

# METHODS IN MOLECULAR BIOLOGY

*Series Editor*

John M. Walker

School of Life and Medical Sciences  
University of Hertfordshire  
Hatfield, Hertfordshire, AL10 9AB, UK

For further volumes:  
<http://www.springer.com/series/7651>

# **Statistical Human Genetics**

**Methods and Protocols**

**Second Edition**

Edited by

**Robert C. Elston**

*Case Western Reserve University, Cleveland, Ohio, USA*



*Editor*

Robert C. Elston  
Case Western Reserve University  
Cleveland, Ohio, USA

ISSN 1064-3745

ISSN 1940-6029 (electronic)

Methods in Molecular Biology

ISBN 978-1-4939-7273-9

ISBN 978-1-4939-7274-6 (eBook)

DOI 10.1007/978-1-4939-7274-6

Library of Congress Control Number: 2017951222

© Springer Science+Business Media LLC 2017, corrected publication 2018

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Printed on acid-free paper

This Humana Press imprint is published by Springer Nature  
The registered company is Springer Science+Business Media LLC  
The registered company address is: 233 Spring Street, New York, NY 10013, U.S.A.

---

## Preface

The recent advances in genetics, especially in the molecular techniques that have over the last quarter of a century spectacularly reduced the cost of determining genetic markers, open up a field of research that is becoming of increasing help in detecting, preventing, and/or curing many diseases that afflict us. This has brought with it the need for novel methods of statistical analysis and the implementation of these methods in a wide variety of computer programs. The aim in this book is to make these methods and programs more easily accessible to the beginner who has data to analyze, whether a student or a senior investigator. Twenty-eight of the thirty original chapters have been retained (all updated) and, in view of changes in the field, three new chapters have been added. Apart from the first chapter, which defines some of the genetic terms used in the book, each chapter of this book takes up a particular topic and illustrates the use of at least one piece of software that the authors have found helpful for the relevant statistical analysis of their own human genetic data. There is often more than one program that performs a particular type of analysis and, once you have used one program for a particular analysis, you may find you prefer another program—and there is a good chance you will find that the same basic analysis is described in more than one chapter of this book. You may therefore wish to browse over several chapters, in the first place restricting your reading to only the introductory sections, which describe the underlying theory. The chapters are ordered in the approximate logical order in which human genetic studies are often conducted; so, if you are new to research in human genetics, this initial reading could serve as an introduction to the subject. For the most part, the second section of each chapter gives you step-by-step instructions for running programs and interpreting the program outputs, with extra notes in the third section. However, although the aim is very much to offer a “do it yourself” manual, there may well be times when you will need to consult a statistical geneticist, especially for the interpretation of computer output.

*Cleveland, Ohio, USA*

*Robert C. Elston*

---

## Contents

<i>Preface</i> .....	v
<i>Contributors</i> .....	ix
1 Statistical Genetic Terminology .....	1
<i>Robert C. Elston, Jaya Satagopan, and Shuying Sun</i>	
2 Identification of Genotype Errors .....	11
<i>Jeffery O'Connell and Yin Yao</i>	
3 Detecting Pedigree Relationship Errors .....	25
<i>Lei Sun</i>	
4 Identifying Cryptic Relationships .....	45
<i>Lei Sun, Apostolos Dimitromanolakis, and Wei-Min Chen</i>	
5 Estimating Allele Frequencies .....	61
<i>Indra Adrianto and Courtney Montgomery</i>	
6 Testing Departure from Hardy-Weinberg Proportions .....	83
<i>Jian Wang and Sanjay Shete</i>	
7 Estimating Disequilibrium Coefficients .....	117
<i>Maren Vens and Andreas Ziegler</i>	
8 Detecting Familial Aggregation .....	133
<i>Adam C. Naj and Terri H. Beaty</i>	
9 Estimating Heritability from Twin Studies .....	171
<i>Katrina L. Grasby, Karin J.H. Verweij, Miriam A. Mosing, Brendan P. Zietsch, and Sarah E. Medland</i>	
10 Estimating Heritability from Nuclear Family and Pedigree Data .....	195
<i>Murielle Bochud</i>	
11 Correcting for Ascertainment .....	211
<i>Warren Ewens and Robert C. Elston</i>	
12 Segregation Analysis Using the Unified Model .....	233
<i>Xiangqing Sun</i>	
13 Design Considerations for Genetic Linkage and Association Studies .....	257
<i>Jérémie Nsengimana and D. Timothy Bishop</i>	
14 Model-Based Linkage Analysis of a Quantitative Trait .....	283
<i>Yeunjoo E. Song, Sunah Song, and Audrey H. Schnell</i>	
15 Model-Based Linkage Analysis of a Binary Trait .....	311
<i>Rita M. Cantor</i>	
16 Model-Free Linkage Analysis of a Quantitative Trait .....	327
<i>Nathan J. Morris and Catherine M. Stein</i>	
17 Model-Free Linkage Analysis of a Binary Trait .....	343
<i>Wei Xu, Jin Ma, Celia M.T. Greenwood, Andrew D. Paterson, and Shelley B. Bull</i>	

