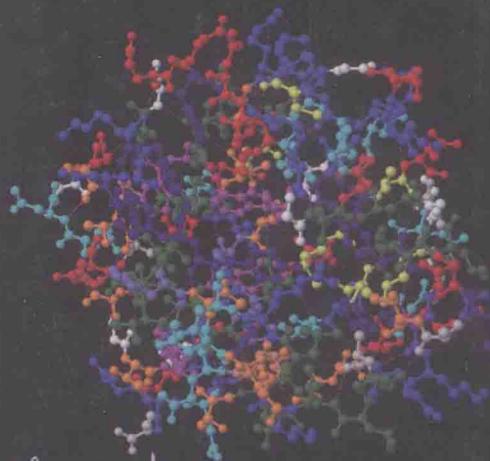




普通高等教育“十一五”国家级规划教材

生物信息学

主编 刘 娟



HTADLSPVLLSUNFPAVWALGSAEELLSVAVLWREKPTNFTLLSTFVCFKAGEK
NYLUSGTVLSSNFGVLAQKSPKREKALLESVIVLWREKPTNFTLLSTFVCFKAGEK
--FALLSFTVLLSSTVFEVHVALEKSPKREKALLESVIVLWREKPTNFTLLSTFVCFKAGEK
SPVHLSFTVLLSSTVFEVHVALEKSPKREKALLESVIVLWREKPTNFTLLSTFVCFKAGEK
SVAQVSPVAVLWREKPTNFTLLSTFVCFKAGEK

高等教育出版社



普通高等教育“十一五”国家级规划教材

内容简介

本书以生物信息学的发展为背景，介绍生物信息学的基本概念、基本理论和基本方法。全书共分10章，主要内容包括：生物信息学概论、生物数据库、生物序列分析、生物结构分析、生物系统分析、生物网络分析、生物数据挖掘、生物信息学应用等。

生物信息学

S H E N G W U X I N X I X U E

主编 刘娟

编者 (按姓氏拼音排序)

段谟杰 李论 刘娟 陆枫

罗飞 马闯 周到

第1次印刷 第1次印刷 第1次印刷 第1次印刷

400-810-0298
http://www.hep.edu.cn
http://www.hep.com.cn
http://www.hep.cn
http://www.hep.com.cn
2014年12月第1版
2014年12月第1次印刷
28 005

高等教育出版社
北京市西城区德胜门内大街2号
100120
三河市河北新华印刷有限公司
787mm×1092mm 1/16
19.2
490千字
2014年12月第1版

高等教育出版社·北京

内容简介

本书围绕目前生物信息学研究与应用的主要内容,以丰富的实例,重点介绍了相关数据库和软件的功能、应用策略和使用方法。具体内容包括:核酸与蛋白质序列数据资源、序列比较与相似序列搜索、分子系统发育分析、基因组结构注释、蛋白质结构分析、蛋白质序列分析、Microarray 基因表达数据分析、蛋白质组数据分析、生物信息学在疾病相关基因与药物发现中的应用,以及生物信息导航资源。本书试图综合介绍生物信息学研究解决的问题、基本方法、现有成果与存在的问题,特别是能使读者把握生物信息学自身的特点和分析解决问题的基本途径,使不同专业背景读者都能有一定的收获。

本书适合作为生命科学、计算机科学等相关专业的教材使用,也可供相关科研人员参考使用。

图书在版编目(CIP)数据

生物信息学 / 刘娟主编. -- 北京:高等教育出版社, 2014.12

ISBN 978-7-04-040975-8

I. ①生… II. ①刘… III. ①生物信息论-高等学校-教材 IV. ①Q811.4

中国版本图书馆 CIP 数据核字(2014)第 204427 号

策划编辑 王莉 责任编辑 孟丽 封面设计 姜磊 责任印制 尤静

出版发行 高等教育出版社
社址 北京市西城区德外大街4号
邮政编码 100120
印刷 三河市华润印刷有限公司
开本 787mm × 1092mm 1/16
印张 19.5
字数 460千字
购书热线 010-58581118

