

Statistics for Biology and Health

Series Editors

W. Wong, M. Gail, K. Krickeberg, A. Tsiatis, J. Samet

Robert Gentleman
Vincent J. Carey
Wolfgang Huber

Rafael A. Irizarry
Sandrine Dudoit

Editors

Bioinformatics and Computational Biology Solutions Using R and Bioconductor

With 128 Illustrations

 Springer

Editors

Robert Gentleman
Program in Computational Biology
Division of Public Health Sciences
Fred Hutchinson Cancer Research Center
1100 Fairview Ave. N, M2-B876
PO Box 19024
Seattle, Washington 98109-1024 USA
rgentlem@fhcrc.org

Vincent J. Carey
Channing Laboratory
Brigham and Women's Hospital
Harvard Medical School
181 Longwood Ave Boston MA 02115 USA
stvjc@channing.harvard.edu

Wolfgang Huber
European Bioinformatics Institute
European Molecular Biology
Laboratory
Cambridge, CB10 1SD UK
huber@ebi.ac.uk

Rafael A. Irizarry
Department of Biostatistics
Johns Hopkins Bloomberg
School of Public Health
615 North Wolfe Street
Baltimore, MD 21205 USA
rafa@jhu.edu

Sandrine Dudoit
Division of Biostatistics
School of Public Health
University of California,
Berkeley
140 Earl Warren Hall, #7360
Berkeley, CA 94720-7360
USA
sandrine@stat.berkeley.edu

Series Editors

Wing Wong
Department of Statistics
Stanford University
Stanford, CA 94305
USA

M. Gail
National Cancer Institute
Rockville, MD 20892
USA

K. Krickeberg
Le Châtelet
F-63270 Manglieu
France

A. Tsiatis
Department of Statistics
North Carolina State University
Raleigh, NC 27695
USA

J. Samet
Department of Epidemiology
School of Public Health
Johns Hopkins University
615 Wolfe Street
Baltimore, MD 21205
USA

Library of Congress Control Number: 2005923843

ISBN-10: 0-387-25146-4

Printed on acid-free paper.

ISBN-13: 978-0387-25146-2

© 2005 Springer Science+Business Media, Inc.

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, Inc., 233 Spring Street, New York, NY 10013, USA), 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 in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed in China. (EVB)

9 8 7 6 5 4 3 2 1

springeronline.com

Preface

During the past few years, there have been enormous advances in genomics and molecular biology, which carry the promise of understanding the functioning of whole genomes in a systematic manner. The challenge of interpreting the vast amounts of data from microarrays and other high throughput technologies has led to the development of new tools in the fields of computational biology and bioinformatics, and opened exciting new connections to areas such as chemometrics, exploratory data analysis, statistics, machine learning, and graph theory.

The Bioconductor project is an open source and open development software project for the analysis and comprehension of genomic data. It is rooted in the open source statistical computing environment R. This book's coverage is broad and ranges across most of the key capabilities of the Bioconductor project. Thanks to the hard work and dedication of many developers, a responsive and enthusiastic user community has formed. Although this book is self-contained with respect to the data processing and data analytic tasks covered, readers of this book are advised to acquaint themselves with other aspects of the project by touring the project web site www.bioconductor.org.

This book represents an innovative approach to publishing about scientific software. We made a commitment at the outset to have a fully *computable book*. Tables, figures, and other outputs are dynamically generated directly from the experimental data. Through the companion web site, www.bioconductor.org/mogr, readers have full access to the source code and necessary supporting libraries and hence will be able to see how every plot and statistic was computed. They will be able to reproduce those calculations on their own computers and should be able to extend most of those computations to address their own needs.

Acknowledgments

This book, like so many projects in bioinformatics and computational biology, is a large collaborative effort. The editors would like to thank the chapter authors for their dedication and their efforts in producing widely used software, and also in producing well-written descriptions of how to use that software.

We would like to thank the developers of R, without whom there would be no Bioconductor project. Many of these developers have provided additional help and engaged in discussions about software development and design. We would like to thank the many Bioconductor developers and users who have helped us to find bugs, think differently about problems, and whose enthusiasm has made the long hours somewhat more bearable.

We would also like to thank Dorit Arlt, Michael Boutros, Sabina Chiaretti, James MacDonald, Meher Majety, Annemarie Poustka, Jerome

Ritz, Mamatha Sauermann, Holger Sülthmann, Stefan Wiemann, and Seth Falcon, who have contributed in many different ways to the production of this monograph. Much of the preliminary work on the `MLInterfaces` package, described in Chapter 16, was carried out by Jess Mar, Department of Biostatistics, Harvard School of Public Health. Ms Mar's efforts were supported in part by a grant from Insightful Corporation.

The Bioconductor project is supported by grant 1R33 HG002708 from the NIH as well as by institutional funds at both the Dana Farber Cancer Institute and the Fred Hutchinson Cancer Research Center. W.H. received project-related funding from the German Ministry for Education and Research through National Genome Research Network (NGFN) grant FKZ 01GR0450.

