

METHODS IN MOLECULAR BIOLOGY

Series Editor
John M. Walker
School of Life and Medical Sciences
University of Hertfordshire
Hatfield, Hertfordshire, UK

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

Statistical Genomics

Methods and Protocols

Edited by

Ewy Mathé

Biomedical Informatics, College of Medicine, Ohio State University, Columbus, OH, USA

Sean Davis

National Institutes of Health, National Cancer Institute, Bethesda, MD, USA

Editors

Ewy Mathé
Biomedical Informatics, College of Medicine
Ohio State University
Columbus, OH, USA

Sean Davis
National Institutes of Health
National Cancer Institute
Bethesda, MD, USA

ISSN 1064-3745

ISSN 1940-6029 (electronic)

Methods in Molecular Biology

ISBN 978-1-4939-3576-5

ISBN 978-1-4939-3578-9 (eBook)

DOI 10.1007/978-1-4939-3578-9

Library of Congress Control Number: 2016933669

© Springer Science+Business Media New York 2016

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

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

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

Printed on acid-free paper

This Humana Press imprint is published by Springer Nature

The registered company is Springer Science+Business Media LLC New York

Preface

Statistical Analysis of Genomic Data is, indeed, a very broad topic. We have attempted in this volume to provide chapters with cross-cutting groundwork materials, public data repositories, common applications of statistical analysis in genomics, and some representative toolsets for operating on genomic data. While we cannot be comprehensive in a single volume, we have tried to provide a breadth of both applications and tools. The authors of the individual chapters have largely focused on practical aspects of their topics, as we feel that application is an integral part of learning about statistical analysis of genomic data.

More specifically, the volume is divided into four parts. In the first part, we have included overview material and resources that can be applied across topics later in the book. In the second part, a couple of prominent public repositories for genomic data are covered in some depth. In the third part, several different biological applications of statistical genomics are presented. In the fourth and last part, software tools that can be used to facilitate ad hoc analysis and data integration are highlighted. Finally, we thank the chapter authors for the generosity of their time and insight in preparing their excellent contributions.

Columbus, OH, USA
Bethesda, MD, USA

Ewy Mathé
Sean Davis

Contents

<i>Preface</i>	<i>v</i>
<i>Contributors</i>	<i>ix</i>
PART I GROUNDWORK	
1 Overview of Sequence Data Formats	3
<i>Hongen Zhang</i>	
2 Integrative Exploratory Analysis of Two or More Genomic Datasets	19
<i>Chen Meng and Aedin Culhane</i>	
3 Study Design for Sequencing Studies	39
<i>Loren A. Honaas, Naomi S. Altman, and Martin Krzywinski</i>	
4 Genomic Annotation Resources in R/Bioconductor	67
<i>Marc R.J. Carlson, Hervé Pagès, Sonali Arora, Valerie Obenchain, and Martin Morgan</i>	
PART II PUBLIC GENOMIC DATA	
5 The Gene Expression Omnibus Database	93
<i>Emily Clough and Tanya Barrett</i>	
6 A Practical Guide to The Cancer Genome Atlas (TCGA)	111
<i>Zhining Wang, Mark A. Jensen, and Jean Claude Zenklusen</i>	
PART III APPLICATIONS	
7 Working with Oligonucleotide Arrays	145
<i>Benilton S. Carvalho</i>	
8 Meta-Analysis in Gene Expression Studies	161
<i>Levi Waldron and Markus Riester</i>	
9 Practical Analysis of Genome Contact Interaction Experiments	177
<i>Mark A. Carty and Olivier Elemento</i>	
10 Quantitative Comparison of Large-Scale DNA Enrichment Sequencing Data	191
<i>Matthias Lienhard and Lukas Chavez</i>	
11 Variant Calling From Next Generation Sequence Data	209
<i>Nancy F. Hansen</i>	
12 Genome-Scale Analysis of Cell-Specific Regulatory Codes Using Nuclear Enzymes	225
<i>Songjoon Baek and Myong-Hee Sung</i>	

