Shiyi Shen

Jack A. Tuszynski

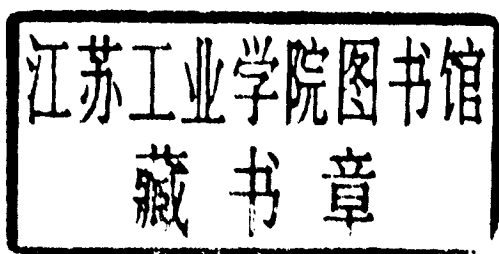# Theory and Mathematical Methods for Bioinformatics

Springer

Shiyi Shen · Jack A. Tuszynski

# Theory and Mathematical Methods for Bioinformatics

With 47 Figures and 59 Tables

Springer

Prof. Shiyi Shen
Nankai University
College of Mathematical Sciences
Tianjin 300071
China
syshen@nankai.edu.cn

Prof. Jack A. Tuszynski
University of Alberta
Department of Physics
Edmonton T6G 2G7, Alberta
Canada
jtus@phys.ualberta.ca

# BIOLOGICAL AND MEDICAL PHYSICS,
# BIOMEDICAL ENGINEERING

# BIOLOGICAL AND MEDICAL PHYSICS, BIOMEDICAL ENGINEERING

The fields of biological and medical physics and biomedical engineering are broad, multidisciplinary and dynamic. They lie at the crossroads of frontier research in physics, biology, chemistry, and medicine. The Biological and Medical Physics, Biomedical Engineering Series is intended to be comprehensive, covering a broad range of topics important to the study of the physical, chemical and biological sciences. Its goal is to provide scientists and engineers with textbooks, monographs, and reference works to address the growing need for information.

Books in the series emphasize established and emergent areas of science including molecular, membrane, and mathematical biophysics; photosynthetic energy harvesting and conversion; information processing; physical principles of genetics; sensory communications; automata networks, neural networks, and cellular automata. Equally important will be coverage of applied aspects of biological and medical physics and biomedical engineering such as molecular electronic components and devices, biosensors, medicine, imaging, physical principles of renewable energy production, advanced prostheses, and environmental control and engineering.

# Preface

Bioinformatics is an interdisciplinary science which involves molecular biology, molecular chemistry, physics, mathematics, computational sciences, etc. Most of the books on biomathematics published within the past ten years have consisted of collections of standard bioinformatics problems and informational methods, and focus mainly on the logistics of implementing and making use of various websites, databases, software packages and serving platforms. While these types of books do introduce some mathematical and computational methods alongside the software packages, they are lacking in a systematic and professional treatment of the mathematics behind these methods.

It is significant in the field of bioinformatics that not only is the amount of data increasing exponentially, but collaboration is also both widening and deepening among biologists, chemists, physicists, mathematicians, and computer scientists. The sheer volume of problems and databases requires researchers to continually develop software packages in order to process the huge amounts of data, utilizing the latest mathematical methods. The intent of this book is to provide a professional and in-depth treatment of the mathematical topics necessary in the study of bioinformatics.

Although there has been great progress in bioinformatics research, many difficult problems are still to be solved. Some of the most well-known include: multiple sequence alignment, prediction of 3D structures and recognition of the eukaryote genes. Since the Human Genome Project (HGP) was developed, the problems of the network structures of the genomes and proteomes, as well as regulation of the cellular systems are of great interest. Although there is still much work to be done before these problems are solved, it is our hope that the key to solving these problems lies in an increased collaboration among the different disciplines to make use of the mathematical methods in bioinformatics.

This book is divided into two parts: the first part introduces the mutation and alignment theory of sequences, which includes the general models and theory to describe the structure of mutation and alignment, the fast algorithms for pairwise and multiple alignment, as well as discussion based on

the output given by fast multiple alignment. Part I contains a fairly advanced treatment, and it demonstrates how mathematics may be successfully used in bioinformatics. The success achieved using fast algorithms of multiple alignment illustrates the important role of mathematics.

