

“十一五”国家重点图书

计算机科学与技术学科前沿丛书

计算机科学与技术学科研究生系列教材（中文版）

---

# 生物信息学导论

## ——面向高性能计算的算法与应用

---

王勇献 王正华 编著

---



清华大学出版社



中国科学院植物研究所 中国科学院昆明植物研究所 昆明植物研究所



# 生物信息学导论

—— 基因组学、蛋白质组学、代谢组学、系统生物学

张松海 王中林 编著



科学出版社

“十一五”国家重点图书

计算机科学与技术学科前沿丛书

计算机科学与技术学科研究生系列教材（中文版）

---

# 生物信息学导论

## ——面向高性能计算的算法与应用

---

王勇献 王正华 编著

---

清华大学出版社

北京

## 内 容 简 介

本书主要针对生物信息学中的典型应用,从计算方法角度介绍相关算法的原理及应用;内容分成生物学及数理基础、生物序列分析、蛋白质组学分析以及大规模生物学网络分析等四个专题,涉及生物分子序列分析、基因发现、分子进化分析、蛋白质结构预测、蛋白质肽测序、生物学网络模块划分等具体问题的求解原理及算法实现。

本书的读者对象是具有现代分子生物学及计算机科学基本知识的研究生及相关科研人员,在附加习题后也可作为生物信息学方面的入门及进阶教材,供生物医学工程、计算机应用等专业学生使用。

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

生物信息学导论:面向高性能计算的算法与应用/王勇献,王正华编著. —北京:清华大学出版社,2011.6

(计算机科学与技术学科前沿丛书 计算机科学与技术学科研究生系列教材(中文版))

ISBN 978-7-302-25022-7

I. ①生… II. ①王… ②王… III. ①生物信息论-研究生-教材 IV. ①Q811.4

中国版本图书馆CIP数据核字(2011)第045333号

责任编辑:焦虹 徐跃进

责任校对:李建庄

责任印制:何芊

出版发行:清华大学出版社

地 址:北京清华大学学研大厦A座

<http://www.tup.com.cn>

邮 编:100084

社 总 机:010-62770175

邮 购:010-62786544

投稿与读者服务:010-62795954,jsjic@tup.tsinghua.edu.cn

质 量 反 馈:010-62772015,zhiliang@tup.tsinghua.edu.cn

印 刷 者:北京富博印刷有限公司

装 订 者:北京市密云县京文制本装订厂

经 销:全国新华书店

开 本:185×260 印 张:32.5 字 数:810千字

版 次:2011年6月第1版 印 次:2011年6月第1次印刷

印 数:1~3000

定 价:49.00元

---

产品编号:040630-01

“十一五”国家重点图书 计算机科学与技术学科前沿丛书

计算机科学与技术学科研究生系列教材

编  
委  
会

■ 名誉主任：陈火旺

■ 主 任：王志英

■ 副 主 任：钱德沛 周立柱

■ 编委委员：（按姓氏笔画为序）

马殿富 李晓明 李仲麟 吴朝晖

何炎祥 陈道蓄 周兴社 钱乐秋

蒋宗礼 廖明宏

# 序

未来的社会是信息化的社会，计算机科学与技术在其中占据了最重要的地位，这对高素质创新型计算机人才的培养提出了迫切的要求。计算机科学与技术已经成为一门基础技术学科，理论性和技术性都很强。与传统的数学、物理和化学等基础学科相比，该学科的教育工作者既要培养学科理论研究和基本系统的开发人才，还要培养应用系统开发人才，甚至是应用人才。从层次上来讲，则需要培养系统的设计、实现、使用与维护等各个层次的人才。这就要求我国的计算机教育按照定位的需要，从知识、能力、素质三个方面进行人才培养。

硕士研究生的教育须突出“研究”，要加强理论基础的教育和科研能力的训练，使学生能够站在一定的高度去分析研究问题、解决问题。硕士研究生要通过课程的学习，进一步提高理论水平，为今后的研究和发展打下坚实的基础；通过相应的研究及学位论文撰写工作来接受全面的科研训练，了解科学研究的艰辛和科研工作者的奉献精神，培养良好的科研作风，锻炼攻关能力，养成协作精神。