Seattle
Boston
Cambridge (UK)
Baltimore
Berkeley

Robert Gentleman
Vincent Carey
Wolfgang Huber
Rafael Irizarry
Sandrine Dudoit
February 2005

J. Gentry, Center for Cancer Research, Massachusetts General Hospital, Boston, MA, USA

F. Hahne, Division of Molecular Genome Analysis, German Cancer Research Center, Heidelberg, FRG

L. Harris, Department of Cancer Biology, Dana Farber Cancer Institute, Boston, MA, USA

T. Hothorn, Institut für Medizininformatik, Biometrie und Epidemiologie, Friedrich-Alexander-Universität Erlangen-Nürnberg, FRG

W. Huber, European Molecular Biology Laboratory, European Bioinformatics Institute, Cambridge, UK

J. Ibrahim, Department of Biostatistics, University of North Carolina, Chapel Hill, NC, USA

J. D. Iglehart, Department of Cancer Biology, Dana Farber Cancer Institute, Boston, MA, USA

R. A. Irizarry, Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

X. Li, Department of Biostatistics and Computational Biology, Dana Farber Cancer Institute, Boston, MA, USA

X. Lu, Department of Biostatistics, Harvard School of Public Health, Boston, MA, USA

A. Miron, Department of Cancer Biology, Dana Farber Cancer Institute, Boston, MA, USA

A. C. Paquet, Department of Biostatistics, University of California, San Francisco, CA, USA

K. S. Pollard, Center for Biomolecular Science and Engineering, University of California, Santa Cruz, USA

D. Scholtens, Department of Preventive Medicine, Northwestern University, Chicago, IL, USA

Q. Shi, Department of Cancer Biology, Dana Farber Cancer Institute, Boston, MA, USA

Contents

I	Preprocessing data from genomic experiments	1
1	Preprocessing Overview	3
	W. Huber, R. A. Irizarry, and R. Gentleman	
1.1	Introduction	3
1.2	Tasks	4
	1.2.1 Prerequisites	5
	1.2.2 Stepwise and integrated approaches	5
1.3	Data structures	6
	1.3.1 Data sources	6
	1.3.2 Facilities in R and Bioconductor	7
1.4	Statistical background	8
	1.4.1 An error model	9
	1.4.2 The variance-bias trade-off	11
	1.4.3 Sensitivity and specificity of probes	11
1.5	Conclusion	12
2	Preprocessing High-density Oligonucleotide Arrays	13
	B. M. Bolstad, R. A. Irizarry, L. Gautier, and Z. Wu	
2.1	Introduction	13
2.2	Importing and accessing probe-level data	15
	2.2.1 Importing	15
	2.2.2 Examining probe-level data	15
2.3	Background adjustment and normalization	18
	2.3.1 Background adjustment	18
	2.3.2 Normalization	20
	2.3.3 vsn	24
2.4	Summarization	25
	2.4.1 <code>expresso</code>	25
	2.4.2 <code>threestep</code>	26
	2.4.3 <code>RMA</code>	27
	2.4.4 <code>GCRMA</code>	27
	2.4.5 <code>affypdnn</code>	28

2.5	Assessing preprocessing methods	29
2.5.1	Carrying out the assessment	30
2.6	Conclusion	32
3	Quality Assessment of Affymetrix GeneChip Data	33
	B. M. Bolstad, F. Collin, J. Brettschneider, K. Simpson, L. Cope, R. A. Irizarry, and T. P. Speed	
3.1	Introduction	33
3.2	Exploratory data analysis	34
3.2.1	Multi-array approaches	35
3.3	Affymetrix quality assessment metrics	37
3.4	RNA degradation	38
3.5	Probe level models	41
3.5.1	Quality diagnostics using PLM	42
3.6	Conclusion	47
4	Preprocessing Two-Color Spotted Arrays	49
	Y. H. Yang and A. C. Paquet	
4.1	Introduction	49
4.2	Two-color spotted microarrays	50
4.2.1	Illustrative data	50
4.3	Importing and accessing probe-level data	51
4.3.1	Importing	51
4.3.2	Reading target information	52
4.3.3	Reading probe-related information	53
4.3.4	Reading probe and background intensities	54
4.3.5	Data structure: the <i>marrayRaw</i> class	54
4.3.6	Accessing the data	56
4.3.7	Subsetting	56
4.4	Quality assessment	57
4.4.1	Diagnostic plots	57
4.4.2	Spatial plots of spot statistics - <code>image</code>	59
4.4.3	Boxplots of spot statistics - <code>boxplot</code>	60
4.4.4	Scatter-plots of spot statistics - <code>plot</code>	61
4.5	Normalization	62
4.5.1	Two-channel normalization	63
4.5.2	Separate-channel normalization	64
4.6	Case study	67
5	Cell-Based Assays	71
	W. Huber and F. Hahne	
5.1	Scope	71
5.2	Experimental technologies	71
5.2.1	Expression assays	72
5.2.2	Loss of function assays	72

5.2.3	Monitoring the response	72
5.3	Reading data	73
5.3.1	Plate reader data	74
5.3.2	Further directions in normalization	76
5.3.3	FCS format	77
5.4	Quality assessment and visualization	79
5.4.1	Visualization at the level of individual cells	79
5.4.2	Visualization at the level of microtiter plates	82
5.4.3	Brushing with Rggobi	83
5.5	Detection of effectors	85
5.5.1	Discrete Response	85
5.5.2	Continuous response	88
5.5.3	Outlook	90

6 SELDI-TOF Mass Spectrometry Protein Data 91

X. Li, R. Gentleman, X. Lu, Q. Shi, J. D. Iglehart, L. Harris, and A. Miron

6.1	Introduction	91
6.2	Baseline subtraction	93
6.3	Peak detection	95
6.4	Processing a set of calibration spectra	96
6.4.1	Apply baseline subtraction to a set of spectra	98
6.4.2	Normalize spectra	99
6.4.3	Cutoff selection	100
6.4.4	Identify peaks	101
6.4.5	Quality assessment	101
6.4.6	Get proto-biomarkers	102
6.5	An example	105
6.6	Conclusion	108

II Meta-data: biological annotation and visualization 111

7 Meta-data Resources and Tools in Bioconductor 113

R. Gentleman, V. J. Carey, and J. Zhang

7.1	Introduction	113
7.2	External annotation resources	115
7.3	Bioconductor annotation concepts: curated persistent packages and Web services	116
7.3.1	Annotating a platform: HG-U95Av2	117
7.3.2	An Example	118
7.3.3	Annotating a genome	119
7.4	The <code>annotate</code> package	119
7.5	Software tools for working with Gene Ontology (GO)	120

