapter 4) and comparing spatial distributions (Chapter 8) are also new. Recent major advances in the analysis of multiple dependent tests are recorded in Chapter 5 on Multivariate Analysis.

My thanks to the many individuals who previewed chapters for this edition, including, in alphabetical order, Brian Cade, Mike Ernst, Barbara Heller, John Kimmel, Patrick Onghena, Fortunato Pesarin, and John Thaden.

Phillip Good  
Huntington Beach, California

# Preface to the First Edition

Permutation tests permit us to choose the test statistic best suited to the task at hand. This freedom of choice opens up a thousand practical applications, including many which are beyond the reach of conventional parametric statistics. Flexible, robust in the face of missing data and violations of assumptions, the permutation test is among the most powerful of statistical procedures. Through sample size reduction, permutation tests can reduce the costs of experiments and surveys.

This text on the application of permutation tests in biology, medicine, science, and engineering may be used as a step-by-step self-guiding reference manual by research workers and as an intermediate text for undergraduates and graduates in statistics and the applied sciences with a first course in statistics and probability under their belts.

Research workers in the applied sciences are advised to read through Chapters 1 and 2 once quickly before proceeding to Chapters 3 through 8 which cover the principal applications they are likely to encounter in practice.

Chapter 9 is a must for the practitioner, with advice for coping with real-life emergencies such as missing or censored data, after-the-fact covariates, and outliers.

Chapter 10 uses practical applications in archeology, biology, climatology, education and social science to show the research worker how to develop new permutation statistics to meet the needs of specific applications. The practitioner will find Chapter 10 a source of inspiration as well as a practical guide to the development of new and novel statistics.

The expert system in Chapter 11 will guide you to the correct statistic for your application. Chapter 12, more “must” reading, provides practical advice on experimental design and shows how to document the results of permutation tests for publication.

Chapter 13 describes techniques for reducing computation time; and a guide to off-the-shelf statistical software is provided in an appendix.

The sequence of recommended readings is somewhat different for the student and will depend on whether he or she is studying the permutation tests by

themselves or as part of a larger course on resampling methods encompassing both the permutation test and the bootstrap resampling method.

This book can replace a senior-level text on testing hypotheses. I have also found it of value in introducing students who are primarily mathematicians to the applications which make statistics a unique mathematical science. Chapters 1, 2, and 14 provide a comprehensive introduction to the theory. Despite its placement in the latter part of the text, Chapter 14, on the theory of permutation tests, is self-standing. Chapter 3 on applications also deserves a careful reading. Here in detail are the basic testing situations and the basic tests to be applied to them. Chapters 4, 5, and 6 may be used to supplement Chapter 3, time permitting (the first part of Chapter 6 describing the Fisher exact test is a must). Rather than skipping from section to section, it might be best for the student to consider one of these latter chapters in depth—supplementing his or her study with original research articles.

My own preference is to parallel discussions of permutation methods with discussion of a second resampling method, the bootstrap. Again, Chapters 1, 2, and 3—supplemented with portions of Chapter 14—are musts. Chapter 7, on tests of dependence, is a natural sequel. Students in statistical computing also are asked to program and test at least one of the advanced algorithms in Chapter 12.

For the reader's convenience, the bibliography is divided into four parts: the first consists of 34 seminal articles; the second of two dozen background articles referred to in the text that are not directly concerned with permutation methods; the third of 111 articles on increasing computational efficiency; and a fourth, principal bibliography of 574 articles and books on the theory and application of permutation techniques.

Exercises are included at the end of each chapter to enhance and reinforce your understanding. But the best exercise of all is to substitute your own data for the examples in the text.

My thanks to Symantek, TSSI, and Perceptronics without whose Grand-View® outliner, Exact® equation generator, and Einstein Writer® word processor this text would not have been possible.

I am deeply indebted to Mike Chernick for our frequent conversations and his many invaluable insights, to Mike Ernst, Alan Forsythe, Karim Hiriji, John Ludbrook, Reza Modarres, and William Schucany for reading and commenting on portions of this compuscript and to my instructors at Berkeley including E. Fix, J. Hodges, E. Lehmann, and J. Neyman.

Phillip Good  
Huntington Beach, California

# Contents

Preface to the Second Edition . . . . .	v
Preface to the First Edition . . . . .	vii
1. A Wide Range of Applications . . . . .	1
1.1. Permutation Tests . . . . .	1
1.1.1. Applications . . . . .	1
1.2. “I Lost the Labels” . . . . .	4
1.3. Five Steps to a Permutation Test . . . . .	6
1.3.1. Analyze the Problem . . . . .	6
1.3.2. Choose a Test Statistic . . . . .	6
1.3.3. Compute the Test Statistic . . . . .	7
1.3.4. Rearrange the Observations . . . . .	7
1.3.5. Make a Decision . . . . .	8
1.4. What’s in a Name? . . . . .	8
1.4.1. Comparison with Other Tests . . . . .	9
1.4.2. Sampling from the Data at Hand . . . . .	9
1.5. History . . . . .	10
1.6. To Learn More . . . . .	12
1.7. Questions . . . . .	12
2. A Simple Test . . . . .	13
2.1. Properties of the Test . . . . .	13
2.2. Fundamental Concepts . . . . .	14
2.2.1. Population and Sample Distributions . . . . .	14
2.2.2. Two Types of Error . . . . .	16
2.2.2.1. Losses and Risk . . . . .	18
2.2.3. Significance Level and Power . . . . .	19
2.2.4. Power and Sample Size . . . . .	20
2.2.5. Power and the Alternative . . . . .	22

2.2.6. Exact, Unbiased Tests . . . . .	23
2.2.7. Exchangeable Observations . . . . .	24
2.3. Which Test? . . . . .	25
2.4. World Views . . . . .	27
2.4.1. Generalized Permutations . . . . .	28
2.5. Questions . . . . .	29
3. Testing Hypotheses . . . . .	31
3.1. One-Sample Tests . . . . .	31
3.1.1. Tests for a Location Parameter . . . . .	31
3.1.2. Properties of the Test . . . . .	33
3.1.3. Exact Significance Levels: A Digression . . . . .	33
3.2. Confidence Intervals . . . . .	34
3.2.1. Comparison with Other Tests . . . . .	35
3.3. Comparison of Locations . . . . .	36
3.3.1. An Example . . . . .	36
3.3.2. Violation of Assumptions . . . . .	37
3.4. Comparing Dispersions . . . . .	38
3.4.1. The Permutation Approach . . . . .	38
3.4.2. The Bootstrap Approach . . . . .	42
3.5. $k$ -Sample Comparisons . . . . .	42
3.5.1. Comparing Location Parameters . . . . .	42
3.5.2. Pitman Correlation . . . . .	45
3.5.3. Effect of Ties . . . . .	47
3.5.4. Cochran–Armitage Test . . . . .	48
3.5.5. Linear Estimation . . . . .	49
3.5.6. A Unifying Theory . . . . .	49
3.6. Blocking . . . . .	49
3.6.1. Extending the Range of Applications . . . . .	51
3.7. Matched Pairs . . . . .	51
3.8. Questions . . . . .	52
4. Experimental Designs . . . . .	54
4.1. Introduction . . . . .	54
4.2. Balanced Designs . . . . .	54
4.2.1. Main Effects . . . . .	55
4.2.2. An Example . . . . .	57
4.2.3. Testing for Interactions . . . . .	59
4.2.3.1. Blocking . . . . .	60
4.2.3.2. Synchronized Permutations . . . . .	61
4.2.3.3. To Learn More . . . . .	62
4.2.4. Designing an Experiment . . . . .	62
4.2.5. Latin Squares . . . . .	63
4.2.6. Other Designs . . . . .	65

4.3. Analysis of Covariance . . . . .	65
4.3.1. Covariates Defined . . . . .	65
4.3.2. Eliminate the Functional Relationship . . . . .	65
4.3.3. Selecting Variables . . . . .	66
4.3.4. Restricted Randomization . . . . .	66
4.4. Unbalanced Designs . . . . .	67
4.4.1. Missing Combinations . . . . .	68
4.4.2. The Boot-Perm Test . . . . .	70
4.5. Clinical Trials . . . . .	70
4.5.1. Avoiding Unbalanced Designs . . . . .	70
4.5.2. Missing Data . . . . .	71
4.6. Sequential Analysis . . . . .	72
4.7. Very Large and Very Small Samples . . . . .	73
4.8. Questions . . . . .	74
 5. Multivariate Analysis . . . . .	76
5.1. Introduction . . . . .	76
5.2. One- and Two-Sample Comparisons . . . . .	76
5.2.1. Hotelling's $T^2$ . . . . .	76
5.2.2. An Example . . . . .	79
5.2.3. Doing the Computations . . . . .	80
5.2.4. Weighting the Variables . . . . .	81
5.2.5. Interpreting the Results . . . . .	82
5.2.6. Alternative Statistics . . . . .	82
5.3. Runs Test . . . . .	83
5.3.1. Which Statistic? . . . . .	84
5.4. Experimental Designs . . . . .	86
5.4.1. Matched Pairs . . . . .	86
5.4.2. Block Effects . . . . .	86
5.5. Multiple Comparisons . . . . .	87
5.5.1. Step Up or Step Down? . . . . .	87
5.5.1.1. Standardized Statistics . . . . .	87
5.5.1.2. Paired Sample Tests . . . . .	88
5.5.2. Fisher's Omnibus Statistic . . . . .	89
5.6. Repeated Measures . . . . .	90
5.6.1. Missing Data . . . . .	91
5.6.2. Bioequivalence . . . . .	91
5.6.3. Omnibus Test . . . . .	92
5.7. Questions . . . . .	92
 6. Categorical Data . . . . .	94
6.1. Fisher's Exact Test . . . . .	94
6.1.1. One-Tailed and Two-Tailed Tests . . . . .	96
6.1.2. The Two-Tailed Test . . . . .	97

6.1.3. Determining the <i>p</i> -Value . . . . .	97
6.1.4. What Is the Alternative? . . . . .	98
6.1.5. Increasing the Power . . . . .	99
6.2. Odds Ratio . . . . .	99
6.2.1. Stratified $2 \times 2$ 's . . . . .	100
6.3. Exact Significance Levels . . . . .	102
6.4. Unordered $r \times c$ Contingency Tables . . . . .	103
6.4.1. Agreement Between Observers . . . . .	106
6.4.2. What Should We Randomize? . . . . .	106
6.4.3. Underlying Assumptions . . . . .	107
6.5. Ordered Contingency Tables . . . . .	109
6.5.1. Ordered $2 \times c$ Tables . . . . .	109
6.5.2. More than Two Rows and Two Columns . . . . .	110
6.5.2.1. Singly Ordered Tables . . . . .	110
6.5.2.2. Doubly Ordered Tables . . . . .	111
6.5.2.3. An Example . . . . .	112
6.6. Covariates . . . . .	113
6.7. Higher Dimensional Tables . . . . .	114
6.8. To Learn More . . . . .	115
6.9. Questions . . . . .	115
7. Dependence . . . . .	118
7.1. The Models . . . . .	118
7.2. Testing for Independence . . . . .	119
7.2.1. Independence . . . . .	119
7.2.2. Quadrant Dependence . . . . .	119
7.3. Testing for Trend . . . . .	120
7.4. Serial Correlation . . . . .	121
7.4.1. An Example . . . . .	122
7.4.2. Trend . . . . .	124
7.4.3. First-Order Dependence . . . . .	124
7.4.4. An Economic Model . . . . .	125
7.5. Known Models . . . . .	126
7.5.1. Testing a Specific Hypothesis . . . . .	126
7.5.2. Testing a General Hypothesis . . . . .	126
7.5.3. Confidence Intervals . . . . .	127
7.6. Multiple Regression . . . . .	128
7.6.1. Eliminating Covariate Effects . . . . .	128
7.6.1.1. Permute the Raw Data . . . . .	129
7.6.1.2. Permute Residuals Under the Full Model . . . . .	129
7.6.1.3. Permute Residuals Under the Reduced Model . . . . .	129
7.6.2. LAD or LSD? . . . . .	130
7.6.3. An Exact Solution . . . . .	130
7.6.4. Testing All Coefficients . . . . .	131

7.7. Single-Case Phase Designs . . . . .	131
7.8. Questions . . . . .	131
8. Clustering in Time and Space . . . . .	134
8.1. The Generalized Quadratic Form . . . . .	134
8.1.1. Mantel's $U$ . . . . .	134
8.1.2. An Example . . . . .	135
8.2. Applications . . . . .	135
8.2.1. The MRPP Statistic . . . . .	136
8.2.2. The BW Statistic of Cliff and Ord [1973] . . . . .	137
8.2.3. Equivalances . . . . .	137
8.2.4. Extensions . . . . .	137
8.2.5. Another Dimension . . . . .	137
8.3. Alternate Approaches . . . . .	138
8.3.1. Quadrant Density . . . . .	138
8.3.2. Nearest-Neighbor Analysis . . . . .	138
8.3.3. Comparing Two Spatial Distributions . . . . .	139
8.4. Questions . . . . .	139
9. Coping with Disaster . . . . .	140
9.1. Missing Data . . . . .	140
9.2. Covariates After the Fact . . . . .	141
9.2.1. Observational Studies . . . . .	142
9.3. Outliers . . . . .	143
9.3.1. Original Data . . . . .	144
9.3.2. Ranks . . . . .	144
9.3.3. Scores . . . . .	145
9.3.4. Robust Transformations . . . . .	146
9.3.5. Use an $L_1$ Test . . . . .	146
9.3.6. Censoring . . . . .	146
9.3.7. Discarding . . . . .	147
9.4. Censored Data . . . . .	147
9.4.1. GAMP Tests . . . . .	147
9.5. Censored Matched Pairs . . . . .	149
9.5.1. GAMP Test . . . . .	149
9.5.2. Ranks . . . . .	150
9.5.3. One-Sample: Bootstrap Estimates . . . . .	151
9.6. Adaptive Tests . . . . .	151
9.7. Questions . . . . .	152
10. Which Statistic? Solving the Insolvable . . . . .	154
10.1. The Permutation Distribution . . . . .	154

10.2. New Statistics . . . . .	154
10.2.1. Nonresponders . . . . .	154
10.2.1.1. Extension to $K$ -Samples . . . . .	155
10.2.2. Animal Movement . . . . .	155
10.2.3. The Building Blocks of Life . . . . .	156
10.2.4. Model Validation . . . . .	157
10.2.5. Structured Exploratory Data Analysis . . . . .	158
10.2.6. Comparing Multiple Methods of Assessment . . . . .	159
10.3. Going Beyond . . . . .	159
10.3.1. Sufficiency . . . . .	162
10.3.2. Invariance . . . . .	163
10.3.3. Applying the Principles . . . . .	163
10.3.4. Losses . . . . .	164
10.4. Likelihood Ratio . . . . .	165
10.4.1. Goodness of Fit and the Restricted Chi-Square . . . . .	166
10.4.2. Censored Data . . . . .	167
10.4.3. Logistic Regression . . . . .	167
10.5. Questions . . . . .	168
11. Which Test Should You Use? . . . . .	170
11.1. Parameters and Parametric Tests . . . . .	170
11.2. Parametric Tests, Permutations, and the Bootstrap . . . . .	171
11.3. What Significance Level Should I Use? . . . . .	172
11.4. A Guide to Selection . . . . .	173
11.4.1. The Data Are in Categories . . . . .	173
11.4.2. The Data Take Discrete Values . . . . .	174
11.4.3. The Data Are Continuous . . . . .	174
11.5. Quick Key . . . . .	176
12. Publishing Your Results . . . . .	179
12.1. Design Methodology . . . . .	179
12.1.1. Randomization in Assignment . . . . .	179
12.1.2. Choosing the Experimental Unit . . . . .	180
12.1.3. Determine Sample Size . . . . .	181
12.2. Statistical Software for Exact Distribution-Free Inference . . . . .	181
12.2.1. Freeware and Shareware . . . . .	181
12.2.2. Commercial Software . . . . .	182
12.3. Preparing Manuscripts for Publication . . . . .	183
13. Increasing Computational Efficiency . . . . .	185
13.1. Seven Techniques . . . . .	185

13.2. Monte Carlo . . . . .	185
13.2.1. Stopping Rules . . . . .	186
13.2.2. Variance of the Result . . . . .	186
13.2.3. Cutting the Computation Time . . . . .	187
13.3. Rapid Enumeration and Selection Algorithms . . . . .	187
13.3.1. Matched Pairs . . . . .	188
13.4. Recursive Relationships . . . . .	189
13.5. Focus on the Tails . . . . .	190
13.5.1. Contingency Tables . . . . .	191
13.5.1.1. Network Representation . . . . .	191
13.5.1.2. The Network Algorithm . . . . .	193
13.6. Gibbs Sampling and a Drunkard's Walk . . . . .	194
13.7. Characteristic Functions . . . . .	195
13.8. Asymptotic Approximations . . . . .	195
13.8.1. A Central Limit Theorem . . . . .	195
13.8.2. Edgeworth Expansions . . . . .	196
13.8.3. Generalized Correlation . . . . .	196
13.9. Confidence Intervals . . . . .	196
13.10. Sample Size and Power . . . . .	197
13.10.1. Simulations . . . . .	197
13.10.2. Network Algorithms . . . . .	198
13.11. Some Conclusions . . . . .	199
13.12. Questions . . . . .	200
14. Theory of Permutation Tests . . . . .	201
14.1. Fundamental Concepts . . . . .	201
14.1.1. Dollars and Decisions . . . . .	201
14.1.2. Tests . . . . .	202
14.1.3. Distribution Functions, Power, Exact, and Unbiased Tests . . . . .	203
14.1.4. Exchangeable Observations . . . . .	203
14.2. Maximizing the Power . . . . .	204
14.2.1. Uniformly Most Powerful Unbiased Tests . . . . .	204
14.2.2. The Fundamental Lemma . . . . .	206
14.2.3. Samples from a Normal Distribution . . . . .	207
14.2.4. Testing the Equality of Variances . . . . .	208
14.2.5. Testing for Bivariate Correlation . . . . .	208
14.3. Confidence Intervals . . . . .	209
14.4. Asymptotic Behavior . . . . .	210
14.4.1. A Theorem on Linear Forms . . . . .	211
14.4.2. Asymptotic Efficiency . . . . .	211
14.4.3. Exchangeability . . . . .	212
14.5. Questions . . . . .	213

Bibliography . . . . .	215
B.1. Permutation Test Articles . . . . .	216
B.2. Computational Methods . . . . .	250
B.3. Seminal Articles . . . . .	255
Author Index . . . . .	257
Subject Index . . . . .	265

## CHAPTER 1

# A Wide Range of Applications

“After all, few statisticians, let alone experimenters in general, will have ready access to an automatic computer.”

M.J.R. Healy, 1952, commenting on a paper by K.D. Tocher.

### 1.1. Permutation Tests

The chief value of permutation tests lies in their wide range of applications.

Permutation tests can be applied to continuous, ordered, and categorical data, and to values that are normal, almost normal, and non-normally distributed.

For almost every parametric and nonparametric test, one may obtain a distribution-free permutation counterpart. The resulting permutation test is usually as or more powerful than alternative approaches. Permutation methods can also sometimes be made to work when other statistical methods fail (see Section 3.4.1 and Chapter 10).

Permutation tests can be applied to homogeneous (textbook) and to heterogeneous (real life) data when subpopulations are mixed together (see Section 10.2.1), when covariates must be taken into account (see Sections 4.3 and 6.6), and when repeated measures on a single subject must be adjusted for (Section 5.6). The ability of permutation methods to be adapted to real-world situations is what led to my writing this book, aimed at the practitioner, as well as upper-division and graduate students in statistics.

#### 1.1.1. Applications

Permutation tests have been applied in

- cluster analysis [Hubert and Levin, 1976];
- Fourier analysis [Friedman and Lane, 1980];
- multivariate analysis [Arnold, 1964; Mielke, 1986];

- single-case analysis [Edgington, 1975A, 1980A,B, 1984; Ferron and Onghena, 1996; Ferron and Ware, 1994, 1995; Kazdin, 1976, but see Kazdin, 1980; McLeod, Taylor, Cohen, and Cullen, 1986; Onghena and Van Damme, 1994].

Its dictionary of applications includes

- agriculture [Eden and Yates, 1933; Higgins and Noble, 1993; Kempthorne, 1952];
- aquatic science [Jackson, 1990; Laroche, Baran, and Rasoanandrasana, 1997; Lorenz and Eiler, 1989; Ponton and Copp, 1997; Quinn, 1987];
- anthropology [Fisher, 1936A; Jorde, Rogers, Bamshad, Watkins, Krakowiak, Sung, Kere, and Harpending, 1997; Konigsberg, 1997; Smith, McDonald, Forster, and Berrington, 1996; Williams-Blangero, 1989];
- archaeology [Berry, Kvamme, and Mielke, 1980, 1983];
- atmospheric science [Adderley, 1961; Gabriel and Hsu, 1983; Mielke, 1979A, 1984; Miller, Shaw, Veitch, and Smith, 1979; Tukey, 1985; Tukey, Brillinger, and Jones, 1978];
- biology [Arnold, Daw, Stenberg, Jayawardene, Srivastava, and Jackson, 1997; Daw, Arnold, Abushullaib, Stenberg, White, Jayawardene, Srivastava, and Jackson, 1998; Howard, 1981];
- biotechnology [Vanlier, 1996];
- botany [Mitchell-Olds, 1986, 1987; Ritland and Ritland, 1989];
- cardiology [Chapelle, Albert, Smeets, Heusghem, and Kulberts, 1982];
- chemistry [vanKeerberghen, Vandenbosch, Smeyers-Verbeke, and Massart, 1991];
- climatology [Robson, Jones, Reed, and Bayliss, 1998];
- clinical trials [Berlin and Ness, 1996; Freedman, 1989; Gail, Tan, and Piantadosi, 1988; Shuster, 1993; Salsburg, 1992; Wei and Lachin, 1988];
- computer science [Yucesan, 1993];
- diagnostic imaging [Arndt, Cizadlo, Andreasen, Heckel, Gold, and Oleary, 1996; Bullmore, Brammer, Williams, Rabehesk, Janot, David, Mekers, Howard, and Slam, 1996];
- ecology [Belyea, 1996; Bersier and Sugihara, 1997; Cade, 1997; Manly, 1983; Meagher and Burdick, 1980; Syrjala, 1996; Saitoh, Stenseth, and Bjornstad, 1997];
- econometrics [Kennedy, 1995; Kim, Nelson, and Startz, 1991; McQueen, 1992];
- education [Manly, 1988];
- endocrinology [O'Sullivan, Whitney, Hinshelwood, and Hauser, 1989];
- entomology [Bryant, 1977; Mackay and Jones, 1989];
- epidemiology [Glass, Mantel, Gunz, and Spears, 1971; Kryscio, Meyers, Prusiner, Heise, and Christine, 1973; Turnbull, Iwano, Burnett, Howe, and Clark, 1990; Wu, Amos, Kemp, Shi, Jiang, Wan, and Spitz, 1998];
- ergonomics [Valdesperez, 1995];
- forensics [Evett, Gill, Scranage, and Weir, 1996; Gastwirth, 1992; Good and Good, 2000; Solomon, 1986; Paternoster, Brame et al., 1998];

- forestry [Magnussen and Boudewyn, 1998];
- genetics [Louis and Dempster, 1987; Levin, 1977; Karlin and Williams, 1984; Rogstad and Pelikan, 1996; Wan, Cohen, and Guerra, 1997; Kidd, Morar et al., 1998];
- geography [Royaltey, Astrachen, and Sokal, 1975];
- geology [Clark, 1989];
- gerontology [Miller, Bookstein, Vandermeulen, Engle, Kim, Mullins, and Faulkner, 1997];
- immunology [Makinodan, Albright, Peter, Good, and Hedrick, 1976; Roper, Doerge, Call, Tung, Hickey, and Teuscher, 1998; Teuscher, Rhein, Livingstone, Paynter, Doerge, Nicholson, and Melvold, 1997];
- library science [Dee, Rankin, and Burns, 1998];
- medicine [Bross, 1964; Feinstein, 1973; McKinney, Young, Hartz, and Bi-Fong Lee, 1989];
- molecular biology [Karlin, Ghandour, Ost, Tauare, and Korph, 1983];
- neurobiology [Edgington and Bland, 1993; Weth, Nadler, and Korschning, 1996];
- neurology [Burgess and Gruzelier, 1997; Diggle, Lange, and Benes, 1991; Faris and Sainsbury, 1990; McIntosh, Bookstein, Haxby, and Grady, 1996; Noble, Gottschalk, Fallon, Ritchie, and Wu, 1997];
- neuropsychopharmacology [Wu, Bell, Najafi, Widmark, Keator, Tang, Klein, Bunney, Fallon, and Bunney, 1997];
- neuropsychology [Stuart, Maruff, and Currie, 1997];
- oncology [Gann, Chatterton, Vogelsong, Dupuis, and Ellman, 1997; Gart, 1986; Hoel and Walburg, 1972; Spitz, Shi, Yang, Hudmon, Jiang, Chamberlain, Amos, Wan, Cinciripini, Hong, and Wu, 1998];
- ornithology [Busby, 1990];
- paleontology [Marcus, 1969; Quinn, 1987; Donnelly and Kramer, 1999];
- pharmacology [Adamson, Hajimohamadenza, Brammer, and Campbell, 1969; Boess, Balasuvramanian, Brammer, and Campbell, 1990; Plackett and Hewlett, 1963];
- physics [Penninckx, Hartmann, Massart, and Smeyersverbeke, 1996];
- physiology [Zempo, Kayama, Kenagy, Lea, and Clowes, 1996];
- psychology [Hubert, 1976, 1978, 1979A,B; Hubert and Baker, 1978; Jennings, McIntosh, Kapur, Tulving, and Houle, 1997; Kelly, 1973; Stilson, 1966];
- reliability [Kalbfleisch and Prentice, 1980; Nelson, 1982; Nelson, 1992];
- sociology [Marascuilo and McSweeney, 1977];
- stochastic processes [Bell, Woodroffe, and Avadhani, 1970; Basawa and Rao, 1980];
- surgery [Majeed, Troy, Nicholl, Smythe, Reed, Stoddard, Peacock, and Johnson, 1996];
- taxonomy [Alroy, 1994; Doolittle, 1981; Gabriel and Sokal, 1969; Klingenberg, Neuenschwander, and Flury, 1996];
- theology [Witztum, Rips, and Rosenberg, 1994];

- toxicology [Cory-Slechta, 1990; Cory-Slechta, Weiss, and Cox, 1989; Farrar and Crump, 1988, 1991];
- virology [Good, 1979];
- vocational guidance [Gliddentracey and Greenwood, 1997; Gliddentracey and Parraga, 1996; Ryan, Tracey, and Rounds, 1996].

Permutation methods are relatively impervious to complications that defeat other statistical techniques. Outliers and “broad tails” may be defended against through the use of preliminary rank or robust transformations (Section 9.3). Missing data are often corrected for automatically. Censored data may affect the power of a permutation test, but not its existence or exactness. A most powerful unbiased permutation test often exists in cases in which a most powerful parametric test fails for lack of knowledge of some yet unknown nuisance parameter (Lehmann, 1986; Good, 1989, 1991, 1992).

A major reason permutation tests have such a wide range of applications is that they require only one or two relatively weak assumptions, e.g., that the underlying distributions are symmetric and the alternatives are simple shifts in value. The permutation test can even be applied to finite populations (see Section 2.4).

Permutation tests have their limitations too, as we shall see and discuss beginning in Section 2.2.7.

## 1.2. “I Lost the Labels”

Shortly after I received my doctorate in statistics, I decided that if I really wanted to help bench scientists apply statistics I ought to become a scientist myself. So back to school I went to learn all about physiology and aging in cells raised in petri dishes.

I soon learned there was a great deal more to an experiment than the random assignment of subjects to treatments. In general, 90% of my effort was spent in mastering various arcane laboratory techniques, 9% in developing new techniques to span the gap between what had been done and what I really wanted to do, and a mere 1% on the experiment itself. But the moment of truth came finally—it had to if I were to publish and not perish—and I succeeded in cloning human diploid fibroblasts in eight culture dishes: Four of these dishes were filled with a conventional nutrient solution and four held an experimental “life-extending” solution to which Vitamin E had been added (see Figure 1.1).

I waited three weeks with my fingers crossed—there is always a risk of contamination with cell cultures—but at the end of this test period three dishes of each type had survived. My technician and I transplanted the cells, let them grow for 24 hours in contact with a radioactive label, and then fixed and stained them before covering them with a photographic emulsion.

Ten days passed and we were ready to examine the autoradiographs. Two years had elapsed since I first envisioned this experiment and now the results were in: I had the six numbers I needed.

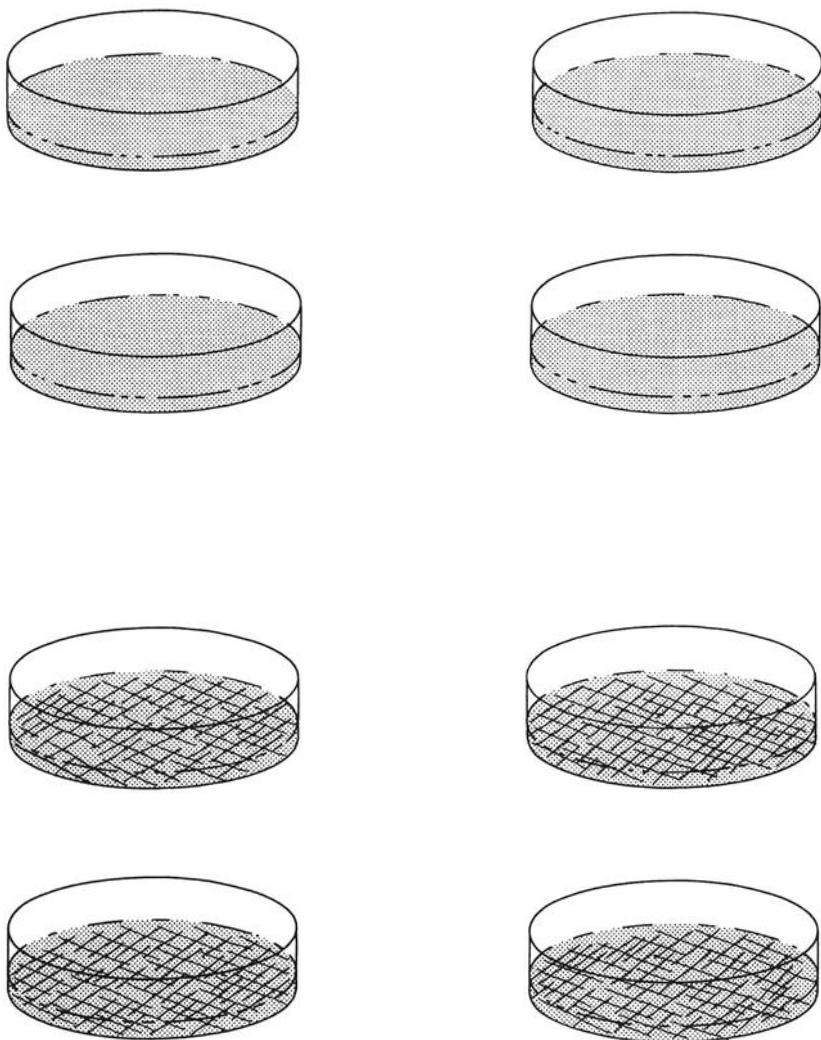


Figure 1.1. Eight petri dishes, 4 containing standard medium, 4 containing standard medium supplemented by Vitamin E. Ten cells inoculated in each dish.

"I've lost the labels," my technician said as he handed me the results.

"What!?" Without the labels, I had no way of knowing which cell cultures had been treated with Vitamin E and which had not.

"121, 118, 110, 34, 12, 22." I read and reread these six numbers over and over again. If the first three counts were from treated colonies and the last three were from untreated, then I had found the fountain of youth. Otherwise, I really had nothing to report.

## 1.3. Five Steps to a Permutation Test

How had I reached that conclusion?

In succeeding chapters, you will learn to apply permutation techniques to a wide variety of testing problems ranging from the simple to the complex. In each case, you will follow the same five-step procedure that we follow in this example.

1. Analyze the problem.
2. Choose a test statistic.
3. Compute the test statistic for the original labeling of the observations.
4. Rearrange (permute) the labels and recompute the test statistic for the rearranged labels. Repeat until you obtain the distribution of the test statistic for all possible permutations.
5. Accept or reject the hypothesis using this permutation distribution as a guide.

### 1.3.1. Analyze the Problem

Let's take a second, more formal look at the problem of the missing labels. First, we identify the hypothesis and alternative of interest.

I wanted to assess the life-extending properties of a new experimental treatment. To do this, I divided my cell cultures into two groups: One grown in a standard medium and one grown in a medium containing Vitamin E. At the conclusion of the experiment and after the elimination of several contaminated cultures, both groups consisted of three independently treated dishes.

My null hypothesis is that the growth potential of a culture will not be affected by the presence of Vitamin E in the media. The alternative of interest is that cells grown in the presence of Vitamin E would be capable of many more cell divisions.

Under the null hypothesis, the labels "treated" and "untreated" provide no information about the outcomes, as the observations are expected to have more or less the same values in each of the two experimental groups. I am free to exchange the labels.

### 1.3.2. Choose a Test Statistic

The next step in the permutation method is to choose a test statistic that discriminates between the hypothesis and the alternative. The statistic I chose was the sum of the counts in the group that had been treated with Vitamin E. If the alternative is true this sum ought to be larger than the sum of the observations in the untreated group. If the null hypothesis is true, that is, if it doesn't make any difference which treatment the cells receive, then the sums of the two groups of observations should be approximately the same. One sum might be smaller or larger than the other by chance, but the two shouldn't be all that different.

### 1.3.3. Compute the Test Statistic

The third step in the permutation method is to compute the test statistic for each of the possible relabelings. But to compute the test statistic for the data as they had been labeled originally, I had to find the labels! Fortunately, I had kept a record of the treatments independent of my technician. In fact, I had deliberately not let my technician know which cultures were which in order to ensure he would give them equal care in handling. As it happened, the first three observations he showed me—121, 118, and 110 were those belonging to the cultures that had received Vitamin E. The value of the test statistic for the observations as originally labelled is  $349 = 121 + 118 + 110$ .

### 1.3.4. Rearrange the Observations

We now rearrange or permute the observations, randomly reassigning the six labels, three “treated” and three “untreated,” to the six observations: for example, treated, 121 118 34, and untreated, 110 12 22. In this rearrangement, the sum of the observations in the first (treated) group is 273. We repeat this step until all  $\binom{6}{3} = \frac{16}{13!(6-3)!} = \frac{6.5.4}{3.2.1} = 20$ , distinct rearrangements have been examined.

	First Group			Second Group			Sum
1.	121	118	110	34	22	12	349
2.	121	118	34	110	22	12	273
3.	121	110	34	118	22	12	265
4.	118	110	34	121	22	12	262
5.	121	118	22	110	34	12	261
6.	121	110	22	118	34	12	253
7.	121	118	12	110	34	22	251
8.	118	110	22	121	34	12	250
9.	121	110	12	118	34	22	243
10.	118	110	12	121	34	22	240
11.	121	34	22	118	110	12	177
12.	118	34	22	121	110	12	174
13.	121	34	12	118	110	22	167
14.	110	34	22	121	118	12	166
15.	118	34	12	121	110	22	164
16.	110	34	12	121	118	22	156
17.	121	22	12	118	110	34	155
18.	118	22	12	121	110	34	152
19.	110	22	12	121	118	34	144
20.	34	22	12	121	118	110	68

**Five Steps to a Permutation Test**

- 1) Analyze the problem.
  - a) What is the hypothesis? What are the alternatives?
  - b) How are the data distributed?
  - c) What losses are associated with bad decisions?
- 2) Choose a statistic and establish a rejection rule that will distinguish the hypothesis from the alternative.
- 3) Compute the test statistic for the original observations.
- 4) Rearrange the observations.
  - a) Compute the test statistic for the new arrangement.
  - b) Compare the new value of test statistic with the value you obtained for the original observations.
  - c) Repeat steps a) and b) until you are ready to make a decision.
- 5) Make a decision.

Reject the hypothesis and accept the alternative if the value of the test statistic for the observations as they were labeled originally is an extreme value with respect to the rejection rule in the permutation distribution of the statistic. Otherwise, accept the hypothesis and reject the alternative.

### 1.3.5. Make a Decision

The sum of the observations in the original Vitamin E treated group, 349, is equaled only once and never exceeded in the 20, distinct random relabelings. If chance alone is operating, then such an extreme value is a rare, only-1-time-in-20 event. I reject the null hypothesis at the 5% (1 in 20) significance level and embrace the alternative that the treatment is effective and responsible for the difference I observed.

In using this decision procedure, I risk making an error and rejecting the null hypothesis when it is true 1 in 20 times. In this case, I made just such an error. I was never able to replicate the observed life-promoting properties of Vitamin E in other repetitions of this experiment. Good statistical methods can reduce and contain the probability of making a bad decision, but they cannot eliminate the possibility.

## 1.4. What's in a Name?

Permutation tests are also known as randomization, rerandomization, and exact tests. Historically, one may distinguish between Fisher's notion of the randomization test applicable only to the samples at hand<sup>1</sup> and Pitman's idea of a permutation test that could be applied inductively to the larger populations from which the samples are drawn,<sup>2</sup> but few research workers honor this distinction today. Gabriel

<sup>1</sup> Fisher [1966, Chapter 3, (first edition 1935)].

<sup>2</sup> Pitman [1937a,b; 1938].

and Hall [1983] use the term “rerandomization” to distinguish between the initial randomization of treatment assignments at the design phase and the subsequent “rerandomizations” that occur during the permutation analysis. When analyzing data, some workers confine themselves to randomizations that might have arisen at the design stage; others make no such distinction.<sup>3</sup> In this book, we shall use the three terms “permutation,” “randomization,” and “rerandomization” interchangeably.

Most permutation tests provide “exact” significance levels. We define “exact,” “significance level,” and other important concepts in Section 2.2 and establish the conditions under which permutation tests are exact and unbiased. We reserve the name “exact test” for the classic Fisher’s test for  $2 \times 2$  tables, studying this test and other permutation tests applied to categorical data in Chapter 6.

The terms “distribution free” and “nonparametric” often arise in connection with the permutation tests. “Distribution free” means that the significance level of the test is independent of the form of the hypothetical infinite population from which the sample is drawn.<sup>4</sup> Permutation tests are almost but not quite “distribution free” in this sense in that only one or two assumptions about the underlying population(s) are required for their application. A preliminary rank transformation often can ensure that tests are distribution free. Bell and Doksum [1967] prove that all distribution-free tests of independence are permutation tests.

“Nonparametric” means that the parametric form of the underlying population distribution is not specified explicitly. It is probably safe to say that 99% of permutation tests are nonparametric and that 99% of common nonparametric tests are permutation tests in which the original observations have been replaced by ranks.

#### 1.4.1. Comparison with Other Tests

When the samples are very large, decisions based on parametric tests like the  $t$ -test and the  $F$  usually agree with decisions based on the corresponding permutation test. With small samples, the parametric test will be preferable IF the assumptions of the parametric test are satisfied completely. The familiar “rank” tests are simply permutation tests applied to the ranks of the observations rather than their original values (see Sections 9.3 and 11.2).

#### 1.4.2. Sampling from the Data at Hand

The two resampling methods—the permutation tests and the bootstrap—have much in common. Both are computer intensive and limited to the data at hand.

<sup>3</sup> Welch [1990] and Romano [1990].

<sup>4</sup> The power of the test will depend on the underlying distribution.

With the permutation test, you recompute the test statistic for all possible relabelings of the combined samples. If the original samples contained the observations 1, 2, 4 and 3, 5, 6 you would consider the relabelings 1, 2, 3 and 4, 5, 6; 1, 2, 5 and 3, 4, 6; and so forth.

With the bootstrap, you recompute the test statistic for each of a series of samples with replacement. These may be taken from the pooled sample or drawn separately from each sample, obtaining in the former case, for example, the bootstrap samples 1, 3, 3 and 3, 4, 6.

For some testing situations and test statistics, the bootstrap and the randomization test are asymptotically equivalent [Romano, 1989; Robinson, 1987], but often they yield quite different results, a point we make at length in Sections 7.2 and 11.2.

When you analyze an experiment or survey with a parametric test—Student's  $t$ , for example—you compare the observed value of the test statistic with the values in a table of its theoretical distribution, for example, in a table of Student's  $t$  with eight degrees of freedom. Analyzing the same experiment with a permutation test, you compare the observed value of the test statistic with the set of what-if values you obtain by rearranging and relabeling the data.

In view of all the necessary computations—the test statistic must be recomputed for each what-if scenario—it is not surprising that the permutation test's revival in popularity parallels the increased availability of high-speed computers. Although the permutation test was introduced by Fisher and Pitman in the 1930's, it represented initially a theoretical standard rather than a practical approach. But with each new quantum leap in computer speed, the permutation test was applied to a wider and wider variety of problems. In earlier eras—the 1950's, the 1960's, and the 1970s—the permutation test's proponents, enthusiastic at first, would grow discouraged as, inevitably, the number of computations proved too demanding for even the largest of the then-available computing machines. But with today's new and more powerful generation of desktops, it is often faster to compute a  $p$ -value for an exact permutation test than to look up an asymptotic approximation in a book of tables.

With both the bootstrap and the permutation test, all significance levels are computed on the fly. The statistician is not limited by the availability of tables, but is free to choose a test statistic exactly matched to hypothesis and alternative [Bradley, 1968].

## 1.5. History

Dependent on a computer technology that was itself under development, the history of permutation tests has been one of constant rediscovery; see, for example, Kempthorne [1952], Bradley [1968], as well as the first, 1980 edition of Edgington [1995].

World War II provided impetus for developing a theoretical basis for parametric procedures that would "serve" in place of the correct but computationally demanding permutations. The theory and certainly the practice of permutation methods

was largely neglected with a few major exceptions, such as Wald and Wolfowitz [1944] and Lehmann and Stein [1949].

In the 1950's and early 1960's, workers in the rapidly burgeoning field of experimental design took pains to demonstrate that their parametric tests were asymptotically equivalent to the correct permutation procedures, as in Kempthorne [1952], Freeman and Halton [1951], Cornfield and Tukey [1956], Barton and David [1961], and Arnold [1964].

Tests based on the permutation of ranks had the advantage that the permutation distribution need be computed only once for each design and combination of sample sizes and thus could be tabulated. Rank tests dominated the nonparametric literature of the 1960's and 1970's; see, for example, Siegel [1956], Hedges and Lehmann [1963], Cox and Kempthorne [1963], Sen [1965, 1967], Bell and Doksum [1965, 1967], Bell and Donoghue [1969], Shane and Puri [1969], Bickel [1969], Puri [1970], Puri and Sen [1971], and Lehmann [1975].

By the middle of the 1970's, most statisticians had a mainframe terminal on their desk and conditions were ripe for the development of a wide variety of resampling procedures, including the bootstrap and density estimation as well as permutation tests; see McCarthy [1969], Hartigan [1969], Efron [1979], Diaconis and Efron [1983], Izenman [1991], and Good [1999].

The late 1960's and 1970's saw the introduction of cluster analysis (Mantel [1967]; Cliff and Ord [1971]; Mielke, Berry, and Johnson [1976]; Mielke [1979]; Mielke, Berry, Brockwell, and Williams [1981]; see also Chapter 8). The late 1970's and early 1980's saw breakthroughs in the analysis of single-case designs (Kazdin, 1976; Edgington, 1980a,b) and directional data (Hubert et al., 1984).

In the 1980's, there were plenty of new technologies to take advantage of the newly enabled permutation tests, as demonstrated in the work of Mann and Hand [1983], who analyzed flow cytometric histograms, and Karlin et al. [1983], who performed computer analysis of DNA sequences.

In the 1950's, it made sense to look for large-sample approximations to the permutation distribution; see, for example, Noether [1950], Hoeffding [1951], Erdos and Renyi [1959], Rosen [1965], Hajek [1968], and Cléoux [1969]. By the 1980's, this effort was questionable because of the development of new and more efficient computational algorithms [Feldman and Kluger, C1963;<sup>5</sup> Rogers, C1964; Sag, C1964; Boothroyd, C1967; Chase, C1970a,b; Minc, C1971; Liu and Tang, C1973; Gentleman, C1975; Bitner et al., C1976; Gail and Mantel, C1975; Ives, C1976; DeCani, C1979; Mehta and Patel, C1980, C1983, C1986; Akl, C1981; Pagano and Halvorsen, 1981; Patefield, C1981; Pagano and Tritchler, C1983; Vitter, C1984, Baglivo et al., C1988; Oden, C1991; and Spino and Pagano, C1991a,b; see also Chapter 13].

Still, the work on asymptotics continued; see, for example, Shapiro and Hubert [1979], Robinson [1980], Fang [1981], Ascher and Bailar [1982], Constanzo et al. [1983], Denker and Puri [1988], and Donegani [1991].

<sup>5</sup> The C in front of a year, as in C1973, refers to a separate bibliography at the end of the text devoted exclusively to computational procedures.

The 1990's saw advances in the testing of simultaneous hypotheses [Westfall and Young, 1993; Troendle, 1995; and Blair et al., 1996] and sequential analysis [Lefebvre, 1982; Lin et al., 1991]. The newest technologies again saw the application of this oldest of statistical methodologies; examples include Bullmore et al. [1996; functional MR image analysis], Jorde et al. [1997; demographics and satellite imagery].

## 1.6. To Learn More

For excellent introductions to the application of permutation tests, see Kempthorne [1952], Feinstein [1973], Howard [1981], Barbellla et al. [1990], and Good [1999]. For a comparison with parametric tests, see Bradbury [1987]; for a comparison with the bootstrap, see Romano [1989], ter Braak [1992], Good, [1999], and Chapter 11. For an introduction to the theory, see Lehmann [1986], Welch [1990], Romano [1990], and Chapter 14.

## 1.7. Questions

Take the time to think about the answers to these questions, even if you don't answer them explicitly. You may wish to return to them after you've read subsequent chapters.

1. In the simple example analyzed in this chapter, what would the result have been if you had used as your test statistic the difference between the sums of the first and second samples? the difference between their means? the sum of the squares of the observations in the first sample? the sum of their ranks?
2. How was the analysis of my experiment affected by the loss of two of the cultures due to contamination? Suppose these cultures had escaped contamination and given rise to the observations 90 and 95; what would be the results of a permutation analysis applied to the new, enlarged data set consisting of the following cell counts:

Treated	121	118	110	90
Untreated	95	34	22	12

3. Read one or two of the articles in your research area that were cited on pages 1 and 2. What led the authors to use a permutation test rather than some other statistical procedure?

## CHAPTER 2

# A Simple Test

“Actually, the statistician does not carry out this very tedious process but his conclusions have no justification beyond the fact they could have been arrived at by this very elementary method.”

R.A. Fisher [1936a]

“Tests of significance in the randomized experiment have frequently been presented by way of normal law theory, whereas their validity stems from randomizaton theory.”  
O. Kempthorne [1955]

## 2.1. Properties of the Test

In this chapter, we consider the assumptions that underlie the permutation test and take a look at some of the permutation test’s formal properties—its significance level, power, and robustness. This first look is relatively nonmathematical in nature. A formal derivation is provided in Chapter 14.

In the example of the missing labels in the preceding chapter, we introduced a statistical test based on the random assignment of labels to treatments. We showed this test provided a significance level of 5%, an *exact* significance level, not an approximation. The test we derived is valid under very broad assumptions. The data could have been drawn from a normal distribution or they could have come from some quite different distribution. All that is required for our permutation test comparing samples from two populations to be valid is that under the null hypothesis the distribution from which the data in the treatment group are drawn be the same as that from which the untreated sample is taken.

This freedom from reliance on numerous assumptions is a big plus. The fewer the assumptions, the fewer the limitations, and the broader the potential applications of a test. But before statisticians introduce a test into their practice, they need to know a few more things about it.

How powerful a test is it? That is, how likely is it to pick up actual differences between treated and untreated populations? Is this test as powerful or more powerful than the test we are using currently?

How robust is the new test? That is, how sensitive is it to violations in the underlying assumptions and the conditions of the experiment?

What if data are missing, as they are in so many of the practical experiments we perform? Will missing data affect the significance level of our test?

What are the effects of extreme values or outliers? In an experiment with only five or six observations, it is obvious that a single extreme value can mislead the experimenter. In Section 9.3 of this text, you will learn techniques for diminishing the effect of extreme values.

Can we extend our results to complex experimental designs in which there are several treatments at several different levels and several simultaneous observations on each subject?

The answer to this last question, as the balance of this book will reveal to you, is yes. For example, you can easily apply permutation methods to studies in which you test a single factor at three or four levels simultaneously (see Chapter 3, Section 3.5). You can also apply permutation methods to experimental designs in which you control and observe the values of multiple variables (Chapters 4 and 5).

The balance of this chapter is devoted to providing a theoretical basis for all the preceding questions and answers.

## 2.2. Fundamental Concepts

Why do we elect to use one statistical procedure rather than another—a permutation test, say, as opposed to a table of chi-square? If you've just completed a course in statistics, you probably already know the answer. If it's been a year or so since you last looked at a statistics text, then you will find this section helpful.

In this section, you are introduced in an informal way to the fundamental concepts of variation, population and sample distributions, Type I and Type II error, significance level, power, and exact and unbiased tests. Formal definitions and derivations are provided in Chapter 14.

### 2.2.1. Population and Sample Distributions

The two factors that distinguish the statistical from the deterministic approach are variation and the possibility of error. The effect of this variation is that a distribution of values takes the place of a single, unique outcome.

I found freshman physics extremely satisfying: Boyle's Law, for example,  $V = KT/P$ , with its tidy relationship between the volume, temperature, and pressure of a perfect gas. The problem was I could never quite duplicate this law in the freshman physics laboratory. Maybe it was the measuring instruments, my lack of familiarity with the equipment, or simple measurement error—but I kept getting different values for the constant  $K$ .

By now, I know that variation is the norm—particularly in the clinical and biological areas. Instead of getting a fixed, reproducible  $V$  to correspond to a

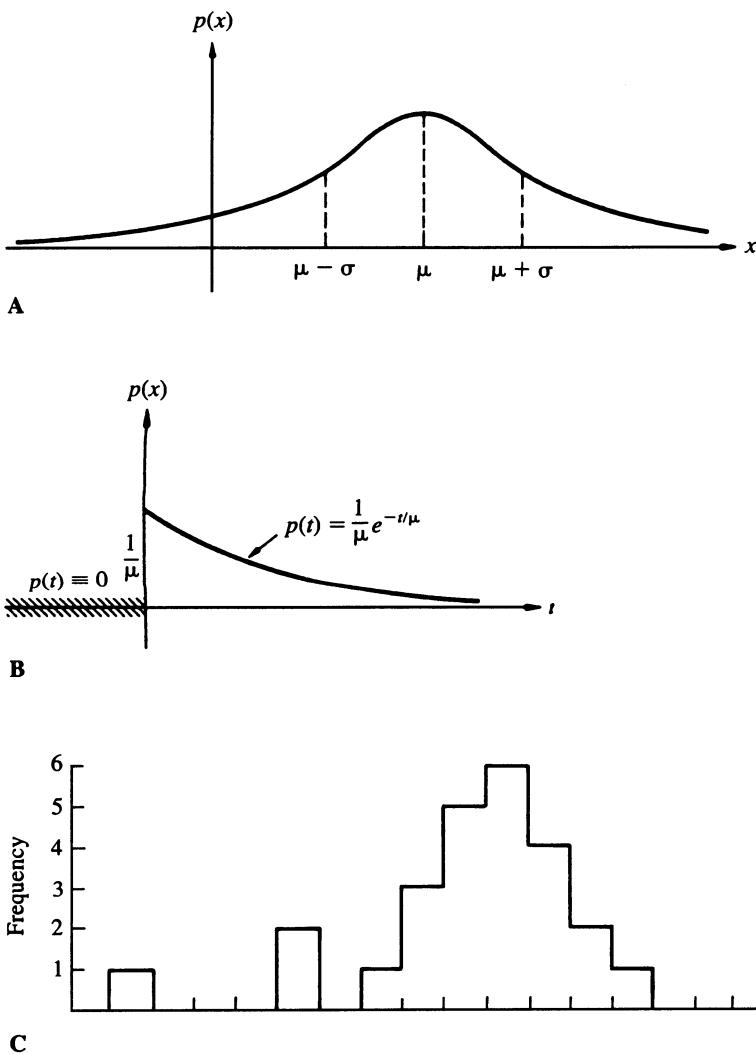


Figure 2.1. Distributions: a) normal distribution, b) exponential distribution, c) distribution of values in a sample taken from a normal distribution.

specific  $T$  and  $P$ , one ends up with a distribution of values instead. But I also know that, with a large enough sample, the mean and shape of this distribution are reproducible.

Figures 2.1a and 2.1b depict two such distributions. The first is a normal distribution. Examining the distribution curve, we see that the normally distributed variable can take all possible values between  $-\infty$  and  $+\infty$ , but most of the time it takes values that are close to its median (and mean)  $\mu$ . The second is an exponential distribution; the exponentially distributed variable only takes positive values; the

majority of the time these values are less than its mean, but on occasion they can be many times larger.

Both of these distributions are limiting cases; they represent the aggregate result of an infinite number of observations; thus, the distribution curves are smooth. The choppy histogram in Figure 2.1c is typical of what one sees with a small, finite sample of observations—in this case, a sample of 25 observations taken from a normal distribution with mean  $\mu$ .

Still, the sample is not the population. To take an extreme, almost offensive example, suppose that every member of a Los Angeles, CA, jury were to be non-white. Could that jury really have been selected at random from the population as the statutes of the State of California requires? Yes; there are court districts in Los Angeles in which fewer than 30% of the population is white; the probability of a jury of 12 individuals containing no whites is approximately 0.7 raised to the 12th power or about 1%. With hundreds of juries being impaneled each week, nonwhite juries are not uncommon; nonetheless, they are not representative.

The good news is that as a sample grows larger, it will more and more closely resemble the population from which it is drawn. How large is large? The answer depends both on the underlying distribution—is it exponential? normal? a mixture of normals?—and the population parameters we wish to approximate.

### 2.2.2. Two Types of Error

It's usually fairly easy to reason from cause to effect—that is, if you have a powerful enough computer. Get the right formula, Boyle's Law, say, plug in enough values to enough decimal places, and out pops the answer. The difficulty with reasoning in the opposite direction, from effect to cause, is that more than one set of causes can be responsible for precisely the same set of effects. We can never be completely sure which set of causes is responsible. Consider the relationship between sex (cause) and height (effect). Boys are taller than girls. Right? So that makes this new 6'2" person in our lives . . . a starter on the women's volleyball team.

In real life, in real populations, there are vast differences from person to person. Some women are tall, and some women are short. In Lake Wobegon MN, all the men are good looking and all the children are above average. But in most other places in the world, there is a wide range of talent and abilities. As a further example of this variation, consider that half an aspirin will usually take care of one of my headaches, while other people can and do take two or three aspirins at a time and get only minimal relief.

Figure 2.2 depicts the results of an experiment in which two groups were each given a “painkiller.” The first group got buffered aspirin, the second group received a new experimental drug. Each of the participants then provided a subjective rating of the effects of the drug. The ratings ranged from “got worse” to “much improved,” depicted on a scale of 0 to 4. Take a close look at Figure 2.2. Does the new drug represent an improvement over aspirin?

Those who took the new experimental drug do seem to have done better on the average than those who took aspirin. Or are the differences we observe in

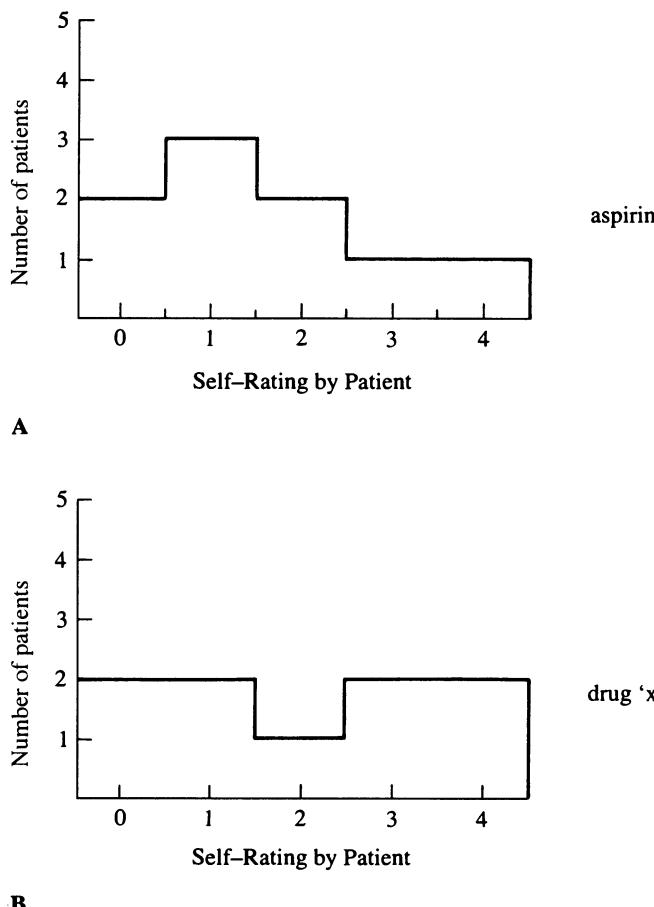


Figure 2.2. Response to treatment: Self-rating by patient, a) aspirin-treated group, b) drug 'x'-treated group.

Figure 2.2 simply the result of chance? If it's just a chance effect and we opt in favor of the new drug, we've made an error. We also make an error if we decide there is no difference and the new drug really is better. These decisions and the effects of making them are summarized in Table 2.1a.

We distinguish the two types of error because they have quite different implications. For example, Fears, Tarone, and Chu [1977] use permutation methods to assess several standard screens for carcinogenicity. Their Type I error, a false positive, consists of labeling a relatively innocuous compound as carcinogenic. Such an action means economic loss for the manufacturer and the denial of the compound's benefits to the public. Neither consequence is desirable. But a false negative, a Type II error, would mean exposing a large number of people to a potentially lethal compound.

Because variation is inherent in nature, we are bound to make the occasional error when we draw inferences from experiments and surveys particularly if, for

Table 2.1a. Decision Making Under Uncertainty

Our Decision		
The Facts	No Difference	Drug is Better
No Difference		Type I error Manufacturer wastes money developing ineffective drug.
Drug is Better	Type II error Manufacturer misses opportunity for profit. Public denied access to effective treatment.	

Table 2.1b. Decision Making Under Uncertainty

Fears, et al.'s Decision		
The Facts	Not a Carcinogen	Compound a Carcinogen
No Effect		Type I error Manufacturer misses opportunity for profit. Public denied access to effective treatment.
Carcinogen	Type II error Patients die; families suffer; Manufacturer sued.	

example, chance hands us a completely unrepresentative sample. When I toss a coin in the air six times, I can get three heads and three tails, but I can also get six heads. This latter event is less probable, but it is not impossible. Does the best team always win?

We can't eliminate the risk in making decisions, but we can contain it by the correct choice of statistical procedure. For example, we can require that the probability of making a Type I error not exceed 5% (or 1% or 10%) and restrict our choice to statistical methods that ensure we do not exceed this level. If we have a choice of several statistical procedures, all of which restrict the Type I error appropriately, we can choose the method which leads to the smallest probability of making a Type II error.

#### 2.2.2.1. Losses and Risk

The preceding discussion is greatly oversimplified. Obviously, our losses will depend not merely on whether we guess right or wrong, but also on how far our guesstimate is off the mark. For example, you've developed a new drug to relieve

anxiety and are investigating its side effects. Does it raise blood pressure? You do a study and find the answer is no. But the truth is your drug raises systolic blood pressure an average of 1 mm. What is the cost to the average patient? Negligible, a bare fraction of the day to day variation.

Now, suppose your new drug actually raises blood pressure an average of 10 mm. What is the cost to the average patient? to the entire potential patient population? to your company in law suits? Clearly, the cost of a Type II error will depend on the magnitude of that error and the nature of the losses associated with it.

Historically, much of the work in testing hypotheses has been limited to 0, 1 loss functions, while that of estimation has focused on losses proportional to the square of the error. The result may have been statistics that were suboptimal in nature with respect to the true, underlying loss [Mielke, 1986; Mielke and Berry, 1997].

Are we more concerned with the losses associated with a specific decision or those we will sustain over time as a result of adhering to a specific decision procedure? Which concerns us and our company the most, reducing average losses over time or avoiding even the remote possibility of a single, catastrophic loss? We return to this topic in Section 14.1.

### 2.2.3. Significance Level and Power

In selecting a statistical method, statisticians work with two closely related concepts, significance level and power. The *significance level* of a test, denoted throughout the text by the Greek letter  $\alpha$  (alpha), is the probability of making a Type I error; that is,  $\alpha$  is the probability of deciding erroneously on the alternative when the hypothesis is true.

To test an hypothesis, we divide the set of possible outcomes into two or more regions. We reject the hypothesis and risk a Type I error when our test statistic lies in the *rejection region*  $R$ ; we accept the hypothesis and risk a Type II error when our test statistic lies in the *acceptance region*  $A$ ; and we may take additional observations when our test statistic lies in the boundary region of *indifference*  $I$ . If  $H$  denotes the hypothesis, then

$$\alpha = \Pr\{A \mid H\}.$$

The *power* of a test, denoted throughout the text by the Greek letter  $\beta$  (beta), is the complement of the probability of making a Type II error; that is,  $\beta$  is the probability of deciding on the alternative when the alternative is the correct choice. If  $K$  denotes the alternative, then

$$\beta = \Pr\{R \mid K\}.$$

The relationship among power, significance level, and effect size for a specific test is summarized in Figure 2.3, provided through the courtesy of Patrick Onghena. For a fixed significance level, the power of a two-tailed test is an increasing function of the absolute effect size. For a fixed effect size, increasing  $\alpha$ , the probability of a Type I error, also increases the power.

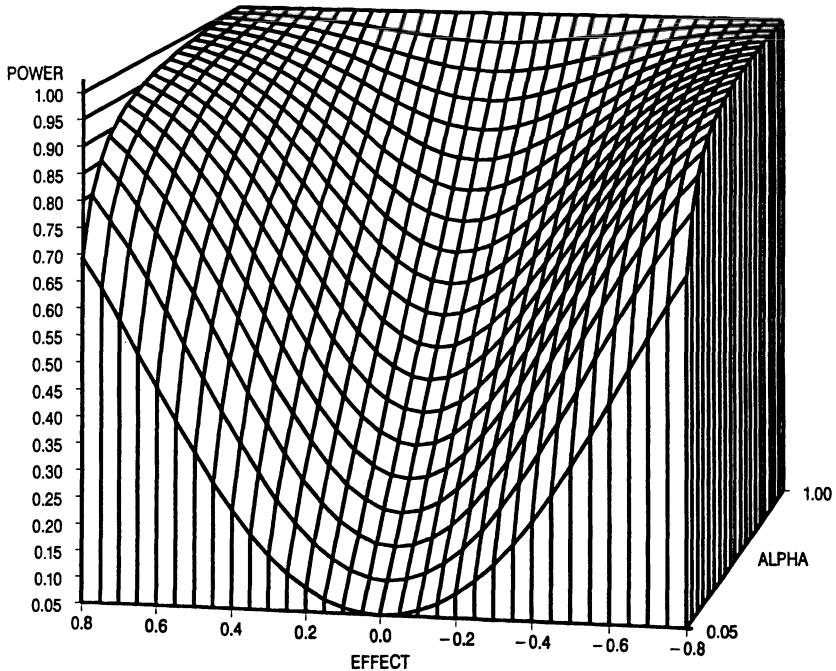


Figure 2.3. Power of the two-tailed  $t$ -test with sample sizes of  $n_1 = n_2 = 20$  as a function of the effect size (EFFECT) and the significance level (ALPHA) under the classical parametric assumptions.

The ideal statistical test would have a significance level  $\alpha$  of zero or 0% and a power  $\beta$  of 1 or 100%. But unless we are all-knowing, this ideal can not be realized. In practice, we fix a significance level  $\alpha > 0$ , where  $\alpha$  is the largest value we feel comfortable with, and choose a statistic that maximizes or comes closest to maximizing the power.

#### 2.2.4. Power and Sample Size

As we saw in Section 2.2.2.1, the greater the discrepancy between the true alternative and our hypothesis, the greater the loss associated with a Type II error. Fortunately, in most practical situations, we can devise a test in which the larger the discrepancy, the greater the power and the less likely we are to make a Type II error.

The relationship among power, effect size, and the number of observations for a specific test is summarized in Figure 2.4, provided through the courtesy of Patrick Onghena.

Figure 2.5 depicts the power as a function of the alternative (effect size) for two tests based on samples of size 6. In the example illustrated, the test  $\phi_1$  is uniformly

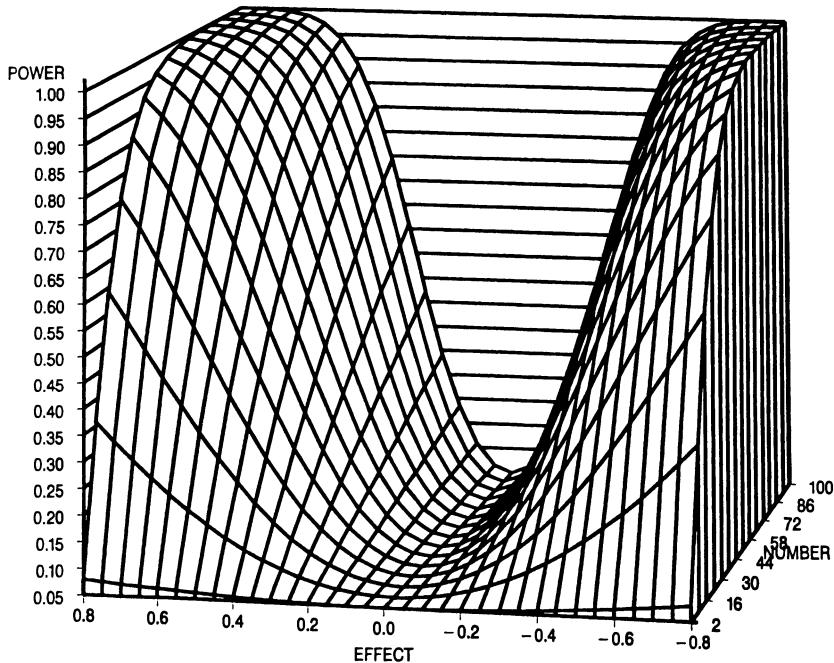


Figure 2.4. Power of the two-tailed  $t$ -test with  $p = 0.05$  as a function of the effect size (EFFECT) and the number of observations (NUMBER,  $n_1 = n_2$ ) under the classical parametric assumptions.

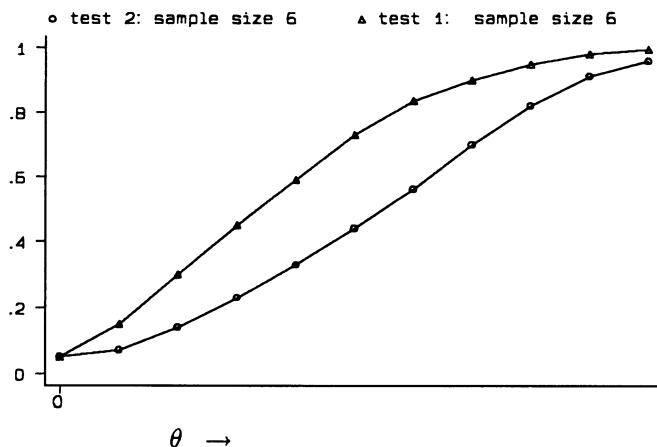


Figure 2.5. Power as a function of the alternative. Both tests have the same sample size.

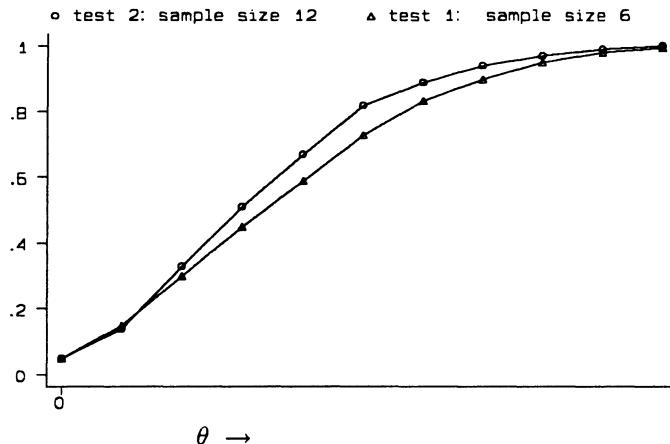


Figure 2.6. Power as a function of the alternative. Tests have different sample sizes.

*more powerful* than  $\phi_2$ ; hence, using  $\phi_1$  in preference to  $\phi_2$  will expose us to less risk.

Figure 2.6 depicts the power curve of these same two tests but using different size samples; the power curve of  $\phi_1$  is still based on a sample of size 6, but that of  $\phi_2$  is now based on a sample of size 12. The two new power curves almost coincide, revealing that the two tests now have equal risks, but we will have to pay for twice as many observations if we use the second test in place of the first.

Moral: *A more powerful test reduces the costs of experimentation along with minimizing the risk.*<sup>1</sup>

### 2.2.5. Power and the Alternative

If a test at a specific significance level  $\alpha$  is more powerful against a specific alternative than all other tests at the same significance level, we term it *most powerful*. But as we see in Figure 2.7, a test that is most powerful for some alternatives may be less powerful for others. When a test at a specific significance level is more powerful against all alternatives than all other tests at the same significance level, we term this test *uniformly* most powerful.

The significance level and power may also depend on how the variables we observe are distributed. Does the population distribution follow a bell-shaped normal curve with the most frequent values in the center as in Figure 2.1a? or is the distribution something quite different? To protect our interests, we may need to require that the Type I error be less than or equal to some predetermined value for all possible distributions. When applied correctly, permutation tests always have

<sup>1</sup> The exception proves the rule. When data gathering is dirt cheap, Lloyd Nelson observes, a less powerful test such as the sign test makes economic sense.

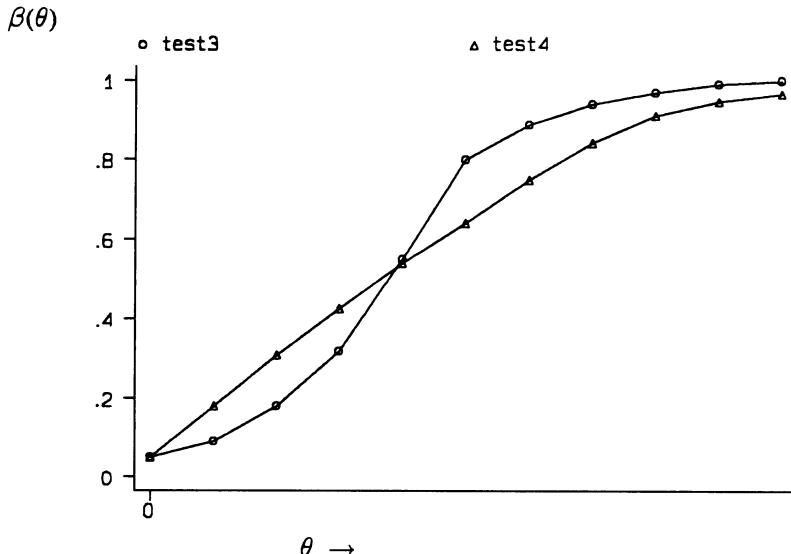


Figure 2.7. Comparing power curves: For near alternatives, with  $\theta$  close to zero, test 4 is the more powerful test; for far alternatives, with  $\theta$  large, test 3 is more powerful. Thus, neither test is uniformly most powerful.

this property. The significance levels of parametric tests and of tests based on the bootstrap are dependent on the underlying distribution.

The power of a test depends on the statistic, the sample size, and the alternative.

### 2.2.6. Exact, Unbiased Tests

In practice, we seldom know the distribution of a variable or its variance. We usually want to test a *compound* hypothesis, such as  $H : X$  has mean 0. This latter hypothesis includes several *simple* hypotheses, such as  $H_1 : X$  is normal with mean 0 and variance 1;  $H_2 : X$  is normal with mean 0 and variance 1.2; and  $H_3 : X$  has a gamma distribution with mean 0 and four degrees of freedom.

A test is said to be *exact* with respect to a compound hypothesis if the probability of making a Type I error is exactly  $\alpha$  for each and every one of the possibilities that make up the hypothesis. A test is said to be *conservative* if the Type I error never exceeds  $\alpha$ . Obviously, an exact test is conservative, though the reverse may not be true.

The importance of an exact test cannot be overestimated, particularly a test that is exact regardless of the underlying distribution. If a test that is nominally at level  $\alpha$  is actually at level  $\chi$ , we may be in trouble before we start: If  $\chi > \alpha$ , the risk of a Type I error is greater than we are willing to bear. If  $\chi < \alpha$ , then our test is

suboptimal, and we can improve on it by enlarging its rejection region. We return to these points again in Chapter 11, on choosing a statistical method.

A test is said to be *unbiased* and of level  $\alpha$ , providing its power function  $\beta$  satisfies the following two conditions:

$\beta$  is conservative; that is,  $\beta_\theta \leq \alpha$  for every  $\theta$  that satisfies the hypothesis; and  
 $\beta_\theta \geq \alpha$  for every  $\theta$  that is an alternative to the hypothesis.

That is, a test is unbiased if using the test you are more likely to reject a false hypothesis than a true one. I find unbiasedness to be a natural and desirable principle, but not everyone shares this view, see, for example, Suissa and Shuster [1984].

Faced with some new experimental situation, our objective is always to derive a uniformly most powerful unbiased test if one exists. But, if we can't derive a uniformly most powerful test, and Figure 2.7 depicts just such a situation, then we will look for a test that is most powerful against those alternatives that are of immediate interest.

### 2.2.7. Exchangeable Observations

A sufficient condition for a permutation test to be exact and unbiased against shifts in the direction of higher values is the *exchangeability* of the observations in the combined sample. The observations  $\{X, Y, \dots, Z\}$  are exchangeable if the probability of any particular joint outcome,  $X + Y + Z = 6$ , for example, is the same regardless of the order in which the observations are considered [Lehmann, 1986, p. 231]. Section 14.1 provides a formal derivation of this fundamental result [see also, Galambos, 1986; Draper, et al., 1993].

Independent, identically distributed observations are exchangeable. So are samples without replacement from a finite population, such as Polya urn models [Koch, 1982]. So are dependent, normally distributed random variables  $X_i$  for which the variance is a constant independent of  $i$  and the covariance of  $X_i$  and  $X_j$  is a constant independent of  $i$  and  $j$ .

When nuisance parameters are involved, variables may not be exchangeable, except with extremely large samples. Practical examples arise when comparing variances when the means are unknown (Section 3.4), testing for interaction in an experimental design when main effects are present (Section 4.2.3), and testing in the presence of covariates (Section 7.6).

Sometimes a simple transformation will ensure that observations are exchangeable. For example, if we know that  $X$  comes from a population with mean  $\mu$  and distribution  $F(x - \mu)$  and an independent observation  $Y$  comes from a population with mean  $\eta$  and distribution  $F(x - \eta)$ , then the independent variables  $X' = X - \mu$  and  $Y' = Y - \eta$  are exchangeable.

In deciding whether your own observations are exchangeable, and a permutation test is applicable, the key question is the one we posed in Chapter 1, Section 1.2:

*Under the null hypothesis of no differences among the various experimental or survey groups, can we exchange the labels on the observations without affecting the results?*

The effect of a “no” answer to this question is discussed in Chapter 9 along with practical guidelines for the design and conduct of experiments and surveys to ensure that the answer is “yes.”

## 2.3. Which Test?

We are now able to make an initial comparison of the four types of statistical tests—permutation, rank, bootstrap, and parametric.

Recall from Chapter 1 that, with a permutation test, we:

1. Choose a test statistic  $S(X)$ .
2. Compute  $S$  for the original set of observations.
3. Obtain the permutation distribution of  $S$  by repeatedly rearranging the observations at random. With two or more samples, we combine all the observations into a single large sample before we rearrange them.
4. Obtain the upper  $\alpha$ -percentage point of the permutation distribution and accept or reject the null hypothesis according to whether  $S$  for the original observations is smaller or larger than this value.

If the observations are exchangeable under the null hypothesis, then the resultant test is exact and unbiased.

As noted in this chapter’s opening quotation from Fisher, although permutation tests were among the very first statistical tests to be developed, they were beyond the computing capacities of the 1930’s. One alternative, which substantially reduces the amount of computation required, is the rank test. To form a rank test (e.g., Mann–Whitney or Friedman’s test), we:

1. Choose a test statistic  $S$ .
2. Replace the original observations  $\{X_{ij}, i = 1, \dots, I, j = 1, \dots, J\}$  by their ranks in the combined sample  $\{R_k, k = 1, \dots, IJ\}$ . As an example, if the original observations are 5.2, 1, and 7, their ranks are 2, 1, and 3. Compute  $S$  for the original set of ranks.
3. Obtain the permutation distribution of  $S$  by repeatedly rearranging the ranks at random and recomputing the test statistic. Or, since ranks always take the same values 1, 2, and so forth, take advantage of a previously tabulated distribution.
4. Accept or reject the hypothesis in accordance with the upper  $\alpha$ -percentage point of this permutation distribution.

In short, a rank test is simply a permutation test applied to the ranks of the observations rather than their original values. If the observations are exchangeable, then the resultant rank test is exact and unbiased. Generally, a rank test is less powerful than a permutation test, but see Section 9.3 for a discussion of the merits and drawbacks of using ranks.

Table 2.2. Comparison of Methods for Testing Equality of Means of Two Populations

Permutation	Distribution-free Methods		
	Rank (e.g., Wilcoxon)	Nonparametric Bootstrap	Parametric (e.g., $t$ -test)
Choose test statistic	Choose test statistic	Choose test statistic	Choose test statistic whose distribution can be derived analytically
(e.g., sum of observations in first sample)	(e.g., sum of ranks in first sample)	(e.g., difference between means of samples)	(e.g., Student's $t$ )
Calculate statistic	Convert to ranks Calculate statistic	Calculate statistic	Calculate statistic
Are observations exchangeable?	Are observations exchangeable?	Are observations independent? With identical parameters of interest?	Are observations independent? Do they follow specified distribution?
Derive permutation distribution from combined sample	Use table of permutation distribution of ranks	Derive bootstrap distribution: resample separately from each sample	Use tabulated distribution
Compare statistic with percentiles of distribution	Compare statistic with percentiles of distribution	Compare statistic with percentiles of distribution	Compare statistic with percentiles of distribution

The bootstrap is a relatively recent introduction (circa 1970), primarily because the bootstrap also is computation intensive. The bootstrap, like the permutation test, requires a minimum number of assumptions and derives its critical values from the data at hand.

To obtain a nonparametric bootstrap, we:

1. Choose a test statistic  $S(X)$ .
2. Compute  $S$  for the original set of observations.
3. Obtain the bootstrap distribution of  $S$  by repeatedly resampling from the observations. We need not combine the samples, but may resample separately from each sample. We resample with replacement.
4. Obtain the upper  $\alpha$ -percentage point of the bootstrap distribution and accept or reject the null hypothesis according to whether  $S$  for the original observations is smaller or larger than this value.

The bootstrap is neither exact nor conservative. Generally, but not always, a nonparametric bootstrap is less powerful than a permutation test. One exception

to the rule is when we compare the variances of two populations (see Section 3.4). If the observations are independent and from distributions with identical values of the parameter of interest, then the bootstrap is asymptotically exact [Liu, 1988]. And it may be possible to bootstrap when no other statistical method is applicable; see Section 4.4.

To obtain a parametric test (e.g., a  $t$ -test or an  $F$ -test), we:

1. Choose a test statistic,  $S$ , whose distribution  $F_s$  may be computed and tabulated independent of the observations.
2. Compute  $S$  for the observations  $X$ .
3. (This step may be skipped as the distribution  $F_s$  is already known and tabulated.)
4. Compare  $S(X)$  with the upper  $\alpha$ -percentage point of  $F_s$  and accept or reject the null hypothesis according to whether  $S(X)$  is smaller or larger than this value.

If  $S$  is distributed as  $F_s$ , then the parametric test is exact and, often, the most powerful test available. In order for  $S$  to have the distribution  $F_s$ , in most cases the observations need to be independent and, with small samples, identically distributed with a specific distribution,  $G_s$ . If  $S$  really has some other distribution, then the parametric test may lack power and may not be conservative. With large samples, the permutation test is usually as powerful as the most powerful parametric test [Bickel and Van Zwet, 1978]. If  $S$  is not distributed as  $F_s$ , it may be more powerful.

## 2.4. World Views

Parametric tests such as Student's  $t$  are based on a *sampling* model. Proponents of this model envision a hypothetical population, infinite in size, whose members take values in accordance with some fixed (if unknown) distribution function. For example, normally distributed observations would be drawn from a population whose values range from minus infinity to plus infinity in accordance with a bell-shaped or normal curve. From this population, proponents claim, we can draw a series of values of independent, identically distributed random variables to form a random sample.

This view of the world is very natural to a trained mathematician, but does it really correspond to the practical reality that confronts the physician, the engineer, or the scientist?

Fortunately, we needn't rely on the existence of a hypothetical infinite population to form a permutation test [Welch, 1937]. The permutation tests make every bit as much sense in a context that Lehmann [1986] terms the *randomization* model in which the results are determined by the specific set of experimental subjects and by how these subjects are assigned to treatment.

Suppose that as a scientist you have done things or are contemplating doing things to the members of some representative subset or sample of a larger population—several cages of rats from the population of all genetically similar rats, several acres of land from the set of all similar acres, several long and twisted rods from the set of all similarly machined rods. Or, as opposed to a sample,

perhaps your particular experiment requires you to perform the same tests on every machine in your factory, or on every available fossil, or on the few surviving members of what was once—before man—a thriving species.

In these experiments, there are two sorts of variation: The variation *within* an experimental subject over which you have little or no control—blood pressure, for example, varies from hour to hour and day to day within a given individual—and the variation *between* subjects over which you have even less control. Observations on untreated subjects take on values that vary about a parameter  $\mu_i$ , which depends on the individual  $i$  who is being examined. Observations on treated subjects have a mean value  $\mu_j + \delta$  where the treatment effect  $\delta$  is confounded with the mean  $\mu_j$  of the  $j$ th experimental subject. How are we to tell if the differences between observations on treated and untreated groups represent a true treatment effect or merely result from differences in the two sets of subjects?

If we assign subjects to treatment categories at *random*, so that every permutation of the labels is equally likely, the joint probability density of the observations is

$$\frac{1}{(n+m)!} \sum_{(j_1, \dots, j_{m+n})} \prod_{i=1}^m f(x_i - \mu_{j_i}) \prod_{i=1}^n f(x_i - \mu_{j_{m+i}} - \delta).$$

Under the null hypothesis of no treatment effect, that is  $\delta = 0$ , this density can be written as

$$\frac{1}{(n+m)!} \sum_{(j_1, \dots, j_{m+n})} \prod_{i=1}^{m+n} f(x_i - \mu_{j_i}).$$

By *randomizing* the assignment of subjects to treatment, we provide a statistical basis for analyzing the results. And we can reduce (but not eliminate) the probability, say, that all the individuals with naturally high blood pressure end up in the treatment group.

Because we know that blood pressure is an important factor, one that varies widely from individual to individual, we could do the experiment somewhat differently, dividing the experimental subjects into blocks so as to randomize separately within a “high” blood pressure group and a “low” blood pressure group. But we may not always know in advance which factors are important. Or, we may not be able to measure these factors until the date of the experiment itself. Fortunately, as we shall see in Sections 4.3 and 9.2, randomizing the assignment of subjects to treatment (or treatments to subject) also ensures that we are in a position to correct for significant cofactors *after* the experiment is completed.

Using a permutation test to analyze an experiment in which we have randomly assigned subjects to treatment is merely to analyze the experiment in the manner in which it was designed.

### 2.4.1. Generalized Permutations

For some, randomization during analysis and randomization during design are irretrievably linked. For others, randomization during analysis is based on the purely

abstract notion of exchangeability. These two world views are often equivalent, but they are not the same.

Consider a balanced two-way design as defined in Section 4.2. Shuffling the residuals under the hypothesis of no treatment effect provides many more sets of residuals than would have been obtained under the alternate design randomizations. Nonetheless, all of these sets are legitimate under the null hypothesis of 1) no main effects, 2) no interaction, and 3) exchangeability of random errors. Exchanging labels is not the same as exchanging designs.

## 2.5. Questions

1. a) *Power.* Sketch the power curve  $\beta(\theta)$  for one or both of the two-sample comparisons described in this chapter. (You already know two of the values for each power curve. What are they?)

- b) Using the same set of axes, sketch the power curve of a test based on a much larger sample.
  - c) Suppose that without looking at the data you
    - i) always reject;
    - ii) always accept; or
    - iii) use a chance device so as to reject with probability  $\alpha$ .

For each of these three tests, determine the power and the significance level. Are any of these three tests exact? Unbiased?

- a) *Decisions.* Suppose you have two potentially different radioactive isotopes with half-life parameters  $\lambda_1$  and  $\lambda_2$ , respectively. You gather data on the two isotopes and, taking advantage of a uniformly-most-powerful-unbiased permutation test, you reject the null hypothesis  $H : \lambda_1 = \lambda_2$  in favor of the one-sided alternative not  $H : \lambda_1 > \lambda_2$ . What are you or the person you are advising going to do about it? Will you need an estimate of  $\lambda_1/\lambda_2$ ? What estimate will you use? (Hint: See Section 3.2 in the next chapter.)

