

数学与现代科学技术丛书 2

生物计算

——生物序列的分析方法与应用

杨晶 胡刚 王奎 沈世镒 编著



 科学出版社
www.sciencep.com

生物计算

生物计算

——生物科学与技术方法与实践

李海、周晓、王伟、周长海、周晓



数学与现代科学技术丛书 2

-21

生物计算

——生物序列的分析方法与应用

杨 晶 胡 刚 王 奎 沈世镒 编著

Q811.4

1202

科学出版社

北京

内 容 简 介

本书介绍生物计算中的几种主要方法，如序列比对、系统发育分析、蛋白质序列的语义分析与结构预测、基因识别与生物芯片的数据分析等，给出它们的基本问题与有关的方法及应用。全书由三部分组成。第一部分介绍这些问题的由来与主要内容，给出它们的基本原理、计算与分析方法及应用意义，同时介绍一些国际上较为通用的软件包。第二部分是生物学备忘录，介绍有关生物学的基础知识。第三部分是数学备忘录，介绍与这些生物计算有关的数学理论与方法。

本书可作为数学、生物、医学、化学等专业的本科生或研究生教材，其中第一部分内容可作为各专业的公共部分，而第二、三部分内容可供各专业适当选用。

图书在版编目(CIP)数据

生物计算：生物序列的分析方法与应用/杨晶等编著。—北京：科学出版社，
2010

(数学与现代科学技术丛书；2)

ISBN 978-7-03-026393-3

I. 生… II. 杨… III. 生物信息论-研究 IV. Q811.4

中国版本图书馆 CIP 数据核字 (2010) 第 007551 号

责任编辑：赵彦超 / 责任校对：鲁 素

责任印制：钱玉芬 / 封面设计：王 浩

科学出版社出版

北京东黄城根北街 16 号

邮政编码：100717

<http://www.sciencep.com>

双青印刷厂印刷

科学出版社发行 各地新华书店经销

*

2010 年 3 月第 一 版 开本：B5(720×1000)

2010 年 3 月第一次印刷 印张：23 3/4

印数：1—3 000 字数：456 000

定价：69.00 元

(如有印装质量问题，我社负责调换〈双青〉)

《数学与现代科学技术丛书》序

当代数学在向纵深发展的同时，被空前广泛地应用于几乎一切领域。一方面，它与其他学科交汇，形成了许许多多交叉学科（例如，信息科学、计算机科学、系统科学、数学物理、数学化学、生物数学、数学语言学、数量经济学、金融数学、复杂性科学、科学计算等）；另一方面，它又被应用于高新技术（如信息安全、信息传输、图像处理、语音识别、网络、海量数据处理、网页搜索、遥测遥感、交通管理、医疗诊断、手术方案、药物检验、商业广告等方面）的开发，成为一些高新技术的核心。应用数学的这种发展趋势急剧地扩展了数学的疆界，也深刻地改变了数学的面貌。

中国的经济正在迅猛发展，其中的科技含量也与日俱增。为了提高自主创新能力，我国已经有不少数学工作者投身于这类应用数学的研究中，还有更多的数学工作者则正在密切关注这方面的进展，看好它的前景。愈来愈多的人希望了解这类应用数学的现状，寻找入门之径。

《数学与现代科学技术丛书》是力图反映这个发展趋势的一套应用数学丛书，它将较全面地向我国读者介绍当今数学在现代科学技术各个领域中应用的状况，通过必要的准备知识，逐步把读者引向相关的研究前沿。

从事交叉学科研究和高新技术开发的应用数学家，除了要精通所需的数学知识外，还必须深入了解其所研究问题的来龙去脉。“建模”是应用数学研究实际问题的关键。这也是一门数学艺术：从复杂的实际问题中抽象出关键的“量的关系”，它既能反映出问题的基本特征，又能用现阶段的数学工具对其加以处理。有鉴于此，这套丛书的一个特点就是不但要介绍有关的数学理论和方法，还必须介绍问题的来源与背景、数学建模以及如何运用数学工具来解决实际问题。

