

BAYESIAN INFERENCE FOR GENE EXPRESSION AND PROTEOMICS

The interdisciplinary nature of bioinformatics presents a research challenge in integrating concepts, methods, software, and multiplatform data. Although there have been rapid developments in new technology and an inundation of statistical methodology and software for the analysis of microarray gene expression arrays, there exist few rigorous statistical methods for addressing other types of high-throughput data, such as proteomic profiles that arise from mass spectrometry experiments. This book discusses the development and application of Bayesian methods in the analysis of high-throughput bioinformatics data that arise from medical, in particular cancer, research, as well as molecular and structural biology. The Bayesian approach has the advantage that evidence can be easily and flexibly incorporated into statistical models.

A basic overview of the biological and technical principles behind multiplatform high-throughput experimentation is followed by expert reviews of Bayesian methodology, tools, and software for single group inference, group comparisons, classification and clustering, motif discovery and regulatory networks, and Bayesian networks and gene interactions.

Kim-Anh Do is a professor in the Department of Biostatistics and Applied Mathematics and the University of Texas M.D. Anderson Cancer Center. Her research interests are in computer-intensive statistical methods with recent focus in the development of methodology and software to analyze data produced from high-throughput technologies.

Peter Müller is also a professor in the Department of Biostatistics and Applied Mathematics and the University of Texas M.D. Anderson Cancer Center. His research interests and contributions are in the areas of Markov chain Monte Carlo posterior simulation, nonparametric Bayesian inference, hierarchical models, mixture models, and Bayesian decision problems.

Marina Vannucci is a professor of Statistics at Texas A&M University. Her research focuses on the theory and practice of Bayesian variable selection techniques and on the development of wavelet-based statistical models and their applications. Her work is often motivated by real problems that need to be addressed with suitable statistical methods.

BAYESIAN INFERENCE FOR GENE EXPRESSION AND PROTEOMICS

Edited by

KIM-ANH DO

University of Texas M.D. Anderson Cancer Center

PETER MÜLLER

University of Texas M.D. Anderson Cancer Center

MARINA VANNUCCI

Texas A&M University



CAMBRIDGE
UNIVERSITY PRESS

Cambridge University Press & Assessment
 978-1-107-63698-9 — Bayesian Inference for Gene Expression and Proteomics
 Edited by Kim-Anh Do, Peter Müller, Marina Vannucci
 Frontmatter
[More Information](#)

CAMBRIDGE
 UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom
 One Liberty Plaza, 20th Floor, New York, NY 10006, USA
 477 Williamstown Road, Port Melbourne, VIC 3207, Australia
 314-321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi - 110025, India
 103 Penang Road, #05-06/07, Visioncrest Commercial, Singapore 238467

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107636989

© Cambridge University Press 2006

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2006

A catalogue record for this publication is available from the British Library

Library of Congress Cataloging in Publication data
 Bayesian inference for gene expression and proteomics /
 edited by Kim-Anh Do, Peter Müller, Marina Vannucci.
 p. cm.

Includes bibliographical references.

ISBN-13: 978-0-521-86092-5 (hardback)

ISBN-10: 0-521-86092-X (hardback)

1. Gene expression – Statistical methods. 2. Proteomics –
 Statistical methods. I. Do, Kim-Anh, 1960– II. Müller, Peter, 1963–
 III. Vannucci, Marina, 1966– IV. Title.

QH450.B39 2006

572.8'6501519542 – dc22 2006005635

ISBN 978-0-521-86092-5 Hardback

ISBN 978-1-107-63698-9 Paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Contents