18	Single Marker Association Analysis for Unrelated Samples . . . . .	375
	<i>Gang Zheng, Ao Yuan, Qizhai Li, and Joseph L. Gastwirth</i>	
19	Single Marker Family-Based Association Analysis Conditional on Parental Information . . . . .	391
	<i>Ren-Hua Chung, Daniel D. Kinnaman, and Eden R. Martin</i>	
20	Single Marker Family-Based Association Analysis Not Conditional on Parental Information . . . . .	409
	<i>Junghyun Namkung and Sungbo Won</i>	
21	Calibrating Population Stratification in Association Analysis . . . . .	441
	<i>Huaizhen Qin and Xiaofeng Zhu</i>	
22	Cross-Phenotype Association Analysis Using Summary Statistics from GWAS . . . . .	455
	<i>Xiaoyin Li and Xiaofeng Zhu</i>	
23	Haplotype Inference . . . . .	469
	<i>Sunah Song, Xin Li, and Jing Li</i>	
24	Multi-SNP Haplotype Analysis Methods for Association Analysis . . . . .	485
	<i>Daniel O. Stram</i>	
25	The Analysis of Ethnic Mixtures . . . . .	505
	<i>Xiaofeng Zhu and Heming Wang</i>	
26	Detecting Multiethnic Rare Variants . . . . .	527
	<i>Weiwei Ouyang, Xiaofeng Zhu, and Huaizhen Qin</i>	
27	Identifying Gene Interaction Networks . . . . .	539
	<i>Danica Wiredja and Gurkan Bebek</i>	
28	Structural Equation Modeling . . . . .	557
	<i>Catherine M. Stein, Nathan J. Morris, Noémi B. Hall, and Nora L. Nock</i>	
29	Mendelian Randomization . . . . .	581
	<i>Sandeep Grover, Fabiola Del Greco M., Catherine M. Stein, and Andreas Ziegler</i>	
30	Preprocessing and Quality Control for Whole-Genome Sequences from the Illumina HiSeq X Platform . . . . .	629
	<i>Marvin N. Wright, Damian Gola, and Andreas Ziegler</i>	
31	Processing and Analyzing Human Microbiome Data . . . . .	649
	<i>Xuan Zhu, Jian Wang, Cielito Reyes-Gibby, and Sanjay Shete</i>	
	Correction to: Processing and Analyzing Human Microbiome Data . . . . .	E1
	<i>Index</i> . . . . .	679

---

## Contributors

INDRA ADRIANTO • *Arthritis and Clinical Immunology Research Program, Oklahoma Medical Research Foundation, Oklahoma City, OK, USA*

TERRI H. BEATY • *Department of Epidemiology, Johns Hopkins University Bloomberg School of Public Health, Baltimore, MD, USA*

GURKAN BEBEK • *Systems Biology and Bioinformatics Graduate Program, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Center for Proteomics and Bioinformatics, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Department of Nutrition, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Department of Electrical Engineering and Computer Science, Case Western Reserve University School of Medicine, Cleveland, OH, USA*