7.5.1	Basics of working with the GO package	121
7.5.2	Navigating the hierarchy	122
7.5.3	Searching for terms	122
7.5.4	Annotation of GO terms to LocusLink sequences: evidence codes	123
7.5.5	The GO graph associated with a term	125
7.6	Pathway annotation packages: KEGG and cMAP	125
7.6.1	KEGG	126
7.6.2	cMAP	127
7.6.3	A Case Study	129
7.7	Cross-organism annotation: the homology packages	130
7.8	Annotation from other sources	132
7.9	Discussion	133
8	Querying On-line Resources	135
	V. J. Carey, D. Temple Lang, J. Gentry, J. Zhang, and R. Gentleman	
8.1	The Tools	135
8.1.1	Entrez	137
8.1.2	Entrez examples	137
8.2	PubMed	138
8.2.1	Accessing PubMed information	139
8.2.2	Generating HTML output for your abstracts	141
8.3	KEGG via SOAP	142
8.4	Getting gene sequence information	144
8.5	Conclusion	145
9	Interactive Outputs	147
	C. A. Smith, W. Huber, and R. Gentleman	
9.1	Introduction	147
9.2	A simple approach	148
9.3	Using the <code>annaffy</code> package	149
9.4	Linking to On-line Databases	152
9.5	Building HTML pages	153
9.5.1	Limiting the results	153
9.5.2	Annotating the probes	154
9.5.3	Adding other data	155
9.6	Graphical displays with drill-down functionality	156
9.6.1	HTML image maps	157
9.6.2	Scalable Vector Graphics (SVG)	158
9.7	Searching Meta-data	159
9.7.1	Text searching	159
9.8	Concluding Remarks	160
10	Visualizing Data	161

W. Huber, X. Li, and R. Gentleman

- 10.1 Introduction 161
- 10.2 Practicalities 162
- 10.3 High-volume scatterplots 163
 - 10.3.1 A note on performance 164
- 10.4 Heatmaps 166
 - 10.4.1 Heatmaps of residuals 168
- 10.5 Visualizing distances 170
 - 10.5.1 Multidimensional scaling 173
- 10.6 Plotting along genomic coordinates 174
 - 10.6.1 Cumulative Expression 178
- 10.7 Conclusion 179

III Statistical analysis for genomic experiments 181

11 Analysis Overview 183

V. J. Carey and R. Gentleman

- 11.1 Introduction and road map 183
 - 11.1.1 Distance concepts 184
 - 11.1.2 Differential expression 184
 - 11.1.3 Cluster analysis 184
 - 11.1.4 Machine learning 184
 - 11.1.5 Multiple comparisons 185
 - 11.1.6 Workflow support 185
- 11.2 Absolute and relative expression measures 185

12 Distance Measures in DNA Microarray Data Analysis. 189

R. Gentleman, B. Ding, S. Dudoit, and J. Ibrahim

- 12.1 Introduction 189
- 12.2 Distances 191
 - 12.2.1 Definitions 191
 - 12.2.2 Distances between points 192
 - 12.2.3 Distances between distributions 195
 - 12.2.4 Experiment-specific distances between genes 198
- 12.3 Microarray data 199
 - 12.3.1 Distances and standardization 199
- 12.4 Examples 201
 - 12.4.1 A co-citation example 203
 - 12.4.2 Adjacency 207
- 12.5 Discussion 208

13 Cluster Analysis of Genomic Data 209

K. S. Pollard and M. J. van der Laan

- 13.1 Introduction 209

13.2	Methods	210
13.2.1	Overview of clustering algorithms	210
13.2.2	Ingredients of a clustering algorithm	211
13.2.3	Building sequences of clustering results	211
13.2.4	Visualizing clustering results	214
13.2.5	Statistical issues in clustering	215
13.2.6	Bootstrapping a cluster analysis	216
13.2.7	Number of clusters	217
13.3	Application: renal cell cancer	222
13.3.1	Gene selection	222
13.3.2	HOPACH clustering of genes	223
13.3.3	Comparison with PAM	224
13.3.4	Bootstrap resampling	224
13.3.5	HOPACH clustering of arrays	224
13.3.6	Output files	226
13.4	Conclusion	228
14	Analysis of Differential Gene Expression Studies	229
	D. Scholtens and A. von Heydebreck	
14.1	Introduction	229
14.2	Differential expression analysis	230
14.2.1	Example: ALL data	232
14.2.2	Example: Kidney cancer data	236
14.3	Multifactor experiments	239
14.3.1	Example: Estrogen data	241
14.4	Conclusion	248
15	Multiple Testing Procedures: the multtest Package and Applications to Genomics	249
	K. S. Pollard, S. Dudoit, and M. J. van der Laan	
15.1	Introduction	249
15.2	Multiple hypothesis testing methodology	250
15.2.1	Multiple hypothesis testing framework	250
15.2.2	Test statistics null distribution	255
15.2.3	Single-step procedures for controlling general Type I error rates $\theta(F_{V_n})$	256
15.2.4	Step-down procedures for controlling the family-wise error rate	257
15.2.5	Augmentation multiple testing procedures for controlling tail probability error rates	258
15.3	Software implementation: R multtest package	259
15.3.1	Resampling-based multiple testing procedures: MTP function	260
15.3.2	Numerical and graphical summaries	262
15.4	Applications: ALL microarray data set	262

