

METHODS IN MOLECULAR BIOLOGY™

Series Editor
John M. Walker
School of Life Sciences
University of Hertfordshire
Hatfield, Hertfordshire, AL10 9AB, UK

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

Statistical Methods for Microarray Data Analysis

Methods and Protocols

Edited by

Andrei Y. Yakovlev

*School of Medicine and Dentistry, Department of Biostatistics and Computational Biology,
University of Rochester, Rochester, NY, USA*

Lev Klebanov

Department of Probability and Statistics, Charles University, Prague, Czech Republic

Daniel Gaile

State University of New York at Buffalo, Buffalo, NY, USA

Editors

Andrei Y. Yakovlev
School of Medicine and Dentistry
Department of Biostatistics and Computational
Biology
University of Rochester
Rochester, NY, USA

Lev Klebanov
Department of Probability and Statistics
Charles University
Prague, Czech Republic

Daniel Gaile
State University of New York at Buffalo
Buffalo, NY, USA

ISSN 1064-3745 ISSN 1940-6029 (electronic)
ISBN 978-1-60327-336-7 ISBN 978-1-60327-337-4 (eBook)
DOI 10.1007/978-1-60327-337-4
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2012956545

© Springer Science+Business Media New York 2013

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. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

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.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Humana Press is a brand of Springer
Springer is part of Springer Science+Business Media (www.springer.com)

Preface

Microarrays for simultaneous measurement of redundancy of RNA species are widely used for fundamental biology research. They are also tested for their use in personalized medicine in disease diagnosis and prognosis. From the point of mathematical statistics, the invention of microarray technology in the mid-1990s allowed the simultaneous monitoring of the expression levels of thousands of genes. Microarrays for simultaneous measurement of redundancy of RNA species are used in fundamental biology as well as in medical research. Microarray may be considered as an observation of very high dimensionality equal to the number of expression levels. There arise some needs to develop new statistical methods to handle the data of such large dimensionality, especially connected to the fact of small number of observations (which is the number of arrays). Because of the small number of observations the standard asymptotic methods of multivariate statistical analysis appear to be inapplicable.

The aim of the book is to familiarize the readers with statistical methods used nowadays in microarray analysis. It is addressed to everybody who is involved or is planning to be involved in statistical data analysis of microarrays, mostly to statisticians but also to biological researchers.

It was impossible to come over all statistical methods published to date in a book. The selection was made on a basis of mathematical correctness of corresponding methods. Some approaches based on intuitive impressions and having no mathematical support were excluded from our consideration. However, the selection criteria were sometimes not scientific. In many cases they reflect personal taste of the editors. Nevertheless, the editors took pains that other valuable methods should be described or mentioned. The editors invite the interested readers to continue their study of the material beyond this book. We are very grateful to all authors for their writing.

Each chapter can be read as a separate entry. The style of the individual writing is essential to show the knowledge and experience of each team that contributed. However, the reader will be guided from microarray technology to statistical problems of corresponding data analysis. Chapters 1 and 2 provide such prolegomena.

In Chapter 3, an introduction to current multiple testing methodology are presented, with the objective of clarifying the methodological issues involved, and hopefully providing the reader with some basis with which to compare and select methods.

Chapter 4 discusses a method of selecting differentially expressed genes based on a newly discovered structure termed as the δ -sequence. Together with the nonparametric empirical Bayes methodology, it leads to dramatic gains in terms of the mean numbers of true and false discoveries, and in the stability of the results of testing. The results of this chapter can be viewed also as a new method of normalization of microarray data.

Chapter 5 is in some sense connected to Chapter 4. It studies different normalization procedures of gene expression levels. Normalization procedures are used for removing systematical variation which affects the measure of expression levels.