<i>List of Contributors</i>	<i>page xi</i>
<i>Preface</i>	xv
1 An Introduction to High-Throughput Bioinformatics Data	
<i>Keith A. Baggerly, Kevin R. Coombes, and Jeffrey S. Morris</i>	1
1.1 Introduction	1
1.2 Microarrays	2
1.3 SAGE	19
1.4 Mass Spectrometry	24
1.5 Finding Data	34
2 Hierarchical Mixture Models for Expression Profiles	
<i>Michael A. Newton, Ping Wang, and Christina Kendziorski</i>	40
2.1 Introduction	40
2.2 Dual Character of Posterior Probabilities	43
2.3 Differential Expression as Independence	45
2.4 The Multigroup Mixture Model	47
2.5 Improving Flexibility	49
3 Bayesian Hierarchical Models for Inference in Microarray Data	
<i>Anne-Mette K. Hein, Alex Lewin, and Sylvia Richardson</i>	53
3.1 Introduction	53
3.2 Bayesian Hierarchical Modeling of Probe Level GeneChip Data	56
3.3 Bayesian Hierarchical Model for Normalization and Differential Expression	67
3.4 Predictive Model Checking	70

vi	<i>Contents</i>	
4	Bayesian Process-Based Modeling of Two-Channel Microarray Experiments: Estimating Absolute mRNA Concentrations <i>Mark A. van de Wiel, Marit Holden, Ingrid K. Glad, Heidi Lyng, and Arnaldo Frigessi</i>	75
4.1	Introduction	75
4.2	The Hierarchical Model	78
4.3	Reparameterization and Identifiability	82
4.4	MCMC-Based Inference	84
4.5	Validation	85
4.6	Illustration	85
4.7	TransCount Web Site and Computing Times	91
4.8	A Statistical Discussion of the Model	91
4.9	Discussion	93
5	Identification of Biomarkers in Classification and Clustering of High-Throughput Data <i>Mahlet G. Tadesse, Naijun Sha, Sinae Kim, and Marina Vannucci</i>	97
5.1	Introduction	97
5.2	Bayesian Variable Selection in Linear Models	100
5.3	Bayesian Variable Selection in Classification	101
5.4	Bayesian Variable Selection in Clustering via Finite Mixture Models	103
5.5	Bayesian Variable Selection in Clustering via Dirichlet Process Mixture Models	106
5.6	Example: Leukemia Gene Expression Data	108
5.7	Conclusion	113
6	Modeling Nonlinear Gene Interactions Using Bayesian MARS <i>Veerabhadran Baladandayuthapani, Chris C. Holmes, Bani K. Mallick, and Raymond J. Carroll</i>	116
6.1	Introduction	116
6.2	Bayesian MARS Model for Gene Interaction	118
6.3	Computation	121
6.4	Prediction and Model Choice	122
6.5	Examples	123
6.6	Discussion and Summary	131
7	Models for Probability of Under- and Overexpression: The POE Scale <i>Elizabeth Garrett-Mayer and Robert Scharpf</i>	137
7.1	POE: A Latent Variable Mixture Model	137
7.2	The POE Model	138
7.3	Unsupervised versus Semisupervised POE	144

<i>Contents</i>		vii
7.4	Using POE Scale	145
7.5	Example: POE as Applied to Lung Cancer Microarray Data	148
7.6	Discussion	152
8	Sparse Statistical Modelling in Gene Expression Genomics <i>Joseph Lucas, Carlos Carvalho, Quanli Wang, Andrea Bild, Joseph R. Nevins, and Mike West</i>	155
8.1	Perspective	156
8.2	Sparse Regression Modelling	157
8.3	Sparse Regression for Artifact Correction with Affymetrix Expression Arrays	162
8.4	Sparse Latent Factor Models and Latent Factor Regressions	167
8.5	Concluding Comments	173
9	Bayesian Analysis of Cell Cycle Gene Expression Data <i>Chuan Zhou, Jon C. Wakefield, and Linda L. Breeden</i>	177
9.1	Introduction	177
9.2	Previous Studies	178
9.3	Data	180
9.4	Bayesian Analysis of Cell Cycle Data	182
9.5	Discussion	197
10	Model-Based Clustering for Expression Data via a Dirichlet Process Mixture Model <i>David B. Dahl</i>	201
10.1	Introduction	201
10.2	Model	203
10.3	Inference	208
10.4	Simulation Study	209
10.5	Example	212
10.6	Conclusion	216
11	Interval Mapping for Expression Quantitative Trait Loci <i>Meng Chen and Christina Kendziorski</i>	219
11.1	Introduction	219
11.2	eQTL Mapping Experiments	221
11.3	QTL Mapping Methods	222
11.4	Currently Available eQTL Mapping Methods	223
11.5	MOM Interval Mapping	225
11.6	Discussion	231
12	Bayesian Mixture Models for Gene Expression and Protein Profiles <i>Michele Guindani, Kim-Anh Do, Peter Müller, and Jeffrey S. Morris</i>	238
12.1	Introduction	238

