



Bioinformatics Algorithms

TECHNIQUES AND APPLICATIONS



EDITED BY

Ion I. Măndoiu and Alexander Zelikovsky

Wiley Series on Bioinformatics: Computational Techniques and Engineering

Yi Pan and Albert Y. Zomaya, Series Editors

Q811.4
B615

BIOINFORMATICS ALGORITHMS

Techniques and Applications

Edited by

Ion I. Măndoiu and Alexander Zelikovsky



**WILEY-
INTERSCIENCE**



E2008001385

A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2008 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey
Published simultaneously in Canada

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400, fax 978-646-8600, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201)-748-6011, fax (201)-748-6008.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services please contact our Customer Care Department within the U. S. at 877-762-2974, outside the U. S. at 317-572-3993 or fax 317-572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print, however, may not be available in electronic format.

Library of Congress Cataloging-in-Publication Data:

Bioinformatics algorithms : techniques and applications / edited by Ion I. Mandoiu and Alexander Zelikovskiy.

p. cm.

ISBN 978-0-470-09773-1 (cloth)

I. Bioinformatics. 2. Algorithms. I. Mandoiu, Ion. II. Zelikovskiy, Alexander.

QH324.2B5472 2008

572.80285-dc22

2007034307

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

BIOINFORMATICS ALGORITHMS

PREFACE

Bioinformatics, broadly defined as the interface between biological and computational sciences, is a rapidly evolving field, driven by advances in high throughput technologies that result in an ever increasing variety and volume of experimental data to be managed, integrated, and analyzed. At the core of many of the recent developments in the field are novel algorithmic techniques that promise to provide the answers to key challenges in postgenomic biomedical sciences, from understanding mechanisms of genome evolution and uncovering the structure of regulatory and protein-interaction networks to determining the genetic basis of disease susceptibility and elucidation of historical patterns of population migration.

This book aims to provide an in-depth survey of the most important developments in bioinformatics algorithms in the postgenomic era. It is neither intended as an introductory text in bioinformatics algorithms nor as a comprehensive review of the many active areas of bioinformatics research—to readers interested in these we recommend the excellent textbook *An Introduction to Bioinformatics Algorithms* by Jones and Pevzner and the *Handbook of Computational Molecular Biology* edited by Srinivas Aluru. Rather, our intention is to make a carefully selected set of advanced algorithmic techniques accessible to a broad readership, including graduate students in bioinformatics and related areas and biomedical professionals who want to expand their repertoire of algorithmic techniques. We hope that our emphasis on both in-depth presentation of theoretical underpinnings and applications to current biomedical problems will best prepare the readers for developing their own extensions to these techniques and for successfully applying them in new contexts.

The book features 21 chapters authored by renowned bioinformatics experts who are active contributors to the respective subjects. The chapters are intended to be largely independent, so that readers do not have to read every chapter nor have to read them in a particular order. The opening chapter is a thought provoking discussion of

the role that algorithms should play in 21st century bioinformatics education. The remaining 20 chapters are grouped into the following five parts:

- Part I focuses on algorithmic techniques that find applications to a wide range of bioinformatics problems, including chapters on dynamic programming, graph-theoretical methods, hidden Markov models, sorting the fast Fourier transform, seeding, and phylogenetic networks comparison approximation algorithms.
- Part II is devoted to algorithms and tools for genome and sequence analysis. It includes chapters on formal and approximate models for gene clusters, and on advanced algorithms for multiple and non-overlapping local alignments and genome things, multiplex PCR primer set selection, and sequence and network motif finding.
- Part III concentrates on algorithms for microarray design and data analysis. The first chapter is devoted to algorithms for microarray layout, with next two chapters describing methods for missing value imputation and meta-analysis of gene expression data.
- Part IV explores algorithmic issues arising in analysis of genetic variation across human population. Two chapters are devoted to computational inference of haplotypes from commonly available genotype data, with a third chapter describing optimization techniques for disease association search in epidemiologic case/control genotype data studies.
- Part V gives an overview of algorithmic approaches in structural and systems biology. First two chapters give a formal introduction to topological and structural classification in biochemistry, while the third chapter surveys protein–protein and domain–domain interaction prediction.

We are grateful to all the authors for their excellent contributions, without which this book would not have been possible. We hope that their deep insights and fresh enthusiasm will help attracting new generations of researchers to this dynamic field. We would also like to thank series editors Yi Pan and Albert Y. Zomaya for nurturing this project since its inception, and the editorial staff at Wiley Interscience for their patience and assistance throughout the project. Finally, we wish to thank our friends and families for their continuous support.

