

前言

本书是为了作为高年级本科生和研究生的生物信息学教材而编写的。书中大多数章节取材于两位作者分别在加拿大滑铁卢大学和新加坡国立大学开设的生物信息学课程的讲义。该书也可以给希望了解生物信息学和计算生物学基础和工具的科研与教学人员作为参考书。

对生物信息学有兴趣的人员往往来自于不同的领域包括数学、计算机、生命科学、甚至生物科技工业。因此，本书的阅读不需要读者具有很多分子生物学方面的知识。在每章节所需要的地方，我们根据情况插入了一些必要的生物学背景知识。我们在写作的时候也不假定读者具有很强的计算机算法和概率方面的背景，尽管对这两方面的基本了解对于阅读这本书的某些章节是很有帮助的。为此，我们使用通俗的语言详细地解释了一些重要的算法概念和概率模型。

为了给学生们提供更多的练习机会，除了在内容上有一些特殊安排外，我们还在每章的最后给出了很多练习题。大多数习题比较容易，可以作为学生作业的题目。个别较难的用星号标注。其中尚未解决的问题特别注明，可以作为研究生的研究课题。

本书共有六章。在第一章里，作者尝试使用比对图通俗地介绍生物序列比对的概念，动态规划算法，和打分矩阵。

第二章讲述有关同源搜索的各种重要的概念和快速比对算法。特别的是我们在该章最后一节系统地介绍了散核技巧。这是近年来提出的一个非常重要的设计比对快速算法的技巧。这是其它生物信息学教科书所不具有的一个特色。

在第三章里，我们主要讲述如何有效地比对多个生物序列。我们介绍了用于表示蛋白结构域序列保守性的各种概念和工具。我们还介绍了各种渐进式比对策略和常用的多序列比对程序。

第四章先介绍隐马尔科夫模型的三大基本问题和相应的算法。隐马尔科夫模型其实就是有限自动机的随机版本。在这一章，我们尝试从形式语言和自动机的角度来理解隐马尔科夫模型。我们希望这对计算机系的学生理解这一重要模型有帮助。然后，我们还讨论了它在蛋白质序列分析以及基因预测中的应用。

在第五章里，我们讲述进化树的基本数学性质、构建进化树的各种算法和两个实际问题。我们还简单地介绍了进化树应用在比较基因组学和综合基因组学 (metagenomics) 中的两个例子：推断基因复制历史和构建一个生物特性的祖先状态。

第六章介绍了计算蛋白质组学的一些基本问题以及相应的算法和统计方法。计算蛋白质组学是生物信息学中近十多年来新兴的一个重要研究领域，主要是通过计算来分析蛋白质的质谱数据，从而推断出蛋白质的序列、修饰、

定量等信息。目前这个领域发展的非常快，而介绍这个领域的教材又非常少，无论是中文还是英文都是如此，所以这一章也是本书的另一个特色。

随着高通量 DNA 测序的成本降为每个人都可负担得起的程度以及蛋白质测序工具的日益成熟，本书所讲述的知识是每个生物信息工作人员所应据备的基础知识。本书大部分章节之间是相互独立的。在讲授生物信息学课程的时候，教师可以按照各自不同需求选择内容。第一至六章的授课时间分别是大约八、八、六、六、十和十小时。

最后，我们提醒读者生物信息学是一门仍在发展的学科。当这本书和读者见面的时候，也许某些章节所讲述的内容已经有了新的突破。为此，我们也特别建立了一个网站 (<http://www.math.nus.edu.sg/~matzlx/BioTextbook>)。在这个网站上，我们会继续补充本书所涉及内容的最新进展，以供读者参考。另外，为方便制作课件，我们也给使用这本书的讲师们提供书中的所有插图。

第六章、§2.3 和 §2.7 由 BM 执笔，其余部分由 LXZ 完成。在长达两年的写作过程中，我们得到了许多朋友和学生们的帮助。在此，我们特别感谢高松、贺琳、林皓、宁康、姚玉华和徐魁。通过认真阅读这本书的初稿，他们帮助我们把它的内容更完美地呈献给读者。尽管我们作了最大的努力，本书难免还有一些由于我们疏忽造成的错误。如果发现任何错误，请你通过上述网站上的联系方式通知我们。除了在上述网站公布外，这些错误也将会在再版时加以更正。

马斌，滑铁卢
张洛欣，新加坡
2014 年 6 月

目 录

前言	iii
第一章 生物序列比对	1
1.1 DNA, RNA, 和蛋白质	1
1.1.1 DNA 分子	1
1.1.2 蛋白质分子	2
1.1.3 RNA 分子	3
1.1.4 从基因到蛋白质的信息传递	4
1.2 比对 – 序列比较的模型	4
1.3 比对图	5
1.3.1 定义	6
1.3.2 双序列之间比对的总数目	7
1.4 比对的记分法则	8
1.5 全序列比对: 动态规划算法	9
1.5.1 基本算法	9
1.5.2 使用仿射空位罚分的算法	11
1.5.3 *全序列比对的 C 语言程序	14
1.6 局部比对: Smith-Waterman 算法	18
1.6.1 Smith-Waterman 算法	18
1.6.2 *局部比对的 C 语言程序	19
1.7 最优占用空间的比对算法	21
1.8 比对蛋白质序列所使用的打分矩阵	24
1.8.1 打分的统计基础	24
1.8.2 BLOSUM 矩阵系列	25
1.9 参考文献	31
1.10 练习题	32
第二章 快速比对方法	35
2.1 同源序列查询和数据库搜索	35
2.2 序列中的字分布	36
2.2.1 DNA 序列的随机模型 I: 一致独立分布	37
2.2.2 DNA 序列的随机模型 II: 马尔科夫链	40
2.3 字匹配的散列表方法	43
2.4 点阵法	45
2.5 *FASTA 程序	47
2.6 BLAST 程序	50