本丛书适用于数学及相关专业的大学生和研究生，以及与数学有关的各专业科技工作者。

张恭庆

2009年11月

前　　言

生物计算中的理论、方法与应用越来越被生物、医学及其他医务工作者所需要与关注，特别是在人类基因组计划实施以来，该学科的发展与研究更凸显出重要的作用。基因、基因组、蛋白质、蛋白质组等生物学信息的数据采集、储存与分析及其生物学意义，是生物计算乃至生物、医学与医药的重点研究内容之一。因此在国内外的许多医科大学均被作为重要课程，与生物信息学和生物计算相关内容的课程不仅是研究生的必修课程或选修课程，也是多个专业本科生的专业必修课程或选修课程。我们先后用了近三年的时间，在开展教学和研究工作的同时编写了本书，目的是为生物学和医学相关专业的本科生与研究生提供一本既通俗易懂，同时又可深入了解相关内容的教材，为该学科的建设与发展服务。

自 2004 年以来，本人有幸多次参加南开大学数学科学学院沈世镒教授主持的“生物信息学”讨论班。在讨论与学习过程中，不仅掌握了一些解决生物序列分析与计算的具体算法，更重要的是学到了解决生物序列分析的一些新方法和新思想。如生物序列的多种比对算法、数据结构中的语义分析及其在蛋白质结构分析中的应用等。这些方法从不同角度对生物计算中的有关问题进行研究与探讨，并在许多方面得到了很好的应用。在学习过程中，与南开大学数学科学学院胡刚、王奎博士等合作，对生物计算中的算法以及相关软件包的使用等问题有了更深入与确切的理解，使本书得以顺利完成。我们希望能将该领域中的主要内容与方法介绍给读者。

“生物计算”与“生物信息学”在本质上无大的区别，国内外的许多院校均把它们看作同一领域的学科。在本书中，我们把“生物计算”看作较偏重于原理与方法，同时注重它们的实现与应用，在介绍国外先进与常用算法的同时，增加了相应软件包的使用与分析等内容。

本书共分为三个部分，第一部分是本书的重点和核心，主要介绍生物信息学中的常用算法与软件，包括它们的基本原理、方法与应用，并介绍国际上常用的计算机软件包，对其性能、使用方法与具体实施加以介绍，以方便读者更快和更有效地实现与应用这些算法。该部分内容尽量简练并具有实用性与可操作性。为便于更多的读者阅读此书，本书第二与第三部分分别介绍了有关的生物学和数学的基础知识。第二部分介绍了与生物计算相关的生物学基础知识，是为具有数学与计算机等专业背景的读者准备的，已经具有生物学与医学等专业背景的读者可以将这部分内容略去不读。第三部分介绍与生物计算相关的数学理论与方法，因为它涉及多个数学学科领域，内容较为广泛，对理解与掌握生物计算的基本方法是有益的，也可供

需要作更深入研究(如自行设计或编写算法程序与软件等)的读者参考。

本书的编写是在沈世镒教授的建议、组织与策划下由多人编撰完成,第5~10章由杨晶编撰,第2~4,13章由胡刚编撰,第1,11,12,14章由王奎编撰,沈世镒审校全书。书中介绍的一些新算法,如序列比对的SPA算法与蛋白质一级结构的语义分析等,是南开大学生物信息学研究组取得的新成果,书中未详细介绍的部分结果可在参考文献[174, 228]中找到。另外,特别感谢天津医科大学生物医学工程系田心教授在科研和教学工作中给予的精心指导和帮助。同时感谢天津医科大学内分泌研究所郭刚教授、张镜宇教授和基础医学院解用虹教授等在分子生物学方面的帮助和指导。感谢天津理工大学计算机学院宋金杰教授及其学生周新博的帮助。

本书的出版获国家自然科学基金项目(20836005; 30870791)、天津市自然科学基金项目(07JCDJC08100)以及天津医科大学和南开大学教务处的支持,特此表示感谢。

本书可作为医学、生物学、计算机和数学等专业本科生或研究生“生物信息学”或“生物计算学”课程的教材或参考书。不同专业对本书内容可有不同的侧重与选择。生物信息学的研究还处在起步阶段,一些理论与方法都在不断发展中,因此有关内容还会不断更新。另外,本书内容涉及医学、生物学、数学与计算机等专业,具有跨学科特征。由于我们的专业知识与工作背景的限制,书中错误或不妥之处在所难免,真诚希望读者批评指正。