15.4.1 ALL data package and initial gene filtering 262

15.4.2 Association of expression measures and tumor cellular subtype: Two-sample t -statistics 263

15.4.3 Augmentation procedures 265

15.4.4 Association of expression measures and tumor molecular subtype: Multi-sample F -statistics . . . 266

15.4.5 Association of expression measures and time to relapse: Cox t -statistics 268

15.5 Discussion 270

16 Machine Learning Concepts and Tools for Statistical Genomics 273

V. J. Carey

16.1 Introduction 273

16.2 Illustration: Two continuous features; decision regions . . 274

16.3 Methodological issues 276

16.3.1 Families of learning methods 276

16.3.2 Model assessment 281

16.3.3 Metatheorems on learner and feature selection . . 283

16.3.4 Computing interfaces 284

16.4 Applications 285

16.4.1 Exploring and comparing classifiers with the ALL data 285

16.4.2 Neural net initialization, convergence, and tuning . 287

16.4.3 Other methods 287

16.4.4 Structured cross-validation support 288

16.4.5 Assessing variable importance 289

16.4.6 Expression density diagnostics 289

16.5 Conclusions 291

17 Ensemble Methods of Computational Inference 293

T. Hothorn, M. Dettling, and P. Bühlmann

17.1 Introduction 293

17.2 Bagging and random forests 295

17.3 Boosting 296

17.4 Multiclass problems 298

17.5 Evaluation 298

17.6 Applications: tumor prediction 300

17.6.1 Acute lymphoblastic leukemia 300

17.6.2 Renal cell cancer 303

17.7 Applications: Survival analysis 307

17.8 Conclusion 310

18 Browser-based Affymetrix Analysis and Annotation 313

C. A. Smith

18.1	Introduction	313
	18.1.1 Key user interface features	314
18.2	Deploying webbioc	315
	18.2.1 System requirements	315
	18.2.2 Installation	315
	18.2.3 Configuration	316
18.3	Using webbioc	317
	18.3.1 Data Preprocessing	317
	18.3.2 Differential expression multiple testing	318
	18.3.3 Linked annotation meta-data	320
	18.3.4 Retrieving results	321
18.4	Extending webbioc	322
	18.4.1 Architectural overview	322
	18.4.2 Creating a new module	324
18.5	Conclusion	326

IV Graphs and networks 327

19 Introduction and Motivating Examples 329

R. Gentleman, W. Huber, and V. J. Carey

19.1	Introduction	329
19.2	Practicalities	330
	19.2.1 Representation	330
	19.2.2 Algorithms	330
	19.2.3 Data Analysis	331
19.3	Motivating examples	331
	19.3.1 Biomolecular Pathways	331
	19.3.2 Gene ontology: A graph of concept-terms	333
	19.3.3 Graphs induced by literature references and citations	334
19.4	Discussion	336

20 Graphs 337

W. Huber, R. Gentleman, and V. J. Carey

20.1	Overview	337
20.2	Definitions	338
	20.2.1 Special types of graphs	341
	20.2.2 Random graphs	343
	20.2.3 Node and edge labeling	344
	20.2.4 Searching and related algorithms	344
20.3	Cohesive subgroups	344
20.4	Distances	346

21 Bioconductor Software for Graphs	347
V. J. Carey, R. Gentleman, W. Huber, and J. Gentry	
21.1 Introduction	347
21.2 The graph package	348
21.2.1 Getting started	349
21.2.2 Random graphs	352
21.3 The RBGL package	352
21.3.1 Connected graphs	355
21.3.2 Paths and related concepts	357
21.3.3 RBGL summary	360
21.4 Drawing graphs	360
21.4.1 Global attributes	363
21.4.2 Node and edge attributes	363
21.4.3 The function <code>agopen</code> and the <code>Ragraph</code> class	365
21.4.4 User-defined drawing functions	366
21.4.5 Image maps on graphs	368
22 Case Studies Using Graphs on Biological Data	369
R. Gentleman, D. Scholtens, B. Ding, V. J. Carey, and W. Huber	
22.1 Introduction	369
22.2 Comparing the transcriptome and the interactome	370
22.2.1 Testing associations	371
22.2.2 Data analysis	373
22.3 Using GO	374
22.3.1 Finding interesting GO terms	375
22.4 Literature co-citation	378
22.4.1 Statistical development	380
22.4.2 Comparisons of interest	382
22.4.3 Examples	382
22.5 Pathways	387
22.5.1 The graph structure of pathways	388
22.5.2 Relating expression data to pathways	390
22.6 Concluding remarks	393
V Case studies	395
23 limma: Linear Models for Microarray Data	397
G. K. Smyth	
23.1 Introduction	397
23.2 Data representations	398
23.3 Linear models	399
23.4 Simple comparisons	400
23.5 Technical Replication	403
23.6 Within-array replicate spots	406

23.7	Two groups	407
23.8	Several groups	409
23.9	Direct two-color designs	411
23.10	Factorial designs	412
23.11	Time course experiments	414
23.12	Statistics for differential expression	415
23.13	Fitted model objects	417
23.14	Preprocessing considerations	418
23.15	Conclusion	420
24	Classification with Gene Expression Data	421
	M. Dettling	
24.1	Introduction	421
24.2	Reading and customizing the data	422
24.3	Training and validating classifiers	423
24.4	Multiple random divisions	426
24.5	Classification of test data	428
24.6	Conclusion	429
25	From CEL Files to Annotated Lists of Interesting Genes	431
	R. A. Irizarry	
25.1	Introduction	431
25.2	Reading CEL files	432
25.3	Preprocessing	432
25.4	Ranking and filtering genes	433
	25.4.1 Summary statistics and tests for ranking	434
	25.4.2 Selecting cutoffs	437
	25.4.3 Comparison	437
25.5	Annotation	438
	25.5.1 PubMed abstracts	439
	25.5.2 Generating reports	441
25.6	Conclusion	442
A	Details on selected resources	443
A.1	Data sets	443
	A.1.1 ALL	443
	A.1.2 Renal cell cancer	443
	A.1.3 Estrogen receptor stimulation	443
A.2	URLs for projects mentioned	444
	References	445
	Index	465