12.2	A Nonparametric Bayesian Model for Differential Gene Expression	240
12.3	A Mixture of Beta Model for MALDI-TOF Data	243
12.4	A Semiparametric Mixture Model for SAGE Data	247
12.5	Summary	250
13	Shrinkage Estimation for SAGE Data Using a Mixture Dirichlet Prior	
	<i>Jeffrey S. Morris, Keith A. Baggerly, and Kevin R. Coombes</i>	254
13.1	Introduction	254
13.2	Overview of SAGE	255
13.3	Methods for Estimating Relative Abundances	257
13.4	Mixture Dirichlet Distribution	260
13.5	Implementation Details	263
13.6	Simulation Study	264
13.7	Conclusion	267
14	Analysis of Mass Spectrometry Data Using Bayesian Wavelet-Based Functional Mixed Models	
	<i>Jeffrey S. Morris, Philip J. Brown, Keith A. Baggerly, and Kevin R. Coombes</i>	269
14.1	Introduction	270
14.2	Overview of MALDI-TOF	270
14.3	Functional Mixed Models	274
14.4	Wavelet-Based Functional Mixed Models	276
14.5	Analyzing Mass Spectrometry Data Using Wavelet-Based Functional Mixed Models	280
14.6	Conclusion	288
15	Nonparametric Models for Proteomic Peak Identification and Quantification	
	<i>Merlise A. Clyde, Leanna L. House, and Robert L. Wolpert</i>	293
15.1	Introduction	293
15.2	Kernel Models for Spectra	294
15.3	Prior Distributions	296
15.4	Likelihood	301
15.5	Posterior Inference	302
15.6	Illustration	303
15.7	Summary	305
16	Bayesian Modeling and Inference for Sequence Motif Discovery	
	<i>Mayetri Gupta and Jun S. Liu</i>	309
16.1	Introduction	309
16.2	Biology of Transcription Regulation	311

<i>Contents</i>		ix
16.3	Problem Formulation, Background, and General Strategies	312
16.4	A Bayesian Approach to Motif Discovery	316
16.5	Extensions of the Product-Multinomial Motif Model	320
16.6	HMM-Type Models for Regulatory Modules	321
16.7	Model Selection through a Bayesian Approach	327
16.8	Discussion: Motif Discovery Beyond Sequence Analysis	329
17	Identification of DNA Regulatory Motifs and Regulators by Integrating Gene Expression and Sequence Data <i>Deukwoo Kwon, Sinae Kim, David B. Dahl, Michael Swartz, Mahlet G. Tadesse, and Marina Vannucci</i>	333
17.1	Introduction	333
17.2	Integrating Gene Expression and Sequence Data	335
17.3	A Model for the Identification of Regulatory Motifs	337
17.4	Identification of Regulatory Motifs and Regulators	340
17.5	Conclusion	344
18	A Misclassification Model for Inferring Transcriptional Regulatory Networks <i>Ning Sun and Hongyu Zhao</i>	347
18.1	Introduction	347
18.2	Methods	348
18.3	Simulation Results	355
18.4	Application to Yeast Cell Cycle Data	360
18.5	Discussion	361
19	Estimating Cellular Signaling from Transcription Data <i>Andrew V. Kossenkov, Ghislain Bidaut, and Michael F. Ochs</i>	366
19.1	Introduction	366
19.2	Bayesian Decomposition	370
19.3	Key Biological Databases	373
19.4	Example: Signaling Activity in <i>Saccharomyces cerevisiae</i>	376
19.5	Conclusion	380
20	Computational Methods for Learning Bayesian Networks from High-Throughput Biological Data <i>Bradley M. Broom and Devika Subramanian</i>	385
20.1	Introduction	385
20.2	Bayesian Networks	387
20.3	Learning Bayesian Networks	389
20.4	Algorithms for Learning Bayesian Networks	391
20.5	Example: Learning Robust Features from Data	395
20.6	Conclusion	398