- b) Review some of the hypotheses you tested in the past. Distinguish your actions after the test was performed from the conclusions you reached. (In other words, did you do more testing? Rush to publication? Abandon a promising line of research?) What losses were connected with your actions? Should you have used a higher/lower significance level? Should you have used a more powerful test or taken more/fewer observations? And, if you used a parametric test like Student's  $t$  or Welch's  $z$ , were all the assumptions for these tests satisfied?

- a) The advertisement reads, "Safe, effective, faster than aspirin." A picture of a happy smiling woman has the caption, "My headache vanished faster than I thought possible." The next time you are down at the pharmacy, the new drug is there at the same price as your favorite headache remedy. Would you buy it? Why or why not? Do you think the ad is telling the truth? What makes you think it is?

- b) In the United States, in early 1995, a variety of government agencies and regulations would almost guarantee the ad is truthful—or, if not, that it would not appear in print a second time. Suppose you are part of the government's regulatory team reviewing the evidence supplied by the drug company. Looking into the claim of safety, you are told only "we could not reject the null hypothesis." Is this statement adequate? What else would you want to know?

- c) Suppose, once again, you are a consumer with a splitting headache, but when you go to buy the new drug, you discover it is twice the price of your favorite remedy. The ad does promise it is faster than aspirin; a footnote to the ad states a statistically significant increase in speed was found in an FDA-approved survey of 100 patients. Would you be willing to pay the difference in price for the new drug? Why or why not?
4. a) Your lab has been gifted with a new instrument offering ten times the precision of your present model. How might this affect the power of your tests? their significance level? the number of samples you'll need to take?
  - b) A directive from above has loosened the purse strings so you now can take larger samples. How might this affect the power of your tests? their significance level? the precision of your observations? the precision of your results?
  - c) A series of law suits over those silicon implants you thought were harmless has totally changed your company's point of view about the costs of sampling. How might this affect the number of samples you'll take? the power of your tests? their significance level? the precision of your observations? the precision of your results?
5. Are the residuals in a regression analysis or an analysis of variance exchangeable?  
If you aren't satisfied with or are uncertain of your answers, you may want to return to these questions as you proceed further into the text.

## CHAPTER 3

# Testing Hypotheses

In this chapter, you learn how to approach and resolve a series of testing problems of increasing complexity, specifically, tests for location and scale parameters in one, two, and  $k$  samples. You learn how to derive confidence intervals for the unknown parameters, and you learn to increase the power of your tests by sampling from blocks of similar composition.

## 3.1. One-Sample Tests

### 3.1.1. Tests for a Location Parameter

One of the simplest testing problems would appear to be that of testing for the value of the location parameter of a distribution using a series of observations  $x_1, x_2, \dots, x_n$  from that distribution. This testing problem is a simple one if we can assume that the underlying distribution is symmetric about the unknown parameter  $\theta$ , that is, if

$$\Pr\{X \leq \theta - x\} = \Pr\{X \geq \theta + x\}, \quad \text{for all } x.$$

The normal distribution with its familiar symmetric bell-shaped curve and the double exponential, Cauchy, and uniform distribution are examples of symmetric distributions. The difference of two independent observations drawn from the same population also has a symmetric distribution, as you will see when we come to consider experiments involving matched pairs in Section 3.7.

Suppose now we wish to test the hypothesis that  $\theta \leq \theta_0$  against the alternative that  $\theta > \theta_0$ . As in Chapter 1, we proceed in five steps.

First, we choose a test statistic that will discriminate between the hypothesis and the alternative. As one possibility, consider the sum of the deviations about  $\theta_0$ . Under the hypothesis, positive and negative deviations ought to cancel and this sum should be close to zero or negative. Under the alternative, positive terms should predominate and this sum should be large. But how large should the sum be for us to reject the hypothesis?

We saw in Chapter 2 that we can use the permutation distribution to obtain the answer, but what should we permute? The principle of sufficiency can help us here.

Suppose we had lost track of the signs (plus or minus) of the deviations. We could attach new signs at random, selecting a plus or a minus with equal probability. If we are correct in our hypothesis that the variables have a symmetric distribution about  $\theta_0$ , the resulting values should have precisely the same distribution as the original observations. The absolute values of the observations are sufficient for regenerating the sample. (You'll find more on the topic of sufficiency in Sections 10.3 and 14.2 with regard to choosing a test statistic.)

Under the alternative of a location parameter larger than  $\theta_0$ , randomizing the signs of the deviations should reduce the sum from what it was originally; as we consider one after another in a series of random reassessments, our original sum should be revealed as an extreme value.

Before implementing this permutation procedure, we note that the sum of just the deviations with plus signs attached is related to the sum of all the deviations by the formula:

$$\sum_{\{x_i > 0\}} x_i = \left( \sum x_i + \sum |x_i| \right) / 2,$$

because the +1's get added twice, once in each sum on the right-hand side of the equation while the -1's and |-1|'s cancel. Thus, we can reduce the number of calculations by summing only the positive deviations.

As an illustration, suppose that  $\theta_0$  is 0 and that the original observations are -1, 2, 3, 1.1, 5. Our first step is to compute the sum of the positive deviations, which is 11.1.

Among the  $2 \times 2 \times 2 \times 2 \times 2$  or  $2^5$  possible reassessments of plus and minus signs are

$$\begin{aligned} &+1, -2, +3, +1.1, +5, \\ &+1, +2, +3, +1.1, +5, \end{aligned}$$

and

$$-1, -2, +3, +1.1, +5.$$

Our third step is to compute the sum of the positive deviations for each rearrangement. For the three rearrangements shown above, this sum would be 10.1, 12.1, and 9.1, respectively.

Our fourth step is to compare the original value of our test statistic with its permutation distribution. Only two of the 32 rearrangements have sums as large as the sum, 11.1, of the original observations. Is  $2/32 = 1/16 = .0625$  statistically significant? Perhaps. It all depends on the relative losses we assign to Type I and Type II error and on the loss function; are small differences of practical as well as statistical significance? Certainly, a significance level of 0.0625 is suggestive. Suggestive enough that in this case we might want to look at additional data or perform additional experiments before accepting the hypothesis that 0 is the true value of  $\theta$ .

### 3.1.2. Properties of the Test

Adopting the sampling model advanced in Section 2.4, we see the preceding permutation test is applicable even if the different observations come from different distributions—provided, that is, that these distributions are all symmetric and all have the same location parameter or median. (If these distributions are symmetric then if the mean exists, it is identical with the median.) If you are willing to specify their values through the use of a parametric model, the medians needn’t be the same (see problem 6).

*Most powerful test.* Against specific normal alternatives, this permutation test provides a most powerful unbiased test of the distribution-free hypothesis  $H : \theta = \theta_0$  [Lehmann, 1986, p. 239]. For large samples, its power is almost the same as Student’s  $t$ -test [Albers, Bickel, and van Zwet, 1976]. We provide proofs of these and related results in Chapter 14.

*Asymptotic consistency.* What happens if the underlying distributions are almost but not quite symmetric? Romano [1990] shows that the permutation test for a location parameter is asymptotically exact, provided the underlying distribution has finite variance. His result applies whether the permutation test is based on the mean, the median, or some statistical functional of the location parameter. If the underlying distribution is almost symmetric, the test will be almost exact even when based on as few as 10 or 12 observations. See Section 13.7 for the details of a Monte Carlo procedure to use in deciding when “almost” means “good enough.”

#### Capsule Summary

ONE-SAMPLE TEST     $H$ : mean/median =  $\theta_0$   
                              $K$ : mean/median  $\neq \theta_0$

##### Assumptions

- 1) exchangeable observations
- 2) distributions  $F_i$  symmetric about median

Transform

Let  $X'_i = X_i - \theta_0$

Test statistic

Sum of nonnegative  $X'_i$

### 3.1.3. Exact Significance Levels: A Digression

Many of us are used to reporting our results in terms of significance levels of 0.01, 0.05, or 0.10, and significance levels of 0.0625 or 0.03125 may seem confusing at first. These “oddball” significance levels often occur with small sample sizes. Five observations means just 32 possibilities and one extreme observation out of 32

corresponds to 0.03125. Things improve as sample sizes get larger. With seven observations, we can test at a significance level of 0.049. Is this close enough to 0.05?

Lehmann [1986] describes a method called “randomization on the boundary” for obtaining a significance level of exactly 5% (or exactly 1%, or exactly 10%). But this method isn’t very practical. In the worst case, “on the boundary,” you must throw a die or use some other chance device to make your decision.

What is the practical solution? We agree with Kempthorne [1975, 1977, 1979]. Forget tradition. There is nothing sacred about a  $p$ -value of 5% or 10%. Report the exact significance level, whether it is 0.065 or 0.049. Let your colleagues reach their own conclusions based on the losses they associate with each type of error.

## 3.2. Confidence Intervals

The method of randomization can help us find a good interval estimate of the unknown location parameter  $\theta$ .

The set of confidence intervals are the duals of the corresponding tests of hypotheses.

In the first step of our permutation test for the location parameter of a single sample, we subtract  $\theta_0$  from each of the observations. We might test a whole series of hypotheses involving different values for  $\theta_0$  until we find a  $\theta_1$  such that as long as  $\theta_0 \geq \theta_1$ , we accept the hypothesis, but if  $\theta_0 < \theta_1$  we reject it. Then an 100  $(1 - \alpha)\%$  confidence interval for  $\theta$  is given by the interval  $\{\theta > \theta_1\}$ .

Suppose the original observations are  $-1, 2, 3, 1.1$ , and  $5$  and we want to find a confidence interval that will cover the true value of the parameter  $\frac{31}{32}$ nds of the time. In the first part of this chapter, we saw that  $\frac{1}{16}$ th of the rearrangements of the signs resulted in samples that were as extreme as these observations. Thus, we would accept the hypothesis that  $\theta \leq 0$  at the  $\frac{1}{16}$ th and any smaller level including the  $\frac{1}{32}$ nd. Similarly, we would accept the hypothesis that  $\theta \leq -0.5$  at the  $\frac{1}{32}$ nd level, or even that  $\theta \leq -1 + \varepsilon$  where  $\varepsilon$  is an arbitrarily small but still positive number. But we would reject the hypothesis that  $\theta \leq -1 - \varepsilon$  as after subtracting  $-1 - \varepsilon$  the transformed observations are  $\varepsilon, 3 + \varepsilon, 4 + \varepsilon, 2.1 + \varepsilon, 6 + \varepsilon$ .

Our one-sided confidence interval is  $\{-1, \infty\}$  and we have confidence that  $\frac{31}{32}$ nds of the time the method we’ve used yields an interval that includes the true value of the location parameter  $\theta$ .

Our one-sided test of a hypothesis gives rise to a one-sided confidence interval. But knowing that  $\theta$  is larger than  $-1$  may not be enough. We may want to pin  $\theta$  down to a more precise two-sided interval, say, that  $\theta$  lies between  $-1$  and  $+1$ .

To accomplish this, we need to begin with a two-sided test. Our hypothesis for this test is that  $\theta = \theta_0$  against the two-sided alternatives that  $\theta$  is smaller or larger than  $\theta_0$ . We use the same test statistic—the sum of the positive observations, that we used in the previous one-sided test. Again, we look at the distribution of our test statistic over all possible assignments of the plus and minus signs to the observations. But this time we reject the hypothesis if the value of the test statistic

for the original observations is either one of the largest or one of the smallest of the possible values.

In our example, we don't have enough observations to find a two-sided confidence interval at the  $\frac{31}{32}$ nd level, so we'll try to find one at the  $\frac{15}{16}$ ths. The lower boundary of the new confidence interval is still  $-1$ . But what is the new upper boundary? If we subtract five from every observation, we would have the values  $-6, -3, -2, -3.9, -0$ ; their sum is  $-14.9$ . Only the current assignment of signs to the transformed values, that is, only one out of the 32 possible assignments, yields this small a sum for the positive values. The symmetry of the permutation test requires that we set aside another  $\frac{1}{32}$ nd of the arrangements at the high end. Thus we would reject the hypothesis that  $\theta = 5$  at the  $\frac{1}{32} + \frac{1}{32}$  or  $\frac{1}{16}$ th level. Consequently, the interval  $\{-1, 5\}$  has a  $\frac{15}{16}$ th chance of covering the unknown parameter value.

These results are readily extended to a confidence interval for a vector of parameters,  $\theta$ , that underlies a one-sample, two-sample, or  $k$ -sample experimental design with single- or vector-valued variables. In each case, the  $100(1 - \alpha)\%$  confidence interval consists of all values of the parameter vector  $\theta$  for which we would accept the hypothesis at level  $\alpha$ . Remember, one-sided tests produce one-sided intervals and two-sided tests produce two-sided confidence intervals.

In deriving a confidence interval, we look first for a *pivotal quantity* or *pivot*,  $Q(X_1, \dots, X_n, \theta)$ , whose distribution is independent of the parameters of the original distribution. One example is  $Q = \bar{X} - v$ , where  $\bar{X}$  is the sample mean, and the  $\{X_i\}_{i=1, \dots, n}$  are independent and identically distributed as  $F(x - v)$ . A second example is  $Q = \bar{X}/\sigma$ , where the  $\{X_i\}$  are independent and identically distributed as  $F(x/\sigma)$ . If the  $\{X_i\}$  are independent with the identical exponential distribution  $1 - \exp[-\lambda t]$  (see problem 2 in Chapter 2), then  $T = 2 \sum t_i / \lambda$  is a pivotal quantity whose distribution does not depend on  $\lambda$ . We can use this distribution to find an  $a$  and  $b$  such that  $\Pr(a < T < b) = 1 - \alpha$ . But then  $\Pr\left\{\frac{1}{2b \sum t_i} < \lambda < \frac{1}{2a \sum t_i}\right\} = 1 - \alpha$ . We use a pivotal quantity in Section 7.5 to derive a confidence interval for a regression coefficient.

For further information on deriving confidence intervals using the randomization approach see Section 14.3, as well as Lehmann [1986, pp. 246–263], Gabriel and Hsu [1983], John and Robinson [1983], Maritz [1981, pp. 7, 25], and Tritchler [1984]. For a discussion of the strengths and weaknesses of pivotal quantities, see Berger and Wolpert [1984].

### 3.2.1. Comparison with Other Tests

When a choice of statistical methods exists, the best method is the one that yields the shortest confidence interval for a given significance level. Robinson [1987] finds approximately the same coverage probabilities for three sets of confidence intervals for the slope of a simple linear regression, based, respectively, on 1) the standardized bootstrap; 2) parametric theory; and 3) a permutation procedure.

### Confidence Intervals and Rejection Regions

There is a close connection between the confidence intervals and the rejection regions we've constructed. If  $A(\theta')$  is a  $1 - \alpha$  level acceptance region for testing the hypothesis  $\theta = \theta'$ , and  $S(X)$  is a  $1 - \alpha$  level confidence interval for  $\theta$  based on the vector of observations  $X$ , then for the confidence intervals defined here,  $S(X)$  consists of all the parameter values  $\theta^*$  for which  $X$  belongs to  $A(\theta^*)$ , while  $A(\theta)$  consists of all the values of the statistic  $x$  for which  $\theta$  belongs to  $S(x)$ .

$$P_{\theta}\{\theta \in S(X)\} = P_{\theta}\{X \in A(\theta)\} \geq 1 - \alpha.$$

In Section 14.3, we show that if  $A(\theta)$  is the acceptance region of an unbiased test, the correct value of the parameter is more likely to be covered by the confidence intervals we've constructed than is an incorrect value.

## 3.3. Comparison of Locations

We tested the equality of the location parameters of two samples in Chapter 1. Recall that we observed 121, 118, and 110 in the treatment group and 34, 12, and 22 in the control group. Our test statistic was the sum of the observations in the first group and we rejected the null hypothesis because the observed value of this statistic, 349, was as large or larger than it would have been in any of the  $\binom{6}{3} = 20$  rearrangements of the data.

In Chapter 14, we show that a permutation test based on this statistic is exact and unbiased against stochastically increasing alternatives of the form  $K : F_2[x] = F_1[x - \delta]$ ,  $\delta > 0$ . In fact, we show that this permutation test is a uniformly most powerful unbiased test of the null hypothesis  $H : F_2 = F_1$  against normally distributed shift alternatives. Against normal alternatives and for large samples, its power is equal to that of the standard  $t$ -test [Bickel and van Zwet, 1978].

The permutation test offers the advantage over the parametric  $t$ -test that it is exact even for very small samples whether or not the observations come from a normal distribution. The parametric  $t$ -test relies on the existence of a mythical infinite population from which all the observations are drawn (see Section 2.4). The permutation test is applicable even to finite populations such as all the machines in a given shop or all the supercomputers in the world.

### 3.3.1. An Example

Suppose we have two samples: The first, control sample takes values 0, 1, 2, 3, and 19. The second, treatment sample takes values 3.1, 3.5, 4, 5, and 6. Does the treatment have an effect?

The answer would be immediate if it were not for the value 19 in the first sample. The presence of this extreme value changes the mean of the first sample from 1.5 to 5. To dilute the effect of this extreme value on the results, we convert all the

data to ranks, giving the smallest observation a rank of 1, the next smallest the rank of 2, and so forth. The first sample includes the ranks 1, 2, 3, 4, and 10 and the second sample includes the ranks 5, 6, 7, 8, and 9. Is the second sample drawn from a different population than the first?

Let's count. The sum of the ranks in the first sample is 20. All the rearrangements with first samples of the form 1, 2, 3, 4,  $k$ , where  $k$  is chosen from {5, 6, 7, 8, 9 or 10} have sums that are as small or smaller than that of our original sample. That's six rearrangements. The four rearrangements whose first sample contains 1, 2, 3, 5, and a fifth number chosen from the set {6, 7, 8, 9} also have smaller sums. That's  $6 + 4 = 10$  rearrangements so far.

Continuing in this fashion—we leave the complete enumeration as an exercise—we find that 19 of the  $\binom{10}{5} = 252$  possible rearrangements have sums that are as small or smaller than that of our original sample. Two samples this different will be drawn from the same population just under 8% of the time by chance.

#### Capsule Summary

##### TWO-SAMPLE TEST FOR LOCATION

- $H_0$ : mean/medians of groups differ by  $d_0$   
 $H_A$ : mean/medians of groups differ by  $d > d_0$

##### Assumptions

- 1) exchangeable observations
- 2)  $F_{1i}(x) = F(x) = F_{2i}(x - d)$

$$\text{Transform} \quad X'_{1i} = X_{1i} - d_0$$

##### Test statistic

Sum of observations in smallest sample

### 3.3.2. Violation of Assumptions

Just because the permutation test is distribution free does not mean that it is assumption free. For example, what if  $F_2$  is not equal to  $F_1$  under the hypothesis? What if only the location parameters of  $F_2$  and  $F_1$  are equal while the other parameters of these two distributions are not the same? Is this permutation test still exact or almost exact? Is it still efficient for testing against normal alternatives?

Romano [1990] shows that the permutation test based on the sum of the observations in the first sample is asymptotically exact for testing whether the first moments of  $F_2$  and  $F_1$  are equal whether  $F_1 = F_2$ , providing both distributions have finite second moments and the samples are of equal size.

But suppose the first sample comes from a population whose members are very close to zero, while the second comes from a population whose members are very large in absolute value. The sum of the observations in the first sample is close to zero in the original sample. When we rerandomize, mixing elements of the second sample with those of the first, the test statistic will either be a very large positive

number or a very large negative one. In this worst case example, we accept the null hypothesis, for  $\alpha < 50\%$  and for any finite sample size, independent of the samples and of parameters to be tested.

The preceding example is a worst case. In a series of Monte Carlo simulations with small sizes, I found that, in many common applications including the Behrens–Fisher problem, the permutation test remains close to exact even for very small equal-sized samples; see Table 3.1. This result is in line with the findings of Box and Andersen [1955] and Brown [1982], but see also, Boik [1987].

A second resampling method, the nonparametric bootstrap, provides asymptotically exact solutions, whether  $F_1 = F_2$  and whether the sample sizes are equal, but see Gine and Zinn [1989]. In a bootstrap, we resample separately with replacement from each of the two original samples. The underlying populations need not be the same, even under the null hypothesis.

The primitive, uncorrected bootstrap is far from exact except for very large samples. If we modify the bootstrap using pivots, studentization, bias, and higher order correction, as in Hall [1992], we can derive an almost exact bootstrap even for samples with as few as eight observations.

This result is not unexpected: Liu [1988] shows the bootstrap test of the hypothesis of equal means retains the second-order convergence properties it has in the case  $F_1 = F_2$ .

## 3.4. Comparing Dispersions

Precision is essential in a manufacturing process. Items that are too far out of tolerance must be discarded and an entire production line brought to a halt if too many items exceed (or fall below) designated specifications. With some testing equipment, such as that used in hospitals, precision can be more important than accuracy. Accuracy can always be achieved through the use of standards with known values, while a lack of precision may render an entire sequence of tests invalid.

### 3.4.1. The Permutation Approach

There is no shortage of parametric methods to test the hypothesis that two samples come from populations with the same inherent variability (see, for example, Conover, Johnson, and Johnson, 1981), but few can be relied on. Many promise an error rate (significance level) of 5%, but in reality make errors as frequently as 8% to 20% of the time. Others have severe restrictions. The  $F$ -ratio test [Fisher, 1924] is exact only if the variates are normally distributed and, unlike the  $t$ -test, is very sensitive to nonnormality.

At first glance, the permutation test for comparing the variances of two populations would appear to be an immediate extension of the test we use for comparing location parameters, in which we use the squares of the observations rather than the observations themselves. But these squares are actually the sum of two

Table 3.1. Significance Level, Effect of Unequal Variance on the Permutation Test, Bootstrap, and Student's for Various Sample Sizes and Distributions. Rejections in 1000 Monte Carlo Simulations.

$\sigma_1/\sigma_2$	1.0					4.0					10.0					
	<i>p</i> -test		boot		<i>t</i> -test	<i>p</i> -test		boot		<i>t</i> -test	<i>p</i> -test		boot		<i>t</i> -test	
Ideal	100	50	100	50	100	50	100	50	100	50	100	50	100	50	100	
Normal	6,6	100	57	108	57	98	45	98	50	84	36	115	58	111	56	83
	8,8	100	47	103	54	102	52	112	67	94	45	117	68	106	64	94
	8,12	99	52	111	47	101	50	69	29	98	48	72	34	51	23	83
	12,8	115	59	103	47	95	37	138	75	80	36	75	37	168	111	95
	12,12	101	53	107	55	100	50	105	62	110	57	87	48	101	43	107
Double Exponential																
	6,6	102	57	130	69	92	44	104	57	119	68	109	61	138	87	148
	8,8	108	45	119	65	116	59	114	62	136	83	132	66	126	70	151
	8,12	89	51	102	48	101	50	61	27	102	44	72	25	65	27	109
	12,8	113	44	129	72	93	28	146	95	121	72	106	52	134	78	99
	12,12	100	44	132	78	80	38	85	46	117	76	85	40	132	77	179

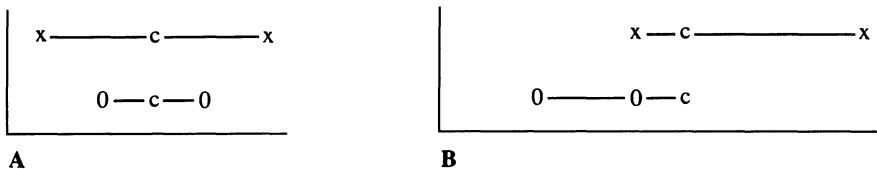


Figure 3.1. Comparison of two samples: a) original data, b) after first sample is shifted to the right. c common center, x—x first sample, 0—0 second sample.

components, one of which depends on the unknown variance, the other on the unknown location parameter. In symbols, where  $EX$  represents the mathematical expectation of  $X$ , that is, the average of a very large number of independent observations,

$$EX^2 = E(X - \mu + \mu)^2 = E(X - \mu)^2 + 2\mu E(X - \mu) + \mu^2 = \sigma^2 + 0 + \mu^2.$$

A permutation test based on the squares of the observations is appropriate only if the location parameters of the two populations are known or known to be equal (Bailer, 1989; see also, Hayes, 1997).

We cannot eliminate the effects of the location parameters by working with the deviations about each sample mean as these deviations are interdependent [Maritz, 1981]. The problem is illustrated in Figure 3.1. In the sketch on the left, the observations in the first sample are both further from the common center than either of the observations in the second sample, and of the four possible rearrangements of four observations between two samples, this arrangement is the most extreme. In the sketch on the right, the observations in the first sample have undergone a shift to the right; this shift has altered the relative ordering of the absolute deviations about the common center, and at least one other rearrangement is more extreme.

Still, we needn't give up; we can obtain an asymptotically exact permutation test with just a few preliminary calculations. First, we compute the median for each sample; next, we replace each of the remaining observations by the square of its deviation about its sample median; last, we discard certain redundant values.

Suppose the first sample contains the observations  $x_1, \dots, x_n$ , whose median is  $\text{mdn}\{x_i\}$ ; we begin by forming the deviates  $x'_j = |x_j - \text{mdn}\{x_i\}|$  for  $j = 1, \dots, n$ .<sup>1</sup> Similarly, we form the set of deviates  $\{y'_j\}$  using the observations in the second sample and their median.

If there is an odd number of observations in the sample, then one of these deviates must be zero. We can't get any information out of a zero, so we throw it away. If there is an even number of observations in the sample, then two of these deviates must be equal. We can't get any information out of the second one that we didn't already get from the first, so we throw it away.

Our test statistic  $S$  is the sum of the remaining deviations in the first sample, that is,  $S = \sum_{j=1}^{n-1} x'_j$ . We obtain its permutation distribution and the cutoff point for

<sup>1</sup> Alternately, we might use the formula  $x'_j = |x_j - \text{mdn}\{x_j\}|^2$ .

the test by considering all possible rearrangements of the deviations that remain in both the first and second samples.

To illustrate the application of this statistic, suppose the first sample consists of the measurements 121, 123, 126, 128.5, 129, and the second sample consists of the measurements 153, 154, 155, 156, 158. Thus,  $x'_1 = 5$ ,  $x'_2 = 3$ ,  $x'_3 = 2.5$ ,  $x'_4 = 3$ , and  $S_0 = 13.5$ . While  $y'_1 = 2$ ,  $y'_2 = 1$ ,  $y'_3 = 1$ ,  $y'_4 = 3$ . There are  $\binom{8}{4} = 70$  arrangements in all, of which only three yield values of the test statistic as or more extreme than our original value.  $3/70 = 0.043$ , and we could conclude that there is a significant difference between the dispersions of the two manufacturing processes at the 5% significance level. Still, with such a small number of rearrangements, the evidence is unconvincing. We should continue to take observations until the evidence is overwhelming in one direction or the other.

There is a weak dependency among these deviates, and thus they are only asymptotically exchangeable. (See Section 14.4.3.) For very small samples, the test is alternately conservative and liberal, see Baker [1995]. For the test to be asymptotically exact, a) the ratio of the sample sizes  $n, m$  must be close to 1, b) the population variances must exist and be equal (as they would be under the null hypothesis), and c) the only other difference between the two populations from which the samples are drawn is that they might have different means.

### Capsule Summary

#### TWO-SAMPLE TEST FOR VARIANCE

$$\begin{aligned} H: & \text{variances of populations are equal} \\ K: & \sigma_2^2 > \sigma_1^2 \end{aligned}$$

#### Assumptions

- 1) independent observations
- 2) continuous observations
- 3)  $F_{1i}[x] = F_{2i}[x - d]$

$$\text{Transform } X'_{ij} = (X_{ij} - M_{di})^2$$

Discard redundant deviate from each sample

#### Test statistic

Sum of  $X'_{ij}$  in smallest sample

The statistic  $\delta$  proposed by Aly [1990]<sup>2</sup> doesn't have even this latter restriction.

$$\delta = \sum_{i=1}^{m-1} (i)(m-i)(X_{(i+1)} - X_{(i)}),$$

where  $X_{(1)} < X_{(2)} < \dots < X_{(m)}$  are the *order statistics* of the first sample.

<sup>2</sup> The choice of test statistic is dictated by the principles of density estimation; see, for example, Izenman [1991].

Applying this statistic to the samples considered above,  $m = 5$ ,  $X_{(1)} = 121$ ,  $X_{(2)} = 123$ , and so forth.

Set  $z_{1i} = X_{(i+1)} - X_{(i)}$  for  $i = 1, \dots, 4$ . In this instance,  $z_{11} = 123 - 121 = 2$ ,  $z_{12} = 3$ ,  $z_{13} = 2.5$ , and  $z_{14} = 0.5$ .

The original value of the test statistic is  $8 + 18 + 15 + 2 = 43$ . To compute the test statistic for other arrangements, we also need to know the differences  $z_{2i} = Y_{(i+1)} - Y_{(i)}$  for  $i = 1, \dots, 4$ .  $z_{21} = 1$ ,  $z_{22} = 1$ ,  $z_{23} = 1$ , and  $z_{24} = 2$ .

One possible rearrangement is  $\{2, 1, 1, 2\}$ , which yields a value of  $S = 8 + 6 + 6 + 8 = 28$ .

There are  $2^4 = 16$  rearrangements in all, of which one  $\{2, 3, 2.5, 2\}$  is more extreme than the original observations. With two out of 16 rearrangements as or more extreme than the original, we accept the null hypothesis.<sup>3</sup>

If our second sample is larger than the first, we have to resample in two stages. First, we select a subset of  $m$  values  $\{Y_i^*, i = 1, \dots, m\}$  without replacement from the  $n$  observations in the second sample, and compute the order statistics  $Y_{(1)}^* < Y_{(2)}^* < \dots < Y_{(m)}^*$  and their differences  $\{z_{2i}\}$ . Last, we examine all possible values of Aly's measure of dispersion for permutations of the combined sample  $\{\{z_{1i}\}, \{z_{2i}\}\}$ , as we did when the two samples were equal in size, and compare Aly's measure for the original observations with this distribution.

A similar procedure, first suggested by Baker [1995], should be used with the first test in the case of unequal samples.

### 3.4.2. The Bootstrap Approach

To obtain a test based on the bootstrap confidence interval for the variance ratio, we resample repeatedly without replacement, drawing independently from the two original samples, until we have two new samples the same size as the originals. Each time we resample, we compute the variances of the two new independent subsamples and calculate their ratio. The resultant bootstrap confidence interval is asymptotically exact [Efron, 1981] and can be made close to exact with samples of as few as eight observations: See Table 3.2a. As Table 3.2b shows, this bootstrap is more powerful than the permutation test we described in the previous section. One caveat also revealed in the table: this bootstrap is still only “almost” exact.

## 3.5. *k*-Sample Comparisons

### 3.5.1. Comparing Location Parameters

Suppose we wanted to assess the effect on crop yield of hours of sunlight, observing the yield  $X_{ij}$  for  $I$  different levels of sunlight,  $i = 1, \dots, I$  with  $n_i$  observations

<sup>3</sup> Again, the wiser course would be to take a few more observations.

Table 3.2a. Significance Level for Variance Comparisons for  $BC_a$  Method, Efron and Tibshirani [1986]. For Various Underlying Distributions by Sample Size. 500 Simulations.

	6, 6	8, 8	8, 12	12, 8	12, 12	15, 15
Ideal	50	50	50	50	50	50
Normal (0, 1)	44	52	53	56	45	49
Double (0, 1)	53	51	63	70	55	54
Gamma (4.1)	48	55	60	65	52	52
Exponential	54	58	56	70	46	63

Table 3.2b. Power as a Function of the Ratio of the Variances. For Various Distributions and Two Samples Each of Size 8. Rejections in 500 Monte Carlo Simulations.

$\phi = \sigma_2/\sigma_1$	Permutation Test					Bootstrap*				
	1.	1.5	2.	3.	4.	1.	1.5	2.	3.	4.
Ideal	50			500	50				500	
Normal	52	185	312	438	483	52	190	329	444	482
Double	55	153	215	355	439	53	151*	250*	379*	433
Gamma	44	158	255	411	462	49	165	288	426	464
Exponential	51	132	224	323	389	54	150*	233*	344*	408

\* Bootstrap intervals shortened so actual significance level is 10%.

at each level. Our model is that

$$X_{ij} = \mu_i + \varepsilon_{ij},$$

where the experimental errors  $\{\varepsilon_{ij}\}$  are exchangeable.

Let  $X_{i\cdot} = \sum_{j=1}^{n_i} X_{ij}/n_i$  and  $X_{..} = \sum_{i=1}^I \sum_{j=1}^{n_i} X_{ij}/N$ , where  $N = \sum_{i=1}^I n_i$ . The sum of squares of the deviations about the grand mean may be analyzed into two sums,

$$\sum_{i=1}^I \sum_{j=1}^{n_i} (X_{ij} - X_{..})^2 = \sum_{i=1}^I \sum_{j=1}^{n_i} (X_{ij} - X_{i\cdot})^2 + \sum_{i=1}^I n_i (X_{i\cdot} - X_{..})^2,$$

the first of which represents the *within-group* sum of squares, and the second represents the *between-group* sum of squares.

Just as Student's  $t$  is the classic parametric statistic for testing the hypothesis that the means of two normal distributions are the same, so the  $F$ -ratio  $W$  of the between-group variance to the within-group variance is the classic parametric statistic for testing the hypothesis that the means of  $k$  normal distributions are the same

$$W = \frac{\sum n_i (X_{i\cdot} - X_{..})^2 / (I - 1)}{\sum \sum (X_{ij} - X_{i\cdot})^2 / (N - I)}.$$

It is easy to see that  $W$  is invariant under transformations of scale or origin. Lehmann [1986, p. 375] shows that against normal alternatives, and among all similarly invariant tests, the parametric test based on  $W$  is a uniformly most powerful procedure for testing the hypothesis  $H: \mu_1 = \dots = \mu_I$  against the alternative  $K: \mu_i \neq \mu_j$  for some pair  $(i, j)$ .

If the  $\{X_{ij}\}$  are normally distributed with a common variance, then under the hypothesis,  $W$  has the  $F$ -distribution with  $I - 1, N - I$  degrees of freedom. But we may not know or not be willing to assume that these observations do come from a normal distribution, but rather from some common but unknown distribution  $G$ . Explicitly, let  $X_{ij}(j = 1, \dots, n_i; i = 1, \dots, I)$  be independently distributed as  $G[x - \mu_i]$  and, thus, exchangeable under the null hypothesis. To obtain the permutation distribution of  $W$ , we examine all possible reassessments of the observations to the various treatment groups subject to the restriction that the number of observations in each of the  $k$  groups remains unchanged. Our analysis is exact if the experimental units were randomly assigned to treatment to begin with.

In a sidebar, we've provided an outline of a computer program that uses a Monte Carlo simulation to estimate the significance level (see Section 13.2). This program is applicable to any of the experimental designs we consider in this chapter and the next. Our one programming trick is to pack all observations into a single linear vector  $X = (X_{11}, \dots, X_{1n_1}, X_{21}, \dots, X_{2n_2}, \dots, X_{In_I})$  and then permute the observations within the vector. If we have  $k$  samples, we only need to select  $k - 1$  of them when we rearrange the data. The  $k$ th sample is left over automatically.

We need to write a subprogram to compute the test statistic, but there's less work involved than the formula for  $W$  would suggest. As is the case with the permutation equivalent of the  $t$ -statistic, we can simplify the calculation of the test statistic by eliminating terms that are invariant under permutation of the subscripts. For example, the *within-group* sum of squares in the denominator of  $W$  may be written as the difference of two sums  $\sum \sum (X_{ij} - X_{..})^2$  and  $\sum n_i (X_{i.} - X_{..})^2$ . The first of these sums is invariant under permutation of the subscripts. The second, the *between-groups* sum of squares, already occurs in the numerator. Our test statistic reduces to the between-groups sum of squares  $\sum n_i (X_{i.} - X_{..})^2$ , which reduces in turn after the elimination of invariants to

$$W' = \sum_i \frac{\left( \sum_j X_{ij} \right)^2}{n_i},$$

with a corresponding reduction in the number of calculations.

The size of this test is unchanged and the power nearly so in the face of violations of the normality assumption, providing that the  $\{X_{ij}; j = 1, \dots, n_i\}$  are samples from distributions  $F[x - \mu_i]$ , where  $F$  is an arbitrary distribution with finite variance [Robinson, 1973, 1983]. The parametric version of this test is almost as robust [Lehmann, 1986, p. 378]. The real value of the permutation approach comes when we realize that we are not restricted to a permutation version of an existing statistic, but are free to choose a test statistic optimal for the problem at hand.

**Sidebar**

Program for estimating permutation significance levels; for tips on optimization, see Chapter 13.

Monte, the number of Monte Carlo simulations; try 400  
 $S_0$ , the value of the test statistic for the unpermuted observations  
 $S$ , the value of the test statistic for the rearranged observations  
 $X[ ]$ , a one-dimensional vector that contain the observations  
 $n[ ]$ , a vector that contains the sample sizes  
 $N$ , the total number of observations

**Main program**

Get data  
 put all the observations into a single linear vector  
 Compute the test statistic  $S_0$   
 Repeat Monte times:  
 Rearrange the observations  
 Recompute the test statistic  $S$   
 Compare  $S$  with  $S_0$   
 Print out the proportion of times  $S$  was equal to or larger than  $S_0$

**Rearrange**

Set  $s$  to the size of the combined sample  
 Start: Choose a random integer  $k$  from 0 to  $s - 1$   
 Swap  $X[k]$  and  $X[s - 1]$ :  
 $\text{temp} = X[k];$   
 $X[k] = X[s - 1];$   
 $X[s - 1] = \text{temp}.$   
 Decrement  $s$  and repeat from start  
 Stop after you've selected all but one of the samples.

**Get data**

This user-written procedure gets all the data and packs it into a single long linear vector  $X$ .

**Compute stat**

This user-written procedure computes the test statistic.

### 3.5.2. Pitman Correlation

The *F*-ratio test and its permutation version offer protection against any and all deviations from the null hypothesis of equality among treatment means. As a result, they may offer less protection against some specific alternative than some other test function(s). When we have a specific alternative in mind, as is so often the case in biomedical research, for example, when we are testing for an ordered dose response, the *F*-ratio may not be the statistic of choice.

Frank, Trzos, and Good [1977] studied the increase in chromosome abnormalities and micronuclei as the dose of various known mutagens was increased. Their object was to develop an inexpensive but sensitive biochemical test for mutagenicity that would be able to detect even marginal effects. Thus they were more than willing to trade the global protection offered by the  $F$ -test for a statistical test that would be sensitive to ordered alternatives.

Fortunately, a most powerful unbiased test (and one that is also most powerful among tests that are invariant to changes in scale) has been known since the late 1930's. Pitman [1937] proposes a test for linear correlation using as test statistic

$$S = \sum_{i=1}^I f[i] \sum_{k=1}^{n_i} X_{ik},$$

where  $f[i]$  is any monotone increasing function. The simplest choice is  $f[i] = i$ .

The permutation distributions of  $S_1$  with  $f[i] = ai + b$  and  $S_2$  with  $f[i] = i$  are equivalent in the sense that, if  $S_{10}$ ,  $S_{20}$  are the values of these test statistics corresponding to the same set of observations  $\{x_i\}$ , then  $\Pr(S_1 > S_{10}) = \Pr(S_2 > S_{20})$ .

Let us apply the Pitman approach to the data collected by Frank et al., shown in Table 3.3. As the anticipated effect is proportional to the logarithm of the dose, we take  $f[\text{dose}] = \log[\text{dose} + 1]$ . (Adding a 1 to the dose keeps this function from blowing up at a dose of zero.) There are four dose groups; the original data for breaks may be written in the form

0 1 1 2      0 1 2 3 5      3 5 7 7      6 7 8 9 9.

As  $\log[0 + 1] = 0$ , the value of the Pitman statistic for the original data is  $0 + 11 * \log[6] + 22 * \log[21] + 39 * \log[81] = 112.1$ . The only larger values are associated with the small handful of rearrangements of the form a statistically significant, ordered dose response ( $\alpha < 0.001$ ) has been detected. The micronuclei also exhibit a statistically significantly dose response when we calculate the permutation distribution of  $S$ , with  $f[i] = \log[\text{dose}_i + 1]$ . To make the calculations, we

Table 3.3. Micronuclei in Polychromatophilic Erythrocytes and Chromosome Alterations in the Bone Marrow of Mice Treated with CY.

Dose (mg/kg)	Number of Animals	Micronuclei per 200 cells	Breaks per 25 cells
0	4	0 0 0 0	0 1 1 2
5	5	1 1 1 4 5	0 1 2 3 5
20	4	0 0 0 4	3 5 7 7
80	5	2 3 5 11 20	6 7 8 9 9

0	0	1	2	1	1	2	3	5	3	5	7	7	6	7	8	9	9
0	0	1	1	1	2	2	3	5	3	5	7	7	6	7	8	9	9
0	0	1	1	1	2	2	3	3	5	5	7	7	6	7	8	9	9
0	0	1	2	1	1	2	3	3	5	5	7	7	6	7	8	9	9
0	1	1	2	0	1	2	3	3	5	5	7	7	6	7	8	9	9
0	1	1	2	0	1	2	3	5	3	5	6	7	7	7	8	9	9
0	0	1	2	1	1	2	3	5	3	5	6	7	7	7	8	9	9
0	0	1	1	1	2	2	3	5	3	5	6	7	7	7	8	9	9
0	0	1	1	1	2	2	3	3	5	5	6	7	7	7	8	9	9
0	0	1	2	1	1	2	3	3	5	5	6	7	7	7	8	9	9
0	1	1	2	0	1	2	3	3	5	5	6	7	7	7	8	9	9

took advantage of the computer program we developed in Section 3.5.1; the only change was in the subroutine used to compute the test statistic.

A word of caution: If we use some function of the dose other than  $f[\text{dose}] = \log[\text{dose} + 1]$ , we might not observe a statistically significant result. Our choice of a test statistic must always make biological as well as statistical sense; see question 3 in Section 3.8.

### 3.5.3. Effect of Ties

Ties can complicate the determination of the significance level. Because of ties, each of the rearrangements noted in the preceding example might actually have resulted from several distinct reassessments of subjects to treatment group and must be weighted accordingly. To illustrate this point, suppose we put tags on the 1's in the original sample

0 1\* 1# 2 0 1 2 3 5 3 5 7 7 6 7 8 9 9.

The rearrangement

0 0 1 2 1 1 2 3 5 3 5 7 7 6 7 8 9 9

corresponds to the three reassessments

0	0	1	2	1*	1#	2	3	5	3	5	7	7	6	7	8	9	9
0	0	1	2	1	1#	2	3	5	3	5	7	7	6	7	8	9	9
0	0	1	2	1	1*	2	3	5	3	5	7	7	6	7	8	9	9

The 18 observations are divided into four dose groups containing 4, 5, 4, and 5 observations, respectively, so that there are  $\binom{18}{4, 5, 4, 5}$  possible reassessments of observations to dose groups. Each reassignment has probability  $\frac{1}{\binom{18}{4, 5, 4, 5}}$  of occurring, so the probability of the rearrangement

0 0 1 2 1 1 2 3 5 3 5 7 7 6 7 8 9 9

is

$$\frac{3}{\binom{18}{4 \ 5 \ 4 \ 5}}.$$

To determine the significance level when there are ties, weight each distinct rearrangement by its probability of occurrence. This weighting is done automatically if you use Monte Carlo sampling methods as is done in the computer program we provide in Section 3.5.1.

#### Capsule Summary

##### *K*-SAMPLE TEST

*H*: all distributions and, in particular,  
all population means the same

*K1*: at least one pair of means differ

*K2*: the population means are ordered

##### Assumptions

- 1) exchangeable observations
- 2)  $F_{ij}(x) = F(x - \mu_i)$

Transform None

##### Test statistic

*K1*:  $\sum n_i(X_i)^2$

*K2*:  $\sum f[i]n_iX_i$

#### 3.5.4. Cochran–Armitage Test

If our observations are binomial response variables, that is, they can only take the values 0 or 1,  $f[i] = d_i$  the magnitude of the *i*th dose, and  $X_i$  denotes the number of responders in the *i*th dose group, then the Pitman correlation,  $\sum d_i X_i$  is more commonly known as the Cochran–Armitage test for trend.

More formally, let  $p_i$  be the probability that an individual in the *i*th dose group will respond to the dose. Our hypothesis *H* is that  $p_1 = p_2 = \dots = p_I$ , and the alternative *K* is that  $p_1 \leq p_2 \leq \dots \leq p_I$  with strict inequality holding in at least one case. If we assume the response probabilities satisfy a logistic regression model, that is,

$$\log \frac{p_i}{1 - p_i} = \gamma + \lambda d_i,$$

then the Cochran–Armitage test is locally most powerful, Armitage [1955].

Bounds on the power can be obtained by the method of Mehta, Patel, and Senchaudhuri [C1998].

### 3.5.5. Linear Estimation

The Pitman correlation may be generalized by replacing the fixed function  $f[i]$  by an estimate  $\hat{\phi}$  derived by a linear estimation procedure, such as least-squares polynomial regression, kernel estimation, local regression, and smoothing splines [Raz, 1990].

Suppose the  $j$ th treatment group is defined by  $x_j$ , a vector-valued design variable ( $x_j$  might include settings for temperature, humidity, and phosphorous concentration). Suppose also that we may represent the  $i$ th observation in the  $j$ th group by a regression model of the form

$$Y_{ji} = \mu(x_j) + e_{ji}, \quad j = 1, \dots, n,$$

where  $e$  is an error variable with mean 0, and  $\mu(x)$  is a smooth regression function (that is, for any  $x$  and  $\varepsilon$  sufficiently small,  $\mu(x + \varepsilon)$  may be closely approximated by the first-order Taylor expansion  $\mu(x) + b\varepsilon$ ).

The null hypothesis is that  $\mu(x) = \mu$ , a constant that does not depend on the design variable  $x$ . As always, we assume that the errors  $e_{ji}$  are exchangeable so that all  $n!$  assignments of the labels to the observations that preserve the sample sizes  $\{n_j\}$  are equally likely.

Raz's test statistic is  $Q = \sum(\hat{\mu}(x_j))^2$  where  $\hat{\mu}$  is an estimate of  $\mu$  derived by a linear estimation procedure such as least-squares polynomial regression, kernel estimation, local regression, and smoothing splines.

This test may be performed using the permutation distribution of  $Q$  or, for large samples, a gamma-distribution approximation. See also Section 7.3.

### 3.5.6. A Unifying Theory

The permutation tests for Pitman correlation and the two-sample comparison of means are really special cases of a more general class of tests that take the form of a dot product of two vectors [Wald and Wolfowitz, 1943; De Cani [1979]. Let  $W = \{W_1, \dots, W_N\}$  and  $Z = \{Z_1, \dots, Z_N\}$  be fixed sets of number and let  $z = \{z_1, \dots, z_N\}$  be a random permutation of the elements of  $Z$ . Then we may use the dot product of the vectors  $z$  and  $W$ ,  $T = \sum z_i w_i$ , to test the hypothesis that the labeling is irrelevant. In the two-sample comparison,  $W$  is a vector of  $m$  1's followed by  $n$  0's. In Pitman correlation,  $W = \{f[1], \dots, f[N]\}$  where  $f$  is a monotone function.

## 3.6. Blocking

Although the significance level of a permutation test may be “distribution-free,” its power strongly depends on the underlying distribution.

Figure 3.2 depicts the effect of a change in the variance of the underlying population on the power of the permutation test for the difference in two means. As

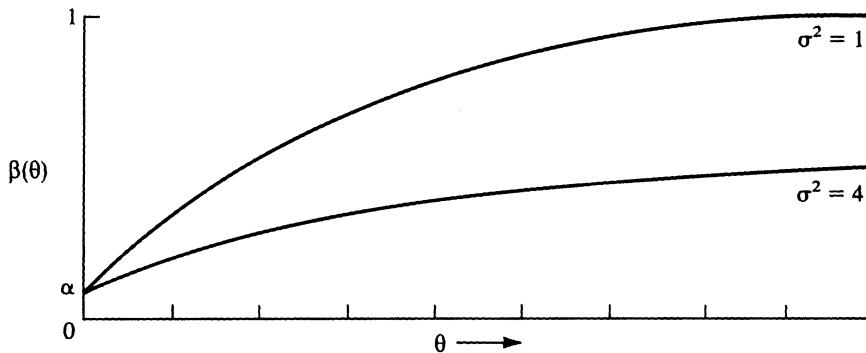


Figure 3.2. Effect of the population variance on the power of a test of two means.  
 $\theta = \theta_1 - \theta_2$ .

the variance increases, the power decreases. *To get the most from your experiments, reduce the variance.*

One way to reduce the variance is to subdivide the population under study into more homogeneous subpopulations and to take separate samples from each. Suppose you were designing a survey on the effect of income level on the respondents' attitudes toward compulsory pregnancy. Obviously, the views of men and women differ markedly on this controversial topic. It would not be prudent to rely on randomization to even out the sex ratios in the various income groups.

The recommended solution is to block the experiment, to interview, and to report on men and women separately. You would probably want to do the same type of blocking in a medical study. Similarly, in an agricultural study, you would want to distinguish among clay soils, sandy, and sandy-loam.

In short, whenever a population can be subdivided into distinguishable subpopulations, you can reduce the variance of your observations and increase the power of your statistical tests by blocking or stratifying your sample.

Suppose we have agreed to divide our sample into two blocks—one for men, one for women. If this is an experiment, rather than a survey, we would then assign subjects to treatments separately within each block.

In a study that involves two treatments and ten experimental subjects, four men and six women, we would first assign the men to treatment and then the women. We could assign the men in any of  $(4)$  = 6 ways and the women in any of  $(6)$  = 20 ways. That is, there are  $6 \times 20 = 120$  possible random assignments in all.

When we come to analyze the results of our experiment, we use the permutation approach to ensure we analyze in the way the experiment was designed. Our test statistic is a natural extension of that used for the two-sample comparison [Lehmann, 1986, pp. 233–4]:

$$S = \sum_{b=1}^B \sum_{j=m_b+1}^{(n_b+m_b)} x_{bj}, \quad (3.6.1)$$

where  $B$  is the number of blocks, two in the present example, and the inner sum extends over the  $n_b$  treated observations  $x_{bj}$  within each block.

We compute the test statistic for the original data. Then, we rearrange the observations at random within each block, subject to the restriction that the number of observations within each treatment category—the pair  $\{n_b, m_b\}$ —remain constant.

We compute  $S$  for each of the 120 possible rearrangements. If the value of  $S$  for the original data is among the  $120\alpha$  largest values, then we reject the null hypothesis; otherwise we accept it.

### 3.6.1. Extending the Range of Applications

The resulting permutation test is exact and most powerful against normal alternatives even if the observations on men and women have different distributions [Lehmann, 1986]. As we saw in Section 2.3, all that is required is that the subsets of errors be exchangeable.

The design need not be balanced. The test statistic  $S$  (equation 3.6.1) is a sum of sums. Unequal sample sizes resulting from missing data or an inability to complete one or more portions of the experiment will affect the analysis only in the relative weights assigned to each subgrouping.

**Warning:** This remark applies only if the data are missing at random. If treatment-related withdrawals are a problem in one of your studies, see Entsuah [1990] for the details of a resampling procedure.

Blocking is applicable to any number of subgroups; in the extreme case, that in which every pair of observations forms a distinct subgroup, we have the case of matched pairs.

## 3.7. Matched Pairs

In a matched pairs experiment, we take blocking to its logical conclusion. Each subject in the treatment group is matched as closely as possible by a subject in the control group. For example, if a 45-year old black male hypertensive is given a blood-pressure lowering pill, then we give a second, similarly built 45-year old black male hypertensive a placebo. One member of each pair is then assigned at random to the treatment group, and the other member is assigned to the controls.

Assuming we have been successful in our matching, we will end up with a series of independent pairs of observations  $(X_i, Y_i)$ , in which the members of each pair have been drawn from the distributions  $F_i[x - \mu_i]$  and  $F_i[x - \mu_i - \delta]$ , respectively. Note that it is the pairs that are independent; the components of each pair  $X_i, Y_i$  have the parameter  $\mu_i$  in common, and it is this commonality we count on to help us reduce the unwanted variability.

Regardless of the form of the unknown distribution  $F$ , the differences  $Z_i = Y_i - X_i$  will be symmetrically distributed about the unknown parameter  $\delta$ . For noting that  $Y' = Y - \mu - \delta$  and  $X' = X - \mu$  have the identical distribution  $F$ , we

see that

$$\begin{aligned}
 \Pr(Z \leq z + \delta) &= \Pr\{Y - X - \delta \leq z\} \\
 &= \Pr\{(Y - \mu - \delta) - (X - \mu) \leq z\} \\
 &= \Pr(Y' - X' \leq z) \\
 &= \Pr(X' - Y' \leq z) \\
 &= \Pr\{(X - \mu) - (Y - \mu - \delta) \leq z\} \\
 &= \Pr\{X - Y + \delta \leq z - \delta\} \\
 &= \Pr\{-Z \leq z - \delta\} \\
 &= \Pr\{Z \geq -z + \delta\}.
 \end{aligned}$$

This is precisely the case we considered at the beginning of this chapter, and the same readily computed permutation test is applicable.

This permutation test has the same properties of exactness, unbiasedness, and sensitivity under the same conditions as the one-sample test with the exception that, if the observation on one member of a pair is missing, then we must discard the remaining observation (see question 9 in Section 3.8).

For an almost most powerful test when one member of the pair is censored, see Section 9.5. For an application of a permutation test to the case in which an experimental subject serves as her own control, see Shen and Quade [1986].

#### Capsule Summary

##### MATCHED-PAIRS

$H$ : distributions and, in particular, means/medians of the members of each pair are the same

$K$ : means/medians of the members of each pair differ by  $d > 0$

##### Assumptions

- 1) independent observations
- 2)  $F_{1i}(x) = F_{2i}(x - d)$

##### Transform

$$z_i = x_{1i} - x_{2i}$$

##### Test statistic

Sum of positive  $z_i$

## 3.8. Questions

1. Show that the following statistics lead to equivalent permutation tests for the equality of two location parameters:
  - a)  $\sum X_{2i}$  (our original choice)
  - b)  $\sum X_{2i}/n_2 - \sum X_{1i}/n_1$  (the difference of the sample means)
  - c)  $\frac{(\sum X_{2i}/n_2 - \sum X_{1i}/n_1)}{\sqrt{(\sum (X_{2i} - \bar{X}_{2.})^2 + \sum (X_{1i} - \bar{X}_{1.})^2)/(m+n-2)}}$  (the  $t$ -statistic).

Hint: The sums ( $\sum X_{2i} + \sum X_{1i}$ ), ( $\sum X_{2i}^2 + \sum X_{1i}^2$ ) and the sample sizes  $n_1, n_2$  are invariant under permutations.

2. How was the analysis of the cell culture experiment described in Section 1.2 affected by the loss of two of the cultures due to contamination? Suppose these cultures had escaped contamination and given rise to the observations 90 and 95; what would be the results of a permutation analysis applied to the new, enlarged data set consisting of the following cell counts:

Treated	121	118	110	90
Untreated	95	34	22	12?

3. In the example of Section 3.3.2, list all rearrangements in which the sum of the ranks in the first sample is less than or equal to the original sum.
4. Use both the  $F$ -ratio and Pitman correlation to analyze the data for micronuclei in Table 3.1. Explain the difference in results.
5. The following vaginal virus titres were observed in mice by H.E. Renis of the Upjohn Company 144 hours after inoculation with Herpes virus Type II; see Good [1979] for complete details.

Saline controls	10000, 3000, 2600, 2400, 1500.
Treated with antibiotic	9000, 1700, 1100, 360, 1.

Is this a one-sample, two-sample, k-sample, or matched pairs study?

Does treatment have an effect?

Most authorities would suggest using a logarithmic transformation before analyzing this data. Repeat your analysis after taking the logarithm of each of the observations. Is there any difference?

Compare your results and interpretations with those of Good [1979].

6. Using the logarithm of the viral titre, determine an approximate 90% confidence interval for the treatment effect. (Hint: Keep subtracting a constant from the logarithms of the observations on saline controls until you can no longer detect a treatment difference.)
7. Suppose you make a series of  $I$ -independent pairs of observations  $\{x_i, y_i; i = 1 \dots I\}$ .  $y_i$  might be tensile strength and  $x_i$  the percentage of some trace metal. You know from your previous work that each of the  $y_i$  has a symmetric distribution.
- How would you test the hypothesis that for all  $i$ , the median of  $y_i$  is  $x_i$ ? (Hint: See Section 3.1.2.)
  - Do you need to assume that the distributions of the  $\{y_i\}$  all have the same shape, i.e., that they are all normal or all double exponential? Are the  $\{y_i\}$  exchangeable? Are the differences  $\{y_i - x_i\}$ ? (We return to these questions in Chapter 7.)
8. Show that the permutation test introduced in Section 3.4. for comparing variances based on deviations about the sample medians is asymptotically exact.
9. Without thinking, you analyze your matched pairs as if they were simply two independent samples and obtain a significant result. Consequently, you decide not to waste time analyzing the data correctly. Right or wrong? [Hint:  $X_i$  and  $Y_i$  are correlated: Why else would you have used matched pairs? Now, suppose some of the data is missing.]

## CHAPTER 4

# Experimental Designs

### 4.1. Introduction

In this chapter, we explore the use of permutation methods for analyzing the results of complex experimental designs that may involve multiple control variables, covariates, and restricted randomization.

### 4.2. Balanced Designs

The analysis of randomized blocks we studied in Chapter 3 can be generalized to very complex experimental designs with multiple control variables and confounded effects. In this section, we consider the evaluation of main effects and interactions in the two- and three-way univariate analysis of variance and in the Latin Square. Only *balanced* designs with the sample sizes equal in all subcategories are considered here. Unbalanced designs are considered in Section 4.4.

What distinguishes the complex experimental design from the simple one-sample, two-sample, and  $k$ -sample experiments we have considered so far is the presence of *multiple* control factors.

For example, we may want to assess the simultaneous effects on crop yield of hours of sunlight and rainfall. We determine to observe the crop yield  $X_{ijm}$  for  $I$  different levels of sunlight,  $i = 1, \dots, I$ , and  $J$  different levels of rainfall,  $j = 1, \dots, J$ , and to make  $M$  observations at each factor combination,  $m = 1, \dots, M$ . We adopt as our model relating the dependent variable, crop yield (the effect) to the independent variables of sunlight and rainfall (the causes)

$$X_{ijm} = \mu + s_i + r_j + (sr)_{ij} + \varepsilon_{ijm}.$$

In this model, terms with a single subscript like  $s_i$ , the effect of sunlight, are called *main effects*. Terms with multiple subscripts like  $(sr)_{ij}$ , the residual and nonadditive effect of sunlight and rainfall, are called *interactions*. The  $\{\varepsilon_{ijm}\}$  represent that portion of crop yield that cannot be explained by the independent variables

alone; these are variously termed the residuals, the errors, or the model errors. To ensure the residuals are exchangeable so that permutation methods can be applied, the experimental units must be assigned at random to treatment (see Section 4.2.4).

If we wanted to assess the simultaneous effect on crop yield of three factors simultaneously—sunlight, rainfall, and fertilizer—we would observe the crop yield  $X_{ijkm}$  for I different levels of sunlight,  $i = 1, \dots, I$ ; J different levels of rainfall,  $j = 1, \dots, J$ ; and K different levels of fertilizer,  $k = 1, \dots, K$ ; and make M observations at each factor combination,  $m = 1, \dots, M$ . Our model would then be

$$X_{ijkm} = \mu + s_i + r_j + f_k + (sr)_{ij} + (sf)_{ik} + (rf)_{jk} + (srf)_{ijk} + \epsilon_{ijkm}.$$

In this model we have three main effects,  $s_i$ ,  $r_j$ , and  $f_k$ , three two-way interactions,  $(sr)_{ij}$ ,  $(sf)_{ik}$ ,  $(rf)_{jk}$ , a single three-way interaction,  $(srf)_{ijk}$ , and the error term  $\epsilon_{ijkm}$ .

Including the additive constant  $\mu$  in the model allows us to define all main effects and interactions so they sum to zero,  $\sum s_i = 0$ ,  $\sum_i (sr)_{ij} = \sum_j (sr)_{ij} = 0$ , and so forth. Thus, under the null hypothesis of no effect of sunlight on crop yield, each of the main effects  $s_1 = s_2 = \dots = s_I = 0$ . Under the alternative, the different terms  $s_i$  represent deviations from a zero average, with the interaction term  $(sr)_{ij}$  representing the deviation from the sum  $s_i + r_j$ .

Clearly, when we have multiple factors, we must also have multiple test statistics. In the preceding example, we require three separate tests and test statistics for the three main effects  $s_i$ ,  $r_j$ , and  $f_k$ , plus four other statistical tests for the three two-way and the one three-way interactions. Will we be able to find statistics that measure a single intended effect without confounding it with a second unrelated effect? Will the several test statistics be independent of one another?

In the permutation analysis of an experimental design as in the parametric analysis of variance, the answer is yes to both questions only if the design is balanced, that is, if there are equal numbers of observations in each subcategory, and if the test statistics are independent of one another.

In a balanced design, the permutation test has a threefold advantage over the parametric ANOVA: It is exact; it is not restricted by an assumption of normality—although, it does require that the experimental errors be exchangeable (see Section 2.2), yet it is as or more powerful than the parametric approach (see Scheffe, 1959; Collier and Baker, 1966; and Bradbury, 1987).

In an unbalanced design, main effects will be confounded with interactions so that the two cannot be tested separately, a topic we return to in Section 4.4.

### 4.2.1. Main Effects

In a  $k$ -way analysis with equal sample sizes  $M$  in each category, we can assess the main effects using essentially the same statistics we would use for randomized blocks. Take sunlight in the preceding example. If we have only two levels of sunlight, then, referring to equation 3.6.1, our test statistic for the effect of

sunlight is

$$S = \sum_{j=1}^J \sum_{k=1}^K \sum_{m=1}^M X_{1jkm}. \quad (4.1)$$

If we have more than two levels of sunlight, our test statistic is

$$F2 = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (X_{ijk\cdot} - X_{\cdot jk\cdot})^2 \quad (4.2)$$

or

$$F1 = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K |X_{ijk\cdot} - X_{\cdot jk\cdot}|. \quad (4.3)$$

The dot  $\cdot$  used as a subscript indicates that we have summed over the corresponding subscript and then taken an average by dividing by the number of terms in that sum; thus,

$$X_{ijk\cdot} = \frac{1}{M} \sum_{m=1}^M X_{ijkm}/M.$$

The statistics F2 and F1 offer protection against a broad variety of shift alternatives, including

$$K_1 : s_1 = s_2 > s_3 = \dots$$

$$K_2 : s_1 > s_2 > s_3 = \dots$$

$$K_3 : s_1 < s_2 > s_3 = \dots$$

As a result, they may not provide a most powerful test for any single one of these alternatives. If we believe the effect to be monotone increasing, then, in line with the thinking detailed in Section 3.5.2, we would use the Pitman correlation statistic

$$R = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K f[i](X_{ijk\cdot} - X_{\cdot jk\cdot}). \quad (4.4)$$

To obtain the permutation distributions of the test statistics  $S$ ,  $F2$ ,  $F1$ , and  $R$ , we permute the observations independently in each of the  $JK$  blocks determined by a specific combination of rainfall and fertilizer. Exchanging observations *within* a category corresponding to a specific level of sunlight leaves the statistics  $S$ ,  $F2$ ,  $F1$ , and  $R$  unchanged. We can concentrate on exchanges *between* levels, and the total number of rearrangements is  $\binom{MI}{M...M}^{JK}$ . As we are only exchanging observations between levels of sunlight, then  $\sum_{j=1}^J \sum_{k=1}^K X_{\cdot jk\cdot}$  and  $\sum_{j=1}^J \sum_{k=1}^K (X_{\cdot jk\cdot})^2$  are

invariant, and after their elimination, F2 reduces to

$$\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left( \sum_{m=1}^M X_{ijkm} \right)^2.$$

Similarly, F1 reduces to

$$\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K \left| \sum_{m=1}^M X_{ijkm} \right|.$$

We compute the test statistic ( $S$ ,  $F1$ , or  $R$ ) for each rearrangement, rejecting the hypothesis that sunlight has no effect on crop yield only if the value of  $S$  (or  $F1$  or  $R$ ) that we obtain using the original arrangement of the observations lies among the  $\alpha$  most extreme of these values.

Of the two  $F$ -statistics, F1 is to be preferred to F2. F1 is as powerful or more powerful for detecting location shifts and more powerful for detecting concentration changes [Mielke and Berry, 1983].

A third alternative to F1 and F2 is

$$F3 = \sum_j \frac{n_j(n_j - 1)(X_{j\cdot} - X_{\cdot\cdot})^2}{\sum_k (X_{jk} - X_{j\cdot})^2}, \quad (4.5)$$

(James, 1951) which Hall [1989] recommends for use with the bootstrap when we cannot be certain that the observations in the various categories all have the same variance. In simulation studies with permutation tests and variances that differed by an order of magnitude, I found F3 was inferior to F1.

A final alternative to the statistics  $S$ ,  $F1$ , and  $F2$  is the standard  $F$ -ratio statistic

$$F = \frac{\sum_{i=1}^I M_i (X_{i\cdot\cdot} - X_{\cdot\cdot\cdot})^2}{(I-1)\hat{\sigma}^2}, \quad (4.6)$$

where  $\hat{\sigma}^2$  is our estimate of the variance of the errors  $\varepsilon_{ijk}$ . But if we use  $F$ , we are forced to consider exchanges between as well as within blocks, thus negating the advantages of blocking as described in Section 3.6.

### 4.2.2. An Example

In this section, we apply the permutation method to determine the main effects of sunlight and fertilizer on crop yield using the data from the two-factor experiment depicted in Table 4.1a. As there are only two levels of sunlight in this experiment, we use  $S$  (equation 4.1) to test for the main effect. For the original observations,  $S = 23 + 55 + 75 = 153$ . One possible rearrangement is shown in Table 4.1b in which we have interchanged the two observations marked with an asterisk, the 5 and 6. The new value of  $S$  is 154.

Table 4.1a. Effect of Sunlight and Fertilizer on Crop Yield

S	Fertilizer			
	LO	MED	HIGH	
u	LO	5	15	21
n		10	22	29
l		8	18	25
i				
g	HI	6	25	55
h		9	32	60
t		12	40	48

Table 4.1b. Effect of Sunlight and Fertilizer. Data Rearranged

	LO	MED	HIGH
LO	6*	15	21
	10#	22	29
	8	18	25
HI	5*	25	55
	9#	32	60
	12	40	48

As can be seen by a continuing series of straightforward hand calculations, the test statistic,  $S$ , for the main effect of sunlight is as small or smaller than it is for the original observations in only 8 out of the  $\binom{6}{3}^3 = 8000$  possible rearrangements. For example, it is smaller when we swap the 9 of the Hi-Lo group for the 10 of the Lo-Lo group (the two observations marked with the pound sign). As a result, we conclude that the effect of sunlight is statistically significant.

The computations for the main effect of fertilizer are more complicated—we must examine  $\binom{9}{3,3,3}^2$  rearrangements, and compute the statistic  $F_1$  for each. We use  $F_1$  rather than  $R$  because of the possibility that too much fertilizer—the “High” level, might actually suppress growth. Only a computer can do this many calculations quickly and correctly, so we adapted our program from Section 3.5 to make them (see Sidebar). The estimated significance level is .001 and we conclude that this main effect, too, is statistically significant.

In this last example, each category held the same number of experimental subjects. If the numbers of observations were unequal, our main effect would have been confounded with one or more of the interactions (see Section 4.5). In contrast to the simpler designs we studied in the previous chapter, missing data will affect our analysis.

**Sidebar**

Program for estimating significance level of the main effect of fertilizer on crop yield in a balanced design

Set aside space for

Monte	the number of Monte Carlo simulations
$S_0$	the original value of test statistic
$S$	test statistic for rearranged data
data	{5, 10, 8, 15, 22, 18, 21, 29, 25, 6, 9, 12, 25, 32, 40, 55, 60, 48}
$n = 3$	number of observations in each category
blocks = 2	number of blocks
levels = 3	number of levels of factor

Main program

Get data

put all the observations into a single linear vector

Compute  $S_0$  for the original observations

Repeat Monte times

for each block

Rearrange the data in the block

Compute  $S$

Compare  $S$  with  $S_0$

Print out the proportion of times  $S$  was larger than  $S_0$

Rearrange

Set  $s$  to the number of observations in the block

Start: Choose a random integer  $k$  from 0 to  $s - 1$

Swap  $X[k]$  and  $X[s - 1]$

Decrement  $s$  and repeat from start

Stop after you've selected all but one of the samples

Get data

user-written procedure gets data and packs them into a two-dimensional array in which each row corresponds to a block

Compute

$$F1 = \sum_{i=1}^I \sum_{j=1}^J |X_{ij}|$$

for each block

calculate the mean of that block

for each level within a block

calculate the mean of that block-level

calculate difference from block mean

### 4.2.3. Testing for Interactions

In the preceding analysis of main effects, we assumed the effect of sunlight was the same regardless of the levels of the other factors. To test the validity of this assumption, we might attempt to eliminate row and column effects by subtracting

the row and column means from the original observations. That is, we set

$$X'_{ijk} = X_{ijk} - \bar{X}_{i..} - \bar{X}_{.j.} + \bar{X}_{...},$$

where, by adding the grand mean,  $\bar{X}_{...}$ , we ensure the overall sum will be zero. In the example of the effect of sunlight and fertilizer on crop yield, we are left with the residuals shown in Table 4.2.

Table 4.2. Effect of Sunlight and Fertilizer on Crop Yield. Testing for Non-additive Interaction

		Fertilizer		
		LO	MED	HIGH
S				
U	LO	4.1	-2.1	-11.2
N		9.1	4.1	-3.2
L		7.1	0.1	-7.2
I				
G	HI	-9.8	-7.7	7.8
H		-7.8	-0.7	12.8
T		-3.8	7.2	0.8

The pattern of plus and minus signs in this table of residuals suggests that fertilizer and sunlight affect crop yield in a superadditive fashion. Note the minus signs associated with the mismatched combinations of a high level of sunlight and a low level of fertilizer and a low level of sunlight with a high level of fertilizer. To encapsulate our intuition in numeric form, we follow Still and White [1981] to sum the deviates from the mean within each cell, square the sum, and then sum the squares to form the test statistic

$$I = \sum_i \sum_j \left( \sum_k X'_{ijk} \right)^2.$$

We compute the deviates and this test statistic for each rerandomization of the residuals. In most, the values of the test statistic are close to zero as the entries in each cell cancel. The value of the test statistic for our original data,  $I = 2127.8$ , stands out as an exceptional value, but how are we to interpret it?

Recall that  $X_{ijk} = \mu + s_i + r_j + (sr)_{ij} + \epsilon_{ijk}$ ; the residual  $X'_{ijk} = \epsilon_{ijk} - \epsilon_{i..} - \epsilon_{.j.} + \epsilon_{...}$ ; and, thus, the residuals are weakly correlated, though the  $\{\epsilon_{ijk}\}$  themselves may be independent. The test based on  $I$  is not exact (but see question 6 at the end of this chapter) and any significance level is at best an approximation.

#### 4.2.3.1. Blocking

Assume now that  $I = 2$  and  $J = 2$  and that we have only a single observation on each treatment combination ( $K = 1$ ), but that we are able to replicate this experiment  $B$

times, for  $b = 1, 2, \dots, B$ . These replications may have been at different times, or on different subjects or plots. Let  $X_{bij}$  denote the result of applying the treatment combination  $i, j$  in the  $b$ th block.

It is easy to see that  $D_{bI} = (X_{b11} - X_{b21}, X_{b12} - X_{b22})$ ,  $D_{bJ} = (X_{b11} - X_{b12}, X_{b21} - X_{b22})$ , and  $D_b = X_{b11} - X_{b12} + X_{b22} - X_{b21}$  provide estimates of the main effects and interaction in the  $b$ th block. Under the various null hypotheses, the distribution of these estimates is symmetric about zero. To obtain a test for the interaction, we may proceed as in Section 3.1 to consider the statistic  $S = \sum_{\{D_b > 0\}} D_b$  for each of the  $2^B$  permutations  $\pm D_1, \dots, \pm D_B$ . If we accept the null hypothesis, we may then proceed to test for the presence of main effects. A test for the row effect might be based on  $S = \sum_{\{D_{bIj} > 0\}} D_{bIj}$  for each of the  $2^{2B}$  permutations  $\pm D_{111}, \pm D_{112}, \dots, \pm D_{B11}, \pm D_{B12}$ .

A formal derivation of these results due to Welch [1990] is provided in Section 10.3.3.

#### 4.2.3.2. Synchronized Permutations

Pesarin [1999] developed an exact test for nonzero interactions using synchronized permutations. We consider here only the simple case of a  $2^2$  factorial design, that is,  $I = J = 2$ . We begin by developing statistics for comparing the first factor at each of the two levels of the second.  $I_{A|1} = \sum_k X_{11k} - \sum_k X_{21k}$  and  $I_{A|2} = \sum_k X_{12k} - \sum_k X_{22k}$ . To obtain a synchronized permutation, we select  $v$  elements at random from the group of  $K$  for which  $i = 1$  and  $j = 1$  and exchange these with  $v$  elements selected at random from the group for which  $i = 1$  and  $j = 2$ , while at the same time exchanging  $v$  elements selected at random from the group for which  $i = 2$  and  $j = 1$ , with  $v$  elements selected at random from the group for which  $i = 2$  and  $j = 2$ .

Recall that

$$X_{ijk} = \mu + s_i + r_j + (sr)_{ij} + \epsilon_{ijk}.$$

As a result, our two statistics take on the values

$$I_{A|1}^* = 4(K - v)(s + (sr)) + 2(K - v)(\epsilon_{11} - \epsilon_{21}).$$

and

$$I_{A|2}^* = 4(K - v)(s - (sr)) + 2(K - v)(\epsilon_{12} - \epsilon_{22}).$$

We set  $I_{AB} = I_{A|1} - I_{A|2}$  to obtain a test of the interaction that is independent of the main effect  $s$  and, as important, that is independent of the permutation test of that main effect based on the statistic  $I_A = I_{A|1} + I_{A|2}$ .

Similarly, we define statistics for the second factor at each of the two levels of the first factor,  $I_{B|1} = \sum_k X_{11r} - \sum_k X_{12r}$  and  $I_{A|2} = \sum_k X_{21r} - \sum_k X_{22r}$ , and obtain test statistics  $I_{BA} = I_{B|1} - I_{B|2}$  for the interaction and  $I_{BA} = I_{B|1} + I_{B|2}$  for the main effect of the second factor. Our test statistic for the interaction is  $I = I_{AB} + I_{BA}$ .

Hi	Med	Lo
Hi	Med	Lo
Hi	Med	Lo
<b>a</b>		
Hi	Med	Lo
Lo	Lo	Med
Hi	Hi	Med
<b>b</b>		
Hi	Med	Lo
Lo	Hi	Med
Med	Lo	Hi
<b>c</b>		

Figure 4.1. a) Systematic assignment of fertilizer levels to plots; b) random assignment of fertilizer levels to plots; c) Latin Square assignment of fertilizer levels to plots.

#### 4.2.3.3. To Learn More

For more information concerning permutation tests for interactions, see Cox [1984], Loughin and Noble [1997], Kennedy and Cade [1996], Manly [1997], Sprent [1998], ter Braak [1992], Welch [1990], and Westfall and Young [1993].

#### 4.2.4. Designing an Experiment

All the preceding results are based on the assumption that the assignment of treatments to plots (or subjects) is made at random. While it might be convenient to fertilize our plots as shown in Figure 4.1a, the result could be a systematic bias, particularly if, for example, there is a gradient in dissolved minerals from east to west across the field.

The layout adopted in Figure 4.1b, obtained with the aid of a computerized random number generator, reduces but does not eliminate the effects of this hypothetical gradient. Because this layout was selected at random, the exchangeability of the error terms and, hence, the exactness of the corresponding permutation test is assured. Unfortunately, the layout of Figure 4.1a with its built-in bias can also result from a random assignment; its selection is neither more nor less probable than any of the other ( ${}^9_{333}$ ) possibilities.

		Factor 1		
		1	2	3
F a c t o r 2	1	A	B	C
	2	B	C	A
	3	C	A	B

Figure 4.2. A Latin Square.

What can we do to avoid such an undesirable event? In the layout of Figure 4.1c, known as a Latin Square, each fertilizer level occurs once and once only in each row and in each column; if there is a systematic gradient of minerals in the soil, then this layout ensures that the gradient will have almost equal impact on each of the three treatment levels. It will have an almost equal impact even if the gradient extends from northeast to southwest rather than from east to west, or north to south. I use the phrase “almost equal” because a gradient effect may still persist. The design and analysis of Latin Squares is described in the next section.

To increase the sensitivity of your experiments and to eliminate any systematic bias, I recommend you use the following three-step procedure during the design phase:

- 1) List all the factors you feel may influence the outcome of your experiment.
- 2) Block all factors which are under your control; this process is described in Section 3.6. You may want to use some of these factors to restrict the scope of your experiment, e.g., eliminate all individuals under 18 and over 60.
- 3) Randomly assign units to treatment within each block. See also, Maxwell and Cole [1991].

#### 4.2.5. Latin Square

The Latin Square considered in Section 4.2.4 is one of the simplest examples of an experimental design in which the statistician takes advantage of some aspect of the model to reduce the overall sample size.

A Latin Square is a three-factor experiment in which each combination of factors occurs once and once only. We can use a Latin Square as in Figure 4.2 to assess the effects of soil composition on crop yield:

In this diagram, Factor 1—gypsum concentration, for example—is increasing from left to right; Factor 2 is increasing from top to bottom (or from North to South); and Factor 3, its varying levels denoted by the capital letters A, B, and C, occurs in combination with the other two in such a way that each combination of factors—row, column, and treatment—occurs once and once only.

Because of this latter restriction, there are only 12 different ways in which we can assign the varying factor levels to form a  $3 \times 3$  Latin Square. Among the other

11 designs are

	1	2	3
1	A	C	B
2	B	A	C
3	C	B	A

and

1	C	B	A
2	B	A	C
3	A	C	B

Let us assume we begin our experiment by selecting one of these twelve designs at random and planting our seeds in accordance with the indicated conditions.

Because there is only a single replication of each factor combination in a Latin Square, we cannot estimate the interactions. Thus, the Latin Square is appropriate *only* if we feel confident in assuming that the effects of the various factors are completely additive, that is, that the interaction terms are zero.

Our model for the Latin Square is

$$X_{ijk} = \mu + s_i + r_j + f_k + \varepsilon_{ijk},$$

where, as always in a permutation analysis, we assume that the errors  $\varepsilon_{ijk}$  are exchangeable. Our null hypothesis is  $H: s_1 = s_2 = s_3$ . If we assume an ordered alternative,  $K: s_1 > s_2 > s_3$ , our test statistic for the main effect is similar to the correlation statistic employed in equation 4.4:

$$R = \sum_{i=1}^3 i(X_{i..} - X_{...})$$

or, equivalently, after eliminating the grand mean  $X_{...}$  which is invariant under permutations,

$$R1 = \sum_{i=-1}^1 iX_{i..} = X_{C..} - X_{A..}$$

We evaluate this test statistic both for the observed design and for each of the twelve possible Latin Square designs that might have been employed in this particular experiment. We reject the hypothesis of no treatment effect only if the test statistic for the original observations is an extreme value.

For example, suppose we employed Design 1 and observed

$$\begin{array}{ccc} 21 & 28 & 17 \\ 14 & 27 & 19 \\ 13 & 18 & 23 \end{array}$$

Then  $3y_{A..} = 58$ ,  $3y_{B..} = 65$ ,  $3y_{C..} = 57$  and our test statistic  $R1 = -1$ . Had we employed Design 2, then  $3y_{A..} = 71$ ,  $3y_{B..} = 49$ ,  $3y_{C..} = 65$ , and our test

statistic  $R1 = -6$ . With Design 3,  $3y_{A..} = 57$ ,  $3y_{B..} = 65$ ,  $3y_{C..} = 58$  and our test statistic  $R1 = +1$ .

We see from the permutation distribution obtained in this manner that the value of our test statistic for the design actually employed in the experiment,  $R1 = -1$ , is an average value, not an extreme one. We accept the null hypothesis and conclude that increasing the treatment level from A to B to C does not significantly increase the yield.

#### 4.2.6. Other Designs

If the three-step rule outlined in Section 4.2.4 leads to a more complex experimental design than those considered here, consult Kempthorne [1955]; Wilk and Kempthorne [1956, 1957]; and Scheffe [1959]. To correct for variables not under your control, see the next section.

### 4.3. Analysis of Covariance

#### 4.3.1. Covariates Defined

Some variables that affect the outcome of an experiment are under our control from the very beginning—e.g., light and fertilizer. But we may only be capable of measuring rather than controlling other equally influential variables, called covariates. Blood chemistry is an example of a covariate in a biomedical experiment. Various factors in the blood can affect an experimental outcome, and most blood factors will be affected by a treatment, but few are under our direct control.

In this section, we will discuss two methods for correcting for the effects of covariates. The first, eliminating the functional relationship, is for use when you know or suspect the nature of the functional relationship between the observables and the covariates. The second method, restricted randomization, is for use when the covariates take only a few discrete values and these values can be used to restrict the randomization.

#### 4.3.2. Eliminate the Functional Relationship

Gail, Tan, and Piantadosi [1988] recommend eliminating the effects of covariates first and then applying permutation methods to the residuals. For example, suppose the observation  $Y$  depends both on the treatment  $\tau_i (i = 1, \dots, l)$  and on the  $p$ -dimensional vector of covariates  $X = (X^1, \dots, X^p)$ , that is

$$Y = \mu + \tau + X\beta + e,$$

where  $Y$ ,  $\mu$ ,  $\tau$ , and  $e$  are  $n \times 1$  vectors of observations, mean values, treatment effects, and errors respectively,  $X$  is an  $n \times p$  matrix of covariate values, and  $\beta$  is a  $p \times 1$  vector of regression coefficients.

We would use least-squares methods to estimate the regression coefficients  $\hat{\beta}$  after which we would apply the permutation methods described in the preceding sections to the residuals  $Z = Y - X\hat{\beta}$ .

We use a similar approach in 4.2.3 in testing a two-factor model for a significant interaction. In that example, as here, we assume that the individual errors are exchangeable. A further assumption in the present case is that both the concomitant variables (the  $X$ 's) and the regression coefficients  $\beta$  are unaffected by the treatment [Kempthorne, 1952, p. 160].

A distribution-free multivariate analysis of covariance in which the effects of the treatments and the covariates are evaluated simultaneously is considered in the next chapter.

### 4.3.3. Selecting Variables

Which covariates should be included in your model? Draper and Stoneman [1966] describe a permutation procedure for selecting covariates using a *forward* stepping rule.

The first variable you select should have the largest squared sample correlation with the dependent variable  $y$ ; thereafter, include the variable with the largest squared partial correlation with  $y$  given the variables that have already been selected. You may use any standard statistics package to obtain these correlations. Equivalently, you may select variables based on the maximum value of the square of the  $t$ -ratio for the regression coefficient of the entering variable, the so-called “ $F$  to enter.” The problem lies in knowing when to stop, that is, in knowing when an additional variable contributes little beyond noise to the model.

Percentiles of the permutation distribution of the  $F$ -to-enter statistic can be used to test whether variables not yet added to the model would be of predictive value. Details for deriving the permutation distribution of this statistic defined in terms of Householder rotations of the permuted variable matrix are given in Forsythe et al. [1973].

### 4.3.4. Restricted Randomization

If the covariates take on only a few discrete values, e.g., smoker vs nonsmoker, or status 0, 1, or 2 we may correct for their effects by restricting the rerandomizations to those whose design matrices match the original [Edgington, 1983].

Consider the artificial data set in Table 4.3 adapted from Rosenbaum [1984, p. 568]. To test the hypothesis that the treatment has no effect on the response, we would use the sum of the observations in the treatment group as our test

Table 4.3. Data for Artificial Example

Subject	Treatment	Result	Covariate
A	1	6	1
B	1	2	0
C	0	5	1
D	0	4	1
E	0	3	1
G	0	1	0
H	0	0	0

statistic. The sum of 8 for the original observations is equaled or exceeded in six of the  $\binom{7}{2} = 21$  possible rerandomizations. This result is not statistically significant.

Now let us take the covariate into consideration. One member of the original treatment group has a covariate value of 0, the other has a covariate value of 1. We limit our attention to the  $12 = \binom{4}{1}\binom{3}{1}$  possible rerandomizations in which the members of the treatment group have similar covariate values. These consist of AB AG AH, CB CG CH, DB DG DH, EB EG EH. With only one of the 12, that of AB the original observations, do we observe a result sum as large as 8. This sum is statistically significant at the 0.1 level. Restricting the randomizations eliminates the masking effect of the covariate and reveals the statistically significant effect of the treatment.

If the covariate varies continuously, it may still be possible to apply the method of restricted randomizations by first subdividing the covariate's range into a few discrete categories. For example, if

$$\begin{aligned}x < -1 &\quad \text{let } x' = 0 \\-1 \leq x < 1 &\quad \text{let } x' = 1 \\1 \leq x &\quad \text{let } x' = 2.\end{aligned}$$

Rosenbaum [1984] suggests that with larger samples one should restrict the randomizations so that a specific mean value of the covariate is attained, rather than a specific set of values.

Subject to certain relatively weak assumptions, the method of restricted randomizations can also be applied to after-the-fact covariates. (See Section 9.2.)