高素质创新型计算机人才应具有较强的实践能力，教学与科研相结合是培养实践能力的有效途径。高水平人才的培养是通过被培养者的高水平学术成果来反映的，而高水平的学术成果主要来源于大量高水平的科研。高水平的科研还为教学活动提供了最先进的高新技术平台和创造性的工作环境，使学生得以接触最先进的计算机理论、技术和环境。高水平的科研也为高水平人才的素质教育提供了良好的物质基础。

为提高高等院校的教学质量，教育部最近实施了精品课程建设工程。由于教材是提高教学质量的关键，必须加快教材建设的步伐。为适应学科的快速发展和培养方案的需要，要采取多种措施鼓励从事前沿研究的学者参与教材的编写和更新，在教材中反映学科前沿的研究成果与发展趋势，以高水平的科研促进教材建设。同时应适当引进国外先进的原版教材，确保所有教学环节充分反映计算机学科与产业的前沿研究水平，并与未来的发展趋势相协调。

中国计算机学会教育专业委员会在清华大学出版社的大力支持下，进行了计算机科学与技术学科硕士研究生培养的系统研究。在此基础上组织来自多所全国重点大学的计算机专家和教授们编写和出版了本系列教材。作者们以自己多年来丰富的教学和科研经验为基础，认真研究和结合我国计算机科学与技术学科硕士研究生教育的特点，力图使本系列教材对我国计算机科学与技术学科硕士研究生的教学方法和教学内容的改革起引导作用。本系列教材的系统性和理论性强，学术水平高，反映科技新发展，具有合适的深度和广度。同时本系列教材两种语种（中文、英文）并存，三种版权（本版、外版、

合作出版)形式并存,这在系列教材的出版上走出了一条新路。

相信本系列教材的出版,能够对提高我国计算机硕士研究生教材的整体水平,进而对我国大学的计算机科学与技术硕士研究生教育以及培养高素质创新型计算机人才产生积极的促进作用。

陈永旺



# 序 言

生物信息学 (又称计算生物学) 是计算机、生物、数学等多学科交叉而新兴的学科分支, 以借助计算机科学、信息科学等领域的算法与工具解决生命科学中的问题为主要特征。随着生命科学研究的深入和后基因组时代的来临, 生物信息学所研究的问题已经发生了巨大的变化, 新的研究越来越需要借助于高性能计算环境的支持, 然而综合生命科学与高性能计算两个领域的知识、有效解决现实中的应用问题并不是一件容易的事, 目前也缺乏系统深入的相关图书资料。编写这本《生物信息学导论》, 作者希望在向读者介绍生物信息学中与高性能计算结合最密切的一些基础性问题, 讨论并总结相关的求解算法与应用技术。

尽管本书定位在生物信息学“导论”的层次, 但是作者并没有打算从生物信息学领域的概念内涵、发展现状入手, 也没有追求内容上的面面俱到、或者对生物信息学领域内容进行泛泛介绍, 而是有所选择地介绍了生物分子序列分析、基因发现、分子进化分析、蛋白质结构预测、蛋白质肽测序、生物学网络模块划分等具体问题的求解算法及原理实现。从这个意义上讲, 也可以将本书看成是生物信息学部分专题内容的汇集。在每部分专题内容中, 既有对经典方法的详细讨论, 也融入了作者及其合作研究者最近几年研究的创新成果, 既注重理论方面的方法 (例如: 利用谱分析挖掘 PPI 网络中的典型模式), 也强调具体应用方面的实现 (例如: 利用 MPI 实现并行计算)。值得说明的是, 结合作者知识背景与研究兴趣, 本书在介绍各类生物信息学问题的求解方法时, 特别关注了如何跟高性能计算技术相结合 (例如: 关于序列比对的并行计算)。

