

王槐春 编著

重序列与核苷酸基础 自读与分析



人民军医出版社



蛋白质与核酸序列分析基础

Essentials of Sequence Analysis

王槐春 编著

人民军医出版社

1994 · 北京 ·

内容提要

蛋白质与核酸序列的计算机分析是分子生物学研究的一项新技术，并逐渐形成一门由分子生物学和计算机信息处理技术相结合的交叉学科——生物信息学。本书介绍了蛋白质和核酸数据库，各种基本的序列分析方法和应用程序，并用大量的图表和实例加以说明，力图使不熟悉计算机的分子生物学家容易接受。本书可供生物化学和分子生物学科研人员及生物技术软件设计人员参考，也可作为大学生物物理、生物化学专业的研究生和高年级本科生的参考书。

责任编辑 苗 芳

*

图书在版编目(CIP)数据

蛋白质与核酸序列分析基础 / 王槐春编著. —北京:人民军医出版社, 1994.12

ISBN 7-80020-424-3

I. 蛋... II. 王... III. ①蛋白质—序列—分析方法 ②核酸—序列—分析方法 IV. ①Q51-34②Q52-34

中国版本图书馆 CIP 数据核字 (94) 第 11968 号

人民军医出版社

(北京复兴路 22 号甲 3 号 邮政编码: 100842)

中国人民解放军第 1201 工厂印刷

新华书店总店科技发行所发行

开本: 850×1168mm 1/32 印张: 8 字数 204 千字

1994 年 10 月第 1 版 1994 年 10 月第 1 次印刷

印数: 1~2000 册 定价: 10.00 元

ISBN 7-80020-424-3 / R · 365

PREFACE

No aspect of modern biology has provided more understanding about life on Earth than the determination of macromolecular sequences. Indeed, the past forty years have witnessed unparalleled advances in our understanding of not only how organisms are related, but, more to the point, how they have come to possess the protein equipment with which they face the world. Thus, the development of life on Earth has largely been the result of a vast succession of gene duplications followed by waves of smaller scale changes in the form of base replacements. The reconstruction of these past events is being made possible by the computer analysis of macromolecular sequences. Such analyses demand not only modern computers, programs and databases, but also informed biologists. It is in the last-named realm that this book makes its important contribution. It is not that the technical aspects of this field are so difficult. Rather, it is simply that they are not familiar to most molecular biologists. It is an area that requires some elementary statistics, on the one hand, and some general background in macromolecular structure, on the other. The author has provided both components in a simple and straightforward way.

It is unusual to be asked to write a preface for a book written in a language that one cannot read. In the case at hand, I have only been able to read a translation of an outline of the contents, but it has been sufficient to convince me that this is a timely and up-to-date volume. Huaichun Wang has gathered together an

impressive assortment of topics. Chinese molecular biologists are fortunate indeed to have these very recent methods gathered into a single volume in their own language. The rapid pace of the biological sciences demands that we all have a basic understanding of these phenomena.

Russell F. Doolittle

序

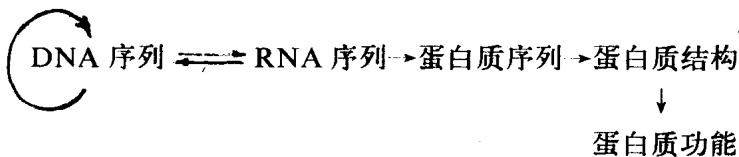
现代生物学的任何一个方面都不能象确定大分子序列那样能提供给我们更多的有关地球上生命的知识。确实，在过去的40年里，我们已经目睹了人类不仅在认识生物体是如何相关的，而且，更确切地说，在了解它们是何以获得使它们面对这个世界的蛋白质方面，已经取得了前所未有的进展。因此，地球上生命的发生大体上是一系列的基因复制伴随着以碱基替换为形式的较小规模的改变。可以用计算机对大分子序列的分析来重建这些过去的事情。这样的分析要求不仅有现代的计算机设备、程序和数据库，而且要有在行的生物学家。正是在后一个意义上，本书作出了它的重要贡献。问题不在于本领域的技术方面是如何的困难，而只是多数分子生物学家对它们不熟悉。它是这样一个领域，一方面需要一些基本的统计学方法，另一方面需要某些大分子结构的背景知识。本书作者以简单和直接的方式提供了这两方面的知识。

要求一个人为一本用他所不懂的语言写作的书来写序言是不寻常的。目前在手边的，我仅能阅读本书目录概要的翻译，但是它足以使我相信这是一本及时的和最新的卷本。王槐春把各类有趣的论题综合在一起。中国的分子生物学家确实是幸运的，能读到一本用他们自己的语言写作的收集了最新方法的书。生命科学的快速发展需要我们都要具备对这些现象的基本认识。