咨询电话 400-810-0598
网 址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>
网上订购 <http://www.landracom.com>
<http://www.landracom.com.cn>
版次 2014年12月第1版
印次 2014年12月第1次印刷
定价 38.00元

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换
版权所有 侵权必究
物料号 40975-00

数字课程 (基础版)

生物信息学

主编 刘娟

登录方法:

1. 访问<http://abook.hep.com.cn/40975>
2. 输入数字课程用户名 (见封底明码)、密码
3. 点击“进入课程”

账号自登录之日起一年内有效, 过期作废
使用本账号如有任何问题
请发邮件至: lifescience@pub.hep.cn



生物信息学

主编 刘娟

用户名 密码 验证码 3 4 5 9 进入课程

使用说明

内容介绍

纸质教材

版权信息

联系方式

围绕本教材知识体系和生物信息学的特点, 立足反映学科快速发展的趋势和成果, 本书配套数字资源涵盖了全书彩图、重要数据库链接和简介、重要生物信息统计分析软件功能简介等。建议教师根据本书资源特点和类型, 引导学生自主学习。



数字课程网站

网址: <http://abook.hep.com.cn/201512>
<http://abook.hep.com.cn/201512>

用户名: 输入教材封底的16位明码; 密码: 刮开“增值服务”涂层, 输入16位暗码; 输入正确的验证码后, 点击“进入课程”开始学习。

Copyright © 2014-2015 高等教育出版社 版权所有

<http://abook.hep.com.cn/40975>

前 言

生物信息学是生命科学、信息科学、统计学等多学科交叉的学科，自人类基因组计划以来迅速发展。促使该学科突飞猛进发展的根本原因就是数据。一方面，海量的实验数据驱动数据的表示、管理和分析方法的发展。另一方面，及时、有效和充分地利用海量的生物信息资源来进行生物机制的理解和生物功能的发现已经成为现代生物、医学、药学等领域研究的必备方法。利用生物信息学的资源和手段开展研究，已经成为多学科研究人员必须要掌握的一种基本技能。鉴于此，我们编写了这本《生物信息学》教材，对生物信息学的常用数据资源、常见的分析手段和工具，以及如何利用数据资源开展应用研究等进行了较为系统的介绍。

全书共分 11 章，各章都配有一定数量的思考题，帮助读者巩固和消化本书内容。第一章简要介绍了生物信息学的发展历史和本书的主要内容。第二章介绍了常用的序列数据库资源及如何从数据库中获取序列数据。第三章比较全面地介绍了序列比对的基本概念、经典的序列比对算法以及有关序列对工具的使用方法等。第四章介绍基因结构的预测方法和 Ensembl 的基因组结构注释流程。第五章系统发生树的基本概念以及分子系统发生树的构建方法。第六章介绍了蛋白质结构的相关知识以及蛋白质结构的预测方法。第七章介绍了蛋白质序列分析方法与蛋白质功能预测方法。第八章介绍了微阵列数据的有关概念和微阵列数据的分析方法。第九章主要介绍了蛋白质质谱数据分析方法和蛋白质-蛋白质互作预测方法等。第十章主要介绍了疾病相关的数据资源和疾病基因发现的生物信息学方法。第十一章简要介绍了 SNP 芯片及深度测序的基本概念。

本书既可以作为生物信息学及相关专业的本科生和研究生的教材，也可以作为生物医学等领域的科研人员有关各种生物信息资源的使用指导书。同时，也可作为生物信息学专业人士的研究参考书。希望不同层次的人员通过阅读此书，各取所需，各有收获。

本书由来自 4 所高校 7 位一线教师编写而成，具体分工如下：第一、二、三章由华中科技大学陆枫编写，第四章由西北农林科技大学马闯编写，第五章由华中科技大学李论、陆枫编写，第六章由华中科技大学段谟杰和武汉大学罗飞编写，第七章由华中科技大学陆枫和武汉大学刘娟编写，第八章由中南民族大学周到、