21	Bayesian Networks and Informative Priors: Transcriptional Regulatory Network Models	
	<i>Alexander J. Hartemink</i>	401
21.1	Introduction	401
21.2	Bayesian Networks and Bayesian Network Inference	403
21.3	Adding Informative Structure Priors	407
21.4	Applications of Informative Structure Priors	409
21.5	Adding Informative Parameter Priors	418
21.6	Discussion	419
21.7	Availability of Papers and Banjo Software	421
21.8	Acknowledgments	421
22	Sample Size Choice for Microarray Experiments	
	<i>Peter Müller, Christian Robert, and Judith Rousseau</i>	425
22.1	Introduction	425
22.2	Optimal Sample Size as a Decision Problem	428
22.3	Monte Carlo Evaluation of Predictive Power	431
22.4	The Probability Model	432
22.5	Pilot Data	435
22.6	Example	435
22.7	Conclusion	436

Contributors

Keith A. Baggerly, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Veerabhadran Baladandayuthapani, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Ghislain Bidaut, *Department of Genetics, The University of Pennsylvania School of Medicine, 1423 Blockley Hall, 423 Guardian Drive, Philadelphia, PA 19104-6021*

Andrea Bild, *Institute for Genome Sciences and Policy, Duke University, Durham, NC 27710*

Linda L. Breeden, *Fred Hutchinson Cancer Research Center, 1100 Fairview Ave. N., Mailstop A2-168, P.O. Box 19024, Seattle, WA 98109-1024*

Bradley M. Broom, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Philip J. Brown, *Institute of Mathematics, Statistics and Actuarial Science, Room E218, Cornwallis Building, University of Kent, Canterbury, Kent CT2 7NF, UK*

Raymond J. Carroll, *Department of Statistics, Texas A&M University, 3143 TAMU, College Station, TX 77843-3143*

Carlos Carvalho, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Meng Chen, *Department of Statistics, The University of Wisconsin-Madison, 1220 Medical Sciences Center, 1300 University Ave., Madison, WI 53703*

Merlise A. Clyde, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Kevin R. Coombes, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

David B. Dahl, *Department of Statistics, Texas A&M University, 3143 TAMU, College Station, TX 77843-3143*

Kim-Anh Do, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Arnoldo Frigessi, *Department of Biostatistics, University of Oslo, P.O. Box 1122, Blindern, 0317 Oslo, Norway*

Elizabeth Garrett-Mayer, *Johns Hopkins Kimmel Cancer Center, Johns Hopkins University, Suite 1103, 550 N. Broadway, Baltimore, MD 21205*

Ingrid K. Glad, *Department of Mathematics, University of Oslo, P.O. Box 1053, Blindern, 0316 Oslo, Norway*

Michele Guindani, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Mayetri Gupta, *Department of Biostatistics, University of North Carolina at Chapel Hill, McGavran, Greenberg Hall B CB#7420, Chapel Hill, NC 27599-7420*

Alexander J. Hartemink, *Department of Computer Science, Duke University, Box 90129, Durham, NC 27708-0129*

Anne-Mette K. Hein, *Department of Epidemiology and Public Health, Imperial School of Medicine, St. Mary's Campus, Norfolk Place, London W2 1PG, UK*

Marit Holden, *Norwegian Computing Center, P.O. Box 114, Blindern, 0314 Oslo, Norway*

Chris C. Holmes, *Department of Statistics, University of Oxford, 1 South Parks Road, Oxford OX1 3TG, UK*

Leanna L. House, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Christina Kendziorski, *Department of Biostatistics and Medical Informatics, 6785 Medical Sciences Center, 1300 University Ave., Madison, WI 53703*