杨 晶

2009年3月于天津

目 录

《数学与现代科学技术丛书》序

前言

第一部分 基本方法

第 1 章 生物序列突变与比对分析	3
1.1 生物序列突变与比对问题	3
1.1.1 生物序列的类型与结构	3
1.1.2 生物序列突变与比对问题的意义与应用	4
1.1.3 生物序列比对的原理与方法	6
1.2 二重序列比对的有关算法	9
1.2.1 关于动态规划算法的一些说明	9
1.2.2 动态规划算法	10
1.2.3 统计判决算法的基本思想	15
1.2.4 BLAST 软件的使用	16
1.3 多重序列的比对问题	19
1.3.1 MSA 的意义与概况	19
1.3.2 MSA 的定义与优化准则	21
1.4 MSA 算法与计算	22
1.4.1 MSA 算法的基本概念	22
1.4.2 MSA 的算法步骤	24
1.4.3 ClustalW 软件的使用	26
1.4.4 关于 MSA 的几点说明	30
1.4.5 几个多重序列比对应用例子	31
1.5 SPA 算法的原理与计算	32
1.5.1 SPA 算法的基本原理	32
1.5.2 SPA 算法的基本步骤	34
1.5.3 SPA 算法源码	36
1.5.4 SPA 算法的有关问题讨论	39
1.5.5 SPA 算法的一个实例计算	41
习题与思考	47
第 2 章 系统发育分析	49
2.1 分子系统发育分析的基本概念	49

2.2 基于距离的方法.....	49
2.2.1 非加权分组平均法.....	49
2.2.2 邻接法.....	52
2.3 基于特征的方法.....	55
2.4 极大似然和 Bayes 方法.....	57
2.4.1 进化的概率论模型.....	58
2.4.2 构建进化树的极大似然方法.....	60
2.4.3 构建进化树的 Bayes 方法.....	62
2.5 构建进化树软件简介.....	63
习题与思考.....	68
第 3 章 蛋白质一级结构的语义分析.....	69
3.1 蛋白质一级结构的信息与统计分析法.....	69
3.1.1 蛋白质一级结构的语义分析简介.....	69
3.1.2 信息、统计分析法的要素与要点.....	70
3.1.3 局部词的定义与判定.....	72
3.1.4 蛋白质一级结构的语义分析.....	74
3.2 蛋白质序列语义结构的组合分析法.....	80
3.2.1 关于组合图论的有关记号.....	81
3.2.2 数据库的复杂度.....	84
3.2.3 数据库的关键词与核心词.....	86
3.2.4 关于组合分析的若干应用问题.....	89
习题与思考.....	92
第 4 章 蛋白质结构预测.....	93
4.1 蛋白质二级结构预测.....	93
4.1.1 蛋白质二级结构预测的评价体系.....	93
4.1.2 Chou-Fasman 方法.....	94
4.1.3 GOR 方法.....	96
4.1.4 PHD 方法.....	98
4.2 蛋白质空间结构预测.....	100
4.2.1 同源序列搜索.....	100
4.2.2 折叠识别方法.....	101
4.2.3 从头预测方法.....	104
4.3 蛋白质结构预测软件简介.....	105
4.3.1 PHD 软件使用简介.....	105
4.3.2 使用 nnpredict 软件预测蛋白质二级结构.....	108

4.3.3 PSIPRED 软件使用简介	109
习题与思考	111
第 5 章 基因识别	112
5.1 绪论	112
5.1.1 原核基因识别	112
5.1.2 真核基因识别	113
5.1.3 常用模式基因组简介	114
5.2 基因序列特征分析	116
5.2.1 内含子与外显子	116
5.2.2 CpG 岛	117
5.2.3 密码子使用偏性	118
5.3 开放阅读框识别	119
5.3.1 开放阅读框特性	119
5.3.2 开放阅读框识别原理	121
5.3.3 开放阅读框识别软件使用	122
5.4 隐 Markov 模型基因识别方法	126
5.4.1 隐 Markov 模型	127
5.4.2 GENSCAN 隐 Markov 模型方法和原理	128
5.4.3 GENSCAN 软件使用	131
5.4.4 基因识别方法评价	134
5.5 其他基因识别方法简介	135
5.5.1 神经网络方法	135
5.5.2 Z 曲线方法	136
习题与思考	138
第 6 章 基因表达数据分析	139
6.1 基因表达序列标签数据分析简介	139
6.1.1 基因表达序列标签的概念	139
6.1.2 基因表达序列标签数据的获取	141
6.1.3 基因表达序列标签数据聚类分析	145
6.1.4 基因表达序列标签的应用	147
6.2 基因芯片数据的获取	147
6.2.1 基本概念	148
6.2.2 基因芯片实验过程	149
6.2.3 基因芯片数据获取	150
6.2.4 基因芯片数据内容	152