## 4.4. Unbalanced Designs

The permutation test is not a panacea. Imbalance in the design will result in the confounding of main effects with interactions. Consider the following two-factor

model for crop yield:

$$X_{ijk} = \mu + s_i + r_j + sr_{ij} + \varepsilon_{ijk}.$$

$$\begin{array}{c|c} N(0,1) & N(2,1) \\ \hline N(2,1) & N(0,1) \end{array}$$

Now suppose that the observations in a two-factor experimental design are normally distributed as in the preceding diagram taken from Cornfield and Tukey [1956]. There are no main effects in this example—both row means and both column means have the same expectations—but there is a clear interaction represented by the two nonzero off-diagonal elements.

If the design is balanced, with equal numbers per cell, the lack of significant main effects and the presence of a significant interaction should and will be confirmed by our analysis. But suppose that the design is not in balance, that for every ten observations in the first column, we have only one observation in the second. Because of this imbalance, when we use the statistic  $S'$  (equation 4.1'), we will uncover a false “row” effect which is actually due to the interaction between rows and columns. The main effect is said to be *confounded* with the interaction.

If a design is unbalanced as in the preceding example, we cannot test for a “pure” main effect or a “pure” interaction. But we may be able to test for the combination of a main effect with an interaction by using the statistic ( $S'$ ,  $F1'$  or  $R'$ ) that we would use to test for the main effect alone. This combined effect will not be confounded with the main effects of other unrelated factors.

For three-factor designs with unequal sample sizes, the test statistics for mixed main/interaction effects are:

$$S' = \sum_{j=1}^J \sum_{k=1}^K \sum_{l=m_{jk}+1}^{n_{jk}+m_{jk}} X_{1jkl} \quad (4.1')$$

$$F1' = \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^I n_{ijk} |X_{ijk\cdot} - X_{\cdot jk\cdot}| \quad (4.3')$$

$$R' = \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^I f[i] n_{ijk} (X_{ijk\cdot} - X_{\cdot jk\cdot}). \quad (4.4')$$

#### 4.4.1. Missing Combinations

If an entire factor combination is missing, we may not be able to estimate or test any of the effects. One very concrete example is an unbalanced design I encountered in the 1970’s.

Makinodan et al. [1976] studied the effects of age on the mediation of the immune response. They measured the anti-SBRC response of spleen cells derived from C57BL mice of various ages. In one set of trials, the cells were derived entirely from the spleens of young mice, in a second set of trials, they came from the spleens of old mice, and in a third, they came from mixtures of the two.

Let  $X_{i,j,k}$  denote the response of the  $k$ th sample taken from a population of type  $i$ ,  $j$  ( $i = 1 = j$ : controls;  $i = 2, j = 1$ : cells from young animals only;  $i = 1, j = 2$ : cells from old animals only;  $i = 2 = j$ : mixture of cells from old and young animals.) We assume that for lymphocytes taken from the spleens of young animals,

$$X_{2,1,k} = \mu + \alpha + e_{2,1,k},$$

for the spleens of old animals,

$$X_{1,2,k} = \mu - \alpha + e_{1,2,k},$$

and for a mixture of  $p$  spleens from young animals and  $(1 - p)$  spleens from old animals, where  $0 \leq p \leq 1$ ,

$$\begin{aligned} X_{2,2,k} &= p(\mu + \alpha) + (1 - p)(\mu - \alpha) - \gamma + e_{2,2,k} \\ &= \mu + (1 - 2p)\alpha - \gamma + e_{2,2,k}, \end{aligned}$$

where the  $e_{i,j,k}$  are independent values.

Makinodan knew in advance of his experiment that  $\alpha > 0$ . He also knew that the distributions of the errors  $e_{i,j,k}$  would be different for the different populations. We can assume only that these errors are independent of one another and that their medians are zero.

Makinodan wanted to test the hypothesis  $\gamma = 0$  as there are immediate biological interpretations for the three alternatives: from  $\gamma = 0$  one may infer independent action of the two cell populations;  $\gamma < 0$  means excess lymphocytes in young populations; and  $\gamma > 0$  suggests the presence of suppressor cells in the spleens of older animals.

But what statistic are we to use to do the test? One possibility is

$$S = |X_{2,2,.} - pX_{1,2,.} - (1 - p)X_{2,1,.}|.$$

If the design were balanced, or we could be sure that the null effect  $\mu = 0$ , this is the statistic we would use. But the design is not balanced, with the result that the main effects (in which we are not interested) are confounded with the interaction (in which we are).

It is small consolation that the standard parametric (ANOVA) approach won't work in this example either. Fortunately, another resampling method, the bootstrap, can provide a solution.

Here is the bootstrap procedure.

Draw an observation at random and with replacement from the set  $\{x_{2,1,k}\}$ ; label it  $x_{2,1,j}^*$ . Similarly, draw the bootstrap observations  $x_{1,2,j}^*$  and  $x_{2,2,j}^*$  from the sets  $\{x_{1,2,k}\}$  and  $\{x_{2,2,k}\}$ .

$$\text{Let } z_j = px_{1,2,j}^* + (1 - p)x_{2,1,j}^* - x_{2,2,j}^*.$$

Repeat this resampling procedure a thousand or more times, obtaining a bootstrap estimate  $z_j$  of the interaction each time you resample. Use the resultant set of bootstrap estimates  $\{z_j\}$  to obtain a confidence interval for  $\gamma$ . If 0 belongs to this confidence interval, accept the hypothesis of additivity; otherwise reject.

One word of caution: unlike a permutation test, a bootstrap is exact only for very large samples. The probability of a Type I error may be greater than the significance level you specify.

### Sidebar

Mean DPFC response. Effect of pooled old BC3FL spleen cells on the anti-SRBC response of indicator pooled BC3FL spleen cells. Data extracted from Makinodan et al (1976). Bootstrap analysis.

Young Cells	Old Cells	$1/2 + 1/2$
5640	1150	7100
5120	2520	11020
5780	900	13065
4430	50	
7230		
Bootstrap sample 1:	$5640 + 900 - 11020$	-4480
Bootstrap sample 2:	$5780 + 1150 - 11020$	-4090
Bootstrap sample 2:	$7230 + 1150 - 7100$	1280
....	....	...
....	....	...
Bootstrap sample 600:	$5780 + 2520 - 7100$	1200

#### 4.4.2. The Boot-Perm Test

The preceding was an extreme example of an unbalanced design. More often, we will have at least a few observations in each category. In those circumstances, we may proceed as follows:

Bootstrap from the original data preserving categories so that the ab's are selected from the ab sample, the aB's from the aB sample, and so forth. Analyze the resultant balanced design using a permutation test. Bootstrap several times, say, five or ten. If you reject every time or accept every time, then draw the corresponding conclusion. Otherwise bootstrap 100 times. Check again. If you still have an ambiguity, then you probably have highly significant interaction and must obtain a balanced design to proceed further.

### 4.5. Clinical Trials

#### 4.5.1. Avoiding Unbalanced Designs

In preceding sections, we tacitly assumed that the assignment of subjects to treatment took place *before* the start of the experiment. We also assumed, tacitly, that the assignment of subjects to treatment was *double blind*, that is, neither the

experimental subject nor the experimenter knew which treatment the subject was receiving. (See Fisher [1951] and Feinstein [1972] for a justification of this double blind approach.) But in a large clinical trial covering several hundreds, even thousands of patients in various treatment categories, not all of the subjects will be available prior to the start of treatment. We even may have tabulated some of the results before the last of the patients have enrolled in the experiment. If we let pure chance determine whether an incoming patient is assigned to treatment or control, the trials may quickly go out of balance and stay out of balance. On the other hand, if we insist on keeping the experiment balanced at each stage, assigning subjects alternately to treatment and placebo, a physician could crack the code, guesstimate the next treatment assignment, and be influenced in her handling of a patient as a result.

One solution [Efron, 1971] is to weight the probability of a particular treatment assignment in accordance with the assignments that have already taken place. For example, if the last subject was assigned to the control group, we might increase the probability of assigning the current subject to the treatment from  $\frac{1}{2}$  to  $\frac{3}{4}$ . The assignment is still random—so no one can crack the code, but there will be a tendency for the two groups—treatment and control—to even out in size. Of course, Efron's biased coin approach is only one of many possible restricted designs. A general form is provided in Smith [1984].

While the numbers of subjects in the various treatment groups will (in theory) even out in the long run, in most cases they will still be unequal when the experiment is completed, taking values  $\{n_i, i = 1, \dots, l\}$ . Fortunately, we may analyze this experiment as if these were the sample sizes we had intended from the very beginning [Cox, 1982].

Following Hollander and Pena [1988], suppose there are  $R$  possible treatments. Let  $T_j = (T_{j1}, \dots, T_{jR-1})'$  be the treatment assignment vector for the  $j$ th patient;  $j = 1, \dots, n$ .  $T_{ji}$  is equal to 1 or 0 according to whether patient  $j$  is or is not assigned to treatment  $i$ . Let  $x_n = (x_1, \dots, x_n)'$  be the vector of patient responses (e.g., time to death, time to relapse). We want to test the null hypothesis that the  $R$  treatments are equivalent. The randomization distribution of the test statistic  $S_n = (T_1, \dots, T_n)x_n$  induced by the randomized treatment allocation grows increasingly more complicated with increasing  $n$ . Nevertheless, it may be determined by recursive means.

Smythe and Wei [1983] show that the permutation method can provide an exact test in the case of two treatments. Their result is extended to  $k$ -treatments by Wei, Smythe, and Smith [1986]. Algorithms for computing the exact distribution of the test statistic, rather than an asymptotic approximation, are provided by Hollander and Pena [1988] and Mehta, Patel, and Wei [1988].

### 4.5.2. Missing Data

A further and as yet unresolved problem in the analysis of clinical trials is the dropping out of patients during the course of the investigation. When such dropouts

occur at random, we still may apply any of the standard permutation methods, that is if we are prepared to deal with confounded effects (see Section 4.4). But what if the dropout rate is directly related to the treatment? In a study of a medication guaranteed to lower cholesterol levels in the blood, a midwest pharmaceutical company found itself without any patients remaining in the treatment group. The medicine, alas, tasted too much like sand coated with slimy egg whites and chalk dust.

In several less extreme cases, Entsuah [1990] shows that permutation methods can be applied even if withdrawal is related to treatment, providing we modify our scoring system to account for the dropouts. Entsuah studies and compares the power of scoring systems based on functions of boundaries, endpoints, and time using either the ranks or the original observations. His results are specific to the applications he studied.

## 4.6. Sequential Analysis

The underlying idea in sequential analysis, like that of meeting and bedding one's favorite movie actress, is simple and straightforward. It is only in the details of its implementation that we encounter difficulties.<sup>1</sup>

To reduce the required number of observations, we gather data in stages, as first suggested by Wald [1950], computing the statistic  $S_k$  at the  $k$ th stage, accepting the hypothesis and terminating the process if  $S_k \in A_k$ , rejecting the hypothesis in favor of the alternative and terminating the process if  $S_k \in R_k$ , and continuing to gather data otherwise. The acceptance regions  $\{A_k\}$  and the rejection regions  $\{R_k\}$  are chosen so that the significance level and power satisfy certain predetermined conditions.

Can we define such a test? Will it reduce the expected number of observations compared to a sample of fixed size? Can we be sure such a test will terminate after a finite number of steps?

As one simple example, taken from Lehmann [1986, p. 7], suppose we draw observations one at a time from the uniform distribution over the interval  $(\theta - \frac{1}{2}, \theta + \frac{1}{2})$  in order to test a hypothesis about  $\theta$ . If  $|X_1 - X_2|$  is sufficiently close to 1, we could quit after taking only a few observations. If after  $n$  steps, however,  $\max |X_i - X_j|$  is close to 0, we would still need to take more observations.

Under the assumption of normality, David and Kruskal [1956] provide a sequential t-test with the desired properties. To obtain a distribution-free sequential permutation test applied to ranks, we need only assume the observations are independent and identically distributed.

<sup>1</sup> Consider my own plan circa 1961 for meeting Miss Kim Novak: Hitchike from Berkeley to Carmel, then on down the coast to Big Sur where I would hang about the Nepenthe Restaurant for several days hoping to meet her. As to what I would say at this meeting and the problem of my appearance after several nights spent sleeping in the woods, consult Good [1994].

Assume, following Slud and Wei [1982], that subjects are introduced into the study in sequence and are assigned to either Treatment A or Treatment B according to some (possibly restricted) randomization rule. At the  $j$ th step, subjects may be introduced in a group of  $t_j$ , with  $n_j$  assigned to Treatment A. The responses that have been observed by the  $i$ th look are converted to ranks; that is,  $r_{kj}^i$  is the rank of the  $k$ th patient in the  $j$ th block at the  $i$ th look (and is defined only if  $i \geq j$ ), and we compute a linear rank statistic  $w_i$  each time we look at the data, based on the ranks of those patients receiving Treatment A.

To test the null hypothesis that the treatments are equivalent against the alternative that Treatment A is superior to Treatment B, we need to determine stopping boundaries  $\{b_i\}$  so that, if at the  $i$ th step,  $w_k < b_k$  for  $k < i$  and  $w_i \geq b_i$ , we stop the study and reject the hypothesis; if  $i < I$ , we continue to observe the subjects; if  $i = I$ , we accept the null hypothesis and stop the study.

If  $\alpha_i$  is the probability of making a Type I error at the  $i$ th look (that is, of rejecting the null hypothesis in error), then  $\sum_{i=1}^I \alpha_i = \alpha$ . The  $\{\alpha_i\}$  may be predetermined or adaptive, as in Lan and DeMets [1983].

Because the number of permutations  $\prod_{j=1}^I \binom{t_j}{n_j}$  is discrete, the best we can do in the early stages of the test when the total sample size is small is to try to get as close to the  $\{\alpha_i\}$  is possible.

$$\Pr(W_1 \geq b_1) = q_1 \leq \alpha_1,$$

$$q_1 + \Pr(W_2 \geq b_2) = q_2 \leq \alpha_2,$$

and so forth.

An efficient method for calculating the boundary values is given by Mehta et al [1994].

To learn more about sequential tests, see Bürlinger, Martin, and Schrieven [1980] and Siegmund [1985].

## 4.7. Very Large And Very Small Samples

When the sample sizes are large, from several dozen to several hundred observations per group, as they often are in clinical trials, the time required to compute a permutation distribution can be prohibitive even if we are taking advantage of one of the optimal computing algorithms described in Chapter 13. Fortunately, when sample sizes are large—and we refer here to the size of the smallest subsample corresponding to a specific factor combination, not to the size of the sample as a whole, we can make use of an asymptotic approximation in place of the exact permutation distribution. A series of papers by Hoeffding [1951], Box and Anderson [1955], and Kempthorne et al. [1961] support the replacement of the permutation distribution of the  $F$ -statistic by the tabulated distribution of the  $F$ -ratio. This approximation can often be improved on if we replace the observed values by their

corresponding ranks or normal scores. Sections 9.3 and 14.4 provide additional discussion of these points.

With very small samples, the permutation distribution is readily calculated. But with few observations, the power of the test may well be too small and we run the risk of overlooking a treatment effect that is of practical significance. A solution in some cases is to take our observations in stages, rejecting or accepting the null hypothesis at each stage only if the  $p$ -value of the data is very large or very small. Otherwise, we continue to take more observations.

## 4.8. Questions

1. Rewrite the computer program in Section 4.2.3 so it will yield the permutation distributions of the three  $k$ -sample statistics:  $F_1$ ,  $F_2$ , and  $R$ . Would you still accept/reject the hypothesis if you used  $F_2$  or  $R$  in place of  $F_1$ ?
2. *Confidence interval.* Derive a 90% confidence interval for the main effect of sunlight using the crop yield data in Table 4.1a. First, restate the model so as to make clear what it is you are estimating:

$$X_{ikl} = \mu + s_i + f_k + s_{fik} + \varepsilon_{ikl},$$

$$\text{with } s_1 = -\delta \text{ and } s_2 = \delta.$$

Recall that we rejected the null hypothesis that  $\delta = 0$ . Suppose you add  $d = 1$  to each of the observations in the low sunlight group and subtract  $d = 1$  from each of the observations in the high sunlight group. Would you still reject the null hypothesis at the 90% level? If your answer is “yes” then  $d = 1$  does not belong to the 90% confidence interval for  $\delta$ . If your answer is “no” then  $d = 1$  does belong. Experiment (be systematic) until you find a value  $\delta_0$  such that you accept the null hypothesis whenever  $d > \delta_0$ .

3. *Covariate analysis.* Suppose your observations obey the model:

$$Y_{ik} = \mu + s_i + bX_k + \varepsilon_{ik},$$

where the errors  $\varepsilon_{ik}$  are exchangeable. What statistic would you use to test if  $b = 0$ ? to test that  $s_i = 0$  for all  $i$ ?

4. *Equality of the slopes of two lines.* Suppose you observed samples from two populations and that

$$Y_{1k} = \mu_1 + b_1 X_k + \varepsilon_{1k},$$

$$Y_{2k} = \mu_2 + b_2 X_k + \varepsilon_{ik},$$

where the errors  $\varepsilon_{ik}$  are exchangeable. What statistic would you use to test that  $b_1 = b_2$ , that is, that the effect of  $X$  on  $Y$  is the same in the two populations? See Chapter 7.

5. Design an experiment. a) List all the factors that might influence the outcome of your experiment. b) Write a model in terms of these factors. c) Which factors are under your control? d) Which of these factors will you use to restrict the scope of the experiment? e) Which of these factors will you use to block? f) Which of the remaining factors will you neglect initially, that is, lump into your error term? g) How will you deal with each of

- the remaining covariates? h) By correction? i) By blocking after the fact? j) How many subjects will you observe in each subcategory? k) Is the subject the correct experimental unit? l) Write out two of the possible assignments of subjects to treatment. m) How many possible assignments are there in all?
6. a. Is the Still–White test for interaction asymptotically exact?  
b. What if we were to generate random relabelings of the original observations and then compute the I statistic. Would the resulting distribution provide an exact test?
  7. It may not always be possible to find a test of fixed size with guaranteed power and significance level. Suppose  $\rho$  is the odds ratio of two binomials,  $\rho = \frac{p_2(1-p_1)}{(1-p_2)p_1}$ . Show that, if  $\alpha < \beta$ , no test with a fixed number of trials exists for testing  $\rho = \rho_0$  against all alternatives with  $\rho = \rho_1$ . Show that, if the observations are taken in pairs, one from each binomial, it is possible to obtain a terminating sequential test with the desired properties.
  8. How many distinct synchronized permutations are there of a  $2^2$  factorial design with K observations per cell?

## CHAPTER 5

# Multivariate Analysis

### 5.1. Introduction

The value of an analysis based on simultaneous observations on several variables—height, weight, blood pressure, and cholesterol level, for example, is that it can be used to detect subtle changes that might not be detectable, except with very large, prohibitively expensive samples, were you to consider only one variable at a time.

Any of the permutation procedures described in Chapters 3 and 4 can be applied in a multivariate setting providing we can find a single-valued test statistic which can stand in place of the multivalued vector of observations.

### 5.2. One- and Two-Sample Comparisons

#### 5.2.1. Hotelling's $T^2$

One example of such a statistic is Hotelling's  $T^2$ , a straightforward generalization of Student's  $t$  to the multivariate case.

Suppose we have made a series of exchangeable vector-valued observations  $\vec{X}_i = \{X_{i1}, X_{i2}, \dots, X_{iJ}\}$ , for  $i = 1, \dots, I$ . Let  $\vec{X}_.$  denote the vector of mean values  $\{X_{.1}, X_{.2}, \dots, X_{.J}\}$ , and  $V$  the  $J \times J$  covariance matrix; that is,  $V_{ij}$  is the covariance of  $X_{ki}$  and  $X_{kj}$ . To test the hypothesis that the midvalue of  $\vec{X}_i = \vec{\xi}$  for all  $i$ , use

$$\text{Hotelling's } T^2 = (\vec{X}_i - \vec{\xi})V^{-1}(\vec{X}_i - \vec{\xi})^T.$$

Loosely speaking, this statistic weighs the contribution of individual variables and pairs of variables in inverse proportion to their covariances. If the variables in each observation vector are independent of one another (a rare case, indeed), Hotelling's  $T^2$  weighs the contributions of the individual variables in inverse proportion to their variances.

The two-sample comparison is only slightly more complicated: Let  $n_1, \vec{X}_{1\cdot}$ ;  $n_2, \vec{X}_{2\cdot}$  denote the sample size and vector of mean values of the first and second samples respectively. We assume under the null hypothesis that the two sets of vector-valued observations  $\{\vec{X}_{1i}\}$  and  $\{\vec{X}_{2i}\}$  come from the same distribution (that is, the sample labels 1 and 2 are exchangeable). Let  $V$  denote the pooled estimate of the common covariance matrix; as in the one-sample case,  $V_{ij}$  denotes the pooled covariance estimate of  $X_{1ki}$  and  $X_{2kj}$ ,

$$(N - 2)V_{ij} = \sum_{m=1}^2 \sum_{k=1}^{n_m} (X_{mki} - X_{m\cdot i})(X_{mkj} - X_{m\cdot j}).$$

To test the hypothesis that the midvalues of the two distributions are the same, we could use the statistic

$$T^2 = (\vec{X}_{1\cdot} - \vec{X}_{2\cdot})V^{-1}(\vec{X}_{1\cdot} - \vec{X}_{2\cdot})^T,$$

but then we would be forced to recompute the covariance matrix  $V$  and its inverse  $V^{-1}$  for each new rearrangement. To reduce the number of computations, Wald and Wolfowitz [1944] suggest a slightly different statistic  $T'$  that is a monotonic function of  $T$  (see Problem 3).

Let

$$U_j = N^{-1} \sum_{i=1}^2 \sum_{k=1}^{n_i} X_{ikj}$$

$$c_{ij} = \sum_{m=1}^2 \sum_{k=1}^{n_m} (X_{mki} - U_i)(X_{mkj} - U_j).$$

Let  $C$  be the matrix whose components are the  $c_{ij}$ . Then

$$T'^2 = (\vec{X}_{1\cdot} - \vec{X}_{2\cdot})C^{-1}(\vec{X}_{1\cdot} - \vec{X}_{2\cdot})^T.$$

As with all permutation tests we proceed in three steps:

- (1) we compute the test statistic for the original observations;
- (2) we compute the test statistic for all relabelings;
- (3) we determine the percentage of relabelings that lead to values of the test statistic that are as extreme or more extreme than the orginal value.

For the purpose of relabeling, each vector of observations on an individual subject is treated as a single indivisible entity. When we relabel, we relabel on a subject-by-subject basis so that all observations on a single subject receive the same new label. If the original vector of observations on subject  $i$  consists of  $k$  distinct observations on  $k$  different variables

$$(x_i^1, x_i^2, \dots, x_i^k),$$

and we give this vector a new label  $p(i)$ , then the individual observations remain together as a unit, each with the new label:

### Sidebar

Calculating the Wald–Wolfowitz variant of Hotelling’s  $T^2$  Blood Chemistry Data from Werner et al [1970]

ID	BC	Albumin	Uric Acid	Mean	Albumin	Uric Acid
2381	N	43	54			
1610	N	41	33	N	41.25	46.25
1149	N	39	50	Y	37.0	52.75
2271	N	42	48	Comb	39.125	49.5
				Y-N	-4.25	6.50
1946	Y	35	72	C		
1797	Y	38	30		8.982	-21.071
575	Y	40	46		-21.071	196.571
39	Y	35	63			
$C^{-1}$						
Hotelling’s $T^2$						
$= (-4.25 \quad 6.50)C^{-1}(-4.25 \quad 6.50)^T$						
$= 2.092$						
.1487						
.01594						
.01594						
.006796						

$$(x_{p(i)}^1, x_{p(i)}^2, \dots, x_{p(i)}^k).$$

This approach to relabeling should be contrasted with the approach we would use if we were testing for independence of the covariates (see Section 7.2).

Hotelling’s  $T^2$  is the appropriate statistic to use if you suspect the data have a distribution that is close to that of the multivariate normal. Under the assumption of multivariate normality, the power of the permutation version of Hotelling’s  $T^2$  converges with increasing sample size to the power of the most powerful parametric test that is invariant under transformations of scale.

The stated significance level of the parametric version of Hotelling’s  $T^2$  cannot be relied on for small samples if the data are not normally distributed [Davis, 1982]. As always, the corresponding permutation test yields an exact significance level even if the errors are not normally distributed, providing that the errors are exchangeable from sample to sample.

Much of the theoretical work on Hotelling’s  $T^2$  has focused on the properties of the *unconditional*<sup>1</sup> permutation test in which the original observations are replaced by ranks. Details of the asymptotic properties and power of the unconditional test are given in Barton and David [1961], Chatterjee and Sen [1964, 1966], and, most recently, Gill and Siotani [1987]. The effect of missing

<sup>1</sup> Recall from our discussion in Section 2.3 that whereas we must compute the permutation distribution anew for each new set of observations, the permutation distribution of a set of ranks is independent or unconditional of the actual values of the observations.

observations on the significance level and power of the test is studied by Servy and Sen [1987].

### 5.2.2. An Example

The following blood chemistry data are taken from Werner et al [1970]. The full data set is included with the BMDP statistical package. An asterisk (\*) denotes missing data.

1	2	3	4	5	6	7	8	9
2381	22	67	144	N	200	43	98	54
1946	22	64	160	Y	600	35	*	72
1610	25	62	128	N	243	41	104	33
1797	25	68	150	Y	50	38	96	30
1149	53	*	178	N	227	39	*	50
575	53	65	140	Y	220	40	107	46
2271	54	66	158	N	305	42	103	48
39	54	60	170	Y	220	35	88	63

The variables are

1. identification number
2. age in years
3. height in inches
4. weight in pounds
5. uses birth control pills?
6. cholesterol level
7. albumin level
8. calcium level
9. uric acid level.

A potential hypothesis of interest is whether birth-control pill usage has any effect on blood chemistries. As the nature of such hypothetical effects very likely depends on age and years of use, before testing this hypothesis using a permutation method, you might want to divide the data into two blocks corresponding to young and old patients.

You could test several univariate hypotheses using the methods of Section 3.5; for example—the hypothesis that using birth control pills lowers the albumin level in blood. You might want to do this now to see if you can obtain significant results. As the sample sizes are small, the univariate observations may not be statistically significant. But by combining the observations that Werner and his colleagues made on several different variables to form a single multivariate statistic, you may obtain a statistically significant result; that is, if taking birth control pills does alter blood chemistries.

**Sidebar**

Program for computing multivariate permutation statistics

```
#define length    119
#define control   60
#define variates   9
```

Set aside space for a multivariate array data [length, variates]; and a vector of sample sizes index[length];

Main program

```
load (data);
compute stat0 (data, index);
repeat Nsim times;
    rearrange data;
    compute stat (data, index);
    record whether stat >= stat0;
print out the significance level of the test;
```

Load

packs the data into a long matrix, each row of which corresponds to  $k$  observations on a single subject; the first  $n$  rows are the control group; the last  $m$  rows are the treatment group. (A second use of this subroutine will be to eliminate variables and subjects that will not be included in the analysis, e.g., to eliminate all records that include missing data, and to define and select specific subgroups.)

Rearrange

randomly rearranges the rows of the data array; the elements in each row are left in the same order.

Compute

```
calculate the mean of each variable for each sample and store the results in a 2 by  $n$  array  $N$ ;
calculate  $n$  by  $n$  array  $V$  of covariances for the combined sample and invert  $V$ ;
matrix mult (mean,  $W, * W$ );
matrix mult ( $W$ , mean);
return  $T^2$ ;
```

### 5.2.3. Doing the Computations

You don't need to use all the dependent variables in forming Hotelling's  $T^2$ . For example, you could just include albumin and uric acid levels as we have in a sidebar. For each relabeling, you would need to compute four sample means corresponding to the two variables and the two treatment groups. And you would need to perform two successive matrix multiplications. I would not attempt these calculations without a computer and the appropriate computer software: Werner et al's full data set includes 188 cases!

As in the univariate examples in Chapters 3 and 4, you need to program and implement three procedures:

- a) one to rearrange the stored data;
- b) one to compute the  $T^2$  statistic;
- c) and one to compute the significance level.

Only the first of these procedures, devoted to rearranging the data, represents a significant change from the simple calculations we performed in the univariate case. In a multivariate analysis, we can't afford to manipulate the actual data; a simple swap could mean the exchange of nine or ten or even 100 different variables; so we rearrange a vector of indices that point to the data instead. Here is a fragment of code in the *C* programming language that does just that:

```
float Data [length, variates];
int index[length];
.....
rearrange (index, length);
.....
for (j = 0; j < ncontrol;j++) Mean [k]+ = Data [index[j], k].
```

#### 5.2.4. Weighting the Variables

With several variables simultaneously lending their support to (or withholding their support from) a hypothesis or an alternative, should some variables be given more weight than others? Or should all variables be given the same importance?

In any multivariate application, whether or not you use Hotelling's  $T^2$  as the test statistic, you may want to "Studentize" the variables by dividing by the covariance matrix before you begin.

Hotelling's  $T^2$  Studentizes the variables in the sense that it weights each variable in inverse proportion to its standard deviation. (This is not quite true if the variables are correlated; see below.) As a result, Hotelling's  $T^2$  is dimensionless; it will not matter if we express a vector of measurements in feet rather than inches or miles. Variables whose values fluctuate widely from observation to observation are given less weight than variables whose values are essentially constant.

When we convert an observation on a variable to the corresponding rank or normal score (see Section 9.3), we are also standardizing it. If we have exactly the same number of observations on each variable—as would be the case, for example, if all the observations on all the variables have been accurately recorded and none are missing, then all ranked variables will have exactly the same variance. The problem of what emphasis to give each individual variable is solved automatically.

Although we have standardized the variances in forming a rank test, we must still divide by the covariances.<sup>1</sup> When we divide by the covariance or, equivalently in the

<sup>1</sup> See Pesarin [1999, p. 140] for an alternative weighting procedure.

case of ranks, by the correlation, we are discounting the importance and reducing the effects of *correlated* or dependent variables. If we have two perfectly correlated variables—the testimony of a ventriloquist and his dummy—then, clearly the second variable (or witness) has no information to contribute beyond that which we have already received from the first.

### 5.2.5. Interpreting the Results

The significance of  $T^2$  or some equivalent multivariate statistic still leaves unanswered the question of *which* variables have led to the rejection of the multivariate hypothesis. For a discussion of the problem of simultaneous inference, see any text on multivariate methods, for example, Morrison [1990]. My own preference on finding a significant result, a preference that reflects my studies under Jerzy Neyman, is to search for a mechanistic, cause-and-effect model that will explain the findings. In Chapters 7 through 10, we consider some of the tests one might perform to verify or disprove such a model.

### 5.2.6. Alternative Statistics

Hotelling's  $T^2$  is designed to test the null hypothesis of no difference between the distributions of the treated and untreated groups against alternatives that involve a shift of the  $k$ -dimensional center of the multivariate distribution. Although Hotelling's  $T^2$  offers protection against a wide variety of alternatives, it is not particularly sensitive to alternatives that entail a shift in just one of the dependent variables.

Boyett and Shuster [1977] show that a more powerful test against such alternatives is based on the permutation distribution of the test statistic

$$\max_{1 \leq j \leq k} \frac{(X_{1\cdot}^j - X_{2\cdot}^j)}{SE^k},$$

a statistic first proposed in a permutation context by Chung and Fraser [1958], where  $SE^k$  is a pooled estimate of the standard error of the mean of the  $k$ th variable.

Let us apply this approach to the subsample of blood chemistry data we studied in Section 5.2.1. We use a two-sided test so we may detect changes up or down. For albumin, the absolute difference in means is 4.25 and the standard error is  $\sqrt{8.982}/2 = 1.498$ ; for uric acid, the difference is 6.5, the standard error is 7.010. Our test statistic is 2.84, the larger of the two weighted differences. To determine whether this value is significant, we need to compute the sample means, their differences, and the maximum difference after weighting by our estimates of the standard errors for each of the  $\binom{8}{4} = 70$  rearrangements of the two samples.

Having to rely on sample-based estimates reduces the power of this test [Pesarin, 1998, p. 193–195]. For alternatives, see Blair et al [1994] and Pesarin [1998].

### 5.3. Runs Test

Friedman and Rafsky (1979) provide a multivariate generalization of the Wald–Wolfowitz and Smirnov distribution-free two-sample tests used for testing  $F_X = F_Y$  against the highly nonspecific alternative  $F_X \neq F_Y$ . In both the univariate and the multivariate versions of these two-sample tests, one measures the degree to which the two samples are segregated within the combined sample. In the univariate version, one forms a single combined sample, sorts and orders it, and then

- a) counts the number of runs in the combined sample, or
- b) computes the maximum difference in cumulative frequency of the two types within the combined sample.

For example, if  $x = (1, 3, 6)$  and  $y = (2, 4, 5)$ , the ordered combined sample is 1, 2, 3, 4, 5, 6, that is, an  $x$  followed by  $y$   $x$   $yy$   $x$  and has five runs.

Highly segregated samples will give rise to a small number of runs (and a large maximum difference in cumulative frequency), while highly interlaced distributions will give rise to a large number of runs (and a very small difference in cumulative frequency). Statistical significance, that is, whether the number of runs is significantly large, can be determined from the permutation distribution of the test statistic.

To create a multivariate version of these tests, we must find a way to order observations that have multiple coordinates. One of the keys to this ordering is the *minimal spanning tree* described by Friedman and Rafsky [1979].

Each point in Figure 5.1a corresponds to a pair of observations, height and weight, say, that were made on a single subject. We build a spanning tree between these data points, as in Figure 5.1b, by connecting the points so that there is exactly one path between each pair of points, and so that no path closes back on itself in a loop. Obviously, we could construct a large number of such trees. A minimal spanning tree is one for which the sum of the lengths of all the paths is a minimum. This tree is unique if there are no ties among the  $N(N - 1)/2$  interpoint distances.

Before computing the test statistic(s) in the multivariate case, we first construct the minimal spanning tree for the combined sample. Once the tree is complete, we can generate the permutation distribution of the runs statistic through a series of random relabelings of the individual data points. After each relabeling, we remove all edges for which the defining nodes originate from different samples. Figure 5.1c illustrates one such result.

Although it can take a multiple of  $N \times N$  calculations to construct the minimal spanning tree for a sample of size  $N$ , each determination of the multivariate runs statistic only takes a multiple of  $N$  calculations. For large samples a normal

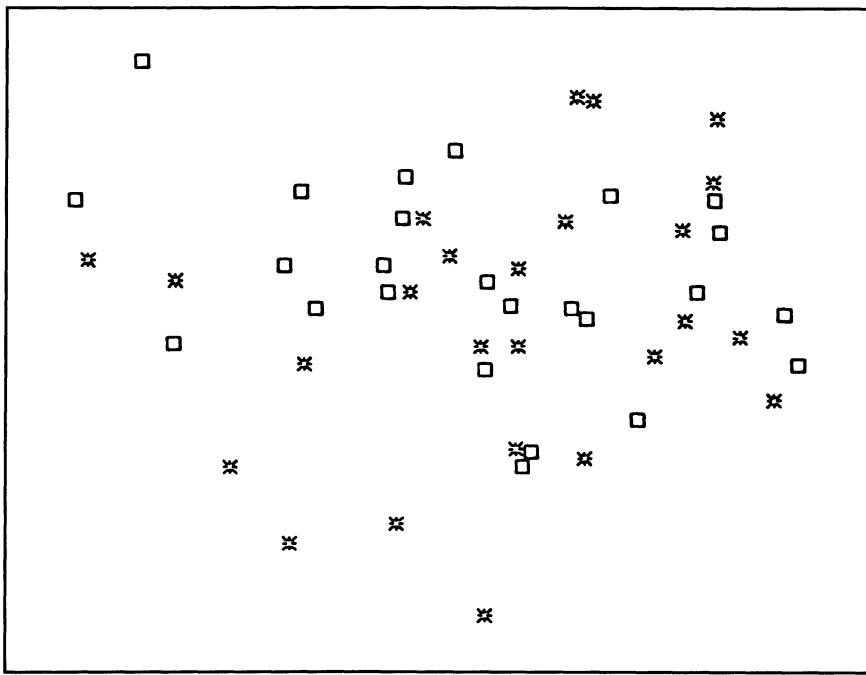
**A**

Figure 5.1. Building a minimal spanning tree.

From "Multivariate generalizations of the Wald–Wolfowitz and Smirnov two-sample test," by J. H. Friedman and L. C. Rafsky, *Annals of Statistics*; 1979; 7: 697–717. Reprinted with permission from the Institute of Mathematical Statistics.

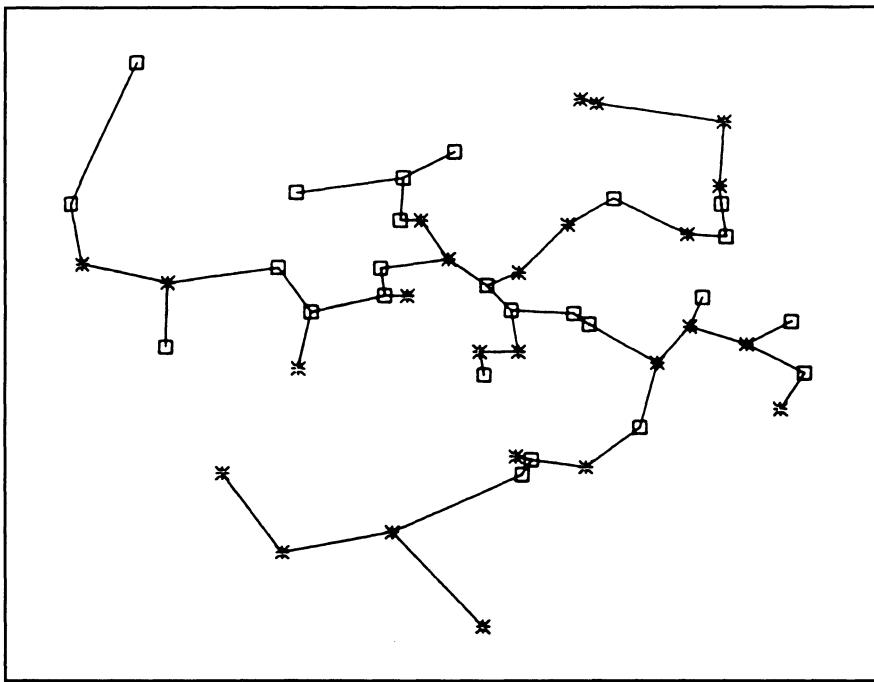
Continued next page.

approximation to the permutation distribution may be used (see Section 14.4); the expected value and variance of the runs statistic are the same as in the univariate case.

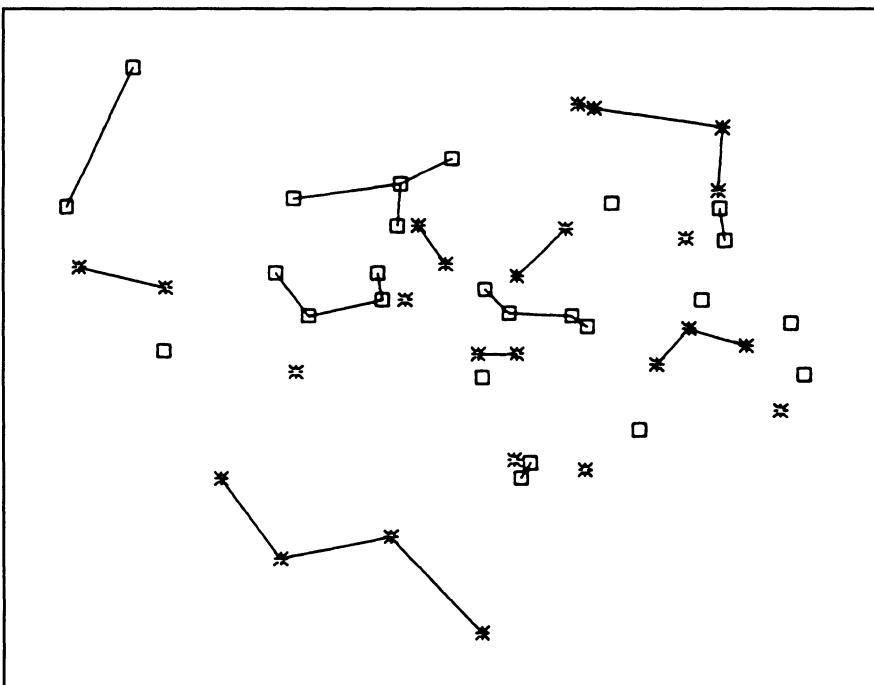
### 5.3.1. Which Statistic?

We've now considered three multivariate test statistics for testing hypotheses based on one or two samples. Which one should we use? To detect a simultaneous shift in the means of several variables, use Hotelling's  $T^2$ ; to detect a shift in *any* of several variables, use the maximum  $t$ ; and to detect an arbitrary change in a distribution (not necessarily a shift) use Friedman and Rafsky's multivariate runs test.

Tests proposed by van-Putten [1987] and Henze [1988] offer advantages over Friedman–Rafsky in some cases.



B



C

Figure 5.1. Continued from previous page.

## 5.4. Experimental Designs

### 5.4.1. Matched Pairs

Puri and Shane [1970] study the multivariate generalization of paired comparisons in an incomplete blocks design (see Sections 3.6 and 3.7). Their procedure is a straightforward generalization of the multivariate one-sample test developed by Sen and Puri [1967]; (see also Sen [1967, 1969]).

For simplicity, suppose we have only a single block. As in Section 3.1, we consider all possible permutations of the signs of the individual multivariate observations. If  $\{\vec{X}_i, \vec{Y}_i\}$  is the  $p$ -dimensional vector of observations on the  $i$ th matched pair, and  $\vec{Z}_i$  is the vector of differences  $(Z^1, \dots, Z^p)$ , then our permutation set consists of vectors of differences of the form  $((-1)^{j_1} Z_1, \dots, (-1)^{j_n} Z_n)$  where  $-Z = (-Z^1, \dots, -Z^p)$ .

Depending on the hypothesis and alternatives of interest, one may want to apply an initial set of linear transformations to each separate coordinate, that is, to replace  $Z^j$  by  $Z^{j'} = a_j + b_j Z^j$ . Puri and Shane studied the case in which the individual variables were replaced by their ranks, with each variable being ranked separately.

### 5.4.2. Block Effects

When we have more than two treatments to compare, an alternative statistic studied by Gerig [1969] is the multivariate extension of Friedman's chi-square test in which ranks take the place of the original observations creating an unconditional permutation test.

The experimental units are divided into  $B$  blocks each of size  $I$  with the elements of each block as closely matched as possible with respect to extraneous variables. During the design phase, one individual from each block is assigned to each of the  $I$  treatments. We assume that  $K$  (possibly) dependent observations are made simultaneously on each subject. To test the hypothesis of identical treatment effects against translation-type alternatives, we first rank each individual variable separately within each block, ranking them from 1 to  $I$  (smallest to largest). The rank totals  $T_{i(k)}$  are computed for each treatment  $i$  and each variable  $(k)$ . The use of ranks automatically rescales each variable so that the variances (but not the covariances) are the same.

Let  $T$  denote the  $I \times K$  matrix whose  $ik$ th component is  $T_{i(k)}$ . Noting that the expected value of  $T_{i(k)}$  is  $(K + 1)/2$ , let  $V$  denote the matrix whose components are the sample covariances

$$V_{st} = \frac{\left[ \sum_{b=1}^B \sum_{i=1}^I T_{bi(s)} T_{bi(t)} - \frac{k(k+1)^2}{4} \right]}{n(k-1)}.$$

By analogy with Hotelling's  $T^2$ , the test statistic is  $TV^{-1}T^T$  [Gerig, 1969]. Gerig [1975] extends these results to include and correct for random covariates.

## 5.5. Multiple Comparisons

One of the difficulties with clinical trials and other large-scale studies is that frequently so many variables are underinvestigation that one or more of them is practically guaranteed to be significant by chance alone. If we perform 20 tests at the 5% level, we expect at least one significant result in 20 on the average.

### 5.5.1. Step Up or Step Down?

One way, and not a very good one, to ensure that the probability of making at least one Type I error is less than some predesignated value  $\alpha$  is to make the  $k$  different comparisons each at level  $\alpha/k$ . A better method, first described by Holm [1979], first orders the p-values from smallest to largest (or the corresponding standardized test statistics from largest to smallest). One begins with the most significant result and decides whether to accept or reject. Obviously, once a hypothesis is accepted, then all hypotheses with larger p-values are accepted as well. If a hypothesis is rejected, then a new critical value is determined and the next p-value inspected. Permutation procedures utilizing this step-down<sup>1</sup> approach were developed independently by Westfall and Young [1993], Blair and Karniski [1994], and Troendle [1995]; a test based on the latter's work is described in the next subsection.

The chief weakness of the step-down procedure is its dependence on the rejection criteria used to test the smallest p-value, normally  $p_{(1)} \leq \alpha/k$ . An alternative developed by Hochberg [1988] begins with the largest p-value at the first step. If a hypothesis is rejected, then all hypotheses with smaller p-values are rejected as well.<sup>2</sup> If a hypothesis is accepted, then a new critical value is determined and the next p-value inspected. Blair, Troendle, and Beck [1996] report that this step-up method, an example of which is provided in Section 5.5.1.2., is slightly more powerful than the step-down.

#### 5.5.1.1. Standardized Statistics

A resampling procedure outlined by Troendle [1995] allows us to work around the dependencies among variables. Suppose we have measured  $k$  variables on each subject and are now confronted with  $k$ -test statistics  $s_1, s_2, \dots, s_k$ . To make these

<sup>1</sup> One steps down among the values of the test statistic.

<sup>2</sup> Again, one steps up among the values of the test statistic.

statistics comparable, we need to standardize them and render them dimensionless, dividing each by its respective  $L_1$  or  $L_2$  norm. For example, if one variable, measured in centimeters, takes values like 144, 150, 156 and the other, measured in meters, takes values like 1.44, 1.50, 1.56, we might set  $t_1 = s_1/4$  and  $t_2 = s_2/0.04$ .

Next, we order the standardized statistics by magnitude so that  $t_{(1)} \leq \dots \leq t_{(k)}$ . Denote the corresponding hypotheses as  $H_{(1)}, \dots, H_{(k)}$ . The probability that at least one of these statistics will be significant by chance alone at the  $\alpha$  level if they are independent is  $1 - (1 - \alpha)^k$  which is approximately  $k\alpha$ . But once we have rejected one hypothesis (assuming it was false), there will only be  $k - 1$  true hypotheses to guard against rejecting.

Begin with  $i = 1$  and

1. Repeatedly resample the data (with or without replacement), estimating the cutoff value  $\phi(\alpha, k - i + 1)$  such that  $\alpha = \Pr\{T(k - i + 1) \leq \phi(\alpha, k - i + 1)\}$ , where  $T(k - i + 1)$  is the largest of the  $k - i + 1$  test statistics  $t_{(1)} \dots t_{(k-i+1)}$  for a given resample.
2. If  $t_{(k-i+1)} \leq \phi(\alpha, k - i + 1)$ , then accept all remaining hypotheses  $H_{(1)} \dots H_{(k-i+1)}$  and STOP.

Otherwise, reject  $H_{(k-i+1)}$ , increment  $i$ , and RETURN to step 1.

### 5.5.1.2. Paired Sample Tests

Suppose we observe before and after values of  $k$  attributes for each of  $N$  subjects and wish to test the  $k$  null hypotheses that the before and after means are the same. We perform a set of  $k$  tests, which may be parametric, permutation, or bootstraps and order the significance levels  $p_{(1)} \geq p_{(2)} \dots \geq p_{(k)}$ . At the  $j$ th step of our multiple comparison procedure, following Blair, Troendle, and Beck [1996], we compute the critical level

$$\gamma_j = \frac{1}{2^N} \sum_{i=1}^{2^N} I \left[ \min_{j \leq i \leq k} p_k(\pi_{ki}) \leq p_j \right],$$

where  $I$  is an indicator function taking the value 1 if its argument is true, and the value 0 otherwise, and  $\pi_{ki}$  denotes the  $i$ th permutation of the before, after data for the  $k$ th variable. Thus, the sum counts the number of permutations for which the inequality is true.

In the step-up procedure, if  $\gamma_j < \alpha$ , we accept the remaining hypotheses and stop. Otherwise, we accept the hypothesis  $H_{(j)}$ , increment  $j$ , and continue.

For example, suppose we have collected the following observations:

subject 1: before =  $\begin{pmatrix} 5 \\ 6 \end{pmatrix}$ , after =  $\begin{pmatrix} 3 \\ 7 \end{pmatrix}$  showing a decline in both variables,  
 subject 2: before =  $\begin{pmatrix} 4 \\ 4 \end{pmatrix}$ , after =  $\begin{pmatrix} 3 \\ 4 \end{pmatrix}$  showing a decline in just one.

A t-test of variable 1 yields a p-value of 10%; a permutation test of variable 2 yields a p-value of 50%.

There are three possible rearrangements of the data in addition to the original, for example, before  $\binom{5}{6}$ , after  $= \binom{3}{7}$ , and before  $= \binom{3}{4}$ , after  $= \binom{4}{4}$ . All rearrangements yield the same p-value for variable 2, but two of the four rearrangements yield p-values of 85% for variable 1.

The step-up test starts with the largest of the original p-values, 50%, for variable 2.  $\gamma_1 = 1$ , and we accept. Thereafter, we ignore variable 2 and focus on the remaining variable.  $\gamma_2 = 0.5$ , and for an experiment-wide error of  $\alpha \leq 0.5$ , we also accept the null hypothesis for variable 1.

### 5.5.2. Fisher's Omnibus Statistic

The methods described in this section due to Pesarin [1990] apply to exchangeable continuous, ordinal, or categorical variables or to any combination thereof. As in Section 5.2.1, suppose we have made a series of exchangeable vector-valued observations on J variables. In the kth vector of observations in the i<sup>th</sup> treatment group,  $\vec{X}_{ik} = \{X_{ik1}, X_{ik2}, \dots, X_{ikJ}\}$ ,  $X_{ik1}$  might be a 0 or 1 according to whether the kth seedling in the i<sup>th</sup> treatment group germinated,  $X_{ik2}$  the height of the i<sup>th</sup> seedling,  $X_{ik3}$  its weight, and so forth. Let  $\vec{T}_o$  denote the vector of single-variable test statistics  $\{T_{o1}, T_{o2}, \dots, T_{oJ}\}$  based on the original unpermuted matrix of observations  $\mathbf{X}_o$ . These might include differences of means or weighted sums of the total number germinated, or any other statistic one might employ when testing just one variable at a time. When we permute the observations vectors among treatments, as in  $\pi(\mathbf{X})$ , we obtain a new vector of single-variable test statistics  $\vec{T}^* = \{T_1^*, T_2^*, \dots, T_J^*\}$ .

To obtain a composite multivariate test, we first must convert these separate and quite distinct statistics to a common basis, in this case, their ranks within their own individual permutation distributions. We proceed as follows:

1. Generate S permutations of  $\mathbf{X}$  and thus obtain S vectors of test statistics that we denote by  $\vec{T}_i^*$ ,  $i = 1, \dots, S$ .

2. Rank the single-variable test statistics;

$R_{ij} = R(T_{ij}^*) = \sum_{h=1}^S I(T_{hj}^* \leq T_{ij}^*)$ , for  $i = 1, \dots, S$ ;  $j = 1, \dots, J$ ; where the indicator function  $I(E)$  takes values 1 or 0 depending on whether E is true or false. The result is J sets of ranks each ranging from 1 to S.

3. Combine the ranks of the J individual tests using Fisher's omnibus method;

$$U_i = - \sum_{j=1}^J \log \left( \frac{S + 0.5 - R_{ij}}{S + 1} \right); \quad i = 1, \dots, S.$$

4. Determine the marginal significance levels of each of the single-variable tests,

$$p_j = \frac{0.5 + \sum I(T_{hj}^* \geq T_{oj})}{S + 1}, \quad j = 1, \dots, J.$$

5. Combine these values  $U_o = - \sum_{j=1}^J \log(p_i) = - \sum_{j=1}^J \log\left(\frac{S+0.5-R_{oj}}{S+1}\right)$ ; note that  $R_{oj}$  can take any value in the range 0 to S.

6. Determine the significance level of the combined test,

$$p = \frac{0.5 + \sum_{i=1}^S I(U_i \geq U_o)}{S+1}.$$

## 5.6. Repeated Measures

In many experiments, we want to study the development of a process over a period of time, such as the growth of a tumor or the gradual progress of a cure. If our observations are made by sacrificing different groups of animals at different periods of time, then time is simply another variable in the analysis which we may treat as a covariate. But if all our observations are made on the same subjects, then the multiple observations on a single individual will be interdependent. And all the observations on a single subject must be treated as a single multivariate vector.

We may ask at least three questions about the response profiles: (1) Are the response profiles the same for the various treatments? (2) Are the response profiles parallel? (3) Are the response profiles at the same level?

A “yes” answer to question (1) implies “yes” answers to questions (2) and (3), but we may get a “yes” answer to (2) even when the answer to (3) is “no”.

One simple test of parallelism entails computing the successive differences  $z_{j,i} = x_{j,i+1} - x_{j,i}$  for  $j = 1, 2; i = 1, \dots, I-1$  and then applying the methods from Sections 5.2 or 5.3 to these differences. Of course, this approach is applicable only if the observations on both treatments were made at identical times.

To circumvent this limitation and to obtain a test of the narrower hypothesis (1), we follow Koziol et al [1981] and suppose there are  $N_i$  subjects in group  $i$ . Let  $X_{tj}^i, t = 1, 2, \dots, T$ ; and  $j = 1, 2, \dots, N_i$  denote the observation on the  $j$ th subject in group  $i$  at time  $t$ . Not all the  $X_{tj}^i$  may be observed in practice; we will only have observations for  $N_{it}$  of the  $N_i$  in the  $i$ th group at time  $t$ . If  $X_{tj}^i$  is observed, let  $R_{tj}^i$  be its rank among the  $N_i$  available values at time  $t$ . Set  $S_{it} = (N_{it})^{-1} \sum R_{tj}^i$ .

If luck is with us so that all subjects remain with us to the end of the experiment, then  $N_{it} = N_i$  for all  $t$  and each  $i$ , and we may adopt as our test statistic  $L_N = \sum N_i \vec{S}_i^T V^{-1} \vec{S}_i$ , where  $\vec{S}_i$  is a  $T \times 1$  vector with components  $(S_{i1}, S_{i2}, \dots, S_{iT})$  and  $V$  is a  $T \times T$  covariance matrix whose  $st$ th component is

$$v_{st} = N^{-1} \sum_{i=1}^I \sum_{j=1}^{N_i} R_{sj}^i R_{tj}^i.$$

This test statistic was proposed and investigated by Puri and Sen [1966, 1969, 1971].

### 5.6.1. Missing Data

If we are missing data, and missing data are almost inevitable in any large clinical study since individuals commonly postpone or even skip follow-up appointments, then no such simplified statistic presents itself. Zerbe and Walker [1977] suggest that each subject's measurements first be reduced to a vector of polynomial regression coefficients with time the independent variable. The subjects needn't have been measured at identical times or over identical periods, nor does each subject need to have the same number of observations. Only the number of coefficients (the rank of the polynomial) needs to be the same for each subject. Thus, we may apply the equations of Koziol et al to these vectors of coefficients though we cannot apply the equations to the original data.

We replace the  $m_k$  observations on the  $k$ th subject,  $\{X_{ki}, i = 1, \dots, m_k\}$  with a set of  $J + 1$  coefficients,  $\{b_{kj}, j = 0, \dots, J\}$ . While the  $m_k$  may vary, the number  $J$  is the same for every subject; of course,  $J < m_k$  for all  $k$ . The  $\{b_{kj}\}$  are chosen so that for all  $k$  and  $i$ ,

$$X_{ki} = b_{k0} + t_{ki}b_{k1} + \dots + t_{ki}^J b_{kJ},$$

where the  $\{t_{ki}, i = 0, \dots, m_k\}$  are the observation times for the  $k$ th subject.

This approach has been adopted by a number of practitioners including Albert et al [1982], Chapelle et al [1982], Goldberg et al [1980], and Hiatt et al [1983]. Multiple-comparison procedures based on it include Foutz et al [1985] and Zerbe and Murphy [1986]. A SAS/IML program to do the calculations is available [Nelson and Zerbe, C1988].

### 5.6.2. Bioequivalence

Zerbe and Walker's solution to the problem of missing data suggests a multivariate approach we may use with any time course data. For example, when we do a bioequivalence study, we replace a set of discrete values with a "smooth" curve. This curve is derived in one of two ways: 1) by numerical analysis and 2) by modeling. The first yields a set of coefficients, the second a set of parameter estimates. Either the coefficients or the estimates may be treated as if they were the components of a multivariate vector and the methods of this chapter applied to them.

Here is an elementary example: Suppose you observe the time course of a drug in the urine over a period for which a linear model would be appropriate. Suppose further that the chief virtue of your measuring system is its low cost; the individual measurements are crude and imprecise. To gain precision, you take a

series of measurements on each patient about half an hour apart and use least-squares methods to derive a best-fitting line for each patient. That is, you replace the set of measurements  $\{X_{ijk}\}$  where  $i = 0$  or  $1$  denotes the drug,  $j = 1, \dots, J$  denotes the subject, and  $k = 1, \dots, K_j$  denotes the observation on a subject, with the set of vectors  $\{\bar{Y}_{ij} = (a_{ij}, b_{ij})\}$  where  $a_{ij}$  and  $b_{ij}$  are the intercept and slope of the regression line for the  $j$ th subject in the  $i$ th treatment group.

Using the computer code in Section 5.2, you calculate the mean vector and the covariance matrix for the  $\{\bar{Y}_{ij}\}$ , and compute Hotelling's  $T^2$  for the original observations and for a set of random arrangements. You use the resultant permutation distribution to determine whether the time courses of the two drugs are similar.

### 5.6.3. Omnibus Test

Yet a fourth way to approach the problem of repeated measures is, after correcting for the baseline so that  $X'_{ij}[t] = X_{ij}[t] - X_{ij}[0]$ , to analyze the differences as we would any set of dependent variables using Fisher's omnibus method, as described in Section 5.5.2. See Pesarin [1997].

## 5.7. Questions

1. You can increase the power of a statistical test in three ways: a) making additional observations, b) making more precise observations, and c) adding covariates. Discuss this remark in light of your own experimental efforts.
2. You are studying a new tranquilizer that you hope will minimize the effects of stress. The peak effects of stress manifest themselves between five and ten minutes after the stressful incident, depending on the individual. To be on the safe side, you've made observations at both the five- and ten-minute marks.

Subject	pre-stress	5-minute	10-minute	Treatment
A	9.3	11.7	10.5	Brand A
B	8.4	10.0	10.5	Brand A
C	7.8	10.4	9.0	Brand A
D	7.5	9.2	9.0	New drug
E	8.9	9.5	10.2	New drug
F	8.3	9.5	9.5	New drug

How would you correct for the pre-stress readings? Is this a univariate or a multivariate problem? List possible univariate and multivariate test statistics. Perform the permutation tests and compare the results.

3. Show that if  $T'$  is a monotonic function of  $T$ , then a test based on the permutation distribution of  $T'$  will accept or reject only if a permutation test based on  $T$  also accepts or rejects.

4. Suppose you wish to test the hypothesis  $\theta_1 \leq 0$  and  $\theta_2 \leq 0$  and  $\cdots \theta_n \leq 0$ . Which of the methods, 5.5.1 or 5.5.2, would be more powerful against the following alternatives:
  - a)  $\theta_1 = \delta; \theta_2 = \delta; \cdots \theta_n = \delta$
  - b)  $\theta_1 = \Delta; \theta_2 = 0; \cdots \theta_n = 0$ ?
5. Consider a repeated-measures problem you encountered in the literature. Which of the methods of Section 5.6 were applicable? Which would be most powerful against the alternatives of interest to you?

# CHAPTER 6

## Categorical Data

In many experiments and in almost all surveys, many if not all the results fall into categories rather than being measurable on a continuous or ordinal scale: e.g., male vs. female, black vs. Hispanic vs. oriental vs. white, in favor vs. against vs. undecided. The corresponding hypotheses concern proportions: “Blacks are as likely to be Democrats as they are to be Republicans.” Or, “the dominant genotype ‘spotted shell’ occurs with three times the frequency of the recessive.” In this chapter, you learn to test hypotheses like these that concern categorical and ordinal data.

### 6.1. Fisher’s Exact Test

As an example, suppose on examining the cancer registry in a hospital, we uncover the following data, which we put in the form of a  $2 \times 2$  *contingency table*

Table 6.1.

	Survived	Died	
Men	9	1	10
Women	4	10	14
	13	11	24

The nine denotes the number of males who survived, the 1 denotes the number of males who died, and so forth. The four marginal totals or *marginals* are 10, 14, 13, and 11. 10 is the total number of men in the study, 14 denotes the total number of women, and so forth.

We see in this table an apparent difference in the survival rates for men and women: Only 1 out of 10 men died following treatment, but 10 of the 14 women failed to survive. Is this difference statistically significant?

The answer is yes. Let's see why, using the same line of reasoning that Fisher advanced at the annual Christmas meeting of the Royal Statistical Society in 1934. After Fisher's talk was concluded, incidentally, a seconding speaker compared Fisher's talk to "the braying of the Golden Ass." I hope you will take more kindly to my own explanation. The preceding contingency table has several fixed elements: The total number of men in the survey, 10, the total number of women, 14, the total number who died, 11, and the total number who survived, 13. These totals are immutable; no swapping of labels will alter the total number of individual men and women or bring back the dead. But these totals do not determine the contents of the table as can be seen from the two tables reproduced below whose marginals are identical with those of our original table.

Table 6.2.

	Survived	Died	
Men	10	0	10
Women	3	11	14
	13	11	24

Table 6.3.

	Survived	Died	
Men	8	2	10
Women	5	9	14
	13	11	24

The first of these tables makes a strong case for the superior fitness of the male, stronger even than our original observations. In the second table, the survival rates for men and women are more alike than they were in our original table.

Fisher would argue that if the survival rates were the same for both sexes, then each of the redistributions of labels to subjects, that is, each of the  $N$  possible contingency tables with these same four fixed marginals, is equally likely, where

$$N = \sum_{x=0}^{10} \binom{13}{x} \binom{11}{10-x} = \binom{13+11}{10}.$$

How did we get this value for  $N$ ? The component terms are taken from the hypergeometric distribution:

$$\sum_{x=0}^t \binom{m}{x} \binom{n}{t-x} / \binom{m+n}{t}, \quad (6.1)$$

where  $n$ ,  $m$ ,  $t$ , and  $x$  occur as the indicated elements in the following  $2 \times 2$  contingency table

	Category 1	Category 2	
Category A	$x$	$t - x$	$t$
Category B	$m - x$	$n - (t - x)$	
	$m$	$n$	$m + n$

If men and women have the same probability of surviving, then all tables with the marginals  $m$ ,  $n$ ,  $t$  are equally likely, and  $\sum_{k=0}^{t-x} \binom{m}{t-k} \binom{n}{k}$  are as or more extreme.

In our example,  $m = 13$ ,  $n = 11$ ,  $x = 9$ , and  $t = 10$ , so that  $\binom{13}{9} \binom{11}{1}$  of the  $N$  tables are as extreme as our original table and  $\binom{13}{10} \binom{11}{0}$  are more extreme.

$11 \binom{13}{9} + \binom{13}{10}$  is a very small fraction of the total, so we conclude that a difference in survival rates as extreme as the difference we observed in our original table is very unlikely to have occurred by chance. We reject the hypothesis that the survival rates for the two sexes are the same and accept the alternative that, in this instance at least, males are more likely to profit from treatment.

### 6.1.1. One-Tailed and Two-Tailed Tests

In the preceding example, we tested the hypothesis that survival rates do not depend on sex against the alternative that men diagnosed as having cancer are likely to live longer than women similarly diagnosed. We rejected the null hypothesis because only a small fraction of the possible tables were as extreme as the one we observed initially. This is an example of a one-tailed test. Or is it? Wouldn't we have been just as likely to reject the null hypothesis if we had observed a table of the following form:

Table 6.4.

	Survived	Died	
Men	0	10	10
Women	13	1	14
	13	11	24

Of course, we would. In determining the significance level in the present example, we must add together the total number of tables that lie in either of the two extremes or tails of the permutation distribution.

McKinney et al. [1989] reviewed some 70 plus articles that appeared in six medical journals. In over half of these articles, Fisher's exact test was applied improperly. Either a one-tailed test had been used when a two-tailed test was

called for or the authors of the paper simply hadn't bothered to state which test they had used.

When you design an experiment, decide at the same time whether you wish to test your hypothesis against a two-sided or a one-sided alternative. A two-sided alternative dictates a two-tailed test; a one-sided alternative dictates a one-tailed test.

As an example, suppose we decide to do a follow-up study of the cancer registry to confirm our original finding that men diagnosed as having tumors live significantly longer than women similarly diagnosed. In this follow-up study, we have a one-sided alternative. Thus, we would analyze the results using a one-tailed test rather than the two-tailed test we applied in the original study.

### 6.1.2. The Two-Tailed Test

Unfortunately, it is not as obvious which tables should be included in the second tail. Is Table 6.4 as extreme as Table 6.2 in the sense that it favors an alternative more than the null hypothesis? One solution is simply to double the p-value we obtained for a one-tailed test. Alternately, we can define and use a test statistic as a basis of comparison. One commonly used measure is the  $\chi^2$  statistic defined for the  $2 \times 2$  contingency table after eliminating invariants as  $(x - t \frac{m}{m+n})^2$ . For Table 6.1, this statistic is 13; for Table 6.4, it is 29. We leave it to you to do the computations to show that Table 6.5 is more extreme than 6.1, but 6.6 is not.

Table 6.5.

	Survived	Died	
Men	1	9	10
Women	12	2	14
	13	11	24

Table 6.6.

	Survived	Died	
Men	2	8	10
Women	11	3	14
	13	11	24

### 6.1.3. Determining the p-Value

A problem with any of the methods we've used so far is that they produce only discrete significance values. We're very unlikely to observe 0.05 exactly; if the number of observations is small, p-values may jump from 0.045 to 0.062

(depending on the table) with no exact matches in between. There are at least five solutions:

- 1) Deliberately err on the conservative side, that is, reject when  $x \leq C(n,m,t)$ , where  $C(n,m,t)$  is the smallest integer for which the proportion of tables with the same marginals  $(n,m,t)$  is less than or equal to the significance level. See Boschloo [1970] and McDonald, Davis, and Miliken [1977] for some slight improvements on this approach.
- 2) Randomize on the boundary. If you get a p-value of 0.062 and the next closest value would have been 0.045, let the computer choose a random number between 0 and 1 for you. If this number is less than  $\frac{0.05 - 0.045}{0.062 - 0.045}$ , reject the hypothesis at the 5% level, accept it otherwise. Don't care to leave your decisions to chance?
- 3) Use the mid-p value. Let  $p$  be equal to half the probability of the table you actually observe plus all of the probability of more extreme results. See Lancaster [1961].
- 4) Present your audience with the data and the p-value you calculated; let them make up their own minds whether it is a significant result. See Section 3.1.3.
- 5) Conduct a sensitivity analysis. Add a single additional case to one of the cells (say, the cell that has the most observations already, so your addition will have the least percentage impact). Does the significance value change appreciably? (Again, leave it to your audience to decide what is an "appreciable" change.) This approach is particularly compelling if you are presenting statistical evidence in a courtroom as it turns impersonal percentages into individuals. See Good and Good [2000].

#### 6.1.4. What Is the Alternative?

In Chapter 2, we noted that every test requires both a hypothesis and an alternative, a concept we formalize in Chapter 14. We stated the null hypothesis for the  $2 \times 2$  contingency table, sort of; but what is the alternative? We need a model.

Let us assume we have taken two independent samples, the two columns, and that each sample consists of independent, identically distributed observations. In the first column, the  $m$  independent binomial trials resulted in  $x$  observations in the first row, the successes. In the second column, the  $n$  independent binomial trials resulted in  $t - x$  successes. As noted in Lehmann [1986, p. 154], the joint probability may be written as

$$\binom{m}{x} p_1^x (1 - p_1)^{m-x} \binom{n}{t-x} p_2^{t-x} (1 - p_2)^{n-(t-x)}$$

or

$$\binom{m}{x} p_1^x (1 - p_1)^{m-x} \binom{n}{t-x} p_2^{t-x} (1 - p_2)^{n-(t-x)}$$

or

$$\binom{m}{x} \binom{n}{t-x} (1-p_1)^m (1-p_2)^n \exp \left[ x \log[-\theta] + t \log \frac{p_2}{(1-p_2)} \right],$$

where  $\theta$  is the odds ratio  $\frac{p_2/(1-p_2)}{p_1/(1-p_1)}$ .

Our null hypothesis of identical distributions in the two columns is that  $p_1 = p_2$  or  $\theta = 1$ , while a one-sided alternative might be that  $\theta > 1$ .

### 6.1.5. Increasing the Power

Providing we are willing to randomize on the boundary, Fisher's exact test based on the conditional distribution of  $x$  given  $m$ ,  $n$ , and  $t$  is uniformly most powerful among all unbiased tests for comparing two binomial populations [Lehmann, 1986, p. 151–162].

It is most powerful under any of the following four world views:

- i) binomial sampling—one set of marginals in the contingency table is random; the other set and the sum  $s = n + m$  are fixed;
- ii) independent Poisson processes—all marginals and  $s$  are random;
- iii) multinomial sampling—all marginals are random and  $s$  fixed;
- iv) an experiment in which sampling is replaced by the random assignment of subjects to treatments—all marginals are fixed.

The power of Fisher's test depends strongly on the composition of the sample. A balanced sample, with equal numbers in each category is the most desirable. If the sample is too unbalanced, for example, if 100 of the observations have the attribute A and only one has the attribute notA, it may not be possible to determine if attribute B is independent of A.

If you have some prior knowledge about the frequency of A and B, then Berkson has suggested and Neyman has proved it is better to select samples of equal size from B and notB provided  $|p_B - \frac{1}{2}| > |p_A - \frac{1}{2}|$ . The “blind faith” method of selecting the sample at random from the population at large is worse than taking equal-sized samples from either A and notA or B and notB.

Studies of the power of Fisher's exact test against various alternatives were conducted by Haber [1987] and Irony and Pereira [1986]. It is easy to see that  $p(f_{11} = x) = \frac{\sum_{u=0}^1 \binom{x}{u} \binom{n-x}{t-u} \theta^u}{\sum_{u=0}^1 \binom{n}{u} \binom{t-u}{t-u} \theta^u}$ , where  $\theta$  is the odds ratio; this is the noncentral hypergeometric distribution [Fisher, 1934; Cornfield, 1956].

## 6.2. Odds Ratio

In most instances, we won't be satisfied with merely rejecting the null hypothesis, but will want to make some more powerful statement like “men are twice as likely

Table 6.7. Transco Employment

Outcome	Young	Old
Fired	1	10
Retained	24	17

as women to get a good-paying job” or “women under thirty are twice as likely as men over 40 to receive an academic appointment.”

In the discrimination case of Fisher vs. Transco Services of Milwaukee [1992], the plaintiffs claimed that Transco was ten times as likely to fire older employees. Can we support this claim with statistics? The Transco data are provided in Table 6.7.

Let  $p_1$  denote the probability of firing a young person, and  $\pi_2$  the probability of firing an older person. We want to go beyond testing the null hypothesis  $p_1 = p_2$  to determine a confidence interval for the odds ratio  $\theta = \frac{p_2/(1-p_2)}{p_1/(1-p_1)}$ .

One can obtain confidence intervals for the odds ratio by iterative methods as described in Section 3.2; see also Cornfield [1956]. Mantel and Hankey [1971], and Thomas [1971]. Baptista and Pike [1977] describe an approach that sometimes gives shorter confidence intervals. We turn for aid to StatXact™, a statistical package whose emphasis is the analysis of categorical and ordinal data. Choosing statistics, two binomials, and CI odds ratio from successive StatXact menus, we obtain the results of Figure 6.1.

Based on these results, we can tell the judge that older workers were fired at a rate at least 1.6 times the rate at which younger workers were discharged.

### 6.2.1. Stratified $2 \times 2$ 's

In trying to develop a cure for a relatively rare disease, we face the problem of having to gather data from a multitude of test centers, each with its own set of

Statistic based on the observed  $2 \times 2$  table:

Binomial proportion for column <young >: pi\_1 = 0.04000  
 Binomial proportion for column <old >: pi\_2 = 0.3704

$$\text{Odds Ratio} = \frac{(pi\_2)/(1-pi\_2)}{(pi\_1)/(1-pi\_1)} = 14.12$$

Results:

p-value(2-sided)      95.00% Confidence Interval  
 0.007145      (1.649,      636.5)

Figure 6.1. StatXact-4 Output Datafile: C:\SX3WIN\EXAMPLES\TRNSCO.CY3 ODDS RATIO OF TWO BINOMIAL PROPORTIONS.

Table 6.8. Sandoz Drug Data

Test Site	New Drug		Control Drug	
	Response	#	Response	#
1	0	15	0	15
2	0	39	6	32
3	1	20	3	18
4	1	14	2	15
5	1	20	2	19
6	0	12	2	10
7	3	49	10	42
8	0	19	2	17
9	1	14	0	15
10	2	26	2	27
11	0	19	2	18
12	0	12	1	11
13	0	24	5	19
14	2	10	2	11
15	0	14	11	3
16	0	53	4	48
17	0	20	0	20
18	0	21	0	21
19	1	50	1	48
20	0	13	1	13
21	0	13	1	13
22	0	21	0	21

procedures and its own way of executing them. Before we can combine the data, we must be sure the odds ratios across the test centers are approximately the same. Consider the set of results in Table 6.8 obtained by the Sandoz drug company and reproduced with permission from the StatXact-3 manual. One of the sites, number 15, stands out from the rest. But is the difference statistically significant?<sup>1</sup>

Zelen [1971] proposed a test based on the number of  $2 \times 2 \times 22$  tables with the same marginals that are as likely or are less likely than the table which was actually observed. With 22 contingency tables, the number of computations needed to examine all rearrangements is in the billions. Fortunately, StatXact utilizes several time-saving algorithms, including the one introduced in Mehta, Patel, and Senchaudhuri [1988] to obtain a Monte Carlo estimate of the significance level. We pull down menus Statistics, Stratified  $2 \times 2$  Tables, Homogeneity of Odds Ratios to obtain the results in Figure 6.2.

The estimated p-value of .013, just a fraction greater than 1%, tells us it would be unwise to combine the results from the different sites.

<sup>1</sup> Similar problems were encountered in a study in which test subjects might use one of several different “identical” machines. I couldn’t combine the results from the different machines or the different technicians who operated them, until I performed an initial test of their equivalence.

[18 2 × 2 informative tables]

Observed Statistics:

BD: Breslow and Day Statistic = 25.78  
ZE: Zelen Statistic = 9.481e - 009

Asymptotic p-value: (based on chi-square distribution with 17 df)  
 $\Pr\{\text{BD.GE. } 25.78\} = 0.0785$

Monte Carlo estimate of p-value:

$\Pr\{\text{ZE.GE. } 9.481\text{e-009}\} = 0.0127$   
99.00% Confidence Interval = (0.0119, 0.0135)

Elapsed Time is 0:16:15.37 (10000 tables sampled; starting seed 85190)

Figure 6.2. StatXact-4 Output Datafile: C:\SX3WIN\EXAMPLES\SANDOZ.CY3 TEST FOR HOMOGENEITY OF ODDS RATIOS.

The output of the StatXact program provides us with one more important findings. Displayed above the Monte Carlo estimate of the exact p-value, 0.01237, is the asymptotic or large-sample approximation based on the chi-square distribution. Its value, 0.0785, is many times larger than the correct value, and relying on this so-called approximation would have led us to a completely different and erroneous conclusion.

### 6.3. Exact Significance Levels

The preceding result is not an isolated one. Asymptotic approximations are to be avoided except with very large samples. Table 6.9 contains data on oral lesions in three regions of India derived from Gupta et al. [1980] by Mehta and Patel. We want

Table 6.9. Oral Lesions in Three Regions of India

Site of Lesion	Kerala	Gujarat	Andh
Labial Mucosa	0	1	0
Buccal Mucosa	8	1	8
Commissure	0	1	0
Gingiva	0	1	0
Hard Palate	0	1	0
Soft Palate	0	1	0
Tongue	0	1	0
Floor of Mouth	1	0	1
Alveolar Ridge	1	0	1

Table 6.10. Three Tests of Independence.

Statistic	$\chi^2$	F-H	LR
Exact p-value	.0269	.0101	.0356
Tabulated p-value	.1400	.2331	.1060

to test the hypothesis that the location of oral lesions is unrelated to geographical region. Possible test statistics include Freeman–Halton  $p$  (see Section 6.4),  $p_\chi$ , and  $p_L$ . This latter statistic is based on the log-likelihood ratio  $\sum \sum f_{ij} \log(f_{ij} f_{..}/f_i f_{.j})$ .

We may calculate the exact significance levels of these test statistics by deriving their permutation distributions or use asymptotic approximations obtained from tables of the chi-square statistic. Table 6.10 taken from the StatXact-3 manual compares the various approaches.

The exact significance level varies from 1% to 3.5%, depending on which test statistic we select. Tabulated p-values based on large-sample approximations vary from 11% to 23%. Using the Freeman–Halton statistic, the permutation test tells us the differences among regions are significant at the 1% level; the large-sample approximation says no, they are insignificant even at the 20% level. The permutation test is correct. The large-sample approximation is grossly in error. With so many near-zero entries in the original contingency table, the chi-square large-sample approximation is not appropriate.<sup>2</sup>

## 6.4. Unordered $r \times c$ Contingency Tables

With a computer at hand, the principal issue in the analysis of a contingency table with  $r$  rows ( $r > 2$ ) and  $c$  columns ( $c > 2$ ) is deciding on an appropriate test statistic. Halter [1969] showed that we can find the probabilities of any individual  $r \times c$  contingency table through a straightforward generalization of the hypergeometric distribution given in equation 6.1. An  $r \times c$  contingency table consists of a set of frequencies  $\{f_{ij}, 1 \leq i \leq r; 1 \leq j \leq c\}$  with row marginals  $\{f_{i.}, 1 \leq i \leq r\}$  and column marginals  $\{f_{.j}, 1 \leq j \leq c\}$ . Suppose once again we have mixed up the labels. To make matters worse, this time every item/subject is to be assigned both a row and a column label from the  $r + c$  stacks of labels,  $f_1$ , of which are labeled row 1,  $f_2$ , of which are labeled row 2, and so forth.

Let  $P$  denote the probability with which a specific table assembled at random will have these exact frequencies.  $P = Q/R$  with<sup>3</sup>

$$Q = \prod_{i=1}^r f_{i.}! \prod_{j=1}^c f_{.j}! f_{..}!$$

<sup>2</sup> See, also, Mudholkar and Hutson [1997].

<sup>3</sup>  $\prod_{i=1}^n f_i! = f_1! f_2! \cdots f_n!$

and

$$R = \prod_{i=1}^r \prod_{j=1}^c f_{ij}!$$

An obvious extension of Fisher's exact test is the Freeman and Halton [1951] test based on the proportion  $p$  of tables for which  $P$  is greater than or equal to  $P_0$  for the original table.

While the extension itself may be obvious, it's not as obvious that this extension offers any protection against the alternatives of interest. Just because one table is less likely than another under the null hypothesis does not mean it is going to be more likely under the alternatives of interest to us. Consider the  $1 \times 3$  contingency table  $\begin{bmatrix} f_1 & f_2 & f_3 \end{bmatrix}$ , which corresponds to the multinomial with probabilities  $p_1 + p_2 + p_3 = 1$ , the table whose entries are 1, 2, 3 argues more in favor of the null hypothesis  $p_1 = p_2 = p_3$  than of the ordered alternative  $p_1 > p_2 > p_3$ .

The classic statistic for independence in a contingency table with  $r$  rows and  $c$  columns is

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c (f_{ij} - Ef_{ij})^2/Ef_{ij},$$

where  $Ef_{ij}$  is the number of observations in the  $ij$ th category one would expect on theoretical grounds.

With very large samples, this statistic has the chi-square distribution with  $(r-1)(c-1)$  degrees of freedom. But in most practical applications, the chi-square distribution is only an approximation and notoriously inexact for small and unevenly distributed samples.

The permutation statistic based on the proportion  $p_\chi$  of tables for which  $\chi^2$  is greater than or equal to  $\chi_0^2$  for the original table provides an exact test and possesses all the advantages of the original chi-square. The distinction between the two approaches, as we observed in Chapter 2, is that with the original chi-square we look up the significance level in a table, while with the permutation statistic, we derive the significance level from the permutation distribution. With large samples, the two approaches are equivalent, as the permutation distribution converges to the tabulated distribution (see Chapter 14 of Bishop, Fienberg, and Holland [1975]).

An alternative is the likelihood ratio test based on the statistic

$$p_L = 2 \sum_{i=1}^r \sum_{j=1}^c f_{ij} \log \left[ \frac{f_{ij}}{f_i f_j / f_{..}} \right].$$

If you wish to compare the strength of the association between the row and column variables across different  $r \times c$  tables having different row and column dimensions, use one of the contingency coefficient's described by Liebetrau [1983]. One example is Cramer's V, which ranges between 0 and 1, with 0 signifying no

association and signifying total dependence:

$$V = \sqrt{\frac{\chi^2}{f_{..}(\min[r, c] - 1)}}.$$

All these permutation tests have one of the original chi-square test's disadvantages: While they offer global protection against a wide variety of alternatives, they offer no particular protection against any single one of them. The statistics  $p$ ,  $p_\chi$ , and  $p_L$  treat row and column categories symmetrically, and no attempt is made to distinguish between cause and effect. To address this deficiency, Goodman and Kruskal [1954] introduce an asymmetric measure of association for nominal scale variables called tau  $\tau$ , which measures the proportional reduction in error obtained when one variable, the "cause" or independent variable, is used to predict the other, the "effect" or dependent variable.

Assuming the independent variable determines the row,

$$\tau = \frac{\sum_j f_{mj} - f_{..}}{f_{..} - f_{m.}},$$

where  $f_{mj} = \max_i f_{ij}$  and  $f_{m.} = \max_i f_{i.}$

$0 \leq \tau \leq 1$ .  $\tau = 0$  when the variables are independent;  $\tau = 1$  when, for each category of the independent variables, all observations fall into exactly one category of the dependent. These points are illustrated in the following  $2 \times 3$  tables:

3	6	9	$\tau = 0$
6	12	18	

18	0	0	$\tau = 1$
0	36	0	

3	6	9	$\tau = 0.166$
12	18	6	

A permutation test of independence is based on the proportion of tables  $p_\tau$  for which  $\tau \geq \tau_0$ .

An alternative is the uncertainty coefficient derived from the likelihood ratio statistic

$$U_{R|C} = \frac{\sum_{i=1}^r \sum_{j=1}^c f_{ij} \log[E f_{ij}]}{\sum_{i=1}^r f_{i.} \log[f_{i.}/f_{..}]}.$$

Like tau, this statistic measures the reduction in error and ranges between 0 and 1 as the association ranges from complete independence to complete dependence of the row, column variables.

### 6.4.1. Agreement Between Observers

Suppose two observers assign the same set to various categories so the results can be put in the form of an  $r \times r$  table. An example would be two teachers assigning letter grades to the same set of students. A permutation test based on Cohen's Kappa, as described in Agresti [1990] and Berry and Mielke [1988], allows us to measure the degree of agreement between the two observers:

$$\kappa = \frac{f_{..} \sum_{i=1}^r (f_{ii} - f_{i.} f_{.i})}{f_{..}^2 - \sum_{i=1}^r f_{i.} f_{.i}}.$$

Note that  $0 \leq \kappa \leq 1$ .

Which Test?

The data are in categories

The categories can't be ordered

There are exactly two rows and two columns

Use Fisher's Exact Test

There are more than two rows and at least two columns

You want to test whether the relative frequencies are the same in each row and in each column

Use the Freedman–Halton Test or use chi-square

You want to test whether the column frequencies depend on the row

Use Tau or the uncertainty coefficient

The number of rows is equal to the number of columns

You want to test whether the row and column classifications are in agreement

Use Kappa

### 6.4.2. What Should We Randomize?

Table 6.11a summarizes Clarke's [1960, 1962] observations on the relation between habitat and the relative frequencies of different varieties of *C. nemoralis* snail. It is tempting to analyze this table using the methods of the preceding section, but before we can analyze a data set, we need to understand *how* it was collected. In this instance, observers went to a series of locations in southern England. At each location, they noted the type of habitat—beechwoods, grasslands, and so forth, and the frequencies of each of 12 different varieties of snail. The original findings are summarized in Table 6.11b, reproduced from Manly [1983]. Note that each row in this table corresponds to a single multivariate observation.

Table 6.11a. Summary of Clarke's [1960, 1962] Data on *C. nemoralis*

Habitat	N1	N2	N3	N4	N5	N6	N7	N8	N9	N10	N11	N12
Beechwoods	9	1	34	26	0	46	8	59	126	6	40	115
Other deciduous	10	1	1	0	0	85	8	13	44	2	1	12
Fens	73	3	8	4	6	89	1	23	21	11	0	22
Hedgerows	76	15	32	19	36	98	3	12	8	14	1	18
Grasslands	49	29	75	7	28	23	17	60	12	14	14	24

Manly computed the chi-square statistic for the original data as summarized in Table 6.11a. Then, using the information in Table 6.11b, he randomly reassigned the location labels to different habitats, preserving the number of locations at each habitat. For example, in one of the rearrangements, Clipper Down Wood, Boarstall Wood, Hatford, and Charlbury Hill and only these four locations were designated as Fens. He formed a summary table similar to 6.11a for each rearrangement and computed the chi-square statistic for that table. He found the original value of the chi-square statistic 1756.9 was greater than any of the values he observed in each of 500 random reassignments and concluded that habitat type has a significant effect on the distribution of the various body types of the *C. nemoralis* snail.