R.F.杜利特尔
加州大学圣迭戈分校
生物学和化学教授

前　　言

运用计算机进行蛋白质和核酸的序列分析是分子生物学研究的一个较新的发展，该项技术已越来越多地用于研究新测定的生物大分子序列或大量积累的序列数据。今天可以说，没有一篇发表关于序列的分子生物学研究论文不具有诸如此类的语句：经检索核酸或蛋白质序列数据库，“该序列与数据库中的某某序列同源”，或者“该序列与序列库中的任何序列未见明显类似”，等等。然而，序列分析不单是序列的类似性检索或比较，它还包含很多内容：二级结构预测、亲疏水性分析、序列模式识别、结构域或编码区识别以及蛋白质家族分类和进化树的构建。因此，大分子序列数据是分子生物学研究的一个重要资源，序列分析已成为当代分子生物学家必须掌握的技术。事实上，有关序列分析方法及其应用的研究已形成一个新的交叉学科——生物信息学(bioinformatics)或计算分子生物学(computational molecular biology)。它的基本出发点是通过计算方法由DNA和蛋白质的序列推导出它们的结构和功能：



序列分析不同于常规的分子生物学方法(如蛋白质分离纯化、分子克隆和序列测定)，它要求研究者具备两方面的知识，其一是基本的数学和统计学方法以及计算机程序应用的知识，其二是分子生物学特别是生物大分子结构的背景知识。由于这第一个因素，许多分子生物学家望而生畏；由于第二个因素，专业计算机软件人员和数学、统计学家难以入门，从而阻碍了序列分析程序的开发和方法的应用。在我国，虽然已经有许多有关分子生

生物学实验技术的书籍，但是迄今为止，尚未见到一本关于核酸和蛋白质序列的计算机分析的专著。为了填补这一空白，本人在多年从事生物大分子数据库和软件开发的基础上，编写了这本《蛋白质与核酸序列分析基础》，期望能促进分子生物学家熟悉和掌握序列分析技术，也为计算机软件人员开发分子生物学软件提供指南。

全书共分七章。第1章介绍序列数据库，第2和第3章介绍序列对数据库的类似性检索、序列比较和同源性构建，第4章描述蛋白质家族、分子钟和进化树的概念和研究方法，第5和第6章分别介绍蛋白质和核酸的信号检索、结构预测和功能分析，第7章介绍利用电子网络(Internet)进行序列分析。本书的最后有9个附录，它们收集了许多有用的信息，包括核苷酸碱基和氨基酸的标准符号、密码子字典、序列数据库和序列分析软件、计算机硬件的选择等。

本书力图从一个生物学家的角度来描述各种基本的序列分析方法和程序，并提供尽量多的图表和实例，而不是过多地去叙述算法和公式。其目的是使读者知道如何分析一个新的DNA或蛋白质序列，如何利用现有的数据库，有哪些程序和方法可以使用，它们能得到哪些结论，这些方法还有哪些限制或不足之处。毫不讳言，本书在较大程度上类似于Doolittle1987年出版的《OF URFS AND ORFS - A Primer on How to Analyze Derived Amino Acid Sequences》和Gribskov与Devereux1990年发表的《Sequence Analysis Primer》，因为它们与本书一样都是序列分析入门的书籍。但是本书也包含了一些新近发展起来的方法，如氨基酸组成分析及其用于纯化蛋白质的鉴定、蛋白质稳定性预测、寡核苷酸引物和探针的检索；另外，本书尽可能地收进我国学者在序列分析方面所开展的工作。

我要特别感谢美国加州大学R. F. Doolittle教授为本书所写的序言，以及允许我使用他的上述著作中的图表资料。作为序列

数据库和序列分析研究的先驱者之一，他用序列比较方法发现了许多序列同源和进化关系，特别是 1983 年他发现了病毒癌基因 v-sis 产物与人血小板源性生长因子同源。他于 1990 年为著名的《酶学方法》系列丛书编辑了第 183 卷《分子进化：蛋白质与核酸序列的计算机分析》。读者将从本书中了解到 Doolittle 的更多的贡献。我还要感谢美国 Brookhaven 国家实验室的 T.Koetzle、美国 NIH 全国生物技术信息中心的 D.Benson、瑞士日内瓦大学的 A.Bairoch、日本东京理科大学的 A.Tsugita 和意大利热那亚国家癌症研究所的 P.Romano 等向我提供了 PDB、Entrez: Sequences、EMBL、ATLAS 和 MPDB 等核酸和蛋白质数据库。军事医学科学院基础医学研究所唐佩弦研究员为本书题写书名，该院情报研究所廖应昌、缪其宏研究员和中国科学院生物物理研究所徐军同志对本书的出版给予很多帮助，在此一并表示感谢。