6.3 基因芯片数据分析	153
6.3.1 基因表达谱芯片数据标准化	154
6.3.2 基因表达谱芯片数据散点图分析	156
6.3.3 基因表达差异显著性分析	157
6.4 基因芯片数据聚类分析	159
6.4.1 基本概念	159
6.4.2 特征描述	160
6.4.3 分层聚类方法	162
6.4.4 模糊聚类方法	167
6.5 其他基因芯片数据分析方法简介	173
6.5.1 支持向量机方法	173
6.5.2 K 均值聚类	173
6.5.3 自组织映射图聚类	174
6.6 基因芯片数据分析软件简介	175
习题与思考	176

第二部分 生物学备忘录

第 7 章 核酸与 DNA	179
7.1 细胞与染色体	179
7.1.1 细胞	179
7.1.2 染色体概念	180
7.1.3 染色体特征	181
7.2 核酸分子与 DNA 结构	182
7.2.1 核酸分子	182
7.2.2 DNA 分子结构	184
7.3 RNA 结构与分类	187
7.3.1 RNA 结构	187
7.3.2 RNA 分类	188
第 8 章 氨基酸与蛋白质	190
8.1 氨基酸	190
8.1.1 氨基酸组成	190
8.1.2 氨基酸符号表示	190
8.1.3 氨基酸分类	192
8.2 肽链	193

8.3 蛋白质	194
8.3.1 蛋白质分类	194
8.3.2 蛋白质一级结构	194
8.3.3 蛋白质空间结构	195
8.3.4 蛋白质功能	196
8.3.5 蛋白质组	197
8.4 中心法则与遗传密码	197
8.4.1 中心法则	197
8.4.2 遗传密码	199
第 9 章 基因与基因组	201
9.1 基因	201
9.1.1 基本概念	201
9.1.2 基因突变	202
9.2 基因组	203
9.2.1 基本概念	203
9.2.2 人类基因组	205
9.2.3 后基因组计划	206
9.3 基因表达与调控	207
9.3.1 基本概念	207
9.3.2 原核生物基因表达与调控	209
9.3.3 真核生物基因表达与调控	211
第 10 章 生物信息数据库	213
10.1 GenBank 数据库	213
10.1.1 数据来源	213
10.1.2 数据内容与类型	213
10.1.3 序列格式	215
10.1.4 数据检索与下载	215
10.1.5 数据提交	216
10.1.6 应用实例	218
10.2 Swiss-Prot 数据库	220
10.2.1 数据来源	221
10.2.2 数据内容	221
10.2.3 序列格式	222
10.2.4 数据检索与下载	222
10.2.5 数据提交	224

10.2.6 应用实例	224
附录 1 GenBank 数据库中的核酸序列记录	228
附录 2 Swiss-Prot 数据库中的蛋白质序列记录	231
 第三部分 数学备忘录	
第 11 章 智能计算理论与算法	237
11.1 智能计算概论与感知器理论	237
11.1.1 感知器模型及其学习算法	237
11.1.2 感知器模型的推广	241
11.1.3 支持向量机	244
11.2 EM 算法	246
11.2.1 EM 算法概论	246
11.2.2 极大似然估计的 EM 算法	247
11.2.3 组合决策中的 EM 计算	250
11.3 EM 算法在其他统计问题中的应用	254
11.3.1 互熵与 Fisher 矩阵	254
11.3.2 混合分布参数估计中的 EM 算法	257
11.3.3 分布族的聚类中的 EM 算法	261
11.4 Weka 软件的使用	267
11.4.1 Weka 的基本工作环境与数据准备	267
11.4.2 Weka 的使用	269
第 12 章 概率、信息与统计	275
12.1 概率与信息	275
12.1.1 随机变量与多重随机变量	275
12.1.2 随机变量的特征数	280
12.1.3 随机变量与概率分布的信息度量	282
12.2 重要随机变量和极限定理	285
12.2.1 几种重要的随机变量及其概率分布	285
12.2.2 随机变量的极限定理	290
12.3 统计分析简介	293
12.3.1 统计分析的基本要素	293
12.3.2 参数的点估计理论	295
12.3.3 参数的区间估计理论	298
12.3.4 其他问题	299