Part II analyzes the protein structures, which includes the semantic and cluster analysis based on the primary structure and the analysis of the 3D structure for main chains and side chains of proteins. The wiggly angle (dihedral angle) was used when analyzing the configuration of proteins, making the description of the configuration more exact. Analyzing the configuration differs from predicting the secondary or 3D structures. We collect all pockets, grooves, and channels in a protein as configuration characteristics, analyze the structure of these characteristics, and give the algorithms to compute them.

Parts I and II offer independent treatments of biology and mathematics. This division is convenient, as the reader may study both separately. In each part we include results and references from our own research experiences. We propose some novel concepts, for example, the modulus structures, alignment space, semantic analysis for protein sequences, and the geometrical methods to compute configuration characteristics of proteins, etc. Study of these concepts is still in its infancy and so there is much to still be explored. It is our hope that these issues continue to be examined mathematically so that they remain at the forefront, both in mathematics and bioinformatics.

We recognize the importance of considering the computational aspect while introducing mathematical theories. A collection of computational examples have been included in this book so that our theoretical results may be tested, and so that the reader may see the corresponding theories illustrated. Additionally, some of these examples have implications which may be applied to biology directly, and may be downloaded from the website [99] as the data is too large to include in this book. (As an example, when examining the alignment of the HIV gene, the size is $m = 405, n = 10{,}000$ bp.)

An understanding of the fundamentals of probability, statistics, combinatorial graph theory, and molecular biology is assumed, as well as programming ability. In order that the reader may solidify their understanding, problems have been added at the end of each chapter, with hints to help get started, It is our hope that this book will be useful in the field of bioinformatics both as a textbook and as a reference book.

# Acknowledgements

Ruan and his graduate students provided the initial translation of the Chinese edition into non-native-speaker English, and later checked the native-speaker English edition to confirm that the meaning of the original was retained. These students are Guangyue Hu, Tuo Zhang, Guoxiang Lu, Jianzhao Gao, Yubing Shang, Shengqi Li, and Hanzhe Chen. The manuscript was rendered into native-speaker English by Jack Tuszynski's assistant, Michelle Hanlon. The authors acknowledge and are grateful to all involved for the hard work that went into this project.

Both authors would like to thank Angela Lahee of Springer-Verlag for her guidance and encouragement.


Tianjin, China                                                  *Shiyi Shen*
Edmonton, Canada                                      *Jack A. Tuszynski*
                                                                *June 2007*

# Contents

## Part II Protein Configuration Analysis

# Outline

This book discusses several important issues including mathematical and computational methods and their applications to bioinformatics.

This book contains two parts. Part I introduces sequence mutations and the theory of data structure used in alignment. Stochastic models of alignment and algebraic theory are introduced as a means to describe data structure. Dynamic programming algorithms and statistical decision algorithms for pairwise sequence alignment, fast alignment algorithms, and the analysis and application of multiple sequence alignment output are also introduced. Part II includes the introduction of frequency analysis, cluster analysis and semantic analysis of the primary sequences, the introduction of the 3D-structure analysis of the main chains and the side chains, and the introduction of configuration analysis.

Many mathematical theories are presented in this monograph, such as stochastic analysis, combinational graph theory, geometry and algebra, informatics, intelligent computation, etc. A large number of algorithms and corresponding software packages make use of these theories, which have been developed to deal with the large amounts of biological and clinical data that play such an important role in the fields of bioinformatics, biomedicine, and so on.

This book has three main goals. The first is to introduce these classical mathematical theories and methods within the context of the current state of mathematics, as they are used in the fields of molecular biology and bioinformatics. The second is to discuss the potential mathematical requirements in the study of molecular biology and bioinformatics, which will drive the development of new theories and methods. Our third goal is to propose a framework within which bioinformatics may be combined with mathematics.

Within each chapter, we have included results and references from our own research experience, to illustrate our points. It is our hope that this book will be useful as a textbook for both undergraduate and graduate students, as well as a reference for teachers or researchers in the field of bioinformatics or mathematics, or those interested in understanding the relation between the two.