Manly's analysis combines multivariate and categorical techniques. It makes optimal use of all the data because it takes into account how the data were collected. Could Manly have used Table 6.11b alone to analyze the data? No, because this table lacks essential information about interdependencies among the various types of snail.

#### 6.4.3. Underlying Assumptions

The assumptions that underlie the analysis of an  $r \times c$  contingency table are the same as those that underlie the analysis of the  $k$ - or  $r$ -sample problem. To see this, note that a contingency table is merely a way of summarizing a set of  $N$  bivariate observations. We may convert from this table to  $r$  distinct samples by using the first or row observation as the sample or treatment label and the second or column observation as the "value." Keeping the marginals fixed while we rearrange the labels ensures that the  $r$  sample sizes and the  $N$  individual values remain unchanged.

As in the  $r$ -sample problem, the labels must be exchangeable under the null hypothesis. This entails two assumptions: First, that the row and column scores are mutually independent, and second, that the observations themselves are independent of one another. We as statisticians can only test the first of these assumptions. We rely on the investigator to ensure that the latter assumption is satisfied. (See question 4 at the end of this chapter.)

Table 6.11b. Clarke's [1960, 1962] Data\* on *C. nemoralis*

Habitat type	Location	N1	N2	N3	N4	N5	N6	N7	N8	N9	N10	N11	N12
Beechwoods	Clipper Down Wood	1	0	0	0	0	8	0	1	12	1	0	0
	Hackpen Wood	0	0	5	4	0	0	0	5	20	0	1	1
	Kingstone Coombes	0	0	0	2	0	4	1	0	0	0	0	2
	Danks Down Wood	0	0	2	0	0	9	0	15	21	0	1	27
	Fawley Bottom Wood	0	1	0	0	0	5	3	0	2	3	0	0
	Maidensgrove Wood	0	0	0	0	0	3	2	0	5	2	0	0
	Aston Rowant Wood	0	0	0	0	0	6	1	0	23	0	0	0
	Rockley Wood	0	0	10	15	0	0	0	4	20	0	0	21
	Manton Wood	0	0	3	1	0	0	1	6	2	0	3	9
	Knoll Down A	3	0	0	0	0	8	0	9	2	0	35	47
	Knoll Down B	0	0	7	4	0	0	0	0	0	0	0	8
	Roundway Wood	5	0	7	0	0	3	0	19	20	0	0	0
Other deciduous woods	Boarstall Wood	0	0	0	0	0	13	0	9	28	1	0	0
	Rockley Copse	9	1	1	0	0	63	8	4	10	0	0	8
	Elsfield Covert	1	0	0	0	0	6	0	0	4	0	0	0
	Uffington Wood 2	0	0	0	0	0	3	0	0	2	1	1	4
Fens	Shippon	54	1	3	3	1	54	0	8	13	7	0	20
	Headington Wick	5	1	3	0	2	14	1	13	4	2	0	0
	Cothill Fen	2	1	1	0	1	3	0	0	1	1	0	0
	Shippon Fen 2	12	0	1	1	2	18	0	2	3	1	0	2
Hedgerows and rough herbage	Hatford	1	1	0	15	0	2	0	1	3	2	0	4
	Shepherd's Rest 1	16	7	9	0	19	11	1	0	0	6	0	0
	Shepherd's Rest 2	13	4	4	0	9	0	1	0	0	1	0	0
	Standford in Vale	5	0	0	0	0	5	0	1	4	0	0	0
	Wootton	2	0	3	0	0	7	0	1	0	0	0	0
	Chiseldon	6	2	0	2	4	9	0	0	0	1	0	1
	Faringdon	18	0	8	0	1	34	0	5	0	0	1	4
	The Ham	8	0	2	1	1	1	0	1	0	0	0	9
	Wanborough Plain	2	0	0	0	0	24	0	0	1	3	0	0
	Watchfield	3	1	0	0	0	2	1	0	0	1	0	0
	Hill Barn Tumulus	1	0	5	0	0	0	0	3	0	0	0	0
	Little Hinton	1	0	1	1	2	3	0	0	0	0	0	0
Crasslands	Charlbury Hill	2	0	5	1	0	1	0	4	7	0	0	5
	White Horse 1	4	10	4	0	3	3	3	7	0	1	2	1
	White Horse 2	6	6	10	0	0	0	0	0	0	0	0	0
	White Horse 3	7	2	12	0	7	7	4	5	0	2	0	0
	White Horse 4	7	0	2	0	2	0	1	1	0	0	0	0
	Dragons Hill 1	2	4	5	0	0	3	4	19	0	5	2	4
	Dragons Hill 2	1	1	6	0	0	0	1	4	0	0	2	2
	Dragons Hill 3	1	2	3	0	2	2	3	12	0	1	0	4
	West Down 1	0	1	4	3	1	0	0	0	0	0	7	2
	West Down 2	0	0	5	3	0	0	0	0	1	1	0	5
	Sparsholt Down	13	1	15	0	6	0	0	0	0	0	0	0
	Little Hinton	5	0	1	0	5	5	0	1	3	1	0	0
	White Horse 5	0	2	2	0	0	1	0	1	0	1	1	0
	Dragons Hill 4	1	0	2	0	2	1	1	6	1	2	0	1

\*The morph types are similar to those for *hortensis*, with up to five bands present. They are: N1, yellow fully banded (Y12345); N2, yellow part-banded (N00345); N3, yellow mid-banded (Y00300); N4, yellow unbanded (Y00000); N5, other yellows; N6, pink fully banded (P12345); N7, pink part-banded (P00345); N8, pink mid-banded (P00300); N9, pink unbanded (P00000); N10, other pinks; N11, brown banded; N12, brown unbanded.

Note: From "Analysis of Polymorphic Variation in Different Types of Habitat" by BFJ Mainly, which appeared in *Biometrics*; 1983; 16: 13–27. Reprinted with permission from the Biometric Society.

## 6.5. Ordered Contingency Tables

When data are measured on a continuous basis, such as 1.12, 1.13, 1.14, ties are a relatively infrequent occurrence, but when we ask someone to provide a self-rating on a discrete ordinal scale, 1 through 5, for example, ties are inevitable, the rule, not the exception, and the methods of this chapter may be more appropriate for analyzing such ordinal data than those of Chapter 3.

Table 6.12. Data Gathered by Graubard and Korn [1987]

	Maternal Alcohol Consumption (drinks/day)					
Malformation	0	<1	1–2	3–5	≥6	Total
Absent	17066	14464	788	126	37	32481
Present	48	38	5	1	1	93
	17114	14502	793	127	38	32574

### 6.5.1. Ordered $2 \times c$ Tables

Our analysis of a  $2 \times c$  ordered contingency table is straightforward and parallels the approach used in Section 3.5.2 for a  $k$ -sample comparison, once we have determined what value to assign each of the ordered categories. We illustrate this dilemma with data gathered by Graubard and Korn [1987].

Recall that our test statistic is  $\sum g[j] f_{1j}$ , where  $g$  is any monotone increasing function. Among the leading choices for a scoring method  $g$  are

- i) the category number: 1 for the 1st category, 2 for the second, and so forth,
- ii) the midrank scores,
- iii) scores determined by the user, the choice we made in Section 3.5.2 when we analyzed the micronuclei data.

Consider the following  $1 \times 2$  contingency table

	Alcohol Consumption	
drinks/day	0	1–2
frequency	3	5

The category or equidistant scores are 1 and 2. The ranks of the eight observations are 1 through 3, and 4 through 8, so that the midrank score of those in the first category is 2, and in the second 6. Our user-chosen scores, corresponding to alcohol consumption, are 0 and 1.5.

Using Testimate™ to analyze the full Gaubard–Korn data set, we obtain p-values that range from the insignificant, 0.29 for midrank scores, to marginally significant, 0.10 for the equidistant scores, to highly-significant, 0.01 for our user-chosen

Table 6.13. Response to Chemotherapy

	None	Partial	Complete
CTX	2	0	0
CCNU	1	1	0
MTX	3	0	0
CTX + CCNU	2	2	0
CTX + CCNU + MTX	1	1	4

scores.<sup>4</sup> A user-chosen score based on the user's knowledge of underlying cause and effect is always recommended as it will be the most effective at distinguishing between hypothesis and alternative.

### 6.5.2. More than Two Rows and Two Columns

Two cases need to be considered: The first when the columns but not the rows of the table may be ordered (the other variable being purely categorical), and the second when both columns and rows can be ordered.

#### 6.5.2.1. Singly Ordered Tables

Several tests have been proposed; see Agresti [1992], Haberman [1974], and Soms [1985]. Our approach parallels that of Section 4.2 in which we describe a two-way analysis of variance. Our test statistic is

$$F_2 = \Sigma(T_i - \bar{T})^2 \quad \text{where} \quad T_i = \Sigma g_j f_{ij}.$$

As in the case of the  $2 \times c$  table, our problem is in deciding on the appropriate scores  $\{g_j\}$ . Among the proposals are ranks, normal scores, and Savage scores.

Table 6.13 provides tumor remission data for five chemotherapy regimes. As partial response corresponds to approximately two years in remission (about 100 weeks) and complete response to an average of three years (150 weeks), we assign scores of 0, 100, and 150 to the ordered response categories.

To use StatXact-3 to analyze the tumor data, we first click on TableData in the main menu, click on Settings, then enter Column Scores. To execute the analysis, we select in turn Statistics, Singly Ordered  $r \times c$  Table, ANOVA with Arbitrary Scores, and Exact test method. The results are in Figure 6.3.

Although our estimated significance level is less than 0.0444, a 99% confidence interval for this estimate, based on a sample of 3200 possible rearrangements, does include values greater than 0.05. We can narrow this confidence interval by sampling additional rearrangements. With a Monte Carlo of 10000, our Monte

<sup>4</sup> The chi-square approximation yields values ranging from 0.017 for the classic Pearson chi-square statistic to 0.19 for the likelihood ratio; Agresti [1992].

Statistic based on the observed data:  
The Observed Statistic = 6.507  
Asymptotic p-value: (based on chi-square distribution with 4 df)  
 $\Pr\{\text{Statistic. GE. } 6.507\} = 0.0747$   
Monte Carlo estimate of p-value:  
 $\Pr\{\text{Statistic. GE. } 6.507\} = 0.0444$   
99.00% Confidence Interval = (0.0350, 0.0538)

Elapsed Time i0:1:20.58 (3200 tables sampled with starting seed 16229)

Figure 6.3. StatXact-4 Output Datafile: C:\SX3WIN\EXAMPLES\TUMOR.CY3 ANOVA TEST [That the 5 rows are identically distributed].

Carlo estimate of p-value is 0.0434 with a 99% confidence interval of (0.0382, 0.0486). Of course, the calculations take three times as long. When I first perform a permutation test, I use as few as 400 to 1600 simulations. If the results are equivocal, as they were in this example, then and only then will I run 10,000 simulations.

#### 6.5.2.2. Doubly Ordered Tables

Our log-linear model is that

$$\log[\mathbb{E} n_{ij}] = \mu + \lambda_i^X + \lambda_j^Y + \lambda_{ij}.$$

In an  $r \times c$  contingency table conditioned on fixed marginal totals, Cornfield [1956] showed that the outcome depends only on the  $(r - 1)(c - 1)$  odds ratios

$$\phi_{ij} = \frac{\pi_{ij}\pi_{i+1,j+1}}{\pi_{i,j+1}\pi_{i+1,j}},$$

where  $\pi_{ij}$  is the probability of an individual being classified in row  $i$  and column  $j$ .

In a  $2 \times 2$  table, conditional probabilities depend on a single odds ratio, and hence, one- and two-tailed tests of association are easily defined. In an  $r \times c$  table, there are potentially  $n = 2(r - 1)(c - 1)$  sets of extreme values, two for each odds ratio. Hence, an omnibus test for no association, e.g.,  $\chi^2$ , might have as many as  $2^n$  tails.

Following Patefield [1982], we consider tests of the null hypothesis of no association between row and column categories  $H: \phi_{ij} = 1$  for all  $i, j$  against the alternative of a positive trend  $K: \phi_{ij} \geq 1$  for all  $i, j$ .

The two principal test statistics considered by Patefield [1982], are

$$\lambda_3 = \sum \sum f_{ij} r_i c_j,$$

Table 6.14. Incidents of Pairs of  
Anencephalic Infants

km apart	Months Apart		
	<1	<2	<4
<1	39	101	235
<5	53	156	364
<25	211	652	1516

where  $\{r_i\}$  and  $\{c_j\}$  are user-chosen row and column scores,<sup>5</sup> and

$$\lambda_2 = \sup \sum \sum f_{ij} x_i y_j,$$

where the supremum is taken over all sets  $\{x_i\}$  and  $\{y_j\}$ , satisfying the conditions  $\sum f_{i..} x_i = 0$ ,  $\sum f_{..j} y_j = 0$ ,  $\sum f_{i..} x_i^2 = n_{..}$ ,  $\sum f_{..j} y_j^2 = f_{..}$ ; and

$$x_1 \leq x_2 \cdots \leq x_r; \quad y_1 \leq y_2 \cdots \leq y_c.$$

Patefield finds that  $\lambda_2$  has higher power than the linear-by-linear association test  $\lambda_3$  when some but not all of the odds ratios  $\phi_{ij}$  are close to unity, whereas  $\lambda_3$  has higher power than  $\lambda_2$  when the odds ratios all have about the same value.

The likelihood ratio behaves like  $\lambda_2$ ; the Goodman and Kruskal test of association behaves like  $\lambda_3$ .

Other possible statistics, including one based on the difference between the numbers of concordant and discordant pairs, are considered by Agresti and Wackerly [1977].

### 6.5.2.3. An Example

An ongoing fear of many parents is that something in the environment— asbestos or radon in the walls of their house, or toxic chemicals in their air and ground water, will affect their offspring. Table 6.14 is extracted from data collected by Siemiatycki and McDonald [1972] on congenital neural-tube defects. Eyeballing the gradient along the diagonal of this table, one might infer that births of anencephalic infants occur in clusters. The question arises as to which measures of distance and time we should employ as weights  $\{r_i\}$  and  $\{c_j\}$ .

Mantel [1967] reports striking differences between one analysis of epidemiologic data in which the coefficients are proportional to the differences in position and a second approach (which he recommends) to the same data in which the coefficients are proportional to the reciprocals of these differences. Using Mantel's approach, a pair of infants born 5 km and three months apart contribute  $\frac{1}{3} * \frac{1}{5} = \frac{1}{15}$  to the correlation. Summing the contributions from all pairs, then repeating the summing process for a series of random rearrangements, Siemiatycki

<sup>5</sup> This statistic is actually just another form of Mantel's U, perhaps the most widely used of all multivariate statistics. See Chapter 9.

and McDonald conclude that the clustering of anencephalic infants is not statistically significant.

Entering the data into StatXact-3, then selecting Statistics, Doubly Ordered Entering the data into StatXact-3, then selecting Statistics, Doubly Ordered  $r \times c$  Table. Linear-by-Linear, and Exact from the menus, and using the weights (2,.67,.33) and (2, .4, .08), We verify that a table this extreme would occur less than 1 in 5,000 times by chance.

## 6.6. Covariates

The presence of a covariate adds a third dimension to a contingency table. Bross [1964] studies the effects of treatment on the survival of premature infants. His results are summarized in Table 6.15.

Table 6.15. Effect of Treatment of Survival of the Premature

	Dead	Recovered	
Placebo	6	5	11
Treatment	2	12	14
	8	17	23

These results, though suggestive, are not statistically significant.

Bross notes that survival is very much a function of a third, concomitant variable—the birth weight of the child. A low birth weight indicates greater prematurity and, hence, greater odds against a child's survival. An analysis of treatment is out of the question unless, somehow, he can correct for the effects of birth weight.

A solution we studied in earlier chapters is to set up an experiment in which we study the effects of treatment in pairs that have been matched on the basis of birth weight. But Bross' study of the premature was not an experiment; he could only observe, not control birth weight.

Table 6.16 depicts his first nine observations, ordered by birth weight. The last two columns of this table deserve explanation. The column headed NI records the number of cases in which a child of lower birth weight treated with ukinase recovered when an untreated child of higher birth weight died. Such a result is to be expected under the alternative of a positive treatment effect, though it would occur only occasionally by chance under the null hypothesis.

The column headed I records the number of cases in which an untreated child of lower birth weight recovered when a child of higher birth weight treated with ukinase died. Such an event or inversion would be highly unlikely under the alternative.

As his test statistic, Bross adopts

$$S = \frac{(NI - I)^2}{NI + I}.$$

Table 6.16. Effect of Treatment and Birth Weight on Survival of the Premature

Weight	Treatment	Outcome	NI	I
			TR/PL	PL/TR
1.08	TR	D		
1.13	TR	R	3	
1.14	placebo	D		
1.20	TR	R	2	
1.30	TR	R	2	
1.40	placebo	D		
1.59	TR	D		
1.69	TR	R	1	
1.88	placebo	D		

Note that  $NI = 8$ ,  $I = 0$ , and  $S = 8$  for the original observations. Bross computes  $S$  for each of the  $\binom{9}{3}$  possible rearrangements of the treatment labels—and only the labels were changed, the pairing of birth weight with outcome was preserved. None of the other rearrangements yield as large a value of  $S$  as the original observations. Bross concludes that the treatment has a statistically significant effect on survival of the premature.

## 6.7. Higher Dimensional Tables

Another way to correct for the effects of a covariate is to divide the observations into blocks, so that the value of the covariate is approximately constant within each block. Under the assumption that the odds ratio is the same for each block, Mehta, Patel, and Gray [1985] provide a method for combining the results from several  $2 \times 2$  contingency tables. To test the assumption of a constant odds ratio,  $\theta_1 = \theta_2 = \dots = \theta_B$ , use Zelen's test (see Section 6.2.1).

On the other hand, it may be that an apparent association between the variables determining the rows and columns of a contingency table is actually the result of an association with a third factor. By separating the data into blocks based on the values of this third factor, we may test this latter assumption,  $\theta_1 = \theta_2 = \dots = \theta_n = 1$ . Lehmann [1986, p. 162–166] showed that a UMPU test exists and is given by rejecting if  $T = \sum_{k=1}^B f_{11k}$  is an extreme value relative to that of other tables with the same fixed marginals.

A two-sided test can be obtained by doubling the exact one-sided p-value, or by specifying that  $|T - ET|$  be less than or equal to the value actually observed.

Birch [1964] showed that this result can be extended to  $B 2 \times c$  contingency tables whose  $c$  columns are ordered, using the linear rank statistic given by

$T = \sum_{k=1}^B \sum_{j=1}^c w_j f_{1jk}$ , where the weights or scores  $w_j$  are selected as described in Section 6.5.1.<sup>6</sup>

Agresti [1992] showed this result may be extended still further to tests of the conditional independence of the row and column variables in an  $r \times c$  table given a third blocking variable. If we can assume that the  $(r - 1)(c - 1)$  odds ratios are identical for all values of the blocking factor, our test statistic is

$$\mathbf{d}' V^{-1} \mathbf{d},$$

where  $\mathbf{d}$  is the matrix with elements  $d_{ij} = \sum_k [f_{ijk} - \frac{f_{i,k} f_{j,k}}{f_{.,k}}]$ ,  $i = 1, \dots, R - 1$ ;  $j = 1, \dots, C - 1$ ; and  $\mathbf{V}$  is the null covariance matrix of  $\mathbf{d}$ .

If we cannot assume the odds ratios are identical under the alternative, then we may still test for conditional independence using the statistic  $\sum \chi_k^2$ , where  $\chi_k^2$  is the chi-square statistic for testing independence of rows and columns within the  $k$ th level of the blocking factor.

## 6.8. To Learn More

Excellent introductions to the analysis of contingency tables may be found in Agresti [1990, 1992], and in the StatXact-4 manual authored by Mehta and Patel. Major advances in analysis by resampling means have come about through the efforts of Gail and Mantel [1977], Mehta and Patel [1983], Mehta, Patel, and Senchaudhuri [1988], Baglivio, Oliver, and Pagano [1988, 1992], and Smith, Forster and McDonald [1996]. Chapter 13 is devoted to this topic.

Berkson [1978], Basu [1979], Garside and Mack [1976], Haber [1987], and Mielke and Berry [1992] examine Fisher's exact test. For some alternatives to the Zelen test statistic for testing for homogeneity among tables, see Jones, O'Gorman, Lemke, and Woolson [1989] and Liang and Self [1985].

The power of the Freeman-Halton statistic in the  $r \times 2$  case is studied by Krewski, Brennan, and Bickis [1984]. For a description of some other, alternative statistics for use in  $r \times c$  contingency tables, see Nguyen [1985]. For a review of the literature on higher dimensional tables, see Agresti [1992].

## 6.9. Questions

1. a) Fisher's original presentation of his "exact test" was marked by acrimony and dissent. You may wonder what all the fuss was about. Fisher's exact test agrees asymptotically with the chi-square test based on one degree of freedom, a fact that is no longer in dispute. But many of the participants at the meeting raged over whether there should

<sup>6</sup> We need to assume the absence of an underlying joint dependence among the three variables.

be three or four degrees of freedom corresponding to the number of marginals or just one degree as Fisher asserted. To find out why, read the discussions following Fisher [1934] as well as Box [1978].

- b) Surprisingly, the controversy is not dead, though it has taken a somewhat different form. The analyst may chose to view the table as we have here, conditional on the margins, or as the result of taking two separate and independent Binomial samples, the so-called unconditional case. The unconditional approach appears to me less desirable as it requires we average into the p-value tables whose marginals did not occur, but see Barnard [1945, 1949, 1979, 1989], Greenland [1991], Haber [1987], Storer and Kim [1990], and Suisa and Shuster [1984, 1985].

2. Referring to the literature of your own discipline, see if you can find a case where a  $2 \times 2$  table with at least one entry smaller than seven gave rise to a border line p-value using the traditional chi-square approximation. Reanalyze this table using Fisher's exact test.

Did the original authors use a one-tailed or a two-tailed test? Was their choice appropriate?

3. Again, refer to the literature of your own discipline for an example in which the chi-square approximation was used with an  $r \times 2$  table. Do you feel the chi-square statistic was appropriate? What statistic would you have used? Reanalyze the table using the statistic you have chosen.

4. If we were to question one respondent in the presence of another, would their answers be independent? If we were to make observations on several individuals in the same household, would these observations be independent? Criticize your own past work.

5. According to the Los Angeles Times, a recent report in the New England Journal of Medicine states that a group of patients with a severe bacterial infection of their blood stream who received a single intravenous does of a genetically altered antibody had a 30% death rate compared with a 49% death rate for a group of untreated patients. How large a sample size would you require using Fisher's exact test to show that such a percentage difference was statistically significant?

Before you start your calculations, determine whether you should be using a one-tailed or a two-tailed test.

6. Your friend rolls a die 120 times. Each time, before she rolls, you concentrate and visualize the die coming up a six. The die actually lands on six a total of 28 times. Are you psychic?

7. Will encouraging your child promote his or her intellectual development? A sample of 100 children and their mothers were observed and the children's IQs tested at 6 and 12 years. Results were as follows:

		Mothers Encourage Schoolwork		
		Rarely	Sometimes	Always
IQ increased	8	15	27	
	30	9	11	

- a) Do you plan to perform a one-tailed or two-tailed test?  
 b) Which test statistic is appropriate?  
 c) What is the significance level of your test?
8. Suppose you observed the following table:

10	90
20	90

Determine the p-value as many different ways as you can. Conduct a sensitivity test by determining the p-values for the table

10	91
20	90

For a discussion of your results, see Dupont [1986].

9. How would you go about obtaining a confidence interval for  $\pi_1 - \pi_2$ ?  $\pi_1/\pi_2$ ? See Santner and Snell [1980].
10. Referring to Table 6.8, if Sandoz excluded site 15 from their calculations could they safely combine the data from the remaining cites?
11. Holmes and Williams [1954] studied tonsil size in children to verify a possible association with the virus S. pyrogenes. Do you feel there is an association? How many rows and columns in the following contingency table? Which, if any, of the variables is ordered?

Tonsil Size by Whether Carrier of S. Pyrogenes			
	Not Enlarged	Enlarged	Greatly Enlarged
Noncarrier	497	560	269
Carrier	19	29	24

## CHAPTER 7

# Dependence

The title of this chapter, “dependence,” reflects our continuing emphasis on the alternative rather than on the null hypothesis. As you discover anew in this chapter, the permutation test is invaluable<sup>1</sup> whether you wish to focus on one or two specific hypotheses of dependence or provide protection against a broad spectrum of alternatives.

In this chapter, we consider five models of dependence and contrast the permutation approach to each with the bootstrap approach. You learn how to apply permutation tests, tracing a real-life regression problem from start to finish. And, of particular interest to economists, you learn methods for testing for first- and higher order correlations in stationary time series.

### 7.1. The Models

We consider five models of dependence in order of increasing complexity.

**Model 1 (Independence):** For all  $i$ , the pairs  $\{X_i, Y_i\}$  are independent and identically distributed with joint probability  $P$ , and  $P_X, P_Y$  are the corresponding marginal distributions. Having observed the pairs  $\{X_i, Y_i; i = 1, \dots, n\}$ , we wish to test the hypothesis that  $P$  is the product probability  $P_X * P_Y$ . Model 1 is the simplest of the five models, requiring the fewest assumptions about the data; its primary interest is theoretical rather than applied.

**Model 2 (Quadrant dependence):** When  $X$  is positive,  $Y$  is more likely to be positive, and vice versa. This model is appropriate when we have categorical or partially ordered data.

**Model 3 (Trend):**  $Y_i = G[X_i] + \zeta_i$  for  $i = 1, \dots, n$ ; where  $G$  is a monotone function of the (single) preset variable  $X$ , and the  $\{\zeta_i\}$ , the errors or residuals after the function  $G$  is used to predict  $Y$ , are exchangeable random variables

<sup>1</sup> And inevitable; see Bell and Donoghue [1969].

with zero expectations.  $G$  is a monotone increasing function of  $X$ , for example, if  $x_1 > x_2$  means that  $G[x_1] > G[x_2]$ . Having observed the pairs  $\{X_i, Y_i; i = 1, \dots, n\}$ , we wish to test the hypothesis that the distribution of  $Y$  is independent of  $X$  versus the alternative that  $Y$  is stochastically increasing in  $X$ . We have already encountered this model in Chapter 3, in testing for a dose response.

**Model 4 (Serial correlation):**  $Y_i = G[X_i] + \zeta_i$   $i = 1, \dots, n$ ; where  $G$  is a continuous function of the (single) preset variable  $X$  in the sense that if  $X_1$  is “close” to  $X_2$  then  $G[X_1]$  is “close” to  $G[X_2]$ , and the  $\zeta_i$  are independent random variables with expectation 0. Having observed the pairs  $\{X_i, Y_i; i = 1, \dots, n\}$ , we wish to test the hypothesis that the distribution of  $Y$  is independent of  $X$  versus the alternative that  $Y$  depends on  $X$  through some unknown  $G$ .

**Model 5 (Known model):**  $Y_i = G[X_i, \beta] + \zeta_i$   $i = 1, \dots, n$  where  $G$  is a known (arbitrary) function of  $X$  a vector of preset values,  $\beta$  is a vector of unknown parameters, and the  $\{\zeta_i\}$  are independent variables symmetrically distributed about 0. Having observed  $\{X_i, Y_i; i = 1, \dots, n\}$ , we wish to test the adequacy of some estimate  $\hat{\beta}$  of  $\beta$ , the true parameter value.

## 7.2. Testing for Independence

### 7.2.1. Independence

For Model 1,  $P$  is the product probability  $P_X * P_Y$ ; distribution-free bootstrap and randomization tests in the spirit of Kolmogorov–Smirnov test statistics are provided by Romano [1989]. Under the assumption that the pairs  $\{Y_i, X_i\}$  are independent and identically distributed, Romano finds that the bootstrap and the rerandomization test lead to almost the same confidence intervals for very large sample sizes.

Against parametric alternatives, the most powerful and/or locally most powerful tests are permutation tests based on the likelihood function [Bell and Doksum, 1967].

### 7.2.2. Quadrant Dependence

In Model 2, no ordinal relationship is implied;  $X$  and  $Y$  may even take categorical values, so that the problem reduces to that of analyzing a  $2 \times 2$  contingency table. The most powerful permutation test and, not incidentally, the most powerful unbiased test is Fisher’s exact test described in Section 6.2.

The bootstrap for the  $2 \times 2$  contingency table may be determined entirely on theoretical grounds without the need to resort to resampling. Estimates of the probabilities  $P\{Y > 0 | X > 0\}$ , and  $P\{Y > 0 | X < 0\}$  are used to obtain a confidence

interval for the odds ratio. If this interval contains unity, we accept the null hypothesis of independence, otherwise we reject it in favor of the alternative of quadrant dependence.

This model is occasionally used in practice while exploring the relationship between  $X$  and  $Y$ , first transforming to the deviations about the sample mean,  $X'_i = X_i - \bar{X}$ ,  $Y'_i = Y_i - \bar{Y}$ .

### 7.3. Testing for Trend

Consider an experiment in which you make two disparate observations on each of a series of experimental subjects, for example, observing the birth weight of an infant and its weight after one year, or the blood pressure and caffeine intake of each of a series of adults. You wish to test the hypothesis that the two variables vary independently against the alternative that there is a positive dependence between them.

More accurately, you wish to test the alternative of positive dependence against the null hypothesis of independence. In formal terms, if  $X$  and  $Y$  are the two variables, and  $Y_x$  is the random variable whose distribution is the conditional distribution of  $Y$  given that  $X = x$ , we want to test the null hypothesis that  $Y_x$  has the same distribution for all  $x$ , against the alternative that if  $x' > x$ , then  $Y_{x'}$  is likely to be larger than  $Y_x$ .

In Section 14.2, we show that Pitman's correlation  $\sum x_{(i)}y_i$ , where  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ , provides a most powerful unbiased test against alternatives with a bivariate normal density. As the sample size increases, the cutoff point for Pitman's test coincides with the cutoff point for the corresponding  $\alpha$  normal test based on the Pearson correlation coefficient.

Let's apply this test to the weight and cholesterol levels taken from a subset of the blood chemistry data collected by Werner et al [1970]; the full data set is included with the BMDP statistical package.

Wt	Chol
144	200
160	600
128	243
150	50
178	227
140	220
158	305
170	220

Is there a trend in cholesterol level by weight? Reordering the data by weight provides a clearer picture.

Wt	Chol
128	243
140	220
144	200
150	50
158	305
160	600
170	220
178	227

The cholesterol level does not appear to be related to weight; or, at least, it is not directly related. Again, we can confirm our intuition by the permutation test based on the statistic  $r$ .

But before we perform the test, what should we do about the subjects who had cholesterol values of 50 and 600? Are these typographical errors, or a misreading of the test results? Should we discard these values completely or perhaps replace them by ranks? Chapter 9 is devoted to a discussion of these and other alternatives for dealing with suspect data. In this chapter, we play the data as they lay. For the original data,  $r = 128 * 243 + \dots + 178 * 227 = 320,200$ , while  $r = 332,476$  for the following worst-case permutation:

Wt	Chol
128	50
140	200
144	220
150	220
158	227
160	243
170	305
178	600

Examining several more rearrangements, we easily confirm our eyeball intuition that cholesterol level is not directly related to weight. The majority of permutations of the data have sample correlations larger and more extreme than that of our original sample. We accept the null hypothesis.

## 7.4. Serial Correlation

For Model 4, advocates of the permutation test can take advantage of the (possible) local association between  $Y$  and  $X$ , reordering the  $X_i$  so that  $X_1 \leq \dots \leq X_n$ ,

and adopting as test statistic  $M = \sum_{i=1}^{n-1} (Y_i - Y_{i+1})^2$  [Wald and Wolfowitz, 1943; Noether, 1950; Maritz, 1981, p. 219]. We reject the null hypothesis if the value of the statistic  $M$  for the original observations is less than the  $\alpha$ th percentile of the permutation distribution of  $M$ . Again, we need not make specific assumptions about the nature of the association. If we can make specific assumptions, then some other permutation test may recommend itself. Ghosh [1954], for example, considers tests against the alternative of periodic fluctuations. Manly [1991] also considers a number of practical examples.

It is not clear what statistic, beyond that proposed by Romano for the simpler Model 1, might be used as the basis of a bootstrap test of Model 4. Of course, if we are prepared to specify the dependence function  $G$  explicitly, as is the case in Model 5, we may apply bootstrap methods to the residuals or to a suitable transformation thereof; see, for example, Stine [1987].

#### 7.4.1. An Example

To see a second illustration of the regression method (Model 3) while making a novel application of the present model, let us consider a second example, this time employing hypothetical data.

In Table 7.1,  $X$  represents the independent variable or cause, and  $Y$  represents the dependent variable or effect. Plotting  $Y$  versus  $X$  as in Figure 7.1 suggests a linear trend, and our permutation test for Model 3 confirms its presence. Our next step is to formulate a specific model and to estimate its parameters. The simplest model is a linear one  $Y = a + bX + \varepsilon$ . We can estimate the coefficients  $a$  and  $b$  using the method of least squares.

$$b = \frac{S_{xy}}{S_{xx}} = 4.5$$

$$a = \bar{y} - b\bar{x} = 3.53,$$

Table 7.1. Exploring a Cause–Effect Relationship

$X$	$Y$	$a + bX$	Residual	Rank	$a + bX + cX^2$	Residual
1	10.56	8.74	1.81	7	11.68	-1.12
2	15.28	14.15	1.12	6	14.57	.70
3	20.13	19.56	.56	5	18.31	1.82
4	22.26	24.98	-2.72	1	22.88	-0.62
5	28.06	30.38	-2.32	2	28.29	-0.23
6	33.61	35.80	-2.18	3	34.53	-0.93
7	41.13	41.20	-0.08	4	41.62	-0.49
8	50.41	46.62	3.79	8	49.55	0.86

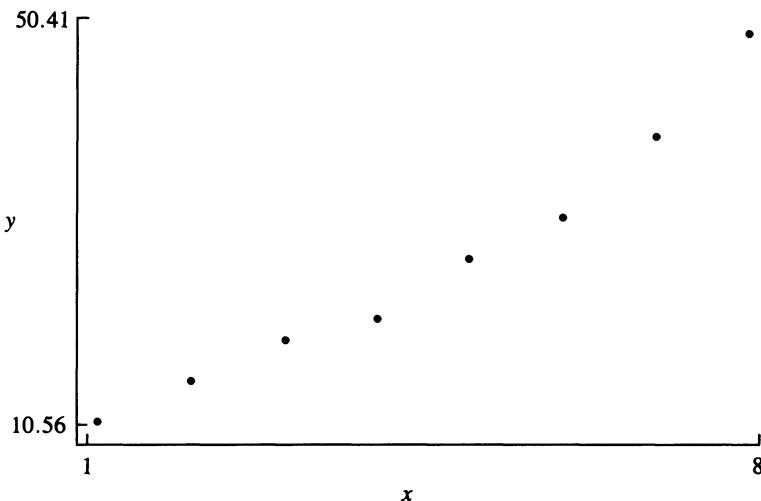


Figure 7.1. Plotting the effect  $Y$  against values of the cause  $X$ .

where

$$\begin{aligned} S_{xy} &= \sum (x_i - \bar{x})(y_i - \bar{y}) \\ S_{xx} &= \sum (x_i - \bar{x})^2. \end{aligned}$$

But a simple model is not always the right model. Suppose we compare our predictions from the linear model with the actual observations as in the second and third columns of our table. The fourth column of this table, which lists the differences between the predicted and observed values, attracts our interest. Is there some kind of trend here also? Examining a plot of the residuals in Figure 7.2, there does appear to be some sort of relationship between the residuals and our variable  $X$ . We can confirm the existence of this relationship by testing for serial correlation among the residuals. As a preliminary aid to the intuition, examine the ranks of the residuals in the fifth column of the table: 7 6 5 1 2 3 4 8. How likely is such a highly organized sequence to occur by chance? The value of  $M$  for the original residuals is 39.45; not one of 400 random rearrangements yields a value of  $M$  this extreme. The permutation test confirms the presence of a residual relationship not accounted for by our initial first-order model.

Let's try a second-order model:  $Y = a + bX + cX^2 + \varepsilon$ ; the least-squares coefficients are  $Y = 9.6 + 1.6X + 0.42X^2$ ; we've plotted the results in the final columns of Table 7.1; note the dramatic reduction in the size of the residuals; the second-order model provides a satisfactory fit to our data.

We could obtain bootstrap estimates of the joint distribution of  $X, Y$  by selecting random pairs, but with far less efficiency. If we are willing and justified in making additional assumptions about the nature of the trend function and the residuals as

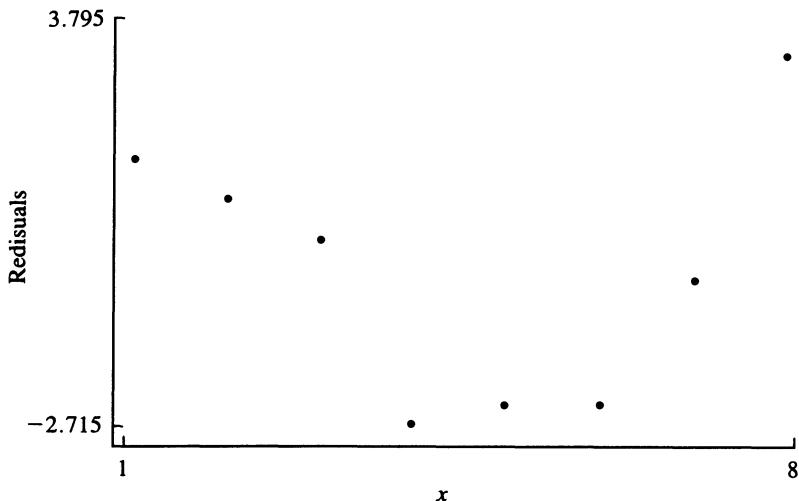


Figure 7.2. Plotting the residuals against values of the cause  $X$  after estimating and subtracting the linear trend.

in Model 5, then a number of more powerful bootstrap tests may be formulated. While we remain in an exploratory phase, our best choice of a test procedure appears to be Pitman's test followed by the Wald and Wolfowitz test for serial correlation among the residuals.

#### 7.4.2. Trend

We can test for a trend over time by using the Pitman correlation  $\sum t X(t)$ , where  $X(t)$  is the value of  $X$  at time  $t$  [Wald and Wolfowitz, 1943]. In the presence of a trend, the value of the test statistic should be more extreme than  $1 - \alpha$  of the values  $\sum \pi(t)X(t)$  where  $\pi$  is a permutation of the index set  $\{t_1, \dots, t_n\}$ .

We reach the same decision—accept or reject—whether we use the original values of the index set, for example, the dates 1985, 1986, 1987... to compute our test statistic or rezero them first as in 0, 1, 2, ... For  $\sum(t - C)X(t) = \sum t X(t) - C \sum X(t)$ , and the latter sum is invariant under permutations of the index set.

#### 7.4.3. First-Order Dependence

In a large number of economic and physical applications, we are willing to accept the existence of a first-order dependence in time (or space) of the form  $X(t + \tau) = f[\tau]X(t) + e_{t+\tau}$  but we would like to be sure that second- and higher order interactions are zero or close to zero. That is, if we are trying to predict  $X(t)$  and already know  $X(t - 1)$ , we would like to be sure that no further information is to be gained from a knowledge of  $X(t - 2)$ ,  $X(t - 3)$ , and so forth.

Gabriel [1962] generalizes the issue as follows: Define  $\{x(i)\}$  to be  $s$ th-degree antedependent if the partial correlation  $\rho$  of  $x(i)$  and  $x(i-s-z-1), \dots, x(1)$  for  $x(i-1), \dots, x(i-s-z)$  held constant is zero for all nonnegative  $z$ .

To test the hypothesis  $H_s$  that the  $\{x(i)\}$  are  $s$ th antedependent against the alternative  $H_{s+1}$ , accept  $H_s$  if

$$R = -N \sum_{i=1}^{p-s-1} \ln(1 - r_{i,i+s+1:i+1,\dots,i+s}^2)$$

is small, where the  $\{r_{i,i+s+1:i+1,\dots,i+s}^2\}$  denote the partial correlations derived from the original or permuted sample. We can assess  $R$  against its permutation distribution.

#### 7.4.4. An Economic Model

Are changes in stock market prices simply the result of a random walk? Or are there long-term serial correlations? More specifically, if a stock shoots up in price this year, must its price inevitably fall over the course of the next several years?

Kim et al [1991] studied the return data. They wanted to compare the financial return in the period  $[t, t + K]$  with that in the prior time interval  $[t - K, t]$ . They assumed that the return in the new period would depend in a linear fashion on the return in the old; in symbols,

$$r_{K,t+K} = \alpha_K + \beta_K r_{K,t} + \epsilon_{K,t+K},$$

where  $\alpha_K, \beta_K$  denote the intercept and slope of the line, and the  $\{\epsilon_{K,t+K}\}$  are a set of independent identically distributed errors.

Ordinary least squares were used to estimate the coefficient  $\beta_K$ ; then they randomized the original return series and obtained a new estimate of  $\beta_K$  for the permuted data. This randomization procedure was repeated 1000 times. If the  $\beta_K$  for the original unpermuted series was an extreme value, they were in a position to assert the presence of a correlation for that time period and that lag  $K$ . The entire procedure was repeated for several time periods and several lags (see question 5 at the end of this chapter).

## 7.5. Known Models

In Model 5, we may be given the vector of parameters or we may need to estimate it. We consider the case of known parameters first. Confidence intervals for the estimated parameters are covered in Section 7.5.3.

### 7.5.1. Testing a Specific Hypothesis

Just knowing there is a trend may not be enough. We may wish to test a specific hypothesis, for example, that  $\text{Chol} = C + 3Wt$ . Under the assumption of independent, identically symmetrically distributed residuals, we may form an unbiased permutation test of the hypothesis  $F = G_{\beta=3}$  by permuting the signs of the deviations  $d_i = \text{Chol}_i - C - 3Wt_i$  to obtain the distribution of the statistic

$$M = \Sigma_+ d_{\pi(i)},$$

where  $\Sigma_+$  ranges over the set for which  $d_{\pi(i)} \geq 0$ .

### 7.5.2. Testing a General Hypothesis

Suppose we have a linear regression model in mind, that is, we suspect that  $Y = \alpha + \beta f[X] + \varepsilon$ , where  $\varepsilon$  is a random variable with zero mean and  $f$  is a known monotone function in  $X$ , but we do not know the values of the parameters  $\alpha$  and  $\beta$ .

Without loss of generality, we assume that  $f[X] = X$  in what follows. To test the hypothesis that the slope is  $\beta_0$ , Oja [1981] proposes the test statistic,

$$\begin{aligned} T(X, \beta_0) &= \sum_{i < j} (e_j - e_i)(x_j - x_i) \\ &= \sum_{i < j} (y_j - \beta_0 x_j - y_i + \beta_0 x_i)(x_j - x_i) \\ &= \sum_{i < j} (y_j - y_i)(x_j - x_i) - \beta_0 \sum_{i < j} (x_j - x_i)(x_j - x_i), \end{aligned}$$

whose permutation distribution may be obtained from

$$\begin{aligned} T(\pi(X), \beta_0) &= \sum_{i < j} (y_j - y_i)(\pi[x_j] - \pi[x_i]) \\ &\quad - \beta_0 \sum_{i < j} (x_j - x_i)(\pi[x_j] - \pi[x_i]). \end{aligned}$$

### 7.5.3. Confidence Intervals

In most cases, it is not enough to know that  $Y$  is dependent on  $X$ , we want to know the specific nature of this dependence. As an example, suppose we have satisfied ourselves that Model 3 (Trend) holds—that is,  $Y_i = G[X_i] + \zeta_i$  for  $i = 1, \dots, n$ , where  $G$  is a monotone function of the (single) preset variable  $X$ ,  $G = a + bX$ , say, and the  $\zeta_i$  are exchangeable random variables with zero expectations. Having decided that  $b \neq 0$ , we would like to obtain a confidence interval for  $b$ . First note

that a permutation test of the hypothesis  $H_0 : b = b_0$  may be based on the Pitman correlation  $\sum x_i \zeta_i^0$ , where  $\zeta_i^0 = y_i - y. - x_i b_0$ ,  $i = 1, \dots, n$  are the deviations about the line whose slope is  $b_0$ .

Let  $\pi$  denote a permutation of the subscripts  $1, \dots, n$  and put  $b^\pi[w] = \sum x_i w_{\pi[i]} / \sum x_i^2$ . For example,

$$b^\pi[\zeta^0] = \sum x_i \zeta_{\pi[i]}^0 / \sum x_i^2.$$

We reject or accept  $H_0$  according to whether  $b'[\zeta^0]$  for the original, unpermuted deviations is an extreme value of the distribution of  $b^\pi[\zeta^0]$ .

We can obtain a confidence interval for  $b$  by following the trial and error procedure described in Section 3.2, but there is a better way, due to Robinson [1987].

The least-squares estimate of  $b$  is  $\hat{b} = \sum x_i y_i / \sum x_i^2$ , so that  $b'[\zeta^0] = \hat{b} - b_0$  for the original, unpermuted deviations. Let  $\hat{\zeta}_i = y_i - y. - x_i \hat{b}$ ;

$$\begin{aligned} \zeta_i^0 &= \hat{\zeta}_i + (\hat{b} - b_0)x_i; \\ b^\pi[\zeta^0] &= b^\pi[\hat{\zeta}] + (\hat{b} - b_0)b^\pi[x]; \\ \{b_0 : b^\pi[\zeta^0] \geq \hat{b} - b_0\} &= \{b_0 : b^\pi[\hat{\zeta}] \geq (\hat{b} - b_0)(1 - b^\pi[x])\} \\ &= \{b_0 : b^\pi[\hat{\zeta}] / (1 - b^\pi[x]) \geq \hat{b} - b_0\}. \end{aligned}$$

$b^\pi[\hat{\zeta}] / (1 - b^\pi[x])$  is a pivotal quantity that does not depend on  $b_0$ . The desired confidence region is the interval between the  $k$ th and the  $(n! - k + 1)$ th-order statistics of this pivotal quantity where  $k = (n!\alpha/2)$ .

## 7.6. Multiple Regression

### 7.6.1. Eliminating Covariate Effects

We easily can modify the definition of Oja's test to encompass multivariate regression. Note that the statistic in Section 7.5.2 may be rewritten as

$$T_1 = \sum_{i < j} \Delta_{ij}^y \Delta_{ij}^{\pi(x)},$$

where

$$\Delta_{ij}^y = y_j - y_i = \begin{vmatrix} 1 & 1 \\ y_i & y_j \end{vmatrix}$$

and similarly for  $\Delta_{ij}^{\pi(x)}$ .

Now, suppose,  $y = \alpha + \beta_{x_1|x_2}x_1 + \beta_{x_2|x_1}x_2 + \varepsilon$  and we wish to test the hypothesis that  $\beta_{x_2|x_1} = 0$  after correcting for the effect of the covariate  $x_1$ . Our test statistic is  $\sum_{i < j} \Delta_{ijk}^y \Delta_{ijk}^{\pi(x_2)}$ , where

$$\Delta_{ijk}^y = \begin{vmatrix} 1 & 1 & 1 \\ y_i & y_j & y_k \\ x_{1i} & x_{1j} & x_{1k} \end{vmatrix}.$$

Our permutations are obtained by holding the couples  $(y_i, x_{1i})$  fixed and randomizing with respect to  $x_2$ . If we knew what the relationship was between Y and  $X_2 (X_3, \dots)$  exactly, that is, if we knew what  $\beta_{x_2|x_1}$  was and so needn't estimate it, we could obtain an exact test of the partial regression coefficient  $\beta_{x_1|x_2}$  of Y on  $X_1$  with ease. But because we must estimate all coefficients, the ones we are not interested in as well as the ones we are, these permutations are not exchangeable, and the resulting test is not exact (though it is asymptotically exact). Nor does this method of shuffling hold constant the collinearity between  $X_1$  and  $X_2$ , a violation of the ancillarity principle (see Welch, 1990). The resulting slope coefficient estimates will not vary as much in repeated samples as the original estimate, causing the Type I error to be larger than the declared value. Kennedy and Cade [1996] propose the use of an asymptotically pivotal test statistic (such as the F statistic one might use to test that the partial regression coefficient is zero) to minimize the discrepancy.

Three other methods for determining the partial regression coefficients by permutation means have been considered. Though none of the three is exact, they are all asymptotically exact and, in one instance, close to exact even for small samples under a wide variety of conditions. As described by Anderson and Legendre [1998], these methods are

- 1) permuting the raw data;
- 2) permuting the residuals under the full model;
- 3) permuting the residuals under the reduced model obtained by first correcting for the covariates.

#### 7.6.1.1. Permute the Raw Data

The following procedure is employed as suggested by Manly [1991] and modified by Anderson and Legendre [1998] to use a pivotal statistic.

- 1) Regress Y on all independent variables simultaneously to obtain an estimate  $b$  of  $\beta_{x_1|x_2\dots}$  and a value  $t_0$  of the t-statistic for testing  $\beta_{x_1|x_2\dots} = 0$  for the real data.
- 2) Permute the  $\{Y_i\}$  to obtain a permuted set of values  $\{Y_i^*\}$ .
- 3) Obtain and estimate  $b^*$  and a t-statistic  $t^*$  for the permuted values. Repeat Steps 2 and 3 until you obtain a permutation distribution of t against which  $t_0$  may be assessed.

### 7.6.1.2. Permute Residuals Under the Full Model

This approach, due to ter Braak [1992], an analog of a bootstrap method proposed by Hall and Titterington [1989], proceeds as follows.

- 1) Regress  $Y$  on all independent variables simultaneously to obtain estimates of all regression coefficients and a value  $t_0$  of the t-statistic for testing  $\beta_{x_1|x_2\dots} = 0$  for the real data.
- 2) Permute the residuals randomly.
- 3) Calculate new values  $\{Y^*\}$  using the permuted residuals and the coefficients obtained at Step 1. (That is,  $Y^* = b_0 + b_{1,2,\dots}X_1 + b_{2,1,\dots}X_2 + \pi[\varepsilon']$ .)
- 4) Regress  $Y^*$  on all independent variables simultaneously to obtain new estimates of all regression coefficients and a value  $t^*$  of the t-statistic for testing  $\beta_{x_1|x_2\dots} = 0$  for the permuted data.

Repeat Steps 2 through 4 until you obtain a permutation distribution of  $t$  against which  $t_0$  may be assessed.

### 7.6.1.3. Permute Residuals Under the Reduced Model

Freedman and Lane [1983] propose the following permutational procedure.

- 1) Regress  $Y$  on all independent variables simultaneously to obtain estimates of all regression coefficients and a value  $t_0$  of the t-statistic for testing  $\beta_{x_1|x_2\dots} = 0$  for the real data.
- 2) Regress  $Y$  on  $X_2$  alone after correcting  $Y$  for the remaining covariates using the regression coefficients obtained in Step 1. (That is,  $Y^{cor} = b_0 + b_1X + \varepsilon'$ .)
- 3) Permute the residuals from the regression in Step 2.
- 4) Obtain new values  $\{Y^*\}$  by adding the permuted values of the residuals to the fitted values obtained in Step 2. (That is,  $Y^* = b_0 + b_2X_2 + \pi[\varepsilon']$ .)
- 5) Regress  $Y^*$  on all independent variables simultaneously to obtain new estimates of all regression coefficients and a value  $t^*$  of the t-statistic for testing  $\beta_{x_1|x_2\dots} = 0$  for the permuted data.

Repeat Steps 2 through 5 until you obtain a permutation distribution of  $t$  against which  $t_0$  may be assessed.

While this latter method requires the greatest number of computations, Anderson and Legendre [1998] found that with it one can obtain almost exact results, even though 1) the data are radically nonnormal, 2) there is at least one very large outlier, and 3) the independent variables are highly colinear.<sup>2</sup>

<sup>2</sup> A variant of the reduced model approach due to Kennedy [1995] proved the worst in their tests.

### 7.6.2. LAD or LSD?

A careful reader of the preceding sections may have noticed that we did not specify *how* the regression coefficients are to be estimated. As the minimization functions used in least absolute deviation regression are pivotal quantities for regression quantiles including the median (Parzen, Wei, and Ying, 1994), the preceding results apply equally to LAD or LSD regression (Cade and Richards, 1996). We need not assume the errors  $\{\epsilon_i\}$  are identically symmetrically distributed about zero, necessary assumptions for ordinary least squares, only that they are identically distributed with median zero. We can choose our regression method based on the loss function (see Section 14.1.1), not because one method lends itself more readily to an analytic solution than another.

### 7.6.3. An Exact Solution

Brown and Maritz [1982] propose that we stratify or block our observations on  $Y$  and  $X_2$  by the values of  $X_1$ . (If  $X_1$  is continuous, then we need to group its values into a small number of ordered categories.) Within each block, we compute the test statistic,

$$T(j, \gamma) = \sum_{k=1}^{s_j} z_{jk}(y_{jk} - \hat{\beta}_j x_{1j} - \gamma z_{jk}),$$

where  $\hat{\beta}_j = \hat{\beta}_j(\gamma)$  is the solution for  $\beta$  of

$$0 = \sum_{k=1}^{s_j} (y_{jk} - \hat{\beta}_j x_{1j} - \gamma z_{jk})$$

and  $z_{j1}, z_{j2}, \dots, z_{js_j}$  are the values taken by  $X_2$  when  $X_1 = x_{1j}$ . Our overall test statistic is  $T(., \gamma) = \sum T(j, \gamma)$ , and to obtain its distribution, we permute *on a block-by-block basis*. The resulting statistic is exact, though with a possible loss of power if we need to group the values of  $X_1$  to form a reduced number of discrete levels.

### 7.6.4. Testing All Coefficients

Alternatively, we may want to test both the partial independence of  $(y, x_1 | x_2)$  and  $(y, x_2 | x_1)$ . Then, following Manly [1997], we would hold the couples  $(x_{1i}, x_{2i})$  fixed and randomize with respect to  $y$  using test statistics of similar form.

These results extend to any number of explanatory variables  $x_1, x_2, \dots, x_n$ .

## 7.7. Single-Case Phase Designs

In single-case designs, the power depends not only on the effect size and number of observations, but also on both the method of randomization and the correlation between successive observations.

Consider an AB design consisting of 30 observations. The initiation of treatment B could be randomly assigned to any of the sixth through 25th observations for a total of 20 possible assignments, or, the experiment could be divided into several phases of fixed length and the treatments, not the points of initiation, randomly assigned. A 30-observation design consisting of six phases of five observations each would also lead to 20 possible assignments.

Successive observations might be uncorrelated, but more often will obey a first-order autoregressive model,  $y_t = \phi y_{t-1} + e_t$ , where  $y_t$  is the deviation of the observation from the mean at time  $t$ ,  $\phi$  is the autoregressive parameter, and the  $\{e_t\}$  are independent errors (Ferron and Ware, 1995; Huitema and McKean, 1991).

Ferron and Onghena [1996] studied the power of single-case designs for both methods of random assignment and for both correlated and uncorrelated observations. Designs based on the random assignment of treatments to phases were substantially more powerful than those based on the random assignment of intervention points. Positive autocorrelation led to greater power with random assignment of treatments and to decreased power with random assignment of intervention points. Negative correlation had the opposite effect in both instances.

## 7.8. Questions

1. a) Are the bootstrap and permutation tests against quadrant dependence equivalent for very large samples?
- b) Suppose you observed the contingency table

	Republican	Democrat
White	8	3
Black	3	8

Is race associated with political preference? Use both the bootstrap and Fisher's Exact test (Section 6.2) to make the inference.

2. In your own area of specialization, there is undoubtedly a controversy about the nature of the association between some pair of variables. Which of the models, 1, 2, ..., 5 would be most appropriate for describing this association?
3. Adding platinum to a metallic coating will increase the mean time between failures. But is it worth it? This will depend on the cost of platinum, the magnitude of the effect, and the cost of a failure. Using the data in the following table and the prediction equation

$MTBF = a + b(PT)$ , obtain a confidence interval for the effect  $b$ . Use both the trial and error method (Section 3.2) and the pivotal quantity developed in Section 7.4.

Table 7.2. Effect of Platinum on MTBF

Grams Platinum per KG	MTBF (Hours)
1	900
2	1000
5	1100
10	1300
15	1600
20	1800

4. a) Table 7.3 records monthly sales for a two-year period, taken from Makridakis, Wheelwright, and McGee [1983]. Is there a seasonal trend?

Table 7.3. Monthly Sales as a Function of  $X$

$t$	$X$	Sales	$t$	$X$	Sales
0	116.44	202.66	12	129.85	260.51
1	119.58	232.91	13	122.65	266.34
2	125.74	272.07	14	121.64	281.24
3	124.55	290.97	15	127.24	286.19
4	122.35	299.09	16	132.35	271.97
5	120.44	296.95	17	130.86	265.01
6	123.24	279.49	18	122.90	274.44
7	127.99	255.75	19	117.15	291.81
8	121.19	242.78	20	109.47	290.91
9	118.00	255.34	21	114.34	264.95
10	121.81	271.58	22	123.72	228.40
11	126.54	268.27	23	130.33	209.33

- b) After eliminating the seasonal trend from the sales data in Table 7.3, is there a significant upward trend in the remaining averages? Your test statistic is what sum?
- c) The “ $X$ ” of Table 7.3 is actually advertising expenditures. Can a knowledge of your advertising expenditures explain part of the trend in sales? What statistic would you use to determine if sales do depend on advertising  $X$ .
- d) Should you test this multivariate regression before eliminating the seasonal trend? Would the sales in month  $i$  depend on the advertising expenditures in month  $i$ ? or the previous month  $i - 1$ ? or on those in several previous months? What statistics would you use to resolve these issues?
5. In the work described in Section 7.4.4, Kim et al [1991] studied three out of four time periods, singly and in combination, and some 20 different lags. They note: “If we have no prior basis for choosing a particular lag we may be overstating significance by focusing on the lag with the lowest p value.” Keeping in mind what you read in Section 12.3, what is your view?

6. Is performance on the LSAT's related to undergraduate GPA? Use the data in the following table to make your determination.

LSAT	GPA	LSAT	GPA
545	2.76	594	2.96
555	3.00	605	3.13
558	2.81	635	3.30
572	2.88	651	3.36
575	2.74	653	3.12
576	3.39	661	3.43
578	3.03	666	3.44
580	3.07		

7. Show that Oja's test 7.5.2 is or is not exact. Show that it is or is not asymptotically exact. Similarly, show that Manley's test 7.6.4 is or is not exact (asymptotically exact).
8. A pivotal statistic has a sampling distribution that is independent of all parameters. Show that if a pivotal statistic is not employed in multiple regression that the degree of collinearity in the original data will affect the outcome of the test (see Kennedy and Cade, 1996).

## CHAPTER 8

# Clustering in Time and Space

In this chapter, you learn how to detect clustering in time and space and to validate clustering models. We use the generalized quadratic form in its several guises including Mantel's  $U$  and Mielke's multiresponse permutation procedure to work through a series of applications in atmospheric science, epidemiology, ecology, and archeology.

## 8.1. The Generalized Quadratic Form

### 8.1.1. Mantel's $U$

Mantel's  $U$  [Mantel, 1967]  $\sum \sum a_{ij} b_{ij}$  is perhaps the most widely used of all multivariate statistics. In Mantel's original formulation,  $a_{ij}$  is a measure of the temporal distance between items  $i$  and  $j$ , while  $b_{ij}$  is a measure of the spatial distance. As an example, suppose the pair  $(t_i, l_i)$  represents the day  $t_i$  on which the  $i$ th individual in a study came down with cholera and  $l_i = (l_{i1}, l_{i2})$  denotes her position in space. For all  $i, j$  set  $a_{ij} = 1/(t_i - t_j)$  and

$$b_{ij} = 1/\sqrt{(l_{i1} - l_{j1})^2 + (l_{i2} - l_{j2})^2}.$$

A large value for  $U$  would support the view that cholera spreads by contagion from one household to the next. How large is large? As always, we compare the value of  $U$  for the original data with the values obtained when we fix the  $i$ 's but permute the  $j$ 's as in  $U' = \sum \sum a_{ij} b_{i\pi(j)}$ .

The generalized quadratic form has seen widespread application in anthropology, archaeology [Klauber, 1971, 1975], ecology [Bryant, 1977; Douglas and Endler, 1982; Highton, 1977; Levin, 1977; Royaltey, Astrachen, and Sokal, 1975; Ryman et al, 1980], education [Schultz and Hubert, 1976], epidemiology [Alderson and Nayak, 1971; Fraumeni and Li, 1969; Glass and Mantel, 1969; Klauber and Mustacchi 1970; Kryscio et al, 1973; Mantel and Bailar, 1970; Merrington and Spicer, 1969; Siemiatycki and McDonald, 1972; Smith and Pike, 1976; geography

**Table 8.1. Incidents of Pairs of Anencephalic Infants by Distance and Time Months Apart**

km apart	<1	<2	<4
<1	39	101	235
<5	53	156	364
<25	211	652	1516

[Cliff and Ord, 1971, 1973, 1981; Hubert, 1978b; Hubert, Golledge, and Costanzo, 1981; Hubert et al, 1984, 1985], management science [Graves and Whinston, 1970], psychology [Hubert and Schultz 1976; Hubert, 1978a, 1979], sociology [Hubert and Baker, 1978], and systematics [Dietz, 1983; Gabriel and Sokal, 1969; Jones, Selander, and Schnell, 1980; Selander and Kaufman, 1975; Sokal, 1979].

### 8.1.2. An Example

An ongoing fear among many parents is that something in their environment — asbestos or radon in the walls of their house, or toxic chemicals in their air and ground water, will affect their offspring. Table 8.1 is extracted from data collected by Siemiatycki and McDonald [1972] on congenital neural-tube defects. Eyeballing the gradient along the diagonal of this table, one might infer that births of anencephalic infants occur in clusters. One could test this hypothesis statistically using the methods of Chapter 6 for ordered categories, but a better approach, since the exact time and location of each event is known, is to use Mantel's  $U$ . The question arises as to which measures of distance and time we should employ. Mantel [1967] reports striking differences between one analysis of epidemiologic data in which the coefficients are proportional to the differences in position and a second approach (which he recommends) to the same data in which the coefficients are proportional to the reciprocals of these differences.<sup>1</sup> Using Mantel's approach, a pair of infants born 5 km and three months apart contribute  $\frac{1}{5} * \frac{1}{3} = 1/15$  to the correlation. Summing the contributions from all pairs, then repeating the summing process for a series of random rearrangements, Siemiatycki and McDonald conclude that the clustering of anencephalic infants is not statistically significant.

## 8.2. Applications

By appropriately restricting the values of  $a_{ij}$  and  $b_{ij}$ , the definition of Mantel's  $U$  can be seen to include several of the standard measures of correlation including

<sup>1</sup> One further caveat: Mantel's  $U$  fails completely if the spatial distribution of the underlying population is also changing with time [Roberson and Fisher, 1986].

those usually attributed to Pearson, Pitman, Kendall, and Spearman [Hubert, 1985]. Mantel's  $U$  has been rediscovered frequently, often without proper attribution; see Whaley [1983]. In this section we consider three diverse approaches to the problem of assessing the presence of clustering in space and time. In each case, the permutation distribution of the quadratic form is used to provide a baseline against which the behavior of the observations may be assessed.

### 8.2.1. The MRPP Statistic

One such variant is the MRPP or multiresponse permutation procedure [Mielke, 1979], which is used in applications as diverse as the weather and the spatial distribution of archaeological artifacts. The MRPP uses the permutation distribution of between-object distances to determine whether a classification structure has a nonrandom distribution in space or time. With large samples, a Pearson Type III curve based on the first three (or four) exact moments may be used in place of the permutation distribution [Mielke, Berry, and Brier, 1981].

An example of the application of the MRPP arises in the assignment of antiquities (artifacts) to specific classes based on their spatial locations in an archaeological dig. Presumably, the kitchen tools of primitive man—woks and Cuisinarts—should be found together, just as a future archaeologist can expect to find TV, VCR, and stereo side by side in a neolithic living room.

Following Berry et al [1980, 1983], let  $\Omega = \{\omega_1, \dots, \omega_N\}$  designate a collection of  $N$  artifacts within a site; let  $X_{1i}, \dots, X_{ri}$  denote the  $r$  coordinates for the site space for artifact  $\omega_i$ ; let  $S_1, \dots, S_{g+1}$  represent an exhaustive partitioning of the  $N$  artifacts into  $g + 1$  disjoint classes (the  $g + 1$ st being reserved for not-yet-classified items); and let  $n_j$  be the number of artifacts in the  $j$ th class.

Define the Euclidian distance between two artifacts,

$$\delta_{ij} = \left[ \sum_{k=1}^r (X_{ki} - X_{kj})^2 \right]^{1/2}.$$

Define the average between-artifact distance for all artifacts within the  $i$ th class,

$$\xi_i = \frac{2}{n_i(n_i - 1)} \sum_{i < j} \delta_{ij} \phi_i(\omega_i) \phi_i(\omega_j),$$

where  $\phi_i(\omega)$  is an indicator function that is 1 if  $\omega \in S_i$  and 0 otherwise.

The test statistic is the weighted within-class average of these distances,

$$\Delta = \sum_{i=1}^g n_i \xi_i / K,$$

where  $K = \sum_{i=1}^g n_i$ .

The permutation distribution associated with  $\Delta$  is taken over all  $\frac{N!}{\prod_{i=1}^{g+1} n_i!}$  allocations of the  $N$  artifacts to the  $g + 1$  classes with the same numbers  $\{n_i\}$  assigned to each class.

Empirical power comparisons between MRPP rank tests and with other rank tests are made by Tracy and Tajuddin [1985] and Tracy and Khan [1990].

### 8.2.2. The BW Statistic of Cliff and Ord [1973]

As a second application of generalized correlation, suppose we want to measure the degree to which the presence of some factor in an area (or time period) increases the chances that this factor will be found in a nearby area.

The BW statistic of Cliff and Ord [1973] is defined as  $\sum \sum \delta_{ij}(x_i - x_j)^2$ , where

$$\begin{aligned} &= 1 && \text{if the } i\text{th area has the characteristic} \\ x_i &= 0 && \text{otherwise} \\ &= 1 && \text{if the } i\text{th and } j\text{th areas are adjacent} \\ \delta_{ij} &= 0 && \text{otherwise.} \end{aligned}$$

### 8.2.3. Equivalances

The generalized quadratic form has been rediscovered and redefined in many different guises. Whaley [1983] shows that Mantel's  $U$  and the BW statistic are equivalent to the MRPP for testing purposes. A third equivalent example is the  $k$ -dimensional runs test of Friedman and Rafsky [1979] studied in Section 5.3.

### 8.2.4. Reduce

Mantel's  $U$  is quite general in its application. The sets of coefficients  $\{a_{ij}\}$  and  $\{b_{ij}\}$  need not represent positions or changes in time and space.

In a completely disparate application in sociology, Hubert and Schultz [1976], observers studied  $k$  distinct variables in each of a large number of subjects. Their object was to test a specific sociological model for the relationships among the variables. This time, the  $\{a_{ij}\}$  in Mantel's  $U$  are elements of the  $k \times k$  sample correlation matrix while the  $\{b_{ij}\}$  are elements of an idealized or theoretical correlation matrix derived from the model. A large value of  $U$  supports the model, a small value rules against it.

### 8.2.5. Another Dimension

Vecchia and Iyer [1989] generalized the MRPP for use in the comparison of several linear models. In the words of these authors, "Regarding algebraic quantities useful

to detect concentrations of points within distinct groups, one might have asked: *When are two points concurrent?*. The answer, that they coincide whenever the *distance between them is zero* motivates the definition of the MRPP statistic in terms of interpoint distance.

“Extending this approach, for example, to the *comparison of straight line relations*, the analogous geometric argument is that three points are colinear only if their triangular *area is zero*.”

The statistic used in Vecchia and Iyer’s new test is a symmetric volume: A real-valued function, symmetric in its  $n + 1$  arguments, that is zero if and only if the Euclidean volume of the simplex formed by the arguments is zero. An immediate application for this statistic is assessing the consistency of multiclinic designs. Some of this statistic’s asymptotic properties are considered in Vecchia and Iyer [1991].

## 8.3. Alternate Approaches

### 8.3.1. Quadrant Density

Following Mead [1974], we overlay an (possibly irregular) area with a grid and divide it into squares. We then group the squares into K regions so there is an equal number of squares in each region. Finally, we count the number  $n_i$  of events (nests, animals) observed within each region and form the test statistic  $S = \sum n_i^2$ .

As we are working with squares, S takes its largest value if the data are clustered by region.

Permute the squares among regions and compute S each time. Accept the alternative that there is clustering if only a small percentage of the permutations yields values of S that are as large as  $S_o$  for the original arrangement of squares.

Suppose we have only eight squares, which we group into regions corresponding to the counts 0, 1, 2, 0 and 3, 4, 5, 2. Clustering is evident.  $S_o = 9 + 196 = 205$ . We rearrange the squares so that the counts within each region are 0, 1, 4, 5 and 0, 2, 2, 3.  $S = 100 + 49 = 149$ . Continuing in this fashion, we see that  $S_o$  is an extreme value and we reject the null hypothesis that the counts are distributed uniformly.

### 8.3.2. Nearest-Neighbor Analysis

Following Ripley [1981], let  $\{p_j\}$  be a set of points in a region where specific events have been observed (cases of leukemia, birds nests, and so forth). Let  $q_i(p_j)$  denote the distance from the point  $p_j$  to its  $i$ th nearest neighbor, and let  $q_i$  denote the mean of these distances. Now, overlay the area with a grid so as to divide it into squares. Permute the squares; determine  $q_i$  for the permutation, and compare it with  $q_{io}$  for the original observations. The question remains as to which of the  $q_i$  to use for testing purposes.

### 8.3.3. Comparing Two Spatial Distributions

Upton [1984] objects to Mead's procedure observing that the result depends strongly on how the regions are defined, particularly for irregular areas and when there are missing data. Syrjala's [1993] use of a cumulative distribution based on the work of Zimmerman [1993] overcomes this objection and, moreover, allows us to extend the procedure to compare two distinct distributions.

Again, we overlay the region with a grid and divide it into squares whose centers are at the points  $(x_k, y_k)$ . Define the density  $d_i(x_k, y_k) = n_i/N$ , where  $N = \sum n_i$ ; In order to compare two populations, normalize the density so that  $\gamma_i(x_k, y_k) = \frac{d_i(x_k, y_k)}{\sum d_j(x_k, y_k)}$  and define

$$F_i[x_k, y_k] = \sum_{S_k} d_i(x, y),$$

where the sum  $\sum_{S_k}$  is taken over the region  $S_k = \{x \leq x_k; y \leq y_k\}$ . Our test statistic  $\Gamma = \sum_k (F_1 - F_2)^2$ , and to obtain its permutation distribution, we evaluate all  $2^K$  permutations of the two species at the  $K$  points of the grid  $\{x_k, y_k\}$ .<sup>1</sup>

## 8.4. Questions

1. Show that Pitman's correlation is a special case of Mantel's U.
2. List at least two applications for Vecchia and Iyer's test.

<sup>1</sup> Since the value of this statistic also depends on the location of the origin, we may define  $F_{j1}, F_{j2}, \Gamma_j$  for  $j = 1, \dots, 4$ , corresponding to the placing of the origin at each of the four corners of a (nearly) rectangular region. Our test statistic then would be the average of the four values,  $\Gamma = \sum \Gamma_i / 4$ .

## CHAPTER 9

# Coping with Disaster

In this chapter, you receive practical guidelines for coping with the many catastrophes that confront the applied statistician:

- \* subjects who miss an appointment,
- \* subjects who disappear completely and mysteriously in the middle of an experiment,
- \* incomplete questionnaires,
- \* covariates after the fact,
- \* outlying observations whose extreme and questionable values suggest they may have been recorded incorrectly,
- \* off-scale and other censored values that cannot be determined with precision,
- \* and even studies that must be brought to a rapid and untimely conclusion well in advance of the scheduled date.

### 9.1. Missing Data

The effects of missing data depend upon the nature of the study. In some instances, for example, in the one-factor,  $k$ -sample comparison, missing data have no effect upon the analysis other than to reduce the power of the test. In other, more complex designs, missing data may result in an unbalanced design in which several factors are confounded with one another. In most, though not all, of these latter cases, no special statistical procedures are required, *providing* we are careful in how we interpret the results. We must identify which effects are confounded with one another, a main effect with an interaction, say. In other studies (and one such example was examined in Section 4.4.2), we may have to abandon permutation procedures altogether and consider using the bootstrap.

The majority of experimental designs belong to the correctable category. We proceed with the permutation analysis using a revised set of marginal constraints that reflect the actual rather than the hoped-for sample sizes. And in analyzing

the results, we acknowledge that one or more higher order interactions may have contaminated the observed effects.

Consider an example we studied in Section 4.2, the effect of sunlight and fertilizer on crop yield. Suppose that one of the observations in the lowsunlight, medium-fertilizer group, the 22 noted in parenthesis in the table below, is missing from the study.

Effect of Sunlight and Fertilizer on  
Crop Yield

		Fertilizer		
		LO	MED	HIGH
LO	LO	5	15	21
	HI	10	(22)	29
	HI	8	18	25
HI	LO	6	25	55
	LO	9	32	60
	LO	12	40	48

The test statistic for the main effect of sunlight  $S = 23 + (15 + 18) + 75 = 131$  for these observations. Such an extremely low value is found in only a small handful of the rearrangements in which we swap observations at random between the low and high groups. The number of rearrangements after correcting for the missing data item is  $\binom{17}{8}$ . The reduction from the hoped for  $\binom{18}{9}$  rearrangements reduced the power of the test. But the reduction is irrelevant in this instance as we are rejecting the hypothesis. Had we accepted the null hypothesis, we would have been forced to consider whether a larger sample size might have enabled us to detect an effect.

A missing data item in only one of the groups means that the main effect of sunlight is partially confounded with the interaction between sunlight and fertilizer. But our common sense strengthened by a glance at the table tells us that the confounding also is irrelevant in this instance.

The preceding discussion was based on the implicit assumption that dropouts occur at random. If the dropout rate is directly related to the treatment, we must either abandon the study or modify our scoring system explicitly to account for the dropout. See, for example, Entsuah [1990].

A further example of using the permutation distribution to cope with missing data is given in Section 10.2.6.

## 9.2. Covariates After the Fact

After World War II, public policymakers in the United States did a slow about-face on the dangers of tobacco smoke. The changes in policy accelerated during

the 1970's. One moment it seemed the cigarette was the ultimate symbol of masculinity, and the next it was the primary cause of emphysema, hypertension, lung cancer, and fetal defects. One month you could design a 400-patient, six-week, 50-variable clinical study with the full support of a Food and Drug Administration panel, and the next the panel would be asking if you'd corrected for the smokers in the control group. Of course you hadn't, not then, not in those days.

Today, we know that smoking is harmful, but "cigarettes smoked per week" is only one of hundreds of possible covariates. Regardless of how many covariates you have controlled or matched in putting together a clinical study, there are sure to be one or two more covariates that you didn't think of, that no one thought of, that no one could have envisioned—that is, until the day after your 300-page report on the study was sent to the printers.

All is not lost, it is still possible to make a comparison among treatment groups using the method of permutations by restricting the rerandomizations to those with specific after-the-fact design matrices.

Using the method due to Rosenbaum [1984], described at length in Section 4.3, we block the data into smokers and nonsmokers (or lemon eaters and non-lemon-eaters), and then randomize separately within each block.

Restricting the number of randomizations may reduce the power of the test. (It may also increase it by eliminating a source of variability; see Section 3.6.) As a result, we may need to add more subjects and an additional clinical center to the study to justify and confirm any negative results.

### 9.2.1. Observational Studies

An extreme example of the use of an after-the-fact covariate comes when we attempt to create matched pairs from two groups that were part of an observational study. In an observational study, the groupings themselves are after the fact. The subjects are not randomly assigned to treatment or control but are merely "observed" to belong to one group or the other. Through the use of after-the-fact covariates, we hope to reduce or eliminate any built-in biases.

An example provided by Rosenbaum [1988] is that of a study in humans of the effect of vasectomy on the risk of myocardial infarction. Obviously, we do not have the luxury (nor the authority, thankfully) to select a random sample of patients for a mandatory vasectomy, but must analyze the data as they lie. We can take advantage of concurrent data on obesity and smoking history (both of which are known to affect the risk of myocardial infarction) to help us block the two samples so as to reduce the between-sample variance. See Rosenbaum [1988] for methods for dealing with imperfect matching.

While no justification for the use of restricted randomization is required when the covariates are built in to the experimental design, formal justification for the use of Rosenbaum's method after the fact requires us to make three assumptions.

First, for all observations, the observed treatment assignment  $z$  ( $z = j$  if the unit is assigned to treatment  $j$ ) and the vector  $r = (r_1, \dots, r_j)$  of potential responses to

treatment of that unit are conditionally independent given the vector of observed covariates. Second, regardless of the values taken by the covariates, all treatment assignments are possible. And third, the conditional probability  $e[X]$  of receiving a particular treatment given a vector of observed covariates  $X$ , follows a logistic model [Cox and Shell, 1989], that is

$$\log \left\{ \frac{e[X]}{(1 - e[X])} \right\} = \beta^T f(X),$$

where  $f(X)$  is a known but arbitrary vector-valued function of  $X$ . Since  $f(X)$  is arbitrary, this latter condition is not particularly restrictive.

*All three of these assumptions are satisfied if the covariates did not affect the treatment assignment.* For example, obesity and smoking history would satisfy these conditions if they were not factors in the patient/physician decision to have or perform a specific treatment.

### 9.3. Outliers

Consider the set of observations 0, 1, 2, 3, 19. Does the 19 represent a genuine response to treatment, the response we have been looking for, or is it an outlier—a typographical error or a bad reading that will only lead us astray? In the first case, we will want to utilize the data just as they are; in the second, we will want to modify or perhaps even to discard the questionable reading.

Shall we deal with such outliers on a one-by-one basis? Or should we establish a policy that will automatically adjust for and diminish the effect of outliers? Ad hoc rejection of suspect data could lead to charges of bias. A systematic policy can be adjusted for sample size and power determinations.

We consider seven policies here:

- 1) preserving the original data;
- 2) using ranks in place of the original observations, thus diminishing the effects of outliers;
- 3) replacing the observations/ranks by scores derived from some standard distribution, e.g., the order statistics of a standardized normal distribution;
- 4) applying a robust tail-compression transformation to all the data;
- 5) using an  $L_1$  test;
- 6) censoring extreme observations;
- 7) deleting extreme observations.

Whichever policy we elect, the permutation method will be more robust to outliers than a test based on a parametric distribution. The influence functions of a two-sample permutation test are always bounded above, even if the influence functions of the corresponding parametric test are unbounded from above and below [Lambert, 1981]. Our only concern need be the selection of a test statistic that is both practical and optimal.

### 9.3.1. Original Data

"The Method of Randomization applied to the original observations produced stunningly efficient tests which were disarmingly impractical." [Bradley, 1968]

Despite these discouraging words from James V. Bradley, I almost always make use of the original observations rather than their transform.

The exception that proves the rule is in my analysis of the Renis data considered in problem 5 of Chapter 3 and in Good [1979]. In that study, I used a preliminary logarithmic transformation, but it was to equalize the variances in the two samples, not to eliminate large values.

The computational difficulties to which Bradley alluded have largely been resolved through advances in computer technology between 1968 and today; the efficiency of the permutation test remains. The power and high relative efficiency of the permutation test comes from its use of exact values. Throw away one of the observations or replace it with its rank or a trimmed value and you reduce the power of the corresponding test. The gain in power is particularly evident when there is a mixture of responders and nonresponders [Good, 1979]; but see Boos and Browne [1986].

On the other hand, a single extreme observation often can have a disproportionate effect. Given the observations 0, 1, 2, 3, 19, would you rather guessimate the population mean as 2 or 2.5 than estimate it using the sample mean of 5? By taking ranks or applying some other tail-compressing transformation to all the observations, we can "democratize" the data so that each data item has a relatively equal influence upon the final calculation. (See also Hampel et al [1986].)

### 9.3.2. Ranks

Suppose we have two samples: The first control sample takes values 0, 1, 2, 3, 15. The second treatment sample takes values 3.1, 3.5, 4, 5, and 6. Does the second sample include larger values than the first?

When we rank the data giving the smallest observation a rank of 1, the next smallest the rank of 2, and so forth, the first sample includes the ranks 1, 2, 3, 4, 10, and the second sample includes the ranks 5, 6, 7, 8, 9. Does the second sample include larger values than the first?

Applying the two-sample comparison described in Section 3.2 to the ranked data, we conclude at the 10% level that the second sample is significantly larger. The sums of the ranks in the original first sample, 20, is as large or larger in just 19 of the  $\binom{10}{5} = 252$  rearrangements.

Obviously, taking ranks diminishes the effects of outliers. Taking ranks has a second advantage from the computational point of view: When we take ranks, the results are unconditionally distribution free. As we are working with the same values—the ranks, over and over regardless of the actual values of the observations,

we can tabulate the significance levels of our test statistics (at least for small samples) and avoid lengthy computations. And we may determine analytically when a sample of ranks is large enough that its permutation distribution may be replaced by an asymptotic approximation. It's not surprising that much of the literature on distribution-free tests is devoted to an analysis of the permutation distributions of ranked data.

The cost of using ranks is a loss of power, that is, a diminished probability of detecting a real difference between the distributions under test. But it is not a great loss. To achieve the same power as the permutation or parametric  $t$ -test with very large samples, the Mann–Whitney test—a two-sample comparison that uses ranks in place of the original observations—requires only 3% or 4% more observations. Cheap, if the units are widgets; expensive, if the units are patients or rare Rhesus monkeys.

### 9.3.3. Scores

If we are testing against normal alternatives, we can improve on the power of the Mann–Whitney test by using normal scores in place of ranks.

In the general case, we replace the rank of the  $i$ th observation,  $r_i$ , say, by the expected value of the  $r_i$ th largest value in a sample of  $n$  values drawn from the distribution  $F$ ,  $F^{-1}[r_i/(n + 1)]$ , where  $F$  is our best guess of how the observations are really distributed; (see also David [1970, p. 65]).

A good guess will produce an optimal test, and, sometimes, even a “bad” guess can be close to optimum. For example, Chernoff and Savage [1958] show that the normal-scores test, in which  $\phi$  is the Gaussian distribution, has a minimum asymptotic efficiency of 1 relative to the usual  $t$ -test regardless of the true underlying distribution.

Bell and Doksum [1965] provide detailed comparisons of the rank and normal scores tests in a variety of settings. In Bell and Doksum [1967] they provide conditions under which the normal-scores test is minimax.

Hajek and Sidak [1967] show that, in general, optimal scores for tests of location are based on the scores

$$a(j) = -\frac{f'(F^{-1}[u])}{f(F^{-1}[u])},$$

where  $u = j(N + 1)$ , and  $f$  and  $F$  are the density and cumulative distribution functions, respectively, of the underlying distribution. For optimal rank tests of scale, the scores are

$$a(j) = 1 - \frac{F^{-1}[u]f'(F^{-1}[u])}{f(F^{-1}[u])}.$$

### 9.3.4. Robust Transformations

A robust transformation preserves sample values at the center of a distribution while shrinking those in the tails. As one example [Maritz, 1981], consider

$$\phi(u) = u/(1 + u^2).$$

For  $u$  small,  $\phi(u)$  is approximately  $u$ . For  $u < 1$ ,  $\phi(u)$  is a slowly increasing function of  $u$ . If we replace  $x_i$  by  $\phi(x_i)$  in computing the mean, then large values will make virtually no contribution to the total.

As a second example Huber [1972], take

$$\phi(u) = (1 - \exp[-u])/(1 + \exp[-u]).$$

Again  $\phi(u)$  is approximately  $u$  for  $u$  small, and is bounded between 0 and 1.

In a complex experimental design, the transformation may be applied to the residual rather than the original observation. For example, to test whether  $Y = bX$ , one would apply  $\phi$  to  $y' = y - bx$ , rather than to  $y$ .

If you are uncertain which transformation to use, you can reduce the effect of extreme values in some cases simply by switching to a statistic based on the absolute differences  $|x_i - y_i|$  in place of the squared differences  $(x_i - y_i)^2$ . The final choice should be dictated by your loss function (see Section 10.4).

If extreme values are unlikely, as is the case with normal alternatives, then a robust transformation will have little or no effect on the power of a test. See Maritz [1981] and Lambert [1985] for further discussion.

### 9.3.5. Use an $L_1$ Test

A test based on the absolute values of the deviations about the median rather than the squares of the deviations about the mean is less likely to be affected by extreme values. Such a test is also the appropriate one to use with a first-order loss function. See Wilson [1978], Mielke [1986], Dodge [1987], Wang and Scott [1994], Cade and Richards [1996], and Mielke and Berry [1997].