全书正文各章节结构如下图所示, 共分为“预备知识篇”、“序列分析篇”、“蛋白质组分析篇”和“生物学网络分析篇”等四部分。

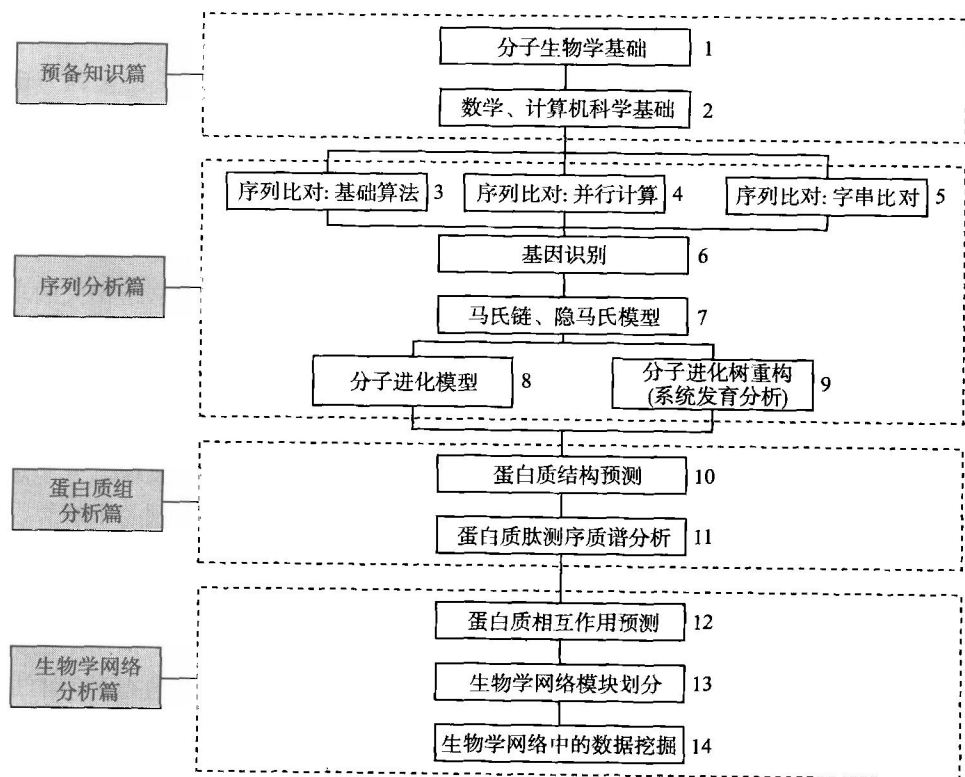
“预备知识篇” (包括第 1 章和第 2 章) 提供了生物信息学分析中涉及的常用分子生物学、线性代数、概率统计以及计算机算法等方面的基础知识, 不同知识背景的读者可以选择自己所需的内容。读者阅读时也可以直接跳过这两章, 在后续章节需要相关知识时, 再回头翻阅这些内容。

“序列分析篇”是全书的重点内容之一, 共包括七章内容 (第 3 章到第 9 章), 主要涵盖生物大分子的序列比对、DNA 序列上的基因识别 (或称基因发现) 和分子系统发育分析等。其中, 序列比对共有三章, 分别从经典基础算法、高性能并行实现以及基于字符串的模式匹配三个层面进行介绍; 分子系统发育分析共有两章, 分别介绍分子进化模型和进化树重构方法等内容; 作为序列分析中最常用的一种模型, 我们在本篇内容中专门介绍了马氏链和隐马氏模型 (第 7 章), 该模型在后续各篇中也有应用。



“蛋白质组分析篇”包括两章内容(第10章和第11章),分别介绍了基于序列预测蛋白质结构的计算分析,以及基于质谱数据分析的蛋白质序列测定方法等内容。

“生物学网络分析篇”是全书的最后一部分,共包含三章内容(第12章到第14章),重点以蛋白质相互作用为例,介绍了蛋白质相互作用预测、蛋白质相互作用网络的模块划分与功能预测等应用中的计算方法,最后一章系统总结了一般生物学网络中的数据挖掘方法与结果。



本书的读者对象是具有现代分子生物学及计算机科学基本知识的研究生及相关科研人员,在附加习题后也可作为生物信息学方面的入门及进阶教材,供分子生物学、计算机应用等专业的读者使用。为了照顾不同学科背景知识读者的需求,本书在开头还简要介绍了阅读后续章节所需要的分子生物学、数学及计算机科学基础知识;全书最后附有详细的主题索引、人名索引及以字母排序的参考文献,每条文献还特意标注了在正文中被引用处的页码,以方便读者检索。

本书是在为国防科技大学硕士研究生开设的“生物信息学导论”课程讲义以及作者从事国家自然科学基金项目研究成果的基础上整理而成。作者感谢2005—2010各学年选修这门课程学习的所有同学,他们参与课程讨论的许多内容构成了本书的基本素材;本书在选题与出版方面得到了国家自然科学基金(60603054)、湖南省自然科学基金(08JJ4021)及国家重点基础研究发展计划课题(2009CB723803)的资助。

