

Rui Jiang  
Xuegong Zhang  
Michael Q. Zhang *Editors*

# 生物信息学课程导引

——生物信息学研究生暑期学校讲义

## Basics of Bioinformatics

Lecture Notes of the Graduate Summer  
School on Bioinformatics of China



清华大学出版社



Springer

TUP-Springer Project

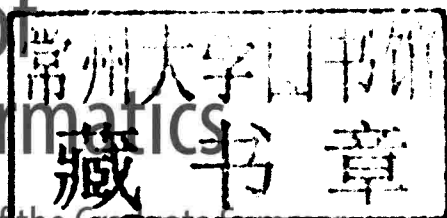
Rui Jiang  
Xuegong Zhang  
Michael Q. Zhang *Editors*

# 生物信息学课程导引

——生物信息学研究生暑期学校讲义

Basics of  
Bioinformatics

Lecture Notes of the Graduate Summer  
School on Bioinformatics of China



清华大学出版社  
北京



Springer

## 内 容 简 介

本书根据清华大学承办的全国生物信息学暑期学校课程,高度概括地介绍了与生物信息学研究紧密相关的 11 门基础课程和 15 个前沿专题报告。全书分 12 章,包括:生物信息学引论、生物信息学中的基础统计、计算基因组学专题、生物信息学中的高级统计、计算生物学算法基础、生物信息学中的多元统计、人类疾病关联研究方法与实例、生物信息学中的数据挖掘与知识发现、生物信息学应用工具、蛋白质结构与功能基础、中医药研究的计算系统生物学方法、生物信息学与计算系统生物学前沿等。本书不仅可以作为生物信息学初学者的入门读物,还可作为生物信息学领域专业研究人员高度概括而又不失系统性的参考书籍。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话:010-62782989 13701121933

### 图书在版编目(CIP)数据

生物信息学课程导引=Basics of Bioinformatics: 英文/江瑞,张学工,张奇伟主编.--北京:清华大学出版社,2014

ISBN 978-7-302-32359-4

I. ①生… II. ①江… ②张… ③张… III. ①生物信息论—高等学校—教学参考资料—英文 IV. ①Q811.4

中国版本图书馆 CIP 数据核字(2013)第 094210 号

责任编辑:王一玲

封面设计:何凤霞

责任校对:白 蕾

责任印制:刘海龙

出版发行:清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址:北京清华大学学研大厦 A 座 邮 编:100084

社总机:010-62770175

邮 购:010-62786544

投稿与读者服务:010-62776969, [c-service@tup.tsinghua.edu.cn](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈:010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 装 者:三河市春园印刷有限公司

经 销:全国新华书店

开 本:155mm×235mm

印 张:26

字 数:462 千字

版 次:2014 年 5 月第 1 版

印 次:2014 年 5 月第 1 次印刷

印 数:1~1000

定 价:99.00 元

---

产品编号:036335-01

*Editors*

Rui Jiang  
Xuegong Zhang  
Department of Automation  
Tsinghua University  
Beijing  
China, People's Republic

Michael Q. Zhang  
Department of Molecular and Cell Biology  
The University of Texas at Dallas  
Richardson, TX, USA

Tsinghua National Laboratory  
for Information Science and Technology  
Tsinghua University  
Beijing, China, People's Republic

ISBN 978-3-642-38950-4      ISBN 978-3-642-38951-1 (eBook)  
DOI 10.1007/978-3-642-38951-1  
Springer Heidelberg New York Dordrecht London

Jointly published with Tsinghua University Press, Beijing  
ISBN: 978-7-302-32359-4 Tsinghua University Press, Beijing

Library of Congress Control Number: 2013950934

© Tsinghua University Press, Beijing and Springer-Verlag Berlin Heidelberg 2013

This work is subject to copyright. All rights are reserved by the Publishers, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publishers' locations, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

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

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

Printed on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

# Foreword

This ambitious volume is the result of the successful 2007 Graduate Summer School on Bioinformatics of China held at Tsinghua University. It is remarkable for its range of topics as well as the depth of coverage. Bioinformatics draws on many subjects for analysis of the data generated by the biological sciences and biotechnology. This foreword will describe briefly each of the 12 chapters and close with additional general comments about the field. Many of the chapters overlap and include useful introductions to concepts such as gene or Bayesian methods. This is a valuable aspect of the volume allowing a student various angles of approach to a new topic.