Sinae Kim, *Department of Biostatistics, School of Public Health, University of Michigan, 1420 Washington Heights, Ann Arbor, MI 48109-2029*

Andrew V. Kossenkov, *Division of Population Science, Fox Chase Cancer Center, 333 Cottman Ave., Philadelphia, PA 19111-2497*

Deukwoo Kwon, *Radiation Epidemiology Branch, National Cancer Institute, 6120 Executive Blvd, MSC 7238, Executive Plaza South, Room 7045, Bethesda, MD 20892-7238*

Contributors

xiii

Alex Lewin, *Department of Epidemiology and Public Health, Imperial School of Medicine, St. Mary's Campus, Norfolk Place, London W2 1PG, UK*

Jun S. Liu, *Statistics Department, Harvard University, Science Center, One Oxford Street, Cambridge, MA 02138-2901*

Joseph Lucas, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Heidi Lyng, *Department of Biophysics, The Norwegian Radium Hospital, Montebello, 0310 Oslo, Norway*

Bani K. Mallick, *Department of Statistics, Texas A&M University, 3143 TAMU, College Station, TX 77843-3143*

Jeffrey S. Morris, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Peter Müller, *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Michael A. Newton, *Department of Statistics, University of Wisconsin-Madison, Medical Sciences Center, RM 1245A, 1300 University Ave., Madison, WI 53706-1532*

Joseph R. Nevins, *Institute for Genome Sciences and Policy, Duke University, Durham, NC 27710*

Michael F. Ochs, *Division of Population Science, Fox Chase Cancer Center, 333 Cottman Ave., Philadelphia, PA 19111-2497*

Sylvia Richardson, *Department of Epidemiology and Public Health, Imperial School of Medicine, St. Mary's Campus, Norfolk Place, London W2 1PG, UK*

Christian Robert, *Ceremade – Université Paris-Dauphine, Bureau C638, Place du Maréchal de Lattre de Tassigny, 75775 Paris Cedex 16, France*

Judith Rousseau, *Ceremade – Université Paris-Dauphine Bureau B638, Place du Maréchal de Lattre de Tassigny, 75775 Paris Cedex 16, France*

Robert Scharpf, *Department of Biostatistics, Bloomberg School of Public Health, Johns Hopkins University, E3034 615 N. Wolfe Street, Baltimore, MD 21205-2179*

Naijun Sha, *Department of Mathematical Science, University of Texas at El Paso, Bell Hall 203, El Paso, TX 79968-0514*

Devika Subramanian, *Department of Computer Science, Rice University, 6100 Main Street, Houston, TX 77005*

Ning Sun, *Division of Biostatistics, Department of Epidemiology and Public Health, Yale University School of Medicine, 60 College Street, P.O. Box 208034, New Haven, CT 06520-8034*

Michael Swartz, *Department of Statistics, Texas A&M University, 3143 TAMU, College Station, TX 77843-3143* and *Department of Biostatistics & Applied Mathematics, University of Texas M.D. Anderson Cancer Center, 1515 Holcombe Blvd., Houston, TX 77030-4075*

Mahlet G. Tadesse, *Department of Epidemiology & Biostatistics, University of Pennsylvania School of Medicine, 918 Blockley Hall, 423 Guardian Drive, Philadelphia, PA 19104-6021*

Mark A. van de Wiel, *Department of Mathematics, Vrije Universiteit, De Boelelaan 1081a, 1081 HV Amsterdam, The Netherlands*

Marina Vannucci, *Department of Statistics, Texas A&M University, 3143 TAMU, College Station, TX 77843-3143*

Jon C. Wakefield, *University of Washington, Department of Biostatistics, Box 357232, Seattle, WA 98195-7232* and *Department of Statistics, Box 354322, Seattle, WA 98195-4322*

Ping Wang, *Department of Statistics, The University of Wisconsin-Madison, B248 Medical Sciences Center, 1300 University Ave., Madison, WI 53703*

Quanli Wang, *Institute for Genome Sciences & Policy, Duke University, Durham, NC 27710* and *Institute of Statistics & Decision Sciences, Duke University, Durham, NC 27708-0251*