在书稿准备出版过程中,清华大学出版社的广大员工给予了大力支持与帮助,在此

一并表示感谢。

全书由王勇献主编，王正华对内容进行了统稿并提出了改进意见。由于作者水平所限，书中还有很多错误和不足之处，希望读者批评指正（作者联系方式：湖南长沙国防科技大学计算机学院，邮编：410073，电子邮箱：yxwang@nudt.edu.cn）。

作者于长沙

2010年12月

# 目 录

<b>第一篇 预备知识篇</b>	<b>1</b>
<b>第 1 章 分子生物学基础</b>	<b>3</b>
1.1 生命的演化与分类	4
1.2 核酸: DNA 与 RNA	5
1.3 蛋白质	7
1.4 DNA 的复制	9
1.5 基因与染色体	10
1.6 基因表达	10
1.6.1 转录	11
1.6.2 遗传密码	11
1.6.3 基因的进化——遗传与变异	14
1.7 现代生物工程技术	16
1.8 现代分子生物学中的经典计算问题	18
<b>第 2 章 数学及计算机科学基础</b>	<b>20</b>
2.1 线性代数理论	20
2.1.1 记号与约定	20
2.1.2 矩阵的范数	20
2.1.3 矩阵的特征值与特征向量	21
2.1.4 矩阵的广义逆	22
2.2 概率论基础知识	22
2.2.1 随机事件	22
2.2.2 概率的三种定义	23
2.2.3 概率的加法原理	24
2.2.4 条件概率	24
2.2.5 全概率公式和 Bayes 公式	24
2.2.6 独立随机试验与贝努利定律	25
2.2.7 随机变量及其分布	26
2.2.8 常用的随机分布	27
2.2.9 概率分布的熵与相对熵	30

2.2.10	随机过程	31
2.2.11	一阶马氏链	31
2.2.12	随机游动	34
2.2.13	高阶马氏链	34
2.2.14	统计推断与假设检验	35
2.3	最优化理论	35
2.3.1	问题描述	35
2.3.2	Lagrange 理论	37
2.4	统计学习理论	41
2.4.1	引言	41
2.4.2	机器学习的基本问题和方法	42
2.4.3	统计学习理论的核心内容	45
2.5	函数增长速度的比较	54
 <b>第二篇 序列分析篇</b>		 <b>57</b>
 <b>第 3 章 序列比对的基本方法</b>		 <b>59</b>
3.1	序列的相似性与同源性	59
3.2	点阵图	60
3.3	两序列比对概述	61
3.4	全局比对的动态规划方法	62
3.5	局部比对的动态规划方法	64
3.6	重叠区域匹配的准全局比对算法	66
3.7	空位罚分模型	68
3.8	仿射空位罚分模型下的全局比对算法	69
3.9	仿射空位罚分模型下的局部比对算法	72
3.10	降价空间存储的两序列比对算法	75
3.10.1	线性空间复杂性算法	75
3.10.2	CheckPoint 算法	77
3.11	降低时间开销的两序列比对算法	82
3.11.1	分块比对算法	82
3.11.2	带状比对算法	83
3.12	比对得分的正则化	85
3.13	启发式的近似寻优比对算法	86
3.13.1	FASTA	86
3.13.2	BLAST	88
3.14	比对得分的统计学显著性	90