Chapter 1, “Basics for Bioinformatics,” defines bioinformatics as “the storage, manipulation and interpretation of biological data especially data of nucleic acids and amino acids, and studies molecular rules and systems that govern or affect the structure, function and evolution of various forms of life from computational approaches.” Thus, the first subject they turn to is molecular biology, a subject that has had an enormous development in the last decades and shows no signs of slowing down. Without a basic knowledge of biology, the bioinformatics student is greatly handicapped. From basic biology the authors turn to biotechnology, in particular, methods for DNA sequencing, microarrays, and proteomics. DNA sequencing is undergoing a revolution. The mass of data collected in a decade of the Human Genome Project from 1990 to 2001 can be generated in 1 day in 2010. This is changing the science of biology at the same time. A 1,000 genome project became a 10,000 genome project 2 years later, and one expects another zero any time now. Chromatin Immunoprecipitation or ChIP allows access to DNA bound by proteins and thus to a large number of important biological processes. Another topic under the umbrella of biological sciences is genetics, the study of heredity and inherited characteristics (phenotypes). Heredity is encoded in DNA and thus is closely related to the goals of bioinformatics. This whole area of genetics beginning with Mendel’s laws deserves careful attention, and genetics is a key aspect of the so-called genetic mapping and other techniques where the chromosomal locations of disease genes are sought.

Chapter 2, “Basic Statistics for Bioinformatics,” presents important material for the understanding and analysis of data. Probability and statistics are basic to bioinformatics, and this chapter begins with the fundamentals including many classical distributions (including the binomial, Poisson, and normal). Usually the observation of complete populations such as “all people in China over 35 years old” is not practical to obtain. Instead random samples of the population of interest are obtained and then inferences about parameters of the population are made. Statistics guides us in making those inferences and gaining information about the quality of the estimates. The chapter describes techniques such as method of moments, maximum likelihood, and Bayesian methods. Bayesian methods have become indispensable in the era of powerful computing machines. The chapter treats hypothesis testing which is less used than parameter estimation, but hypothesis testing provides understanding of  $p$ -values which are ubiquitous in bioinformatics and data analysis. Classical testing situations reveal useful statistics such as the  $t$ -statistic. Analysis of variance and regression analysis are crucial for testing and fitting large data sets. All of these methods and many more are included in the free open-source package called *R*.

Chapter 3, “Topics in Computational Genomics,” takes us on a tour of important topics that arise when complete genome information is available. The subject did not begin until nearly 2000 when complete genome sequences became a possibility. The authors present us with a list of questions, some of which are listed next. What are the genes of an organism? How are they turned off and on? How do they interact with each other? How are introns and exons organized and expressed in RNA transcripts? What are the gene products, both structure and function? How has a genome evolved? This last question has to be asked with other genomes and with members of the population comprising the species. Then the authors treat some of the questions in detail. They describe “finding protein coding genes,” “identifying promoters,” “genomic arrays and a CGH/CNP analysis,” “modeling regulatory elements,” “predicting transcription factor binding sites,” and motif enrichment and analysis. Within this last topic, for example, various word counting methods are employed including the Bayesian methods of expectation maximization and Gibbs sampling.