### 9.3.6. Censoring

Lambert [1985] offers a two-sample test that is both robust and powerful. First, we order the data, so that

$$X_{(1)} < \cdots < X_{(n)} \quad \text{and} \quad Y_{(1)} < \cdots < Y_{(m)}.$$

To test against the alternative that the  $Y$ 's are larger on the average than the  $X$ 's, we replace each  $X_i$  and  $Y_j$  that is less than  $k_1 = X_{(n\beta_1)}$  by  $k_1$  and each  $X_i$  and  $Y_j$  that is greater than  $k_2 = Y_{(n\beta_2)}$  by  $k_2$  and then carry out the usual permutation test based on the sum of the observations in the first sample. Note that the censoring

values are determined by the data itself. Unfortunately, there can be more than one “right” choice for  $\beta_1$  and  $\beta_2$ , and the computations are far from straightforward. One possible compromise is to let  $k_1 = X_{(2)}$  and  $k_2 = Y_{(m-1)}$  for samples of 15 or less.

### 9.3.7. Discarding

The most extreme method of dealing with outliers is to discard them. Although Welch and Guitierrez [1988] obtain narrower confidence intervals in matched-pairs designs through the use of permutation applied to trimmed means, there are two objections to this method. First, the resultant test is unlikely to be exact (Theorem 3.3, [Romano, 1990]). Second, discarding data reduces the power of the test. In Good [1991], I improve on the power of the Welch–Guitierrez test by treating the outliers as if they were censored. My approach is described in more detail in the next section.

## 9.4. Censored Data

We may not be able to make all our measurements with the same precision.

In a radioimmune assay, for example, the typical concentration curve has a sigmoidal shape with flat regions at the two extremes. In the lower, flat region of the curve, estimation is difficult, if not impossible. While binding values elsewhere may be determined to one part in a billion, in this region they merely are recorded as “below minimum.”

Here is a second example: In many clinical studies, it is neither possible nor desirable to follow all patients to the end of their lifespans. Limiting the duration of the study cuts the costs of observation and puts promising new materials and processes into immediate service. But while some lifespans will be known with precision, others can be noted only as “exceeded treatment period.”

In each of these examples some of the data have been censored.

### 9.4.1. GAMP Tests

When observations are censored, the most powerful test typically depends on the alternative, so that it is not possible to obtain a uniformly most powerful test.

Good [1989, 1991, 1992] found that by establishing a region of indifference, it may be possible to obtain a permutation test that is close to the most powerful test, “almost most powerful,” regardless of the underlying parameter values.

Suppose we wish to perform a test of a hypothesis  $F$  against a series of alternatives  $F_1, F_2, \dots$ . To obtain a test that is globally almost most powerful (GAMP), we proceed in three stages.

First, we use the likelihood ratio to obtain a locally most powerful unbiased  $\alpha$ -level test of the hypothesis  $F$  against the alternative  $F_1$ . We repeat this procedure for each alternative  $F_i$  to obtain a family of rejection regions  $\{R_i\}$ .

Next, we form two regions: (i) A rejection region  $R \subseteq \cap R_i$  that contains only events common to all the rejection regions of the preceding family; and (ii) an acceptance region  $A$  that contains only events common to all the acceptance regions.

Last, we construct a permutation test whose  $p$ -value is determined by assigning each rearrangement of the data to one of three regions: Rejection ( $R$ ), acceptance ( $A$ ), or indifference ( $I$ ). While we cannot determine the  $p$ -value of their new test exactly, we can bound it:

$$\Pr\{R | X\} \leq p \leq 1 - \Pr\{A | X\}.$$

In Good [1992], I showed that GAMP's exist when the joint log likelihood of the observations takes the particularly simple form  $S_U * f(\theta) + N_C * g(\theta)$  where  $S_U$  and  $N_C$  are the sum of the uncensored observations and the number of censored observations in the treatment sample, respectively, and  $f$  and  $g$  are monotone functions of  $\theta$ . Examples include normally distributed, exponentially distributed, and gamma distributed random variables subject to Type I censoring.

A permutation (or rerandomization) approach is utilized.

There are two distinct cases, which I term left- and right-censoring respectively, though the actual directions—left or right—will depend upon the alternative. To fix ideas, suppose we have samples from two populations and are testing a null hypothesis  $H: F_2 = F_1$ , against stochastically larger alternatives,  $K: F_2(x) = F_1(x - \delta)$ . With left-censoring, we can assign  $x$  a precise value only if  $x \geq c$ ; for example, radioimmune assay involves left-censoring. With right-censoring, we can assign  $x$  a precise value only if  $x \leq c$ ; for example, reliability studies usually involve right-censoring.

To eliminate any dependence on the zero point of the underlying scale, we transform the data before we derive the permutation distribution; from each of the original observations we subtract  $\bar{X}_U$  the mean of the uncensored observations in the sample taken from  $G$ ;  $X'_{ij} = X_{ij} - \bar{X}_U$ , for  $i = 1, 2, j = 1, \dots, n_i$ ; and  $S'_{U_0} = 0$  and the transformed observations are censored at  $c' = c - \bar{X}_U$ . Next, we compute  $S_{U_0}$  and  $N_{C_0}$  for the original treatment sample; and permute repeatedly, computing  $S_U$  and  $N_C$  for each permuted sample.

With left-censoring, we assign a permutation to the rejection region  $R$  if  $S_U \geq S_{U_0}$  and  $N_C \geq N_{C_0}$ . We assign it to the acceptance region  $A$  if  $S_U < S_{U_0}$  and  $N_C \leq N_{C_0}$ . We assign it to the indifference region otherwise.

With right-censoring, we impute the value  $c$  to the censored observations. Let  $k = N_C - N_{C_0}$ . We assign a permutation to the rejection region  $R$  if  $S_U + kc \geq S_{U_0}$ . We assign it to the acceptance region  $A$  if  $S_U + kc < S_{U_0}$ . We assign it to the indifference region otherwise.

The indifference region is small enough in most instances to permit effective decision making [Good, 1989]. As the sample size increases, the GAMP test converges in probability to a UMP unbiased test [Good, 1992]. In the rare case

where the result does lie in the indifference region, I recommend taking additional observations.

The application of permutation methods to censored data was first suggested by Kalbfleisch and Prentice [1980], who sampled from the permutation distribution of censored data to obtain estimates in a process akin to bootstrapping.

For a survey of other permutation tests that have been applied to censored data, see Schemper [1984]. Conditional rank tests for randomly censored survival data are described by Andersen et al [1982] and Janssen [1991].

## 9.5. Censored Matched Pairs

As we showed in Section 3.7, the sensitivity of an experiment can be increased through the use of matched pairs. But it may happen that an exact observation cannot be made for one or more subjects, the only available information being that the required measurement is greater or less than some known value. Often this censoring process is accidental, but in many toxicology studies and reliability trials, it is a matter of deliberate design: The experimenter trades the cost of enrolling a larger number of subjects at the onset of the experiment for a shortened study period.

Suppose  $z = y - x$  is the difference between the (transformed) observations on the two members of a pair, and that observations are not recorded if they exceed  $C$  on the (transformed) scale. As noted by Sampford and Taylor [1959], any pair provides information on the distribution of  $z$  in one of the following four forms:

- (i) both  $y$  and  $x$  are observed, so that  $z$  is determined exactly;
- (ii)  $x$  is observed, but we only know that  $y$  exceeds  $C$ ; that is  $z > C - x$ , so we say  $z$  is upper censored;
- (iii)  $y$  is observed, but we only know that  $x$  exceeds  $C$ ; that is  $z < y - C$ , so we say  $z$  is lower censored;
- (iv) both  $x$  and  $y$  exceed  $C$ , so that no information is available on  $z$  for this pair; the sample size is effectively reduced.

While cases (ii) and (iii) provide less information than case (i), they are not uninformative, and a variety of hypothesis testing methods have been proposed for capitalizing on the information they provide. Recently [Good, 1991], I developed an “almost” most powerful distribution-free method based strictly on the data at hand. To see how this method is applied, assume that the first observation in each pair has the distribution  $F$  and the second has the distribution  $G$ . The hypothesis, unless stated to the contrary, is that  $F \geq G$ . The alternative is that  $F < G$ .

### 9.5.1. GAMP Test

The GAMP test for matched pairs represents a simple extension of the GAMP test for two independent samples derived in Good [1989, 1992]. Record  $U$ , the number

of upper censored pairs in the original sample, and  $Z$ , the sum of the uncensored  $z$ 's in the original sample. Randomize the observations, permuting the treatment labels within each pair, and let  $U'$  and  $Z'$  be the corresponding statistics for the permuted sample.

If  $U' \geq U$  and  $Z' \geq Z$ , then assign the permuted sample to the rejection region  $R$ .

If  $U' \leq U$  and  $Z' < Z$ , then assign the permuted sample to the acceptance region  $A$ .

Otherwise, assign the permuted sample to a region of indifference.

Repeat the randomization process for all possible permutations (or for a suitably large number  $N$  of randomly selected permutations) and let  $f_R$ ,  $f_A$ , and  $f_I$  be the frequency with which permutations are assigned to the rejection, acceptance, and indifference regions, respectively.

This method of construction ensures that the acceptance region  $A$  of the GAMP test is contained in the acceptance regions of each of the most powerful  $\alpha$ -level permutation tests of a simple hypothesis  $G = F = F^*$  against the simple alternative  $G^* = G > F = F^*$ . Similarly, the rejection region  $R$  of the GAMP test is contained in the rejection regions of each of the most powerful  $\alpha$ -level permutation tests.

$f_R \leq p \leq N - f_A$ , where  $p$  is the significance level of any member of the family of most powerful permutation tests of a simple hypothesis against a simple alternative. Thus, a test of the composite hypothesis  $F \leq G$  against the composite alternative  $F > G$  based on the bounds defined by  $A$  and  $R$  is globally almost most powerful, or GAMP.

In practice, an investigator using a GAMP will elect one of three courses of action: 1) accept the null hypothesis, noting the bounds on the  $p$  level; 2) reject the hypothesis in favor of a stochastically larger alternative; or, 3) in order that  $p$  might be known with greater certainty, elect to take additional observations. If you require exact significance levels to make power comparisons with other tests, you must randomize on the indifference region as follows.

If  $f_R$  is greater than the desired  $\alpha$ -level, accept the null hypothesis. If  $N - f_A$  is less, reject. If neither condition holds, choose a random number  $Z = U(0, 1)$  and reject the hypothesis if  $Z \leq (N\alpha - f_R)/(N - f_R - f_A)$ , accepting it otherwise.

### 9.5.2. Ranks

When data are heavily censored, you can improve on this method by replacing the original observations with ranks. Two approaches suggest themselves: In the first, which I term “post-ranking,” compute the differences,  $z$ , for each pair, then rank these differences in absolute value, dividing the highest ranks among the censored observations. Denote by  $Z$  the sum of the ranks which correspond to those pairs in which  $y$  is known to be larger than  $x$ . As in the GAMP test, now randomize the observations, permuting the treatment labels within each pair, and denote by  $Z'$

the new rank sum. Assign this randomization to  $R$ ,  $I$ , or  $A$  according to whether  $Z' >$ ,  $=$ , or  $<$  than  $Z$ . As with the GAMP test, reject  $H$  in favor of  $K$  if only a small proportion of rerandomizations are assigned to  $R$ ; randomize on the indifference region  $I$  to obtain a test at a specific significance level  $p$ .

Post-ranking has the drawback that if, say, 2 is the censoring point, the difference “censored–1.99” is automatically assigned a higher rank than the difference “1.99 – 0.” To avoid this difficulty, in a second approach, which I term preranking, first rank the individual observations, again dividing the highest ranks among the censored observation. Next, compute the differences of the ranks within each pair, and, as a third and final step, rank the absolute values of the differences. The drawbacks of this second, preredanked approach are computational: you must rank the data twice and you must correct for ties during the second ranking.

When the underlying distribution is normal and censoring is heavy, the preredanked permutation test provides the greatest sensitivity [Good, 1991].

When the underlying distribution is normal and censoring is light, or when the underlying distribution is exponential, the GAMP test is preferable.

The strength of the GAMP lies in its use of exact values rather than ranks—thus its effectiveness with heavy-tailed distributions, like the exponential, which have many extreme values. The GAMP is also the most readily computed. Its weakness lies in its dependence on a region of indifference whose size varies from sample to sample.

### 9.5.3. One-Sample: Bootstrap Estimates

If you are willing to assume the underlying distribution(s) are symmetric, then these methods for paired comparisons may also be applied to hypotheses based on a single sample. If censoring is one-sided, we are forced to censor on the opposite side in order to obtain an exact test. If you are unwilling to assume symmetry, and/or to throw away data through censoring, have 15 or more observations (30 would be better) and are willing to assume that all observations are drawn from the same distribution, then you may apply Efron’s [1981] bootstrap method of extending the Kaplan–Meir estimates.

## 9.6. Adaptive Tests

In an adaptive test [Hogg and Lenth, 1984], we compute several different test statistics, but make use only of the one we estimate to be the most powerful. For example, we could compute both a  $t$ -test and a robust test based on an  $M$ -estimate and, after the fact, use the one which seems best suited to the data. With some adaptive methods, the frequency of Type I error may increase as a result of this selection procedure. But with Donegani’s method [1991] applied to

two permutation tests, we can obtain a single test that is both exact and equal in power asymptotically to the most powerful of the two tests.

Let  $T_1$ , and  $T_2$  be the two tests and let  $c_1$ , and  $c_2$ , the “criteria”, be two positive real functions defined on the vector of observations  $X$  such that if  $c_1(X) < c_2(X)$ , then  $T_1$  is preferable to  $T_2$ . Suppose that large values of either test statistic indicate a departure from the null hypothesis. Proceed in four steps as follows.

1. Evaluate  $c_1(X)$ ,  $c_2(X)$  and let ‘opt’ refer to the index of the criterion having the smaller value.
2. Partition the set,  $P$ , of all possible rearrangements of the data into two sets

$$P_1 = \{\pi : c_1(\pi X) < c_2(\pi X)\}$$

$$P_2 = \{\pi : c_1(\pi X) > c_2(\pi X)\}.$$

3. Let  $H_{opt}$  be the randomization distribution obtained by evaluating the optimal test statistic  $T_{opt}$  on each element of the set that contains the original rearrangement.
4. Reject the null hypothesis at the level  $\alpha$  if  $T_{opt}$  exceeds the 100- $\alpha$ th percentile of  $H_{opt}$ . In other words, if  $c_1(x) < c_2(X)$  restrict attention to those rearrangements that are in  $P_1$ .

Let  $N_i$  denote the number of rearrangements in  $P_i$ . Let  $C_i$  denote the choice of the statistic  $T_i$ . Then

$$\begin{aligned} P\{R | H\} &= P\{R | H, C_1\}P\{C_1 | H\} + P\{R | H, C_2\}P\{C_2 | H\} \\ &= \alpha(N_1/(N_1 + N_2)) + \alpha(N_2/(N_1 + N_2)) \\ &= \alpha. \end{aligned}$$

Donegani [1991] shows that this adaptive procedure is asymptotically optimal and, in the case of matched pairs, that it is optimal with as few as nine pairs of observations.

## 9.7. Questions

1. Prove that ranking the data will eliminate any distortions brought about by a nonlinear measuring device. That is, prove that the ranks of the observations are invariant under any continuous, strictly increasing transformation. (We take advantage of this result in a multivariate analysis in which we use ranks to bring several disparate variables together on a single common scale; see Section 5.2.)
2. Show that an exact one-sample permutation test for singly censored data can exist only if you deliberately censor the data from the other side.
3. Let  $x_1, \dots, x_n$  be a sample from the exponential distribution with density  $\frac{1}{b}e^{-x/b}$ ,  $b > 0$ . If you have a scintillation counter at hand, you can generate just

such a sample by recording the time elapsed between counts. Alternately, you may stand on a street corner or at night club entrance and record the number of seconds before the next redhead or the next BMW goes by. If you have access to a computer, use its random number generator and take the logarithms of the random numbers you generate. Guesstimate the mean waiting time,  $b$ , before you start. Test your guesstimate (see Section 3.1) using a) the original observations, b) ranks, c) normal scores, and d) the data remaining after you've thrown out all observations that are three times the guesstimated value. Compare your results with the different statistical procedures for samples of size 5, 6, and 7.

## CHAPTER 10

# Which Statistic? Solving the Insolvable

### 10.1. The Permutation Distribution

Many common statistical problems defy conventional parametric analysis simply because the distributions of the resultant test statistics are not well tabulated, or, worse, we settle for a less-than-optimal statistic simply because a table for the less-than-optimal statistic is readily available—the chi-square statistic (Section 6.3) and its misapplication to sparse contingency tables is one obvious example.

We need not settle for less than the best. Given a sufficiently powerful computer and the time needed to perform the necessary calculations, we can always obtain the permutation distribution of the statistic that best separates the hypothesis from the alternative.

The freedom of choice provided by permutation methods creates its own new set of problems. Given complete freedom in the selection of a test statistic, which statistic are we to choose?

The purpose of this chapter is twofold: 1) To describe a number of practical applications in animal behavior, atmospheric science, education, epidemiology, molecular genetics, and sociology where permutation distributions have provided new and more powerful solutions; and 2) to provide some general rules to use in the derivation of test statistics for your own demanding applications.

### 10.2. New Statistics

#### 10.2.1. Nonresponders

In this section, we consider several new statistics designed specifically for use in a permutation test. An elementary example is a statistic Good [1979] proposed for use when there is a response threshold, a common occurrence in pharmacological studies.

We assume that the control observations  $\{X_i\}$  are independent and identically distributed with distribution  $F$ , while responders in the treatment group are

independent and identically distributed as  $G[x] = F[x - \delta]$ . Unfortunately, not every member of the treatment group is capable of responding to the treatment, with the result that we are forced to test the hypothesis  $G = F$  against contaminated alternatives of the form

$$G = \pi F[x - \delta] + (1 - \pi)F[x], \quad \text{with } 0 < \pi \leq 1.$$

The conventional statistics for the two-sample comparison—Student's t and the Wilcoxon test—are subject to a loss of power in the presence of nonresponders. This reduction in power of the t-test is due to two factors: 1) A decrease in the absolute difference between the means of the two testing groups and 2) an increase in the variance of the treatment sample. This last change is the key to the selection of a new test statistic:

$$\nu(p) = p' \frac{nm}{n+m} (X - Y)^2 + (1-p)S_y^2.$$

This new statistic has two components: The first is proportional to the difference  $(X - Y)$  in the means of the two samples, and the second to  $S_y$ , the variance of the treatment sample.

Barring the availability of an independent test for response, the  $p$  used in the equation for  $\nu$  is at best only a guess as to the true value of  $\pi$ . Good [1979] found that using a value of  $p = 0.67$  appears to offer relatively good protection against a broad range of values of  $\pi$ . Boos and Browne [1986] question whether the gain in power is really worth all the extra computation. An increase in power can mean a decrease in sample size with fewer experimental subjects placed at risk and a shortened study time with more rapid dissemination of important results. An increase in computation time puts the strain where it belongs—on the computer.

### 10.2.1.1. Extension to K-Samples

Mielke and Berry [1994] have extended Good's result to  $k$  samples, choosing as their test statistic

$$S = \sum_{k=1}^K \frac{n_k}{N} \binom{n_k}{2}^{-1} \sum_{i < j} \Delta_{ij} \phi_k(i) \phi_k(j),$$

where  $n_k$  is the number of observations in the  $k$ th sample,  $\phi_k(j)$  is one if  $j$  belongs to the  $k$ th sample and zero otherwise, and  $\Delta_{ij} = |x_i - x_j|^m$ ,  $m > 0$ ; typically,  $m = 1$  or  $m = 2$ .

### 10.2.2. Animal Movement

Let  $\{(w_i, x_i), i = 1, \dots, n\}$  denote a series of paired observations on the successive positions of two organisms in space. We would like to know if the movements of the two organisms are independent or coordinated. The ecological literature favors

a test of independence based on the ratio of the actual distance travelled to the distance from the starting point:

$$R_1 = \frac{\sum \{(w_{i+1} - w_i)^2 + (x_{i+1} - x_i)^2\}}{\sum \{w_i^2 + x_i^2\}}.$$

Our own intuition suggests a more powerful test of the hypothesis of independence would result from using either

$$R_2 = \frac{\sum (w_i - x_i)^2}{\sum \{w_i^2 + x_i^2\}},$$

the ratio of the successive distances of the two organisms from each other and from the starting point, or

$$R_3 = \frac{\sum (w_{i+1} - w_i)(x_{i+1} - x_i)}{\sum \{w_i^2 + x_i^2\}},$$

the traditional measure of correlation.

We also favor  $R_2$  and  $R_3$  on the grounds of simplicity. To compute the permutation distribution of  $R_1$ , we need to rearrange both sets of movements  $\{w_i\}$  and  $\{x_i\}$ . To compute the permutation distribution of  $R_2$  or  $R_3$ , we only need to rearrange one set of movements. Whatever statistic we choose, we may use its permutation distribution to obtain a test of statistical significance.

### 10.2.3. The Building Blocks of Life

In a fascinating state-of-the-art biological application, DNA sequencing, Karlin et al [1983] use permutation methods to assess the significance of certain repeated patterns of nucleic acids in several viruses.

DNA, the self-replicating molecule that is the basis of life on Earth, is assembled from four specific nitrogenous bases—adenine, guanine, thymine, and cytosine. The sequence in which these bases occur in the DNA molecule determines the structure of the organism. The triplet of deoxyribonucleotides guanine-adenine-cytosine leads to the production of the amino acid asparagine, for example. At issue is whether certain repeated patterns involving multiple copies of lengthy nucleotide sequences is also significant or merely the result of chance. Studying the distribution of repeated patterns that result when one randomly reassigns the labels on the nucleotides while preserving the total numbers of each label, Karlin et al conclude that the observed patterns are statistically significant. Hasegawa, Krishino, and Yano [1988] approach an analogous problem in DNA sequencing using bootstrap methods. The unraveling of the biological significance of the patterns continues to be an important research problem.

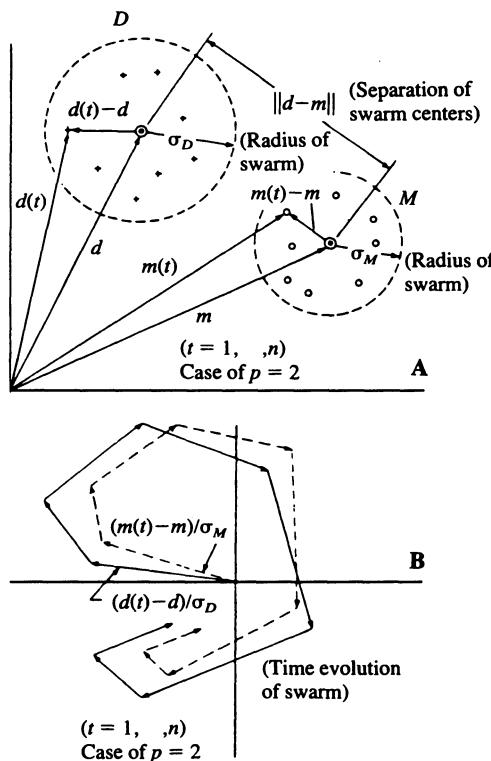


Figure 10.1. The geometric meaning of the trinity statistics SITES, SPRED, and SHAPE. The statistic SITES is essentially a dimensionless measure of the separation of data swarm centroids, while SPRED is a dimensionless measure of the differences in the root-mean-square radii of the swarms. The statistic SHAPE is a combined measure of the time evolution of the data swarms (and their associated maps). Note: From “The numerical model/reality intercomposition tests using small-sample statistics,” by R.W. Priesendorfer and T.P. Barnett, which appeared in *Journal of the Atmospheric Sciences*; 1983; 40: 1884–96. Reprinted with permission from the American Meteorological Society.

#### 10.2.4. Model Validation

The general circulation models of the Earth’s atmosphere and oceans used in weather and current prediction are of mind-boggling complexity, while the available data are all too finite. Priesendorfer and Barnett [1983] confront the problem of model-reality comparison studies for general circulation models head on by developing their own triple of metrics. In Figure 10.1a, and b which illustrates some of their concepts, the set  $D$  represents actual on-site data while  $M$  corresponds to a computer-generated model.

Rerandomization is accomplished in two steps. First, the data from  $D$  and  $M$  are combined into a single data set. Then, this combined set is repeatedly subdivided at random into sets of the same size as the original  $D$  and  $M$ . The resultant reference

distributions for each of the three metrics are used to assess the agreement of the model with reality.

How good is the Priesendorfer–Barnett test? The answer to this question illustrates the value of the permutation approach to the scientist and engineer whose primary training is not in statistics. For the answer does not depend on the abilities of Priesendorfer and Barnett as statisticians—the calculations in their test are straightforward—but on their abilities as meteorologists and oceanographers. Their test of statistical significance will be a good one, *if* they have selected the appropriate metric and the appropriate variables.

### 10.2.5. Structured Exploratory Data Analysis

A further illustration of this principle is given by Karlin and Williams [1984] in their use of permutation methods in a structured exploratory data analysis (SEDA) of familial traits. A SEDA has four principal steps.

- 1) The data are examined for heterogeneity, discreteness, outliers, and so forth, after which they may be adjusted for covariates (as in Section 4.3) and the appropriate transform applied (as in Section 9.3).
- 2) A collection of summary SEDA statistics are formed from ratios of functionals.
- 3) The SEDA statistics are computed for the original family trait values and for reconstructed family sets formed by permuting the trait values within or across families.
- 4) The values of the SEDA statistics for the original data are compared with the resulting permutation distributions.

As one example of a SEDA statistic, consider the OBP, the Offspring-Between-Parent SEDA statistic:

$$\frac{\sum_i^K \sum_j^{K_i} |O_{ij} - (M_i + F_i)/2|}{\sum_i^N |F_i - M_i|}.$$

In family  $i = 1, \dots, I$ ,  $F_i$  and  $M_i$  are the trait values of the father and mother (the cholesterol levels in the blood of the father and mother, for example), while  $O_{ij}$  is the trait value of the  $j$ th child,  $j = 1, \dots, K_i$ .

To evaluate the permutation distribution of the OBP, we consider all permutations in which the children are kept together in their respective family units, while we either:

- a) randomly assign to them a father and (separately) a mother; or
- b) randomly assign to them an existing pair of spouses. The second of these methods preserves the spousal interaction. Which method we choose will depend upon the alternative(s) of interest.

It would be difficult to establish the distribution of these measures or any other SEDA statistics analytically. To obtain the permutation distribution for the OBP

statistic, we merely substitute its formula (10.6) in place of the compute subroutine in our sample program (Section 4.2).

### 10.2.6. Comparing Multiple Methods of Assessment

We are often forced to combine several methods of assessment; one obvious example is in quality control; another is in grading students: Is an “A” in statistics equivalent to an “A” in Spanish? Direct comparisons are difficult, if not impossible, when students are free to choose their own courses. Table 10.1, reproduced with permission from Manly [1988], illustrates some of the problems associated with free choice: Missing data are one obvious problem. A second, hidden problem is that there is no guarantee that a student who is good in statistics will do equally well in Spanish.

The solution to both problems is to develop some kind of aggregate measure, compute this measure separately for each course, and then check to see how the distribution of this measure is affected by random relabelings of the students.

Table 10.2, also taken from Manly, illustrates the computation of just such a measure for the course in  $F$ . (The names of the actual courses have been changed to letters to protect the identities of overly-generous and overly-stingy graders.) The students are arranged in Table 10.2 in order of increasing mean grade. Each student’s mark in course  $F$  is subtracted from that student’s mean grade and the differences are cumulated.

If the *marks* in the various subjects are comparable, then each random rearrangement of an individual student’s marks is equally likely. For example, under the null hypothesis, student 6, who we see from Table 10.1 received marks of 75, 46, 45, and 64 in subjects  $A$ ,  $C$ ,  $E$ , and  $F$  might just as easily have received marks of 64, 45, 75, and 46 in those same subjects. Had this been the case, the CUMSUM score for subject  $F$  would have been 67.2 rather than 85.2. By looking at all possible arrangements of each student’s marks, we obtain a permutation distribution against which the CUMSUM score for the original arrangement can be assessed.

If the original score does not represent an extreme value, we conclude that the marking for subject  $F$  is consistent with the marking for the other subjects.

If, on the other hand the original CUMSUM score does represent an extreme value, our next step is to rescale the marks for subject  $F$ , subtracting and/or dividing by a constant. We repeat the test procedure using the rescaled values. And, in a manner akin to the way in which we derive a confidence interval (see Section 3.2), we continue testing and rescaling until all the marks in all the courses have been brought into alignment. Then, we may safely combine the assessments.

## 10.3. Going Beyond

At this point, you may already be thinking about several problems of your own for which you would like to develop an optimal test statistic. The purpose of this

Table 10.1. Examination Results from Seven Examinations (Subjects A-G) for 64 Students<sup>†</sup>

21	62	—	40	—	42	80	—	56.0	53	65	—	39	—	—	—	—	—	—	52.0
22	78	—	48	—	40	—	—	66	63.0	54	60	—	—	—	—	—	—	—	60.0
23	—	—	72	—	—	65	80	—	72.3	56	90	—	52	—	—	—	—	—	74.0
24	—	—	45	—	—	—	—	56	50.5	57	91	63	84	83	—	—	—	—	71.0
25	—	—	60	—	—	54	—	—	57.0	58	90	—	92	—	—	—	—	—	80.3
26	—	—	78	—	—	70	—	—	74.0	59	64	—	41	—	—	—	—	—	52.5
27	—	—	—	—	—	35	67	—	51.0	60	20	—	1	—	—	—	—	—	10.5
28	—	—	81	—	—	74	—	79	78.0	61	45	—	26	—	—	—	—	—	35.5
29	—	—	64	—	32	—	—	—	48.0	62	91	75	79	82	—	—	—	—	81.8
30	96	—	91	—	—	—	—	93.5	63	60	—	56	—	—	—	—	—	—	58.0
31	70	—	65	—	—	—	—	67.5	64	—	—	—	66	—	—	92	—	—	79.0

† A dash indicates that the student concerned did not sit the examination: *Sd.*, student number; *M*, student mean mark.

Note: From Manly [1988]. Reprinted with permission from the Royal Statistical Society.

Table 10.2. CUMSUM Calculations for the Subject  
*F* Marks of Table 1\*

Student	<i>F</i> mark	Mean	Difference	CUMSUM
38	42	32.5	9.5	9.5
28	67	51.0	16.0	25.5
21	80	56.0	24.0	49.5
6	64	57.5	6.5	56.0
33	51	58.5	-7.5	48.5
24	80	72.3	7.7	56.2
55	90	74.0	16.0	72.2
64	92	79.0	13.0	85.2
19	92	—	—	—

\*Student 19 only took subject *F*. There is therefore no comparison possible with other subjects and no contribution to the CUMSUM.  
 Note: From “The comparison and scaling of student assessment marks in several subjects” by B.F.J. Manly which appeared in Applied Statistics; 1988; 37: 385–95.

Note: Reprinted with permission from the Royal Statistical Society

last section of this chapter is to provide you with the basic principles of selection. While in Chapter 14 we consider a number of formal derivations based on the likelihood ratio, our approach in this chapter is more intuitive. The three essential concepts we consider are sufficiency, invariance, and loss.

### 10.3.1. Sufficiency

A statistic  $T(X)$  is *sufficient* for a parameter  $\theta$  (or a set of parameters  $\{\theta_i\}$ ) if the conditional distribution of  $X$  given  $T$  is independent of  $\theta$ . Once we have calculated the value of a sufficient statistic or statistics, we may be able to throw away the original observations, for frequently, a sufficient statistic(s) can provide us with all the information a sample has to offer.

An example we have already encountered is that of the order statistics  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ . If we know these order statistics, we know as much about the unknown distribution as we would if we had the original observations in hand.

Another commonly encountered example is that of the mean of a sample of independent, identically Poisson-distributed random variables, a statistic which is sufficient for the mean of the underlying Poisson distribution. Likewise the mean of a sample of normally distributed random variables is sufficient for the mean of the underlying normally distributed population. But there is distinction: In the first example, the Poisson, the sample mean possesses all the information the sample has to offer with regard to the underlying single-parameter distribution. A normal distribution depends on two parameters, the population mean and the population variance. We need to compute both the sample mean and the sample variance to obtain all the information a sample from a normal distribution has to offer.

In selecting a statistic to test a hypothesis about a population parameter  $\theta$ , look first at those statistics which are sufficient for  $\theta$ .

### 10.3.2. Invariance

If your measurements are made in feet, would you expect to reach the same conclusions as you would if your measurements were made in inches? What if you discover *after* you report your results that you forgot to rezero the measurement device so that each of your readings is off by exactly 0.0123 grams. Would you still believe that your decision to accept the hypothesis is correct? If your answers to both these questions is an unconditional “yes,” then you are already applying the principle of invariance, implicitly if not explicitly.

Many statistical problems involve symmetries. In the examples we’ve considered so far, the observations are exchangeable, so that the order in which we made these observations is irrelevant. Our test statistic(s) should and do reflect this same symmetry. The sample mean and sample variance are good examples of statistics that are symmetric in the underlying variables. Symmetry and invariance are related. The mathematical expression of symmetry is invariance under a suitable group of transformations. In generating an optimal test, look for test statistics that preserve the structure and symmetry of a problem.

### 10.3.3. Applying the Principles

Recall that in Section 4.2.3.1 we studied tests of the various hypotheses connected with the model

$$X_{ijb} = \mu + s_i + r_j + f_b + (sr)_{ij} + (sf)_{ib} + (rf)_{jb} + \epsilon_{ijb},$$

where  $i = 1, 2$ ;  $j = 1, 2$ ;  $b = 1, \dots, B$  replicates; by convention,  $s_1 + s_2 = 0$ ;  $r_1 + r_2 = 0$ ;  $\sum \sum (sr)_{ij} = 0$ ;  $\sum f_b = 0$ ;  $\sum \sum (sf)_{ib} = 0$ ;  $\sum \sum (rf)_{jb} = 0$ . One such hypothesis was that of zero interaction between the two factors of primary interest, in symbols,

$$(sr)_{11} - (sr)_{12} - (sr)_{21} + (sr)_{22} = 0.$$

Following Welch [1990], we see that invariance of the model with respect to  $g(X_{ijb}) = X_{ijb} + (sf)_{ib} + (rf)_{jb}$  eliminates the nuisance parameters for the block effects and all main effects. One such maximal invariant is that adopted in Section 4.2.3.1,  $D = (D_1, \dots, D_B)$ , where  $D_b = X_{b11} - X_{b12} + X_{b22} - X_{b21}$ .

We argue now in terms of sufficiency. Replacing the  $X_{bij}$  by the corresponding terms in the model, we see that  $D_b = (sr) + (\epsilon)$ , where  $(sr) = (sr)_{11} - (sr)_{12} - (sr)_{21} + (sr)_{22}$  and  $(\epsilon) = \epsilon_{11b} - \epsilon_{12b} - \epsilon_{21b} + \epsilon_{22b}$ . Because of the assumption of exchangeability within blocks, we only need examine permutations within blocks. The only permutations of the data within block  $b$  that leave the likelihood of  $D$  unchanged are permuting  $X_{b11}$  with  $X_{b22}$  and/or permuting  $X_{b12}$  with  $X_{b21}$ ; the

nuisance interactions still cancel. As these permutations do not change  $D$ , there is no reduction of  $D$  by sufficiency. Under the hypothesis,  $D_b = (\epsilon)$  and the likelihood of  $D$  remains constant if the data in the  $b$ th block are rearranged by independent permutations of the  $i$  and  $j$  subscripts. These permutations just change the sign of  $D_b$ , and a sufficient statistic is  $|D_1|, \dots, |D_B|$ .

#### 10.3.4. Losses

A statistical problem is defined by three elements:

- 1) the class  $P = (P_\theta, \theta \in \Omega)$  to which the probability distribution of the observations is assumed to belong;
- 2) the set  $D$  of possible decisions  $\{d\}$  one can make on observing  $X = (X_1, \dots, X_n)$ ;
- 3) the loss  $L(d, \theta)$ , expressed in dollars, men's lives, or some other quantifiable measure, that results when we make the decision  $d$  when  $\theta$  is true.

When you and I differ in our assessment of the loss function, we are likely to differ in our assessment of the practical significance of Type I and Type II error and, hence, in our choice of test statistic.

The loss function should be a key factor in the selection of a statistical test. Even when we don't know the exact values taken by a loss function, we have some idea about its form. In many testing situations, for example, in the analysis of variance and in some matched pair applications, the traditional test statistic (or discrepancy measure in Mehta and Patel's terminology) is a function of the square of the distance between the observed or estimated values and the hypothesis. Yet the natural measure is the distance itself. A statistical procedure that minimizes the expected value of the one may not minimize the expected value of the other [Mielke and Berry, 1982, 1983].

The principal reason for using the square of the distance is that it yields a maximum likelihood solution when the underlying distribution is normal. An assumption of normality may or may not be justified while maximum likelihood itself can only be justified on the grounds of convenience.

A second and more compelling reason for using the square of the distance in the data space would be that the loss function, a discrepancy measure in the parameter space, is also proportional to the square. But if we are uncertain about the form of the loss function, wouldn't it be more natural to utilize a test statistic that is linear in both the data and parameter spaces? A first-order statistic will be more robust than a second-order statistic in the face of questionably large deviations [Dodge, 1987].

The permutation approach frees us to choose the test statistic that is best suited to the problem at hand. If a second-order statistic is called for, we may use it, and if a first-order statistic is more appropriate, we may take advantage of it instead. Through the use of resampling methods we are free to choose the statistic best suited to the problem.

Recall from Section 4.2 that if we have more than two levels of a factor, we have a choice of at least three test statistics:

$$F2 = \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^I n_{ijk} (X_{ijk\cdot} - X_{\cdot jk\cdot})^2,$$

a second-order statistic;

$$F1 = \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^I n_{ijk} |X_{ijk\cdot} - X_{\cdot jk\cdot}|,$$

a first-order statistic; and

$$R = \sum_{j=1}^J \sum_{k=1}^K \sum_{i=1}^I n_{ijk} f[i] (X_{ijk\cdot} - X_{\cdot jk\cdot}).$$

With the permutation approach, we are free to select the optimal statistic in accordance with both the alternatives of interest and the underlying loss function.

## 10.4. Likelihood Ratio

As we shall see in Section 14.2, the primary criteria for selecting a test statistic is the likelihood ratio. We assign to our acceptance region those values of our test statistic for which the likelihood under the hypothesis is much greater than it is under the alternative and to the rejection region those values which are much more likely under the alternative than they are under the hypothesis.

To see this intuitively, suppose the variables can take only a countable number of values,  $P_i\{X = x\} = p_i(x)$  for  $i = 0, 1, \dots$ .

The optimal test is obtained by finding a set of values  $S$  to form the rejection region for which the significance level

$$\sum_{x \in S} p_0(x) \leq \alpha$$

and the power

$$\sum_{x \in S} p_1(x) \text{ is a maximum.}$$

Which values of  $x$  should we include in  $S$ ? Clearly, we should include those values which contribute the least to the significance level while contributing the most to the power. In other words, we should include those values of  $x$  with the largest values of the likelihood ratio

$$r(x) = \frac{p_1(x)}{p_0(x)}.$$

We extend this result to continuous distribution functions in Section 14.2 with the fundamental lemma of Neyman and Pearson.

The cutoff—that is, the precise definition of “largest” values—is determined by the significance level. Using the likelihood ratio, we show in Chapter 14 that the same criteria which led to the  $t$ -statistic and the  $F$ -ratio for the parametric analysis of two and  $k$ -samples, respectively, leads to the use of statistics equivalent to the  $t$  and the  $F$ -ratio for the corresponding permutation analyses. In Section 6.1, the likelihood ratio is used to derive Fisher’s exact test and to show that it is the most powerful unbiased test we can use with a  $2 \times 2$  contingency table.

### 10.4.1. Goodness of Fit and the Restricted Chi-Square

In the next example, that of an  $r \times 1$  contingency table, we cannot derive a most powerful test that will protect us against all alternatives, but we can use the likelihood ratio to derive a most powerful test against those alternatives which are of immediate interest. The approach lends itself to any set of data for which we have knowledge of an underlying model.

Suppose the hypothesis to be tested is that certain events (births, deaths, accidents) occur randomly over a given time interval. If we divide this time interval into  $m$  equal parts and  $p_i$  denotes the probability of an event in the  $i$ th subinterval, the null hypothesis becomes  $H : p_i = 1/m$  for  $i = 1, \dots, m$ . Our test statistic is

$$\chi^2 = mn \sum_{i=1}^m \left( v_i - \frac{1}{m} \right)^2,$$

where  $v_i$  is the relative frequency of occurrence in the  $i$ th interval.

0	1	2	3	$m - 1$
$v_0$	$v_1$			$v_{m-1}$

To determine whether this test statistic is large, small, or merely average, we examine the distribution of  $\chi^2$  for all sets of frequencies  $\{v_i\}$  that satisfy the two conditions

- 1)  $v_i \geq 0 \quad i = 1, \dots, m$ ; and
- 2)  $\sum v_i = 1$ .

We reject the hypothesis if the fraction of tables for which  $\chi^2 \leq \chi_0^2$  is less than  $\alpha$ .

We can obtain a still more powerful test when we know more about the underlying model and, thus, are able to focus on a narrower class of alternatives.

Suppose, in contrast to the previous example, that we use the  $m$  categories to record the results of  $n$  repetitions of a series of  $m - 1$  trials, that is, we let the  $i$ th category correspond to the number of repetitions which result in exactly  $i - 1$  successes. If our hypothesis is that the probability of success is .5 in each individual trial, then the expected number of repetitions resulting in exactly  $k$  successes is  $\pi_k[.5] = n(\frac{m}{k})(.5)^m$ .

If we proceed as we did in the preceding example, then our test statistic would be

$$S_1 = \chi^2 = n \sum_{k=1}^m \frac{(v_k - \pi_k[.5])^2}{\pi_k[.5]}.$$

Such a test provides us with protection against a wide variety of alternatives. But from the description of the problem we see that we can restrict ourselves to alternatives for which

$$\pi_k[p] = n \binom{m}{k} (p)^k (1-p)^{m-k}.$$

Fix, Hodges, and Lehmann [1959] show that a more powerful test statistic against such alternatives is

$$S = S_1 - S_2,$$

where

$$S_2 = \min_p \sum_{i=1}^m \frac{(v_i - p_i[p])^2}{\pi_i[p]}.$$

The parametric form of the distribution of  $S$  is difficult if not impossible to obtain analytically, except for very large sample sizes; as always, we can approximate the permutation distribution by Monte Carlo means, assigning the  $\Sigma v_i$  items to the  $m$  categories at random and computing  $S_2$  for each such rerandomization.

### 10.4.2. Censored Data

In Section 9.5, we use the likelihood ratio to derive a globally almost powerful test for use with censored data.

Kalbfleisch and Prentice [1980] also use the likelihood ratio to obtain tests for use against highly specific alternatives when the underlying distributions are censored. The calculations are complex, so these authors suggest *bootstrapping* from the permutation distribution as a computational shortcut. Their test is appropriate when the parameters of the alternative are known with some precision. Against global and unspecified alternatives, the GAMP test described in Section 9.5 is to be preferred.

### 10.4.3. Logistic Regression

Finally, we use the likelihood ratio to derive a procedure which is of inestimable value in the analysis of epidemiological data. One of the earliest applications of logistic regression is that of Pike, Casagrande, and Smith [1975]. For each subject, we have a pair of observations,  $x_i$  the length of exposure and  $y_i$  the apparent effect, where  $y_i$  may be a vector of several variables. To eliminate extraneous variation, we divide the data into blocks based on age, duration of residence, marital

status, and so forth. Each block may be further subdivided into two not necessarily equal-sized groups—cases and controls. We would like to know if the risk of exposure is the same for each group and to estimate the relative risk.

Following Breslow and Day [1980, 1987], we condition the likelihood of  $x$  given  $y$  on the set of exposures without regard to which are cases and which are controls.

$$\frac{\prod_{j=1}^n L(x_j | y_j = 1) \prod_{j=1}^m L(x_j | y_j = 0)}{\sum_{\pi \in R} \prod_{j=1}^n L(x_{\pi(j)} | y_j = 1) \prod_{j=1}^m L(x_{\pi(j)} | y_j = 0)},$$

where  $R$  is the set of  $\binom{n+m}{n}$  possible reassessments  $\pi$  of case labels to subjects and the likelihood

$$L(x | y) = \frac{pr(y | x)pr(x)}{pr(y)}.$$

Assume that within a block, the observations satisfy the logistic regression model, so that

$$pr\{y | x\} = \frac{\exp[\alpha + \beta x]}{1 + \exp[\alpha + \beta x]}.$$

The conditional likelihood (10.1) reduces to

$$\frac{\prod_{j=1}^n \exp \left[ \sum_{k=0}^1 \beta_k x_{\pi(j)^k} \right]}{\sum_{\pi \in R} \prod_{j=1}^n \exp \left[ \sum_{k=0}^1 \beta_k x_{\pi(j)^k} \right]},$$

an expression which depends only on the relative risk parameters  $\beta_0$  and  $\beta_1$ .

## 10.5. Questions

1. Suppose you wish to compare two groups of observations. Would it be better to compare them using the two-sample comparison of Section 3.3 or the matched pairs technique of Section 3.6? Is your decision rule an “always . . .” or does it depend on how the observations are dispersed and the relative importance of the covariates used to do the matching?
2. Suppose you have discarded the  $n$  original observations in the sample, keeping only the  $n$  order statistics, when you obtain independent evidence that the data are normally distributed, can you still compute the sample mean and variance?
3. Suppose you have multiple observations on each subject, some in feet, some in inches, some in pounds. Should they all be transformed to a common unit of reference before you begin your multivariate analysis? What transformation(s) should you use?

4. What statistic(s) remain invariant under an arbitrary monotone increasing transformation of the observations? Is this result relevant to the preceding question?
5. Ninety-nine percent of all scientists ignore the loss function and make do with a pre-designated significance level and a minimum power level against one or two selected alternatives. Reconsider the statistical analyses you performed recently. What was the loss function in each instance? Were the test statistics you selected appropriate for this loss function?
6.
  - a) Can the four  $k$ -sample statistics,  $F_1$ ,  $F_2$ ,  $F_3$ , and  $R$  introduced in Section 4.2.2 be made equivalent to one another if we eliminate terms that are invariant under permutations?
  - b) If your answer to the previous question is “no,” will there be data sets for which tests based on  $F_1$ ,  $F_2$ , and  $R$  lead to different conclusions?
  - c) How would you decide which of these three statistics to use?
  - d) Are you free to compute the permutation distributions of  $F_1$ ,  $F_2$ , and  $R$  for a specific data set and then choose the statistic that does the best job of proving your point?
  - e) Suppose you were an examiner at the FDA; how would you react to a submission in which the authors had done just that?
  - f) If you were one of those authors, how would you justify your choice of test statistic to an examiner at the FDA?
  - g) Throughout this text, we have tried to justify our choice of statistic on the grounds that the resultant test was a) unbiased, b) most powerful, c) minimized losses, or d) invariant under transformations of location and scale. Do these criteria satisfy your own instincts? What other criteria can you suggest?
7. Are the two tests described in Section 10.2.1 equivalent when comparing two samples? If not, is one uniformly more powerful than another?

## CHAPTER 11

# Which Test Should You Use?

In this chapter, we provide you with an expert system to use in choosing an appropriate testing technique. Your expert system comes to you in two versions—a professional’s handbook with detailed explanations of the choices and a short “quick-reference” version at the end of the chapter.

### 11.1. Parameters and Parametric Tests

To perform a parametric test, we must assume the observations come from a probability distribution that has a specific parametric form. For example, an observation  $X$  has the Poisson distribution with parameter  $\lambda$  if the probability that  $X = k$  is  $\lambda^k \exp[-\lambda]/k!$  for  $k = 0, 1, 2, \dots$ . An observation  $X$  has the normal distribution with location parameter  $\mu$  and scale parameter  $\sigma$  if the probability density  $h(x)$  is

$$\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right].$$

While various techniques for verifying whether a set of observations does or does not have a Poisson or normal distribution exist, the following heuristic definitions have proven of great value in practice.

An observation has the Poisson distribution if it is the cumulative result of a large number of opportunities, each of which has only a small chance of occurring. For example, if we seed a small number of cells into a petri dish that is divided into a large number of squares, the distribution of cells per square follows the Poisson.

An observation has the Gaussian or normal distribution if it is the sum of a large number of factors, each of which makes a very small contribution to the total. This explains why the mean of a sufficiently large number  $N$  of observations,  $X = \sum_{i=1}^N X_i/N$ , will be normally distributed even if the individual observations  $X_i$  come from quite different distributions.

By contrast, proportions and the ratios of variables or sums of variables seldom have a normal distribution.

In many applications in economics and pharmacology where changes are often best expressed in percentages, a variable may be the product of a large number of variables, each of which makes only a near unit contribution to the total. Such a variable has the log-normal distribution, and, because  $\log[\Pi x_i] = \Sigma \log[x_i]$ , its logarithm has a normal distribution.

The normal distribution is easy to recognize. It is symmetrically distributed about the mean and falls off rapidly in the tails, so there is only a small probability of observing extremely large or extremely small values.

Many other distributions one encounters in practice—chi-square, Beta, Student's t, and the F-ratio are all examples—may be derived from variables that have the normal distribution. For example, if X has the normal distribution with mean 0 and variance  $\sigma^2$ , then  $Y = (X/\sigma)^2$  has the chi-square distribution with one degree of freedom.

Gamma distributions for which the density

$$f(x|a, b) = \frac{1}{\Gamma[a]b^a} x^{a-1} \exp[-x/b]$$

come into existence in complex systems, where the failure of several simple parallel components is necessary before the system fails to function.

The literature is replete with methods for determining whether observations have this or that distribution. My own preference is to use a permutation test whenever there is the slightest doubt as to the nature of the underlying population.

## 11.2. Parametric Tests, Permutations, and the Bootstrap

Of course, one should always use a parametric test

- 1) when you have a large number of observations ( $\geq 30$ ) in each category;
- 2) when you have a very small number of observations in each category ( $\leq 5$ ) and if the assumptions underlying the corresponding parametric test may be relied on. For example, if we have only three observations with which to test the hypothesis that the mean of a symmetric distribution is zero, the sample space for the permutation test is limited to  $2^3$  or eight rerandomizations. As a result, we must randomize on the boundary except for significance levels that are multiples of 1/8th. At all significance levels, a more powerful parametric test, and, if we may rely on the normality of the observations, a uniformly most powerful unbiased parametric test, may be obtained directly from tables of the t-statistic.

These caveats aside, in most practical testing situations, we would advise the reader to use a permutation test or, at least, to use the permutation distribution in place of the parametric distribution.

- o The permutation test is exact under relatively nonstringent conditions: In the one-sample problem, the variables must have symmetric distributions; in

the two- and k-sample problem, the variables must be exchangeable among the samples.

- The permutation test provides protection against deviations from parametric assumptions, yet it is usually as powerful as the corresponding unbiased parametric test even for small samples.

- With two binomial or two Poisson populations, the most powerful unbiased permutation test and the most powerful parametric test coincide. With two normal populations, the most powerful unbiased permutation test and the most powerful unbiased parametric test are asymptotically equivalent.

- Using the permutation test means you can choose the statistic that is best adapted to your problem and to the alternatives of interest. While most parametric tests rely on  $L_2$  norms, a reasonable procedure when the loss function is based on squared error, permutation tests permit us to use either the  $L_1$  or  $L_2$  norm, determined, as it should be, by the loss function.

Permutations test hypotheses concerning distributions; bootstraps test hypotheses concerning parameters. As a result, the bootstrap entails less-stringent assumptions. For example, in Chapter 3, we saw that, to perform a one-sample test of the median using a permutation test, we had to assume that all observations came from the same distribution, while with the bootstrap we only had to assume that all observations had the same median.

Nonetheless, consider a permutation test before you turn to a bootstrap. The bootstrap is not exact except for quite large samples and, often, is not very powerful. But the bootstrap can sometimes be applied when the permutation test fails: An example is interaction in an unbalanced design (Chapter 4.4) for which neither an exact parametric test nor an exact permutation test can be formulated.

### 11.3. What Significance Level Should I Use?

Choice of significance level and power are determined by the environment in which you work.

Most scientists simply report the observed p-value, leaving it to their peers to decide for themselves what weight should be given the results.

A manufacturer preparing to launch a new product line or a pharmaceutical company conducting a research for promising compounds typically adopts a three-way decision procedure: If the observed p-value is less than 1%, they go forward with the project. If the p-value is greater than 20%, they abandon it. And if the p-value lies in the gray area in between, they arrange for additional surveys or experiments.

A regulatory commission like the FDA that is charged with oversight responsibility must work at a fixed significance level, typically,  $\alpha = 0.05$  or  $0.10$ . The choice of a fixed significance level ensures consistency in both result and interpretation as the agency reviews the findings from literally thousands of tests. They specify a minimum power level of 90% or 95%, or even 99% against biologically significant alternatives if potentially dangerous side-effects are under investigation.

## 11.4. A Guide to Selection

The initial division of this guide is into three groupings: categorical data, discrete data, and continuous data.

### 11.4.1. The Data Are in Categories.

Examples include men vs. women, white vs. black vs. Hispanic vs. other, and much improved vs. improved vs. no change vs. worse vs. much worse.

Only a single factor is involved.

You are testing the goodness of fit of a specific model. See Section 10.4.1.

Only two factors are involved. For example, sex vs. political party.

Each factor is at exactly two levels:

There is a single table.

Use Fisher's Exact Test (See Section 6.1).

There are several  $2 \times 2$  tables.

Use Zelen's exact test (See Section 6.2.1).

One factor is at three or more levels.

This factor is not ordered as would be the case with a factor like race.

You want a test that provides protection against a broad variety of alternatives.

Use the permutation distribution of  $\chi^2$  (Section 6.4).

You wish to test against the alternative of a cause-effect dependence.

Use the permutation distribution of P [Freeman and Halton, 1958]  
(see Section 6.4 for other possible tests).

This factor can be ordered.

Use Pitman correlation (see Section 3.5).

Both factors are at three or more levels.

Neither factor can be ordered.

The alternative is the first factor is caused or affected by the other.

Use the permutation distribution of Kendall's tau or Cochran's Q  
(see Section 6.4).

A cause and effect relationship is not suspected.

Use the permutation distribution of  $\chi^2$  (see Section 6.4).

One factor can be ordered.

Assign scores to this factor based on your best understanding of its effects on the second variable (see Section 6.5.1).

Both factors can be ordered.

All the odds ratios are approximately equal.

Use  $\lambda_3$  or the Goodman–Kruskall test (see Section 6.5.2.2).

Some but not all of the odds ratios are close to one.

Use  $\lambda_2$  or the likelihood ratio test (see Section 6.5.2.2).

A third covariate factor is present.

Use the method of Bross [1964]. See Sections 6.6 and 6.7.

### 11.4.2. The Data Take Discrete Values

Each sample consists of a fixed number of independent, identically distributed observations that can be either zero or one. (A set of trials in which each may result in a success or a failure is one example.)

Only one or two samples are involved.

Use the parametric test for the binomial. See, for example, Lehmann [1986, pp. 81 and 154].

More than two samples, but only one factor is involved.

Analyze as indicated above under categorical data.

More than one factor is involved.

Transform the data to equalize the variances. For each factor combination, take the arcsin of the square root of the proportion of observations that take the value 1. Analyze the results as indicated below under continuous data.

Each sample consists of a set of independent, identically distributed Poisson observations.

Only one or two samples are involved.

Use the parametric test for the Poisson. In the two-sample case, note that the UMPU test uses the binomial distribution. See for example, Lehmann [1986, pp. 81 and 15]).

More than two samples are involved.

Transform the data to equalize the variances by taking the square root of each observation. Analyze as indicated below under continuous data.

Each sample consists of a set of exchangeable observations whose distribution is unknown.

There is a single sample.

The data may be assumed to come from a symmetric distribution. Use the permutation test for a location parameter that is described in Section 3.1.

The data may not be assumed to come from a symmetric distribution.

Use the bootstrap described in Section 3.4. If you have only a few subjects, consider using a multivariate approach (see Chapter 5).

There is more than one sample.

Use one of the permutation tests designed for data with continuous distributions, which is described in Chapters 3 and 4. Treat tied observations as separate distinct observations when you form rearrangements. Be cautious in interpreting a negative finding; the significance level may be too large simply because the test statistic can take on too few distinct values.

### 11.4.3. The Data Are Continuous

How precise do our measurements have to be so that we may categorize them as “continuous” rather than discrete? Should they be accurate to two decimal places

as in 1.02? or four as in 1.0203? To apply statistical procedures for continuous variables, the observations need only be precise enough so that there are no or only a very few ties.

If you can be sure the data has the normal distribution,

a parametric test like Student's t or the F-ratio may be applicable, but you can protect yourself against deviations from normality by making use of a permutation test based on the t statistic or the F.

You have only a single sample.

You want to test that the location parameter has a specific value, and you feel safe in assuming that the underlying distribution is symmetric about the location parameter.

Use the procedure described in Section 3.1.

If the distribution is not symmetric,

but has a known parametric form,

apply the corresponding parametric test.

If the distribution does not have a known parametric form,

consider applying an initial transformation that will symmetrize the data. For example, take the logarithm of data that undergoes percentage changes. Be warned that such a transformation affects the form of the loss function

and/or bootstrap; see Good [1999, Chapter 3]

You want to test that the scale parameter has a specific value.

First, divide each observation by the hypothesized value of the scale parameter. Then, apply one of the procedures noted above for testing a location parameter.

You have two samples.

You want to test for equality of the scale parameters of the two populations.

You know the means/medians of the two populations or you know they are equal.

Use the permutation distribution of the F-ratio based on the sample variances (see Section 3.4.1).

You have no information about the means/medians of the populations.

Use one of the methods of Section 3.4.

You want to test for equality of the location parameters of the populations.

If changes are proportional rather than additive, consider working with the logarithms of the observations.

If the data are censored or you suspect outliers, see Chapter 9.

Each sample consists of measures taken on different subjects.

Use the two-sample comparison described in Section 3.3.

Two observations were made on each subject.

Use the matched-pair comparison described in Section 3.7.

You have more than two samples.

If changes are proportional rather than additive, consider working with the logarithms of the observations.

If the data are censored or you suspect outliers, see Chapter 9.

A single factor distinguishes the various samples.

You can't take advantage of other factors to block the samples.

The factor levels are not ordered.

Use the permutation distribution of an F-ratio (see Section 3.5).

The factor levels are ordered.

Use Pitman's correlation (Section 3.5.2).

You can take advantage of other factors to block the samples.

Rerandomize on a block-by-block basis, then apply one of the techniques described in Sections 3.6 and 3.7.

Multiple factors are involved.

All observations are exchangeable.

The experimental design is balanced.

All the factors are under your control.

Use one of the permutation techniques described in Section 4.2.

Not all the factors are under your control.

First, correct for the functional relationship among factors or use restricted randomization as described in Section 4.3. Then, use one of the permutation techniques described in Section 4.2.

The experimental design is not balanced.

Some factors will be confounded. A book on experimental design, such as that of Kempthorne [1952], can help you determine which factors.

Consider the bootstrap (see Section 4.4.2) or the bootstrap-permutation test (see Section 4.4.3).

Multiple variables are involved.

Look first for a cause–effect model to be tested as a whole.

See Sections 5.1–5.4.

The variables are not all continuous; see Section 5.5.2.

Repeated measurements were made over time.

Treat the repeated measurements as components of a single multivariate vector. See Section 5.6.

## 11.5. Quick Key

Categorical data

Single factor,  $r = 1$ .

Goodness of fit, 10.4.1

Two factors,  $r = 2$ .

$c = 2$

single table

Use Fisher's Exact Test, 6.1

several  $2 \times 2$  tables

Use Zelen's exact test, 6.2.1

$c > 2$

- not ordered
  - use  $\chi^2$  or  $\tau$ , 6.4
- ordered
  - Use Pitman correlation, 3.5
- Two factors,  $r > 2$ ,  $c > 2$ 
  - not ordered
    - Use  $\tau$ , Q,  $\chi^2$ , 6.4
  - ordered
    - Use  $\lambda_2$  or  $\lambda_3$  6.5.2.2
- With covariate
  - Use Bross method, 6.6, 6.7.
- Discrete data
  - Binomial data
    - one factor, one or two samples
      - See Lehmann [1986, pp. 81 and 154].
    - one factor, more than two samples
      - See categorical data.
    - more than one factor
      - See continuous data.
  - Poisson data
    - one or two samples
      - See Lehmann [1986, pp. 81 and 152].
    - more than two samples
      - See under continuous data.
  - Other exchangeable observations
    - one sample.
      - symmetric distribution, see 3.1.
      - not symmetric, use bootstrap
      - more than one sample
        - transform data; see under continuous data
- Continuous data
  - one sample
    - test of location parameter
      - symmetric distribution, see 3.1
      - not symmetric
        - attempt to transform to known parametric or symmetric form.
        - test of scale parameter
        - rescale and test as for location parameter.
  - two samples
    - test equality of scale parameters
      - means/medians are known or are known to be equal
        - F-ratio of the sample variances, 3.4.1
        - otherwise permute or bootstrap, 3.4
    - test equality of location parameters

- samples not matched
  - two-sample comparison, 3.3.
- samples are matched
  - matched-pair comparison, 3.7
- more than two samples
  - single factor
    - no blocking
      - levels not ordered, F-ratio, 3.5
      - levels ordered, Pitman's correlation, 3.5.2
  - blocks
    - resample block by block, 3.6, 3.7, Chapter 4
  - multiple factors
    - balanced design
      - all factors under your control, 4.2
      - otherwise, correct as in 4.3, then apply 4.2
  - unbalanced design
    - consult text on experimental design; consider bootstrap, 4.4

## CHAPTER 12

# Publishing Your Results

McKinney et al [1989] report that more than half the published articles that apply Fisher's exact test do so improperly. Our own survey of some 50 biological and medical journals supports their findings. This chapter provides you with a positive prescription for the successful application and publication of the results of resampling procedures. First, we consider the rules you must follow to ensure that your data can be analyzed by statistical and permutation methods. Then, we describe two commercially available computer programs that can perform a wide variety of permutation analyses. Finally, we provide you with five simple rules to prepare your report for publication.

### 12.1. Design Methodology

It's never too late to recheck your design methodology. Recheck it now in the privacy of your office rather than before a large and critical audience. All hypothesis-testing methods rely on the independence and/or the exchangeability of the observations. Were your observations independent of one another? What was the experimental unit? Were your subjects/plots assigned at random to treatment? If not, how was randomization restricted? With complex multifactor experiments, you need to list the blocking variables and describe your randomization scheme.

#### 12.1.1. Randomization in Assignment

Are we ever really justified in exchanging labels among observations? Consider an experiment in which we give six different animals exactly the same treatment. Because of inherent differences among the animals, we end up with six different measurements, some large, some small, some in between. Suppose we arbitrarily label the first three measurements as "controls" and the last three as "treatment." These arbitrary labels are exchangeable and thus the probability is one in 20 that the three "control" observations will all be smaller than the three "treatment." Now

suppose we repeat the experiment, only this time we give three of the animals an experimental drug and three a saline solution. To be sure of getting a positive result, we give the experimental drug to those animals who got the three highest scores in the first experiment. Not fair, you say. Illegal! Illegitimate! No one would ever do this in practice.

In the very first set of clinical data that was brought to me for statistical analysis, a young surgeon described the problems he was having with his chief of surgery. "I've developed a new method for giving arteriograms which I feel can cut down on the necessity for repeated amputations. But my chief will only let me try out the technique on patients that he feels are hopeless. Will this affect my results?" It would and it did. Patients examined by the new method had a very poor recovery rate. But, of course, the only patients who'd been examined by the new method were those with a poor prognosis. The young surgeon realized that he would not be able to test his theory until he was able to assign patients to treatment at random.

Not incidentally, it took us three more tries until we got this particular experiment right. In our next attempt, the chief of surgery—Mark Craig of St Eligius in Boston—announced that he would do the "random" assignments. He finally was persuaded to let me make the assignment using a table of random numbers. But then he announced that he, and not the younger surgeon, would perform the operations on the patients examined by the traditional method to make sure "they were done right." Of course, this turned a comparison of methods into a comparison of surgeons and intent.

In the end, we were able to create the ideal "double blind" study: The young surgeon performed all the operations, but the incision points were determined by his chief after examining one or the other of the two types of arteriogram.

### 12.1.2. Choosing the Experimental Unit

The exchangeability of the observations is a sufficient condition for a permutation test to be exact. It is also a necessary condition for the application of any statistical test.

Suppose you were to study several pregnant animals that had been inadvertently exposed to radiation (or acid rain or some other undesirable pollutant) and examine their offspring for birth defects. Let  $X_{ij} i = 1, \dots, l; j = 1, \dots, n_i$  denote the number of defects in the  $j$ th offspring of the  $i$ th parent; let  $Y_i = \sum_{j=1}^{n_i} X_{ij} i = 1, \dots, l$  denote the number of defects in the  $i$ th litter. The  $\{Y_i\}$  may be exchangeable; (we would have to know more about how the data were collected). The  $\{X_{ij}\}$  are not; the observations within a litter are interdependent; what affects a parent affects all her offspring. In this experiment, the litter is the correct experimental unit.

In a typical toxicology study, a pathologist may have to examine three to five slides at each of 15 to 20 sites in each of three to five animals just to get a sample size of *one*.

### 12.1.3. Determine Sample Size

As noted in Chapter 2, the number of observations must be large enough that the resultant hypothesis test will have sufficiently high probability (power) of detecting effects that are of scientific and/or practical interest. Before you start, specify the significance level, the minimum effect of interest, and the desired power for that effect, then, use one of the methods described in Section 13.8 to determine the appropriate sample size.

You may need to conduct your experiment in several stages, using your initial efforts as a basis for estimating the population parameters needed in the power calculations.

## 12.2. Statistical Software for Exact Distribution-Free Inference

### 12.2.1. Freeware and Shareware

A no-name package of basic resampling procedures (bootstrap and permutation tests for a single sample, permutation tests for two samples, the one-way layout and  $2^2$  factorials) for use on PC compatibles can be downloaded from the author's website at

<http://users.oco.net/drphilgood>.

Program source in both C and Fortran 77 is available by anonymous FTP from <ftp.salford.ac.uk>. Use the following commands:

```
ftp          ftp.salford.ac.uk
ftp>user      anonymous
ftp>password  guest
ftp>          cd misc/perm
ftp>          get read.me
ftp>          get [put filename here]
```

Software for PC compatibles cross-referenced to the corresponding program listings in Edgington [1998] may be obtained from

<ftp.acs.ucalgary.ca>. Use the following commands:

```
ftp          ftp.acs.ucalgary.ca
ftp>user      anonymous
ftp>password  guest
ftp>          cd pub/private_group_info/randibm
ftp>          get readme.doc
ftp>          binary
ftp>          get randibm.exe
```

SMP provides six different methods for deriving the unconditional p-values for comparing two different proportions. Obtain from Antonio Lopes, [alopes@itqb.unl.pt](mailto:alopes@itqb.unl.pt).

*Blossom Statistical Analysis Package* This interactive program for analyzing data utilizing multiresponse permutation procedures (MRPP) includes statistical procedures for grouped data, agreement of model predictions, circular distributions, goodness of fit, least absolute deviation and quantile regression. Programmed by Brian Cade at the U.S. Geological Survey, Midcontinent Ecological Science Center. PC compatibles with online manual in HTML Frames format, may be downloaded from <http://www.mesc.usgs.gov/blossom/blossom.html>.

*NPSTAT* carries out parametric and randomization tests on one factor designs for independent groups (equal or unequal n), repeated measures, correlations, and Fisher's exact test. *NPFACt* carries out both parametric and randomization tests for 2-factor and 3-factor designs in independent, dependent, and mixed designs. For PC compatibles, 1217 Brickyard Road # 106, Salt Lake City, UT 84106 801/463-1839 [rmay@utahinter.net](mailto:rmay@utahinter.net). Download from <http://home.utahinter.net/rmay/npstat.html>.

*PERMUSTAT* computes exact and approximate p-values for k-samples and k-proportions. For Macintosh System 7; Andrew F. Hayes, Dartmouth College, Business Administration, Hanover, NH 03755. [Andrew.F.Hayes@dartmouth.edu](mailto:Andrew.F.Hayes@dartmouth.edu). Download from <http://mba.tuck.dartmouth.edu/pages/faculty/Andrew.Hayes/pstat.htm>.)

### 12.2.2. Commercial Software

*StatXact* uses the algorithms developed by Mehta, Patel, and their colleagues to help you perform a wide variety of permutation tests for one, two and k samples, ordered and unordered R × C contingency tables, and stratified 2 × 2 and 2 × C contingency tables. The two-sample procedures include stratified linear rank test, Wilcoxon-Mann-Whitney test, logrank and Wilcoxon-Gehan tests for censored survival data, normal scores test, and trend test with equally spaced scores. Use *StatXact* to help determine power and sample size for Fisher's Exact test, trend tests on k binomial samples, and linear rank tests on two multinomial samples. The manual incorporates many excellent examples from the literature and is a textbook in itself. For IBM-PC compatible microcomputers from Cytel Software, 137 Erie St, Cambridge MA 02139. [cytel.com](http://cytel.com) 617/661-2011. Also available as an add-on module for SAS and SPSS.

*LogXact* performs exact logistic regressions as described in Cox and Shell [1989]. (For IBM-PC compatible microcomputers from Cytel Software, 137 Erie St, Cambridge MA 02139. 617/661-2011.)

*RT* performs permutation tests on one- and two-samples (though fewer than for StatXact), plus analysis of variance, regression analysis, matrix randomization tests, tests on spatial data, time series analysis, and multivariate analysis using Wilk's lambda statistic and Romesburg's sum of squares statistic E. Applications are drawn from Manly [1996]. For IBM-PC compatible microcomputers from West, 1406 South Greeley Highway, Cheyenn WY 82007. 307/634-1756.

*Statistical Calculator* (SC) is an extensible statistical environment, supplied with over 1200 built-in (compiled C) and external (written in SC's C-like language) routines. Permutation-based methods for contingency tables (chi-square, likelihood, Kendall S, Theil U, kappa, tau, odds ratio), one and two-sample inference for both means and variances, correlation, and multivariate analysis (MV runs test, Boyett-Schuster, Hoetelling's T are available). Also includes ready-made bootstrap routines for testing homoscedacity, detecting multimodality, plus general bootstrapping and jack-knifing facilities. For PC compatibles, Unix, and T800 transputer. Tony Dusoir, Mole Software, 23 Cable Rd., Whitehead, Co. Antrim BT38 9PZ, N. Ireland; e-mail: fbgj23@ujvax.ulster.ac.uk. phone/fax: (+44) (0)960 378988.

*Testimate* lacks the pull-down menus and prompts that Windows users have grown to rely on. Too few statistical routines for its price tag. PC compatibles. idv: Datenanalyse und Versuchsplanung, Wessobrunner Strabe 6, D-82131 Gauting, Munich, Germany; phone: 89.850.8001; also, SciTech International, Inc., 2231 North Clybourn Avenue, Chicago, IL 60614-3011, USA; phone: 800/622-3345.

## 12.3. Preparing Manuscripts for Publication

You've laid the groundwork. You've done the experiment. You've completed the analysis. Five simple rules can help you prepare your article for publication.

1. State the test statistic explicitly. Reproduce the formulae. If you cite a text, for example, Good [1993], please include the page number(s) on which the statistic you are using is defined.
2. State your assumptions. Are your observations independent? exchangeable? Is the underlying distribution symmetric? Contrary to statements that have appeared in several recent journal articles—we withhold the names to protect the guilty—permutation tests cannot be employed without one or both of these essential assumptions. See Draper et al [1993], Gastwirth and Rubin [1971], and Hettmansperger [1984] for discussions of this point.
3. State which labels you are rearranging. Provide enough detail that any interested reader can readily reproduce your results. In other words, report your statistical procedures in the same detail you report your other experimental and survey methodologies.

4. State whether you are using a one-tailed or a two-tailed test. See Section 6.1.1 for help in making a decision.
5. a) If you detect a statistically significant effect, then provide a confidence interval (see Section 3.2). Remember, an effect can be statistically significant without being of practical or biological significance.  
b) If you do not detect a statistically significant effect, could a larger sample or a more sensitive experiment have detected one? Consider reporting the power of your test. See Sections 2.2 and 13.9.

## CHAPTER 13

# Increasing Computational Efficiency

### 13.1. Seven Techniques

With today's high-speed computers, drawing large numbers of subsamples with replacement (the bootstrap) or without (the permutation test) is no longer a problem, unless or until the entire world begins computing resampling tests. To prepare for this eventuality, and because computational efficiency is essential in the search for more powerful tests, a primary focus of research in resampling today is the development of algorithms for rapid computation.

There are seven main computational approaches, several of which may be and usually are employed in tandem, as follows:

1. the *Monte Carlo*, in which a sample of the possible rearrangements is drawn at random and these samples are used in place of the complete permutation distribution;
2. *rapid enumeration and selection algorithms*, whose object is to provide a rapid transition from one rearrangement to the next;
3. *recursive relationships* reduce the number of computations;
4. *branch and bound algorithms* that eliminate the need to evaluate each individual rearrangement;
5. *Gibbs sampling*;
6. solutions through *characteristic functions and fast Fourier transforms*;
7. *asymptotic approximations*, for use with sufficiently large samples.

In the following sections, we consider each of these approaches in turn.

### 13.2. Monte Carlo

Instead of examining all possible rearrangements, we can substantially reduce the computations required by examining only a small but representative random sample [Dwass, 1957; Barnard, 1963]. In this process, termed a Monte Carlo, we proceed in stages: 1) We rearrange the data at random; 2) we compute the test

statistic for the rearranged data and compare its value with that of the statistic for the original sample; and 3) we apply a stopping rule to determine whether we should continue sampling, or whether we are already in a position to accept or reject.

The program fragments reproduced in Chapters 3 through 5 of this text use the Monte Carlo approach. In the, not necessarily optimal, computer algorithm introduced in those chapters, all observations in all subsamples are loaded into a single linear vector  $X = \{X[0], X[1], \dots, X[N-1]\}$ . Then, a random number is chosen from the set of integers  $0, 1, \dots, I$  with  $I = N-1$  initially. If the number we choose is  $i$ ,  $X[i]$  is swapped with  $X[I]$  in a three-step process:

```
temp := X[i];
X[i] := X[N - 1];
X[N - 1] := temp;
```

and  $I$  is decremented. This process is repeated until we have rearranged the desired number of observations and are ready to compute the test statistic for the new rearrangement.

We don't always need to reselect all  $N$  observations. For example, in a two-sample comparison of means, with  $N = n + m$ , our test statistic only makes use of the last  $m$  observations. Consequently, we only need to choose  $m$  random numbers each time.

After we obtain the new value of the test statistic, we compare it with the value obtained for the original data. We continue until we have examined  $N$  random rearrangements and  $N$  values of the test statistic. Typically,  $N$  is assigned a value between 100 and 1600, depending on the precision that is desired (see Section 13.2.2 and Marriott, 1979). Through the use of a Monte Carlo, even the most complicated multivariate experimental design can be analyzed in less than a minute on a desktop computer.

### 13.2.1. Stopping Rules

If a simple accept/reject decision is required, we needn't perform all  $N$  calculations, but can stop as soon as it is obvious that we must accept or reject the hypothesis at a specific level. In practice, I use a one-sided stopping rule based on the 10% level. Suppose in the first  $n$  rearrangements we observe a fraction  $H_n$  with a value of the test statistic that is as or more extreme than the value for the original observations. If  $H_n > 0.1N$ , then we accept the hypothesis at the 10% level. Otherwise, we continue until  $n = N$  and report the exact percentage of rejections. Besag and Clifford [1991] and Lock [1991] describe two-sided sequential procedures in which the decision to accept, reject, or continue is made after each rearrangement is examined.

### 13.2.2. Variance of the Result

The resultant estimated significance level  $\hat{p}$  is actually a binomial random variable  $B(N, p)$ , where  $N$  is the number of random rearrangements and  $p$  is the true but still

unknown value of the significance level. The variance of  $\hat{p}$  is  $p(1-p)/N$ . If  $p$  is 10%, then using a sample of 81 randomly selected rearrangements provides a standard deviation for  $\hat{p}$  of 1%. A sample of 364 reduces the standard deviation to 0.25%.

The use of a variable in place of a fixed significance level results in a minor reduction in the power of the test particularly with near alternatives [Dwass, 1957]. In most cases, this reduction does not appear to be of any practical significance; see Vadiveloo [1983], Jockel [1986], Bailer [1989], Edgington [1987], and Noreen [1989].

In a Monte Carlo variant called *importance sampling*, the rearrangements are drawn with weights chosen so as to minimize the variance. In some instances, when combined with branch and bound techniques, as in Mehta, Patel, and Senchaudhuri [1988], importance sampling can markedly reduce the number of samples that are required. See also Besag and Clifford [1989].

### 13.2.3. Cutting the Computation Time

The generation of random rearrangements creates its own set of computational problems.

Each time a data element is selected for use in the test statistic, two computations are required: 1) A random number is selected and 2) two elements in the combined sample are swapped.

The ideal futuristic computer will have a built-in random number generator—for example, it might contain a small quantity of a radioactive isotope, with the random intervals between decays producing a steady stream of random numbers. This futuristic computer might also have a butterfly network that would randomly swap 10 or 100 elements of an array in a single pass.

Today, in the absence of such technology, any improvements in computation speed must be brought about through software. Little direct research has been done in the area, although recently Baglivo et al [1992] reported on techniques for doing many of the repetitive computations in parallel. I did some preliminary work in which I considered a sort of drunkard's walk through the set of rearrangements: The first rearrangement was chosen at random; thereafter, the program stumbled from rearrangement to rearrangement swapping exactly two data elements at random each time. The results were disappointing. Any savings in computation time per rearrangement were more than offset by the need to sample four or five times as many rearrangements to achieve the same precision in the result. I did achieve a substantial increase in efficiency by selecting several, separated random bits from each random number.

## 13.3. Rapid Enumeration and Selection Algorithms

If we are systematic and proceed in an orderly fashion from one rearrangement to the next, we can substantially reduce the time required to examine a series of

rearrangements. The literature on the generation of permutations is so extensive that we include a separate bibliography (C) on this and related computational topics. Optimal algorithms for generating sequences of rearrangements are advanced by Walsh [C1957], Boothroyd [C1967], Plackett [C1968], Yanagimoto and Okamoto [1969], Boulton [C1974], Hancock [C1974], Bitner, Ehrlich, Rheingold [C1976], Akl [C1981], and Bissell [C1986]. See, for example, the review by Wright [C1984]. Recent minimal change algorithms include those of Berry [C1982], Lam and Sotchen [C1982], Nigam and Gupta [1984], and Marsh [C1987].

### 13.3.1. Matched Pairs

Sometimes we can reduce the number of computations that are required by taking advantage of the way we label or identify individual permutations. In the case of paired comparisons, we readily enumerate each possible combination by running through the binary numbers from 0 to  $2^n - 1$ , letting the zero's and one's in each number (obtained via successive right shifts, a single machine-language instruction in most computers) correspond to positive and negative paired differences, respectively.

The shift algorithm, introduced in this context by Baker and Tilbury [C1993] for use with discrete data, avoids the need for assembler level programming. The

test statistic  $T_k \sum_{i=1}^k |X_i| - \sum_{i=1}^k X_i$  is calculated one variable at time, and an array of counters or bar chart  $N[ ]$  incremented appropriately. At step 0,  $T_0 = 0$ ; we initialize the array of counters, so that  $N[0] = 1$  and all other elements are zero. At step 1, we add  $X_1$  to  $T_0$ ; as  $X_1$  could be either positive or negative, we increment  $N[|X_1|]$  by  $N[0]$  so that both  $N[0]$  and  $N[|X_1|]$  are now equal to 1. At step 2, we add  $X_2$  to  $T_1$  and increment  $N[|X_1| + |X_2|]$  by  $N[|X_1|]$  and  $N[|X_2|]$  by  $N[0]$ . Note that if  $X_1 = X_2$ ,  $N[|X_1|] = 2$ . We continue in this fashion, so that at the  $k$ th step, we increment  $N[j]$  by  $N[j - |X_k|]$  for  $j = \sum_{i=1}^k |X_i|, \dots, |X_k|$ .

Censoring actually reduces the time required for enumeration. For if there are  $n_c$  censored pairs, then enumeration need only extend over the  $2^{n-n_c}$  values that might be assumed by the uncensored pairs. In computing the GAMP test for paired comparisons, it is easy to see that

$$\Pr\{U' \geq U \text{ AND } S' \geq S\} = \Pr\{U' \geq U\} * \Pr\{S' \geq S\};$$

$$\Pr\{U' \geq U\} = \frac{1}{2^{U+L}} \sum_{k=U}^{U+L} \binom{U+L}{k}.$$

The remaining probability,  $\Pr\{S' \geq S\}$ , may be obtained by enumeration and inspection.

## 13.4. Recursive Relationships

Although tables for determining the significance level of Fisher's Exact test are available, in Finney [1948] and Latscha [1953], for example, these are restricted to a few discrete p-values. Today, it is usually much faster to compute a significance level than it is to look it up in tables. Beginning with Leslie [1955], much of the subsequent research on Fisher's Exact test has been devoted to developing algorithms that would speed or reduce the number of computations required to obtain a significance level.

As one rapid alternative to equation (6.1), we may use the recursive relationship provided by Feldman and Kluger [1963]: With table entries  $(a_o, b_o, c_o, d_o)$ , define

$$p = \frac{(a_o + b_o)!(a_o + c_o)!(d_o + b_o)!(d_o + c_o)!}{N!a_o!b_o!c_o!d_o!}.$$

It is easy to see that

$$p_{i+1} = \frac{a_i d_i}{b_{i+1} c_{i+1}} p_i,$$

where  $a_i = a_o - i$ .

We may speed the computations of the statistics for unordered  $r \times c$  contingency tables considered in Section 6.4 by noting that  $Q$  is invariant under permutations that leave the marginals intact. Thus, we may neglect the numerator  $Q$  in calculating the permutation distribution and focus on the denominator  $R$  March [C1972].

We may use a recursive algorithm developed by Gail and Mantel [C1977] to speed the computations for  $r \times 2$  contingency tables. If  $N_i(\mathbf{f}_{.1}; f_1, f_2, \dots, f_n)$  denotes the number of tables with the indicated marginals, then

$$N_{i+1}(\mathbf{f}_{.1}; f_1, f_2, \dots, f_n) = \sum_j N_i(\mathbf{f}_{.1} - j; f_1, f_2, \dots, f_n).$$

The algorithms we developed in Chapters 3 and 4 are much too slow, since they treat each observation as an individual value.

Algorithms for speeding the computations of the Freeman–Halton statistic in the general  $r \times c$  case are given in March [C1972], Gail and Mantel [C1977], Mehta and Patel [C1983, 1986a, 1986b], and Pagano and Halvorsen [C1981]. Details of the Mehta and Patel approach are given in Section 13.4. An efficient method for generating  $r \times c$  tables with given row and column totals is provided by Patefield [1981]. See also Agresti, Wackerly and Boyett [1979] and Streitberg and Rohmed [C1986].

The power of the Freeman–Halton statistic in the  $r \times 2$  case is studied by Krewski, Brennan, and Bickis [1984].

### 13.5. Focus on the Tails

We can avoid examining all  $N!$  rearrangements, if we focus on the tails, using the internal logic of the problem to deduce the number of rearrangements that yield values of the test statistic that are as or more extreme than the original.

Consider the shift algorithm introduced in the preceding section. Suppose that  $T^o$  is the test statistic for the original data; as  $T_k$  is nondecreasing, we only need to keep track of individual values of  $T_k$  that are less than  $T^o$ . Our modified procedure at the  $k$ th step is as follows:

if  $\sum_{i=1}^k |X_i| < T^o$  increment  $N[j]$  by  $N[j - |X_k|]$  for  $j = \sum_{i=1}^k |X_i|, \dots, |X_k|$ ;

otherwise, set  $N[T^o] = 2N[T^o]$ , then increment  $N[j]$  by  $N[j - |X_k|]$  for  $j = T^o, \dots, |X_k|$ .

Of course, if  $N[T^o] > \alpha n$ , we would terminate the procedure and accept the null hypothesis.

Green [C1977] was the first to suggest a branch and bound method for use in two-sample tests and correlation. Our description of Green's method is based on De Cani [C1979].

In the two-sample comparison described in Section 3.2, suppose our test statistic  $T = \Sigma x_{\pi(i)}$  and that the observed value is  $T_o$ . We seek  $P(T \geq T_o)$ , the probability under the null hypothesis that a random value of  $T$  equals or exceeds  $T_o$ .

Assume that the combined observations are arranged in descending order  $X_{(1)} \geq X_{(2)} \geq \dots \geq X_{(N)}$ . To simplify the notation, let  $Z_i$  denote the  $i$ th order statistic  $X_{(i)}$ . If the labels (subscripts) on the  $X$ 's really are irrelevant (as they would be under the null hypothesis) then  $T$  can be regarded as a random sample of  $m$  of the observations selected at random without replacement from the  $\{Z_i\}$ .

Suppose we have selected  $k$  such values,  $Z_{l_1}, \dots, Z_{l_k}$ ,  $k < m$ . The maximum attainable value of  $T$  is obtained by adding to  $Z_{l_1} + \dots + Z_{l_k}$  the  $m - k$  largest of the  $N - k$  remaining elements. Call this maximum  $T(l_1, \dots, l_k)$ . Similarly, the minimum attainable value of  $T$  is obtained by adding to  $Z_{l_1} + \dots + Z_{l_k}$  the  $m - k$  smallest of the  $N - k$  remaining elements. Call this minimum  $t(l_1, \dots, l_k)$ . Given  $l_1, \dots, l_k$ , we can bound  $T$ :

$$t(l_1, \dots, l_k) \leq T \leq T(l_1, \dots, l_k).$$

There are  $\binom{N-k}{m-k}$  sets of  $m$  elements of  $Z$  whose totals lie between the given bounds.