List of Contributors

B. M. Bolstad, Department of Statistics, University of California, Berkeley, CA, USA

J. Brettschneider, Department of Statistics, University of California, Berkeley, CA, USA

P. Buhlmann, Swiss Federal Institute of Technology, Zürich, CH

V. J. Carey, Channing Laboratory, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

F. Collin, Department of Statistics, University of California, Berkeley, CA, USA

L. Cope, Division of Oncology Biostatistics, The Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins Medical School, Baltimore, MD, USA

M. Dettling, Division of Oncology and Biostatistics, The Sidney Kimmel Comprehensive Cancer Center, Johns Hopkins Medical School, Baltimore, MD, USA

B. Ding, Medical Affairs Biostatistics, Amgen Inc., Thousand Oaks, CA, USA

S. Dudoit, Department of Biostatistics, University of California, Berkeley, CA, USA

L. Gautier, Independent investigator, Copenhagen, DK

R. Gentleman, Program in Computational Biology, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

K. Simpson, The Walter and Eliza Hall Institute of Medical Research, Melbourne, Australia

C. A. Smith, Department of Molecular Biology, The Scripps Research Institute, La Jolla, CA, USA

G. K. Smyth, The Walter and Eliza Hall Institute of Medical Research, Melbourne, Australia

T. P. Speed, Department of Statistics, University of California, Berkeley, CA, USA

D. Temple Lang, Department of Statistics, University of California, Davis, CA, USA

M. J. van der Laan, Department of Biostatistics, University of California, Berkeley, CA, USA

A. von Heydebreck, Global Technologies, Merck KGaA, Darmstadt, FRG

Z. Wu, Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

Y. H. Yang, Department of Biostatistics, University of California, San Francisco, CA, USA

J. Zhang, Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

Part I

Preprocessing data from genomic experiments

1

Preprocessing Overview

W. Huber, R. A. Irizarry, and R. Gentleman

Abstract

In this chapter, we give a brief overview of the tasks of microarray data preprocessing. There are a variety of microarray technology platforms in use, and each of them requires specific considerations. These will be described in detail by other chapters in this part of the book. This overview chapter describes relevant data structures, and provides with some broadly applicable theoretical background.

1.1 Introduction

Microarray technology takes advantage of hybridization properties of nucleic acid and uses complementary molecules attached to a solid surface, referred to as *probes*, to measure the quantity of specific nucleic acid transcripts of interest that are present in a sample, referred to as the *target*. The molecules in the target are labeled, and a specialized scanner is used to measure the amount of hybridized target at each probe, which is reported as an intensity. Various manufacturers provide a large assortment of different platforms. Most manufacturers, realizing the effects of optical noise and non-specific binding, include features in their arrays to directly measure these effects. The raw or *probe-level* data are the intensities read for each of these components. In practice, various sources of variation need to be accounted for, and these data are heavily manipulated before one obtains the genomic-level measurements that most biologists and clinicians use in their research. This procedure is commonly referred to as *preprocessing*.

The different platforms can be divided into two main classes that are differentiated by the type of data they produce. The *high-density oligonucleotide array* platforms produce one set of probe-level data per microarray with some probes designed to measure specific binding and others to measure non-specific binding. The two-color spotted platforms produce two sets

of probe-level data per microarray (the red and green channels), and local background noise levels are measured from areas in the glass slide not containing probe.

Despite the differences among the different platforms, there are some tasks that are common to all microarray technology. These tasks are described in Section 1.2. The data structures needed to effectively preprocess microarray data are described in Section 1.3. In Section 1.4 we present statistical background that serves as a mathematical framework for developing preprocessing methodology. Detailed description of the preprocessing tasks for this platforms are described in Chapters 2 and 3. The specifics for the two-color spotted platforms are described in Chapter 4. Chapters 5 and 6 describe preprocessing methodology for related technologies where similar principles apply.

1.2 Tasks

Preprocessing can be divided into 6 tasks: image analysis, data import, background adjustment, normalization, summarization, and quality assessment. *Image analysis* permits us to convert the pixel intensities in the scanned images into probe-level data. Flexible data import methods are needed because data come in different formats and are often scattered across a number of files or database tables from which they need to be extracted and organized. *Background adjustment* is essential because part of the measured probe intensities are due to non-specific hybridization and the noise in the optical detection system. Observed intensities need to be adjusted to give accurate measurements of specific hybridization. Without proper *normalization*, it is impossible to compare measurements from different array hybridizations due to many obscuring sources of variation. These include different efficiencies of reverse transcription, labeling, or hybridization reactions, physical problems with the arrays, reagent batch effects, and laboratory conditions. In some platforms, *summarization* is needed because transcripts are represented by multiple probes. For each gene, the background adjusted and normalized intensities need to be summarized into one quantity that estimates an amount proportional to the amount of RNA transcript. *Quality assessment* is an important procedure that detects divergent measurements beyond the acceptable level of random fluctuations. These data are usually flagged and not used, or down weighted, in subsequent statistical analyses.

The complex nature of microarray data and data formats makes it necessary to have flexible and efficient statistical methodology and software. This part of the book describes what Bioconductor has to offer in this capacity. In the rest of this section, we describe prerequisites necessary to perform these tasks and two general approaches to preprocessing.

1.2.1 Prerequisites

A number of important steps are involved in the generation of the raw data. The experimental design includes the choice and collection of samples (tissue biopsies or cell lines exposed to different treatments); the choice of probes and array platform; the choice of controls, RNA extraction, amplification, labeling, and hybridization procedures; the allocation of replicates; and the scheduling of the experiments. The experimental design must take into account technical, logistic, and financial boundary conditions. Its quality determines to a large extent the utility of the data. A fundamental guideline is the avoidance of *confounding* between different biological factors of interest or between a biological factor of interest and a technical factor that is anticipated to affect the measurements. The experiment then has to be carried out, which requires great skill and expertise.