An alert reader will have noticed the prominence of Bayesian methods in the preceding paragraphs. Chapter 4, “Statistical Methods in Bioinformatics,” in this collection focuses on this subject. There is a nice discussion of statistical modeling and then Bayesian inference. Dynamic programming, a recursive method of optimization, is introduced and then employed in the development of Hidden Markov Models (HMMs). Of course the basics of Markov chains must also be covered. The Metropolis-Hastings algorithm, Monte Carlo Markov chains (MCMC), and Gibbs sampling are carefully presented. Then these ideas find application in the analysis of microarray data. Here the challenging aspects of multiple hypothesis testing appear, and false discovery rate analysis is described. Hierarchical clustering and bi-clustering appear naturally in the context of microarray analysis. Then the issues of sequence analysis (especially multiple sequence analysis) are approached using these HMM and Bayesian methods along with pattern discovery in the sequences.

Discovering regulatory sequence patterns is an especially important topic in this section. The topics of this chapter appear in computer science as “machine learning” or under “data mining”; here the subject is called statistical or Bayesian methods. Whatever it is named, this is an essential area for bioinformatics.

The next chapter (Chap. 5), “Algorithms in Computational Biology,” takes up the formal computational approach to our biological problems. It should be pointed out that the previous chapters contained algorithmic content, but there it was less acknowledged. It is my belief that the statistical and algorithmic approaches go hand in hand. Even with the Euclid’s algorithm example of the present chapter, there are statistical issues nearby. For example, the three descriptions of Euclid’s algorithm are analyzed for time complexity. It is easy to ask how efficient the algorithms are on randomly chosen pairs of integers. What is the expected running time of the algorithms? What is the variance? Amazingly these questions have answers which are rather deep. The authors soon turn to dynamic programming (DP), and once again they present clear illustrative examples, in this case Fibonacci numbers. Designing DP algorithms for sequence alignment is covered. Then a more recently developed area of genome rearrangements is described along with some of the impressive (and deep) results from the area. This topic is relevant to whole genome analysis as chromosomes evolve on a larger scale than just alterations of individual letters as covered by sequence alignment.

In Chap. 6, “Multivariate Statistical Methods in Bioinformatics Research,” we have a thorough excursion into multivariate statistics. This can be viewed as the third statistical chapter in this volume. Here the multivariate normal distribution is studied in its many rich incarnations. This is justified by the ubiquitous nature of the normal distribution. Just as with the bell-shaped curve which appears in one dimension due to the central limit theorem (add up enough independent random variables and suitably normalized, one gets the normal under quite general conditions), there is also a multivariate central limit theorem. Here detailed properties are described as well as related distributions such as the Wishart distribution (the analog of the chi-square). Estimation is relevant as is a multivariate  $t$ -test. Principal component analysis, factor analysis, and linear discriminant analysis are all covered with some nice examples to illustrate the power of approaches. Then classification problems and variable selection both give platforms to further illustrate and develop the methods on important bioinformatics application areas.

Chapter 7, “Association Analysis for Human Diseases: Methods and Examples,” gives us the opportunity to look more deeply into aspects of genetics. While this chapter emphasizes statistics, be aware that computational issues also drive much of the research and cannot be ignored. Population genetics is introduced and then the important subjects of genetic linkage analysis and association studies. Genomic information such as single-nucleotide polymorphisms (SNPs) provide voluminous data for many of these studies, where multiple hypothesis testing is a critical issue.

Chapter 8, “Data Mining and Knowledge Discovery Methods with Case Examples,” deals with the area of knowledge discovery and data mining. To quote the authors, this area “has emerged as an important research direction for extracting useful information from vast repositories of data of various types. The basic

concepts, problems and challenges deals with the area of knowledge discovery and data mining that has emerged as an important research direction for extracting useful information from vast repositories of data of various types. The basic concepts, problems and challenges are first briefly discussed. Some of the major data mining tasks like classification, clustering and association rule mining are then described in some detail. This is followed by a description of some tools that are frequently used for data mining. Two case examples of supervised and unsupervised classification for satellite image analysis are presented. Finally an extensive bibliography is provided.”