If  $t(l_1, \dots, l_k) \geq T_0$ , then

$$P(T \geq T_0) \geq \binom{N-k}{m-k} / \binom{N}{m}.$$

If  $T_0 > T(l_1, \dots, l_k)$ , then

$$P(T \geq T_0) \leq 1 - \binom{N-k}{m-k} / \binom{N}{m}.$$

If  $T_0$  lies between the bounds, or if we require an improved bound on  $P(T \geq T_0)$ , then we can add a  $k+1$ th element to the index set.

Our results apply equally to any test statistic of the form  $\sum_{i=1}^m f[x_{\pi(i)}]$ , where  $f$  is a monotone increasing function. Examples of such monotone functions include the logarithm (when applied to positive values), ranks, and any of the other robust transformations described in Chapter 9.

### 13.5.1. Contingency Tables

A large number of authors have joined in the search for a more rapid method for enumerating the tail probabilities for Fisher's Exact test, including Leslie [1955], Feldman and Kluger [1963], Good [1976], Gail and Mantel [1977], Pagano and Halvorsen [1981], and Patefield [1981]. See, for example, the review by Agresti [1993]. A quantum leap toward a more rapid method took place with the publication of the network approach of Mehta and Patel [1980]. Their approach is widely applicable, as we shall see below. It has three principal steps:

1. representation of each contingency table as a path through a directed acyclic network with nodes and arcs;
2. an algorithm with which to enumerate the paths in the tail of the distribution without tracing more than a small fraction of those paths;
3. determination of the smallest and largest path lengths at each node.

Only the last of these steps is application specific.

Network algorithms have been developed for all of the following:

$2 \times c$  contingency tables [Mehta and Patel, 1980];

$r \times c$  contingency tables [Mehta and Patel, 1983];

the common odds ratio in several  $2 \times 2$  contingency tables [Mehta, Patel, and Gray, 1985];

logistic regression [Hirji, Mehta, and Patel, 1987];

restricted clinical trials [Mehta, Patel, and Wei, 1988];

linear rank tests and the Mantel-Haenszel trend test [Mehta, Patel, and Senchaudhuri, 1988].

For simplicity, we focus in what follows on the  $2 \times c$  contingency table.

#### 13.5.1.1. Network Representation

Define the reference set  $\Gamma$  to be all possible  $2 \times k$  contingency tables (see Chapter 6) with row marginals  $(m, n)$  and column marginals  $(t_1, t_2, \dots, t_k)$ . Thus each

table,  $x \in \Gamma$ , is of the form

$$\begin{array}{ccccc} x_1 & x_2 & \dots & x_k & m \\ x'_1 & x'_2 & \dots & x'_k & n \\ t_1 & t_2 & \dots & t_k & N. \end{array}$$

For each table  $x \in \Gamma$ , we may define a discrepancy measure

$$d(x) = \sum_{i=1}^k a_i(m_{i-1}, x_i)$$

and a probability

$$h(x) = C^{-1} \prod_{i=1}^k \lambda_i(m_{i-1}, x_i),$$

where the partial sum  $m_j = \sum_{i=1}^j x_i$ , and the normalizing constant  $C = \sum_{x \in \Gamma} \prod_{i=1}^k \lambda_i(m_{i-1}, x_i)$ . Important special cases of  $d(x)$  and  $h(x)$  are

$$d(x) = \prod_{i=1}^k a_i x_i$$

for linear rank tests and

$$h(x) = \prod_{i=1}^k \binom{t_i}{x_i} / \binom{N}{m}$$

for unordered contingency tables.

As in Section 6.3, our object is to compute the one-sided significance level  $p = \sum_R h(x)$ , where  $R$  is the set on which  $d(X) \geq d_0$ .

First, we represent  $\Gamma$  as a directed acyclic network of nodes and arcs. Following Mehta and Patel [1983], the network is constructed recursively in  $k + 1$  stages labeled  $0, 1, 2, \dots, k$ . The nodes at the  $j$ th stage are ordered pairs  $(j, m_j)$  whose first element is  $j$  and whose second is the partial sum of the frequencies in the first  $j$  categories of the first row. If there is a total of two observations in the first category, then there will be three nodes at the first stage  $(1, 0), (1, 1), (1, 2)$ —corresponding to the three possible distributions of elements in this category.

Arcs emanate from the node  $(j, m_j)$ ; each arc is connected to exactly one successor node. Each path linking  $(0, 0)$  with the terminal node  $(k, m)$  corresponds to a unique contingency table. For example, the path

$$(0, 0) \rightarrow (1, 0) \rightarrow (2, 2) \rightarrow (3, 4) \rightarrow (4, 4)$$

corresponds to the table

0	2	2	0	4
2	0	0	2	4
2	2	2	2	

The total number of paths in the network corresponds to the total number of tables. We could count the total number of tables by tracing each of the individual paths, but we can do better.

### 13.5.1.2. The Network Algorithm

Our goal in network terms is to quickly identify and sum all paths whose lengths do not exceed  $d \cdot h$ : for the original unpermuted table. Let  $\Gamma_j = \Gamma(j, m_j)$  denote the set of all paths from any node  $(j, m_j)$  to the terminal node  $(k, m)$ . In other words,  $\Gamma_j$  represents all possible completions of those tables in  $\Gamma$  for which the sum of the first  $j$  cells of row 1 is  $m_j$ . Define the shortest path length

$$SP(j, m_j) = \min_{x \in \Gamma_j} \sum_{i=j+1}^k a_i(m_{i-1}, x_i)$$

and the longest path length

$$LP(j, m_j) = \max_{x \in \Gamma_j} \sum_{i=j+1}^k a_i(m_{i-1}, x_i).$$

Let  $L(PAST)$  denote the length of a path from  $(0, 0)$  to  $(j, m_j)$ . If this path is such that

$$L(PAST) + LP(j, m_j) \leq d \cdot h,$$

then all similar subpaths from  $(0, 0)$  to  $(j, m_j)$  of equal or smaller length contribute to the  $p$  value. This number can be determined by induction—the details depend on the actual form of  $d$  and  $h$ , and thus we need not enumerate the tables explicitly. If this path is such that

$$L(PAST) + SP(j, m_j) \geq d \cdot h,$$

then we can ignore it and all similar paths of equal or greater length—again, without actually enumerating them.

If the path satisfies neither condition, then we extend it to a node at the  $j + 1$ th stage, compute the new shortest and longest path lengths, and repeat the calculation.

The shortest and longest path lengths may be determined by dynamic programming in a single backward pass through the network. Dynamic programming is used by Mehta and Patel [1980] in their first seminal paper. Their original approach can be improved upon in three ways:

- 1) by taking advantage of the structure of the problem;
- 2) by a Monte Carlo, randomly selecting the successor node at each stage;

- 3) by a Monte Carlo utilizing importance sampling, that is, weighting the probabilities with which an available node is selected so as to reduce the variance of the resultant estimate of  $p$ .

The three approaches can be combined: A highly efficient two-pass algorithm for importance sampling using backward induction followed by forward induction was developed by Mehta, Patel, and Senchaudhuri [1988]. Their new algorithm guarantees that all rearrangements sampled will lie inside the critical region. A result of Joe [1988] also represents a substantial increase in computational efficiency.

### 13.6. Gibbs Sampling and a Drunkard's Walk

Suppose we have a  $2 \times 2$  table with entries  $f_{11}, f_{12}, f_{21}, f_{22}$ ; assume this results from a sequence of random variables  $X_0, Y_0, X_1, Y_1, \dots$ , each taking the value 0 or 1, where the estimated conditional probabilities of  $Y|X$  and  $X|Y$  can be expressed in the two matrices

$$A_{y|x} = \begin{pmatrix} \frac{f_{11}}{f_{11}+f_{21}} & \frac{f_{21}}{f_{11}+f_{21}} \\ \frac{f_{12}}{f_{12}+f_{22}} & \frac{f_{22}}{f_{12}+f_{22}} \end{pmatrix}, \quad A_{x|y} = \begin{pmatrix} \frac{f_{11}}{f_{11}+f_{12}} & \frac{f_{12}}{f_{11}+f_{12}} \\ \frac{f_{21}}{f_{21}+f_{22}} & \frac{f_{22}}{f_{21}+f_{22}} \end{pmatrix}.$$

Using these matrices, generate a single couple  $y, x$ . Modify the table (preserving the marginals) to provide for this new entry; if it is not possible to preserve the marginals, do not modify the table. Compute the test statistic, and compare with the original value of the test statistic. Modify the transition matrices to reflect the change, and repeat the procedure.

A similar procedure but one guaranteed to converge to the correct result is based on the Gibbs sampler, a technique for generating random variables from a (marginal) distribution indirectly, without having to calculate the density (see Casella and George, 1992). At each step, we draw from the hypergeometric distribution produced by taking a binomial ( $p, f_{11} + f_{21}$ ) and an independent binomial ( $p, f_{12} + f_{22}$ ) conditional on the sum of the two binomials being  $f_{11} + f_{12}$ . To obtain a new table, we let the computer pick a uniformly distributed random number between 0 and 1, and evaluate the hypergeometric quantile at this number. Methods for its rapid calculation are described in Kolassa and Tanner [1994].

By taking advantage of a second result of Kolassa and Tanner [1994], we can extend the preceding to contingency tables with  $r$  rows and  $c$  columns.

Let  $\{f_{ij}\}$  be an  $r \times c$  contingency table with independent entries.

If  $i < r$  and  $j < c$ , then the distribution of element  $f_{ij}$  conditional on all other elements except for those in the last row and column, and conditional on all marginals, is the same as the distribution of the first element in the  $2 \times 2$  table with elements  $f_{ij}, f_{ic}, f_{rj}, f_{rc}$  conditional on all marginals.

Thus, we may proceed from cell to cell, drawing Gibbs samples as described above. Ambiguities arise in how we are to balance the marginals. We may balance cell by cell as we go or we may keep a running tabulation and balance only when sampling is complete. The overage or discrepancy may be assigned in a number

of ways: To an adjacent cell, to a following cell, to a cell in the same column or row chosen at random, or to a cell further along in the same column or row chosen at random.

## 13.7. Characteristic Functions

As the sample size increases, the number of possible rearrangements increases exponentially. For example, in the one-sample test of a location parameter based on  $n$  observations, there are  $2^n$  possible rearrangements. When finding the permutation distribution of a statistic that is a linear combination of some function of the original observations, Pagano and Tritchler [1983] show we can reduce the computation time from  $C_1 2^n$  to  $C_2 n^c$  where  $c$  is, we hope, much less than  $n$ .

Their technique requires two steps: In the first, they determine the characteristic function of the permutation distribution through a set of difference equations. This step requires  $2Qm(m + n)$  complex multiplications and additions to find the characteristic function at  $Q$  points. In the second, they use the basic theorem in Fourier series to invert the characteristic function and determine or approximate the permutation distribution at  $U < Q$  different points. This step requires  $2Q \log Q$  calculations.  $Q$  is normally chosen to be a power of two (e.g., 256 or 512) so that one can take advantage of a fast Fourier transform; the exact number will depend on the precision with which one wants to estimate the significance level.

This method is chiefly of historic interest; branch and bound algorithms offer greater computational efficiency, particularly when coupled with importance sampling. Vollset, Hirji, and Elashoff [1991] found that the fast Fourier transform method can result in considerable loss of numerical accuracy.

## 13.8. Asymptotic Approximations

### 13.8.1. A Central Limit Theorem

The fundamental asymptotic result for the permutation distribution of the two-sample test statistic for a location parameter was first stated by Madow [1948] and formalized by Hoeffding [1951, 1952], who demonstrates convergence of the distribution of the Studentized test statistic under the alternative as well as under the null hypothesis.

Let  $T_n = T(X_{(1)}, \dots, X_{(n)})$  be the test statistic, and let  $\mu_n$  and  $\sigma_n^2$  be its first and second moments, respectively. Then the permutation distribution  $F_n$  of  $Z_n = \frac{T_n - \mu_n}{\sigma_n}$  obtained by randomly rearranging the subscripts of the arguments of  $T_n$  converges to  $\Phi$ , the Gaussian (normal) distribution function.

This result means that for sufficiently large samples, we can give our computers a rest, at least temporarily, and approximate the desired  $p$ -value with the aid of tables of the normal distribution. To use these tables, we need to know the first and second moments of the permutation distribution. Occasionally, with samples

of moderate size, we may also need to know and use the third and higher moments in order to obtain an accurate approximation. Moments for the randomized block design are given by Pitman [1937] and Welch [1937]; for the Latin Square by Welch [1937]; for the balanced incomplete block by Mitra [1961]; and for the completely randomized design by Robinson [1983], and Bradbury [1988].

Extensions to, and refinements of, Hoeffding's work are provided by Silvey [1954, 1956], Dwass [1955], Motoo [1957], Erdos and Renyi [1959], Hajek [1960, 1961], and Kolchin and Christyakov [1973]. Asymptotic results for rank tests are given in Jogdeo [1968] and Tardif [1981]. For further details of the practical application of asymptotic approximations to the analysis of complex experimental designs, see Lehmann [1986], Kempthorne, Zyskind, Addelman, Throckmorton, and White [1961], and Ogawa [1963].

### 13.8.2. Edgeworth Expansions

While the Gaussian distribution may provide a valid approximation to the *center* of the permutation distribution, it is the tails (and the *p*-values of the tails) with which we are primarily concerned. Edgeworth expansions give good approximations to the tails in many cases. Edgeworth expansions for the distribution function under both the alternative and the null hypothesis have been obtained by Albers, Bickel, and Van Zwet [1976], Bickel and Van Zwet [1978], Robinson [1978], and John and Robinson [1983].

Saddlepoint methods and large deviation results give still better approximations in the tails. Saddlepoint approximations for the one- and two-sample tests of location as suggested by Daniels [1955, 1958] are derived by Robinson [1982]. Saddlepoint approximations for use with general linear models for both the permutation distribution and the bootstrap are given by Booth and Butler [1990].

### 13.8.3. Generalized Correlation

Test statistics for location parameters are almost always linear or first-order functions of the observations. By contrast, test statistics for scale parameters, the chi-square statistic, and the Mantel–Valand statistic for generalized correlation are quadratic or second-order functions of the observations. Their limiting distributions are not Gaussian but chi-square or a Pearson Type III distribution [Berry and Mielke, 1984, 1986; Mielke and Berry, 1985]. Other asymptotic approximations for second-order statistics are given by Shapiro and Hubert [1979], O'Reilly and Mielke [1980], and Ascher and Bailar [1982].

## 13.9. Confidence Intervals

The trial-and-error method of determining confidence intervals described in Section 3.2 is time consuming and confusing and entails a seemingly unending

number of calculations. The stepwise approach suggested by Garthwaite [1996] is both systematic and efficient with the need for only a single permutation at each step.

Let  $T_o(U)$  be the value of the statistic used to test the hypothesis  $\theta = U$ , obtained for the actual sample data. Let  $T_i(U)$  be the value of the statistic obtained for a random permutation  $\pi_i$  of the data.

Observing  $T(U_i)$ , we update our estimate of the upper limit at the  $i$ th step as follows:

If  $T_i(U) > T_o(U)$ , set  $U_{i+1} = U_i - c\alpha/i$ .

Otherwise, set  $U_{i+1} = U_i + c\alpha/i$ ,

where  $\alpha$  is the significance level and  $c$  is known as the step-length constant.

We continue in this fashion generating exactly one new permutation and evaluating and comparing exactly two values of the test statistic at each step. It is easy to see that the process converges.

$$c = k(U_1 - \hat{\theta}), \quad \text{where } k = \frac{2}{2\pi z_\alpha - \frac{1}{2} \exp(-z_\alpha^2/2)}.$$

One possible starting guess is  $U_1 = \hat{\theta} + z_\alpha s$ , where  $s$  is the sample standard deviation.

## 13.10. Sample Size and Power

Suppose we are in the design stages of a study and we intend to use a permutation test for the analysis. How large should our sample sizes be? Our answer will depend on three things:

- o the alternative(s) of interest;
- o the power desired at these alternatives;
- o the significance level.

A not unrelated question arises if we conclude an analysis by accepting the null hypothesis. Does this mean the alternative is false or that we simply did not have a large enough sample to detect the deviation from the null hypothesis? Again, we must compute the power of the test for several alternatives before we are able to reach a decision.

### 13.10.1. Simulations

One way to estimate the power is by drawing a series of  $K$  (simulated) random samples from a distribution similar to that which would hold under the alternative. For each sample, we perform the permutation test at the stated significance level and record whether we accept or reject the null hypothesis. The proportion of rejections becomes our estimate of the power of the test.

This proportion is a random variable with the binomial distribution with  $K$  trials and a probability  $\beta$  of success in each trial, where  $\beta$  is the (unknown) power of the test to be estimated.

When designing a study, I use  $K = 100$  until I am ready to fine tune the sample size, when I switch to  $K = 400$ . I also study (estimate) the power for at least two distinct alternatives.

For example, when testing the hypothesis that the observations are normal with mean 0 against the alternative that they have a mean of at least 1, one might sample from alternatives with at least two different variances, say, one with variance equal to unity and one with variance equal to 2, where 1 is our best guess of the unknown variance, and 2 is a worst-case possibility.

When doing an after-the-fact analysis of the power, use estimates of the parameters based on the actual data. If the pooled sample variance is 1.5, then use a best guess of 1.5 and a worst case of 3 or even 4. (The use of a single estimate alone would be misleading; see Zumbo and Hubley, 1998). The final result may require  $8KN$  computations, where  $N$  is the average number of resamplings required each time we perform the test.

With such a large number of calculations, it is essential that I take advantage of one or more of the computational procedures described in Sections 2 through 6 of this chapter. Oden [C1991] offers several recommendations. Gabriel and Hsu [1983] describe an application-specific method for reducing the number of computations required to estimate the power and determine the appropriate sample size.

### 13.10.2. Network Algorithms

The same network algorithms that we used to determine significance level can also be used to calculate power, providing we can determine the probability of each specific permutation under the alternative; see, for example, Hilton and Mehta [1993], and Mehta, Patel, and Senchaudhuri [1998]. For example, Mehta, Patel, and Senchaudhuri [1998] studied the Cochran–Armitage test for trend, as described in Section 3.5.2.1, for which the test statistic is  $T = \sum_{j=1}^J d_j X_j$ , where random variable  $X_j$  denotes the integer number of responders among the  $n_j$  subjects treated at dose  $d_j$ , and assumes the value  $x_j$ . The reference set for permutations is

$$\Gamma_m = \left\{ x : \sum_{j=1}^J x_j = m \right\},$$

and its critical region is

$$\Gamma_m(t) = \left\{ x \in \Gamma_m : \sum_{j=1}^J d_j x_j \geq t \right\}.$$

For a given significance level  $\alpha$ , let  $t_\alpha(m)$  be the smallest possible cutoff value such that  $\Pr\{T \geq t_\alpha(m) | m, H_0\} \leq \alpha$ . This cutoff is data dependent through the

marginal  $m$  and  $\beta(m) = \Pr\{T \geq t_\alpha(m) | m, K\}$ . So that the unconditional power is  $\sum \beta(m) \Pr\{m | K\}$ .

Again, we represent the permutation reference set as a network of nodes and arcs constructed in  $J + 1$  stages. At any stage  $j$ , the network contains a set of nodes of the form  $(j, m_j)$ , where  $j$  represents the  $j$ th of the  $J$  binomial populations and  $m_j$  is one possible value of the partial sum of responses from the first  $j$  populations. Arcs emanate from each node, and each such arc is connected to a successor node  $(j + 1, m_{j+1})$  at stage  $j + 1$ . When the network is complete, it will terminate with single node  $(J, m)$ , and each path from  $(0, 0)$  to  $(J, m)$  represents one and only one response vector (permutation, rearrangement) in  $\Gamma_m$ .

The arc connecting the nodes  $(j, m_j)$  with its successor is assigned a rank length based on the Cochran–Armitage statistic,  $r_{j+1} = d_{j+1}(m_{j+1} - m_j)$ , and two probability lengths  $p_{H,j+1}$  and  $p_{K,j+1}$  based on their likelihoods under the hypothesis and alternative, respectively. By specifying a path through the network, we automatically know the corresponding response vector  $x$ , its test statistic  $t(x) = \sum r_j$ , and its unnormalized probability under the null hypothesis and alternative, respectively. Any method we use to generate an estimate of the significance level will provide us with an estimate of the power at the same time.

## 13.11. Some Conclusions

In the Monte Carlo, we compute the test statistic for a sample of the possible rearrangements, and use the resultant sampling distribution and its percentiles in place of the actual permutation distribution and its percentiles. The drawback of this approach is that the resultant significance level  $p'$  may differ from the significance level  $p$  of a test based on the entire permutation distribution.  $p'$  is a consistent estimate of  $p$  with a standard deviation on the order of  $Np(1 - p)$  where  $n$  is the number of rearrangements considered in the Monte Carlo.

In the original Monte Carlo, the rearrangements are drawn with equal probability. In a variant called *importance sampling*, the rearrangements are drawn with weights chosen so as to minimize the variance. In some instances, when combined with branch and bound techniques as in Mehta, Patel, and Senchaudhuri [1988], importance sampling can markedly reduce the number of samples that are required. (See also Besag and Clifford [1989].)

A second drawback of the Monte Carlo is that selecting a random arrangement is itself a time-consuming operation that can take several multiples of the time required to compute the sample statistic. A current research focus is on rapid enumeration and selection algorithms that can provide a fast transition from one rearrangement to the next. To date, all solutions have been highly application-specific.

Branch and bound algorithms eliminate the need to evaluate each rearrangement individually. The network approach advanced by Mehta and Patel can cut computation time by several orders of magnitude.

Solutions through characteristic functions are seldom of practical interest. When subsamples are large—and it is the size of the subsample or block, not the sample

as a whole, that is the determining factor—an asymptotic approximation should be considered. In my experience as an industrial statistician with the pharmaceutical and energy industries, the opportunity to take advantage of an asymptotic approximation seldom arises. In preclinical work, one seldom has enough observations. And in a clinical trial, though the sample size is large initially, one is usually forced to divide the sample again and again to correct for covariates. In practice, contingency tables always have one or two empty cells. The errors in significance level that can result from an inappropriate application of an asymptotic approximation are amply illustrated in Table 6.4.

If you are one of the favored few able to take advantage of an asymptotic approximation, you first will need to compute the mean and variance of the permutation distribution. In some cases, you will also need to calculate and use the third and fourth moments to increase the accuracy of the approximation. The calculations are different for each test; for details, consult the references in the corresponding sections of this text.

### 13.12. Questions

1. Most microcomputer-based random number generators use multiplicative congruence to produce a 16-bit unsigned integer between zero and  $2^{15}$ . Yet in the two-sample comparison, for example, we only use one of the 15 bits, the least significant bit, in selecting items for rearrangement. Could we use more of the bits? That is, are some or all of the bits independent of one another? Write algorithm(s) that take advantage of multiple bits.
2. Apply the Mehta and Patel approach to the following  $3 \times 2$  contingency table:

3	1	0
1	2	1

Compute the marginals for this table. Draw a directed graph in which each node corresponds to a  $3 \times 2$  table whose marginals are the same as those of the proceeding table. Choose a test statistic (see Section 6.3). Identify those nodes which give rise to a value of the test statistic less than that of the original table.

3. Suppose you are interested in the theoretical alternative

4/6	1/6	1/6
1/6	4/6	1/6

How big a sample size would you need to insure that the probability of detecting this alternative was 80% at the 10% significance level? (Hint: use a six-sided die to simulate the drawing of samples.)

## CHAPTER 14

# Theory of Permutation Tests

In this chapter, we establish the underlying theory of permutation tests. The content is heavily mathematical, in contrast to previous chapters, and a knowledge of calculus is desirable.

## 14.1. Fundamental Concepts

In this section, we provide formal definitions for some of the concepts introduced in Chapter 2, including *distribution*, *power*, *exact*, *unbiased*, and the permutation test itself.

### 14.1.1. Dollars and Decisions

A statistical problem is defined by three elements:

- 1) the class  $F = (F_\theta, \theta \in \Omega)$  to which the probability distribution of the observations belongs; for example, we might specify that this distribution is unimodal, or symmetric, or normal;
- 2) the set  $D$  of possible decisions  $\{d\}$  one can make on observing  $X = (X_1, \dots, X_n)$ ;
- 3) the loss  $L(d, \theta)$ , expressed in dollars, men's lives, or some other quantifiable measure, that results when we make the decision  $d$  when  $\theta$  is true.

A problem is a statistical one when the investigator is not in a position to say that  $X$  will take on exactly the value  $x$ , but only that  $X$  has some probability  $P\{A\}$  of taking on values in the set  $A$ .

In this text, we've limited ourselves to two-sided decisions in which either we accept a hypothesis,  $H$ , and reject an alternative,  $K$ , or we reject the hypothesis,  $H$ , and accept the alternative,  $K$ .

One example is:

$$H: \theta \leq \theta_0$$

$$K: \theta > \theta_0.$$

In this example, we would probably follow up our decision to accept or reject with a confidence interval for the unknown parameter  $\theta$ . This would take the form of an interval  $(\theta_{\min}, \theta_{\max})$  and a statement to the effect that the probability that this interval covers the true parameter value is not less than  $1 - \alpha$ . This use of an interval can rescue us from the sometimes undesirable “all or nothing” dichotomy of hypothesis vs. alternative.

Another hypothesis/alternative pair which we considered in Section 3.6, under “testing for a dose response,” is

$$\begin{aligned} H: \theta_1 &= \dots = \theta_J \\ K: \theta_1 &< \dots < \theta_J. \end{aligned}$$

In this example, we might want to provide a confidence interval for  $\max_j \theta_j - \min_j \theta_j$ . Again, see Sections 3.2 and 7.4.

Typically, losses,  $L$ , depend on some function of the difference between the true (but unknown) value  $\theta$  and our best guess  $\theta^*$  of this value;  $L(\theta, \theta^*) = |\theta - \theta^*|$  for example. In the first of the preceding examples, we might have

$$\begin{aligned} L(\theta, d) &= \theta - \theta_0 && \text{if } \theta \in K \text{ and } d = H, \\ L(\theta, d) &= 10 && \text{if } \theta \in H \text{ and } d = K, \\ L(\theta, d) &= 0 && \text{otherwise.} \end{aligned}$$

Our objective is to come up with a decision rule,  $D$ , such that when we average out over *all* possible sets of observations  $X$ , we minimize the associated risk or expected loss,

$$R(\theta, D) = E L(\theta, D(X)).$$

Unfortunately, a testing procedure that is optimal for one value of the parameter,  $\theta$ , might not be optimal for another. This situation is illustrated in Chapter 2, in Figure 2.4, with two decision curves that cross over each other. The risk,  $R$ , depends on  $\theta$  and we don’t know what the true value of  $\theta$  is! How are we to choose the best decision?

This problem is complex with philosophical as well as mathematical overtones; we refer the interested reader to the discussions in the first chapter of Erich Lehmann’s book, *Testing Statistical Hypotheses* [1986]. Our own solution in selecting an optimal test is to focus on the principle of unbiasedness, which is discussed below in 14.1.3.

### 14.1.2. Tests

A test,  $\phi$ , is simply a decision rule that takes values between 0 and 1. When  $\phi(x) = 1$ , we reject the hypothesis and accept the alternative; when  $\phi(x) = 0$ , we accept the hypothesis and reject the alternative; and when  $\phi(x) = p$ , with  $0 < p < 1$ , we flip a coin that has been weighted so that the probability is  $p$  that it will come up heads,

whence we reject the hypothesis, and  $1 - p$  that it will come up tails, whence we accept the hypothesis.

An  $\alpha$ -level permutation test consists of a vector of  $N$  observations  $z$ , a statistic  $T[z]$ , and an acceptance criterion  $A : R \times R \rightarrow [0, 1]$ , such that for all  $z$ ,  $\phi(z) = 1$  if and only if

$$W(z) = \sum_{\pi \in \Pi} A(T[z], T[\pi z]) \leq \alpha N!,$$

where  $\Pi$  is the set of all possible rearrangements of the  $n + m$  observations.

### 14.1.3. Distribution Functions, Power, Exact, and Unbiased Tests

The *distribution function*  $F(x) = \Pr\{X \leq x\}$ ;  $F(x)$  is nondecreasing on the real line and  $0 \leq F(x) \leq 1$ . If  $F$  is continuous and differentiable, then it has a density  $f(x)$  such that  $\int_{-\infty}^{\infty} f(z) dz = F(x)$ .

We define the *power*  $\beta_\phi$  of a test  $\phi$  based on a statistic  $X$  as the expectation of  $\phi$ :  $\beta_\phi(\theta) = E^\theta \phi(X) = \int_{-\infty}^{\infty} f dF_\theta$ , where  $F_\theta$  is the distribution of  $X$ . Note that  $\beta_\phi$  is a function of the unknown parameter  $\theta$  (and, possibly, of other, nuisance parameters as well). For the majority of the tests in this book,  $\beta_\phi(\theta) = \Pr\{\phi = 1 | \theta\}$ .

If  $\theta$  satisfies the hypothesis, then  $\beta_\phi(\theta)$  is the probability of making a Type I error if  $\theta$  is true.

If  $\theta$  satisfies the alternative, then  $1 - \beta_\phi(\theta)$  is the probability of making a Type II error if  $\theta$  is true.

A test,  $\phi$ , is said to be *exact* with respect to a set  $\omega$  of hypotheses if  $E^H \phi = \alpha$ , for all  $H \in \omega$ . A test is conservative under the same circumstances if  $E^H \phi \leq \alpha$ , for all  $H \in \omega$ . The use of an exact, and thus conservative, test guarantees that the Type I error will be held at or below a predetermined level.

A test,  $\phi$ , is said to be *unbiased* and of level  $\alpha$  providing that its power function  $\beta$  satisfies

$$\begin{aligned}\beta_\phi(\theta) &\leq \alpha \text{ if } \theta \text{ satisfies the hypothesis} \\ \beta_\phi(\theta) &\geq \alpha \text{ if } \theta \text{ satisfies one of the alternatives.}\end{aligned}$$

That is, using an unbiased test,  $\phi$ , you are more likely to reject a false hypothesis than a true one.

### 14.1.4. Exchangeable Observations

Suppose that  $X_1, \dots, X_n$  are distributed as  $F(x)$ , while  $Y_1, \dots, Y_n$  are distributed as  $F(x - \delta)$ , and that  $F$  has probability density  $f$ .

A sufficient condition for a permutation test to be exact is the *exchangeability* of the observations [Lehmann, 1986, p. 231]. Let  $S(z)$  be the set of points obtained from  $z = (x_1, \dots, x_m, y_1, \dots, y_n)$  by permuting the coordinates of  $z$  in all  $(n+m)!$  possible ways.

**Theorem 1.** If  $F$  is the family of all  $(n + m)$ -dimensional distributions with probability densities,  $f$ , that are integrable and symmetric in their arguments, and we wish to test alternatives of the form  $f(x_1, \dots, x_m, y_1 - \delta, \dots, y_n - \delta)$  against the hypothesis that  $\delta = 0$ , a test  $\phi$  is unbiased for all  $f \in F$  if and only if

$$\sum_{z' \in S(z)} \phi(z') = \alpha(n + m)! \text{ a.e.}$$

The proof of this result relies on the fact that the set of order statistics constitute a complete sufficient statistic for  $F$ . See, for example, Lehmann [1986, pp. 45–6, 143–4, 231]. Also see problem 2 in this chapter. For more on exchangeability, see Koch [1982] and Romano [1990].

## 14.2. Maximizing the Power

In this section, we set about deriving the most powerful unbiased test for the two-sample testing problem. We will show that the two-sample test for a location parameter is unbiased against stochastically increasing alternatives. We define the likelihood ratio and restate, without proof, the fundamental theorem of Neyman and Pearson. We apply this theorem to show that the two-sample permutation test based on the sum of the observations is uniformly most powerful among unbiased tests against normal alternatives. Finally, we establish the intimate interdependence of confidence intervals and hypothesis tests. We follow closely the derivations provided in Lehmann [1986].

### 14.2.1. Uniformly Most Powerful Unbiased Tests

A family of cumulative distribution functions is said to be *stochastically increasing* if the distributions are distinct and if  $\theta < \theta'$  implies  $F_\theta(x) \geq F'_{\theta'}(x)$  for all  $x$ . One example is the location parameter family for which  $F_\theta(x) = F(x - \theta)$ . If  $X$  and  $X'$  have distributions  $F_\theta$  and  $F'_{\theta'}$ , then  $P\{X > x\} \leq P(X' > x)$ ; that is,  $X'$  tends to have larger values than  $X$ . Formally, we say that  $X'$  is *stochastically larger* than  $X$ .

**Lemma 1.**  $F_1(x) \leq F_0(x)$  for all  $x$  only if there exist two nondecreasing functions  $f_0$  and  $f_1$  and a random variable  $V$  such that  $f_0 \leq f_1$  for all  $v$  and the distributions of  $f_0$  and  $f_1$  are  $F_0$  and  $F_1$ , respectively.

**Proof.** Set  $f_i(y) = \inf\{x : F_i(x - 0) \leq y \leq F_i(x)\}$ ,  $i = 0, 1$ . These functions are nondecreasing and for  $f_i = f$ ,  $F_i = F$  satisfy  $f[F(x)] \leq x$  and  $F[f(y)] \geq y$  for all  $x$  and  $y$ . Thus,  $y \leq F(x_0)$  implies  $f(y) \leq f[F(x_0)] \leq x_0$  and  $f(y) \leq x_0$  implies  $F[f(y)] \leq F(x_0)$  implies  $y \leq F(x_0)$ .

Let  $V$  be uniformly distributed on  $(0, 1)$ . Then  $P\{f_i(V) \leq x\} = P\{V \leq F_i(x)\} = F_i(x)$  which completes the proof.  $\square$

We can apply this result immediately.

**Lemma 2.** Let  $X_1, \dots, X_m; Y_1, \dots, Y_n$  be samples from continuous distributions  $F, G$ , and let  $\phi[X_1, \dots, X_m; Y_1, \dots, Y_n]$  be a test such that a) whenever  $F = G$ , its expectation is  $\alpha$ ; and b)  $y_i \leq y'_i$  for  $i = 1, \dots, n$  implies  $\phi[x_1, \dots, x_m; y_1, \dots, y_n] \leq \phi[x_1, \dots, x_m; y'_1, \dots, y'_n]$ . Then the expectation of  $\phi$  is greater than or equal to  $\alpha$  for all pairs of distributions for which  $Y$  is stochastically larger than  $X$ .

**Proof.** From our first lemma, we know there exist functions,  $f$  and  $g$ , and independent random variables,  $V_1, \dots, V_{m+n}$ , such that the distributions of  $f(V_i)$  and  $g(V_i)$  are  $F$  and  $G$ , respectively, and  $f(z) \leq g(z)$  for all  $z$ .

$$E\phi[f(V_1), \dots, f(V_m); f(V_1), \dots, f(V_n)] = \alpha$$

and

$$E\phi[f(V_1), \dots, f(V_m); g(V_1), \dots, g(V_n)] = \beta.$$

From condition b) of the lemma, we see that  $\beta > \alpha$  as was to be proved.  $\square$

We are now in a position to state the principal result of this section.

**Theorem 2. (Unbiased).** Let  $X_1, \dots, X_m; Y_1, \dots, Y_n$  be samples from continuous distributions  $F, G$ . Let  $\beta(F, G)$  be the expectation of the critical function  $\phi$  defined in (14.1); that is,  $\phi[X_1, \dots, X_m; Y_1, \dots, Y_n] = 1$  only if  $\sum Y_j$  is greater than the equivalent sum in  $\alpha$  of the  $\binom{n+m}{n}$  possible rearrangements. Then  $\beta(F, F) = \alpha$  and  $\beta(F, G) \geq \alpha$  for all pairs of distributions for which  $Y$  is stochastically larger than  $X$ ;  $\beta(F, G) \leq \alpha$  if  $X$  is stochastically larger than  $Y$ .

**Proof.**  $\beta(F, F) = \alpha$ , follows from Theorem 1 and the definition of  $\phi$ . We can apply our lemmas and establish that the two-sample permutation test is unbiased if we can show that  $y_j \leq y'_j$  for  $j = 1, \dots, n$  implies

$$\phi[x_1, \dots, x_m; y_1, \dots, y_n] \leq \phi[x_1, \dots, x_m; y'_1, \dots, y'_n].$$

$\phi = 1$  if sufficiently many of the differences

$$d(\pi) = \sum_{i=m+1}^{m+n} z_i - \sum_{i=m+1}^{m+n} z_{j_i}$$

are positive. For a particular permutation  $\pi = (j_1, \dots, j_{m+n})$ ,

$$d(\pi) = \sum_{i=1}^p z_{s_i} - \sum_{i=m+1}^p z_{r_i},$$

where  $r_1 < \dots < r_p$  denote those of the integers  $j_{m+1}, \dots, j_{m+n}$  that are less than or equal to  $m$ , and  $s_1 < \dots < s_p$  denote those of the integers  $m+1, \dots, m+n$  that are not included in the set  $(j_{m+1}, \dots, j_{m+n})$ .

If  $\sum z_{s_i} - \sum z_{r_i}$  is positive and  $y_i \leq y'_i$ , that is  $z_i \leq z'_i$  for  $i = m+1, \dots, m+n$ , then the difference  $\sum z'_{s_i} - \sum z_{r_i}$  is also positive, so that  $\phi(z') \geq \phi(z)$ . But then

we may apply the lemmas to obtain the desired result. The proof is similar for the case in which  $X$  is stochastically larger than  $Y$ .  $\square$

### 14.2.2. The Fundamental Lemma

In Section 10.4, we showed that if the variables take only a countable number of values, then the most powerful test of a simple hypothesis  $P_0$  against a simple alternative  $P_1$  rejects the hypothesis in favor of the alternative only for those values of  $x$  with the largest values of the likelihood ratio

$$r(x) = \frac{p_1(x)}{p_0(x)}.$$

We can extend this result to continuous distribution functions with the aid of the fundamental lemma of Neyman and Pearson.

**Theorem 3.** *Let  $P_0$  and  $P_1$  be probability distributions possessing densities  $p_0$  and  $p_1$ , respectively.*

a) *There exists a test  $\phi$  and a constant  $k$  such that*

$$E_0\phi(X) = \alpha$$

*and*

$$\phi(x) = \begin{cases} 1 & \text{when } p_1(x) > kp_0(x) \\ 0 & \text{when } p_1(x) \leq kp_0(x). \end{cases}$$

b) *A test that satisfies these conditions for some  $k$  is most powerful for testing  $p_1$  against  $p_0$  at level  $\alpha$ .*

c) *If  $\phi$  is most powerful for testing  $p_1$  against  $p_0$  at level  $\alpha$ , then for some  $k$  it satisfies these conditions (except on a set that is assigned probability zero by both distributions and unless there exists a test at a smaller significance level whose power is one).*

A proof of this seminal lemma is given in Lehmann [1986, p. 74].

Let  $z$  denote a vector of  $n + m$  observations, and let  $S(z)$  be the set of points obtained from  $z$  by permuting the coordinates  $z_i$  ( $i = 1, \dots, n+m$ ) in all  $(n+m)!$  possible ways.

Among all the unbiased tests of the null hypothesis that two sets of observations come from the same distribution, which satisfy the permutation condition

$$\sum_{z' \in S(z)} \phi(z') = \alpha(n+m)!,$$

which is the most powerful?

Let  $t = T(z)$  denote the corresponding set of order statistics ( $z_{(1)} < z_{(2)} < \dots < z_{(n+m)}$ ). Lehmann [1986, p. 232] showed that the problem of maximizing the power

of a test  $\phi$  subject to the permutation condition against an alternative with arbitrary fixed density,  $h$ , reduces to maximizing

$$\sum_{z \in S(t)} \phi(z) \frac{h(z)}{\sum_{z' \in S(t)} h(z')}.$$

By the fundamental lemma of Neyman and Pearson, this expression is maximized by rejecting the hypothesis and setting  $\phi(z) = 1$  for those points  $z$  of  $S(t)$  for which the ratio

$$h(z) / \sum_{z' \in S(t)} h(z') \quad (14.1)$$

is largest. The most powerful unbiased  $\alpha$ -level test is given by rejecting when  $h(z) > C[T(z), \alpha]$  and accepting when  $h(z) < C[T(z), \alpha]$ . To achieve  $\alpha$  exactly, it may also be necessary to use a chance device and to reject with some probability,  $\gamma$ , if  $h(z) = C[T(z), \alpha]$ .

To carry out this test, we order the permutations according to the values of the density  $h$ . We reject the hypothesis for the  $k$  largest of the values, where

$$k \leq \alpha(n + m)! \leq k + 1.$$

The critical value  $C$  depends on the sample through its order statistics  $T$  and on the density  $h$ . Thus different distributions for  $X$  will give rise to different optimal tests.

### 14.2.3. Samples from a Normal Distribution

In what follows, we consider three applications of the likelihood ratio: Testing for the equality of the location parameters in two populations, testing for the equality of the variances, and testing for bivariate correlation.

Suppose that  $Z_1, \dots, Z_m$  and  $Z_{m+1}, \dots, Z_{n+m}$  are independent random samples from normal populations  $N(\eta, \sigma^2)$  and  $N(\eta + \delta, \sigma^2)$ . Then

$$\begin{aligned} h(z) &= (2\pi\sigma)^{-N/2} \exp \left[ -\frac{1}{2\sigma^2} \left( \sum_{j=1}^m (z_j - \eta)^2 + \sum_{j=m+1}^{n+m} (z_j - \eta - \delta)^2 \right) \right] \\ &= (2\pi\sigma)^{-N/2} \exp \left[ -\frac{1}{2\sigma^2} \left( \sum_{j=1}^{n+m} (z_j - \eta)^2 - 2\delta \sum_{j=m+1}^{n+m} (z_j - \eta) + n\delta^2 \right) \right]. \end{aligned}$$

Before substituting this expression in our formula, 14.1, we may eliminate all factors which remain constant under permutations of the subscripts. These include  $(2\pi\sigma)^{-(n+m)/2}$ ,  $n\delta(\delta + \eta)$  and  $\sum_{j=1}^{n+m} (z_j - \eta)^2$ . The resulting test rejects when

$\exp \left[ \delta \sum_{j=m+1}^{n+m} z_j \right] > C[T(z), \alpha]$  or, equivalently, when the sum of the observations in the treated sample  $\sum z_j$  is large. This sum can take at most  $(n+m)!$  possible values and our rejection region consists of the  $\alpha(n+m)!$  largest.

This permutation test is the same whatever the unknown values of  $\eta$  and  $\sigma$  and thus is uniformly most powerful against normally distributed alternatives among all unbiased tests of the hypothesis that the two samples come from the same population.

#### 14.2.4. Testing the Equality of Variances

As a second and elementary illustration of the likelihood ratio approach, suppose we are given that  $z_1, \dots, z_n$  are independent and identically normally distributed with mean 0 and variance  $\sigma^2$ ,  $N(0, \sigma^2)$ , and that  $z_{n+1}, \dots, z_{m+n}$  are independent and identically normally distributed with mean 0 and variance  $\tau^2$ ,  $N(0, \tau^2)$ . We wish to test the hypothesis that  $\sigma^2 = \tau^2$  against the alternative that  $\sigma^2 < \tau^2$ .

Let  $\theta = \tau^2/\sigma^2$  and note that hypothesis and alternative may be rewritten as  $H:\theta = 1$  vs  $K:\theta > 1$ .

Then

$$\begin{aligned} h(z) &= (2\pi)^{-(n+m)/2} \sigma^{-m} \tau^{-n} \exp \left[ -\frac{1}{2\sigma^2} \sum_{j=1}^n z_j^2 + \frac{1}{2\tau^2} \sum_{j=m+1}^{n+m} z_j^2 \right] \\ &= (2\pi\sigma)^{-(n+m)/2} \theta^{-n/2} \exp \left[ -\frac{1}{2\tau^2} \left( \theta \sum_{j=1}^n z_j^2 + \sum_{j=m+1}^{n+m} z_j^2 \right) \right] \\ &= (2\pi\sigma)^{-(n+m)/2} \theta^{-n/2} \exp \left[ -\frac{1}{2\tau^2} \left( (\theta - 1) \sum_{j=1}^n z_j^2 + \sum_{j=1}^{n+m} z_j^2 \right) \right]. \end{aligned}$$

Eliminating terms that are invariant under permutations of the combined sample, such as the sum of the squares of all  $n+m$  observations, we are left with the expression

$$\exp \left[ -\frac{1}{2\tau^2} (\theta - 1) \sum_{j=1}^n z_j^2 \right].$$

Our test statistic is the sum of the squares of the observations in the first sample.

#### 14.2.5. Testing for Bivariate Correlation

Suppose we have made  $N$  simultaneous observations on the pair of variables  $X, Y$  and wish to test the alternative of positive dependence of  $Y$  on  $X$  against the null

hypothesis of independence. In formal terms, if  $Y_x$  is the random variable whose distribution is the conditional distribution of  $Y$  given that  $X = x$ , we want to test the null hypothesis that  $Y_x$  has the same distribution for all  $x$ , against the alternative that if  $x' > x$ , then  $Y_{x'}$  is likely to be larger than  $Y_x$ .

To find a most powerful test of this hypothesis that is unbiased against alternatives with probability density  $h(z)$ , we need to maximize the expression

$$\sum_{z \in S(t)} \phi(z) \frac{h(z)}{\sum_{z' \in S(t)} h(z')}.$$

For bivariate normal alternatives,

$$h(z) = (2\pi\sigma\tau\sqrt{1-\rho^2})^{-n} \exp\left[-\frac{A}{2(1-\rho^2)}\right],$$

$$\text{where } A = \frac{1}{2\sigma^2} \sum_{j=1}^n (x_j - \eta)^2 + \frac{2\rho}{\sigma\tau} \sum_{j=1}^n (x_j - \eta)(y_j - \nu) + \frac{1}{2\tau^2} \sum_{j=1}^n (y_j - \nu)^2.$$

Many of the sums that occur in this expression are invariant under permutations of the subscripts  $j$ . These include the four sums  $\sum x_j$ ,  $\sum y_j$ ,  $\sum x_j^2$ ,  $\sum y_j^2$ . Eliminating all these invariant terms leaves us with the test statistic  $r = \sum x_j y_{\pi(j)}$ .

We evaluate this statistic both for the original data and for all  $n!$  permutations of the subscripts of the  $y$ 's, keeping the subscripts on the  $x$ 's fixed. We reject the null hypothesis in favor of the alternative of positive dependence only if the original value of the test statistic exceeds all but  $\alpha\%$  of the values for the rearrangements.

Reordering the  $x$ 's so that  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ , we see that this test is equivalent to using Pitman correlation (Section 3.5) to test the hypothesis of the randomness of the  $y$ 's against the alternative of an upward trend.

## 14.3. Confidence Intervals

Let  $x = \{X_1, X_2, \dots, X_n\}$  be an exchangeable sample from a distribution,  $F_\theta$ , which depends upon a parameter,  $\theta \in \Omega$ . A family of subsets,  $S(x)$ , of the parameter space  $\Omega$  is said to be a family of confidence sets for  $\theta$  at level  $1 - \alpha$  if

$$P_\theta\{\theta' \in S(X)\} \geq 1 - \alpha \quad \text{for all } \theta \in H(\theta').$$

The family is said to be unbiased if

$$P_\theta\{\theta' \in S(X)\} \leq 1 - \alpha \quad \text{for all } \theta \in \Omega - H(\theta').$$

The construction of a confidence set from a family of acceptance regions is described in Chapter 3. The following theorem shows us this construction can proceed in either direction.

**Theorem 4.1.** For each  $\theta' \in \Omega$ , let  $A(\theta')$  be the acceptance region of the level- $\alpha$  test for  $H(\theta'): \theta = \theta'$ , and for each sample point,  $x$ , let  $S(x)$  denote the set of parameter values  $\{\theta: x \in A(\theta), \theta \in \Omega\}$ . Then  $S(x)$  is a family of confidence sets for  $\theta$  at confidence level  $1 - \alpha$ .

**Theorem 4.2.** If for all  $\theta'$ ,  $A(\theta')$  is UMPU for testing  $H(\theta')$  at level  $\alpha$  against the alternatives  $K(\theta')$ , then for each  $\theta'$  in  $\Omega$ ,  $S(X)$  minimizes the probability

$$P_\theta\{\theta' \in S(X)\} \quad \text{for all } \theta \in K(\theta')$$

among all unbiased level  $1 - \alpha$  family of confidence sets for  $\theta$ .

**Proof 4.1.** By definition,  $\theta \in S(x)$  if and only if  $x \in A(\theta)$ ; hence,  $P_\theta\{\theta \in S(X)\} = P_\theta\{X \in A(\theta)\} \geq 1 - \alpha$ .

**Proof 4.2.** If  $S^*(x)$  is any other family of unbiased confidence sets at level  $1 - \alpha$  and if  $A^*(\theta) = \{x: \theta \in S^*(x)\}$ ; then

$$P_\theta\{X \in A^*(\theta')\} = P_\theta\{\theta' \in S^*(x)\} \geq 1 - \alpha \quad \text{for all } \theta \in H(\theta'),$$

and

$$P_\theta\{X \in A^*(\theta')\} = P_\theta\{\theta' \in S^*(x)\} \leq 1 - \alpha \quad \text{for all } \theta \in \Omega - H(\theta'),$$

so that  $A^*(\theta')$  is the acceptance region of a level- $\alpha$  unbiased test of  $H(\theta')$ . Since  $A$  is UMPU,

$$P_\theta\{X \in A^*(\theta')\} \geq P_\theta\{X \in A(\theta')\} \quad \text{for all } \theta \in \Omega - H(\theta');$$

hence,  $P_\theta\{\theta' \in S^*(x)\} \geq P_\theta\{\theta' \in S(x)\}$  for all  $\theta \in \Omega - H(\theta')$ , as was to be proved.  $\square$

## 14.4. Asymptotic Behavior

A major reason for the popularity of the permutation tests is that with very large samples their power is almost indistinguishable from that of the most powerful parametric tests. To establish this result, we need to know something about the distribution of the permutation statistics as the sample size increases without limit. Two sets of results are available to us. The first, due to Wald and Wolfowitz [1947] and Hoeffding [1953] provides us with conditions under which the limiting distribution is normal under the null hypothesis; the second, due to Albers, Bickel, and Van Zwet [1976] and Bickel and Van Zwet [1978] provides conditions under which this distribution is normal for near alternatives.

### 14.4.1. A Theorem on Linear Forms

Let  $S_N = (s_{N1}, s_{N2}, \dots, s_{NN})$  and  $U_N = (u_{N1}, u_{N2}, \dots, u_{NN})$  be sequences of real numbers, and let  $s_{N\cdot} = \sum s_{Nj}/N$ ;  $u_{N\cdot} = \sum u_{Nj}/N$ .

The sequences  $S_N$  satisfy the condition  $W$ , if for all integers  $r > 2$ ,

$$W(S_N, r) = \frac{\left| \sum (s_{Nj} - s_{N\cdot})^r \right|}{\sum [(s_{Nj} - s_{N\cdot})^2]^{r/2}} \text{ is bounded above for all } n.$$

The sequences  $S_N, U_N$  jointly satisfy the condition  $H_1$ , if for all integers  $r > 2$ ,

$$\lim_N N^{r/2-1} W(S_N, r) W(U_N, r) = 0.$$

The sequences  $S_N, U_N$  jointly satisfy the condition  $H_2$ , if for all integers  $r > 2$ ,

$$\lim_N N \frac{\max_j (s_{Nj} - s_{N\cdot})^r}{\sum (s_{Nj} - s_{N\cdot})^r} \frac{\max_j (u_{Nj} - u_{N\cdot})^r}{\sum (u_{Nj} - u_{N\cdot})^r}.$$

For any value of  $N$  let  $X = (x_1, x_2, \dots, x_N)$  be a chance variable whose possible values correspond to the  $N!$  permutations of the sequence  $A_N = (a_1, a_2, \dots, a_N)$ . Let each permutation of  $A_N$  have the same probability  $1/N!$ , and let  $E(Y)$  and  $SD(Y)$  denote the expectation and standard deviation of the variable  $Y$ .

**Theorem 5.** *Let the sequences  $A_N = (a_1, a_2, \dots, a_N)$  and  $D_N = (d_1, d_2, \dots, d_N)$  for  $N = 1, 2, \dots$ , satisfy any of the three conditions  $W$ ,  $H_1$ , and  $H_2$ . Let the chance variable  $L_N$  be defined as  $L_N = \sum d_i x_i$ . Then, as  $N \rightarrow \infty$ ,  $\Pr\{L_N - E(L_N) < t SD(L_N)\} \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^t e^{-x^2/2} dx$ .*

A proof of this result for condition  $W$  is given in Wald and Wolfowitz [1944]. The proof for conditions  $H_1$  and  $H_2$  is given in Hoeffding [1953].

This theorem applies to the majority of the tests we have already considered, including:

- 1) Pitman's correlation  $\sum d_i a_i$ ;
- 2) the two-sample test with observations  $a_1, \dots, a_{m+n}$ , and  $d_i$  equal to one if  $i = 1, \dots, m$  and zero otherwise;
- 3) Hotelling's  $T$  with  $\{a_{1j}\}$  and  $\{a_{2j}\}$  the observations—both sequences must separately satisfy the conditions of the theorem, and  $d_i = 1/m$  for  $i = 1, \dots, m$  and  $d_i = -1/n$  for  $i = m+1, \dots, m+n$ .

### 14.4.2. Asymptotic Efficiency

In this section, we provide asymptotic expansions to order  $N^{-1}$  for the power of the one- and two-sample permutation tests and compare them with the asymptotic

expansions for the most powerful parametric unbiased tests. The general expansion takes the form

$$b_N = c_0 + c_1 N^{-1/2} + c_{2,N} N^{-1} + o(N^{-1}),$$

where the coefficients depend on the form of the distribution, the significance level, and the alternative—but in both the one- and two-sample cases, the expansions for the permutation test and the  $t$ -test coincide for all terms through  $N^{-1}$ . The underlying assumptions are: 1) The observations are independent; 2) within each sample they are identically distributed; and 3) the two populations differ at most by a shift,  $G(x) = F(x - \delta)$ , where  $\delta \geq 0$ .  $\beta(p, F, \delta)$  and  $\beta(t, F, \delta)$  are the power functions of the permutation test and the parametric  $t$ -test, respectively (see Section 2.3). The theorem's other restrictions are technical in nature and provide few or no limitations in practice; e.g., the significance level must lie between 0 and 1, and the distribution must have absolute moments of at least ninth order. We state the theorem for the one-sample case only.

**Theorem 6.** Suppose the distribution  $F$  is continuous and that positive numbers  $C, D$ , and  $r > 8$  exist such that  $\int |x|^r dF[x] \leq C$  and  $0 \leq \delta \leq DN^{-1/2}$ ; then if  $\alpha$  is neither 0 nor 1, there exists a  $B > 0$  depending on  $C$  and  $D$ , and a  $b > 0$  depending only on  $r$  such that  $|\beta(p, F, \delta) - \beta(t, F, \delta)| \leq BN^{-1/b}$ .

Proof of this result and details of the expansion are given in Bickel and Van Zwet [1976]. The practical implication is that for large samples the permutation test and the parametric  $t$ -test make equally efficient use of the data.

Robinson [1989] finds approximately the same coverage probabilities for three sets of confidence intervals for the slope of a simple linear regression based, respectively, on 1) the standardized bootstrap, 2) parametric theory, and 3) a permutation procedure. Under the standard parametric assumptions, the coverage probabilities differ by  $o(n^{-1})$ , and the intervals themselves differ by  $O(n^{-1})$  on a set of probability  $1 - O(n^{-1})$ .

### 14.4.3. Exchangeability

The requirement that the observations be exchangeable can be relaxed at least asymptotically for some one-sample and two-sample tests. Let  $X_1, \dots, X_n$  be a sample from a distribution  $F$  that may or may not be symmetric. Let  $R_n(x, \Pi_n)$  be the permutation distribution of the statistic  $T_n(X_1, \dots, X_n)$ , and let  $r_n$  denote the critical value of the associated permutation test; let  $J_n(x, F)$  be the unconditional distribution of this same statistic under  $F$ , and let  $\Phi$  denote the standard normal distribution function.

**Theorem 7.** If  $F$  has mean 0 and finite variance  $\sigma^2 > 0$ , and  $T_n = n^{1/2}\bar{X}$ , then as  $n \rightarrow \infty$ ,

$$\sup_x |R_n(x, \Pi_n) - J_n(x, F)| \rightarrow 0 \text{ with probability 1,}$$

and  $\sup_x |R_n(x, \Pi_n) - \Phi(x/\sigma)| \rightarrow 0$  with probability 1. Thus  $r_n \rightarrow \sigma z_\alpha$ , with probability 1 and  $E_F[\phi(R_n)] \rightarrow \alpha$ .

A proof of this one-sample result is given in Romano [1990]; a similar one-sample result holds for a permutation test of the median subject to some mild continuity restrictions in the neighborhood of the median.

The two-sample case is quite different. Romano [1990] shows that if  $F_X$  and  $F_Y$  have common mean  $\mu$  and finite variances  $\sigma_X^2$  and  $\sigma_Y^2$ , respectively,  $T_{m,n} = n^{1/2}(\bar{X} - \bar{Y})$ , and  $m/n \rightarrow \lambda$  as  $n \rightarrow \infty$ , the unconditional distribution of  $T_{m,n}$  is asymptotically Gaussian with mean 0 and variance  $\sigma_X^2 + (1-\lambda)\sigma_Y^2/\lambda$ , while the permutation distribution of  $T_{m,n}$  is asymptotically Gaussian with mean 0 and variance  $\sigma_Y^2 + (1-\lambda)\sigma_X^2/\lambda$ . Thus, the two asymptotic distributions are the same only if either a) the variances of the two populations are the same or b) the sizes of the two samples are equal (whence  $\lambda = 1$ ).

Romano also shows that whatever the sample sizes, a permutation test for the difference of the medians of two populations will not be exact, even asymptotically (except in rare circumstances) unless the underlying distributions are the same.

## 14.5. Questions

1. **Unbiased.** The test  $\phi \equiv \alpha$  is a great timesaver; you don't have to analyze the data; you don't even have to gather data! All you have to do is flip a coin.
  - a) Prove that this test is unbiased.
  - b) Prove that a biased test cannot be uniformly most powerful.
2. **Sufficiency.** A statistic,  $T$ , is said to be sufficient for a family of distributions  $P = \{P_\theta, \theta \in \Omega\}$  (or sufficient for  $\theta$ ) if the conditional probability of an event given  $T = t$  is independent of  $\theta$ .
  - a) Let  $x_1, \dots, x_n$  be independent, identically distributed observations from a continuous distribution  $F_\theta$ . Show that the set of order statistics  $T = \{x_{(1)} < \dots < x_{(n)}\}$  is sufficient for  $\theta$ .
  - b) Let  $x_1, \dots, x_n$  be a sample from a uniform distribution  $U(0, \theta)$ , with density  $h(x) = 1/\theta$ , that is,  $P(x \leq u) = u/\theta$  for  $0 \leq u \leq \theta$ . Show that  $T = \max(x_1, \dots, x_n)$  is sufficient for  $\theta$ .
  - c) Let  $x_1, \dots, x_n$  be a sample from the exponential distribution with density  $\frac{1}{b}e^{-(x-a)/b}$ ,  $b > 0$ . Show that the pair  $\{\min(x_1, \dots, x_n), \sum x_i\}$  is sufficient for  $a, b$ .
3. **Likelihood ratio.**
  - a) Suppose that  $\{X_i, i = 1, \dots, n\}$  is  $N(\mu, \sigma^2)$  and  $\{Y_i, i = 1, \dots, m\}$  is  $N(\mu, \tau^2)$ . Derive the most powerful unbiased permutation test for testing  $H: \tau^2/\sigma^2 = 1$  against *not*  $H: \tau^2/\sigma^2 = 2$ .
  - b) The times between successive decays of a radioactive isotope are said to follow the exponential distribution, that is, the probability that an atom will not decay until after an interval of length  $t$  is  $1 - \exp[-t/\lambda]$ . (A similar formula provides a first-order approximation to the time,  $t$ , you will spend waiting for the next bus.) Suppose you had two potentially different isotopes with parameters  $\lambda_1$  and  $\lambda_2$ ,

respectively. Derive a *UMPU* permutation test for testing  $H: \lambda_1 = \lambda_2$ , against *not H*:  $\lambda_1 > \lambda_2$ .

- c) More generally, suppose that an item is reliable for a fixed period,  $b$ , after which its reliability decays at a constant rate  $\lambda$ . Then its lifetime has the exponential density  $\lambda^{-1} \exp[x - b]/\lambda$ . What statistic would you use for testing that  $H: \lambda_1 = \lambda_2$ , against *not H*:  $\lambda_1 > \lambda_2$ ? Is your answer the same as in 2b)? Why not? (Hint: Look for sufficient statistics. Note that the problem remains invariant under an arbitrary scale transformation applied to both sets of data. See Section 3.4).

#### 4. Asymptotic exchangeability.

- a) Show that the test for comparing variances provided in Section 3.4.1. is asymptotically exact.
- b) Show that the test for interactions based on the residuals described in Section 4.2.3. is asymptotically exact.
- c) Which, if any, of the tests of the multiple regression coefficients described in Section 7.6.1. are asymptotically exact?

# Bibliography

For your convenience, this bibliography is divided into three parts.

The first, main bibliography is of the research literature on permutation tests from the introduction of this straightforward approach to hypothesis testing by E.J.G. Pitman and R.A. Fisher in the mid-1930's to the present date.

This bibliography also includes articles we have cited in the text but that are not articles on permutation per se.

Since so much of today's research on permutation tests focuses on methods of rapid computation, we include a second, separate bibliography on computational methods.

A final and third bibliography consists of those few papers that we consider seminal both to an understanding of permutation tests and to the development of the subsequent vast wealth of articles on the topic. We hope every reader will select readings from this latter bibliography along with articles that are specific to their own interests.

In forming these bibliographies, we restricted ourselves to material on permutations and permutation tests that was directly related to hypothesis testing and estimation. Although, strictly speaking, every rank test is a permutation test, we did not include articles on rank tests in the bibliography unless, as is the case with some seminal work on multivariate analysis, the material is essential to an understanding of all permutation tests. Conference proceedings are excluded, the expected exception being a seminal paper by John Tukey that is available in no other form.

We have tried to be comprehensive, yet selective, and have personally read all but three of the articles in the bibliography. We hope you will find this bibliography of value in your work. We would appreciate your drawing to our attention articles on the theory and application of permutation tests that we may have excluded inadvertently.

## Bibliography Part 1:

# Permutation Test Articles

- Adams DC; Anthony CD. Using randomization techniques to analyse behavioural data. *Animal Behav.* 1996; 51: 733–738.
- Adamson P; Hajimohamadenza I; Brammer M; and Campbell IC. Intrasynatosomal free calcium concentration is increased by phobol esters via a 14-dihydropyridine-sensitive I-type CA2+. *European J. Pharm.* 1989; 162: 59–66.
- Adderley EE. Nonparametric methods of analysis applied to large-scale seeding experiments. *J Meteor.* 1961; 18: 692–694.
- Agresti A. *Categorical Data Analysis*. New York: Wiley; 1990.
- Agresti A. A survey of exact inference for contingency tables. *Statist. Sci.* 1992; 7: 131–177.
- Agresti A; Lang JB; and Mehta C. Some empirical comparisons of exact, modified exact, and higher-order asymptotic tests of independence for ordered categorical variables. *Commun. Statist. Simul.* 1993; 22: 1–18.
- Agresti A; Mehta CR; and Patel NR. Exact inference for contingency tables with ordered categories. *JASA*. 1990; 85: 453–458.
- Agresti A; Wackerly D. Some exact conditional tests of independence for  $R \times C$  cross-classification tables. *Psychometrika*. 1977; 42: 111–126.
- Agresti A; Wackerly D; and Boyett JM. Exact conditional tests for cross-classifications: Approximations of attained significance levels. *Psychometrika*. 1979; 44: 75–83.
- Agresti A; Yang M. An empirical investigation of some effects of sparseness in contingency tables. *Comput. Statist. Data Anal.* 1987; 5: 9–21.
- Ahmad IA; Kochar SC. Testing for dispersive ordering. *Statist. Probab. Ltr.* 1989; 8: 179–185.
- Albers W; Bickel PJ; and Van Zwet WR. Asymptotic expansions for the power of distribution-free tests in the one-sample problem. *Ann. Statist.* 1976; 4: 108–156.
- Albert A; Chapelle JP; Huesghem C; Kulbertus GE; and Harris EK. Evaluation of risk using serial laboratory data in acute myocardial infarction. In *Advanced Interpretation of Clinical Laboratory Data*. C Huesghem; A. Albert; ES Benson, Eds. New York: Marcel-Dekker; 1982.
- Alderson MR; Nayak RA. Study of space-time clustering in Hodgkin's disease in the Manchester Region. *British J. Preventive Social Med.* 1971; 25: 168–173.
- Alroy J. Permutation tests for the presence of phylogenetic structure: An editorial. *Systematics Biol.* 1994; 43: 430–437.
- Altham PME. Exact Bayesian analysis of a  $2 \times 2$  contingency table and Fisher's 'exact' significance test. *J. Roy. Statist. Soc. B* 1969; 31: 261–269.
- Altham PME. Improving the precision of estimates by fitting a model. *J. Roy. Statist. Soc. B* 1984; 46: 118–119.
- Andersen AH. Multidimensional contingency tables. *Scand. J. Statist.* 1974; 1: 115–127.

- Andersen PK; Borgan O; Gill R; and Keiding N. Linear nonparametric tests for comparison of counting processes with applications to censored survival data. *Int. Stat. Rev.* 1982; 50: 219–258.
- Anderson MJ; Legendre P. An empirical comparison of permutation methods for tests of partial regression coefficients in a linear model. *J. Statist. Comp. Simul.* 1999; 62: 271–303.
- Andres AM. A review of classic non-asymptotic methods for comparing two proportions by means of independent samples. *Commun. Statist. Simul.* 1991; 20: 551–583.
- Armitage P. Test for linear trend in proportions and frequencies. *Biometrics.* 1955; 11: 375–386.
- Arndt S; Cizadlo T; Andreasen NC; Heckel D; Gold S; and Oleary DS. Tests for comparing images based on randomization and permutation methods. *J. Cerebral Blood Flow Metabolism.* 1996; 16: 1271–1279.
- Arnold HJ. Permutation support for multivariate techniques. *Biometrika.* 1964; 51: 65–70.
- Arnold JT; Daw NC; Stenberg PE; Jayawardene D; Srivastava DK; and Jackson CWA. Single injection pegylated murine megakaryocyte growth and development factor (mgdf) into mice is sufficient to produce a profound stimulation megakaryocyte frequency size and ploidization. *Blood.* 1997; 89: 823–833.
- Ascher S; Bailar J. Moments of the Mantel-Valand procedure. *J. Statist. Comput. Simul.* 1982; 14: 101–111.
- Ayala G; Simo A. Stochastic labelling of biological images. *Statist. Neerl.* 1998; 52: 141–152.
- Baglivio J; Olivier D; and Pagano M. Methods for the analysis of contingency tables with large and small cell counts. *JASA.* 1988; 83: 1006–1013.
- Bahadur RR; Raghavachari M. Some asymptotic properties of likelihood ratios on general sample spaces. *Proceedings 6th Berkeley Symposium of Mathematical Statistics and Probability.* L. LeCam; J. Neyman Eds. Berkeley, CA: University of California Press; 1970: 129–152.
- Bailer AJ. Testing variance equality with randomization tests. *J. Statist. Comput. Simul.* 1989; 31: 1–8.
- Bailey RA. Randomization constrained. In *Encyclopedia of Statistical Sciences.* S. Kotz; Johnson NL, Eds. New York: Wiley. 1986; 7: 524–530.
- Bakeman R; Robinson BF; and Quera V. Testing sequential association: Estimating exact p values using sampled permutations. *Psych. Meth.* 1996; 1: 4–15.
- Baker FB; Collier RO. Some empirical results on variance ratios under permutation in the completely randomized design. *JASA.* 1966; 61: 813–820.
- Baker FB; Hubert LJ. Inference procedures for ordering theory. *J. Educ. Statist.* 1977; 2: 217–233.
- Baker RD. Two permutation tests of equality of variance. *Statist. Comput.* 1995; 5(4): 289–296.
- Baker RJ. Exact distributions derived from two-way tables. *J. Roy. Statist. Soc. C.* 1977; 26: 199–206.
- Ballin M; Pesarin F. Una procedura di ricampionamento e di combinazione non parametrica per il problema di Behrens-Fisher multivariato. *Proc. It. Statist. Soc.* 1990; 2: 351–358.
- Baptista J; Pike MC. Exact two-sided confidence limits for the odds ratio in a  $2 \times 2$  table. *J. Roy. Statist. Soc. C.* 1977; 26: 214–220.
- Barbella P; Denby L; and Glandwehr JM. Beyond exploratory data analysis: The randomization test. *Math. Teacher.* 1990; 83: 144–149.
- Barnard GA. A new test for  $2 \times 2$  tables (letter to the editor). *Nature.* 1945; 156: 177.
- Barnard GA. Statistical inference. *J. Roy. Statist. Soc. B.* 1949; 11: 115–139.
- Barnard GA. Discussion of paper by MS Bartlett. *J. Roy. Statist. Soc. B.* 1963; 25: 294.
- Barnard GA. In contradiction to J. Berkson's dispraise: Conditional tests can be more efficient. *J. Statist. Plan. Infer.* 1979; 3: 115–139.

- Barnard GA. Conditionality versus similarity in the analysis of  $2 \times 2$  tables. In *Essays in Honor of C.R. Rao*. Amsterdam: North Holland; 1982: 59–65.
- Barnard GA. On alleged gains in power from lower p-values. *Statist. Med.* 1989; 8: 1469–1477.
- Barnard GA. Must clinical trials be large? The interpretation of p-values and the combination of test results. *Statist. Med.* 1990; 9: 601–614.
- Barton DE; David FN. Randomization basis for multivariate tests. *Bull. Int. Statist. Inst.* 1961; 39(2): 455–467.
- Barton DE; David FN. The random intersection of two graphs. In *Research Papers in Statistics*. FN David, Ed. New York: Wiley; 1966.
- Barton DE; David FN; Fix E; Merrington M; and Mustacchi P. Tests for space-time interaction and a power transformation. *Proceedings 5th Berkeley Symposium on Mathematical Statistics and Probability*. Lucien LeCam, Ed. Berkeley, CA: University of California Press; 1967; IV.
- Basawa IV; Rao Praska. *BLS Statistical Inference for Stochastic Processes*. New York: Academic Press; 1980.
- Basu D. On the relevance of randomization in data analysis (with discussion). In *Survey Sampling and Measurement*. NK Namboodiri, Ed. New York: Academic Press; 1978: 267–339.
- Basu D. Discussion of Joseph Berkson's paper "In dispraise of the exact test". *J. Statist. Plan. Infer.* 1979; 3: 189–192.
- Basu D. Randomization analysis of experimental data: The Fisher randomization test. *JASA* 1980; 75: 575–582.
- Bayer L; Cox C. Exact test of significance in binary regression models. *J. Roy. Statist. Soc. C* 1979; 28: 319–324.
- Bedrick KE; Hill JR. Outlier tests for logistic regression: a conditional approach. *Biometrika* 1990; 77: 815–827.
- Bell CB; Doksum KA. Some new distribution free statistics. *Ann. Math. Statist.* 1965; 36: 203–214.
- Bell CB; Doksum KA. Distribution-free tests of independence. *Ann. Math. Statist.* 1967; 38: 429–446.
- Bell CB; Donoghue JF. Distribution-free tests of randomness. *Sankhya A* 1969; 31: 157–176.
- Bell CP; Sen PK. Randomization Procedures. In *Nonparametric Methods*. PR Krishnaiah; PK Sen, Eds. Amsterdam: North Holland; 1984; 4: 1–30.
- Bell CB; Woodroofe M; and Avadhani TV. Some nonparametric tests for stochastic processes. In *Nonparametric Techniques in Statistical Inference*. ML Puri, Ed. Cambridge, U.K.: Cambridge University Press; 1970.
- Belyea LR. Separating the effects of litter quality and microenvironment on decomposition rates in a patterned peatland. *Oikos*. 1996; 77: 529–539.
- Berger RL; Boos DD. P values maximized over a confidence set for a nuisance parameter. *JASA*. 1994; 89: 1012–1016.
- Berger JO; Wolpert RW. *The Likelihood Principle*. IMS Lecture Notes-Monograph Series. Hayward, CA: IMS; 1984.
- Berkson J. In dispraise of the exact test. *J. Statist. Plan. Inf.* 1978; 2: 27–42.
- Berkson J. Do the marginals of the  $2 \times 2$  table contain relevant information respecting the table proportions? *J. Statist. Plan. Inf.* 1979; 3: 193–197.
- Berlin JA; Ness RB. Randomized clinical-trials in the presence diagnostic uncertainty-implications for measures of efficacy and sample-size. *Controlled Clinical Trials*. 1996; 17: 191–200.
- Berry KJ; Kvamme KL; and Mielke PW Jr. Permutation techniques for the spatial analysis of the distribution of artifacts into classes. *Amer. Antiquity*. 1980; 45: 55–59.
- Berry KJ; Kvamme KL; and Mielke PW Jr. Improvements in the permutation test for the spatial analysis of the distribution of artifacts into classes. *Amer. Antiquity*. 1983; 48: 547–553.

- Berry KJ; Mielke PW. Computation of finite population parameters and approximate probability values of multi-response permutation procedures (MRPP). *Commun Statist B*. 1983; 12: 83–107.
- Berry KJ; Mielke PW Jr. Computation of exact probability values for multi-response permutation procedures (MRPP). *Commun Statist B*. 1984; 13: 417–432.
- Berry KJ; Mielke PW. Computation of exact and approximate probability values for a matched-pairs permutation test. *Commun Statist B*. 1985; 14: 229–248.
- Berry KJ; Mielke PW. A generalization of Cohen's kappa agreement measure to interval measurement and multiple raters. *Educ. Psych. Measurement*. 1988; 48: 921–933.
- Berry KJ; Mielke PW Jr. A family of multivariate measures of association for nominal independent variables. *Educ. Psych. Measurement*. 1992; 52: 97–101.
- Berry KJ; Mielke PW Jr. A measure of association for nominal independent variables. *Educ. Psych. Measurement*. 1992; 52: 895–898.
- Berry KJ; Mielke PW. Least sum of absolute deviations regression: Distance, leverage and influence. *Perceptual Motor Skills*. 1998; 86: 1063–1070.
- Berry KJ; Mielke PW Jr; and Helmericks SG. Exact confidence limits for proportions. *Educ. Psych. Measurement*. 1988; 48: 713–716.
- Berry KJ; Mielke PW; and Wary RKW. Approximate MRPP p-values obtained from four exact moments. *Commun Statist B*. 1986; 15: 581–589.
- Bersier LF; Sugihara G. Species abundance patterns—the problem of testing stochastic models. *J. Animal Ecol.* 1997; 66: 769–774.
- Bertacche R; Pesarin F. Treatment of missing data in multidimensional testing problems for categorical variables. *Metron*. 1997; 55: 135–149.
- Besag JE. Some methods of statistical analysis for spatial data. *Bull. Int. Statist. Inst.* 1978; 47: 77–92.
- Besag J; Clifford P. Generalize Monte Carlo significance tests. *Biometrika*. 1989; 76: 633–642.
- Besag J; Diggle PJ. Simple Monte Carlo tests for spatial pattern. *Appl. Statist.* 1977; 25: 327–333.
- Bickel PJ. A distribution free version of the Smirnov two-sample test in the multivariate case. *Ann. Math. Statist.* 1969; 40: 1–23.
- Bickel PM; Van Zwet WR. Asymptotic expansion for the power of distribution free tests in the two-sample problem. *Ann. Statist.* 1978; 6: 987–1004 (corr. 1170–1171).
- Birch MW. The detection of partial association. *J. Roy. Statist. Soc. B*. 1964; 26/27: I 313–324, II 1–124.
- Birnbaum ZW. Computers and unconventional test statistics. In *Reliability and Biometry*. F Presham; RJ Serfling Eds. Philadelphia, PA: SIAM; 1974.
- Bishop YMM; Fienberg SE; and Holland PW. *Discrete Multivariate Analysis: Theory and Practice*. Cambridge, MA: MIT Press; 1975.
- Blair C; Higgins JJ; Karinski W; Krom R; and Rey JD. A study of multivariate permutation tests which may replace Hotellings T-test in prescribed circumstances. *Multivar. Beh. Res.* 1994; 29: 141–163.
- Blair C; Karinski W. Distribution-free statistical analyses of surface and volumetric maps. In *Functional Neuroimaging: Technical Foundations*. Tatcher RW; Hallett M; Zeffiro T; John ER; and Huerta M, Eds., New York: Academic Press; 1994.
- Blair RC; Troendle JF; and Beck RW. Control of familywise errors in multiple endpoint assessments via stepwise permutation tests. *Statist. Med.* 1996; 15: 1107–1121.
- Boess FG; Balasuramanian R; Brammer MJ; and Campbell IC. Stimulation of muscarinic acetylcholine receptors increases synaptosomal free calcium concentration by protein kinase-dependent opening of L-type calcium channels. *J. Neurochem.* 1990; 55: 230–236.
- Boik RJ. The Fisher-Pitman permutation test: A non-robust alternative to the normal theory F-test when variances are heterogeneous. *Brit. J. Math. Stat. Psych.* 1987; 40: 26–42.
- Bookstein FL. Shape and the information in medical images—a decade of the morphometric synthesis. *Computer Vision Image Understanding*. 1997; 66: 97–118.

- Bookstein FL. A hundred years of morphometrics. *Acta Zool. Acad. Sci. H.* 1998; 44: 7–59.
- Boos DD; Browne C. Testing for a treatment effect in the presence of nonresponders. *Biometrics*. 1986; 42: 191–197.
- Booth JG; Butler RW. Randomization distributions and saddlepoint approximations in general linear models. *Biometrika*. 1990; 77: 787–796.
- Boschloo RD. Raised conditional level of significance for the  $2 \times 2$  table when testing the equality of two probabilities. *Statist. Neer.* 1970; 24: 1–35.
- Box GEP; Anderson SL. Permutation theory in the development of robust criteria and the study of departures from assumptions. *J. Roy. Statist. Soc. B.* 1955; 17: 1–34 (with discussion).
- Box JF. *The Life of a Scientist*. New York: Wiley; 1978.
- Boyd MN; Sen PK. Union intersection tests for ordered alternatives in ANOCOVA. *JASA*. 1986; 81: 526–532.
- Boyett JM; Shuster JJ. Nonparametric one-sided tests in multivariate analysis with medical applications. *JASA*. 1977; 72: 665–668.
- Bradbury IS. Analysis of variance vs randomization tests: a comparison (with discussion by White and Still). *Brit. J. Math. Statist. Psychol.* 1987; 40: 177–195.
- Bradbury IS. Approximations to permutation distributions in the completely randomized design. *Commun Statist. T-M A.* 1988; 17: 543–555.
- Bradley JV. *Distribution Free Statistical Tests*. Englewood Cliffs, NJ: Prentice-Hall; 1968.
- Bradley RA; Scott E. Perspectives from a weather modification experiment. *Commun Statist. A.* 1980; 9: 1941–1961.
- Breslow NE; Day NE. *I. Analysis of Case Control Studies. II. Design and Analysis of Cohort Studies*. New York: Oxford University Press; 1980, 1987.
- Brockwell PJ; Mielke PW Jr. Asymptotic distributions of matched-pair permutation statistics based on distance measures. *Australian J. Statist.* 1984; 26: 30–38.
- Brockwell PJ; Mielke PW; and Robinson J. On non-normal invariance principles for multi-response permutation procedures. *Australian J. Statist.* 1982; 24: 33–41.
- Bross IDJ. Taking a covariate into account. *JASA*. 1964; 59: 725–736.
- Brown BM. Robustness against inequality of variances *Australian J. Statist.* 1982; 24: 283–295.
- Brown BM. Cramer-von Mises distributions and permutation tests. *Biometrika*. 1982; 69: 619–624.
- Brown BM; Hettmansperger TP. Affine invariant rank methods in the bivariate location model. *J. Roy. Statist. Soc. B.* 1987; 49: 301–310.
- Brown BM; Maritz JS. Distribution-free methods in regression. *Australian J. Statist.* 1982; 24: 318–331.
- Brown CC; Fears TR. Exact significance levels for multiple binomial testing with applications to carcinogenicity screens. *Biometrics*. 1981; 37: 763–774.
- Bryant EH. Morphometric adaptation of the housefly, *Musa domestica* L; in the United States. *Evolution*. 1977; 31: 580–596.
- Buckland ST; Garthwaite PH. Quantifying precision of mark-recapture estimates using the bootstrap and related methods. *Biometrics*. 1991; 47: 225–268.
- Bullmore E; Brammer M; Williams SC; Rabehesk S; Janet N; David A; Mekers J; Howard R; and Slam P. Statistical methods for estimation and inference for functional MR image analysis. *Magn. Res. M.* 1996; 35(2): 261–277.
- Buonaccorsi JP. A note on confidence intervals for proportions in finite populations. *Amer. Statist.* 1987; 41: 215–218.
- Burgess AP; Gruzelier JH. Short-duration synchronization human theta rhythm during recognition memory. *Neuroreport*. 1997; 8: 1039–1042.
- Büringer H; Martin H; and Schrieven K-H. *Nonparametric Sequential Selection Procedures*. Boston, Ma: Birkhauser; 1980.
- Busby DG. Effects of aerial spraying of fenithrothion on breeding white-throated sparrows. *J. Appl. Ecol.* 1990; 27: 745–755.

- Cade B. Comparison of tree basal area and canopy cover in habitat models: Subalpine forest. *J. Alpine Mgmt.* 1997; 61: 326–335.
- Cade BS; Hoffman RW. Winter use of Douglas-fir forests by blue grouse in Colorado. *J. Wildlife Mgmt.* 1990; 27: 743–755.
- Cade B; Hoffman H. Differential migration of blue grouse in Colorado. *Auk.* 1993; 110: 70–77.
- Cade B; Richards L. Permutation tests for least absolute deviation regression. *Biometrics.* 1996; 52: 886–902.
- Cardwell KF; Wehrly TE. A rank test for distinguishing environmentally and genetically induced disease resistance in plant varieties. *Biometrics.* 1997; 53: 195–206.
- Casagrande JJ; Pike MC; and Smith PG. The power function of the exact test for comparing two binomial distributions. *Appl. Statist.* 1978; 27: 176–181.
- Casella G; George EI. Explaining the Gibbs Sampler. *JASA.* 1992; 46: 167–174.
- Chapelle JP; Albert A; Smeets JP; Heusghem C; and Kulberts HE. Effect of the hypotoglobin phenotype on the size of a myocardial infarct. *New England J. Med.* 1982; 307: 457–463.
- Chapman JW. A comparison of the chi-square,  $-2\log R$ , and multinomial probability criteria for significance tests when expected frequencies are small. *JASA.* 1976; 71: 854–863.
- Chatterjee SK; Sen PK. Nonparametric tests for the bivariate two-sample location problem. *Calcutta Statist. Ass. Bull.* 1964; 13: 18–58.
- Chatterjee SK; Sen PK. Nonparametric tests for the multivariate, multisample location problem. In *Essays in Probability and Statistics in memory of SN Roy.* RC Bose, et al., Eds. Chapel Hill, NC: Univ. North Carolina Press; 1966.
- Chatterjee SK; Sen PK. *Calcutta Statist. Ass. Bull.* 1973; 22: 13–50.
- Chen XR. A two-sample permutation test with heterogeneous blocks. *Wuhan Daxue Xuebao.* 1980; 4: 1–14.
- Chen XR. Two problems of linear permutation statistics. *Acta Math Appl Sinica.* 1981; 4: 342–355.
- Chen XR. Large sample theory of permutation tests in the case of a randomized block design. *J. Wuhan University, Natural Science Edition.* 1983; 4: 1–12.
- Chernoff H; Savage IR. Asymptotic normality and efficiency of certain nonparametric test statistics. *Ann. Math. Statist.* 1929.
- Chhikara RK. State of the art in credit evaluation. *Amer. J. Agric. Econ.* 1989; 71: 1138–1144.
- Chung JH; Fraser DAS. Randomization tests for a multivariate two-sample problem. *JASA.* 1958; 53: 729–735.
- Clark RM. A randomization test for the comparison of ordered sequences. *Math. Geology.* 1989; 21: 429–442.
- Clarke B. Divergent effects of natural selection on two closely-related polymorphic snails. *Heredity.* 1960; 14: 423–443.
- Clarke B. Natural selection in mixed populations of two closely-related polymorphic snails. *Nature.* 1962; 16: 319–345.
- Cléoux R. First and second moments of the randomization test in two associate PBIB designs. *JASA.* 1969; 64: 1224–1235.
- Cliff AD; Ord JK. Evaluating the percentage points of a spatial autocorrelation coefficient. *Geog. Anal.* 1971; 3: 51–62.
- Cliff AD; Ord JK. *Spatial Processes: Models and Applications.* London, U.K.: Pion Ltd; 1981.
- Cohen A. Unbiasedness of tests for homogeneity. *Ann. Statist.* 1987; 15: 805–816.
- Collier RO Jr; Baker FB. The randomization distribution of F-ratios for the split-plot design—an empirical investigation. *Biometrika.* 1963; 50: 431–438.
- Collier RO Jr; Baker FB. Some Monte Carlo results on the power of the F-test under permutation in the simple randomized block design. *Biometrika.* 1966; 53: 199–203.