In the *image analysis* step, we extract probe intensities out of the scanned images containing pixel-level data. The arrays are scanned by the detector at a high spatial resolution to produce a digitized image in which each probe is represented by dozens of pixels. To obtain a single overall intensity value for each probe, the associated pixels need to be identified (segmentation) and their intensities summarized (quantification). In addition to the overall probe intensity, further auxiliary quantities may be calculated, such as an estimate of apparent unspecific “local background” intensity, or spot quality measures. Various software packages offer different segmentation and quantification methods. They differ in their robustness against irregularities and in the amount of human interaction that they require. The different platforms present different problems which implies that the types of image analysis algorithms used are quite different. Currently, Bioconductor does not offer image processing software. Thus, the user will need alternative software to process the image pixel-level data. However, import functions that are compatible with most of the existing image analysis products are available. For an evaluation of image analysis methods for two-color spotted arrays see, for example, the study of Yang et al. (2002a). Details on image analysis methodology for high-density oligonucleotide arrays were described by Schadt et al. (2001).

1.2.2 Stepwise and integrated approaches

The *stepwise* approach to microarray data preprocessing starts with probe-level data as input, performs the tasks sequentially and produces an *expression matrix* as output. In this matrix, rows correspond to gene transcripts, columns to conditions, and each element represents the abundance or relative abundance of a transcript. Subsequent biological analyses work off the expression matrix and generally do not consider the statistical manipulations performed on the probe-level data. The preprocessing task are divided into a set of sequential instructions: for example, subtract

the background, then normalize the intensities, then summarize replicate probes, then summarize replicate arrays. The modularity of this approach allows us to structure the analysis work-flow. Software, data structures, and methodology can be easily re-used. For example, the same machine learning algorithm can be applied to an expression matrix irrespective of whether the data were obtained on high-density oligonucleotide chips or two-color spotted arrays. A potential disadvantage of the stepwise approach is that each step is independently optimized without considering the effect of previous or subsequent steps. This could lead to sub-optimal bottom-line results.

In contrast, *integrated* approaches solve specific problems by carrying out the analysis in one unified estimation procedure. This approach has the potential of using the available data more efficiently. For example, rather than calculating an expression matrix, one might fit an ANOVA-type linear model to the probe-level data, which includes both technical covariates, such as dye and sample effects, and biological covariates, such as treatment effects (Kerr et al., 2000). In the *affyPLM* package, the weighting and summarization of the multiple probes per transcript on Affymetrix chips is integrated with the detection of differential expression. Another example is the *vsn* method (Huber et al., 2002), which integrates background subtraction and normalization in a non-linear model.

Stepwise approaches are often presented as modular data processing pipelines; integrated approaches are motivated by statistical models with parameters representing quantities of interest. In practice, data analysts will often choose to use a combination of both approaches. For example, a researcher may start with the stepwise approach and do a first round of high-level analyses that motivates an integrated approach that is applied to obtain final results. Bioconductor software allows users to explore, adapt, and combine stepwise and integrated methods.

1.3 Data structures

1.3.1 Data sources

The basic data types that we deal with in microarray data preprocessing are probe and background intensities, probe annotations, array layout, and sample annotations. Typically, they come in the form of rectangular tables, stored either in flat files or in a database server. The probe intensities are the result of image processing. The format in which they are reported varies between different vendors of image processing software. Examples are discussed in Sections 2 and 4.

The probe annotations are usually provided by the organization that selected the probes for the array. This may be a commercial vendor, another laboratory, or the experimenters themselves. For high-density oligonucleotide arrays, the primary annotation is the sequence. In addition, there

may be a database identifier of the gene transcript that the probe is intended to match and possibly the exact location. Often, the probe sequences are derived from cDNA sequence clusterings such as Unigene (Pontius et al., 2003). For spotted cDNA arrays, the primary probe identifier is often a clone ID in a nucleotide sequence database. The largest public nucleotide sequence databases are EMBL in Europe, DDBJ in Japan, and Genbank in the United States. Through a system of cross-mirroring, their contents are essentially equivalent. These databases contain full or partial sequences of a large number of expressed sequences. Their clone identifiers can be mapped to genomic databases such as Entrez Gene, H-inv, or Ensembl. Further annotations of the genes that are represented by the probes are provided by various genomic database, for example genomic locus, disease associations, participation in biological processes, molecular function, cellular localization. This will be discussed in Part II of the book.

The array layout is provided by the organization that produced the array. As a minimum, the layout specifies the physical position of each probe on the array. In principle, this can be done through its x - and y -coordinates. For spotted arrays, it is customary to specify probe coordinates through three coordinates: *block*, *row*, and *column*, where the *block* coordinate addresses a particular sub-sector of the array, and the *row* and *column* coordinates address the probe within that sub-sector. Details are discussed in Sections 2 and 4.

The sample annotations describe the labeled cDNA that has been hybridized to the array. This includes technical information on processing protocols (e.g., isolation, amplification, labeling, hybridization) as well as the biologically more interesting covariates such as treatment conditions and harvesting time points for cell lines or histopathological and clinical data for tissue biopsies and the individuals that the biopsies originated from. A table containing this information can sometimes be obtained from the laboratory information management system (LIMS) of the lab that performed the experiments. Sometimes, it is produced *ad hoc* with office spreadsheet software.

1.3.2 Facilities in R and Bioconductor

Specific data structures and functions for the import and processing of data from different experimental platforms are provided in specialized packages. We will see a number of examples in the subsequent sections. A more general-purpose data structure to represent the data from a microarray experiment is provided by the class *exprSet* in the package Biobase.

The design of the *exprSet* class supports the stepwise approach to microarray preprocessing, as discussed in Section 1.2. This class represents a self-documenting data structure, with data separated into logically distinct but substantively interdependent components. Our primary motivation was to link together the large expression arrays with the phenotypic data in such

a way that it would be easy to further process the data. Ensuring correct alignment of data when subsets are taken or when resampling schemes are used should be left to well-designed computer code and generally should not be done by hand.