The valuable chapter on Applied Bioinformatics Tools (Chap. 9) provides a step-by-step description of the application tools used in the course and data sources as well as a list of the problems. It should be strongly emphasized that no one learns this material without actually having hands-on experience with the derivations and the applications. This is not a subject for contemplation only!

Protein structure and function is a vast and critically important topic. In this collection it is covered by Chap. 10, “Foundations for the Study of Structure and Function of Proteins.” There the detailed structure of amino acids is presented with their role in the various levels of protein structure (including amino acid sequence, secondary structure, tertiary structure, and spatial arrangements of the subunits). The geometry of the polypeptide chain is key to these studies as are the forces causing the three-dimensional structures (including electrostatic and van der Waals forces). Secondary structural units are classified into  $\alpha$ -helix,  $\beta$ -sheets, and  $\beta$ -turns. Structural motifs and folds are described. Protein structure prediction is an active field, and various approaches are described including homology modeling and machine learning.

Systems biology is a recently described approach to combining system-wide data of biology in order to gain a global understanding of a biological system, such as a bacterial cell. The science is far from succeeding in this endeavor in general, let alone having powerful techniques to understand the biology of multicellular organisms. It is a grand challenge goal at this time. The fascinating chapter on Computational Systems Biology Approaches for Deciphering Traditional Chinese Medicine (Chap. 11) seeks to apply the computational systems biology (CSB) approach to traditional Chinese medicine (TCM). The chapter sets up parallel concepts between CSB and CTM. In Sect. 11.3.2 the main focus is “on a CSB-based case study for TCM ZHENG—a systems biology approach with the combination of computational analysis and animal experiment to investigate Cold ZHENG and Hot ZHENG in the context of the neuro-endocrine-immune (NEI) system.” With increasing emphasis on the so-called nontraditional medicine, these studies have great potential to unlock new understandings for both CSB and TCM.

Finally I close with a few remarks about this general area. Biology is a major science for our new century; perhaps it will be the major science of the twenty-first century. However, if someone is not excited by biology, then they should find a subject that does excite them. I have almost continuously found the new discoveries such as introns or microRNA absolutely amazing. It is such a young science when such profound wonders keep showing up. Clearly no one analysis subject can

solve all the problems arising in modern computational molecular biology. Statistics alone, computer science alone, experimental molecular biology alone, none of these are sufficient in isolation. Protein structure studies require an entire additional set of tools such as classical mechanics. And as systems biology comes into play, systems of differential equations and scientific computing will surely be important. None of us can learn everything, but everyone working in this area needs a set of well-understood tools. We all learn new techniques as we proceed, learning things required to solve the problems. This requires people who evolve with the subject. This is exciting, but I admit it is hard work too. Bioinformatics will evolve as it confronts new data created by the latest biotechnology and biological sciences.

University of Southern California  
Los Angeles, USA  
March 2, 2013

Michael S. Waterman

# Contents

<b>1</b>	<b>Basics for Bioinformatics</b>	<b>1</b>
	Xuegong Zhang, Xueya Zhou, and Xiaowo Wang	
1.1	What Is Bioinformatics	1
1.2	Some Basic Biology	2
1.2.1	Scale and Time	3
1.2.2	Cells	3
1.2.3	DNA and Chromosome	5
1.2.4	The Central Dogma	6
1.2.5	Genes and the Genome	7
1.2.6	Measurements Along the Central Dogma	10
1.2.7	DNA Sequencing	10
1.2.8	Transcriptomics and DNA Microarrays	13
1.2.9	Proteomics and Mass Spectrometry	16
1.2.10	ChIP-Chip and ChIP-Seq	17
1.3	Example Topics of Bioinformatics	18
1.3.1	Examples of Algorithmatic Topics	19
1.3.2	Examples of Statistical Topics	20
1.3.3	Machine Learning and Pattern Recognition Examples	21
1.3.4	Basic Principles of Genetics	21
	References	26
<b>2</b>	<b>Basic Statistics for Bioinformatics</b>	<b>27</b>
	Yuanlie Lin and Rui Jiang	
2.1	Introduction	27
2.2	Foundations of Statistics	27
2.2.1	Probabilities	27
2.2.2	Random Variables	30
2.2.3	Multiple Random Variables	32
2.2.4	Distributions	34