2.6.1	基本算法：连续核的概念	50
2.6.2	E- 值的计算公式	51
2.6.3	BLAST 程序系列	53
2.7	散核方法	55
2.7.1	散核模型	55
2.7.2	散核的优化	59
2.7.3	基于散核的相似性查找的程序实现	61
2.7.4	多散核	62
2.7.5	*其它有关散核的研究	63
2.8	参考文献	64
2.9	练习题	66
第三章	多序列比对	69
3.1	为什么需要比对多个生物序列?	69
3.2	模体、谱、共识序列	69
3.3	Logo: 一个序列保守区域的可视化方法	72
3.4	多序列比对的 SP 分数	73
3.5	多序列比对的复杂性	75
3.5.1	动态规划算法	75
3.5.2	NP- 难解性	76
3.6	渐进式比对	79
3.6.1	渐进式的基本策略	79
3.6.2	Feng-Doolittle 比对算法	80
3.7	近似算法	82
3.7.1	序列编辑距离	82
3.7.2	星型比对算法	83
3.8	多序列比对实用程序	87
3.8.1	ClustalW	87
3.8.2	MUSCLE	90
3.8.3	其它多序列比对程序	92
3.9	*基因组的比对	93
3.10	参考文献	95
3.11	练习题	97
第四章	隐马尔科夫模型及基因序列的识别	101
4.1	隐马尔科夫模型	101
4.1.1	隐马尔科夫模型的定义	101
4.1.2	隐马尔科夫模型的基本问题	102
4.2	基本算法	103

4.2.1	前向算法和后向算法	103
4.2.2	Viterbi 算法	106
4.2.3	建模算法	107
4.3	蛋白质簇的隐马尔科夫链模型	110
4.3.1	谱 HMM	110
4.3.2	从多序列比对到谱 HMM	111
4.3.3	从谱 HMM 到多序列比对	113
4.3.4	Pfam 数据库	114
4.4	GENSCAN: 预测人基因组中的全基因结构程序	115
4.4.1	真核生物基因的结构	115
4.4.2	半 HMM	116
4.4.3	基因的 Burge-Karlin 模型	117
4.4.4	自动识别人基因组中的基因序列	120
4.5	参考文献	121
4.6	练习题	123
第五章	分子进化树分析	125
5.1	达尔文的进化树	125
5.2	进化树的数学性质	126
5.2.1	基本概念	126
5.2.2	进化树的个数	128
5.2.3	常见的无根进化树变换	130
5.2.4	进化树之间的距离	133
5.2.5	二叉树和多叉树	134
5.3	构建分子进化树 I: Parsimony 方法	135
5.3.1	Fitch 算法	138
5.3.2	寻找简约进化树	139
5.4	构建分子进化树 II: 基于距离的方法	140
5.4.1	加权进化树和距离矩阵	140
5.4.2	计算序列间的距离	141
5.4.3	Neighbor-Joining 算法	143
5.4.4	UPGMA 算法	146
5.5	构建分子进化树 III: 最大似然法和贝叶斯方法	148
5.5.1	最大似然法	148
5.5.2	*贝叶斯方法	150
5.6	*构建分子进化树的两个实际问题	151
5.6.1	一致性和长分枝相吸现象	151
5.6.2	Bootstrap 分析	153
5.7	祖先状态的推断	153

5.7.1	问题的定义	153
5.7.2	Sankoff 算法	154
5.7.3	最大似然法	155
5.7.4	推断方法的准确率	156
5.8	基因树和物种树的融合	158
5.8.1	基因簇和基因树	158
5.8.2	基因树和物种树的融合的定义	159
5.8.3	推断基因复制事件	160
5.9	参考文献	162
5.10	练习题	164
第六章	计算蛋白质组学	169
6.1	基础知识	169
6.1.1	氨基酸和肽的质量	169
6.1.2	质谱仪和质谱	169
6.1.3	同位素峰、误差和噪音	171
6.1.4	连续质谱	174
6.1.5	复杂蛋白样本的处理	175
6.1.6	肽鉴定的基本方式	177
6.2	肽从头测序	178
6.2.1	打分函数	178
6.2.2	PEAKS 算法	179
6.2.3	谱图算法	185
6.3	搜库及其统计学验证	188
6.3.1	打分函数	188
6.3.2	对结果的质控	190
6.4	翻译后修饰	192
6.5	其它研究课题	193
6.5.1	定量分析	194
6.5.2	糖鉴定	194
6.5.3	新型肽鉴定方法	194
6.5.4	其它分子的鉴定	195
6.6	参考文献	195
6.7	练习题	198
索 引		199
英汉术语对照		205