3.15 多序列比对 . . . . .	90
3.15.1 MSA . . . . .	93
3.15.2 渐进式比对 . . . . .	94
3.15.3 Gibbs 采样方法 . . . . .	97
3.15.4 启发式多序列比对软件与工具 . . . . .	98
3.16 氨基酸替换矩阵 . . . . .	99
3.16.1 PAM 氨基酸替换矩阵 . . . . .	99
3.16.2 BLOSUM 氨基酸替换矩阵 . . . . .	101
3.17 小结 . . . . .	102
<b>第 4 章 序列比对的并行计算 . . . . .</b>	<b>103</b>
4.1 并行编程模型 . . . . .	103
4.1.1 并行计算的粒度 . . . . .	103
4.1.2 进程间的通信 . . . . .	104
4.2 并行计算机系统结构 . . . . .	105
4.2.1 通用并行计算机系统 . . . . .	105
4.2.2 专用并行处理硬件 . . . . .	105
4.3 序列比对及其并行化方案 . . . . .	106
4.4 Smith-Waterman 算法的细粒度并行实现 . . . . .	107
4.4.1 SWMMX 并行算法 . . . . .	108
4.4.2 SWSSE2 并行算法 . . . . .	109
4.4.3 条带型并行算法 . . . . .	110
4.4.4 基于分块分治策略的并行算法 . . . . .	111
4.4.5 其他并行算法 . . . . .	115
4.5 序列数据库搜索的粗粒度并行算法 . . . . .	116
4.5.1 并行 FASTA . . . . .	116
4.5.2 TurboBLAST . . . . .	117
4.5.3 mpiBLAST . . . . .	117
4.6 多序列比对的并行算法 . . . . .	118
4.6.1 HMMER 及其并行算法 . . . . .	118
4.6.2 ClustalW . . . . .	119
4.6.3 ClustalW-MPI . . . . .	121
4.6.4 并行 ClustalW、HT Clustal 和 MULTICLUSTAL . . . . .	121
4.7 基于专用硬件 FPGA 的序列比对 . . . . .	123
4.7.1 FPGA 硬件设备 . . . . .	123
4.7.2 FPGA 并行计算 . . . . .	124

<b>第 5 章 基于字符串精确匹配的序列比较</b>	<b>127</b>
5.1 模式的精确匹配与非精确匹配	127
5.2 朴素的模式匹配算法	128
5.3 线性时间的字符串搜索算法	128
5.4 基于关键字树的模式集合匹配算法	130
5.5 后缀树	132
5.6 后缀树的构造	134
5.7 后缀数组	135
5.8 基因组中的重复序列	136
5.9 后缀树用于搜索重复子串和独特子串	136
5.10 最长重复序列的搜索算法	137
5.11 广义后缀树	138
5.12 最长公共子串问题	138
5.13 $k$ 次失配问题	139
5.14 小结	141
<b>第 6 章 基因识别</b>	<b>142</b>
6.1 基因识别与预测的计算方法	142
6.2 预测算法的准确性度量	144
6.3 独立识别法	145
6.3.1 用于基因识别的常用序列信号	146
6.3.2 阅读框的相位及基因中的外显子类型	146
6.3.3 密码子使用偏好	147
6.3.4 用序列特征图寻找剪接位点	149
6.3.5 外显子链问题	151
6.4 基于比较的基因识别方法	153
<b>第 7 章 马氏链与隐马氏模型</b>	<b>156</b>
7.1 马尔可夫链	156
7.2 隐马尔可夫模型	159
7.3 计算全概率的正向算法	162
7.4 计算全概率的反向算法	164
7.5 解码问题的 Viterbi 算法	166
7.5.1 各时间点独立考虑的最可能路径	166
7.5.2 各时间点综合考虑的最可能路径	167
7.6 模型参数的估计	169
7.6.1 已知路径时的参数重估	169
7.6.2 Baum-Welch 方法	170

7.6.3	Baum-Welch 算法的推导	173
7.6.4	参数重估的 Baldi-Chauvin 梯度下降法	174
7.6.5	Baldi-Chauvin 梯度下降法的推导	175
7.6.6	Mamitsuka 算法	177
7.6.7	Mamitsuka 参数重估算法的推导	177
7.7	带有哑状态的 HMM	178
7.8	谱 HMM	181
7.9	采用谱 HMM 进行多序列比对建模	183
7.10	利用 HMM 对基因识别问题进行建模	184
<b>第 8 章</b>	<b>序列进化的基本模型</b>	<b>186</b>
8.1	核苷酸替代的进化模型	186
8.2	连续时间下的进化模型	190
8.2.1	Jukes-Cantor 进化模型	190
8.2.2	Kimura 进化模型	191
8.2.3	Felsenstein 进化模型	192
8.2.4	HKY 进化模型	192
8.3	离散时间下的进化模型	193
8.3.1	Jukes-Cantor 进化模型	193
8.3.2	Kimura 进化模型	194
8.3.3	Felsenstein 进化模型	196
8.3.4	HKY 进化模型	197
<b>第 9 章</b>	<b>分子进化树的重构</b>	<b>198</b>
9.1	进化树的概念与术语	198
9.1.1	二叉树	198
9.1.2	树的标度	198
9.1.3	有根树与无根树	199
9.1.4	树的定根方法	199
9.1.5	物种树与基因树	200
9.1.6	分歧经历的时间	202
9.1.7	树的文本表示法	202
9.1.8	进化树拓扑结构的计数	202
9.1.9	不同树之间的拓扑距离	204
9.1.10	一致树	206
9.1.11	分子进化树重构的基本流程	207
9.2	进化树重构的简约类方法	208
9.3	进化树重构的距离类方法	213