王槐春

1994 年 4 月 30 日

目 录

第 1 章 序列数据库	(1)
1.1 发展历史	(1)
1.2 数据收集和数据发行	(3)
1.2.1 数据收集	(3)
1.2.2 数据发行	(4)
1.3 序列检索与序列数据格式	(7)
1.4 序列数据库与序列分析软件	(12)
1.5 其它分子生物学数据库	(18)
小结	(20)
第 2 章 序列对数据库的类似性检索	(22)
2.1 序列检索和比较应考虑的问题	(23)
2.1.1 氨基酸的替换与记分方法	(23)
2.1.2 空位罚分	(26)
2.1.3 基因重组和序列内部重复	(28)
2.2 数据库检索的算法	(29)
2.2.1 片段覆盖法	(29)
2.2.2 K-tuple 匹配法	(30)
2.3 应用举例: 执行一次数据库检索	(30)
2.4 检索短序列	(34)
2.5 序列类似性检索的发现	(36)
小结	(38)
第 3 章 序列比较与同源性	(39)
3.1 双重序列比较	(40)
3.1.1 动态程序对准	(40)

3.1.2 点矩阵作图	(45)
3.1.3 序列同源与类似	(58)
3.2 序列类似的显著性评价.....	(59)
3.2.1 Monte Carlo 模拟法	(59)
3.2.2 类似性积分的分布状况	(63)
3.2.3 Karlin-Altschul 概率法.....	(66)
3.2.4 经验判断法	(66)
3.2.5 短序列对准的显著性评价	(68)
3.2.6 核酸序列比较的显著性评价	(70)
3.3 多序列比较.....	(72)
3.3.1 Feng 和 Doolittle 的方法	(73)
3.3.2 Schuler, Altschul 和 Lipman 的方法	(77)
小结	(82)
第 4 章 蛋白质家族与分子进化	(84)
4.1 蛋白质家族和超家族.....	(84)
4.2 分子进化钟.....	(88)
4.3 进化树.....	(93)
4.3.1 序列进化树	(93)
4.3.2 结构进化树	(97)
小结	(105)
第 5 章 蛋白质结构预测与功能分析	(107)
5.1 物理化学特性	(107)
5.1.1 由蛋白质序列推算蛋白质的分子量、 氨基酸组成、等电点和消光系数	(107)
5.1.2 氨基酸组成的统计分析	(110)
5.1.3 根据蛋白质的理化特性快速鉴定蛋白质 ...	(120)
5.1.4 蛋白质稳定性预测	(120)

5.2 一级结构特征	(126)
5.2.1 内部重复单位	(126)
5.2.2 序列模式和位点	(129)
5.2.3 序列结构域与模式匹配方法	(134)
5.3 二级结构特征	(148)
5.3.1 亲疏水性特征与作图	(148)
5.3.2 变异矩与螺旋外表面的检测	(157)
5.3.3 二级结构预测	(158)
小结	(167)
 第 6 章 核酸的信号检索与结构预测	(169)
6.1 限制性酶切位点和固定序列模式	(170)
6.1.1 限制性酶切位点检索	(170)
6.1.2 固定序列模式的检索	(172)
6.2 核酸序列的特殊信号检索	(175)
6.2.1 细菌启动子	(177)
6.2.2 mRNA 剪接位点	(178)
6.3 基因编码区鉴定与翻译	(180)
6.3.1 开放阅读框架分析	(180)
6.3.2 编码区鉴定	(181)
6.3.3 基因序列翻译	(191)
6.4 RNA 二级结构预测	(194)
6.5 寡核苷酸探针和引物设计	(199)
6.5.1 分子探针数据库	(199)
6.5.2 计算机辅助的引物和探针设计	(201)
小结	(205)
 第 7 章 电子网络与序列分析	(208)
7.1 使计算机联网	(209)

7.2 通过电子邮件开展序列分析	(209)
7.2.1 电子邮件与电子地址	(209)
7.2.2 电子邮件服务器	(212)
7.3 文件传送协议	(214)
7.4 网络信息工具	(215)
7.5 电子公告牌	(217)
小结	(219)
 参考文献	(220)
 附录 I 核苷酸符号	(230)
附录 II 遗传密码字典	(231)
附录 III 氨基酸符号	(232)
附录 IV 蛋白质的“平均”组成	(233)
附录 V Dayhoff 突变数值矩阵(PAM 250)	(234)
附录 VI 一些基因组的大小	(235)
附录 VII 序列和结构数据库	(236)
附录 VIII 硬件和软件的选择	(239)
附录 IX 军事医学科学院生物技术软件库	(243)