D. TIMOTHY BISHOP • *Section of Epidemiology and Biostatistics, Leeds Institute of Cancer and Pathology, University of Leeds, Leeds, UK*

MURIELLE BOCHUD • *Institute of Social and Preventive Medicine, Lausanne University Hospital, Lausanne, Switzerland*

SHELLEY B. BULL • *Lunenfeld-Tanenbaum Research Institute, Sinai Health System, Toronto, ON, Canada; Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada*

RITA M. CANTOR • *Department of Human Genetics, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA; Center for Neurobehavioral Genetics, Department of Psychiatry, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA*

WEI-MIN CHEN • *Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA; Department of Public Health Sciences, University of Virginia, Charlottesville, VA, USA*

REN-HUA CHUNG • *Division of Biostatistics and Bioinformatics, Institute of Population Health Sciences, National Health Research Institutes, Miaoli County, Taiwan*

FABIOLA DEL GRECO M. • *Center for Biomedicine, EURAC Research, Bolzano, Italy*

APOSTOLOS DIMITROMANOLAKIS • *Department of Statistical Sciences, Faculty of Arts and Science, Toronto, ON, Canada; Lunenfeld-Tanenbaum Research Institute, Mount Sinai Hospital, Toronto, ON, Canada*

ROBERT C. ELSTON • *Case Western Reserve University, Cleveland, OH, USA*

WARREN EWENS • *Department of Biology, University of Pennsylvania, Philadelphia, PA, USA*

JOSEPH L. GASTWIRTH • *Department of Statistics, George Washington University, Washington, DC, USA*

DAMIAN GOLA • *Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein - Campus Lübeck, Lübeck, Germany*

KATRINA L. GRASBY • *Genetic Epidemiology, QIMR Berghofer Medical Research Institute, Brisbane, Australia*

CELIA M.T. GREENWOOD • *Lady Davis Research Institute, Jewish General Hospital, Montréal, QC, Canada; Department of Oncology and Department of Epidemiology, Biostatistics & Occupational Health, McGill University, QC, Canada*

SANDEEP GROVER • *Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Lübeck, Germany*

- NOÉMI B. HALL • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- DANIEL D. KINNAMON • *Division of Human Genetics, Department of Internal Medicine, The Ohio State University Wexner Medical Center, Columbus, OH, USA*
- QIZHAI LI • *Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing, China*
- XIAOYIN LI • *Department of Population and Quantitative Health Sciences, School of Medicine, Case Western Reserve University, Cleveland, OH, USA*
- XIN LI • *Department of Electrical Engineering and Computer Science, Case Western Reserve University, Cleveland, OH, USA*
- JING LI • *Department of Electrical Engineering and Computer Science, Case Western Reserve University, Cleveland, OH, USA*
- JIN MA • *Lunenfeld-Tanenbaum Research Institute, Sinai Health System, Toronto, ON, Canada*
- EDEN R. MARTIN • *John P. Hussman Institute for Human Genomics, Leonard M. Miller School of Medicine, University of Miami, Miami, FL, USA*
- SARAH E. MEDLAND • *Genetic Epidemiology, QIMR Berghofer Medical Research Institute, Brisbane, Australia*
- COURTNEY MONTGOMERY • *Arthritis and Clinical Immunology Research Program, Oklahoma Medical Research Foundation, Oklahoma City, OK, USA*
- NATHAN J. MORRIS • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- MIRIAM A. MOSING • *Department of Neuroscience, Karolinska Institute, Stockholm, Sweden; Department of Medical Epidemiology and Biostatistics, Karolinska Institute, Stockholm, Sweden*
- ADAM C. NAJ • *Department of Biostatistics, Epidemiology, and Informatics, University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, USA; Department of Pathology and Laboratory Medicine, University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, USA; Center for Clinical Epidemiology and Biostatistics, University of Pennsylvania Perelman School of Medicine, Philadelphia, PA, USA*
- JUNGHYUN NAMKUNG • *Molecular Diagnostics Team, IVD Business Unit, SK Telecom, Seoul, South Korea*
- NORA L. NOCK • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- JÉRÉMIE NSENGIMANA • *Section of Epidemiology and Biostatistics, Leeds Institute of Cancer and Pathology, University of Leeds, Leeds, UK*
- JEFFERY O'CONNELL • *University of Maryland, Baltimore, MD, USA*
- WEIWEI OUYANG • *Department of Global Biostatistics and Data Science, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA, USA*
- ANDREW D. PATERSON • *Genetics and Genome Biology, The Hospital for Sick Children, Toronto, ON, Canada; Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada*
- HUAIZHEN QIN • *Department of Global Biostatistics and Data Science, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA, USA; Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- CIELITO REYES-GIBBY • *Department of Emergency Medicine, The University of Texas MD Anderson Cancer Center, Houston, TX, USA*