Mike West, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Robert L. Wolpert, *Institute of Statistics and Decision Sciences, Box 90251, Duke University, Durham, NC 27708-0251*

Hongyu Zhao, *Division of Biostatistics, Department of Epidemiology and Public Health, Yale University School of Medicine, 60 College Street, P.O. Box 208034, New Haven, CT 06520-8034*

Chuan Zhou, *Department of Biostatistics, Vanderbilt University, S-2323 Medical Center North, Nashville, TN 37232-2158*

Preface

Recent rapid technical advances in genome sequencing (genomics) and protein identification (proteomics) have given rise to research problems that require combined expertise from statistics, biology, computer science, and other fields. The interdisciplinary nature of bioinformatics presents many research challenges related to integrating concepts, methods, software, and multiplatform data. In addition to new tools for investigating biological systems via high-throughput genomic and proteomic measurements, statisticians face many novel methodological research questions generated by such data. The work in this book is dedicated to the development and application of Bayesian statistical methods in the analysis of high-throughput bioinformatics data that arise from problems in medical research, in particular cancer research, and molecular and structural biology. This book does not aim to be comprehensive in all areas of bioinformatics. Rather, it presents a broad overview of statistical inference problems related to three main high-throughput platforms: microarray gene expression, serial analysis gene expression (SAGE), and mass spectrometry proteomic profiles. The book's main focus is on the design, statistical inference, and data analysis, from a Bayesian perspective, of data sets arising from such high-throughput experiments.

Chapter 1 provides a detailed introduction to the three main data platforms and sets the scene for subsequent methodology chapters. This chapter is mainly aimed at nonbiologists and covers elementary biological concepts, details the unique measurement technology with associated idiosyncrasies for the different platforms, and generates an overall outline of issues that statistical methodology can address.

Subsequent chapters focus on specific methodology developments and are grouped approximately by the main bioinformatics platform, with several chapters discussing the integration of at least two platforms. The central statistical topics addressed include experimental design, single group inference, group comparisons, classification and clustering, motif discovery and regulatory networks, and Bayesian networks and gene interactions. The general theme of each

chapter is to review existing methods, followed by a specific novel method developed by the author(s). Results are often demonstrated on simulated data and/or a real application data set. Additionally, relevant software may be discussed.

Chapters 2 through 11 are concerned with Bayesian inference for gene expression, focusing on microarray data. Chapter 2 discusses inference about differential expression based on hierarchical mixture models, including a discussion of more than two patterns of differential expression. Inference is based on conjugate parametric models, with empirical Bayes estimation of hyperparameters. Chapter 3 explores the use of Bayesian hierarchical models for an integrated approach to the analysis of microarray data that includes flexible model-based normalization. The chapter defines a Bayesian gene expression (BGE) index as a gene-specific mean parameter in a hierarchical model. In Chapter 4, a model that mimics the detailed experimental process, from gene preparation to image analysis, is developed. The detailed process-based model allows to estimate absolute and relative mRNA concentrations. The use of Bayesian variable selection methods for biomarker selection in classification and clustering problems is reviewed in Chapter 5. The model for classification includes a mixture prior that specifies a positive probability of a given gene not being included in the model. In the clustering setting, the group structure of the data is uncovered by specifying mixture models where the random inclusion of genes can be interpreted as an attribute selection. Chapter 6 applies multivariate adaptive regression splines (MARS) to define a flexible model for the relationship between gene expression and disease status. For a binary classification problem the MARS model is defined on the logistic transformation of the class probability. The resulting classification boundaries are highly nonlinear and account for gene interactions. Chapter 7 reviews the popular probability of expression (POE) model for differential gene expression. The model postulates a mixture of a uniform submodel for under- and overexpression, and a central normal for typical expression. Posterior probabilities of mixture indicators in this model define the POE scale. Chapter 8 explores the use of sparsity priors in multivariate regression and latent factor regression, applied to inference for gene expression. The sparsity prior is a variation of a mixture prior on regression coefficients, with a positive point mass at zero- and a second-level mixture allowing for gene- and covariate-specific relative weights in the mixture. Chapter 9 develops a model for cell cycle gene expression, using a first-order Fourier model. The set of all genes is partitioned into subsets of different frequency and time-dependent amplitude, including a zero class of not-cycle-dependent genes. The prior on the random partitioning is defined as a Dirichlet distribution for cluster membership indicators. Chapter 10 defines a semiparametric Bayesian model for gene expression. The model exploits the clustering that is implicitly defined by the Dirichlet process prior to define subsets of genes based on gene-specific mean and sampling precision. Chapter 11 reviews the