- Collins MF. A permutation test for planar regression. *Australian J. Statist.* 1987; 29: 303–308.
- Conover WJ; Johnson ME; and Johnson MM. Comparative study of tests for homogeneity of variances, with applications to the outer continental shelf bidding data. *Technometrics*. 1981; 23: 351–361.
- Constanzo CM; Hubert LJ; and Golledge RG. A higher moment for spatial statistics. *Geog. Anal.* 1983; 15: 347–351.
- Cormack RS. The meaning of probability in relation to Fisher's exact test. *Metron.* 1986; 44: 1–30.
- Cormack RS; Mantel N. Fisher's exact test: The marginal totals as seen from two different angles. *Statistician.* 1991; 40: 27–34.
- Cornfield J. On samples from finite populations. *JASA.* 1944; 39: 236–239.
- Cornfield J. A statistical problem arising from retrospective studies. In *Proceedings of 3rd Berkeley Symposium of Mathematical Statistics and Probability* J. Neyman, Ed. Berkeley, CA: University of California Press; 1956; 4: 135–138.
- Cornfield J; Tukey JW. Average values of mean squares in factorials. *Ann. Math. Statist.* 1956; 27: 907–949.
- Cory-Slechta DA. Exposure duration modalities, the effects of low-level lead on fixed interval performances. *Neurotoxicology.* 1990; 11: 427–442.
- Cory-Slechta DA; Weiss B; and Cox C. Tissue distribution of Pb in adult vs. old rats: A pilot study. *Toxicology.* 1989; 59: 139–150.
- Cotton JW. Even better than before. *Contemp Psychol.* 1973; 18: 168–169.
- Cox DF; Kempthorne O. Randomization tests for comparing survival curves. *Biometrics.* 1963; 19: 307–317.
- Cox DR. A note on weighted randomization. *Ann. Math. Statist.* 1956; 27: 1144–1150.
- Cox DR. Interpretation of nonadditivity in Latin Square. *Biometrika.* 1958; 45: 69–73.
- Cox DR. Regression models and life tables (with discussion). *J. Roy. Statist. Soc B.* 1972; 34: 187–220.
- Cox DR. Partial likelihood. *Biometrika.* 1975; 62: 269–276.
- Cox DR. A remark on randomization in clinical trials. *Utilitas Math.* 1982; 21A: 242–252.
- Cox DR. Interaction. *Int. Statist. Rev.* 1984; 52: 1–31.
- Cox DR; Shell EJ. *Applied Statistics. Principles and Examples.* London, U.K.: Chapman and Hall; 1981.
- Cox DR; Shell EJ. *Analysis of Binary Data* 2 Ed. London: Chapman-Hall; 1989.
- Cox MAA; Plackett RL. Small samples in contingency tables. *Biometrika.* 1980; 67: 1–13.
- Crisp MD; Linder HP; and Weston PH. Cladistic biogeography of plants in Australia and New Guinea: Congruent pattern reveals two endemic tropical tracks. *Syst. Biol.* 1995; 44: 457–473.
- Crump KS; Howe RB; and Kodell RL. Permutation tests for detecting teratogenic effects. In *Statistics in Toxicology*. Krewski D; Franklin C, Eds. New York: Gordon and Breach Science Publishers; 1990: 347–375.
- D'Abadie C; Proschan F. Stochastic versions of rearrangement inequalities. In *Inequalities in Statistics and Probability*. YL Tong, Ed. Hayward, CA: IMS; 1984: 4–12.
- D'Agostino RB; Chase W; and Belanger A. The appropriateness of some common procedures for testing the equality of two independent binomial populations. *Amer. Statist.* 1988; 42: 198–202.
- Daniels HE. Relation between measures of correlation in the universe of sample permutations. *Biometrika.* 1944; 33: 129–135.
- Dansie BR. A note on permutation probabilities. *J. Roy. Statist. Soc. B.* 1983; 45: 22–24.
- David FN. Measurement of diversity. In *Proceedings of 6th Berkeley Symposium of Mathematical Statistics and Probability.* 1971; 1: 631–648.
- David FN; Barton DE. Two space-time interaction tests for epidemicity. *Brit. J. Prev. Soc. Med.* 1966; 20: 44–48.
- David HA. *Order Statistics.* New York: Wiley; 1970.

- David HT; Kruskal WH. The Wagr sequential t-test reaches a decision with probability one. *Ann. Math. Statist.* 1956; 27: 797–805, 1958; 29: 936.
- Davis AW. On certain ratio statistics in weather modification experiments. *Technometrics*. 1979; 21: 283–290.
- Davis AW. On the effects of moderate nonnormality on Roy's largest root test. *JASA*. 1982; 77: 896–900.
- Davis AW; Speed TP. An Edgeworth expansion for the distribution of the F-ratio under the randomization model for the randomized block design. In *Statistical Decision Theory and Related Topics*. SS Gupta; JO Berger, Eds. New York: Springer-Verlag; 1988; 2: 119–129.
- Davis LJ. Exact tests for  $2 \times 2$  contingency tables. *Amer. Statist.* 1986; 40: 139–140.
- Daw NC; Arnold JT; Abushullah BA; Stenberg PE; White MM; Jayawardene D; Srivastava DK; and Jackson CWA. Single intravenous dose murine megakaryocyte growth and development factor potently stimulates platelet production challenging the necessity for daily administration. *Blood*. 1998; 91: 466–474.
- Dee CR; Rankin JA; and Burns CA. Using scientific evidence to improve hospital library services: Southern Chapter Medical Library Association journal usage study. *B. Med. Libr. Assoc.* 1998; 86: 301–306.
- Denker M; Puri ML. Asymptotic behavior of multiresponse permutation procedures. *Adv. Appl. Math.* 1988; 9: 200–210.
- Dennis AS; Miller JR; Cain DE; and Schwaller RL. Evaluation by Monte Carlo tests of effects of cloud seeding on growing season rainfall in North Dakota. *J. Appl. Meter.* 1975; 14: 959–964.
- Deutsch SJ; Schmeiser BW. The power of paired sample t-tests. *Naval Res. Logistics Quart.* 1982; 29: 635–649.
- Deutsch ST; Schmeiser BW. The computation of the component randomization test for paired comparisons. *J. Quality Technology*. 1983; 15: 94–98.
- Devroye L; Gyorfi L. *Nonparametric Density Estimation: The L1 View*. New York: Wiley; 1985.
- Diaconis P. Statistical problems in ESP research. *Science*. 1978; 201: 131–136.
- Diaconis P; Efron B. Computer intensive methods in statistics. *Scientific American*. 1983; 248: 116–130.
- Diaconis P; Graham RL. Spearman's footrule as a measure of disarray. *J. Roy. Statist. Soc. B*. 1972; 39: 262–268.
- DiCiccio TJ; Romano J. A review of bootstrap confidence intervals (with discussions). *J. Roy. Statist. Soc. B*. 1988; 50: 163–170.
- Dietz EJ. Permutation tests for the association between two distance matrices. *Systemic Zoology*. 1983; 32: 21–26.
- Diggle PJ; Lange N; and Benes FM. Analysis of variance for replicated spatial point patterns in clinical neuroanatomy. *JASA*. 1991; 86: 618–625.
- Dodge Y. ed. *Statistical Data Analysis Based on the L1 Norm and Related Methods*. Amsterdam: North Holland; 1987.
- Donegani M. An adaptive and powerful test. *Biometrika*. 1991; 78: 930–933.
- Donegani M. Asymptotic and approximate distribution of a statistic obtained by resampling with and without replacement. *Statist. Prob. Letters*. 1991; 11: 181–183.
- Donnelly SM; Kramer A. Testing for multiple species in fossil samples: An evaluation and comparison of tests for equal relative variation. *Amer. J. Phys. Anthro.* 1999; 108: 507–529.
- Donner A. Odds ratio inference with dependent data: A relationship between two procedures. *Biometrika*. 1974.
- Doolittle RF. Similar amino acid sequences: Chance or common ancestry. *Science*. 1981; 214: 149–159.
- Douglas ME; Endler JA. Quantitative matrix comparisons in ecological and evolutionary investigations. *J. Theoret. Biol.* 1982; 99: 777–795.

- Draper D; Hodges JS; Mallows CL; and Pregibon D. Exchangeability and data analysis (with discussion). *J. Roy. Statist. Soc. A.* 1993; 156: 9–28.
- Draper NR; Stoneman DM. Testing for the inclusion of variables in linear regression by a randomization technique. *Technometrics*. 1966; 8: 695–699.
- Duffy DE; Quinoz AJ. Permutation-based algorithm for block clustering. *J. Classification*. 1991; 8: 65–91.
- Dupont WD. Sensitivity of Fisher's exact test to minor perturbations in  $2 \times 2$  contingency tables. *Statist. Med.* 1986; 5: 629–635.
- Dwass M. On the asymptotic normality of some statistics used in nonparametric tests. *Ann. Math. Statist.* 1955; 26: 334–339.
- Dwass M. Modified randomization tests for nonparametric hypotheses. *Ann. Math. Statist.* 1957; 28: 181–187.
- Easterling RG. Randomization and statistical inference. *Commun. Statist.* 1975; 4: 723–735.
- Edelman D. Bounds for a nonparametric t-table. *Biometrika*. 1986; 73: 242–243.
- Eden T; Yates F. On the validity of Fisher's z test when applied to an actual sample of nonnormal data. *J. Agricultural Sci.* 1933; 23: 6–16.
- Edgington ES. Randomization tests. *J. Psych.* 1964; 57: 445–449.
- Edgington ES. Statistical inference and nonrandom samples. *Psychol. Bull.* 1966; 66: 485–487.
- Edgington ES. Approximate randomization tests. *J. Psych.* 1969; 72: 143–149.
- Edgington ES. Hypothesis testing without fixed levels of significance. *J. Psych.* 1970; 76: 109–115.
- Edgington ES. Randomization tests for one-subject operant experiments. *J. Psych.* 1975; 90: 57–68.
- Edgington ES. Randomization tests for predicted trends. *Canadian Psych. Review*. 1975; 16: 49–53.
- Edgington ES. Validity of randomization tests for one-subject experiments. *J. Educ. Statist.* 1980; 5: 235–251.
- Edgington ES. Overcoming obstacles to single-subject experimentation. *J. Educ. Statist.* 1980; 5: 261–267.
- Edgington ES. The role of permutation groups in randomization tests. *J. Educ. Statist.* 1983; 8: 121–145.
- Edgington ES. Statistics and single-case analysis. In *Progress in Behavior Modification*. Miltersen; RM Eisler; and PM Miller. New York: Academic Press; 1984: 16.
- Edgington ES. *Randomization Tests*. 3rd ed. New York: Marcel-Dekker; 1995.
- Edginton ES. Randomized single-subject experimental designs. *Beh. Res. Therapy*. 1996; 34: 567–574.
- Edginton ES; Bland BH. Randomization tests: Application to single-cell and other single-unit neuroscience experiments. *J. NeuroSci. Meth.* 1993; 47: 169–177.
- Edginton ES; Ezinga G. Randomization tests and outlier scores. *J. Psych.* 1978; 99: 259–262.
- Edginton ES; Gore AP. Randomization tests for censored survival distributions. *Biometrical J.* 1986; 28: 673–681.
- Edwards D. Exact simulation-based inference: A survey, with additions. *J. Statist. Comp. Simul.* 1980; 22: 307–326.
- Edwards D. On model prespecification in confirmatory randomized studies. *Statist. Med.* 1999; 18: 771–785.
- Efron B. Forcing sequential experiments to be balanced. *Biometrika*. 1971; 58: 403–417.
- Efron B. Bootstrap methods: Another look at the jackknife. *Ann. Statist.* 1979; 7: 1–26.
- Efron B. Censored data and the bootstrap. *JASA*. 1981; 76: 312–319.
- Efron B. Better bootstrap confidence intervals (with discussion) *JASA*. 1987; 82: 171–200.
- Efron B. Three examples of computer intensive statistical inference. *Sankhya A*. 1988; 50: 338–362.

- Efron B; Johnstone I. Fisher's information in terms of the hazard rate. *Ann. Statist.* 1990; 18: 38–62.
- Efron B; Tibshriani R. Bootstrap methods for standard errors, confidence intervals and other measures of scientific accuracy. *Statist. Sci.* 1986; 1: 54–77.
- Efron B; Tibshriani R. Statistical data analysis in the computer age. *Science.* 1991; 253: 390–395.
- Eisenhart C. Assumptions underlying the analysis of variance, *Biometrics.* 1947; 3: 1–21.
- Entsuah AR. Randomization procedures for analyzing clinical trend data with treatment related withdrawls. *Commun. Statist. A.* 1990; 19: 3859–3880.
- Erdos P; Renyi A. On a central limit theorem for samples from a finite population. *Publ. Math. Inst. Hung. Acad. Sci.* 1959; 4: 49–61.
- Evett IW; Gill PD; Scranage JK; and Weir BS. Establishing the robustness short-tandem-repeat statistics for forensic applications. *Amer. J. Human Genetics.* 1996; 58: 398–407.
- Fan CT; Muller ME; and Rezucha I. Development of sampling plans by using sequential (item by item) selection techniques and digital computers. *JASA.* 1962; 57: 387–402.
- Fang KT. The limit distribution of linear permutation statistics and its applications. *Acta Math. Appl. Sinica.* 1981; 4: 69–82.
- Faris PD; Sainsbury RS. The role of the Pontis Oralis in the generation of RSA activity in the hippocampus of the guinea pig. *Psych. Beh.* 1990; 47: 1193–1199.
- Farrar DA; Crump KS. Exact statistical tests for any carcinogenic effect in animal assays. *Fund. Appl. Toxicol.* 1988; 11: 652–663.
- Farrar DA; Crump KS. Exact statistical tests for any carcinogenic effect in animal assays. II age adjusted tests. *Fund. Appl. Toxicol.* 1991; 15: 710–721.
- Fears TR; Tarone RE; and Chu KC. False-positive and false-negative rates for carcinogenicity screens. *Cancer Res.* 1977; 37: 1941–1945.
- Feinstein AR. Clinical biostatistics XXIII. The role of randomization in sampling, testing, allocation, and credulous idolatry (part 2). *Clinical Pharm.* 1973; 14: 989–1019.
- Ferron J; Onghena P. The power of randomization tests for single-case phase designs. *J. Exper. Edu.* 1996; 64: 231–239.
- Ferron J; Ware W. Using randomization tests with responsive single-case designs. *Beh. Res. Therapy* 1994; 32: 787–791.
- Ferron J; Ware W. Analyzing single-case data: The power of randomization tests. *J. Exper. Edu.* 1995; 63: 167–178.
- Festinger LC; Carlsmith JM. Cognitive consequences of forced compliance. *J. Abnorm. Soc. Psych.* 1959; 58: 203–210.
- Finch PD. Description and analogy in the practice of statistics (with disc). *Biometrika.* 1979; 66: 195–205.
- Finch PD. Randomization-I. In *Encyclopedia of Statistical Sciences.* S. Kotz; Johnson NL, Eds. New York: Wiley; 1986: 7: 516–519.
- Finney DJ. Fisher-Yates test of significance in  $2 \times 2$  contingency table. *Biometrika.* 1948; 35: 145–156.
- Fishbein M; Venable DL. Diversity and temporal change in the effective pollinators of *asclepias tuberosa*. *Ecology.* 1996; 77: 1061–1073.
- Fisher NI; Hall P. On bootstrap hypothesis testing. *Australian J. Statist.* 1990; 32: 177–190.
- Fisher RA. *Statistical Methods for Research Workers.* 1st ed. Edinburgh, Scotland: Oliver and Boyd; 1925.
- Fisher RA. The logic of inductive inference (with discussion). *J. Roy. Statist. Soc A.* 1934; 98: 39–54.
- Fisher RA. *Design of Experiments.* New York: Hafner; 1935.
- Fisher RA. Coefficient of racial likeness and the future of craniometry. *J. Royal Anthropol. Soc.* 1936; 66: 57–63.

- Fisher RA. The use of multiple measurements in taxonomic problems. *Ann. Eugenics*. 1936; 7: 179–188.
- Fix E; Hodges JL Jr; and Lehmann EL. The restricted chi-square test. In *Studies in Probability and Statistics Dedicated to Harold Cramer*. Stockholm, Sweden: Almqvist and Wiksell; 1959.
- Folks JL. Use of randomization in experimental research. *Experimental design, Statistical Models, and Genetic Statistics; Essays in Honor of Oscar Kempthorne*. In Hinkelman K, Ed. New York: Marcel-Dekker; 1995: 17–22.
- Ford RD; Colom LV; and Bland BH. The classification of medial septum-diagonal band cells as theta-on or theta-off in relation to hippo campal EEG states. *Brain Res*. 1989; 493: 269–282.
- Forster JJ; McDonald JW; and Smith PWF. Monte Carlo exact conditional tests for log-linear and logistic models. *J. Roy. Statist. Soc. B*. 1996; 58: 445–453.
- Forsythe AB; Engleman L; and Jennrich R. A stopping rule for variable selection in multivariate regression. *JASA*. 1973; 68: 75–77.
- Forsythe AB; Frey HS. Tests of significance from survival data. *Comp. Biomed. Res*. 1970; 3: 124–132.
- Forsythe AB; Haritagn JA. Efficiency of confidence intervals generated by repeated subsample calculations. *Biometrika*. 1970; 57: 629–639.
- Foutz RV. A method for constructing exact tests from test statistics that have unknown null distributions. *J. Statist. Comp. Simul*. 1980; 10: 187–193.
- Foutz RV. Simultaneous randomization tests. *Biometrical J*. 1984; 26: 655–663.
- Foutz RN; Jensen DR; and Anderson GW. Multiple comparisons in the randomization analysis of designed experiments with growth curve responses. *Biometrics*. 1985; 41: 29–37.
- Frank D; Trzos RJ; and Good P. Evaluating drug-induced chromosome alterations. *Mutation Res*. 1978; 56: 311–317.
- Fraser DW. Clustering of disease in population units: an exact test and its asymptotic version. *Amer. J. Epidemiol*. 1983; 118: 732–739.
- Fraumeni JF; Li FP. Hodgkin's disease in childhood: an epidemiological study. *J. Nat. Cancer Inst*. 1969; 42: 681–691.
- Freedman DA. Bootstrap regression models. *Ann. Statist*. 1981; 9: 118–128.
- Freedman D; Lane D. The empirical distribution of Fourier coefficients. *Ann. Statist*. 1980; 8: 1244–1251.
- Freedman D; Lane D. Nonstochastic interpretation of reported significance levels. *J. Bus. Econ. Stat*. 1983; 1: 292–298.
- Freedman L. Using permutation tests and bootstrap confidence limits to analyze repeated events data. *Contr. Clin. Trials*. 1989; 10: 129–141.
- Freeman GH; Halton JH. Note on an exact treatment of contingency, goodness of fit, and other problems of significance. *Biometrika*. 1951; 38: 141–149.
- Friedman M. The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *JASA*. 1937; 32: 675–701.
- Friedman JH; Rafsky LC. Multivariate generalizations of the Wald-Wolfowitz and Smirnov two-sample test. *Ann. Statist*. 1979; 7: 697–717.
- Gabriel KR. Ante-dependent analysis of an ordered set of variables. *Annal. Math. Statist*. 1962; 33: 201–212.
- Gabriel KR. Some statistical issues in weather experimentation. *Commun. Statist. A*. 1979; 8: 975–1015.
- Gabriel KR; Feder P. On the distribution of statistics suitable for evaluating rainfall simulation experiments. *Technometrics*. 1969; 11: 149–160.
- Gabriel KR; Hall WJ. Rerandomization inference on regression and shift effects: Computationally feasible methods. *JASA*. 1983; 78: 827–836.
- Gabriel KR; Hsu CF. Evaluation of the power of rerandomization tests, with application to weather modification experiments. *JASA*. 1983; 78: 766–775.

- Gabriel KR; Sokal RR. A new statistical approach to geographical variation analysis. *Systematic Zoology*. 1969; 18: 259–270.
- Gail MH; Gart JJ. The determination of sample size for use with the exact conditional test in  $2 \times 2$  comparative trials. *Biometrics*. 1973; 29: 441–448.
- Gail MH; Mantel N. Counting the number of  $r \times c$  contingency tables with fixed marginals. *JASA*. 1977; 72: 859–862.
- Gail MH; Tan WY; and Piantadosi S. Tests for no treatment effect in randomized clinical trials. *Biometrika*. 1988; 75: 57–64.
- Galambos J. Exchangeability. In *Encyclopedia of Statistical Sciences*. S. Kotz; Johnson NL, Eds. New York: Wiley; 1986; 7: 573–577.
- Gann P; Chatterton R; Vogelsong K; Dupuis J; and Ellman A. Mitogenic growth-factors in breast fluid obtained from healthy women—evaluation of biological and extraneous sources variability. *Cancer Epidemiology Biomarkers Prevention*. 1997; 6: 421–428.
- Gans LP; Robertson CA. Distributions of Goodman and Kruskal's Gamma and Spearman's Rho in  $2 \times 2$  tables for small and moderate sample sizes. *JASA*. 1981; 76: 942–946.
- Garside GR; Mack C. Actual type I error probabilities for various tests in the homogeneity case of the  $2 \times 2$  contingency table. *Amer. Statist.* 1976; 30: 18–21.
- Gart J. Point and interval estimation of the common odds ratio in the combination of  $2 \times 2$  tables with fixed marginals. *Biometrika*. 1970; 57: 471–475.
- Gart JJ. *Statistical Methods in Cancer Res.*, Vol III—The design and analysis of long term animal experiments. Lyon: IARC Scientific Publications; 1986.
- Garthwaite PH. Confidence intervals from randomization tests. *Biometrics*. 1996; 52: 1387–1393.
- Garthwaite PH; Buckland ST. Generating Monte Carlo confidence intervals by the Robbins-Monro process. *Appl. Statist.* 1992; 41: 159–171.
- Gastwirth JL. Statistical reasoning in the legal setting. *Amer. Statist.* 1992; 46: 55–69.
- Gastwirth JL; Rubin H. Effects of dependence on the level of some one-sample tests. *JASA*. 1971; 66: 816–820.
- Geary RC. Some properties of correlation and regression in a limited universe. *Metron*. 1927; 7: 83.
- Geisser S. The predictive sample reuse method with applications. *JASA*. 1975; 70: 320–328.
- Gerig TM. A multivariate extension of Friedman's chi-square test. *JASA*. 1969; 64: 1595–1608.
- Gerig TM. A multivariate extension of Friedman's chi-square test with random covariates. *JASA*. 1975; 70: 443–447.
- Ghosh MN. Asymptotic distributions of serial statistics and applications to nonparametric tests of hypotheses. *Ann. Math. Statist.* 1954; 25: 218–251.
- Gibbons JD. Permutation tests. In *Encyclopedia of Statistical Sciences*. S. Kotz; Johnson NL, Eds. New York: Wiley; 1986; 6: 690, 740–790.
- Gill DS; Siotani M. On randomization in the multivariate analysis of variance. *J. Statist. Plan. Infer.* 1987; 17: 217–226.
- Gine E; Zinn J. Necessary conditions for a bootstrap of the mean. *Ann. Statist.* 1989; 17: 684–691.
- Glass AG; Mantel N. Lack of time-space clustering of childhood leukemia, Los Angeles County 1960–64. *Cancer Res.* 1969; 29: 1995–2001.
- Glass AG; Mantel N; Gunz FW; and Spears GFS. Time-space clustering of childhood leukemia in New Zealand. *J. Nat. Cancer Inst.* 1971; 47: 329–336.
- Glick BJ. Tests for space-time clustering used in Cancer Research. *Geographical Analysis*. 1979; 11: 202–208.
- Gliddentracey C; Greenwood AK. A validation study of the Spanish self directed search using back translation procedures. *J. Career Assess.* 1997; 5: 105–113.
- Gliddentracey CE; Parraga MI. Assessing the structure of vocational interests among Bolivian university students. *J. Vocational Beh.* 1996; 48: 96–106.

- Goldberg P; Leffert F; Gonzales M; Gorgenola I; and Zerbe GO. Intravenous amino-phylline in asthma: A comparison of two methods of administration in children. *Amer. J. Diseases Children.* 1980; 134: 12–18.
- Gonzales; Manly BJ. Analysis of variance by randomization with small data sets. *Environmetrics.* 1998; 9: 53–65.
- Good IJ. On the analysis of symmetric Dirichlet distributions and their mixtures to contingency tables. *Ann. Statist.* 1976; 4: 1159–1189.
- Good IJ. On the exact distribution of Pearson's chi-square for the lady tasting beer. *J. Statist. Comput. Simul.* 1990; 36: 177–179.
- Good PI. Detection of a treatment effect when not all experimental subjects respond to treatment. *Biometrics.* 1979; 35(2): 483–489.
- Good PI. Review of Edgington's Randomization Tests. *J. Statist. Comput. Simul.* 1980; 11: 157–160.
- Good PI. Almost most powerful tests for composite alternatives. *Comm. Statist.—Theory Methods.* 1989; 18(5): 1913–1925.
- Good PI. Most powerful tests for use in matched pair experiments when data may be censored. *J. Statist. Comp. Simul.* 1991; 38: 57–63.
- Good PI. Globally almost powerful tests for censored data. *Nonpar. Statist.* 1992; 1: 253–262.
- Good PI. Shave and a haircut. *Echoes.* 1994; 1(3): 43–61.
- Good PI. *Resampling Methods.* New York: Birkhauser; 1999.
- Good P; Good S. *Legal Applications of Statistics.* New York: Wiley; to be Published.
- Good P; Kemp P. Almost most powerful test for censored data. *Randomization.* 1969; 2: 25–33.
- Goodall DW. Contingency tables and computers. *Praximetric.* 1968; 9: 113–119.
- Goodman LA; Kruskal WH. *Measures of Association for Cross-Classifications.* New York: Springer-Verlag; 1979.
- Graubard BI; Korn EL. Choice of column scores for testing independence in ordered 2 by K contingency tables. *Biometrics.* 1987; 43: 471–476.
- Graves TS; Pazdan JL. A permutation test analogue to Tarone's test for trend in survival analysis. *J. Statist. Compu. Simul.* 1995; 53: 79–89.
- Green BF. Review of Edgington's *Randomization Tests.* *JASA.* 1981; 76: 495.
- Greenland S. On the logical justification of conditional tests for  $2 \times 2$  contingency tables. *Amer. Statist.* 1991; 45: 248–251.
- Grubbs G. Fiducial bounds on reliability for the two-parameter negative exponential distribution. *Technometrics.* 1971; 13: 873–876.
- Haber M. A comparison of some conditional and unconditional exact tests for  $2 \times 2$  contingency tables. *Commun. Statist. A.* 1987; 18: 147–156.
- Haberman SJ. Log-linear models for frequency tables with ordered classifications. *Biometrics.* 1974; 30: 589–600.
- Hack HRB. An empirical investigation into the distribution of the F-ratio in samples from two non-normal populations. *Biometrika.* 1958; 45: 260–265.
- Hajek J. Limiting distributions in simple random sampling from a finite population. *Publ. Math. Inst. Hung. Acad. Sci.* 1960; 5: 361–374.
- Hajek J. Some extensions of the Wald-Wolfowitz-Noether theorem. *Ann. Math. Statist.* 1961; 32: 506–523.
- Hajek J. Asymptotic normality of simple linear rank statistics under alternatives. *Ann. Math. Statist.* 1968; 39: 325–346.
- Hajek J; Sidak Z. *Theory of Rank Tests.* New York: Academic Press; 1967.
- Hall P. On efficient bootstrap simulation. *Biometrika.* 1989; 76: 613–617.
- Hall P. *The Bootstrap and Edgeworth Expansion.* New York: Springer-Verlag; 1992.
- Hall P; Padmanabhan AR. Adaptive inference for the 2-sample scale problem. *Technometrics.* 1997; 39: 412–422.

- Hall P; Titterington M. The effect of simulation order on level accuracy and power of Monte Carlo tests. *J. Roy. Statist. Soc. B.* 1989; 51: 459–467.
- Hall P; Wilson SR. Two guidelines for bootstrap hypothesis testing. *Biometrics*. 1991; 47: 757–762.
- Halperin M; Ware JH; Byar DP; Mantel N; Brown CC; Koziol J; Gail M; and Green SB. Testing for interaction in an  $I \times J \times K$  contingency table. *Biometrika*. 1977; 64: 271–275.
- Halter JH. A rigorous derivation of the exact contingency formula. *Proc. Cambridge Phil. Soc.* 1969; 65: 527–530.
- Hamilton MA; Collings BJ. Determining the appropriate sample size for nonparametric tests for location shift. *Technometrics*. 1991; 3: 327–337.
- Hampel FR; Ronchetti EM; Rousseeuw PJ; and Stahel WA. *Robust Statistics: The Approach Based on Influence Functions*. New York: Wiley; 1986.
- Hartigan JA. Using subsample values as typical values. *JASA*. 1969; 64: 1303–1317.
- Hasegawa M; Kishino H; and Yano T. Phylogenetic inference from DNA sequence data. In *Statistical Theory and Data Analysis*. K. Matusita, Ed. Amsterdam: North Holland; 1988.
- Hayes AF. Permutation test is not distribution-free: Testing  $H_0: \rho = 0$ . *Psych. Meth.* 1996; 1: 184–198.
- Hayes AF. Cautions in testing variance equality with randomization tests. *J. Statist. Comp. Simul.* 1997; 59: 25–31.
- Healy MJR. Comments on a paper by KD Tocher. *J. Roy. Statist. Soc. B.* 1952; 14: 92.
- Henery RJ. Permutation probabilities for gamma random variables. *J. Appl. Probab.* 1983; 20(4): 822–834.
- Henze N. A multivariate two-sample test based on the number of nearest neighbor coincidence. *Ann. Statist.* 1988; 16: 772–783.
- Hess JC; Elsner JB. Extended range hindcasts of tropical-origin Atlantic hurricane activity. *Geophys. Res. Lett.* 1994; 21: 365–368.
- Hettmansperger TP. *Statistical Inference Based on Ranks*. New York: Wiley; 1984.
- Hiatt WR; Fradl DC; Zerbe GO; Byyny RL; and Niels AS. Comparative effects of selective and nonselective beta blockers on the peripheral circulation. *Clinical Pharmacology Therapeutics*. 1983; 35: 12–18.
- Higgins JJ; Bain PT. Nonparametric tests for lattice-ordered alternatives in unreplicated two-factor experiments. *J. Nonpar. Statist.* 1998; XX: 1–12.
- Higgins JJ; Noble W. A permutation test for a repeated measures design. *Appl. Statist. Agriculture*. 1993; 5: 240–255.
- Higgins JJ; Tashtoush S. An aligned rank transform test for interaction. *Nonlin. World*. 1994; 1: 201–11.
- Highton R. Comparison of microgeographic variation in morphological and electrophoretic traits. In *Evolutionary Biology*. Hecht MK; Steer WC; and B Wallace, Eds. New York: Plenum; 1977; 10: 397–436.
- Hilton JF; Mehta CR. Power and sample size for exact conditional tests with ordered categorical data. *Biometrics*. 1993; 49: 609–616.
- Hilton JF; Mehta CR. Exact power of conditional and unconditional tests. Going beyond the  $2 \times 2$  contingency table. *JASA*. 1993; 47: 91–98.
- Hinkley DV. Comment on article by D Basu. *JASA*. 1980; 75: 582–584.
- Hirji KF. A comparison of exact, mid-P, and score tests for matched case-control studies. *Biometrics*. 1991; 47: 487–496.
- Hirji KF; Mehta CR; and Patel NR. Computing distributions for exact logistic regression. *JASA*. 1987; 82: 1110–1117.
- Hirji KF; Mehta CR, and Patel NR. Exact inference for matched case-control studies. *Biometrics*. 1988; 44: 803–814.

- Hiriji KF; Tan SJ; and Elashoff RM. A quasi-exact test for comparing two binomial proportions. *Statist. Medicine*. 1991; 10: 1137–1153.
- Ho ST; Chen LHY. An  $L_p$  bound for the remainder in a combinatorial central limit theorem. *Annals Probability*. 1978; 6: 231–249.
- Hochberg Y. A sharper Bonferroni procedure for multiple tests of significance. *Biometrika*. 1988; 75: 800–803.
- Hodges JL; Lehmann EL. Testing the approximate validity of statistical hypotheses. *J. Roy. Statist. Soc. B*. 1954; 16: 261–268.
- Hodges JL; Lehmann EL. Estimates of location based on rank tests. *Ann. Math. Statist.* 1963; 34: 598–611.
- Hoeffding W. Combinatorial central limit theorem. *Ann. Math. Statist.* 1951; 22: 556–558.
- Hoeffding W. The large-sample power of tests based on permutations of observations. *Ann. Math. Statist.* 1952; 23: 169–192.
- Hoel DG; Walburg HE. Statistical analysis of survival experiments. *J. Nat. Cancer Inst.* 1972; 49: 361–372.
- Hogg RV; Lenth RV. A review of some adaptive statistical techniques. *Commun. Statist.* 1984; 13: 1551–1579.
- Hollander M; Pena E. Nonparametric tests under restricted treatment assignment rules. *JASA*. 1988; 83(404): 1144–1151.
- Hollander M; Sethuraman J. Testing for agreement between two groups of judges. *Biometrika*. 1978; 65: 403–412.
- Hollander M; Wolfe DA. *Nonparametric Methods in Statistics*. New York: Wiley; 1973.
- Holm S. A simple sequentially rejective multiple test procedure. *Scand. J. Statist.* 1979; 6: 65–70.
- Holmes MC; Williams REO. The distribution of carriers of streptococcus pyrogenes among 2413 healthy children. *J. Hyg. Camb.* 1954; 52: 165–179.
- Hope ACA. A modified Monte Carlo significance test procedure. *J. Roy. Statist. Soc. B*. 1968; 30: 582–598.
- Howard M (pseud for Good P). Randomization in the analysis of experiments and clinical trials. *Amer. Lab.* 1981; 13: 98–102.
- Huber PJ. A robust version of the probability ratio test. *Ann. Math. Statist.* 1965; 36: 1753–1758.
- Hubert LJ. Evaluating object set partitions. *J. Verbal Learn. Behav.* 1976; 15: 459–470.
- Hubert LJ. Seriation using asymmetric proximity measures. *Brit. J. Math Statist. Psych.* 1976; 29: 32–52.
- Hubert LJ. Generalized proximity function comparisons. *Brit. J. Math. Statist. Psych.* 1978A; 31: 179–192.
- Hubert LJ. Nonparametric tests for patterns in geographic variation: possible generalizations. *Geographical Analysis*. 1978B; 10: 86–88.
- Hubert LJ. Matching methods in the analysis of cross-classification. *Psychometrika*. 1979; 44: 21–41.
- Hubert LJ. Generalized concordance. *Psychometrika*. 1979; 44: 3–20.
- Hubert LJ. Combinatorial data analysis: Association and partial association. *Psychometrika*. 1985; 50: 449–467.
- Hubert LJ. *Assignment Methods in Combinatorial Data Analysis*. New York: Marcel-Dekker; 1987.
- Hubert LJ; Baker FB. Data analysis by single-link and complete-link hierarchical clustering. *J. Educ. Statist.* 1976; 1: 87–111.
- Hubert LJ; Baker FB. Analyzing distinctive features confusion matrix. *J. Educ. Statist.* 1977; 2: 79–98.
- Hubert LJ; Baker FB. The comparison and fitting of given classification schemes. *J. Math. Psychol.* 1977; 16: 233–53.

- Hubert LJ; Baker FB. Evaluating the conformity of sociometric measurements. *Psychometrika*. 1978; 43: 31–42.
- Hubert LJ; Golledge RG; and Costanzo CM. Generalized procedures for evaluating spatial autocorrelation. *Geographical Analysis*. 1981; 13: 224–233.
- Hubert LJ; Golledge RG; and Costanzo CM. Analysis of variance procedures based on a proximity measure between subjects. *Psych. Bull.* 1982; 91: 424–430.
- Hubert LJ; Golledge RG; Costanzo CM; Gale N; and Halperin WC. Nonparametric tests for directional data. In *Recent Developments in Spatial Analysis: Methodology, Measurement, Models*. Bahrenberg G; Fischer M; and Nijkamp P, Eds. Aldershot, U.K.: Gower; 1984: 171–190.
- Hubert LJ; Golledge RG; Costanzo CM; and Gale N. Measuring association between spatially defined variables: An alternative procedure. *Geographical Analysis*. 1985; 17: 36–46.
- Hubert LJ; Levin JR. General statistical framework for assessing categorical clustering in free recall. *Psych. Bull.* 1976; 83: 1072–1080.
- Hubert LJ; Levin JR. Inference models for categorical clustering. *Psych. Bull.* 1976; 83: 878–887.
- Hubert LJ; Schultz JR. Maximum likelihood paired comparison ranking and quadratic assessment. *Biometrika*. 1975; 62: 655–660.
- Hubert LJ; Schultz J. Quadratic assignment as a general data analysis strategy. *Brit. J. Math. Statist. Psych.* 1976; 29: 190–241.
- Huitema BE; McKean JW. Autocorrelation estimation and inference with small samples. *Psych. Bull.* 1991; 291–304.
- Hunter MA; May R. Some myths concerning parametric and nonparametric tests. *Canadian Psychol.* 1993; 34: 384–389.
- Ingenbleek JF. Tests simultanés de permutation des rangs pour bruit-blanc multivariable. *Statist. Anal. Données*. 1981; 6: 60–65.
- Irony TZ; Pereira CAB. Exact tests for equality of two proportions: Fisher vs Bayes. *J. Statist. Comput. Simul.* 1986; 25: 83–114.
- Iyer HK; Berry KJ; and Mielke PW. Computation of finite population parameters and approximate probability values for multi-response randomized block permutation procedures (MRPP). *Commun. Statist. B*. 1983; 12: 479–499.
- Izenman AJ. Recent developments in nonparametric density estimation. *JASA*. 1991; 86(413): 205–224.
- Jackson DA. Ratios in aquatic sciences: Statistical shortcomings with mean depth and the morphoedaphic index. *Canadian J. Fisheries Acquatic Sciences*. 1990; 47: 1788–1795.
- James GS. The comparison of several groups of observations when the ratios of the population variances are unknown. *Biometrika*. 1950; 38: 324–329.
- Janssen A. Conditional rank tests for randomly censored data. *Ann. Statist.* 1991; 19: 1434–1456.
- Jennings JM; McIntosh AR; Kapur S; Tulving E; and Houle S. Cognitive subtractions may not add up—the interactions between semantic processing and response-mode. *Neuroimage*. 1997; 5: 229–239.
- Jennrich RI. A note on the behaviour of the log rank permutation test under unequal censoring. *Biometrika*. 1983; 70: 133–137.
- Jennrich RI. Some exact tests for comparing survival curves in the presence of unequal right censoring. *Biometrika*. 1984; 71: 57–64.
- Jin MZ. On the multisample permutation test when the experimental units are nonuniform and random experimental errors exist. *J. System Sci. Math. Sci.* 1984; 4: 117–27, 236–243.
- Jockel KH. Finite sample properties and asymptotic efficiency of Monte Carlo tests. *Ann. Statist.* 1986; 14: 336–347.

- Joe H. Extreme probabilities for contingency tables under row and column independence with applications to Fisher's exact test. *Commun. Statist.—Theory Methods.* 1988; 17: 3677–3685.
- John RD; Robinson J. Significance levels and confidence intervals for randomization tests. *J. Statist. Comput. Simul.* 1983; 16: 161–173.
- John RD; Robinson J. Edgeworth expansions for the power of permutation tests. *Ann. Statist.* 1983; 11: 625–631.
- Jones HL. Investigating the properties of a sample mean by employing random subsample means. *JASA.* 1956; 51: 54–83.
- Jones JS; Selander RK; and Schnell GD. *Biol. J. Linnean Society;* 1980; 14: 359.
- Jones MP; O'Gorman TW; Lemke JH; and Woolson RF. A Monte Carlo investigation of homogeneity tests of the odds ratio under various sample sizes and configurations. *Biometrics.* 1989; 45: 171–181.
- Jorde LB; Rogers AR; Bamshad M; Watkins WS; Krakowiak P; Sung S; Kere J; and Harpending HC. Microsatellite diversity and the demographic history of modern humans. *Proc. Nat. Acad. Sci.* 1997; 94: 3100–3103.
- Kalbfleisch JD. Likelihood methods and nonparametric tests. *JASA.* 1978; 73: 167–170.
- Kalbfleisch JD; Prentice RL. *The Statistical Analysis of Failure Time Data.* New York: Wiley; 1980.
- Kaplan EL; Meier P. Non-parametric estimation from incomplete observations. *JASA.* 1958; 53: 457–81, 562–563.
- Karlin S; Ghandour G; Ost F; Tauare S; and Korph K. New approaches for computer analysis of DNA sequences. *Proc. Nat. Acad. Sci. USA.* 1983; 80: 5660–5664.
- Karlin S; Williams PT. Permutation methods for the structured exploratory data analysis (SEDA) of familial trait values. *Amer. J. Human Genetics.* 1984; 36: 873–898.
- Kazdin AE. Statistical analysis for single-case experimental designs. In *Strategies for Studying Behavioral Change.* M Hersen; DH Barlow, Eds. New York: Pergamon; 1976.
- Kazdin AE. Obstacles in using randomization tests in single-case experiments. *J. Educ. Statist.* 1980; 5: 253–260.
- Keller-McNulty S; Higgins JJ. Effect of tail weight and outliers on power and type I error of robust permutation tests for location. *Commun. Statist.—Theory Methods.* 1987; 16: 17–35.
- Kelly FP; Vonder Haar TH; and Mielke PW. Imagery randomized block analysis (IRBA) applied to the verification of cloud-edge detectors. *J. Atmos. Oceanic. Tech.* 1989; 6: 671–679.
- Kelly ME. Application of the theory of combinatorial chance to the estimation of significance of clustering in free recall. *Brit. J. Math. Statist. Psych.* 1973; 26: 270–280.
- Kempthorne O. *Design and Analysis of Experiments.* New York: Wiley; 1952.
- Kempthorne O. The randomization theory of experimental inference. *JASA.* 1955; 50: 946–967.
- Kempthorne O. Some aspects of experimental inference. *JASA.* 1966; 61: 11–34.
- Kempthorne O. Inference from experiments and randomization. In *A Survey of Statistical Design and Linear Models.* JN Srivastava, Ed. Amsterdam: North Holland; 1975: 303–332.
- Kempthorne O. Why randomize? *J. Statist. Plan. Infer.* 1977; 1: 1–26.
- Kempthorne O. In dispraise of the exact test: reactions. *J. Statist. Plan. Infer.* 1979; 3: 199–213.
- Kempthorne O. Comments on paper by PD Frich. *Biometrika.* 1979; 66: 206–207.
- Kempthorne O; Doerfler TE. The behavior of some significance tests under experimental randomization. *Biometrika.* 1969; 56: 231–248.
- Kempthorne O; Zyskind G; Addelman S; Throckmorton T; and White R. Analysis of variance procedures. Aeronautical Research Laboratory, USAF; 1961.
- Kendall MG; Stuart A; and Ord JK. *Advanced Theory of Statistics.* London, U.K.: Charles Griffin and Co; 1977.

- Kennedy PE. Randomization tests in econometrics. *J. Bus. Econ. Statist.* 1995; 13: 85–95.
- Kennedy PE; Cade BS. Randomization tests for multiple regression. *Comm. Statist.—Simul. Comp.* 1996; 25: 923–936.
- Khan KA; Tracy DS. Fourth exact moment results for MRBP tests with 2 or 3 treatments. *Commun. Statist A.* 1991; 20: 3863–3877.
- Kidd KK; Morar B; Castiglione CM, et al. A global survey of haplotype frequencies and linkage disequilibrium at the DRD2 locus. *Human Genetics.* 1998; 103: 211–227.
- Kim MJ; Nelson CR; and Startz R. Mean revision in stock prices? A reappraisal of the empirical evidence. *Rev. Econ. Stud.* 1991; 58: 515–528.
- Klauber MR. Two-sample randomization tests for space-time clustering. *Biometrics.* 1971; 27: 129–142.
- Klauber MR. Space-time clustering tests for more than two samples. *Biometrics.* 1975; 31: 719–726.
- Klauber MR; Mustacchi A. Space-time clustering of childhood leukemia in San Francisco. *Cancer Res.* 1970; 30: 1969–1973.
- Kleinbaum DB; Kupper LL; and Chambliss LE. Logistic regression analysis of epidemiologic data: theory and practice. *Commun. Statist. A.* 1982; 11: 485–547.
- Klingenberg CP. Individual variation ontogenies—a longitudinal-study of growth and timing. *Evolution.* 1996; 50: 2412–2428.
- Klingenberg CP; Ekau WA. Combined morphometric and phylogenetic analysis an eco-morphological trend-pelagization in antarctic fishes (Perciformes Nototheniidae). *Biological J.—Linnean Soc.* 1996; 59: 143–177.
- Klingenberg CP; Leigh RHB; Keddie BA; and Spence JR. Influence of gut parasites on growth-performance in the water strider gerris-buenoi (Hemiptera Gerridae). *Ecography.* 1997; 20: 29–36.
- Klingenberg CP; Neuenschwander BE; and Flury BD. Ontogeny and individual variation—analysis patterned covariance matrices with common principal components. *Systematic Biol.* 1996; 45: 135–150.
- Knight K. On the bootstrap of the sample mean in the infinite variance case. *Ann. Statist.* 1989; 17: 1168–1173.
- Koch G, Ed. *Exchangeability in Probability and Statistics.* Amsterdam: North Holland; 1982.
- Kolchin VF; Christyakov VP. On a combinatorial central limit theorem. *Theor. Prob. Appl.* 1973; 18: 728–739.
- Kolassa JE; Tanner MS. Approximate conditional inference in exponential families via the Gibbs sampler. *JASA.* 1994; 89: 697–702.
- Konigsberg LW. Comments on matrix permutation tests in the evaluation of competing models for modern human origins. *J. Human Evolution.* 1997; 32: 479–488.
- Koziol JA; Maxwell DA; Fukushima M; Colmer A; and Pilch YHA. Distribution-free test for tumor-growth curve analyses with applications to an animal tumor immunotherapy experiment. *Biometrics.* 1981; 37: 383–390.
- Krehbiel K. Are Congressional committees composed of preference outliers? *Amer. Poli. Sci. Rev.* 1990; 84: 149–163.
- Krewski D; Brennan J; and M Bickis. The power of the Fisher permutation test in 2 by k tables. *Commun. Statist. B.* 1984; 13: 433–448.
- Kryscio RJ; Meyers MH; Prusiner SI; Heise HW; and BW Christine. The space-time distribution of Hodgkin's disease in Connecticut, 1940–1969. *J. Nat. Cancer Inst.* 1973; 50: 1107–1110.
- Lachin JM. Properties of sample randomization in clinical trials. *Contr. Clin. Trials.* 1988; 9: 312–326.
- Lachin JM. Statistical properties of randomization in clinical trials. *Contr. Clin. Trials.* 1988; 9: 289–311.
- Lambert D. Influence functions for testing. *JASA.* 1981; 76: 649–657.
- Lambert D. Qualitative robustness of tests. *JASA.* 1982; 77: 352–357.

- Lambert D. Robust two-sample permutation tests. *Ann. Statist.* 1985; 13: 606–625.
- Lambert D; Hall WJ. Asymptotic lognormality of p-values. *Ann. Statist.* 1982; 10: 44–64.
- Lan K; DeMets DL. Discrete sequential boundaries for clinical trials. *Biometrika*. 1983; 70: 663.
- Lancaster HO. Significance tests in discrete distributions. *JASA*. 1961; 56: 223–234.
- Laroche J; Baran E; and Rasoanandrasana NB. Temporal patterns in a fish assemblage of a semiarid mangrove zone in Madagascar. *J. Fish Biol.* 1997; 51: 3–20.
- Latscha R. Tests of significance in a rxr contingency table; extension of Finney's table. *Biometrika*. 1953; 40: 74–86.
- Lee MLT. Tests of independence against LR dependence in ordered contingency tables. In *Topics in Statistical Dependence*. HW Block; AR Sampson; TH Savits, Eds. Hayward, CA: IMS; 1990; 16: 351–357.
- Lee TJ; Pielke RA; and Mielke PW. Modelling the clear-sky surface energy budget during FIFE 1987. *J. Geophys. Res.* 1995; 100: 25585–25593.
- Leemis LM. Relationships among common univariate distributions. *Amer. Statist.* 1986; 40: 143–146.
- Lefebvre M. Une application des methodes sequentielles aux tests de permutations. *Canadian J. Statist.* 1982; 10: 173–180.
- Legendre P; Legendre L. *Numerical Ecology*. 2<sup>nd</sup> English ed. Amsterdam: Elsevier Science; 1998.
- Lehmann EL. Consistency and unbiasedness of certain nonparametric tests. *Ann. Math. Statist.* 1951; 22: 165–179.
- Lehmann EL. Nonparametric confidence intervals for a shift parameter. *Ann. Math. Statist.* 1963; 34: 1507–1512.
- Lehmann EL. Some concepts of dependence. *Ann. Math. Statist.* 1966; 37: 1137–1153.
- Lehmann EL. *Non-parametrics: Statistical Methods Based on Ranks*. San Francisco, CA: Holden-Day; 1975.
- Lehmann EL. *Testing Statistical Hypotheses*. 2<sup>nd</sup> ed. New York: Wiley; 1986.
- Lehmann EL; Stein C. On the theory of some nonparametric hypotheses. *Ann. Math. Statist.* 1949; 20: 28–45.
- Leonard T. Density estimation, stochastic processes, and p-information (with discussion). *J. Roy. Statist. Soc. B*. 1978; 40: 113–146.
- Leslie PH. A method of calculating the exact probabilities in  $2 \times 2$  contingency tables with small marginal totals. *Biometrika*. 1955; 42: 522–523.
- Levin B; Robbins H. Urn models for regression analysis with application to employment discrimination studies. *Law Contemporary Problems*. 1983; 46: 246–267.
- Levin DA. The organization of genetic variability in Phlox drummondii. *Evolution*. 1977; 31: 477–494.
- Levin JM; Marascuilo LA; and Hubert LJM. Nonparametric randomization tests. In *Single Subject Research: Strategies for Evaluating Change*. JR Kratochwill, Ed. New York: Academic Press; 1978: 167–196.
- Liang KY; Self SG. Test for homogeneity of the odds ratio when the data are sparse. *Biometrika*. 1985; 72: 353–358.
- Liebtrau AM. *Measures of Association*. Newhall, CA: Sage Publications; 1983.
- Light JR; Margolin BH. An analysis of variance of categorical data. *JASA*. 1971; 66: 534–544.
- Lin DY; Wei IJ; and DeMets DL. Exact statistical inference for group sequential trials. *Biometrics*. 1991; 47: 1399–1408.
- Lindsey JK. Likelihood analyses and test for binary data. *Appl. Statist.* 1975; 241: 1–16.
- Liu RY. Bootstrap procedures under some non i.i.d. models. *Ann. Statist.* 1988; 16: 1696–1788.
- Livezey RE. Statistical analysis of general circulation model climate simulation: Sensitivity and prediction experiments. *J. Atmospheric Sci.* 1985; 42: 1139–1149.

- Livezey RE; Chen W. Statistical field significance and its determination by Monte Carlo techniques. *Monthly Weather Review*. 1983; 111: 46–59.
- Lock RH. A sequential approximation to a permutation test. *Commun. Statist. Simul.* 1991; 20: 341–363.
- Loh WY. Estimating an endpoint of a distribution with resampling methods. *Ann. Statist.* 1984; 12: 1543–1550.
- Lorenz J; Eiler JH. Spawning habitat and characteristics of sockeye salmon in the Glacial Taker River, British Columbia and Alaska. *Trans. Amer. Fisheries Soc.* 1989; 18: 495–502.
- Loughin TM; Noble W. A permutation test for effects in an unreplicated factorial design. *Technometrics*. 1997; 39: 180–190.
- Loughin TM; Scherer PN. Testing for association in contingency tables with multiple column responses. *Biometrics*. 1998; 54: 630–637.
- Louis EJ; Dempster ER. An exact test for Hardy-Weinberg and multiple alleles. *Biometrics*. 1987; 43: 805–811.
- Ludbrook J; Dudley H. Why permutation tests are superior to t and F tests in biomedical research. *Amer. Statist.* 1998; 52: 127–132.
- Lunneborg CE. Estimating the correlation coefficient: The bootstrap approach. *Psychol. Bull.* 1985; 98: 209–215.
- Mackay DA; Jones RE. Leaf-shape and the host-finding behavior of two ovipositing monophagous butterfly species. *Ecol. Entom.* 1989; 14: 423–431.
- Macuson R; Nordbrock E. A multivariate permutation test for the analysis of arbitrarily censored survival data. *Biometrical J.* 1981; 23: 461–465.
- Madow WG. On the limiting distribution of estimates based on samples from finite universes. *Ann. Math. Statist.* 1948; 19: 534–545.
- Magnussen S; Boudevyn P. Derivations of stand heights from airborne laser scanner data with canopy-based quantile estimators. *Canadian J. Forest Res.* 1998; 28: 1016–1031.
- Majeed AW; Troy G; Nicholl JP; Smythe A; Reed MWR; Stoddard CJ; Peacock J; and Johnson AG. Randomized prospective single-blind comparison laparoscopic versus small-incision cholecystectomy. *Lancet*. 1996; 347: 989–994.
- Makinodan T; Albright JW; Peter CP; Good PI; and Hedrick ML. Reduced humoral activity in long-lived mice. *Immunology*. 1976; 31: 400–408.
- Makridakis S; Wheelwright SC; and McGee VE. *Forecasting Methods and Applications*. New York: Wiley; 1983.
- Makuch RW; Parks WP. Response to serum antigen level to AZT for the treatment of AIDS. *AIDS Research Human Retroviruses*. 1988; 4: 305–316.
- Manly BFJ. Analysis of polymorphic variation in different types of habitat. *Biometrics*. 1983; 39.
- Manly BFJ. The comparison and scaling of student assessment marks in several subjects. *Appl. Statist.* 1988; 37: 385–395.
- Manly BFJ. *Randomization, Bootstrap and Monte Carlo Methods in Biology*. 2nd ed. London, U.K.: Chapman and Hall; 1997.
- Manly BFJ; McAlevey L; and Stevens D. A randomization procedure for comparing group means in multiple measurements. *Brit. J. Math. Statist. Psychol.* 1986; 39: 183–189.
- Mann HB. Nonparametric tests against trend. *Econometrica*. 1945; 13: 245–259.
- Mann RC; Hand RE Jr. The randomization test applied to flow cytometric histograms. *Computer Programs Biomedicine*. 1983; 17: 95–100.
- Mantel N. The detection of disease clustering and a generalized regression approach. *Cancer Res.* 1967; 27: 209–220.
- Mantel N. Re: “Clustering of disease in population units: an exact test and its asymptotic version.”. *Amer. J. Epidemiol.* 1983; 118: 628–629.
- Mantel N; Bailar JC. A class of permutational and multinomial tests arising in epidemiological research. *Biometrics*. 1970; 26: 687–700.

- Mantel N; Hankey BJ. Programmed analysis of a  $2 \times 2$  contingency table. *Amer. Statist.* 1971; 25: 40–44.
- Mantel N; Valand RS. A technique of nonparametric multivariate analysis. *Biometrics*. 1970; 26: 547–558.
- Marascuilo LA; McSweeney M. *Nonparametric and Distribution-Free Methods for the Social Sciences*. Monterey, CA: Brooks/Cole; 1977.
- Marcus LF. Measurement of selection using distance statistics in prehistoric orang-utan pongo pygamous palaeosumativens. *Evolution*. 1969; 23: 301.
- Mardia KV; Kent JT; and Bibby JM. *Multivariate Analysis*. New York: Academic Press; 1979.
- Maritz JS. A permutation test allowing for missing values. *Australian J. Statist.* 1995; 37: 153–159.
- Maritz JS. *Distribution Free Statistical Methods*. 2nd ed. London, U.K.: Chapman and Hall; 1996.
- Marriott FHC. Barnard's Monte Carlo tests: How many simulations? *Appl. Statist.* 1979; 28: 75–77.
- Marshall-Olds T. Analysis of local variation in plant size. *Ecology*. 1987; 68: 82–87.
- Martin AA; Silva MA. Choosing the optimal unconditioned test for comparing independent proportions. *Comput. Statist. Data Anal.* 1994; 17: 555–574.
- Martin AA; Sanchez S; and Silva MA. Fisher's mid-p-value arrangement in  $2 \times 2$  comparative trials. *Comput. Statist. Data Anal.* 1998; 29: 107–115.
- Martin RL. On the design of experiments under spatial correlation. *Biometrika*. 1986; 73: 247–277.
- Martin-Lof P. Exact tests, confidence regions and estimates. In *Proceeding of the Conference of Foundational Questions in Statistical Inference*. Barndorff-Nielsen O; Blasild P; and Schow G, Eds. Aarhus: Institute of Mathematics, University of Aarhus; 1974; 1: 121–138.
- Maxwell SE; Cole DA. A comparison of methods for increasing power in randomized between-subjects designs. *Psych. Bull.* 1991; 110: 328–337.
- May RB; Hunter MA. Some advantages of permutation tests. *Canadian Psychol.* 1993; 34: 401–407.
- May RB; Masson MEJ; and Hunter MA. Randomization tests: viable alternatives to normal curve tests. *Beh. Res. Meth. Instr. Comp.* 1989; 21: 482–483.
- May RB; Masson MEJ; and Hunter MA. *Application of Statistics in Behavioral Research*. New York: Harper and Row; 1990.
- McCarthy MD. On the application of the z-test to randomized blocks. *Ann. Math. Statist.* 1937; 10: 337–359.
- McCarthy PJ. Psuedo-replication: Half samples. *Rev. Int. Statist. Inst.* 1969; 37: 239–264.
- McDonald LL; Davis BM; and Miliken GA. A nonrandomized unconditional test for comparing two proportions in  $2 \times 2$  contingency tables. *Technometrics*. 1977; 19: 145–158.
- Mcintosh AR; Bookstein FL; Haxby JV; and Grady CL. Spatial pattern-analysis functional brain images using partial least-squares. *Neuroimage*. 1996; 3: 143–157.
- Mckenzie DP; Mackinnon AJ; Peladeau N; Ongena P; Bruce PC; Clarke DM; Harrigan S; and McGorry PD. Comparing correlated kappas by resampling—is one level of agreement significantly different from another? *J. Psychiatric Res.* 1996; 30: 483–492.
- McKinney PW; Young MJ; Hartz A; and Bi-Fong Lee M. The inexact use of Fisher's exact test in six major medical journals. *J. Amer. Med. Assoc.* 1989; 261: 3430–3433.
- McLeod RS; Taylor DW; Cohen A; and Cullen JB. Single patient randomized clinical trials; its use in determining optimal treatment for patients with inflammation of a Kock continent ileostomy reservoir. *Lancet*. 1986; 29: 726–728.
- McQueen G. Long-horizon mean-reverting stock prices revisited. *J. Financial Quant. Anal.* 1992; 27: 1–17.

- Mead R. A test for spatial pattern at several scales using data from a grid of contiguous quadrats. *Biometrics*. 1974; 30: 295–307.
- Meagher TR; Burdick DS. The use of nearest neighbor frequency analysis in studies of association. *Ecology*. 1980; 61: 1253–1255.
- Mehta CR. An interdisciplinary approach to exact inference for contingency tables. *Statist. Sci.* 1992; 7: 167–170.
- Mehta CR; Hilton JF. Exact power of conditional and unconditional tests going beyond the  $2 \times 2$  contingency table. *Amer. Statist.* 1993; 47: 91–98.
- Mehta CR; Patel NR. A network algorithm for the exact treatment of the  $2 \times K$  contingency table. *Commun. Statist. B*. 1980; 9: 649–664.
- Mehta CR; Patel NR. A network algorithm for performing Fisher's exact test in  $r \times c$  contingency tables. *JASA*. 1983; 78: 427–434.
- Mehta CR; Patel NR. Exact inferences for categorical data. In *Encyclopedia of Biostatistics*. New York: Wiley. 1998.
- Mehta CR; Patel NR; and Gray R. On computing an exact confidence interval for the common odds ratio in several  $2 \times 2$  contingency tables. *JASA*. 1985; 80: 969–973.
- Mehta CR; Patel NR; and Senchaudhuri P. Approximately Exact Inference for the Common Odds Ratio in Several  $2 \times 2$  Tables: Comment. *JASA*. 1998; 93(444): 1313–1316.
- Mehta CR; Patel NR; and Senchaudhuri P. Efficient Monte Carlo Methods for Conditional Logistic Regression. *JASA*, March 2000.
- Mehta CR; Patel NR; and Senchaudhuri P. Importance sampling for estimating exact probabilities in permutational inference. *JASA*. 1988; 83: 999–1005.
- Mehta CR; Patel NR; Senchaudhuri P; and Tsiatis AA. Exact permutational tests for group sequential clinical trials. *Biometrics*. 1994; 50: 1042–1053.
- Mehta CR; Patel NR; and Tsiatis AA. Exact significance testing to establish treatment equivalence with ordered categorical data. *Biometrics*. 1984; 40: 819–825.
- Mehta CR; Patel NR; and Wei LJ. Computing exact permutational distributions with restricted randomization designs. *Biometrika*. 1988; 75: 295–302.
- Mehta CR; Walsh SJ. Comparison of exact, mid-P, and Mantel-Haenszel confidence intervals for the common odds ratio across several  $2 \times 2$  contingency tables. *Amer. Statist.* 1992; 46: 146–150.
- Melia KF; Ehlers CL. Signal detection analysis of ethanol effects on a complex conditional discrimination. *Pharm. Biochem Behavior*. 1989; 19(3): 581–584.
- Merrington M; Spicer CC. Acute leukemia in New England. *Brit. J. Prevent. Soc. Med.* 1969; 23: 124–127.
- Micceri T. The unicorn, the normal curve, and other improbable creatures. *Psychol. Bull.* 1989; 105: 156–166.
- Mielke PW. Clarification and appropriate inferences for Mantel and Valand's nonparametric multivariate analysis technique. *Biometrics*. 1978; 34: 277–282.
- Mielke PW. Some parametric, nonparametric and permutation inference procedures resulting from weather modification experiments. *Commun. Statist. A*. 1979; 8: 1083–1096.
- Mielke PW. On asymptotic nonnormality of null distributions of MRPP statistics. *Commun. Statist. A*. 1979; 8: 1541–1550.
- Mielke PW. Meterological applications of permutation techniques based on distance functions. In *Handbook of Statistics*. Krishnaiah PR; Sen PK, Eds. Amsterdam: North Holland; 1984; 4: 813–830.
- Mielke PW Jr. Geometric concerns pertaining to applications of statistical tests in the atmospheric sciences. *J. Atmospheric Sci.* 1985; 42: 1209–1212.
- Mielke PW. Non-metric statistical analysis: Some metric alternatives. *J. Statist. Plan. Infer.* 1986; 13: 377–387.
- Mielke PW.  $L_1$ ,  $L_2$ , and  $L_\infty$  regression models: Is there a difference? *J. Statist. Plan. Infer.* 1986; 13: 430.

- Mielke PW Jr. The application of multivariate permutation methods based on distance functions in the earth sciences. *Earth-Science Rev.* 1991; 31: 55–71.
- Mielke PW; Berry KJ. An extended class of permutation techniques for matched pairs. *Commun. Statist.-Theory Methods.* 1982; 11: 1197–1207.
- Mielke PW; Berry KJ. Asymptotic clarifications, generalizations, and concerns regarding an extended class of matched pairs tests based on powers of ranks. *Psychometrika.* 1985; 48: 483–485.
- Mielke PW; Berry KJ. Non-asymptotic inferences based on the chi-square statistic for  $r \times c$  contingency tables. *J. Statist. Plan. Infer.* 1985; 12: 41–45.
- Mielke PW Jr; Berry KJ. Fisher's exact probability test for cross-classification tables. *Educ. Psych. Measurement.* 1992; 52: 97–101.
- Mielke PW Jr; Berry KJ. Permutation tests for common locations among samples with unequal variances. *J. Educ. Behav. Statist.* 1994; 19: 217–236.
- Mielke PW Jr; Berry KJ. Exact probabilities for first-order and second-order interactions of  $2 \times 2 \times 2$  contingency tables. *Educ. Psych. Measurement.* 1996; 56.
- Mielke PW Jr; Berry KJ. Permutation covariate analysis of residuals based on Euclidean distance. *Psych. Rep.* 1997; 81: 795–802.
- Mielke PW Jr; Berry KJ. Permutation-based multivariate regression analysis: The case for least sum of absolute deviations regression. *Ann. Operat. Res.* 1997; 74: 259–268.
- Mielke PW Jr; Berry KJ. Multivariate tests for correlated data in completely randomized designs. *J. Educ. Behav. Statist.* 1999; 24.
- Mielke PW; Berry KJ; and Brier GW. Application of multiresponse permutation procedures for examining seasonal changes in monthly mean sea-level pressure patterns. *Monthly Weather Rev.* 1981; 109: 120–126.
- Mielke PW; Berry KJ; Brockwell PJ; and Williams JS. A class of nonparametric tests based on multiresponse permutation procedures. *Biometrika.* 1981; 68: 720–724.
- Mielke PW; Berry KJ; and Johnson ES. Multiresponse permutation procedures for a priori classifications. *Commun. Statist.* 1976; A5(14): 1409–1424.
- Mielke PW; Berry KJ; and Medina J. Climax I and II: distortion resistant residual analysis. *J. Appl. Meterology.* 1982; 21: 788–792.
- Mielke PW; Iyer HK. Permutation techniques for analyzing multi-response data from randomized block experiments. *Commun. Statist. A.* 1982; 11: 1427–1437.
- Mielke PW; Sen PK. On asymptotic non-normal null distributions for locally most powerful rank tests statistics. *Commun. Statist. A.* 1981; 10: 1079–1094.
- Miller AJ; Shaw DE; Veitch LG; and Smith EJ. Analyzing the results of a cloud-seeding experiment in Tasmania. *Commun. Statist. A.* 1979; 8: 1017–1047.
- Miller RA; Bookstein F; Vandermeulen J; Engle S; Kim J; Mullins L; and Faulkner J. Candidate biomarkers aging-age-sensitive indexes of immune and muscle function covary in genetically heterogeneous mice. *J. Gerontology A—Biol. Sci. Med. Sci.* 1997; 52: B39–B47.
- Mitchell-Olds T. Quantitative genetics of survival and growth in *Impatiens. capensis*. *Evolution.* 1986; 40.
- Mitchell-Olds T. Analysis of local variation in plant size. *Ecology.* 1987; 68: 82–87.
- Mitra SK. On the F-test in the intrablock analysis of a balanced incomplete block design. *Sankhya.* 1961; 22: 279–284.
- Morgan WM; Blumenstein BA. Exact conditional tests for hierarchical models in multidimensional contingency tables. *J. Roy. Statist. Soc. C.* 1991; 40: 435–442.
- Morrison DF. *Multivariate Statistical Methods.* New York: McGraw-Hill; 1990.
- Motoo M. On the Hoeffding's combinatorial central limit theorem. *Ann. Inst. Statist. Math.* 1957; 8: 145–154.
- Mueller LD; Altenberg L. Statistical inference on measures of niche overlap. *Ecology.* 1985; 66: 1204–1210.
- Mudholkar GS; Hutson AD. Continuity corrected approximations for and “exact” inference with Pearson's chi-square. *J. Statist. Plan. Infer.* 1997; 23: 61–78.

- Mukhopadhyay I. Nonparametric tests for multiple regression under permutation symmetry. *Calcutta Statist. Assoc. Bull.* 1989; 38: 93–114.
- Murphy BP. Comparison of some two-sample tests by means of simulation. *Commun. Statist. Simul.* 1976; B5: 23–32.
- Nelson LS. A randomization test for ordered alternatives. *J. Quality Tech.* 1992; 24: 51–53.
- Nelson LS. Randomization test for linear correlation/regression. *J. Quality Tech.* 1997; 29: 354–356.
- Nelson W. *Applied Life Data Analysis*. New York: Wiley 1982.
- Neyman J. *First Course in Probability and Statistics*. New York: Holt; 1950.
- Neyman J; Scott E. Field galaxies: luminosity, redshift, and abundance of types. Part I: Theory. *Proceedings of the 4th Berkeley Symposium of Mathematical Statistics and Probability*. 1960; 3: 261–276.
- Nguyen TT. A generalization of Fisher's exact test in pxq contingency tables using more concordant relations. *Commun. Statist. B*. 1985; 14: 633–645.
- Nicholson TAJ. A method for optimizing permutation probabilities and its industrial applications. In *Combinatorial Mathematics and its Applications* PJA Welsh, Ed. New York: Academic Press; 1971: 201–217.
- Noble EP; Gottschalk LA; Fallon JH; Ritchie TL; and Wu JC. D-2 dopamine-receptor polymorphism and brain regional glucose-metabolism. *Amer. J. Med. Genetics*. 1997; 74: 162–166.
- Noether GE. On a theorem by Wald and Wolfowitz. *Ann. Math. Statist.* 1949; 20.
- Noether GE. Asymptotic properties of the Wald-Wolfowitz test of randomness. *Ann. Math. Statist.* 1950; 21: 231–246.
- Noether GE. Distribution-free confidence intervals. *Statistica Neer.* 1978; 32: 104–122.
- Noether GE. Distribution-free confidence intervals. *Amer. Statist.* 1972; 26: 39–41.
- Noether GE. Some Simple distribution-free confidence intervals for the center of a symmetric distribution. *JASA*. 1973; 68: 184–188.
- Noether GE. Sample size determination for some common nonparametric tests. *JASA* 1987; 82: 645–647.
- Noreen E. *Computer Intensive Methods for Testing Hypotheses*. New York: Wiley; 1989.
- Oden A; Wedel H. Arguments for Fisher's permutation test. *Ann. Statist.* 1975; 3: 518–520.
- Ogawa J. Effect of randomization on the analysis of a randomized block design. *Ann. Inst. Statist. Math. Tokyo*. 1961; 13: 105–117.
- Ogawa J. On the null distribution of the F-statistic in a randomized block under the Neyman model. *Ann. Math. Statist.* 1963; 34: 1558.
- Ogawa J. *Statistical Theory of the Analysis of Experimental Designs*. New York: Marcel-Dekker; 1974.
- Ogawa J. Exact and approximate sampling distribution of the F-statistic under the randomization procedure. In *A Modern Course on Statistical Distributions in Scientific Work*. GP Patil; S. Kotz; and JK Ord, Eds. Dordret-Holland: Reidel Publishing Company; 1975.
- Ogbonmwan E; Wynn A. Resampling generalized likelihoods. In *Statistical Decision Theory and Related Topics*. SS Gupta; JO Berger. New York: Springer-Verlag; 1988; 1: 133–147.
- Oja H. On permutation tests in multiple regression and analysis of covariance problems. *Australian J. Statist.* 1987; 29: 91–100.
- Onghena P; Edgington ES. Randomization tests for restricted alternating treatments designs. *Beh. Res. Therapy*. 1994; 32: 783–786.
- Onghena P; May RB. Pitfalls in computing and interpreting randomization test p-values: A commentary on Chen and Dunlap. *Beh. Res. Meth. Instr. Comp.* 1995; 27: 408–411.
- Onghena P; Van Damme G. SCRT 1 Single case randomization tests. *Beh. Res. Meth. Instr. Comp.* 1994; 26: 369.
- Onukogu IB. An analysis of variance of nominal data. *Statistica*. 1984; 64: 87–96.

- Opdal SH; Rognum TO; Vege A, et al. Increased number of substitutions in the D-loop of mitochondrial DNA in the sudden infant death syndrome. *Acta Paediatr* 87: 1039–1044.
- O'Reilly FJ; Mielke PW. Asymptotic normality of MRPP statistics from invariance principles of U-statistics. *Commun. Statist. A*. 1980; 9: 629–637.
- O'Sullivan F; Whitney P; Hinshelwood MM; and Hauser ER. Analysis of repeated measurement experiments in endocrinology. *J. Animal Science*. 1989; 59: 1070–1079.
- Pallini A. Estimating probabilities from invariant permutation distributions. *J. Ital. Statist. Soc.* 1994; 3: 77–91.
- Pallini A; Pesarin F. A class of combinations of dependent tests by a resampling procedure. In *Bootstrapping and Related Techniques*. KH Jockel; G Rothe; and W Sendler, Eds. 79–86. Berlin: Springer-Verlag; 1992.
- Parkhurst DF. Arithmetic versus geometric means for environmental concentration data. *Environ. Sci. Tech.* 1998; 32: A92–A98.
- Parzen MI; Wei LJ; and Ying Z. A resampling method based on pivotal estimating functions. *Biometrika*. 1994; 81: 341–350.
- Passing H. Exact simultaneous comparisons with controls in an  $r \times c$  contingency table. *Biometrical J.* 1984; 26: 643–654.
- Patefield WM. Exact tests for trends in ordered contingency tables. *Appl. Statist.* 1982; 31: 32–43.
- Paternoster R; Brame R; Piquero A, et al. The forward specialization coefficient: Distributional properties and subgroup differences. *J. Quant. Criminol.* 1998; 14: 133–154.
- Patil CHK. Cochran's Q test: exact distribution. *JASA*. 1975; 70: 186–189.
- Pearson ES. Some aspects of the problem of randomization. *Biometrika*. 1937; 29: 53–64.
- Pecaric JE; Proschan F; and Tong YL. *Convex Functions, Partial Orderings, and Statistical Applications*. Boston, MA: Academic Press; 1992.
- Penninckx W; Hartmann C; Massart DL; and Smeystersverbeke J. Validation of the calibration procedure in atomic absorption spectrometric methods. *J. Anal. Atomic Spectrometry*. 1996; 11: 237–246.
- Pesarin F. On a nonparametric combination method for dependent permutation tests with applications. *Psychotherapy Psychosomatics*. 1990; 54: 172–179.
- Pesarin F. Multidimensional testing for mixed variables by resampling procedures. *Statistica Applicata*. 1990; 2: 395–406.
- Pesarin F. Some multidimensional testing problems for missing values via a resampling procedures. *Statistica Applicata*. 1991; 3: 567–577.
- Pesarin F. A resampling procedure for a nonparametric combination method of several dependent permutation tests. *J. Ital. Statist. Soc.* 1992; 1: 87–101.
- Pesarin F. Goodness of fit testing for ordered discrete distributions by resampling techniques. *Metron*. 1994; 52: 57–71.
- Pesarin F. A new solution for the generalized Behrens-Fisher Problem. *Statistica*. 1995; 55: 131–146.
- Pesarin F. An almost exact solution for the multi-variate Behrens-Fisher Problem. *Metron*. 1997; 55: 85–100.
- Pesarin F. A nonparametric combination method for dependent permutation tests with application to some problems with repeated measures. In *Industrial Statistics*, Kitsos CP; Edler L, Eds. Heidelberg: Physics-Verlag, 1997: 259–268.
- Pesarin F. *Permutation Testing of Multidimensional Hypotheses*. Cleup Edtrice: Padova; 1997.
- Peritz E. Exact tests for matched pairs: studies with covariates. *Commun. Statist. A*. 1982; 11: 2157–2167 (errata 12: 1209–1210).
- Peritz E. Modified Mantel-Haenszel procedures for matched pairs. *Commun. Statist. A*. 1985; 14: 2263–2285.

- Peto R; Peto J. Asymptotically efficient rank invariant test procedures. *J. Roy. Statist. Soc. A.* 1972; 135: 185–206.
- Petrondas DA; Gabriel RK. Multiple comparisons by rerandomization tests. *JASA*. 1983; 78: 949–957.
- Picard R. Randomization and design. In *Fisher RA, An Appreciation*. Fienberg SE; Hinkley DV, Eds. New York: Springer-Verlag; 1980: 208–213.
- Pike MC; Smith PG. A case-control approach to examine disease for evidence of contagion including diseases with long latent periods. *Biometrics*. 1974; 30: 263–279.
- Pitman EJG. Significance tests which may be applied to samples from any population. *J. Roy. Statist. Soc. Suppl.* 1937; 4: 119–30, 225–232.
- Pitman EJG. Significance tests which may be applied to samples from any population. Part III. The analysis of variance test. *Biometrika*. 1938; 29: 322–335.
- Plackett RL. Random permutations. *J. Roy. Statist. Soc. B*. 1968; 30: 517–534.
- Plackett RL. *Analysis of Categorical Data*. London, U.K.: Griffin; 1974.
- Plackett RL. Analysis of permutations. *Appl. Statist.* 1975; 24: 163–171.
- Plackett RL; Hewlett PS. A unified theory of quantal responses to mixtures of drugs. The fitting to data of certain models for two non-interactive drugs with complete positive correlation of tolerances. *Biometrics*. 1963; 19: 517–531.
- Pollard E; Lakhand KH; and Rothrey P. The detection of density dependence from a series of annual censuses. *Ecology*. 1987; 68: 2046–2055.
- Ponton D; Copp GH. Early dry-season community structure and habitat use of young fish in tributaries of the river sinnamary (French-Guiana South-America) before and after hydrodam operation. *Environ. Biol. Fishes*. 1997; 50: 235–256.
- Prager MH; Hoenig JM. Superposed epoch analysis: A randomization test of environmental effects on recruitment with application to chub mackrel. *Trans. Amer. Fisheries Soc.* 1989; 18: 608–619.
- Praska Rao BLS. *Nonparametric Functional Estimation*. New York: Academic Press; 1983.
- Pratt JW; Gibbons JD. *Concepts of Nonparametric Theory*. New York: Springer; 1981.
- Priesendorfer RW; Barnett TP. Numerical model/reality intercomparison tests using small-sample statistics. *J. Atmospheric Sci.* 1983; 40: 1884–1896.
- Puri ML; Sen PK. A class of rank order tests for a general linear hypothesis. *Ann. Math. Statist.* 1969; 40: 1325–1343.
- Puri ML; Sen PK. On a class of multivariate, multisample rank-order tests. *Sankhya A*. 1966; 28: 353–376.
- Puri ML; Sen PK. *Nonparametric Techniques in Multivariate Analysis*. New York: Wiley; 1971.
- Puri M; Sen PK. *Nonparametric Methods in General Linear Models*. New York: Wiley; 1985.
- Puri ML; Shane HD. Statistical inference in incomplete blocks design. In *Nonparametric Techniques in Statistical Inference*. ML Puri, Ed. Cambridge, U.K.: Cambridge University Cambridge Press; 1970: 131–155.
- Putter J. Treatment of ties in some nonparametric tests. *Ann. Math. Statist.* 1955; 26: 368–386.
- Pyhel N. Distribution-free r-sample tests for the hypothesis of parallelism of response profiles. *Biometric J*. 1980; 22: 703–714.
- Quinn JF. On the statistical detection of cycles in extinctions in the marine fossil record. *Paleobiology*. 1987; 13: 465–478.
- Randles RH; Wolfe DA. *Introduction to the Theory of Nonparametric Statistics*. New York: Wiley, 1979.
- Rao JNK; Bellhouse DR. Optimal estimation of a finite population mean under generalized random permutation models. *J. Statist. Plan. Infer.* 1978; 2: 125–141.
- Rao TJ. Some aspects of random permutation models in finite population sampling theory. *Metrika*. 1984; 31: 25–32.

- Rasmussen J. Estimating correlation coefficients: bootstrap and parametric approaches. *Psych. Bull.* 1987; 101: 136–139.
- Ray W. Logic for a rank test. *Behav. Science*. 1966; 11: 405.
- Raz J. Analysis of repeated measurements using nonparametric smoothing and randomization tests. *Biometrics*. 1989; 45: 851–871.
- Raz J. Testing for no effect when estimating a smooth function by nonparametric regression: a randomization approach. *JASA*. 1990; 85: 132–138.
- Recchia M; Recchetti M. The simulated randomization test. *Computer Programs Biomedicine*. 1982; 15: 111–116.
- Rice WR. A new probability model for determining exact p-values for  $2 \times 2$  contingency tables when comparing binomial proportions. *Biometrics*. 1988; 44: 1–14.
- Richards LE; Byrd J. AS 304: Fishers randomization test for two small independent samples. *Appl. Statist.* 1996; 45: 394–398.
- Ripley BD. *Spatial Statistics*. Wiley: New York; 1981.
- Ritland C; Ritland K. Variation of sex allocation among eight taxa of the *Minimuls guttatus* species complex (Scrophulariaceae). *Amer. J. Botany*. 1989; 76.
- Roberson P; Fisher L. Lack of robustness in time-space disease clustering. *Commun. Statist B: Simulation Computing*. 1986; 12: 11–22.
- Robert V; Awono-Ambene HP; and Thioulouse J. Ecology of larval mosquitoes, with special reference to *Anopheles arabiensis* (Diptera Culicidae) in market-garden wells in urban Dakar, Senegal. *J. Med. Entomol.* 1998; 35: 948–955.
- Robinson J. A converse to a combinatorial central limit theorem. *Ann. Math. Statist.* 1972; 43: 2055–2057.
- Robinson J. The large-sample power of permutation tests for randomization models. *Ann. Statist.* 1973; 1: 291–296.
- Robinson J. On the test for additivity in a randomized block design. *JASA*. 1975; 70: 184–194.
- Robinson J. Large deviation probabilities for samples from a finite population. *Ann. Prob.* 1977; 5: 913–925.
- Robinson J. An asymptotic expansion for samples from a finite population. *Ann. Statist.* 1978; 6: 1005–1011.
- Robinson J. An asymptotic expansion for permutation tests with several samples. *Ann. Statist.* 1980; 8: 851–864.
- Robinson J. Saddlepoint approximations for permutation tests and confidence intervals. *J. Roy. Statist. Soc. B*. 1982; 44: 91–101.
- Robinson J. Approximations to some test statistics for permutation tests in a completely randomized design. *Australian J. Statist.* 1983; 25: 358–369.
- Robinson J. Nonparametric confidence intervals in regression: The bootstrap and randomization methods. In *New Perspectives in Theoretical and Appl. Statistics*. Puri M; Vilaplana JP; Wertz W., Eds. New York: Wiley; 1987: 243–256.
- Robson AJ; Jones TK; Reed DW; and Bayliss AC. A study of national trend and variation in UK floods. *Int. J Climatology*. 1998; 18: 165–182.
- Romano JP. Bootstrap and randomization tests of some nonparametric hypotheses. *Ann. Statist.* 1989; 17: 141–159.
- Romano JP. On the behavior of randomization tests without a group invariance assumption. *JASA*. 1990; 85(411): 686–692.
- Romesburg HC. Exploring, confirming and randomization techniques. *Computers Geosciences*. 1985; 11: 19–37.
- Roper RJ; Doerge RW; Call SB; Tung KSK; Hickey WF; and Teuscher C. Autoimmune orchitis epididymitis and vasitis are immunogenetically distinct lesions. *Amer. J. Path.* 1998; 152: 1337–1345.
- Rosen B. Limit theorems for sampling from a finite population. *Arkiv fur Matematik*. 1965; 5: 383–424.

- Rosenbaum PR. Conditional permutation tests and the propensity score in observational studies. *JASA*. 1984; 79: 565–574.
- Rosenbaum PR. Sensitivity analysis for certain permutation tests in matched observational studies. *Biometrika*. 1987; 74: 13–26.
- Rosenbaum PR. Permutation tests for matched pairs with adjustments for covariates. *Appl. Statist.* 1988; 37(3): 401–411.
- Rosenbaum PR. Sensitivity analysis for matching with multiple controls. *Biometrika*. 1988; 75: 577–581.
- Rosenbaum PR. On permutation tests for hidden biases in observational studies: an application of Holley's inequality to the Savage lattice. *Ann. Statist.* 1989; 17: 643–653.
- Rosenbaum PR. Sensitivity analysis for matched observational studies with many ordered treatments. *Scandinavian J Statist.* 1989; 16: 227–236.
- Rosenbaum PR. Sensitivity analysis for matched case-control studies. *Biometrics*. 1991; 47: 87–100.
- Rosenbaum PR; Krieger AM. Sensitivity analysis of two-sample permutation inferences in observational studies. *JASA*. 1990; 85: 493–498.
- Rounds J; Tracey TJ. Cross cultural structural equivalence of riasec models and measures. *J. Counseling Psych.* 1996; 43: 310–329.
- Royalté HH; Astrachan E; and Sokal RR. Tests for patterns in geographic variation. *Geographic Anal.* 1975; 7: 369–395.
- Rubin DB. Bayesian inference for causal effects: the role of randomization. *Ann. Statist.* 1978; 6: 34–58.
- Runger GC; Eaton MI. Most powerful invariant tests. *J. Multiv. Anal.* 1992; 42: 202–209.
- Ryan JM; Tracey TJG; and Rounds J. Generalizability of Holland's structure of vocational interests across ethnicity, gender, and socioeconomic status. *J. Counseling Psych.* 1996; 43: 330–337.
- Ryman N; Reuterwall C; Nygren K; and Nygren T. Genetic variation and differentiation in Scandinavian moose (*Alces Alces*): Are large mammals monomorphic? *Evolution*. 1980; 34: 1037–1049.
- Saitoh T; Stenseth NC; and Bjornstad ON. Density dependence in fluctuating grey-sided vole populations. *J. Animal Ecol.* 1997; 66: 14–24.
- Salsburg DS. *The Use of Restricted Significance Tests in Clinical Trials*. New York: Springer-Verlag; 1992.
- Sampford MR; Taylor J. Censored observations in randomized block experiments. *J. Roy. Statist. Soc. B*. 1959; 21: 214–237.
- Sandelius M. A simple randomization procedure. *J. Roy. Statist. Soc. B*. 1962; 24: 472–481.
- Santner TJ; Snell MK. Small-sample confidence intervals for  $p_1 - p_2$  and  $p_1/p_2$  in  $2 \times 2$  contingency tables. *JASA*. 1980; 75: 386–394.
- Sawilowsky SS. Nonparametric test of interaction in experimental design. *Rev. Educ. Res.* 1990; 60: 91–126.
- Scheffe H. Statistical inference in the nonparametric case. *Ann. Math. Statist.* 1943; 14: 305–332.
- Scheffe H. *Analysis of Variance*. New York: Wiley; 1959.
- Schemper M. A survey of permutation tests for censored survival data. *Commun. Statist. A*. 1984; 13: 433–448.
- Schemper M. A generalization of the intraclass tau correlation for tied and censored data. *Biometrical J.* 1984; 26: 609–617.
- Schmeiser B; Deutsch SJ. Computation of the component randomization test for paired comparisons. *J. Quality Tech.* 1983; 15: 94–98.
- Schmid PE; Schmidaraya JM. Predation on meiobenthic assemblages-resource use of a tanypod guild (Chironomidae diptera) in a gravel stream. *Freshwater Biol.* 1997; 38: 67–91.