The general premise is that there is an array, or a set of arrays, that are of interest. The *exprSet* structure imposes an order on the sample-specific expression measures in the set, provides convenient access to probe and sample identifier codes, allows coordinated management of standard errors of expression, and couples to this expression information sample- and experiment-level information, following the MIAME standard (Brazma et al., 2001). This data structure is straightforwardly employed with data from single-channel experiments, for ratio quantities derived from double-channel experiments, and for protein mass-spectrometry data. It can be extended, using formal inheritance infrastructure, to accommodate other output formats. One advantage to the use of `exprSets` is demonstrated in Chapter 16 where we describe the use of a uniform calling sequence for many machine learning algorithms (package `MLInterfaces`). This greatly simplifies individual users' interactions and will simplify the design and construction of graphical user interfaces. Establishment of a standardized calling paradigm is most simply accomplished when there are structural standards for the inputs. Both users and developers will profit from closer acquaintance with the `exprSet` structure, especially those who are contemplating complex downstream workflows.

1.4 Statistical background

The purpose of this section is to provide a general statistical framework for the following components of preprocessing: background adjustment, normalization, summarization, and quality assessment. More specific issues relating to the individual technological platforms will be discussed in Chapters 2–4.

With a microarray experiment, we aim to make statements about the abundances of specific molecules in a set of biological samples. However, the quantities that we measure are the fluorescence intensities of the different elements of the array. The measurement process consists of a cascade of biochemical reactions and an optical detection system with a laser scanner or a CCD camera. Biochemical reactions and detection are performed in parallel, allowing up to a million measurements on one array. Subtle variations between arrays, the reagents used, and the environmental conditions lead to slightly different measurements even for the same sample.

The effects of these variations may be grouped in two classes: *systematic effects*, which affect a large number of measurements (for example, the measurements for all probes on one array; or the measurements from one probe

across several arrays) simultaneously. Such effects can be estimated and approximately removed. Other kinds of effects are completely random, with no well-understood pattern. These effects are commonly called *stochastic components* or *noise*.

Stochastic models are useful for preprocessing because they permit us to find *optimal* estimates of the systematic effects. We are interested in estimates that are precise and accurate. However, given the noise structure of the data, we sometimes have to sacrifice accuracy for better precision and *vice versa*. An appropriate stochastic model will aid in understanding the accuracy-precision, or bias-variance, trade-off.

Stochastic models are also useful for construction of inferential statements about experimental results. Consider an experiment in which we want to compare gene expression in the colons of mice that were treated with a substance and mice that were not. If we have many measurements from two populations being compared, we can, for example, perform a Wilcoxon test to obtain a p -value for each transcript of interest. But often it is not possible, too expensive, or unethical, to obtain so many replicate measurements for all genes and for all conditions of interest. Often, it is also not necessary. Models that provide good approximations of reality can add power to our statistical results.

Quality assessment is yet another example of the usefulness of stochastic models: if the distribution of a new set of data greatly deviates from the model, this may direct our attention to quality issues with these data. Chapter 3 demonstrates an example of the use of models for quality assessment.

1.4.1 An error model

A generic model for the value of the intensity y of a single probe on a microarray is given by

$$Y = B + \alpha S \quad (1.1)$$

where B is a random quantity due to *background noise*, usually composed of optical effects and non-specific binding, α is a gain factor, and S is the amount of measured specific binding. The signal S is considered a random variable as well and accounts for measurement error and probe effects. The measurement error is typically assumed to be multiplicative so we write:

$$\log(S) = \theta + \phi + \varepsilon. \quad (1.2)$$

Here θ represents the logarithm of the true abundance, ϕ is a probe-specific effect, and ε accounts for measurement error. This is the *additive-multiplicative error model* for microarray data, which was first proposed by Rocke and Durbin (2001) and in a closely related form by Ideker et al. (2000).

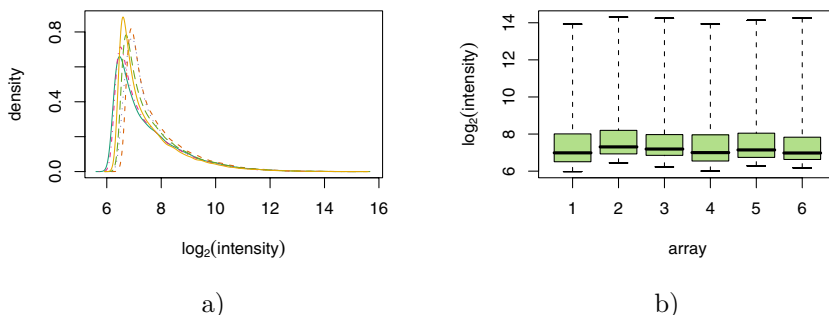


Figure 1.1. a) Density estimates of data from six replicate Affymetrix arrays. The x -axis is on a logarithmic scale (base 2). b) Box-plots.

Different arrays will have different distributions of B and different values of α , resulting in quite different distributions of the values of Y even if S is the same. To see this, let us look at the empirical distribution of six replicate Affymetrix arrays.

```
> library("affy")
> library("SpikeInSubset")
> data("spikein95")
> hist(spikein95)
> boxplot(spikein95)
```

The resulting plots are shown in Figure 1.1.