ION I. MĂNDOIU AND ALEXANDER ZELIKOVSKY

CONTRIBUTORS

Sudha Balla, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA

Sergey Bereg, Department of Computer Science, University of Texas at Dallas, Dallas, TX, USA

Anne Bergeron, Comparative Genomics Laboratory, Université du Québec à Montréal, Canada

Paola Bonizzoni, Dipartimento di Informatica, Sistemistica e Comunicazione, Università degli Studi di Milano-Bicocca, Milano, Italy

Broňa Brejová, Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, NY, USA

Dumitru Brinza, Department of Computer Science, Georgia State University, Atlanta, GA, USA

Daniel G. Brown, Cheriton School of Computer Science, University of Waterloo, Waterloo, Ontario, Canada

Zhipeng Cai, Department of Computing Science, University of Alberta, Edmonton, Alberta, Canada

Cedric Chauve, Department of Mathematics, Simon Fraser University, Vancouver, Canada

Bhaskar DasGupta, Department of Computer Science, University of Illinois at Chicago, Chicago, IL, USA

Sérgio A. de Carvalho Jr., Technische Fakultät, Bielefeld University, D-33594 Bielefeld, Germany

- Jaime Davila**, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA
- Gianluca Della Vedova**, Dipartimento di Statistica, Università degli Studi di Milano-Bicocca, Milano, Italy
- Riccardo Dondi**, Dipartimento di Scienze dei Linguaggi, della Comunicazione e degli Studi Culturali, Università degli Studi di Bergamo, Bergamo, Italy
- Laurent Essioux**, Hoffmann-La Roche Ltd, Basel, Switzerland
- Bruce Futcher**, Department of Molecular Genetics and Microbiology, Stony Brook University, Stony Brook, NY, USA
- Yannick Gingras**, Comparative Genomics Laboratory, Université du Québec à Montréal, Canada
- Daniel Gusfield**, Department of Computer Science, University of California, Davis, CA, USA
- Robert W. Harrison**, Department of Computer Science, Georgia State University, Atlanta, GA, USA
- Jingwu He**, Department of Computer Science, Georgia State University, Atlanta, GA, USA
- Raja Jothi**, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA
- Ming-Yang Kao**, Department of Electrical Engineering and Computer Science, Northwestern University, Evanston, IL, USA
- Gunnar W. Klau**, Mathematics in Life Sciences Group, Department of Mathematics and Computer Science, University Berlin, and DFG Research Center MATHEON “Mathematics for Key Technologies,” Berlin, Germany
- Mikko Koivisto**, Department of Computer Science and HIIT Basic Research Unit, University of Helsinki, Finland
- Kishori M. Konwar**, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA
- Guohui Lin**, Department of Computing Science, University of Alberta, Edmonton, Alberta, Canada
- Sebastien Lissarrague**, Genset SA, Paris, France
- Ion I. Măndoiu**, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA
- Heikki Mannila**, Department of Computer Science and HIIT Basic Research Unit, University of Helsinki, Finland
- Giancarlo Mauri**, Dipartimento di Informatica, Sistemistica e Comunicazione, Università degli Studi di Milano-Bicocca, Milano, Italy

Steven Hecht Orzack, Fresh Pond Research Institute, Cambridge, MA, USA

Pavel Pevzner, Department of Computer Science and Engineering, University of California, San Diego, CA, USA

Teresa M. Przytycka, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

Saumyadipta Pyne, The Broad Institute of MIT and Harvard, Cambridge, MA, USA

Sven Rahmann, Bioinformatics for High-Throughput Technologies, Department of Computer Science 11, Technical University of Dortmund, Dortmund, Germany

Sanguthevar Rajasekaran, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA

Pasi Rastas, Department of Computer Science and HIIT Basic Research Unit, University of Helsinki, Finland

Alexander C. Russell, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA

Yi Shi, Department of Computing Science, University of Alberta, Edmonton, Alberta, Canada

Alexander A. Shvartsman, Department of Computer Science and Engineering, University of Connecticut, Storrs, Connecticut, USA

Steve Skiena, Department of Computer Science, Stony Brook University, Stony Brook, NY, USA

Lakshman Subrahmanyam, University of Massachusetts Medical School, Worcester, MA, USA

Sing-Hoi Sze, Departments of Computer Science and of Biochemistry and Biophysics, Texas A&M University, College Station, Texas, USA

Haixu Tang, School of Informatics and Center for Genomic and Bioinformatics, Indiana University, Bloomington, IN, USA