Content

Preface by R. F. Doolittle

Forword

Chapter 1 Sequence Databases.....	(1)
1.1 Historical Overview	(1)
1.2 Data Collection and Distribution	(3)
1.2.1 Data Collection	(3)
1.2.1 Data Distribution.....	(4)
1.3 Sequence Retrieval and Database Format	(7)
1.4 Sequence Database and Sequence Analysis Software	(12)
1.5 Other Molecular Biology Databases	(18)
Summary	(20)
Chapter 2 Searching through Database	(22)
2.1 Basic Principals.....	(23)
2.1.1 Scoring System.....	(23)
2.1.2 Gap Panelty	(26)
2.1.3 Gene Rearrangement and Internal Sequence Repeat	(28)
2.2 Algorithm for Fast Database Search	(29)
2.3 Example: Making a Database Search	(30)
2.4 Searching for Short Sequence.....	(34)
2.5 Discoveries of Computer Searching	(36)
Summary	(38)

Chapter 3 Sequence Comparison and Homology	(39)
3.1 Pairwise Sequence Comparison	(40)
3.1.1 Dynamic Programming Method	(40)
3.1.2 Dot Matrix Method	(45)
3.2 Significance Evaluation	(59)
3.2.1 Monte Carlo Simulation.....	(59)
3.2.2 Distribution of Similarity Scores	(63)
3.2.3 Karlin– Altschul Probability	(66)
3.2.4 Doolittle's NAS Curve	(66)
3.2.5 Significance of Short Sequence Alignment	(68)
3.2.6 Significance of Nucleic Acid Sequence Alignment.....	(70)
3.3 Multiple Sequence Alignment.....	(72)
3.3.1 The Method of Feng and Doolittle	(73)
3.3.2 The Method of Schuler, Altschul and Lipman	(77)
Summary	(82)
 Chapter 4 Protein Families and Molecular Evolution	(84)
4.1 Protein Families and Superfamilies	(84)
4.2 Molecular Clock	(88)
4.3 Evolutionary Tree	(93)
4.3.1 Sequential Evolutionary Tree.....	(93)
4.3.2 Structural Evolutionary Tree.....	(97)
Summary	(105)
 Chapter 5 Proteins: Structure Prediction and Function Analysis	(107)

5.1 Physicochemical Properties	(107)
5.1.1 Derivation of Molecular Weight, Amino Acid Composition, Isoelectric Point and Extinction Coefficient from a Protein Sequence	(107)
5.1.2 Statistical Analysis of Amino Acid Composition	(110)
5.1.3 Rapid Protein Identification through its Physicochemical Properties	(120)
5.1.4 Protein Stability Prediction	(120)
5.2 Primary Structure Properties	(126)
5.2.1 Internal Repeating Unit	(126)
5.2.2 Sequence Pattern and Site	(129)
5.2.3 Sequential Domain and Pattern Matching Method	(134)
5.3 Secondary Structure Properties	(148)
5.3.1 Hydropathy Plot	(148)
5.3.2 Variation Moment and Detection of the Outside Face of Helices	(157)
5.3.3 Secondary Structure Prediction	(158)
Summary	(167)

Chapter 6 Nucleic Acids: Signal Search and Structure Prediction	(169)
6.1 Restriction Endonuclease Site and Fixed Sequence Pattern	(170)
6.1.1 Restriction Enzyme Database and Site Search	(170)
6.1.2 Fixed Sequence Pattern Search	(172)
6.2 Pattern Matching Analysis for Special	

Signal Search	(175)
6.2.1 Bacterial Promoter	(177)
6.2.2 mRNA Splicing Site	(178)
6.3 Coding Region Identification and Sequence	
Translation	(180)
6.3.1 Open Reading Frame Analysis	(180)
6.3.2 Coding Region Identification	(181)
6.3.3 Sequence Translation	(191)
6.4 RNA Secondary Structure Prediction	(194)
6.5 Oligonucleotide Probe and Primer Design	(199)
6.5.1 Molecular Probe Data Base (MPDB)	(199)
6.5.2 Primer and Probe Design	(201)
Summary	(205)
 Chapter 7 Electronic Network and Sequence Analysis	(208)
7.1 Networking Computer	(209)
7.2 Sequence Analysis by Electronic Mail	(209)
7.2.1 E-mail and Electronic Address	(209)
7.2.2 E-mail Server	(212)
7.3 File Transfer Protocol	(214)
7.4 Network Information Tools	(215)
7.5 Electronic Bulletin Board	(217)
Summary	(219)
 References	(220)
 Appendices	
I Nucleic Acid Codes	(230)
II Dictionary of Genetic Codes	(231)