- JAYA SATAGOPAN • *Memorial Sloan Kettering Cancer Center, New York, NY, USA*
- AUDREY H. SCHNELL • *Cardiovascular Imaging Core Laboratory, Harrington Heart & Vascular Institute, University Hospitals Cleveland Medical Center, Cleveland, OH, USA*
- SANJAY SHETE • *Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston, TX, USA; Department of Epidemiology, The University of Texas MD Anderson Cancer Center, Houston, TX, USA*
- YEUNJOO E. SONG • *Department of Population and Quantitative Health Sciences, Case Western Reserve University, Cleveland, OH, USA*
- SUNAH SONG • *Case Western Reserve University, Cleveland, OH, USA*
- CATHERINE M. STEIN • *Center for Proteomics and Bioinformatics, Case Western Reserve University, Cleveland, OH, USA*
- DANIEL O. STRAM • *Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA*
- SHUYING SUN • *Texas State University, San Marcos, TX, USA*
- LEI SUN • *Department of Statistical Sciences, Faculty of Arts and Science, University of Toronto, Toronto, ON, Canada; Division of Biostatistics, Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada*
- XIANGQING SUN • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- MAREN VENS • *Institut für Medizinische Biometrie und Epidemiologie, Universitätsklinikum Hamburg-Eppendorf, Hamburg, Germany*
- KARIN J.H. VERWEIJ • *Genetic Epidemiology, QIMR Berghofer Medical Research Institute, Brisbane, Australia; Behavioural Science Institute, Radboud University, HR Nijmegen, The Netherlands*
- JIAN WANG • *Department of Biostatistics, University of Texas MD Anderson Cancer Center, Houston, TX, USA*
- HEMING WANG • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Division of Sleep and Circadian Disorders, Brigham and Women's Hospital, Boston, MA, USA; Division of Sleep Medicine, Harvard Medical School, Boston, MA, USA*
- DANICA WIREDJA • *Systems Biology and Bioinformatics Graduate Program, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Center for Proteomics and Bioinformatics, Case Western Reserve University School of Medicine, Cleveland, OH, USA; Department of Nutrition, Case Western Reserve University School of Medicine, Cleveland, OH, USA*
- SUNGHO WON • *Department of Public Health Science, Graduate School of Public Health, Seoul National University, Seoul, South Korea*
- MARVIN N. WRIGHT • *Institut für Medizinische Biometrie und Statistik, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein - Campus Lübeck, Lübeck, Germany*
- WEI XU • *Department of Biostatistics, Princess Margaret Cancer Centre, University Health Network, Toronto, ON, Canada; Dalla Lana School of Public Health, University of Toronto, Toronto, ON, Canada*
- YIN YAO • *Unit of Genomic Statistics, Intramural Research Program, National Institute of Mental Health, Bethesda, MD, USA*
- AO YUAN • *Department of Biostatistics, Bioinformatics and Biomathematics, Georgetown University, Washington, DC, USA*
- XIAOFENG ZHU • *Department of Population and Quantitative Health Sciences, Case Western Reserve University School of Medicine, Cleveland, OH, USA*

XUAN ZHU • *Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston, TX, USA*

ANDREAS ZIEGLER • *Institut für Medizinische Biometrie und Statistik & Zentrum für klinische Studien, Universität zu Lübeck, Universitätsklinikum Schleswig-Holstein, Lübeck, Germany*

BRENDAN P. ZIETSCH • *School of Psychology, University of Queensland, Brisbane, Australia*