Esko Ukkonen, Department of Computer Science and HIIT Basic Research Unit, University of Helsinki, Finland

Tomáš Vinař, Department of Biological Statistics and Computational Biology, Cornell University, Ithaca, NY, USA

Patra Volarath, Department of Computer Science, Georgia State University, Atlanta, GA, USA

Hao Wang, Department of Computer Science, Georgia State University, Atlanta, GA, USA

Yuzhen Ye, The Burnham Institute for Medical Research, San Diego, CA, USA

Alexander Zelikovsky, Department of Computer Science, Georgia State University, Atlanta, GA, USA

Elena Zotenko, National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA and Department of Computer Science, University of Maryland, College Park, MD, USA

CONTENTS

<i>Preface</i>	ix
<i>Contributors</i>	xi
1 Educating Biologists in the 21st Century: Bioinformatics Scientists versus Bioinformatics Technicians	1
<i>Pavel Pevzner</i>	
PART I TECHNIQUES	7
2 Dynamic Programming Algorithms for Biological Sequence and Structure Comparison	9
<i>Yuzhen Ye and Haixu Tang</i>	
3 Graph Theoretical Approaches to Delineate Dynamics of Biological Processes	29
<i>Teresa M. Przytycka and Elena Zotenko</i>	
4 Advances in Hidden Markov Models for Sequence Annotation	55
<i>Broňa Brejová, Daniel G. Brown, and Tomáš Vinař</i>	
5 Sorting- and FFT-Based Techniques in the Discovery of Biopatterns	93
<i>Sudha Balla, Sanguthevar Rajasekaran, and Jaime Davila</i>	

6	A Survey of Seeding for Sequence Alignment	117
	<i>Daniel G. Brown</i>	
7	The Comparison of Phylogenetic Networks: Algorithms and Complexity	143
	<i>Paola Bonizzoni, Gianluca Della Vedova, Riccardo Dondi, and Giancarlo Mauri</i>	
PART II GENOME AND SEQUENCE ANALYSIS		175
8	Formal Models of Gene Clusters	177
	<i>Anne Bergeron, Cedric Chauve, and Yannick Gingras</i>	
9	Integer Linear Programming Techniques for Discovering Approximate Gene Clusters	203
	<i>Sven Rahmann and Gunnar W. Klau</i>	
10	Efficient Combinatorial Algorithms for DNA Sequence Processing	223
	<i>Bhaskar DasGupta and Ming-Yang Kao</i>	
11	Algorithms for Multiplex PCR Primer Set Selection with Amplification Length Constraints	241
	<i>K.M. Konwar, I.I. Măndoiu, A.C. Russell, and A.A. Shvartsman</i>	
12	Recent Developments in Alignment and Motif Finding for Sequences and Networks	259
	<i>Sing-Hoi Sze</i>	
PART III MICROARRAY DESIGN AND DATA ANALYSIS		277
13	Algorithms for Oligonucleotide Microarray Layout	279
	<i>Sérgio A. De Carvalho Jr. and Sven Rahmann</i>	
14	Classification Accuracy Based Microarray Missing Value Imputation	303
	<i>Yi Shi, Zhipeng Cai, and Guohui Lin</i>	
15	Meta-Analysis of Microarray Data	329
	<i>Saumyadipta Pyne, Steve Skiena, and Bruce Fletcher</i>	

PART IV GENETIC VARIATION ANALYSIS	353
16 Phasing Genotypes Using a Hidden Markov Model	355
<i>P. Rastas, M. Koivisto, H. Mannila, and E. Ukkonen</i>	
17 Analytical and Algorithmic Methods for Haplotype Frequency Inference: What Do They Tell Us?	373
<i>Steven Hecht Orzack, Daniel Gusfield, Lakshman Subrahmanyam, Laurent Essioux, and Sebastien Lissarrague</i>	
18 Optimization Methods for Genotype Data Analysis in Epidemiological Studies	395
<i>Dumitru Brinza, Jingwu He, and Alexander Zelikovsky</i>	
PART V STRUCTURAL AND SYSTEMS BIOLOGY	417
19 Topological Indices in Combinatorial Chemistry	419
<i>Sergey Bereg</i>	
20 Efficient Algorithms for Structural Recall in Databases	439
<i>Hao Wang, Patra Volarath, and Robert W. Harrison</i>	
21 Computational Approaches to Predict Protein–Protein and Domain–Domain Interactions	465
<i>Raja Jothi and Teresa M. Przytycka</i>	
<i>Index</i>	493

1

EDUCATING BIOLOGISTS IN THE 21ST CENTURY: BIOINFORMATICS SCIENTISTS VERSUS BIOINFORMATICS TECHNICIANS¹

PAVEL PEVZNER

*Department of Computer Science and Engineering, University of California, San Diego,
CA, USA*

For many years algorithms were taught exclusively to computer scientists, with relatively few students from other disciplines attending algorithm courses. A biology student in an algorithm class would be a surprising and unlikely (though not entirely unwelcome) guest in the 1990s. Things have changed; some biology students now take some sort of *Algorithms 101*. At the same time, curious computer science students often take *Genetics 101*.