2.2.5	Random Sampling .....	37
2.2.6	Sufficient Statistics .....	39
2.3	Point Estimation .....	40
2.3.1	Method of Moments .....	41
2.3.2	Maximum Likelihood Estimators .....	41
2.3.3	Bayes Estimators .....	42
2.3.4	Mean Squared Error .....	44
2.4	Hypothesis Testing .....	44
2.4.1	Likelihood Ratio Tests .....	45
2.4.2	Error Probabilities and the Power Function .....	46
2.4.3	$p$ -Values .....	48
2.4.4	Some Widely Used Tests .....	50
2.5	Interval Estimation .....	52
2.6	Analysis of Variance .....	54
2.6.1	One-Way Analysis of Variance .....	55
2.6.2	Two-Way Analysis of Variance .....	59
2.7	Regression Models .....	61
2.7.1	Simple Linear Regression .....	62
2.7.2	Logistic Regression .....	65
2.8	Statistical Computing Environments .....	66
2.8.1	Downloading and Installation .....	66
2.8.2	Storage, Input, and Output of Data .....	67
2.8.3	Distributions .....	67
2.8.4	Hypothesis Testing .....	68
2.8.5	ANOVA and Linear Model .....	68
	References .....	68
<b>3</b>	<b>Topics in Computational Genomics .....</b>	<b>69</b>
	Michael Q. Zhang and Andrew D. Smith	
3.1	Overview: Genome Informatics .....	69
3.2	Finding Protein-Coding Genes .....	71
3.2.1	How to Identify a Coding Exon? .....	72
3.2.2	How to Identify a Gene with Multiple Exons? .....	72
3.3	Identifying Promoters .....	73
3.4	Genomic Arrays and aCGH/CNP Analysis .....	75
3.5	Introduction on Computational Analysis of Transcriptional Genomics Data .....	76
3.6	Modeling Regulatory Elements .....	77
3.6.1	Word-Based Representations .....	77
3.6.2	The Matrix-Based Representation .....	78
3.6.3	Other Representations .....	79
3.7	Predicting Transcription Factor Binding Sites .....	79
3.7.1	The Multinomial Model for Describing Sequences .....	80
3.7.2	Scoring Matrices and Searching Sequences .....	81

3.7.3	Algorithmic Techniques for Identifying High-Scoring Sites .....	82
3.7.4	Measuring Statistical Significance of Matches .....	83
3.8	Modeling Motif Enrichment in Sequences .....	84
3.8.1	Motif Enrichment Based on Likelihood Models.....	84
3.8.2	Relative Enrichment Between Two Sequence Sets .....	86
3.9	Phylogenetic Conservation of Regulatory Elements .....	88
3.9.1	Three Strategies for Identifying Conserved Binding Sites .....	88
3.9.2	Considerations When Using Phylogenetic Footprinting .....	90
3.10	Motif Discovery .....	91
3.10.1	Word-Based and Enumerative Methods .....	92
3.10.2	General Statistical Algorithms Applied to Motif Discovery .....	93
3.10.3	Expectation Maximization .....	94
3.10.4	Gibbs Sampling .....	95
	References .....	96