PART IV TOOLS

13	NGS-QC Generator: A Quality Control System for ChIP-Seq and Related Deep Sequencing-Generated Datasets	243
	<i>Marco Antonio Mendoza-Parra, Mohamed-Ashick M. Saleem, Matthias Blum, Pierre-Etienne Cholley, and Hinrich Gronemeyer</i>	
14	Operating on Genomic Ranges Using BEDOPS	267
	<i>Shane Neph, Alex P. Reynolds, M. Scott Kuehn, and John A. Stamatoyannopoulos</i>	
15	GMAP and GSNAP for Genomic Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality	283
	<i>Thomas D. Wu, Jens Reeder, Michael Lawrence, Gabe Becker, and Matthew J. Brauer</i>	
16	Visualizing Genomic Data Using Gviz and Bioconductor	335
	<i>Florian Habne and Robert Ivanek</i>	
17	Introducing Machine Learning Concepts with WEKA	353
	<i>Tony C. Smith and Eibe Frank</i>	
18	Experimental Design and Power Calculation for RNA-seq Experiments	379
	<i>Zhijin Wu and Hao Wu</i>	
19	It's DE-licious: A Recipe for Differential Expression Analyses of RNA-seq Experiments Using Quasi-Likelihood Methods in edgeR	391
	<i>Aaron T.L. Lun, Yunshun Chen, and Gordon K. Smyth</i>	
	<i>Index</i>	417

Contributors

- NAOMI S. ALTMAN • *Department of Statistics and Huck Institutes of Life Sciences, The Pennsylvania State University, PA, USA*
- SONALI ARORA • *Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA*
- SONGJOON BAEK • *Laboratory of Receptor Biology and Gene Expression, National Cancer Institute, National Institutes of Health, Bethesda, MD, USA*
- TANYA BARRETT • *National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA*
- GABRIEL BECKER • *Genentech, South San Francisco, CA, USA*
- MATTHIAS BLUM • *Equipe Labellisée Ligue Contre le Cancer, Department of Functional Genomics and Cancer, Institut de Génétique et de Biologie Moléculaire et Cellulaire (IGBMC)/CNRS/INSERM/Université de Strasbourg, Illkirch Cedex, France*
- MATTHEW BRAUER • *Genentech, South San Francisco, CA, USA*
- MARC R.J. CARLSON • *Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA; Seattle Children's Research Institute, Seattle, WA, USA*
- MARK A. CARTY • *Institute for Computational Biomedicine, Weill Cornell Medical College, New York, NY, USA; Memorial Sloan Kettering Cancer Center, New York, NY, USA*
- BENILTON S. CARVALHO • *Brazilian Institute of Neuroscience and Neurotechnology (BRAINN) and Department of Statistics, University of Campinas, Campinas, São Paulo, Brazil*
- LUKAS CHAVEZ • *German Cancer Research Center, Heidelberg, Germany*
- YUNSHUN CHEN • *Walter and Eliza Hall Institute of Medical Research, Parkville, VIC, Australia*
- PIERRE-ETIENNE CHOLLEY • *Equipe Labellisée Ligue Contre le Cancer, Department of Functional Genomics and Cancer, Institut de Génétique et de Biologie Moléculaire et Cellulaire (IGBMC)/CNRS/INSERM/Université de Strasbourg, Illkirch Cedex, France*
- EMILY CLOUGH • *National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA*
- AEDIN CULHANE • *Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA, USA; Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA*
- OLIVIER ELEMENTO • *Institute for Computational Biomedicine, Weill Cornell Medical College, New York, NY, USA*
- EIBE FRANK • *Department of Computer Science, University of Waikato, Hamilton, New Zealand*
- HINRICH GRONEMEYER • *Equipe Labellisée Ligue Contre le Cancer, Department of Functional Genomics and Cancer, Institut de Génétique et de Biologie Moléculaire et Cellulaire (IGBMC)/CNRS/INSERM/Université de Strasbourg, Illkirch Cedex, France*
- FLORIAN HAHNE • *Novartis Institute for Biomedical Research, Basel, Switzerland*
- NANCY F. HANSEN • *National Human Genome Research Institute, Rockville, MD, USA*
- LOREN A. HONAAS • *Department of Biology, The Pennsylvania State University, Wenatchee, WA, USA*