12.4 多元统计中的几个典型问题	299
12.4.1 多元统计分析的基本数学模型	299
12.4.2 聚类分析	300
12.4.3 主成分分析与因子分析	303
12.4.4 因子分析	306
12.4.5 判别分析	307
12.5 R 统计软件包简介	309
12.5.1 R 系统初览	309
12.5.2 R 的数据读入	311
12.5.3 使用 R 做统计分析	312
第 13 章 随机过程	314
13.1 随机过程的一般理论	314
13.1.1 随机过程的基本概念	314
13.1.2 独立随机序列	315
13.1.3 Poisson 过程与可加过程	317
13.2 Markov 过程	321
13.2.1 Markov 过程的基本概念	321
13.2.2 Markov 过程的生成算子	324
13.3 隐 Markov 模型	327
13.3.1 隐 Markov 模型的基本概念	327
13.3.2 HMM 的状态估计	328
13.3.3 HMM 的 EM 学习算法	331
第 14 章 有关图与树的基本知识	334
14.1 图的基本概念与结构	334
14.1.1 图的一般定义与记号	334
14.1.2 树图与系统树	336
14.2 组合空间与 de Bruijn-Good 图	337
14.3 序列与数据库的复杂度理论	340
14.3.1 复杂度的定义	340
14.3.2 复杂度的计算算法	341
14.3.3 算法的改进	342
参考文献	344
索引	357
《数学与现代科学技术丛书》已出版书目	362

第一部分

基本方法

第1章 生物序列突变与比对分析

1.1 生物序列突变与比对问题

1.1.1 生物序列的类型与结构

1. 生物序列的定义与记号

生物序列一般指 DNA 序列、RNA 序列或蛋白质序列，它们是具有活性的生物大分子，在一定的环境条件下具有新陈代谢与特定的生物学功能，都是由大量生物小分子按一定的顺序排列而成。其中 DNA 与 RNA 序列由不同的核苷酸排列组成，而蛋白质是由不同的氨基酸排列而成。如果把核苷酸或氨基酸看作一个基本单元，那么生物序列就是由这些基本的分子单元的排列组成。

生物序列的结构有多种表示方式，最常见的是它们的一级结构与空间结构表达。所谓一级结构就是由这些基本分子单元按一定的顺序排列而成的序列，因此可用以下序列记号表示：

$$A = (a_1, a_2, \dots, a_n), \quad B = (b_1, b_2, \dots, b_m), \quad (1.1.1)$$

其中 A, B 等英文大写字母表示序列， a_i, b_i 表示在第 i 位置上的基本分子单元，它们取自某一个特定的分子集合 $V_q = \{0, 1, \dots, q-1\}$ ，当 A, B 是 DNA 或 RNA 序列时，取 $q = 4$ ，这时 $V_4 = \{a, c, g, t\}$ （或 $\{a, c, g, u\}$ ）表示四种不同的核苷酸。当 A, B 是蛋白质序列时， V_q 就是氨基酸的集合，常用的氨基酸有 20 种，这时 $q = 20$ ， V_{20} 就表示这 20 种常用的氨基酸。

在 (1.1.1) 式中，下标 n, m 分别是 A, B 的序列长度。对于一般的多重序列（序列组），记之为

$$\mathcal{A} = \{A_1, A_2, \dots, A_s, \dots, A_m\}, \quad (1.1.2)$$

其中每个 A_s 是 V_q 上的序列，记之为

$$A_s = (a_{s,1}, a_{s,2}, \dots, a_{s,n_s}), \quad s = 1, 2, \dots, m, \quad (1.1.3)$$

而 n_s 是序列 A_s 的长度，称 m 为该多重序列的重数。

2. 生物序列的类型

生物序列有多种类型，从分子结构类型来分，有核苷酸（DNA 与 RNA）序列与氨基酸序列。从数据内容来分，有一级结构与空间结构等，一级结构一般指生物大