**4 Statistical Methods in Bioinformatics .....** 101

Jun S. Liu and Bo Jiang		
4.1	Introduction .....	101
4.2	Basics of Statistical Modeling and Bayesian Inference .....	102
4.2.1	Bayesian Method with Examples .....	102
4.2.2	Dynamic Programming and Hidden Markov Model ....	104
4.2.3	Metropolis–Hastings Algorithm and Gibbs Sampling ..	107
4.3	Gene Expression and Microarray Analysis .....	109
4.3.1	Low-Level Processing and Differential Expression Identification.....	110
4.3.2	Unsupervised Learning .....	113
4.3.3	Dimension Reduction Techniques .....	117
4.3.4	Supervised Learning .....	119
4.4	Sequence Alignment .....	126
4.4.1	Pair-Wise Sequence Analysis .....	126
4.4.2	Multiple Sequence Alignment .....	129
4.5	Sequence Pattern Discovery .....	133
4.5.1	Basic Models and Approaches.....	133
4.5.2	Gibbs Motif Sampler .....	136
4.5.3	Phylogenetic Footprinting Method and the Identification of <i>Cis</i> -Regulatory Modules .....	138
4.6	Combining Sequence and Expression Information for Analyzing Transcription Regulation .....	140
4.6.1	Motif Discovery in ChIP-Array Experiment.....	140
4.6.2	Regression Analysis of Transcription Regulation .....	141
4.6.3	Regulatory Role of Histone Modification .....	143

4.7	Protein Structure and Proteomics .....	144
4.7.1	Protein Structure Prediction .....	145
4.7.2	Protein Chip Data Analysis .....	146
	References .....	147
<b>5</b>	<b>Algorithms in Computational Biology .....</b>	<b>151</b>
	Tao Jiang and Jianxing Feng	
5.1	Introduction .....	151
5.2	Dynamic Programming and Sequence Alignment .....	153
5.2.1	The Paradigm of Dynamic Programming .....	153
5.2.2	Sequence Alignment .....	155
5.3	Greedy Algorithms for Genome Rearrangement .....	157
5.3.1	Genome Rearrangements .....	157
5.3.2	Breakpoint Graph, Greedy Algorithm and Approximation Algorithm .....	159
	References .....	161
<b>6</b>	<b>Multivariate Statistical Methods in Bioinformatics Research .....</b>	<b>163</b>
	Lingsong Zhang and Xihong Lin	
6.1	Introduction .....	163
6.2	Multivariate Normal Distribution .....	163
6.2.1	Definition and Notation .....	163
6.2.2	Properties of the Multivariate Normal Distribution .....	164
6.2.3	Bivariate Normal Distribution .....	165
6.2.4	Wishart Distribution .....	167
6.2.5	Sample Mean and Covariance .....	167
6.3	One-Sample and Two-Sample Multivariate Hypothesis Tests ....	168
6.3.1	One-Sample $t$ Test for a Univariate Outcome .....	168
6.3.2	Hotelling's $T^2$ Test for the Multivariate Outcome .....	169
6.3.3	Properties of Hotelling's $T^2$ Test .....	170
6.3.4	Paired Multivariate Hotelling's $T^2$ Test .....	171
6.3.5	Examples .....	172
6.3.6	Two-Sample Hotelling's $T^2$ Test .....	174
6.4	Principal Component Analysis .....	178
6.4.1	Definition of Principal Components .....	178
6.4.2	Computing Principal Components .....	179
6.4.3	Variance Decomposition .....	179
6.4.4	PCA with a Correlation Matrix .....	180
6.4.5	Geometric Interpretation .....	181
6.4.6	Choosing the Number of Principal Components .....	183
6.4.7	Diabetes Microarray Data .....	184
6.5	Factor Analysis .....	187
6.5.1	Orthogonal Factor Model .....	187
6.5.2	Estimating the Parameters .....	188
6.5.3	An Example .....	190