Part of the task of preprocessing is to eliminate the effect of background noise. Notice in Figure 1.1 that the smallest values attained are around 64, with slight differences between the arrays. We know that many of the probes are not supposed to be hybridizing to anything (as not all genes are expressed), so many measurements should indeed be 0. A bottom line effect of not appropriately removing background noise is that estimates of differential expression are biased. Specifically, the ratios are attenuated toward 1. This can be seen using the Affymetrix spike-in experiment, where genes were spiked in at known concentrations. Figure 1.2a shows the observed concentrations versus nominal concentrations of the spiked-in genes. Measurements with smaller nominal concentrations appear to be affected by attenuation bias. To see this, notice that the curve has a slope of about 1 for high nominal concentrations but gets flat as the nominal concentration gets closer to 0. This is consistent with the additive background noise model (1.1). Mathematically, it is easy to see that if s_1/s_2 is the true ratio and b_1 and b_2 are approximately equal positive numbers, then $(s_1 + b_1)/(s_2 + b_2)$ is closer to 1 than the true ratio, and the more so the smaller the absolute values of the s_i are compared to the b_i .

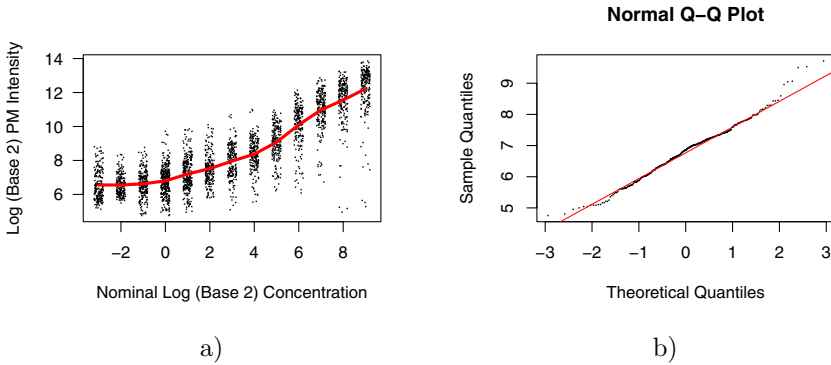


Figure 1.2. a) Plot of observed against nominal concentrations. Both axes are on the logarithmic scale (base 2). The curve represents the average value of all probes at each nominal concentration. Nominal concentrations are measured in picomol. b) Normal quantile-quantile plot of the logarithmic (base 2) intensities for all probes with the same nominal concentration of 1 picomol.

Figure 1.2b shows a normal quantile-quantile plot of logarithmic intensities of probes for genes with the same nominal concentration. Note that these appear to roughly follow a normal distribution. Figure 1.2 supports the multiplicative error assumption of model 1.1.

1.4.2 The variance-bias trade-off

A typical problem with many preprocessing algorithms is that much precision is sacrificed to remove background effects and improve accuracy. Model (1.1) can be used to show that subtracting unbiased estimates of background effects leads to exaggerated variance for genes with small values of a . In fact, background estimates that are often used in practice, such as the “local background values” from many image analysis programs for two-color spotted arrays and the mismatch (MM) value from Affymetrix arrays, tend to be *over*-estimates, which makes the problem even worse.

Various researchers have used models similar to Equation (1.1) to develop preprocessing algorithms that improve both accuracy and precision in a balanced way. Some of these methods propose variance stabilizing transformations (Durbin et al., 2002; Huber et al., 2002, 2004), others use estimation procedures that improve mean squared error (Irizarry et al., 2003b). Some examples will be provided in Chapters 2 and 4.

1.4.3 Sensitivity and specificity of probes

The probes on a microarray are intended to measure the abundance of the particular transcript that they are assigned to. However, probes may

differ in terms of their sensitivity and specificity. This fact is represented by the existence of ϕ in model (1.2). Here, sensitivity means that a probe's fluorescence signal indeed responds to changes in the transcript abundance; specificity, that it does not respond to other transcripts or other types of perturbations.

Probes may lack sensitivity. Some probes initially identified with a gene do not actually hybridize to any of its products. Some probes will have been developed from information that has been superseded. In some cases, the probe may correspond to a different gene or it may in fact not represent any gene. There is also the possibility of human error (Halgren et al., 2001; Knight, 2001).

A potential problem especially with short oligonucleotide technology is that the probes may not be specific, that is, in addition to matching the intended transcript, they may also match some other gene(s). In this case, we expect the observed intensity to be a composite from all matching transcripts. Note that here we are limited by the current state of knowledge of the human transcriptome. As our knowledge improves, the information about sensitivity of probes should also improve.

1.5 Conclusion

Various academic groups have demonstrated that the use of modern statistical methodology can substantially improve accuracy and precision of bottom-line results, relative to *ad hoc* procedures introduced by designers and manufacturers of the technology. In the following chapters, we provide some details of how Bioconductor tools can be used to do this, not only in microarray platforms, but also in other related technologies.

2

Preprocessing High-density Oligonucleotide Arrays

B. M. Bolstad, R. A. Irizarry, L. Gautier, and Z. Wu

Abstract

High-density oligonucleotide expression arrays are a widely used microarray platform. Affymetrix GeneChip arrays dominate this market. An important distinction between the GeneChip and other technologies is that on GeneChips, multiple short probes are used to measure gene expression levels. This makes preprocessing particularly important when using this platform. This chapter begins by describing how to import probe-level data into the system and how these data can be examined using the facilities of the *AffyBatch* class. Then we will describe background adjustment, normalization, and summarization methods. Functionality for GeneChip probe-level data is provided by the `affy`, `affyPLM`, `affycomp`, `gcrma`, and `affypdnn` packages. All these tools are useful for preprocessing probe-level data stored in an *AffyBatch* object into expression-level data stored in an *exprSet* object. Because there are many competing methods for this preprocessing step, it is useful to have a way to assess the differences. In Bioconductor, this can be carried out using the `affycomp` package, which we discuss briefly.

2.1 Introduction

The most popular microarray application is measuring genome-wide expression levels. High-density oligonucleotide expression arrays are a commonly used technology for this purpose. Affymetrix GeneChip arrays dominate this market. In this platform, the choice of preprocessing method can have enormous influence on the quality of the ultimate results. Many preprocessing methods have been proposed for high-density oligonucleotide array data. In this chapter, we discuss methodology and Bioconductor tools