Chapter 6 is dedicated to constructing of multivariate prognostic gene signatures with censored survival data. Modern high-throughput technologies allow us to simultaneously measure the expressions of a huge number of candidate predictors, some of which are likely to be associated with survival. One difficult task is to search among an enormous number of potential predictors and to correctly identify most of the important ones, without mistakenly identifying too many spurious associations. Mere variable selection is insufficient, however, for the information from the multiple predictors must be intelligently combined and calibrated to form the final composite predictor. Many commonly used procedures overfit the training data, miss many important predictors, or both. Author proposes a method that offers a middle ground where some partial multivariate adjustments can be made in an adaptive fashion, regardless of the number of candidate predictors. He demonstrates the performance of our proposed procedure in a simulation study within the Cox proportional hazards regression framework, and applies this new method to a publicly available data set to construct a novel prognostic gene signature for breast cancer survival.

Chapter 7 considers clustering problems for gene-expression data.

There are two distinct but related clustering problems with microarray data. One problem concerns the clustering of the tissue samples (gene signatures) on the basis of the genes; the other concerns the clustering of the genes on the basis of the tissues (gene profiles). The clusters of tissues so obtained in the first problem can play a useful role in the discovery and understanding of new subclasses of diseases. The clusters of genes obtained in the second problem can be used to search for genetic pathways or groups of genes that might be regulated together. Authors focus here on mixtures of normals to provide a model-based clustering of tissue samples (gene signatures) and of gene profiles.

Network-based analysis of multivariate gene expression data is given in Chapter 8. Such data are collected to study genomic responses under special conditions for which the expression levels of given genes are expected to be dependent. One important question from such multivariate gene expression experiments is to identify genes that show different expression patterns over treatment dosages or over time; these genes can also point to the pathways that are perturbed during a given biological process. Several empirical Bayes approaches have been developed for identifying the differentially expressed genes in order to account for the parallel structure of the data and to borrow information across all the genes.

In Chapter 9, author discusses the statistical problem, termed oncogene outlier detection, and discusses a variety of proposals to this problem. A statistical model in the multi-class situation is described; links with multiple testing concepts are established. Some new nonparametric procedures are described and compared to existing methods using simulation studies.

Data quality is intrinsically influenced by design, technical, and analytical parameters. Quality parameters have not yet been well defined for gene expression analysis by microarrays, though ad interim, following recommended good experimental practice guidelines should ensure generation of reliable and reproducible data. In Chapter 10 author summarizes essential practical recommendations for experimental design, technical considerations, feature annotation issues, and standardization efforts.

Inferring gene regulatory networks from microarray data has become a popular activity in recent years, resulting in an ever increasing volume of publications. There are many pitfalls in network analysis that remain either unnoticed or scantily understood. A critical discussion of such pitfalls is long overdue. Chapter 11 discuss one feature of microarray data the investigators need to be aware of when embarking on a study of putative associations between elements of networks and pathways.

Finally, Chapter 12 considers the problem of normality of logs of gene expression levels. In the literature there is no unique point on the fact of normality (or nonnormality) of the distribution of gene expression levels. This chapter discusses different approaches to testing of normality in this situation.

The editors and the contributors assume from the reader a basic knowledge of biological concepts of gene expression and statistical methods of gene expression analysis.

Editors are grateful to all contributors.