9.3.1	距离 . . . . .	213
9.3.2	邻居加入方法 . . . . .	215
9.3.3	UPGMA 方法 . . . . .	223
9.3.4	误差平方和最小方法 . . . . .	226
9.4	进化树重构的统计类方法 . . . . .	228
9.4.1	树的似然度 . . . . .	229
9.4.2	Horner 规则与修剪算法 . . . . .	230
9.4.3	算法加速的策略 . . . . .	232
9.4.4	时间可逆性、树的根结点及分子钟树间的关联性 . . . . .	233
9.4.5	数据缺失及比对空位的处理 . . . . .	234
9.4.6	进化速率关于位点可变的建模方法 . . . . .	235
9.5	树拓扑空间的搜索技术 . . . . .	238
9.5.1	最近邻居交换法 . . . . .	238
9.5.2	子树剪枝嫁接法 . . . . .	239
9.5.3	分支界限法 . . . . .	240
9.6	似然度最大化的数值算法 . . . . .	240
9.6.1	一元函数优化问题 . . . . .	241
9.6.2	多变量优化问题 . . . . .	242
9.6.3	进化树分析中参数估计的应用问题 . . . . .	244
9.7	模型选择与假设检验问题 . . . . .	245
9.7.1	似然比检验 . . . . .	245
9.7.2	Akaike 信息准则方法 . . . . .	246
9.7.3	Bayes 信息准则方法 . . . . .	246
9.8	进化树拓扑结构的建模、估计与检验 . . . . .	246
9.8.1	估计与假设检验 . . . . .	246
9.8.2	Bootstrap 方法 . . . . .	247
9.8.3	内部分支检验法 . . . . .	252
9.8.4	KH 检验与修正 . . . . .	253
9.8.5	简约类方法中的指标 . . . . .	254
<b>第三篇 蛋白质组学分析篇</b>		<b>255</b>
<b>第 10 章 蛋白质的结构预测</b>		<b>257</b>
10.1	蛋白质的层次性结构 . . . . .	257
10.2	常见的二级结构单元 . . . . .	258
10.2.1	螺旋结构 . . . . .	259
10.2.2	$\beta$ 折叠结构 . . . . .	261

10.2.3 $\beta$ 转角结构	262
10.3 蛋白质二级结构检测	263
10.4 蛋白质二级结构预测的计算方法	265
10.4.1 早期的预测方法	266
10.4.2 判别分析法	266
10.4.3 基于神经网络的预测算法	270
10.4.4 最近邻居法	272
10.4.5 基于谱 HMM 的结构预测	273
10.4.6 结构预测的线索化方法	273
10.4.7 结构预测的分子动力学方法	274
10.4.8 蛋白质折叠预测的格子化 HP 模型	276
10.5 蛋白质二级结构预测算法的性能评价	277
10.5.1 问题描述	278
10.5.2 蛋白质结构预测算法性能评估指标	279
10.5.3 性能评估指标对结构预测建模的指导作用	283
10.5.4 各评估指标的比较及使用原则	285
10.6 蛋白质结构的比对方法	286
10.6.1 肽链局部结构特征的提取	286
10.6.2 结构特征的规范化及广义后缀树的构建	288
10.6.3 蛋白质结构的比较与搜索	289
<b>第 11 章 蛋白质序列鉴定的质谱分析</b>	<b>291</b>
11.1 质谱技术	291
11.1.1 质谱仪的基本工作原理	291
11.1.2 串联质谱仪	292
11.2 质谱数据分析	292
11.2.1 串联质谱中的离子类型	292
11.2.2 质谱图	294
11.2.3 碎片离子质量与母离子质量的关系	295
11.2.4 理论质谱与实验质谱	296
11.3 实验质谱数据的预处理	297
11.3.1 噪声过滤的基线确定方法	298
11.3.2 同位素峰识别方法	299
11.4 质谱比较的非概率型打分方法	299
11.4.1 基于单峰或区间匹配的打分	299
11.4.2 基于向量夹角余弦的打分	299
11.4.3 基于信号互相关性的打分	300