6.6	Linear Discriminant Analysis .....	193
6.6.1	Two-Group Linear Discriminant Analysis .....	194
6.6.2	An Example .....	198
6.7	Classification Methods .....	200
6.7.1	Introduction of Classification Methods .....	200
6.7.2	$k$ -Nearest Neighbor Method .....	202
6.7.3	Density-Based Classification Decision Rule .....	205
6.7.4	Quadratic Discriminant Analysis .....	208
6.7.5	Logistic Regression .....	212
6.7.6	Support Vector Machine .....	214
6.8	Variable Selection .....	219
6.8.1	Linear Regression Model .....	220
6.8.2	Motivation for Variable Selection .....	221
6.8.3	Traditional Variable Selection Methods .....	222
6.8.4	Regularization and Variable Selection .....	223
6.8.5	Summary .....	231
	References .....	231
<b>7</b>	<b>Association Analysis for Human Diseases: Methods and Examples .....</b>	<b>233</b>
	Jurg Ott and Qingrun Zhang	
7.1	Why Do We Need Statistics? .....	233
7.2	Basic Concepts in Population and Quantitative Genetics .....	234
7.3	Genetic Linkage Analysis .....	236
7.4	Genetic Case-Control Association Analysis .....	237
7.4.1	Basic Steps in an Association Study .....	238
7.4.2	Multiple Testing Corrections .....	239
7.4.3	Multi-locus Approaches .....	241
7.5	Discussion .....	241
	References .....	241
<b>8</b>	<b>Data Mining and Knowledge Discovery Methods with Case Examples .....</b>	<b>243</b>
	S. Bandyopadhyay and U. Maulik	
8.1	Introduction .....	243
8.2	Different Tasks in Data Mining .....	245
8.2.1	Classification .....	245
8.2.2	Clustering .....	248
8.2.3	Discovering Associations .....	252
8.2.4	Issues and Challenges in Data Mining .....	254
8.3	Some Common Tools and Techniques .....	256
8.3.1	Artificial Neural Networks .....	256
8.3.2	Fuzzy Sets and Fuzzy Logic .....	258
8.3.3	Genetic Algorithms .....	258

8.4	Case Examples .....	259
8.4.1	Pixel Classification .....	260
8.4.2	Clustering of Satellite Images .....	262
8.5	Discussion and Conclusions .....	267
	References .....	267
<b>9</b>	<b>Applied Bioinformatics Tools .....</b>	<b>271</b>
	Jingchu Luo	
9.1	Introduction .....	271
9.1.1	Welcome .....	271
9.1.2	About This Web Site .....	273
9.1.3	Outline .....	274
9.1.4	Lectures .....	275
9.1.5	Exercises .....	276
9.2	Entrez .....	277
9.2.1	PubMed Query .....	277
9.2.2	Entrez Query .....	278
9.2.3	My NCBI .....	278
9.3	ExPASy .....	278
9.3.1	Swiss-Prot Query .....	278
9.3.2	Explore the Swiss-Prot Entry HBA_HUMAN .....	279
9.3.3	Database Query with the EBI SRS .....	279
9.4	Sequence Alignment .....	280
9.4.1	Pairwise Sequence Alignment .....	280
9.4.2	Multiple Sequence Alignment .....	281
9.4.3	BLAST .....	281
9.5	DNA Sequence Analysis .....	282
9.5.1	Gene Structure Analysis and Prediction .....	282
9.5.2	Sequence Composition .....	283
9.5.3	Secondary Structure .....	283
9.6	Protein Sequence Analysis .....	283
9.6.1	Primary Structure .....	283
9.6.2	Secondary Structure .....	283
9.6.3	Transmembrane Helices .....	284
9.6.4	Helical Wheel .....	284
9.7	Motif Search .....	284
9.7.1	SMART Search .....	284
9.7.2	MEME Search .....	284
9.7.3	HMM Search .....	285
9.7.4	Sequence Logo .....	285
9.8	Phylogeny .....	285
9.8.1	Protein .....	285
9.8.2	DNA .....	286