Here comes an important question of how to teach bioinformatics in the 21st century. Will we teach bioinformatics to future biology students as a collection of cookbook-style recipes or as a computational science that first explain ideas and builds on applications afterward? This is particularly important at the time when bioinformatics courses may soon become *required* for all graduate biology students in leading universities. Not to mention that some universities have already started undergraduate bioinformatics programs, and discussions are underway about adding new computational courses to the standard undergraduate biology curriculum—a dramatic paradigm shift in biology education.

¹Reprinted from *Bioinformatics* 20:2159–2161 (2004) with the permission of Oxford University Press.

Since bioinformatics is a computational science, a bioinformatics course should strive to present the principles and the ideas that drive an algorithm's design or explain the crux of a statistical approach, rather than to be a stamp collection of the algorithms and statistical techniques themselves. Many existing bioinformatics books and courses reduce bioinformatics to a compendium of computational protocols without even trying to explain the computational ideas that drove the development of bioinformatics in the past 30 years. Other books (written by computer scientists for computer scientists) try to explain bioinformatics ideas at the level that is well above the computational level of most biologists. These books often fail to connect the computational ideas and applications, thus reducing a biologist's motivation to invest time and effort into such a book. We feel that focusing on ideas has more intellectual value and represents a long-term investment: protocols change quickly, but the computational ideas don't seem to. However, the question of how to deliver these ideas to biologists remains an unsolved educational riddle.

Imagine Alice (a computer scientist), Bob (a biologist), and a chessboard with a lonely king in the lower right corner. Alice and Bob are bored one Sunday afternoon so they play the following game. In each turn, a player may either move a king one square to the left, one square up, or one square "north-west" along the diagonal. Slowly but surely, the king moves toward the upper left corner and the player who places the king to this square wins the game. Alice moves first.

It is not immediately clear what the winning strategy is. Does the first player (or the second) always have an advantage? Bob tries to analyze the game and applies a reductionist approach, and he first tries to find a strategy for the simpler game on a 2×2 board. He quickly sees that the second player (himself, in this case) wins in 2×2 game and decides to write the recipe for the "winning algorithm:"

If Alice moves the king diagonally, I will move him diagonally and win. If Alice moves the king to the left, I will move him to the left as well. As a result, Alice's only choice will be to move the king up. Afterward, I will move the king up again and will win the game. The case when Alice moves the king up is symmetric.

Inspired by this analysis Bob makes a leap of faith: the second player (i.e., himself) wins in any $n \times n$ game. Of course, every hypothesis must be confirmed by experiment, so Bob plays a few rounds with Alice. He tries to come up with a simple recipe for the 3×3 game, but there are already a large number of different game sequences to consider. There is simply no hope of writing a recipe for the 8×8 game since the number of different strategies Alice can take is enormous.

Meanwhile, Alice does not lose hope of finding a winning strategy for the 3×3 game. Moreover, she understands that recipes written in the cookbook style that Bob uses will not help very much: recipe-style instructions are not a sufficiently expressive language for describing algorithms. Instead, she begins by drawing the following table that is filled by the symbols \uparrow , \leftarrow , \nwarrow , and $*$. The entry in position (i, j) (that is, the i th row and the j th column) describes the move that Alice will make in the $i \times j$ game. A \leftarrow indicates that she should move the king to the left. A \uparrow indicates that she should move the king up. A \nwarrow indicates that she should move the king diagonally, and $*$

indicates that she should not bother playing the game because she will definitely lose against an opponent who has a clue.

	0	1	2	3	4	5	6	7	8
0		←	*	←	*	←	*	←	*
1	↑	↖	↑	↖	↑	↖	↑	↖	↑
2	*	←	*	←	*	←	*	←	*
3	↑	↖	↑	↖	↑	↖	↑	↖	↑
4	*	←	*	←	*	←	*	←	*
5	↑	↖	↑	↖	↑	↖	↑	↖	↑
6	*	←	*	←	*	←	*	←	*
7	↑	↖	↑	↖	↑	↖	↑	↖	↑
8	*	←	*	←	*	←	*	←	*

For example, if she is faced with the 3×3 game, she finds a ↖ in the third row and third column, indicating that she should move the king diagonally. This makes Bob take the first move in a 2×2 game, which is marked with a *. No matter what he does, Alice wins using instructions in the table.