expression quantitative trait (eQTL) mapping problem and commonly used approaches. A new method is proposed to facilitate eQTL interval mapping that can account for multiplicities across transcripts. The problem is to match microarray gene expression (phenotype) with a set of genetic markers (genetic map). The mixture over markers (MOM) model defines a mixture model for gene expression, with the mixture being defined over submodels corresponding to the transcript mapping to one of the considered markers. An extension to interval mapping is discussed.

Chapters 12 through 15 discuss statistical inference for protein spectrometry. Chapter 12 reviews the use of semiparametric mixture models for inference on differential gene expression, for protein mass/charge spectra, and for SAGE data. The underlying models are Dirichlet process mixtures of normals, mixtures of beta kernels, and Dirichlet process mixtures of Poissons, respectively. Specific focus on SAGE data alone is detailed in Chapter 13, highlighting the two main characteristics of such data: skewness in the distribution of relative abundances, and small sample size relative to the dimension. A new Bayesian procedure based on the mixture Dirichlet prior is reviewed and specific properties depicted in terms of efficiency advantages over existing methods. Chapters 14 and 15 present two different Bayesian approaches of analyzing MALDI-TOF mass spectrometry. Chapter 14 generalizes the linear mixed model to the case of functional data by using a wavelet-based functional mixed model. In contrast, Chapter 15 presents model-based inference by focusing on nonparametric Bayesian models; a sum of kernel functions is chosen as basis functions for modeling spectral peaks.

Chapters 16 through 21 review motif discovery and regulatory networks. During the process of gene transcription, proteins (transcription factors) interact with control points of DNA sequences known as cis-acting regulatory sequences, called motifs. Chapter 16 focuses on the problem of locating these short sequence patterns in the DNA. The chapter provides an extensive explanation of the biological background and then describes a Bayesian framework for gene regulatory binding site discovery. Alternative models are also reviewed. Recent efforts in the development of methods that aid the detection of transcription factor binding sites have attempted to integrate sequence data with gene expression. In Chapter 17 the authors describe a regression model and a Bayesian variable selection formulation that helps refine the search for candidate motifs by selecting those that correlate with the gene expression. The authors also propose a model extension that includes gene regulators. Data integration is also the focus of Chapter 18, where gene expression and protein–DNA binding information, obtained from ChIP-Chip data, are integrated via a linear classification model to aid the reconstruction of gene regulatory networks. Chapter 19 deals with models that exploit the links between transcription factors and signaling pathways via gene ontologies and annotation databases. Bayesian decomposition models allow the extraction of overlapping

transcriptional signatures. Popular approaches to infer gene regulatory networks are those that utilize probabilistic network models. Bayesian networks, in particular, have received much attention. Chapter 20 provides an extensive survey of computational techniques for learning Bayesian networks from gene expression data. State-of-the art and open questions are also addressed. Chapter 21 discusses the developments of Bayesian networks and dynamic networks. Additional data are used to derive informative prior structures, in particular combining gene expression data with protein–DNA binding information.

The final chapter, Chapter 22, addresses the choice of the sample size for microarray experiments. The authors take a decision theory point of view that attempts to minimize a conditional expected loss. The method is exemplified using a mixture Gamma/Gamma model.

We thank our friends and collaborators for contributing their ideas and insights toward this collection. We are excited by the continuing opportunities for statistical challenges in the area of high-throughput bioinformatics data. We hope our readers will join us in being engaged with changing technologies and statistical development.

Kim-Anh Do
Peter Müller
Marina Vannucci