- Schrage C. Evaluation of permutation tests by means of normal approximation or Monte Carlo methods. *Comput. Statist. Quart.* 1984; 1: 325–332.
- Schulman RS. Ordinal data: an alternative distribution. *Psychometrika*. 1979; 44: 3–20.
- Schultz JR; Hubert L. A nonparametric test for the correspondence between two proximity matrices. *J. Educ. Statist.* 1976; 1: 59–67.
- Scott RJ. The influence of parental care behaviour on localized nest spacing in smallmouth bass micropterus dolomieu. *Environ. Biol. Fishes*. 1996; 46: 103–107.
- Selander RK; Kaufman DW. Genetic structure of populations of the brown snail (*Helix aspersa*). I Microgeographic variation. *Evolution*. 1975; 29: 385–401.
- Sen PK. On some permutation tests based on U-statistics. *Bull. Calcutta Statist. Assoc.* 1965; 14: 106–126.
- Sen PK. On some multisample permutation tests based on a class of U-statistic. *JASA*. 1967; 62: 1201–1213.
- Sen PK. Nonparametric tests for multivariate inter-changeability. Part 1: Problems of location and scale in bivariate distributions. *Sankhya A*. 1967; 29: 351–372.
- Sen PK. Nonparametric tests for multivariate inter-changeability. Part 2. The problem of MANOVA in two-way layouts. *Sankhya A*. 1969; 31: 145–156.
- Sen PK. On permutational central limit theorems for general multivariate linear statistics. *Sankhya A*. 1983; 45: 141–149.
- Sen PK; Puri ML. On the theory of rank order tests for location in the multivariate one sample problem. *Ann. Math. Statist.* 1967; 38: 1216–1228.
- Senchaudhuri P; Mehta CR; and Patel NT. Estimating exact p-values by the method of control variates or Monte Carlo rescue. *JASA*. 1995; 90: 640–648.
- Servy EC; Sen PK. Missing variables in multi-sample rank permutation tests for MANOVA and MANCOVA. *Sankhya A*. 1987; 49: 78–95.
- Shane HD; Puri ML. Rank order tests for multivariate paired comparisons. *Ann. Math. Statist.* 1969; 40: 2101–2117.
- Shapiro CP; Hubert L. Asymptotic normality of permutation probabilities derived from the weighted sums of bivariate functions. *Ann. Statist.* 1979; 7: 788–794.
- Shen CD; Quade D. A randomization test for a three-period three-treatment crossover experiment. *Commun. Statist. B*. 1986; 12: 183–199.
- Shiraishi TA. Studentized robust statistics for main effects in a two-factor manova. *Commun. Statist.—Theory Methods*, 1999; 28: 809–823.
- Shiraishi TA. Studentized robust statistics in multi-variate randomized block design. *J. Nonpar. Statist.* 1999; 10: 95–110.
- Shorack G. Testing and estimating ratios of scale parameters. *JASA*. 1969; 64: 999–1013.
- Shuster JJ. *Practical Handbook of Sample Size Guidelines for Clinical Trials*. Boca Raton, FL: CRC Press; 1993.
- Shuster JJ; Boyett JM. Nonparametric multiple comparison procedures. *JASA*. 1979; 74: 379–382.
- Siegmund H. *Sequential Analysis: Tests and Confidence Intervals*. New York: Springer-Verlag; 1985.
- Siemiatycki J. Mantel's space-time clustering statistic: computing higher moments and a comparison of various data transforms. *J. Statist. Comput. Simul.* 1978; 7: 13–31.
- Siemiatycki J; McDonald AD. Neural tube defects in Quebec: A search for evidence of ‘clustering’ in time and space. *Brit. J. Prev. Soc. Med.* 1972; 26: 10–14.
- Siegel S. *Practical Nonparametric Statistics*. New York: Wiley; 1956.
- Silvey SD. The equivalence of asymptotic distributions arising under randomization and normal theories. *Proc. Glasgow Math. Assoc.* 1954; 1: 139–147.
- Silvey SD. Asymptotic distributions of statistics arising in certain nonparametric tests. *Proc. Glasgow Math. Assoc.* 1956; 2: 47–51.
- Simon JL. *Basic Research Methods in Social Science*. New York: Random House; 1969.
- Simon R. Restricted randomization designs in clinical trials. *Biometrics*. 1979; 35: 503–512.

- Simonsen KL; Kaplan NL; and Martin ER. A Monte-Carlo permutation approach to choosing an affection status model for bipolar affective-disorder. *Genetic Epidemiol.* 1997; 14: 681–686.
- Slud E; Wei LJ. Two-sample repeated significance tests based on the modified Wilcoxon statistic. *JASA.* 1982; 77: 862–868.
- Smirnov NV. Estimate of deviation between empirical distribution functions in two independent samples. *Bull. Moscow U.* 1939; 2: 3–16.
- Smith PG; Pike MC. Generalization of two tests for the detection of household aggregation of disease. *Biometrics.* 1976; 32: 817–828.
- Smith PL; Johnson LR; Priegnitz DL; Boew BA; and Mielke PW. An exploratory analysis of crop-hail insurance data for evidence of cloud-seeding effects in North Dakota. *J. Appl. Meteor.* 1997; 36: 463–473.
- Smith PWF; Forster JJ; and McDonald JW. Monte Carlo exact tests for square contingency tables. *J. Royal Statist. Soc. A.* 1996; 159: 309–321.
- Smith PWF; McDonald JW; Forster JJ; and Berrington AM. Monte Carlo exact methods used for analysing inter-ethnic unions in Great Britain. *Appl. Statist.* 1996; 45: 191–202.
- Smith RL. Sequential treatment allocation using biased coin designs. *J. Roy. Statist. Soc. B.* 1984; 46: 519–543.
- Smith RL. Properties of biased coin designs in sequential clinical trials. *Ann. Statist.* 1984; 12: 1018–1034.
- Smouse PE; Long JC; and Sokal RR. Multiple regression and correlation extensions of the Mantel test of matrix correspondence. *Syst. Zool.* 1968; 35: 627–632.
- Smythe RT. Conditional inference for restricted randomization designs. *Ann. Math. Statist.* 1988; 16: 1155–1161.
- Smythe RT; Wei LJ. Significance tests with restricted randomization design. *Biometrika.* 1983; 70: 496–500.
- Sokal RR. Testing statistical significance in geographical variation patterns. *Syst. Zool.* 1979; 28: 227–232.
- Sokal RR; Rohlf FJ. *Biometry.* San Francisco, CA: Freeman; 1981.
- Solomon H. Confidence intervals in legal settings. In *Statistics and The Law.* DeGroot MH; Fienberg SE; and Kadane JB, Eds. New York: Wiley; 1986: 455–473.
- Solow AR. A randomization test for independence of animal locations. *Ecology.* 1989; 70.
- Solow AR. A randomization test for misclassification problems in discriminatory analysis. *Ecology.* 1990; 71: 2379–2382.
- Soms AP. Permutation tests for k-sample binomial data with comparisons of exact and approximate P-levels. *Commun. Statist. A.* 1985; 14: 217–233.
- Soms AP. Exact confidence intervals, based on the z statistic, for the differences between two proportions. *Commun. Statist. Simul. Comput.* 1989; 18: 1325–1341.
- Soms AP. Some recent results for exact confidence intervals for the differences between two proportions. *Commun. Statist. Simul. Comput.* 1989; 18: 1343–1357.
- Spino C; Pagano M. Efficient calculation of the permutation distribution of trimmed means. *JASA.* 1991; 86: 729–737.
- Spino C; Pagano M. Efficient calculation of the permutation distribution of robust two-sample statistics. *Comp. Statist. Data Anal.* 1991; 12: 349–365.
- Spitz MR; Shi HH; Yang F; Hudmon KS; Jiang H; Chamberlain RM; Amos CI; Wan Y; Cinciripini P; Hong WK; and Wu XF. Case-control study the d2-dopamine-receptor gene and smoking status in lung-cancer patients. *J. Nat. Cancer Inst.* 1998; 90: 358–363.
- Sprent P. *Data Driven Statistical Methods.* London, U.K.: Chapman and Hall; 1998.
- Steyn HS; Stumpf RH. Exact distributions associated with an h<sub>x</sub>k contingency table. *South African Statist. J.* 1984; 18: 135–159.
- Still AW; White AP. The approximate randomization test as an alternative to the F-test in the analysis of variance. *Brit. J. Math. Statist. Psych.* 1981; 34: 243–252.

- Stilson DW. *Psychology and Statistics in Psychological Research and Theory*. San Francisco, CA: Holden Day; 1966.
- Stine RA. Estimating properties of autoregressive forecasts. *JASA*. 1987; 82: 1072–1078.
- Stone M. Cross-validation choice and assessment of statistical predictions. *JASA*. 1974; B36: 111–147.
- Storer BE; Kim C. Exact properties of some exact test statistics for comparing two binomial populations. *JASA*. 1990; 85: 146–155.
- Stuart GW; Maruff P; and Currie J. Object-based visual-attention in luminance increment detection. *Neuropsychologia*. 1997; 35: 843–853.
- Stucky W; Vollmar J. Ein verfahren zur exakten awwertung von  $r \times c$ -haufigkeitstatein. *Biom. Zeit.* 1975; 17: 147–162.
- Stumpf RH; Steyn HS. Exact distributions associated with an  $I \times J \times K$  contingency table. *Commun. Statist.—Theory Methods*. 1986; 15: 1889–1904.
- Suissa S; Shuster JJ. Exact unconditional sample sizes for the  $2 \times 2$  binomial trial. *J. Roy. Statist. Soc. A*. 1985; 148: 317–327.
- Suissa S; Shuster JJ. Are uniformly most powerful unbiased tests really best? *Amer. Statist.* 1984; 38: 204–206.
- Swofford DL; Thorne JL; Felsenstein J; and Wiegmann BM. The topology-dependent permutation test for monophyly does not test for monophyly. *Systematic Biol.* 1996; 45: 575–579.
- Syrjala SE. A statistical test for a difference between the spatial distributions of two populations. *Ecology*. 1996; 77: 75–80.
- Takauechi KEI. Asymptotically efficient tests for location: nonparametric and asymptotically nonparametric. In *Nonparametric Techniques in Statistical Inference*. ML Puri, Ed. Cambridge, U.K.: Cambridge University Press; 1970: 131–155.
- Tardif S. On the almost sure convergence of the permutation distribution for aligned rank test statistics in randomized block designs. *Ann. Statist.* 1981; 9: 190–193.
- Tarter ME; Lock MD. *Model-Free Curve Estimation*. New York: Chapman and Hall; 1993.
- ter Braak CJF. Permutation versus bootstrap significance tests in multiple regression and ANOVA. In *Bootstrapping and Related Techniques*. KH Jockel, G Rothe; and W Sendler Eds. Berlin: Springer-Verlag; 1992: 79–86.
- Teuscher C; Rhein DM; Livingstone KD; Paynter RA; Doerge RW; Nicholson SM; and Melvold RW. Evidence that tmevd2 and eae3 may represent either a common locus or members a gene-complex controlling susceptibility to immunologically mediated demyelination in mice. *J Immunol.* 1997; 159: 4930–4934.
- Thomas DG. Exact confidence limits for the odds ratio in a  $2 \times 2$  table. *J. Roy. Statist. Soc. C*. 1971; 20: 105–110.
- Thornett ML. The role of randomization in model-based inference (with discussion). *Australian J. Statist.* 1982; 24: 137–150.
- Titterington DM; Murray GD; Spiegelhalter DJ; Skene AM; Habbema JDF; and Gelke GJ. Comparison of discrimination techniques applied to a complex data set of head-injured patients. *J. Roy. Statist. Soc. A*. 1981; 144: 145–175.
- Tocher KD. Extension of the Neyman-Pearson theory of tests to discontinuous variates. *Biometrika*. 1950; 37: 1301–1444.
- Tracy DS; Khan KA. Comparison of some MRPP and standard rank tests for three equal sized samples. *Commun. Statist. B*. 1990; 19: 315–333.
- Tracy DS; Khan KA. Fourth exact moment results for MRPB and related power performance. *Commun. Statist. A*. 1991; 20: 2701–2718.
- Tracy DS; Tajuddin IH. Extended moment results for improved inferences based on MRPP. *Commun. Statist. A*. 1985; 14: 1485–1496.
- Tracy DS; Tajuddin IH. Empirical power comparisons of two MRPP rank tests. *Commun. Statist. A*. 1986; 15: 551–570.

- Tritchler D. On inverting permutation tests. *JASA*. 1984; 79: 200–207.
- Troendle JF. A stepwise resampling method of multiple hypothesis testing. *JASA*. 1995; 90: 370–378.
- Tsutakawa RK; Yang SL. Permutation tests applied to antibiotic drug resistance. *JASA*. 1974; 69: 87–92.
- Tucker DF; Mielke PW; and Reiter ER. The verification of numerical models with multivariate randomized block procedures. *Meteorol. Atmosph. Phys.* 1989; 40: 181–188.
- Tukey JW. Dyadic ANOVA: An analysis of variance for vectors. *Human Biol.* 1949; 21: 65–110.
- Tukey JW. Improving crucial randomized experiments—especially in weather modification by double randomization and rank combination. In *Proceeding of the Berkeley Conference in Honor of J Neyman and J Kiefer*. LeCam L; Binckley P, Eds. Heyward, CA: Wadsworth; 1985; 1: 79–108.
- Tukey JW; Brillinger DR; and Jones LV. *Management of Weather Resources: Vol II: The Role of Statistics in Weather Resources Management*. Washington, DC: Department of Commerce, US Government Printing Office; 1978.
- Turnbull BW; Iwano EJ; Burnett WS; Howe HL; and Clark LC. Monitoring for clusters of disease: applications to leukemia incidence in upstate New York. *Amer. J. Epidemiology*. 1990; 132: S136–S143.
- Upton GJG. A comparison of alternative tests for the  $2 \times 2$  comparative trial. *J. Roy. Statist. Soc. A*. 1982; 145: 86–105.
- Upton GJG. On Mead's test for pattern. *Biometrics*. 1984; 40: 759–766.
- Upton GJG; Brook D. Determination of the optimum position on a ballot paper. *Appl. Statist.* 1975; 24: 279–287.
- Vadiveloo J. On the theory of modified randomization tests for nonparametric hypothesis. *Commun. Statist.—Theory Methods*. 1983; 12: 1581–1596.
- Valdes-Perez RE. Some recent human-computer studies in science and what accounts for them. *AI Magazine*. 1995; 16: 37–44.
- Valdes-Perez RE; Stone CA. Systematic detection of subtle spatio-temporal patterns in time-lapse imaging: II. Particle migrations. *Bioimaging*. 1998; 6: 71–78.
- van der Voet H. Comparing the predictive accuracy of models using a simple randomization test. *Chemometrics Intelligent Lab. Sys.* 1994; 25: 313–323.
- van Eltern P; Veerbeek A. De correlatie-coëfficiënt en permutatietoetsen. *Vereniging voor Statistiek Bull.* 1979; 12: 35–44.
- vanKeerberghen P; Vandebosch C; Smeyers-Verbeke J; and Massart DL. Some robust statistical procedures applied to the analysis of chemical data. *Chemometrics Intelligent Lab. Sys.* 1991; 12: 3–13.
- Vanlier JB. Limitations of thermophilic anaerobic waste-water treatment and the consequences for process design. *Antonie Van Leeuwenhoek Int. J. Gen. Molec. Microbiol.* 1996; 69: 1–14.
- van-Putten B. On the construction of multivariate permutation tests in the multivariate two-sample case. *Statist. Neer.* 1987; 41: 191–201.
- Vecchia DF; Iyer HK. Exact distribution-free tests for equality of several linear models. *Commun. Statist. A*. 1989; 18: 2467–2488.
- Vecchia DF; Iyer HK. Moments of the quartic assignment statistic with an application to multiple regression. *Commun. Statist.—Theory Methods*. 1991; 20: 3253–3269.
- Vollset SE; Hirji KF; and Afifi AA. Evaluation of exact and asymptotic interval estimators in logistic analysis of matched case-control studies. *Biometrics*. 1991; 47: 1311–1325.
- Wald A. *Statistical Decision Functions*. New York: Wiley; 1950.
- Wald A; Wolfowitz J. An exact test for randomness in the nonparametric case based on serial correlation. *Ann. Math. Statist.* 1943; 14: 378–388.
- Wald A; Wolfowitz J. Statistical tests based on permutations of the observations. *Ann. Math. Statist.* 1944; 15: 358–372.

- Walker DD; Loftis JC; and Mielke PW. Permutation methods for determining the significance of spatial dependence. *Math. Geol.* 1997; 29: 1011–1024.
- Wampold BE; Furlong MJ. Randomization tests in single-subject designs: illustrative examples. *J. Behav. Assess.* 1981; 3: 329–341.
- Wampold BE; Worsham NL. Randomization tests for multiple-baseline designs. *Behav. Assess.* 1986; 8: 135–143.
- Wan Y; Cohen J; and Guerra R. A permutation test for the robust sib-pair linkage method. *Ann. Human Genetics.* 1997; 61: 79–87.
- Wang FT; Scott DW. The L<sub>1</sub> method for robust nonparametric regression. *JASA.* 1994; 89: 65–76.
- Watson GS. Sufficient statistics, similar regions, and distribution-free tests. *J. Roy. Statist. Soc. B.* 1957; 19: 262–267.
- Wei LJ. Exact two-sample permutation tests based on the randomized play-the-winner rule. *Biometrika.* 1988; 75: 603–605.
- Wei LJ; Lachin JM. Properties of urn-randomization in clinical trials. *Controlled Clinical Trials.* 1988; 9: 345–364.
- Wei LJ; Smythe RT; Lin DY; and Park TS. Statistical inference with data-dependent treatment allocation rules. *JASA.* 1990; 85: 156–162.
- Wei LJ; Smythe RT; and Smith RL. K-treatment comparisons in clinical trials. *Ann. Math. Statist.* 1986; 14: 265–274.
- Welch BL. On the z-test in randomized blocks and Latin squares. *Biometrika.* 1937; 29: 21–52.
- Welch BL. On tests for homogeneity. *Biometrika.* 1938; 30: 149–158.
- Welch WJ. Rerandomizing the median in matched-pairs designs. *Biometrika.* 1987; 74: 609–614.
- Welch WJ. Construction of permutation tests. *JASA.* 1990; 85(411): 693–698.
- Welch WJ; Guitierrez LG. Robust permutation tests for matched pairs designs. *JASA.* 1988; 83: 450–461.
- Wellner JA. Permutation tests for directional data. *Ann. Statist.* 1979; 7: 929–943.
- Werner M; Tolls R; Hultin J; and Mellecker J. Sex and age dependence of serum calcium, inorganic phosphorous, total protein, and albumin in a large ambulatory population. In *Fifth Technical International Congress of Automation, Advances in Automated Analysis.* Mount Kisco, NY: Future Publishing; 1970.
- Westfall DH; Young SS. *Resampling-Based Multiple Testing: Examples and Methods for p-value Adjustment.* New York: Wiley; 1993.
- Weth F; Nadler W; and Korschning S. Nested expression domains for odorant receptors in zebrafish olfactory epithelium. *Proc. Nat. Acad. Sci. USA.* 1996; 93: 13321–13326.
- Wetherill GB. The Wilcoxon test and nonnull hypothesis. *J. Roy. Statist. Soc. B.* 1960; 2: 402–418.
- Whaley FS. The equivalence of three individually derived permutation procedures for testing the homogeneity of multidimensional samples. *Biometrics.* 1983; 39: 741–745.
- White AP; Still AW. Monte Carlo analysis of variance. In *Proceedings of the 6th Symposium in Computational Statistics.* Havranek P; Sidak Z; and Novak M. Wien: Physica-Verlag; 1984.
- Wilk MB. The randomization analysis of a generalized randomized block design. *Biometrika.* 1955; 42: 70–79.
- Wilk MB; Kempthorne O. Some aspects of the analysis of factorial experiments in a completely randomized design. *Ann. Math. Statist.* 1956; 27: 950–984.
- Wilk MB; Kempthorne O. Nonadditivities in a Latin square design. *JASA.* 1957; 52: 218–236.
- Williams-Blangero S. Clan-structured migration and phenotypic differentiation in the Jirels of Nepal. *Human Biol.* 1989; 61: 143–157.
- Willmes K. A comparison between the Lehmacher and Wall rank tests and Pyhel's permu-

- tation test for the analysis of  $r$  independent samples of response curves. *Biometrical J.* 1982; 24: 717–722.
- Wilson HG. Least squares versus minimum absolute deviations estimation in linear models. *Decision Sci.* 1978; 9: 322–335.
- Witting H. On the theory of nonparametric tests. In *Nonparametric Techniques in Statistical Inference*. ML Puri, Ed. Cambridge, U.K.: Cambridge University Press; 1970: 41–51.
- Witztum D; Rips E; and Rosenberg Y. Equidistant letter sequences in the Book of Genesis. *Statist. Science*. 1994; 89: 768–776.
- Wong RKW; Chidambaram N; and Mielke PW. Applications of multi-response permutation procedures and median regression for covariate analyses of possible weather modification effects on hail responses. *Atmosphere Ocean*. 1983; 21: 1–13.
- Wu JC; Bell K; Najafi A; Widmark C; Keator D; Tang C; Klein E; Bunney BG; Fallon J; and Bunney WE. Decreasing striatal 6-fdopa uptake with increasing duration cocaine withdrawal. *Neuropsychopharmacology*. 1997; 17: 402–409.
- Wu XF; Amos CI; Kemp BL; Shi HH; Jiang H; Wan Y; and Spitz MR. Cytochrome-p450 2E1 Drai polymorphisms in lung-cancer in minority populations. *Cancer Epidemiology Biomarkers Prevention*. 1998; 7: 13–18.
- Yanagimoto T; Okamoto M. Partial orderings for permutations and monotonicity of a rank correlation statistic. *Inst. Statist. Math. Ann.* 1969; 21: 489–506.
- Yates F. Tests of significance for  $2 \times 2$  contingency tables (with discussion). *J. Roy. Statist. Soc. A*. 1984; 147: 426–463.
- Young A. Conditional data-based simulations. Some examples from geometric statistics. *Int. Statist. Rev.* 1986; 54: 1–13.
- Yucesan E. Randomization tests for initialization bias in simulation output. *Naval Res. Logistics*. 1993; 40: 643–663.
- Zelen M. The analysis of several  $2 \times 2$  contingency tables. *Biometrika*. 1971; 58: 129–137.
- Zelen M. Exact significance tests for contingency tables embedded in a  $2^n$  classification. In L. LeCam; J. Neyman Eds. *Proc 6th Berkeley Symp. Math. Statist. Probab.* Berkeley, CA: University of California Press; 1972; 1: 737–757.
- Zelterman D. Goodness-of-fit for large sparse multinomial distributions. *JASA*. 1987; 82: 624–629.
- Zelterman D; Chan IS-F; and Mielke PW Jr. Exact tests of significance in higher dimensional tables. *The American Statistician*. 1995; 49: 357–361.
- Zempo N; Kayama N; Kenagy RD; Lea HJ; and Clowes AW. Regulation of vascular smooth-muscle-cell migration and proliferation in vitro and in injured rat arteries by a synthetic matrix metalloproteinase inhibitor. *Art. Thromb. V*. 1996; 16: 28–33.
- Zerbe GO. Randomization analysis of the completely randomized design extended to growth and response curves. *JASA*. 1979; 74: 215–221.
- Zerbe GO. Randomization analysis of randomized block design extended to growth and response curves. *Comm. Statist. A*. 1979; 8: 191–205.
- Zerbe GO; Murphy JR. On multiple comparisons in the randomization analysis of growth and response curves. *Biometrics*. 1986; 42: 795–804.
- Zerbe GO; Walker SH. A randomization test for comparison of groups of growth curves with different polynomial design matrices. *Biometrics*. 1977; 33: 653–657.
- Zimmerman DL. A bivariate Cramer-vonMises type of test for spatial randomness. *Appl. Statist.* 1993; 42: 43–54.
- Zimmerman GM; Goetz H; and Mielke PW Jr. Use of an improved statistical method for group comparisons to study effects of prairie fire. *Ecology*. 1985; 66: 606–611.
- Zumbo BD. Randomization test for coupled data. *Perception Psychophysics*. 1996; 58: 471–478.
- Zumbo BD; Hubley AM. A note on misconceptions concerning prospective and retrospective power. *Statistician*. 1998; 47: 385–388.

## Bibliography Part 2:

# Computational Methods

- Abramson M; Moser WJ. Arrays with fixed row and column sums. *Discrete Math.* 1973; 6: 1–14.
- Agresti A; Mehta CR; and Patel NR. Exact inference for contingency tables with ordered categories. *JASA.* 1990; 85: 453–458.
- Akl SG. A comparison of combination generation methods. *ACM Trans. Math. Software.* 1981; 7: 42–45.
- Amana IA; Koch GG. A macro for multivariate randomization analysis of stratified sample data. *SAS Sugi.* 1980; 5: 134–144.
- Arbuckle J; Astler LS. A program for Pitman's permutation test for differences in location. *Behav. Res. Meth. Instr.* 1975; 7: 381.
- Baglivo J; Oliver D; and Pagano M. Methods for the analysis of contingency tables with large and small cell counts. *JASA.* 1988; 83: 1006–1013.
- Baglivo J; Olivier D; and Pagano M. Methods for exact goodness-of-fit tests. *JASA.* 1992; 87: 464–469.
- Baker FB; Collier RO. Analysis of experimental designs by means of randomization, a Univac 1103 program. *Behav. Science.* 1961; 6: 369.
- Baker RD; Tilbury JB. Algorithm AS 283: Rapid computation of the permutation paired and grouped t-tests. *Appl. Statist.* 1993; 42: 432–441.
- Baker RJ. Exact distributions derived from two-way tables. *Appl. Statist.* 1977; 26: 199–206.
- Balmer DW. Recursive enumeration of  $r \times c$  tables for exact likelihood evaluation. AS 236. *Appl. Statist.* 1988; 37: 290–301.
- Bebington AC. A simple method of drawing a sample without replacement. *Appl. Statist.* 1975; 24: 136.
- Bernard A; Van Efferen P. A generalization of the method of m ranking. *Proc. Kon. Ned. Akad. Wefensch.* 1953: A56.
- Berry KJ. AS179 Enumeration of all permutations of multi-sets with fixed repetition numbers. *Appl. Statist.* 1982; 31.
- Berry K; Mileke PW Jr. A measure of association for nominal independent variables. *Educ. Psych. Measurement.* 1992; 52: 895–898.
- Besag J; Clifford P. Sequential Monte Carlo p-values. *Biometrika.* 1991; 78: 301–304.
- Bissell AF. Ordered random selection without replacement. *Appl. Statist.* 1986; 35.
- Bitner JR; Ehrlich G; and Rheingold E. Efficient generation of the reflected Gray Code and its applications. *Commun. ACM.* 1976; 19: 517–521.
- Booth JG; Butler RW. Random distributions and saddlepoint approximations in general linear models. *Biometrika.* 1990; 77: 787–796.
- Boothroyd J. Algorithm 246. Gray code. *Commun. ACM.* 1964; 7: 701.

- Boothroyd J. Algorithm 29, permutation of the elements of a vector. *Computer J.* 1967; 60: 311.
- Boulton DM. Remarks on Algorithm 434. *Commun. ACM.* 1974; 17: 326.
- Boyett JM. Random  $r \times c$  tables with given row and column totals (algorithm AS 144). *Appl. Statist.* 1979; 28.
- Bratley P. Algorithm 306, Permutations with repetitions. *Commun. ACM.* 1967; 10: 450–451.
- Buckland ST; Garthwaite PH. Algorithm AS 259: Estimating Monte Carlo confidence intervals by the Robbins-Monro search process. *Appl. Statist.* 1990; 39: 413–424.
- Chase PJ. Algorithm 382. Combinations of M out of N objects. *Commun. ACM.* 1970A; 13: 368.
- Chase PJ. Algorithm 383, Permutations of a set with repetitions. *Commun. ACM.* 1970B; 13.
- Chen RS; Dunlap WP. SAS procedures for approximate randomization tests. *Beh. Res. Meth. Instr. Comp.* 1993; 25: 406–409.
- Dallal GE. Pitman: A Fortran program for exact randomization tests. *Comp. Biomed. Res.* 1988; 21: 9–15.
- Daniels HE. Discussion of paper by GEP Box and SL Anderson. *J. Roy. Statist. Soc. B.* 1955; 17: 27–28.
- Daniels HE. Discussion of paper by DR Cox. *J. Roy. Statist. Soc. B.* 1958; 20: 236–238.
- Davison AC; Hinkley DV. Saddlepoint approximations in randomization methods. *Biometrika.* 1988; 75: 417–431.
- DeCani J. An algorithm for bounding tail probabilities for two-variable exact tests. *Randomization.* 1979; 2: 23–24.
- Dershowitz N. A simplified loop-free algorithm for generating permutations. *BIT.* 1975; 15: 158–164.
- Durstenfield R. Random permutations. *Commun. ACM.* 1964; 7: 420.
- Edgington ES; Strain AR. Randomization tests: Computer time requirements. *J. Psychol.* 1973; 85: 89–95.
- Edgington ES; Strain AR. A computer program for randomization tests for predicted trends. *Beh. Res. Meth. Instr.* 1976; 8: 470.
- Edginton ES; Khuller PLV. A randomization test computer program for trends in repeated-measures data. *Educ. Psych. Measurement.* 1992; 52: 93–96.
- Ehrlich G. Algorithm 466. Four combinatorial algorithms. *Commun. ACM.* 1973; 16: 690–691.
- Feldman SE; Kluger E. Shortcut calculations to Fisher-Yates “exact tests”. *Psychometrika.* 1963; 2: 289–291.
- Fike CT. A permutation generation method. *Computer J.* 1975; 18: 21–22.
- Fleishman AI. A program for calculating the exact probabilities along with explorations of  $m$  by  $n$  contingency tables. *Educ. Psych. Measurement.* 1977; 33: 798–803.
- Gail M; Mantel N. Counting the number of  $r \times c$  contingency tables with fixed marginals. *JASA.* 1977; 72: 859–862.
- Gentleman JF. Algorithm AS88. Generation of all  $nCr$  combinations by simulating nested Fortran DO loops. *Appl. Statist.* 1975; 24: 374–376.
- Goetheluck P. Computing binomial coefficients. *Amer. Math. Monthly.* 1987; 94: 360–365.
- Graves GW; Whinston AB. An algorithm for the quadratic assignment probability. *Mgmt. Science.* 1970; 17: 453–471.
- Green BF. A practical interactive program for randomization tests of location. *Amer. Statist.* 1977; 31: 37–39.
- Gregory RJ. A Fortran computer program for the Fisher exact probability test. *Educ. Psych. Measurement.* 1973; 33: 697–700.
- Hancock TW. Remark on Algorithm 434. *Commun. ACM.* 1974; 18: 117–119.
- Hayes AF. PERMUSTAT: Randomization tests for the Macintosh. *Beh. Res. Meth. Instr. Comp.* 1996; 28: 473–475.

- Hayes AF. SPSS procedures for approximate randomization tests. *Beh. Res. Meth. Instr. Comp.* 1998; 30: 536–543.
- Hayes JE. Fortran program for Fisher's exact test. *Beh. Meth. Instr.* 1975; 7: 481.
- Hilton JF; Mehta CR; and Patel NR. A algorithm for conducting exact Smirnov tests. *Comput. Statist. Data Anal.* 1994; 17: 351–361.
- Hirji KF; Mehta CR; and Patel NR. Computing distributions for exact logistic regression. *JASA*. 1987; 82: 1110–1117.
- Howell DC; Gordon LR. Computing the exact probability of an r by c contingency table with fixed marginal totals. *Behav. Res. Meth. Instr. Comp.* 1976; 8: 317.
- Hull ID; Peto R. Alg AS35 probabilities derived from finite populations. *Appl. Statist.* 1971; 20: 99–105.
- Ives FM. Permutation enumeration: Four new permutation algorithms. *Commun. ACM.* 1976; 19: 68–70.
- Joe H. An ordering of dependence for contingency tables., *Linear Algebra Applic.* 1985; 70: 89–103.
- Joe H. Extreme probabilities for contingency tables under row and column independence with applications to Fisher's exact test. *Commun. Statist. A.* 1988; 17: 3677–3685.
- Knott GD. A numbering system for permutations of combinations. *Commun. ACM.* 1976; 19: 355–356.
- Knuth DE. *The Art of Computer Programming*, V 2. Reading, MA: Addison-Wesley; 1973.
- Kreiner S. Analysis for multidimensional contingency tables by exact conditional frequencies: Techniques and strategies. *Scand. J. Statist.* 1987; 14: 97–112.
- Kromrey JD; Chason WM; and Blair RC. PERMUTE: A SAS algorithm for permutation testing. *Appl. Psych. Measure.* 1992; 16: 64.
- Kurtzburg J. Algorithm 94. Combination. *Commun. ACM.* 1962; 5: 344.
- Lam CWH; Sothen LH. Three new combination algorithms with the minimal-change property. *Commun. ACM.* 1982; 25: 555–559.
- Liu CH; Tang DT. Algorithm 452. Enumerating combinations of m out of n objects. *Commun. ACM.* 1973; 16: 485.
- Mackenzie G; O'Flaherty M. Direct simulation of nested Fortran DO loops. *Appl. Statist.* 1982; 31: 71–74.
- March DL. Exact probabilities for r × c contingency tables. *Commun. ACM.* 1972; 15: 991–992.
- Marsh NWA. Efficient generation of all binary patterns by Gray Code. *Appl. Statist.* 1987; 36: 245–249.
- Mehta CR; Patel NR. A network algorithm for the exact treatment of the 2 × K contingency table. *Commun. Statist. B.* 1980; 9: 649–664.
- Mehta CR; Patel NR. A network algorithm for performing Fisher's exact test in r × c contingency tables. *JASA*. 1983; 78: 427–434.
- Mehta CR; Patel NR. A hybrid algorithm for Fisher's exact test in unordered r × c contingency tables. *Commun. Statist.* 1986; 15: 387–403.
- Mehta CR; Patel NR. FEXACT: A Fortran subroutine for Fisher's exact test on unordered r × c contingency tables. *ACM Trans. Math. Software.* 1986; 12: 154–161.
- Mehta CR; Patel NR; and Senchaudhuri P. Exact power and sample size computations for the Cochran Armitage trend test. *Biometrics.* 1998; 54: 1615–1621.
- Mielke PW Jr; Berry KJ. Fisher's exact probability test for cross-classification tables. *Educ. Psych. Measurement.* 1992; 52: 97–101.
- Minc H. Rearrangements. *Trans. Amer. Math. Soc.* 1971; 159: 497–504.
- Nelson DE; Zerbe GO. A SAS/IML program to execute randomization of response curves with multiple comparisons. *Amer. Statist.* 1988; 42: 231–232.
- Nigam AK; Gupta VK. A method of sampling with equal or unequal probabilities without replacement. *Appl. Statist.* 1984; 33.
- Nijenhuis A; Wilf HS. *Combinatorial Algorithms*. New York: Academic Press; 1978.

- Oden NE. Allocation of effort in Monte Carlo simulation for power of permutation tests. *JASA*. 1991; 86: 1074–1076.
- Ord-Smith RJ. Generation of permutation sequences: Part 1. *Computer J.* 1970; 13: 152–155.
- Ord-Smith RJ. Generation of permutation sequences: Part 2. *Computer J.* 1971; 14: 136–139.
- Pagano M; Halvorsen K. An algorithm for finding the exact significance levels of  $r \times c$  contingency tables. *JASA*. 1981; 76.
- Pagano M; Trichler D. On obtaining permutation distributions in polynomial time. *JASA*. 1983; 78: 435–441.
- Page ES. Note on generating random parameters. *Appl. Statist.* 1967; 16: 273–274.
- Page ES; Wilson LB. *An Introduction to Combinatorial Combinations*. Cambridge, U.K.: Cambridge University Press; 1979.
- Patefield WM. An efficient method of generating  $r \times c$  tables with given row and column totals (algorithm AS 159). *Appl. Statist.* 1981; 30: 91–97.
- Payne WH; Ives FM. Combination generators. *ACM Trans. Math. Software*. 1979; 5: 163–172.
- Portnoy S; Koenker R. The Gaussian hare and the Laplacian tortoise: Computability of squared-error versus absolute error estimators. *Statist. Sci.* 1997; 12: 279–300.
- Rabinowitz M; Berenson ML. A comparison of various methods of obtaining random order statistics for Monte Carlo computations. *Amer. Statist.* 1974; 28: 27–29.
- Radlow R; Alf EF. An alternate marginal assessment of the accuracy of the chi-square test of goodness of fit. *JASA*. 1975; 70: 811–813.
- Rao CR. Generation of random permutations of given number of elements using random sampling numbers. *Sankhya*. 1961; 33: 305–307.
- Richards LE; Byrd J. AS 304: Fisher's randomization test for two small independent samples. *Appl. Statist.* 1996; 45: 394–398.
- Robertson WH. Programming Fisher's exact method of comparing two percentages. *Technometrics*. 1960; 2: 103–107.
- Robinson J. Saddlepoint approximations to permutation tests and confidence intervals. *J. Roy. Statist. Soc. B*. 1982; 44: 91–101.
- Rogers MS. A Philco 2000 program to exhibit distinguishably different permutations. *Behav. Sci.* 1964; 9: 289.
- Rogstad SH; Pelikan S. Gelstats—A computer-program for population-genetics analyses using VNTR multilocus probe data. *Biotechniques*. 1996; 21: 1128–1131.
- Rohl JS. Generating permutations by choosing. *Computer J.* 1978; 21: 302–305.
- Romesburg HC; Marshall R and Mauk TP. FITEST-A computer program for “exact chi-square” goodness of fit tests. *Computers Geosciences*. 1981; 7: 457–458.
- Roy MK. Evaluation of permutation algorithms. *Computer J.* 1978; 21: 296–301.
- Sag TW. Algorithm 242. Permutation with a set of repetitions. *Commun. ACM*. 1964; 7: 585.
- Saunders IW. Enumeration of  $r \times c$  tables with repeated row totals. *Appl. Statist.* 1984; 33: 340–352.
- Shamos MI. Geometry and statistics: Problems at the interface. In *Algorithms and Complexity, New Directions and Recent Results*. JF Traub, Ed. New York: Academic Press; 1976: 251–279.
- Silva Mato A; Martin Andres A. Simplifying the calculation of the P-value for Barnard's test and its derivatives. *Statist. Comput.* 1997; 72: 137–143.
- Soms AP. An algorithm for the discrete Fisher's permutation tests. *JASA*. 1977; 72: 662–664.
- Spino C; Pagano M. Efficient calculation of the permutation distribution of robust two-sample statistics. *Comput. Statist. Data Anal.* 1991; 12: 349–368.
- Spino C; Pagano M. Efficient calculation of the permutation distribution of trimmed means. *JASA*. 1991; 86: 729–737.

- Streitberg B; Rohmed R. Exact distributions for permutation and rank tests: an introduction to some recently published algorithms. *Statist. Software Newsletter*. 1986; 12: 10–17.
- Streitberg B; Rohmel J. Exact distributions for rank- and permutation-tests in the general c-sample problem. *EDV Medizin Biologie*. 1987; 18: 12–19.
- Sunter AB. List sequential sampling with equal or unequal probabilities without replacement. *Appl. Statist.* 1977; 26: 261–268.
- Thomas D. Exact and asymptotic methods for the combination of  $2 \times 2$  tables. *Comput. Biomed. Res.* 1975; 8: 423–446.
- Tracey TJG. RANDALL: A Microsoft Fortran program for a randomization test of hypothesized order relations. *Educ. Psych. Measurement*. 1997; 57: 164–168.
- Tritchler D. An algorithm for exact logistic regression. *JASA*. 1984; 79: 709–711.
- Tritchler DL; Pedrini DT. A computer program for Fisher's exact test. *Educ. Psych. Measurement*. 35: 717–720.
- van den Brink WP; van den Brink SGJ. A modified approximate permutation test procedure. *Comp. Statist. Q.* 1990; 3: 241–247.
- Verbeek A; Kroonenberg PM. A survey of algorithms for exact distribution of test statistics in  $r \times c$  tables with fixed marginals. *Comput. Statist. Data Anal.* 1985; 3: 159–185.
- Vitter JS. Faster methods for random sampling. *Commun. ACM*. 1984; 27: 703–718.
- Vollset SE; Hirji KF. A microcomputer program for exact and asymptotic analysis of several  $2 \times 2$  tables. *Epidemiology*. 1991; 2: 217–220.
- Vollset SE; Hirji KF; and Elashoff RM. Fast computation of exact confidence limits for the common odds ratio in a series of  $2 \times 2$  tables. *JASA*. 1991; 86: 404–409.
- Walsh JE. An experimental method for obtaining random digits and permutations. *Sankhya*. 1957; 17: 355–360.
- Wells MB. *Elements of Combinatorial Computing*. Oxford, U.K.: Pergamon; 1971.
- Wichmann BA; Hill ID. Algorithm AS 183: An efficient and portable pseudo-random number generator. *Appl. Statist.* 1982; 31: 188–190.
- Woodhill AD. Generation of permutation sequences. *Computer J.* 1977; 20: 346–349.
- Wright T. A note on Pascal's triangle and simple random sampling. *College Math. J.* 1984; 20: 59–66.
- Zar JH. A fast efficient algorithm for the Fisher exact test. *Behav. Res. Meth. Instr. Comp.* 1987; 19: 413–414.
- Zimmerman H. Exact calculations of permutation distributions for r dependent samples. *Biometrical J.* 1985; 27: 349–352.
- Zimmerman H. Exact calculations for permutation distribution for r independent samples. *Biometrical J.* 1985; 27: 431–434.

## Bibliography Part 3:

# Seminal Articles

- Agresti A; Wackerly D; and Boyett JM. Exact conditional tests for cross-classifications: Approximations of attained significance levels. *Psychometrika*. 1979; 44: 75–83.
- Albers W; Bickel PJ; and Van Zwet WR. Asymptotic expansions for the power of distribution-free tests in the one-sample problem. *Ann. Statist.* 1976; 4: 108–156.
- Barton DE; David FN. Randomization basis for multivariate tests. *Bull. Int. Statist. Inst.* 1961; 39(2): 455–467.
- Basu D. Randomization analysis of experimental data: The Fisher randomization test. *JASA*. 1980; 75: 575–582.
- Bell CB; Doksum KA. Some new distribution free statistics. *Ann. Math. Statist.* 1965; 36: 203–214.
- Bickel PM; Van Zwet WR. Asymptotic expansion for the power of distribution free tests in the two-sample problem. *Ann. Statist.* 1978; 6: 987–1004 (corr. 1170–1171).
- Box GEP; Anderson SL. Permutation theory in the development of robust criteria and the study of departures from assumption. *J. Roy. Statist. Soc. B*. 1955; 17: 1–34 (with discussion).
- Boyett JM; Shuster JJ. Nonparametric one-sided tests in multivariate analysis with medical applications. *JASA*. 1977; 72: 665–668.
- Bradley JV. *Distribution Free Statistical Tests*. Englewood Cliffs, NJ: Prentice-Hall; 1968.
- Daniels HE. Relation between measures of correlation in the universe of sample permutations. *Biometrika*. 1944; 33: 129–135.
- Dwass M. Modified randomization tests for non-parametric hypotheses. *Ann. Math. Statist.* 1957; 28: 181–187.
- Edgington ES. Randomization tests for one-subject operant experiments. *J. Psych.* 1975; 90: 57–68.
- Fisher RA. *The Design of Experiments*. London, U.K.: Oliver and Boyd; 1935.
- Friedman M. The use of ranks to avoid the assumption of normality implicit in the analysis of variance. *JASA*. 1937; 32: 675–701.
- Gabriel KR; Hsu CF. Evaluation of the power of rerandomization tests, with application to weather modification experiments. *JASA*. 1983; 78: 766–775.
- Good PI. Globally almost powerful tests for censored data. *Nonpar. Statist.* 1992; 1: 253–262.
- Hilton JF; Mehta CR. Power and sample size for exact conditional tests with ordered categorical data. *Biometrics*. 1993; 49: 609–616.
- Hoeffding W. Combinatorial central limit theorem. *Ann. Math. Statist.* 1951; 22: 556–558.
- Kempthorne O. *Design and Analysis of Experiments*. New York: Wiley; 1952.
- Kempthorne O; Doerfler TE. The behavior of some significance tests under experimental randomization. *Biometrika*. 1969; 56: 231–248.

- Lehmann EL; Stein C. On the theory of some nonparametric hypotheses. *Ann. Math. Statist.* 1949; 20: 28–45.
- Mantel N. The detection of disease clustering and a generalized regression approach. *Cancer Res.* 1967; 27: 209–220.
- Mehta CR; Patel NR. A network algorithm for the exact treatment of the  $2 \times K$  contingency table. *Commun. Statist. B.* 1980; 9: 649–664.
- Mielke PW; Berry KJ; and Johnson ES. Multiresponse permutation procedures for a priori classifications. *Commun. Statist.* 1976; A5(14): 1409–1424.
- Noether GE. Distribution-free confidence intervals. *Amer. Statist.* 1972; 26: 39–41.
- Oden A; Wedel H. Arguments for Fisher's permutation test. *Ann. Statist.* 1975; 3: 518–520.
- Ogawa J. Effect of randomization on the analysis of a randomized block design. *Ann. Inst. Statist. Math. Tokyo.* 1961; 13: 105–117.
- Pearson ES. Some aspects of the problem of randomization. *Biometrika.* 1937; 29: 53–64.
- Pesarin F. On a nonparametric combination method for dependent permutation tests with applications. *Psychotherapy Psychosomatics.* 1990; 54: 172–179.
- Pitman EJG. Significance tests which may be applied to samples from any population. *Roy. Statist. Soc. Suppl.* 1937; 4: 119–130, 225–232.
- Pitman EJG. Significance tests which may be applied to samples from any population. Part III. The analysis of variance test. *Biometrika.* 1938; 29: 322–335.
- Plackett RL. Random permutations. *J. Roy. Statist. Soc. B.* 1968; 30: 517–534.
- Robinson J. A converse to a combinatorial central limit theorem. *Ann. Math. Statist.* 1972; 43: 2055–2057.
- Romano JP. On the behavior of randomization tests without a group invariance assumption. *JASA.* 1990; 85(411): 686–692.
- Rosenbaum PR. Conditional permutation tests and the propensity score in observational studies. *JASA.* 1984; 79: 565–574.
- Tukey JW; Brillinger DR; and Jones LV. *Management of Weather Resources: Vol II: The Role of Statistics in Weather Resources Management.* Washington, DC: Department of Commerce, US Government Printing Office; 1978.
- Wald A; Wolfowitz J. Statistical tests based on permutations of the observations. *Ann. Math. Statist.* 1944; 15: 358–372.
- Watson GS. Sufficient statistics, similar regions, and distribution-free tests. *J. Roy. Statist. Soc. B.* 1957; 19: 262–267.
- Welch WJ. Construction of permutation tests. *JASA.* 1990; 85: 693–698.

# Author Index

- Abushullaih BA, 2  
Adamson P, 3  
Addelman S, 196  
Adderley EE, 2  
Agresti A, 106, 110, 115, 189, 191  
Akl SG, 1, 187  
Albers W, 33, 196, 210  
Albert A, 2, 91  
Albright JW, 3  
Alderson MR, 134  
Alf EF, 41  
Alroy J, 3  
Amos CI, 2, 3  
Andersen PK, 149  
Anderson MJ, 128, 129  
Anderson SL, 38, 73  
Andreasen NC, 2  
Armitage P, 48  
Arndt S, 2  
Arnold HJ, 2, 11  
Arnold JT, 2  
Ascher S, 11, 196  
Astrachen E, 3, 134  
Avadhani TV, 3  
  
Baglivo J, 11, 115, 187  
Bailer AJ, 40, 187  
Bailar JC, 11, 134, 196  
Baker FB, 3, 55, 135  
Baker RD, 41, 42, 188  
Balasuvramanian R, 3  
  
Bamshad M, 2  
Baptista J, 100  
Baran E, 2  
Barbella P, 12  
Barnard GA, 116, 185  
Barnett TP, 157  
Barton DE, 10, 18  
Basawa IU, 3  
Basu D, 115  
Bayliss AC, 2  
Beck RW, 87  
Bell CB, 3, 9, 11, 118, 145  
Bell K, 3  
Belyea LR, 2  
Benes FM, 3  
Berger RL, 35  
Berkson J, 115  
Berlin JA, 2  
Berrington AM, 2  
Berry KJ, 2, 11, 19, 57, 106, 115, 136, 146, 155, 164, 187, 196  
Bersier LF, 2  
Besag JE, 186, 199  
Bickel PJ, 11, 27, 33, 36, 196, 210, 212  
Bickis M, 115, 189  
Bi-Fong Lee M, 3  
Birch MW, 114  
Bishop YMM, 104  
Bissell AF, 187  
Bitner JR, 11, 187  
Bjornstad ON, 2  
Blair C, 87

- Blair RC, 12, 83, 87, 88  
Bland BH, 3  
Boess FG, 3  
Boik RJ, 38  
Bookstein FL, 3  
Boos DD, 155  
Booth JG, 196  
Boothroyd J, 11, 188  
Boschloo RD, 98  
Boudewyn P, 3  
Boulton DM, 187  
Box GEP, 38, 73  
Box JF, 116  
Boyett JM, 82, 189  
Bradbury IS, 12, 55, 196  
Bradley JV, 144  
Bradley RA, 10  
Bramme R, 2  
Brammer M, 2  
Brammer MJ, 3  
Brennan J, 115, 189  
Breslow NE, 168  
Brier GW, 136  
Brillinger DR, 2  
Brockwell PJ, 11  
Bross IDJ, 3, 113, 173  
Brown BM, 38, 130  
Browne C, 155  
Bryant EH, 2, 134  
Bullmore E, 2, 12  
Bunney BG, 3  
Bunney WE, 3  
Burdick DS, 2  
Burgess AP, 3  
Büringer H, 73  
Burnett WS, 2  
Burns CA, 3  
Busby DG, 3  
Butler RW, 196  
  
Cade B, 2, 62, 128, 130, 133, 146  
Call SB, 3  
Campbell IC, 3  
Casagrande JJ, 167  
Casella G, 194  
Chamberlain RM, 3  
Chapelle JP, 2, 91  
Chase PJ, 11  
  
Chatterjee SK, 78  
Chatterton R, 3  
Chernoff H, 145  
Christine BW, 2  
Christyakov VP, 196  
Chu KC, 17  
Chung JH, 82  
Cinciripini P, 3  
Cizadlo T, 2  
Clark LC, 2  
Clark RM, 3  
Clarke B, 106  
Cléoux R, 11  
Cliff AD, 11, 135  
Clifford P, 186, 199  
Clowes AW, 3  
Cohen A, 2  
Cohen J, 3  
Cole DA, 63  
Conover WJ, 38  
Constanzo CM, 11  
Copp GH, 2  
Cornfield J, 11, 99, 100  
Cory-Slechta DA, 4  
Costanzo CM, 135  
Cox C, 4, 62  
Cox DF, 11  
Cox DR, 143, 182  
Crump KS, 4  
Cullen JB, 2  
Currie J, 3  
  
Daniels HE, 196  
David A, 2  
David FN, 11, 78  
David HA, 145  
David HT, 72  
Davis BM, 98  
Daw NC, 2  
Day NE, 168  
DeCani J, 11, 49, 189  
Dee, CR, 3  
DeMets DL, 73  
Dempster ER, 3  
Denker M, 11  
Diaconis P, 11  
Dietz EJ, 135  
Diggle PJ, 3

- Dodge Y, 146  
Doerge RW, 3  
Doksum KA, 9, 11, 145  
Donegani M, 11, 151  
Donnelly SN, 3  
Donoghue JF, 11, 118  
Doolittle RF, 3  
Douglas ME, 134  
Draper D, 24, 66, 183  
Dupuis J, 3  
Dwass M, 185, 187, 196
- Eden T, 2  
Edgington ES, 2, 3, 10, 11, 65, 187  
Efron B, 11, 42, 43, 71, 151  
Ehrlich G, 187  
Eiler JH, 2  
Elashoff RM, 195  
Ellman A, 3  
Endler JA, 134  
Engle S, 3  
Entsuah AR, 72, 141  
Erdos P, 11, 196  
Evett IW, 2
- Fallon JH, 3  
Fang KT, 11  
Faris PD, 3  
Farrar DA, 4  
Faulkner J, 3  
Fears TR, 17  
Feinstein AR, 3, 12, 71  
Feldman SE, 11, 189, 191  
Ferron J, 2, 131  
Fienberg SE, 104  
Finney DJ, 189  
Fisher L, 135  
Fisher RA, 2, 8, 13, 38, 71, 99, 116  
Fix E, 167  
Flury BD, 3  
Forster JJ, 2, 115  
Forsythe AB, 66  
Foutz RN, 91  
Frank D, 46  
Fraser DAS, 82  
Fraumeni JF, 134  
Freedman D, 129
- Freedman L, 2  
Freeman GH, 11, 173  
Friedman JH, 1, 83, 137
- Gabriel KR, 2, 3, 8, 35, 125, 135, 198  
Gail MH, 2, 11, 65, 115, 189  
Galambos J, 24  
Gann P, 3  
Garside GR, 115  
Gart JJ, 3  
Garthwaite PH, 197  
Gastwirth JL, 2, 183  
Gentleman JF, 11  
George EI, 194  
Gerig TM, 86, 87  
Ghandour G, 3  
Ghosh MN, 122  
Gill DS, 78  
Gill PD, 2  
Glass AG, 2, 134  
Gliddentracey C, 4  
Gold S, 2  
Goldberg P, 91  
Good PI, 2, 3, 4, 11, 12, 46, 72, 98, 144,  
147, 149, 175  
Goodman LA, 115  
Gottschalk LA, 3  
Grady CL, 3  
Graubard BI, 109  
Gray R, 114  
Green BF, 190  
Greenland S, 116  
Greenwood AK, 4  
Gruzelier JH, 3  
Guerra R, 3  
Gutierrez LG, 147  
Gunz FW, 2  
Gupta VK, 102, 187
- Haber M, 98, 116  
Haberman SJ, 110  
Hajek J, 11, 196  
Hajimohamadenza I, 3  
Hall P, 38, 57, 129  
Hall WJ, 8  
Halter JH, 103  
Halton JH, 11, 173

- Halvorsen K, 11, 189  
 Hand RE Jr, 11  
 Harpending HC, 2  
 Hartigan JA, 11  
 Hartmann C, 3  
 Hartz A, 3  
 Hasegawa M, 156  
 Hauser ER, 2  
 Haxby JV, 3  
 Hayes AF, 40  
 Heckel D, 2  
 Hedrick ML, 3  
 Heise HW, 2  
 Henze N, 84  
 Hettmansperger TP, 183  
 Heusghem C, 2  
 Hewlett PS, 3  
 Hiatt WR, 91  
 Hickey WF, 3  
 Higgens JJ, 2  
 Highton R, 134  
 Hilton JF, 198  
 Hinshelwood MM, 2  
 Hirji KF, 195  
 Hochberg Y, 87  
 Hodges JL Jr, 11, 167  
 Hoeffding W, 11, 73, 195, 210, 211  
 Hoel DG, 3  
 Hogg RV, 151  
 Holland PW, 104  
 Holm S, 87  
 Holmes MC, 117  
 Hong WK, 3  
 Houle S, 3  
 Howard M, (see Good P,) 2, 12  
 Howard R, 2  
 Howe HL, 2  
 Hsu CF, 2, 35, 198  
 Huber PJ, 146  
 Hubert LJ, 1, 3, 11, 134–7  
 Hubley AM, 198  
 Hudmon KS, 3  
 Huitema BE, 131  
 Hutson AD, 103  
 Irony TZ, 98  
 Ives FM, 11  
 Iwano EJ, 2  
 Iyer HK, 137–8  
 Izenman AJ, 11, 41  
 Jackson CWA, 2  
 Jackson DA, 2  
 James GS, 57  
 Janot N, 2  
 Janssen A, 149  
 Jayawardene D, 2  
 Jennings JM, 3  
 Jiang H, 2, 3  
 Jockel KH, 187  
 Joe H, 194  
 Jogdeo K, 196  
 John RD, 35, 196  
 Johnson AG, 3  
 Johnson ME, 38  
 Johnson MM, 38  
 Jones JS, 135  
 Jones LV, 2, 135  
 Jones MP, 115  
 Jones TK, 2  
 Jorde LB, 2, 12  
 Kalbfleisch JD, 3, 149  
 Kapur S, 3  
 Karinski W, 87  
 Karlin S, 3, 11, 156, 158  
 Kaufman DW, 135  
 Kayama N, 3  
 Kazdin AE, 2, 11  
 Keator D, 3  
 Kelly ME, 3  
 Kemp P, 2  
 Kemphorne O, 2, 10, 11, 12, 13, 34, 65,  
     66, 73, 176, 196  
 Kenagy RD, 3  
 Kennedy PE, 2, 62, 128, 133  
 Kere J, 2  
 Khan KA, 137  
 Kidd KK, 3  
 Kim C, 116  
 Kim J, 3  
 Kim MJ, 2, 125, 132  
 Klauber MR, 134

- Klein E, 3  
Klingenberg CP, 3  
Kluger E, 11, 189, 191  
Koch GG, 24, 204  
Kolassa JE, 194  
Kolchin VF, 196  
Konigsberg LW, 2  
Korn EL, 109  
Korschning S, 3  
Koziol JA, 90  
Krakowiak P, 2  
Kramer A, 3  
Krewski D, 189  
Krishino H, 156  
Kruskal WH, 72, 105  
Kryscio RJ, 134  
Kulberts HE, 2  
Kvamme KL, 2
- Lachin JM, 2  
Lam CWH, 184  
Lambert D, 143, 146  
Lan K, 73  
Lane D, 1, 129  
Lange N, 3  
Laroche J, 2  
Latscha R, 189  
Lea HJ, 3  
Lefebvre M, 12  
Legendre P, 128, 129  
Lehmann EL, 4, 11, 12, 24, 27, 33–5, 44, 50, 72, 98, 114, 167, 174, 196, 202  
Lemke JH, 115  
Lenth RV, 151  
Leslie PH, 189, 191  
Levin DA, 3  
Levin JR, 1  
Li FP, 134  
Liang KY, 115  
Liebtrau AM, 104  
Lin DY, 12  
Liu CH, 11  
Livingstone KD, 3  
Lock RH, 186  
Lorenz J, 2  
Loughlin TM, 62  
Louis EJ, 3
- Mack C, 115  
Mackay DA, 2  
Madow WG, 195  
Magnussen S, 3  
Majeed AW, 3  
Makinodan T, 3, 68  
Makridakis S, 132  
Manly BFJ, 2, 62, 106, 122, 128, 130, 159, 183  
Mantel N, 2, 11, 100, 112, 134, 189, 191  
Marascuilo LA, 3  
March DL, 189  
Marcus LF, 3  
Maritz JS, 35, 40, 122, 130, 146  
Marriott FHC, 186  
Marsh NWA, 188  
Martin H, 73  
Maruff P, 3  
Massart DL, 2  
Maxwell SE, 63  
McCarthy PJ, 11  
McDonald AD, 112, 134  
McDonald JW, 2, 115  
McDonald LL, 98  
McGee VE, 132  
Mcintosh AR, 3  
McKean JW, 131  
McKinney PW, 3, 96, 179  
McLeod RS, 2  
McQueen G, 2  
McSweeney M, 3  
Mead R, 138  
Meagher TR, 2  
Mehta CR, 48, 71, 101, 114, 115, 187, 189, 191, 198  
Mekers J, 2  
Melvold RW, 3  
Merrington M, 134  
Meyers MH, 2  
Mielke PW Jr, 2, 11, 19, 57, 106, 115, 136, 146, 155, 164, 196  
Miliiken GA, 98  
Miller AJ, 2  
Miller RA, 3  
Minc H, 11  
Mitchell-Olds T, 2  
Mitra SK, 196  
Morar B, 3

- Morrison DF, 82  
 Motoo M, 196  
 Mullins L, 3  
 Murphy BP, 91  
 Mustacchi A, 134  
 Nadler W, 3  
 Nayak RA, 134  
 Nelson CR, 2  
 Nelson DE, 91  
 Nelson LS, 3  
 Nelson W, 3  
 Ness RB, 2  
 Neuenschwander BE, 3  
 Nguyen TT, 115  
 Nicholl JP, 3  
 Nicholson SM, 3  
 Nigam AK, 188  
 Noble EP, 3  
 Noble W, 2, 62  
 Noether GE, 11, 122  
 Noreen E, 187  
**O'Gorman TW**, 115  
 Oden NE, 11, 198  
 Ogawa J, 196  
 Oja H, 126  
 Okamoto M, 188  
 Oleary DS, 2  
 Oliver D, 115  
 Onghena P, 2, 19, 20, 131  
 Ord JK, 11  
 O'Reilly FJ, 196  
 O'Sullivan F, 2  
**Pagano M**, 11, 115, 189, 191  
 Parraga MI, 3  
 Parzen MI, 130  
 Patefield WM, 11, 111, 189, 191  
 Patel NR, 48, 71, 101, 114, 115,  
     187, 189, 191, 198  
 Paternoster R, 2  
 Paynter RA, 3  
 Peacock J, 3  
 Pelikan S, 3  
 Pena E, 71  
 Penninckx W, 3  
 Pereira CAB, 99  
 Pesarin F, 61, 81, 83, 89, 92  
 Piantadosi S, 265  
 Pike MC, 100, 134, 164  
 Pitman EJG, 8, 46, 196  
 Plackett RL, 3, 188  
 Praska Rao BLS, 3  
 Prentice RL, 3, 149, 167  
 Priesendorfer RW, 157  
 Prusiner SI, 2  
 Puri ML, 11, 86, 91  
**Quade D**, 52  
 Quinn JF, 2, 3  
**Rabehesk S**, 2  
 Rafsky LC, 83, 137  
 Rankin JA, 3  
 Raz J, 49  
 Reed DW, 2  
 Reed MWR, 3  
 Renyi A, 11, 196  
 Rhein DM, 3  
 Rheingold E, 188  
 Richards L, 130, 146  
 Ripley BD, 138  
 Rips E, 3  
 Ritchie TL, 3  
 Ritland C, 2  
 Ritland K, 2  
 Roberson P, 135  
 Robinson J, 11, 35, 44, 127, 196, 212  
 Robson AJ, 2  
 Rogers AR, 2  
 Rogers MS, 11  
 Rogstad SH, 3  
 Rohmed R, 189  
 Romano J, 9, 12, 33, 37, 119, 147, 204,  
     213  
 Roper RJ, 3  
 Rosen B, 11  
 Rosenbaum PR, 66, 142  
 Rosenberg Y, 3  
 Rounds J, 4  
 Royaltey HH, 3, 134  
 Rubin H, 183

- Ryan JM, 4  
Ryman N, 134
- Sag TW, 11  
Sainsbury RS, 3  
Saitoh T, 2  
Salsburg DS, 2  
Sampford MR, 149  
Santner TJ, 117  
Scheffe H, 55, 65  
Schemper M, 149  
Schrieven K-H, 73  
Schultz JR, 134, 135  
Scott Elizabeth,  
Scott Elton, 146  
Scranage JK, 2  
Selander RK, 135  
Self SG, 115  
Sen PK, 11, 78, 79, 91  
Senchaudhuri P, 48, 101, 115, 187,  
    191, 198  
Servy EC, 79  
Shane HD, 11, 86  
Shapiro CP, 11, 196  
Shaw DE, 2  
Shell EJ, 143, 182  
Shen CD, 52  
Shi HH, 2, 3  
Shuster JJ, 2, 82, 116  
Sidak Z, 145  
Siegel S, 11  
Siegmund H, 73  
Siemiatycki J, 112, 134  
Silvey SD, 196  
Siotani M, 78  
Slam P, 2  
Slud E, 73  
Smeets JP, 2  
Smeyersverbeke J, 3  
Smeyers-Verbeke J, 2  
Smith EJ, 2  
Smith PG, 134, 164  
Smith PWF, 2, 115  
Smith RL, 71  
Smythe A, 3  
Smythe RT, 71  
Snell MK, 117, 135  
Sokal RR, 3, 134, 135
- Solomon H, 2  
Soms AP, 110  
Sotchen LH, 184  
Spears GFS, 2  
Spicer CC, 134  
Spino C, 11  
Spitz MR, 2, 3  
Sprent P, 62  
Srivastava DK, 2  
Startz R, 2  
Stein C, 11  
Stenberg PE, 2  
Stenseth NC, 2  
Still AW, 60  
Stilson DW, 3  
Stine RA, 122  
Stoddard CJ, 3  
Stoneman DM, 66  
Storer BE, 116  
Streitberg B, 189  
Stuart GW, 3  
Sugihara G, 2  
Suissa S, 116  
Sung S, 2  
Syrjala SE, 2, 139
- Tajuddin IH, 137  
Tan WY, 2, 65  
Tang C, 3  
Tang DT, 11  
Tanner MS, 194  
Tardif S, 196  
Tarone RE, 17  
Taylor DW, 2  
Taylor J, 149  
ter Braak CJF, 12, 62, 129  
Teuscher C, 3  
Thomas DG, 100  
Throckmorton T, 196  
Tibshriani R, 43  
Tilbury JB, 188  
Titterington M, 129  
Tracey TJG, 3  
Tracy DS, 137  
Tritchler D, 11, 35, 195  
Troendle JF, 12, 87  
Trzos RJ, 46  
Tukey JW, 2, 11

- Tulving E, 3  
Tung KSK, 3  
Turnbull BW, 2  
  
Vadiveloo J, 187  
Valdesperez RE, 2  
Van Damme G, 2  
Van Zwet WR, 27, 33, 36, 196, 210  
Vandenbosch C, 2  
Vandermeulen J, 3  
vanKeerberghen P, 2  
Vanlier JB, 2  
van-Putten B, 84  
Vecchia DF, 137  
Veitch LG, 2  
Vitter JS, 11  
Vogelsong K, 3  
Vollset SE, 195  
  
Wackerly D, 112, 189  
Walburg HE, 3  
Wald A, 11, 49, 72, 77, 122, 210  
Walsh JE, 188  
Wan Y, 2, 3  
Wang FT, 146  
Ware JH, 2, 131  
Watkins WS, 2  
Wei LI, 71  
Wei LJ, 2, 73, 130, 191  
Weir BS, 2  
Weiss B, 4  
Welch BL, 27, 196  
Welch WJ, 9, 12, 62, 127, 147  
Werner M, 79, 120  
Westfall DH, 12, 62, 87  
Weth F, 3  
Whaley FS, 136  
Wheelwright SC, 132  
  
Whinston AB, 135  
White AP, 60  
White R, 196  
Whitney P, 2  
Widmark C, 3  
Wilk MB, 65  
Williams JS, 11  
Williams PT, 3, 158  
Williams REO, 117  
Williams SC, 2  
Williams-Blangero S, 2  
Wilson HG, 146  
Witztum D, 3  
Wolfowitz J, 11, 49, 77, 122, 210  
Wolpert RW, 35  
Woodrooffe M, 3  
Woolson RF, 115  
Wright T, 184  
Wu JC, 3  
Wu XF, 2  
  
Yanagimoto T, 188  
Yang F, 3  
Yano T, 156  
Yates F, 2  
Ying Z, 130  
Young MJ, 3  
Young SS, 12, 62, 87  
Yucesan E, 2  
  
Zelen M, 101  
Zempo N, 3  
Zerbe GO, 91  
Zimmerman DL, 139  
Zinn J, 38  
Zumbo BD, 198  
Zyskind G, 196

# Subject Index

Acceptance region, 19, 148, 210  
Acyclic network, 192  
Adaptive test, 151  
Additivity, 54, 69  
Agriculture, 1, 50, 55, 63  
Algorithms,  
    branch and bound, 190  
    characteristic functions, 195  
    Dynamic programming, 160  
    for clinical trials, 71  
    for contingency tables, 191  
    Edgeworth expansions, 196  
    enumeration and selection, 101, 187  
    Gibbs sampling, 194  
    network, 191, 198  
    recursive, 189  
    saddlepoint, 196  
    shift, 190  
Alternative  
    global, 104  
    monotone increasing, 56  
    normal, 51, 145  
    one-sided vs two-sided, 34  
    ordered, 64  
    shift, 35  
    stochastically increasing, 35, 204  
    versus hypothesis, 6, 8, 31, 118, 201  
Analysis of variance (ANOV), 55  
Anthropology, 1  
Aquatic science, 1  
Archaeology, 1, 136  
Assumptions, 13, 29, 37, 67, 107, 183

Asymptotic  
    approximation, 11, 195, 200  
    consistency, 33  
    distribution, 2, 10  
    efficiency, 211  
Atmospheric science, 1, 157  
  
Behrens–Fisher problem, 38  
Bias, 38, 62  
Biased coin approach, 71  
Binomial  
    comparing two Poissons, 174  
    distribution, 98, 198  
    trials, 98  
Bioequivalence, 91  
Biology, 1  
Biotechnology, 1  
Birth defects, 112, 135  
Blocking, 49, 60, 63, 74, 86, 114, 163  
Bootstrap  
    comparison with permutation test, 10,  
        25, 35, 43, 131, 212  
    DNA sequencing, 156  
    estimate, 69, 123  
    nonparametric, 38, 119  
    paired comparisons, 151  
    resampling procedure, 9  
    tests for dependence, 122  
    test for location, 172  
    test for variances, 42  
    test procedure, 26, 38

- Botany, 1  
 Branch and bound (*see* Algorithms)
- Case controls (*see* Observational study)  
 Categorical data, 89, 94  
 Cell culture, 4, 12, 53  
 Censored data, 4, 47, 167  
 Censoring, 146  
 Chemistry, 1  
 Chi-square statistic  
     definition, 104  
     drawbacks, 102  
     restricted, 166  
 Climatology, 1  
 Clinical trial, 1, 70, 72, 87  
 Cluster analysis, 1, 11, 134  
 Cochran–Armitage, 48  
 Computers, 10  
 Computer  
     programs, (*see* Software)  
     science, 1  
 Confidence intervals, 34, 36, 42, 74, 111,  
     127, 196, 202, 209  
 Confound, 55, 67, 72  
 Conservative test, 23, 24, 27, 41  
 Contingency tables  
      $2 \times 2$ , 94, 191  
      $r \times c$ , ordered, 109  
      $r \times c$ , unordered, 103, 191, 200  
     odds ratio test, 99  
     with covariates, 100, 113  
 Control,  
     factor, 52, 54  
     group, 4, 154, 179  
 Correlation  
     bivariate, 208  
     Kendall, 136  
     linear, 45  
     Mantel's U, 134,  
     Pearson, 120  
     Pitman, 45, 120, 124  
     serial, 121  
     Spearman, 136  
 Covariance, 76, 115  
 Covariate  
     after the fact, 141  
     analysis, 24, 65, 74  
 Critical value, 10, 97, 207
- Data  
     categorical, 84, 94, 119, 172  
     discrete, 109, 174  
     ordered, 109  
     continuous, 130, 174  
 Decisions, 18, 29, 164, 201  
 Degrees of freedom, 44  
 Density (*see* Probability density)  
 Dependence,  
     first-order, 124  
     models, 118  
     quadrant, 119  
     serial correlation, 121  
     trend, 120  
 Deterministic, 14  
 Diagnostic imaging, 1  
 Directional data, 11  
 Distribution  
     Beta, 171  
     binomial, 75, 98, 172, 194  
     Cauchy, 31  
     chi-square, 171  
     double exponential, 31  
     exponential, 15, 35, 148, 151  
     F-ratio, 171  
     function defined, 203  
     gamma, 49, 148, 171  
     hypergeometric, 95, 99, 194  
     lognormal, 171  
     normal, 15, 44, 148, 151, 162, 170, 207  
     Poisson, 99, 162, 170, 172  
     Stochastically increasing, 204  
     Student's t, 171  
     symmetric, 31, 61  
     uniform, 31, 72  
 Distribution-free test, 1, 9  
 DNA sequencing, 11, 156  
 Dose response, 46  
 Double blind study, 70, 180  
 Drop outs, (*see* Withdrawals)
- Ecology, 1, 106, 155  
 Economics, 1, 124, 125, 133, 171  
 Education, 1, 106, 133, 159  
 Effect size, 19, 181  
 Efficiency, 145  
 Endocrinology, 1  
 Entomology, 1

- Epidemiology, 1, 112, 117, 134  
Ergonomics, 1  
Errors, Type I and II, 17, 19, 70  
Euclidian distance, 136, 156  
Exact test, 8, 23, 36, 102, 203  
    asymptotically exact, 24, 30, 51, 64, 89, 179, 214  
Exchangeable, 24, 30, 51, 64, 89, 179, 204  
Experimental design  
    balanced, 55  
    balanced incomplete block, 196  
    double blind, 70, 180  
    k-sample, 42  
    Latin square, 63  
    randomized block, 54, 196  
    sequential, 72  
    unbalanced, 58, 67  
Experimental unit, 55, 75, 180
- F**-Statistic, 9, 43, 57, 73  
False positive, 17  
Finite populations, 4, 36  
Fisher's exact test, 94, 116, 179  
Fisher's omnibus statistic, 89, 92  
Forensics, 1, 100  
Fourier analysis, 1  
Fundamental lemma, 165, 206
- G**AMP test, 147, 149  
Genetics, 1, 158  
Geography, 1  
Geology, 1  
Gerontology, 1  
Gibbs sampling, 194  
Goodness-of-fit, 166  
Gradient (*see* Bias)
- H**otelling's T, 76  
Hypothesis  
    null, 25, 64, 82, 104  
    simple vs compound, 23, 206  
    vs alternative, 6, 8, 31, 118
- I**mmunology, 1, 68, 116  
Importance sampling, 187
- Imputed values, 146  
Independent observations  
    exchangeable, 24  
    test for, 104, 119  
Indifference region, 19, 148  
Interaction, 24, 59, 62, 75  
Invariance  
    of a problem, 44, 163, 169, 214  
    of a test, 44, 163, 169  
    under permutations, 53, 64, 169
- k**-sample comparisons, 42, 48  
Kernel estimation, 49
- L**abelings, 6  
LAD regression, 130  
Latin Square, 63  
Least-squares, 66, 125  
Likelihood  
    conditional, 168  
    ratio, 164, 213  
Linear  
    correlation (*see* Correlation)  
    estimation, 49  
    form, 211  
    regression (*see* Regression)  
    statistic, 165  
    transformation, 86  
Logistic regression, 143, 167  
Losses, 18, 146, 164, 169, 201  
 $L_1$  norm, 56, 146, 172
- M**ain effect, 55  
Management science, 132  
Marginals, 94, 107, 189, 194  
Mantel's U  
    applications, 134  
    equivalences, 137, 139  
Matched pairs, 51, 86, 149, 168, 188  
Maximum likelihood, 164  
Median, 40, 146, 172  
Medical applications, 1, 19, 50, 51  
Metallurgy, 53, 131  
Mid-p, 98  
Minimal spanning tree, 85  
Missing data, 4, 58, 71, 91, 140

- Model**  
 randomization, 27  
 sampling, 27
- Molecular biology**, 1
- Monotone function**, 46, 92, 109, 126, 191
- Monte Carlo**, 33, 38, 48, 111, 185, 199
- Most powerful test**, 22
- Multiple comparisons**  
 omnibus statistic, 89, 92  
 step-up, step-down, 87
- Multivariate tests**, 1  
 one-sample, 76  
 two-sample, 77  
 repeated measures, 90
- Nearest-neighbor**, 138
- Neurobiology**, 1
- Neurology**, 1
- Neuropsychopharmacology**, 1
- Neuropsychology**, 1
- Network representation**, 192
- Neyman–Pearson lemma**, 165, 206
- Node**, 192
- Nonlinear device**, 152
- Normal distribution**, 13, 44, 148, 151, 162, 170, 207
- Normal scores**, 74, 110
- Nonparametric test**, 9
- Nonresponders**, 144, 154
- Null hypothesis**, 25, 104
- Observational study**, 142
- Odds ratio**, 99, 111
- Oncology**, 1
- One-sample problem**, 31, 151
- One- vs. two-sided test**, 96, 184
- Order statistics**, 168, 206
- Ordered alternatives**, 45
- Ornithology**, 1
- Outliers**, 4, 14, 129, 143
- Paleontology**, 1
- Parameter**  
 location, 31, 36, 170  
 nuisance, 24, 40, 163  
 scale, 38, 170  
 space, 164
- Parametric test**, 10  
 vs. permutation test, 1, 11, 13, 25, 35, 44, 55, 172
- Partial regression coefficients**, 128
- Percentile**, 66
- Permutation distribution**, 25, 158, 195
- Permutation test**  
 definition, 203  
 properties, 171  
 vs. bootstrap, 10, 25, 35, 43, 131, 212
- p-value** (*see* Critical value)
- Pharmacology**, 3, 16, 29, 79, 154, 171
- Physics**, 3, 14
- Physiology**, 3, 120
- Pitman correlation**, 45, 124, 139
- Pivotal quantity**, 35, 38, 128, 133
- Poisson** (*see* Distributions)
- Population** (*see* Sampling model)
- Power**  
 as function of  
 alternative, 23  
 distribution, 49, 51  
 variance, 50, 92  
 defined, 203  
 effect of  
 blocking, 50  
 covariates, 92  
 design, 131  
 sample size, 23, 181, 197  
 for large samples, 78  
 for variance comparisons, 43  
 maximizing, 204  
 relation to Type II error, 19  
 vs. cost of sampling, 155
- Power curve**, 23
- Precision**, 30, 91
- Probability**  
 conditional, 143  
 distribution (*see* distribution)  
 density, 204  
 of false coverage (*see* Type I and Type II error)
- Proportion**, 171
- Psychology**, 3, 161
- Quadrant density**, 138
- Quadrant dependence**, 119
- Quadratic form**, 134

- Radioactivity**, 29, 152  
**Random**  
    assignments (*see* randomization)  
    integer, 62  
    number, 62, 200  
    rearrangement, 78  
    variable, 24, 27, 204  
    vector, 77  
**Randomization**  
    in experiments, 28  
    model, 27  
    on the boundary, 34, 98  
    rerandomization, 8, 9  
    restricted, 66, 73  
    test (*see* Permutation test)  
**Rank test**, 11, 25, 196  
**Ranks**, 37, 74, 78, 86, 89, 110, 121, 144, 150  
**Ratio**, 171  
**Rearrange**, 7, 77  
**Regression**  
    coefficients, 66  
    compare slopes, 74, 90  
    forward stepping rule, 66  
    LAD vs LSD, 130  
    linear model, 65, 122, 212  
    logistic, 143, 167  
    multivariate, 128  
    polynomial, 49  
**Rejection region**, 19, 24, 36, 148  
**Reliability**, 1  
**Repeated measures**, 90  
**Rerandomize** (*see* Permutation)  
**Resampling methods**  
    with replacement, 26  
    without replacement, 6  
**Residuals**, 55, 60, 65, 129  
**Response profile**, 90  
**Risk**, 22  
**Robust**  
    test, 13  
    transformation, 146  
**Runs test**, 83  
  
**Sample**  
    control sample 4, 154, 179  
    distribution, 15  
    size, 181, 197  
    unrepresentative, 16  
**Sampling model**, 27  
**Scores**, 109, 110, 145  
**Screen**, 15  
**Sensitivity of a test**, 63, 98  
**Sequential**  
    procedure, 12  
    test, 72  
**Serial correlation**, 121, 124, 131  
**Shape**, 15  
**Shift**, 40, 84  
**Significance**  
    choice of level, 172  
    level, 13, 19, 33, 43, 58, 79, 89  
    statistical vs. practical, 84  
**Simultaneous inference**, 82  
**Single-case analysis**, 1, 11, 131  
**Sociology**, 1  
**Software**, 45, 59, 80, 181  
**Spatial distribution**, 136–9, 156, 182  
**Statistic**,  
    first- vs second-order, 164  
**Step-down, Step-up**, 87  
**Stepwise solution**, 66  
**Stochastic**, 1  
**Stochastically increasing**, 204  
**Stopping rule**, 186  
**Stratify**, 50  
**Student's t**, 9, 33, 43, 145  
**Studentize**, 81  
**Study time**, 155  
**Sufficiency**, 32, 162, 213  
**Surgery**, 4  
**Survival data**, 95, 114  
**Symmetry**, 31  
**Synchronized permutations**, 61, 75  
**Systematics**, 135  
**Swap**, 186  
  
**Tail probabilities**, (*see* Critical value)  
**Taxonomy**, 4  
**t-test** (*see* Student's t)  
**Test**,  
    adaptive, 151  
    Boot-Perm, 70  
    chi-square, 102  
    choosing, 170–178  
    Cramer's V, 104

- Test (*cont.*)  
 definition, 202  
 F-test, 27, 38, 45  
 Freeman and Halton, 103  
 Fisher's exact, 94, 116, 179  
 Fisher's omnibus, 89, 92  
 GAMP, 147, 149  
 Goodman-Kruskal tau, 105  
 Hotelling's t, 76  
 Kappa, 106  
 likelihood ratio, 103  
 linear by linear, 112  
 location, for, 31, 172  
 Mann-Whitney, 25, 145, 155  
 Mantel's U, 134  
 MRPP, 136  
 nonparametric, 9  
 odds ratio, 99, 111  
 one-sample, 31, 76, 151  
 one- vs. two-tailed, 96  
 parametric vs. non-parametric, 1, 11, 13,  
   25, 35, 44, 55, 172  
 Pitman correlation, 45, 124  
 restricted chi-square, 166  
 runs, 83  
 scale, for, 38, 208  
 shift, for, 36, 207  
 spatial distribution, 138-9  
 statistic, 6  
 Student's t, 27, 35, 155  
 symmetric volume, 138  
 tau, 105  
 UMPU, 33, 46, 99  
 Wilcoxon (*see* Mann-Whitney)  
 Zelen, 115
- Ties, 47, 109  
 Transformation  
   arcsin, 174
- linear, 86  
 logarithmic, 53, 144, 175  
 rank, 37, 145  
 rescale, 44  
 robust, 146  
 square root, 174  
 Toxicology, 2, 46, 180  
 Tree, 83  
 Trend, 120, 123  
 Two-sample problem, 36, 77, 146, 168  
 Type I and II errors, 19, 70
- UMP test, 22  
 UMPU test, 33, 46, 99, 204  
 Unbalanced designs, 58, 67  
 Unbiased, 24, 203, 213  
 Unconditional test, 80  
 Univariate hypothesis, 79
- Variables, dependent, 82  
 Variance  
   between vs within, 28, 43, 44  
   effect on power, 50  
   reducing between, 49  
   testing equality, 38, 208  
 Variation, 17  
 Vector of observations, 76  
 Virology, 2, 53, 117  
 Vocational guidance, 2
- Weighting variables, 81  
 Withdrawals, 51, 141
- Zero point, 148

# **Springer Series in Statistics**

---

(continued from p. ii)

*Ramsay/Silverman:* Functional Data Analysis.

*Rao/Toutenburg:* Linear Models: Least Squares and Alternatives.

*Read/Cressie:* Goodness-of-Fit Statistics for Discrete Multivariate Data.

*Reinsel:* Elements of Multivariate Time Series Analysis, 2nd edition.

*Reiss:* A Course on Point Processes.

*Reiss:* Approximate Distributions of Order Statistics: With Applications to Non-parametric Statistics.

*Rieder:* Robust Asymptotic Statistics.

*Rosenbaum:* Observational Studies.

*Rosenblatt:* Gaussian and Non-Gaussian Linear Time Series and Random Fields

*Särndal/Swensson/Wretman:* Model Assisted Survey Sampling.

*Schervish:* Theory of Statistics.

*Shao/Tu:* The Jackknife and Bootstrap.

*Siegmund:* Sequential Analysis: Tests and Confidence Intervals.

*Simonoff:* Smoothing Methods in Statistics.

*Singpurwalla and Wilson:* Statistical Methods in Software Engineering:  
Reliability and Risk.

*Small:* The Statistical Theory of Shape.

*Stein:* Interpolation of Spatial Data: Some Theory for Kriging.

*Tanner:* Tools for Statistical Inference: Methods for the Exploration of Posterior  
Distributions and Likelihood Functions, 3rd edition.

*Tong:* The Multivariate Normal Distribution.

*van der Vaart/Wellner:* Weak Convergence and Empirical Processes: With  
Applications to Statistics.

*Weerahandi:* Exact Statistical Methods for Data Analysis.

*West/Harrison:* Bayesian Forecasting and Dynamic Models, 2nd edition.