Impressed by the table, Bob learns how to use it to win the 8×8 game. However, Bob does not know how to construct a similar table for the 20×20 game. The problem is not that Bob is stupid (quite the opposite, a bit later he even figured out how to use the symmetry in this game, thus eliminating the need to memorize Alice's table) but that he has not studied algorithms. Even if Bob figured out the logic behind 20×20 game, a more general $20 \times 20 \times 20$ game on a three-dimensional chessboard would turn into an impossible conundrum for him since he never took *Algorithms 101*.

There are two things Bob could do to remedy this situation. First, he could take a class in algorithms to learn how to solve puzzle-like combinatorial problems. Second, he could memorize a suitably large table that Alice gives him and use that to play the game. Leading questions notwithstanding, what would you do as a biologist?

Of course, the answer we expect to hear is "Why in the world do I care about a game with a lonely king and two nerdy people? I'm interested in biology, and this game has nothing to do with me." This is not actually true: the chess game is, in fact, the ubiquitous *sequence alignment* problem in disguise. Although it is not immediately clear what DNA sequence alignment and our chess game have in common, the computational idea used to solve both problems is the same. The fact that Bob was not able to find the strategy for the game indicates that he does not understand how alignment algorithms work either. He might disagree if he uses alignment algorithms or BLAST on a daily basis, but we argue that since he failed to come up with a strategy, he will also fail when confronted with a new flavor of an alignment problem or a particularly complex bioinformatics analysis. More troubling to Bob, he may find it difficult to compete with the scads of new biologists and computer scientists who think algorithmically about biological problems.

Many biologists are comfortable using algorithms such as BLAST or GenScan without really understanding how the underlying algorithm works. This is not substantially different from a diligent robot following Alice's table, but it does have an important consequence. BLAST solves a particular problem only approximately and it has certain systematic weaknesses (we're not picking on BLAST here). Users that do not know how BLAST works might misapply the algorithm or misinterpret the results it returns (see Iyer et al. Quoderat demonstrandum? The mystery of experimental validation of apparently erroneous computational analyses of protein sequences. *Genome Biol.*, 2001, 2(12):RESEARCH0051). Biologists sometimes use bioinformatics tools simply as computational protocols in quite the same way that an uninformed mathematician might use experimental protocols without any background in biochemistry or molecular biology. In either case, important observations might be missed or incorrect conclusions drawn. Besides, intellectually interesting work can quickly become mere drudgery if one does not really understand it.

Many recent bioinformatics books cater to a protocol-centric pragmatic approach to bioinformatics. They focus on parameter settings, application-specific features, and other details without revealing the *computational ideas* behind the algorithms. This trend often follows the tradition of biology books to present material as a collection of facts and discoveries. In contrast, introductory books in algorithms and mathematics usually focus on ideas rather than on the details of computational recipes. In principle, one can imagine a calculus book teaching physicists and engineers how to take integrals *without* any attempt to explain *what is* integral. Although such a book is not that difficult to write, physicists and engineers somehow escaped this curse, probably because they understand that the recipe-based approach to science is doomed to fail. Biologists are less lucky and many biology departments now offer recipe-based bioinformatics courses without first sending their students to *Algorithms 101* and *Statistics 101*. Some of the students who take these classes get excited about bioinformatics and try to pursue a research career in bioinformatics. Many of them do not understand that, with a few exceptions, such courses prepare *bioinformatics technicians* rather than *bioinformatics scientists*.

Bioinformatics is often defined as "applications of computers in biology." In recent decades, biology has raised fascinating mathematical problems, and reducing bioinformatics to "applications of computers in biology" diminishes the rich intellectual content of bioinformatics. Bioinformatics has become a part of modern biology and often dictates new fashions, enables new approaches, and drives further biological developments. Simply using bioinformatics as a toolkit without understanding the main computational ideas is not very different than using a PCR kit without knowing how PCR works.

Bioinformatics has affected more than just biology: it has also had a profound impact on the computational sciences. Biology has rapidly become a large source for new algorithmic and statistical problems, and has arguably been the target for more algorithms than any of the other fundamental sciences. This link between computer science and biology has important educational implications that change the way we teach computational ideas to biologists, as well as how applied algorithms are taught to computer scientists.