Prague, Czech Republic

Lev Klebanov

Contents

<i>Preface</i>	<i>v</i>
<i>Contributors</i>	<i>xi</i>
1 What Statisticians Should Know About Microarray Gene Expression Technology	1
<i>Stephen Welle</i>	
2 Where Statistics and Molecular Microarray Experiments Biology Meet	15
<i>Diana M. Kelmansky</i>	
3 Multiple Hypothesis Testing: A Methodological Overview	37
<i>Anthony Almudevar</i>	
4 Gene Selection with the δ -Sequence Method.	57
<i>Xing Qiu and Lev Klebanov</i>	
5 Using of Normalizations for Gene Expression Analysis	73
<i>Peter Bubeliny</i>	
6 Constructing Multivariate Prognostic Gene Signatures with Censored Survival Data	85
<i>Derick R. Peterson</i>	
7 Clustering of Gene Expression Data Via Normal Mixture Models	103
<i>G.J. McLachlan, L.K. Flack, S.K. Ng, and K. Wang</i>	
8 Network-Based Analysis of Multivariate Gene Expression Data	121
<i>Wei Zhi, Jane Minturn, Eric Rappaport, Garrett Brodeur, and Hongzhe Li</i>	
9 Genomic Outlier Detection in High-Throughput Data Analysis	141
<i>Debashis Ghosh</i>	
10 Impact of Experimental Noise and Annotation Imprecision on Data Quality in Microarray Experiments.	155
<i>Andreas Scherer, Manhong Dai, and Fan Meng</i>	
11 Aggregation Effect in Microarray Data Analysis.	177
<i>Linlin Chen, Anthony Almudevar, and Lev Klebanov</i>	
12 Test for Normality of the Gene Expression Data	193
<i>Bobosharif Shokirov</i>	
<i>Index</i>	<i>209</i>

Contributors

- ANTHONY ALMUDEVAR • *Department of Biostatistics and Computational Biology, University of Rochester, Rochester, NY, USA*
- GARRETT BRODEUR • *Department of Pediatrics, Children's Hospital of Philadelphia, Philadelphia, PA, USA*
- PETER BUBELÍNY • *Department of Probability and Statistics, Charles University, Prague, Czech Republic*
- LINLIN CHEN • *School of Mathematical Sciences, Rochester Institute of Technology, Rochester, NY, USA*
- MANHONG DAI • *Psychiatry Department and Molecular and Behavioral Neuroscience Institute, University of Michigan, University of Michigan, Ann Arbor, MI, USA*
- L.K. FLACK • *Department of Rheumatology, University of NSW, Ryde, NSW, Australia*
- DEBASHIS GHOSH • *Departments of Statistics and Public Health Sciences, Penn State University, DuBios, PA, USA*
- DIANA M. KELMANSKY • *Instituto de Cálculo, Ciudad Universitaria, Buenos Aires, Argentina*
- LEV KLEBANOV • *Department of Probability and Statistics Charles University Prague, Czech Republic*
- HONGZHE LI • *Department of Biostatistics and Epidemiology, University of Pennsylvania School of Medicine, Philadelphia, PA, USA*
Department of Pediatrics, Children's Hospital of Philadelphia, Philadelphia, PA, USA
- G.J. McLACHLAN • *Department of Mathematics, University of Queensland, Brisbane, Australia*
- FAN MENG • *Psychiatry Department and Molecular and Behavioral Neuroscience Institute, University of Michigan, University of Michigan, Ann Arbor, MI, USA*
- JANE MINTURN • *Department of Pediatrics, Children's Hospital of Philadelphia, Philadelphia, PA, USA*
- S.K. NG • *School of Medicine, Griffith University, Meadowbrook, Australia*
- DERICK R. PETERSON • *Department of Biostatistics and Computational Biology, University of Rochester, Rochester, NY, USA*
- XING QIU • *Department of Biostatistics and Computational Biology, University of Rochester, Rochester, NY, USA*
- ERIC RAPPAPORT • *Children's Hospital of Philadelphia, Philadelphia, PA, USA*
- ANDREAS SCHERER • *Genomics, Biomarker Development, Spheromics, Kontiolahdi, Joensuu, Finland*
- BOBOSHARIF SHOKIROV • *Department of Probability and Statistics, MFF, Charles University, Prague, Czech Republic*
- K. WANG • *University of Queensland, Melbourne, Australia*
- STEPHEN WELLE • *Functional Genomics Center, University of Rochester, Rochester, NY, USA*
- WEI ZHI • *Department of Biostatistics and Epidemiology, New Jersey Institute of Technology, Newark, NJ, USA*