9.9	Projects .....	286
9.9.1	Sequence, Structure, and Function Analysis of the Bar-Headed Goose Hemoglobin.....	286
9.9.2	Exercises .....	287
9.10	Literature .....	287
9.10.1	Courses and Tutorials .....	287
9.10.2	Scientific Stories .....	288
9.10.3	Free Journals and Books .....	288
9.11	Bioinformatics Databases .....	289
9.11.1	List of Databases .....	289
9.11.2	Database Query Systems.....	289
9.11.3	Genome Databases .....	290
9.11.4	Sequence Databases.....	291
9.11.5	Protein Domain, Family, and Function Databases .....	292
9.11.6	Structure Databases .....	293
9.12	Bioinformatics Tools.....	294
9.12.1	List of Bioinformatics Tools at International Bioinformatics Centers.....	295
9.12.2	Web-Based Bioinformatics Platforms .....	295
9.12.3	Bioinformatics Packages to be Downloaded and Installed Locally .....	295
9.13	Sequence Analysis .....	296
9.13.1	Dotplot .....	296
9.13.2	Pairwise Sequence Alignment .....	296
9.13.3	Multiple Sequence Alignment .....	296
9.13.4	Motif Finding .....	297
9.13.5	Gene Identification .....	297
9.13.6	Sequence Logo .....	297
9.13.7	RNA Secondary Structure Prediction .....	297
9.14	Database Search.....	298
9.14.1	BLAST Search .....	298
9.14.2	Other Database Search .....	298
9.15	Molecular Modeling .....	299
9.15.1	Visualization and Modeling Tools.....	299
9.15.2	Protein Modeling Web Servers .....	300
9.16	Phylogenetic Analysis and Tree Construction .....	300
9.16.1	List of Phylogeny Programs .....	300
9.16.2	Online Phylogeny Servers .....	300
9.16.3	Phylogeny Programs .....	301
9.16.4	Display of Phylogenetic Trees .....	301
	References.....	301

<b>10</b>	<b>Foundations for the Study of Structure and Function of Proteins</b> ....	303
	Zhirong Sun	
10.1	Introduction .....	303
10.1.1	Importance of Protein .....	303
10.1.2	Amino Acids, Peptides, and Proteins.....	304
10.1.3	Some Noticeable Problems .....	306
10.2	Basic Concept of Protein Structure .....	306
10.2.1	Different Levels of Protein Structures.....	306
10.2.2	Acting Force to Sustain and Stabilize the High-Dimensional Structure of Protein .....	308
10.3	Fundamental of Macromolecules Structures and Functions .....	310
10.3.1	Different Levels of Protein Structure.....	310
10.3.2	Primary Structure.....	311
10.3.3	Secondary Structure .....	312
10.3.4	Supersecondary Structure .....	314
10.3.5	Folds .....	319
10.3.6	Summary .....	321
10.4	Basis of Protein Structure and Function Prediction .....	322
10.4.1	Overview .....	322
10.4.2	The Significance of Protein Structure Prediction .....	322
10.4.3	The Field of Machine Learning.....	323
10.4.4	Homological Protein Structure Prediction Method .....	331
10.4.5	Ab Initio Prediction Method .....	334
	Reference .....	336
<b>11</b>	<b>Computational Systems Biology Approaches for Deciphering Traditional Chinese Medicine</b> .....	337
	Shao Li and Le Lu	
11.1	Introduction .....	337
11.2	Disease-Related Network.....	338
11.2.1	From a Gene List to Pathway and Network.....	338
11.2.2	Construction of Disease-Related Network .....	340
11.2.3	Biological Network Modularity and Phenotype Network.....	346
11.3	TCM ZHENG-Related Network .....	349
11.3.1	“ZHENG” in TCM .....	350
11.3.2	A CSB-Based Case Study for TCM ZHENG.....	352
11.4	Network-Based Study for TCM “Fu Fang”.....	358
11.4.1	Systems Biology in Drug Discovery .....	358
11.4.2	Network-Based Drug Design .....	359
11.4.3	Progresses in Herbal Medicine .....	360
11.4.4	TCM <i>Fu Fang</i> (Herbal Formula) .....	361
11.4.5	A Network-Based Case Study for TCM <i>Fu Fang</i> .....	361
	References.....	364