- ROBERT IVANEK • *Department of Biomedicine, University of Basel, Basel, Switzerland*
- MARK A. JENSEN • *Research Administration Directorate, Leidos Biomedical Research Inc., Frederick National Laboratory for Cancer Research, Frederick, MD, USA*
- MARTIN KRZYWINSKI • *Canada's Michael Smith Genome Sciences Centre, Vancouver, BC, Canada*
- M. SCOTT KUEHN • *Opower Inc., San Francisco, CA, USA*
- MICHAEL LAWRENCE • *Genentech, South San Francisco, CA, USA*
- MATTHIAS LIENHARD • *Max Planck Institute for Molecular Genetics, Berlin, Germany*
- AARON T. L. LUN • *Walter and Eliza Hall Institute of Medical Research, Parkville, VIC, Australia*
- MARCO ANTONIO MENDOZA-PARRA • *Equipe Labellisée Ligue Contre le Cancer, Department of Functional Genomics and Cancer, Institut de Génétique et de Biologie Moléculaire et Cellulaire (IGBMC)/CNRS/INSERM/Université de Strasbourg, Illkirch Cedex, France*
- CHEN MENG • *Chair of Proteomics and Bioanalytics, Technische Universität München, Freising, Germany*
- MARTIN MORGAN • *Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA*
- SHANE NEPH • *Department of Genome Sciences, Altius Institute for Biomedical Sciences, Seattle, WA, USA*
- VALERIE OBENCHAIN • *Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA*
- HERVÉ PAGÈS • *Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA*
- JENS REEDER • *Genentech, South San Francisco, CA, USA*
- ALEX P. REYNOLDS • *Department of Genome Sciences, Altius Institute for Biomedical Sciences, Seattle, WA, USA*
- MARKUS RIESTER • *Novartis Institutes for BioMedical Research (NIBR), Cambridge, MA, USA*
- MOHAMED-ASHICK M. SALEEM • *Equipe Labellisée Ligue Contre le Cancer, Department of Functional Genomics and Cancer, Institut de Génétique et de Biologie Moléculaire et Cellulaire (IGBMC)/CNRS/INSERM/Université de Strasbourg, Illkirch Cedex, France*
- TONY C. SMITH • *Department of Computer Science, University of Waikato, Hamilton, New Zealand*
- GORDON K. SMYTH • *Walter and Eliza Hall Institute of Medical Research, Parkville, VIC, Australia*
- JOHN A. STAMATOYANNOPOULOS • *Altius Institute for Biomedical Sciences, Seattle, WA, USA; Department of Medicine, University of Washington, Seattle, WA, USA; Department of Genome Sciences, University of Washington, Seattle, WA, USA*
- MYONG-HEE SUNG • *Laboratory of Receptor Biology and Gene Expression, National Cancer Institute, National Institutes of Health, Bethesda, MD, USA; Laboratory of Molecular Biology and Immunology, National Institute on Aging, National Institutes of Health, Baltimore, MD, USA*
- LEVI WALDRON • *Department of Epidemiology and Biostatistics, City University of New York, School of Public Health, New York, NY, USA*
- ZHINING WANG • *Center for Cancer Genomics, National Cancer Institute, National Institutes of Health, Bethesda, MD, USA*

THOMAS D. WU • *Genentech, South San Francisco, CA, USA*

ZHIJIN WU • *Department of Biostatistics, Brown University, Providence, RI, USA*

HAO WU • *Department of Biostatistics and Bioinformatics, Rollins School of Public Health,
Emory University, Atlanta, GA, USA*

JEAN CLAUDE ZENKLUSEN • *Center for Cancer Genomics, National Cancer Institute,
National Institutes of Health, Bethesda, MD, USA*

HONGEN ZHANG • *Center for Cancer Research, National Institutes of Health, National
Cancer Institute, Bethesda, MD, USA*