华中科技大学陆枫和武汉大学刘娟编写，第九章由武汉大学刘娟、罗飞和华中科技大学陆枫编写，第十章由武汉大学刘娟、罗飞和华中科技大学李论、陆枫编写，第十一章由武汉大学罗飞编写。全书所有章节由刘娟统稿并编撰习题。刘娟和罗飞对本书全部内容进行最后校对审核。

高等教育出版社的王莉和孟丽编辑为本书的出版付出了辛勤的劳动，在此一并感谢！同时，我们也感谢在科研工作中给予我们帮助的同仁，以及在本书撰写和出版过程中给我们提供各种帮助的单位和个人！

本书中也引用了生物信息学的国内外文献资料，在此对这些参考文献的作者致以衷心的感谢！由于作者学科背景不同、学术水平有限、写作经验不足，书中不足之处在所难免，敬请读者批评指正。

刘娟

2014年5月

1	绪言	1
1.1	1.1 生物信息学的发展历史	1
1.2	1.2 本书内容简介	3
1.3	1.3 贯穿本书的例子	4
2	序列数据资源	6
2.1	2.1 分子生物学数据库	6
2.2	2.2 序列数据存储格式	9
2.3	2.3 核酸序列数据库	14
2.3.1	2.3.1 GenBank 数据库	14
2.3.2	2.3.2 RefSeq 数据库	20
2.3.3	2.3.3 EPD 数据库	21
2.4	2.4 蛋白质序列数据库	22
2.4.1	2.4.1 UniProt 简介	22
2.4.2	2.4.2 UniProtKB 数据库	22
2.5	2.5 基因组数据资源	27
2.5.1	2.5.1 基础知识	27
2.5.2	2.5.2 不同物种的基因组数据库	29
2.5.3	2.5.3 人类基因组数据库	33
2.6	2.6 数据的检索与获取	44
2.6.1	2.6.1 检索工具	44
2.6.2	2.6.2 获取序列数据的例子	46
	思考题	50
3	序列比对与比对搜索	51
3.1	3.1 基本概念	51
3.1.1	3.1.1 比对序列的选择：核苷酸序列还是蛋白质序列	51
3.1.2	3.1.2 同源性、相似性和一致性	51
3.1.3	3.1.3 空位	54
3.1.4	3.1.4 多序列比对	54
3.2	3.2 Dayhoff 模型：可接受点突变	57
3.2.1	3.2.1 PAM1 矩阵	59
3.2.2	3.2.2 PAM250 和其他 PAM 矩阵	60
3.2.3	3.2.3 从突变概率矩阵到对数比值打分矩阵	61
3.2.4	3.2.4 双序列比对中 PAM 矩阵的实际有用性	63

3.2.5	PAM 矩阵的重要替代者: BLOSUM 打分矩阵	64
3.2.6	双序列比对和检测限度	65
3.3	比对算法: 全局和局部	66
3.3.1	全局序列比对: Needleman - Wunsch 算法	67
3.3.2	局部比对: Smith - Waterman 算法	71
3.3.3	Smith - Waterman 算法的快速和启发式版本	73
3.4	双序列比对的显著性	74
3.4.1	双序列比对统计显著性检验	75
3.4.2	全局比对的统计显著性	75
3.4.3	局部比对的统计显著性	76
3.5	局部比对搜索基本工具 BLAST	76
3.5.1	BLAST 搜索的关键步骤	77
3.5.2	BLAST 算法: 列表、扫描、延伸	80
3.5.3	BLAST 算法的统计学和 E 值	82
3.5.4	BLAST 的各类分值	83
3.5.5	BLAST 搜索示例: 应用搜索原则	85
3.5.6	BLAST 搜索示例: 多结构域蛋白的搜索	90
3.5.7	BLAST 搜索示例: 改变打分矩阵	94
3.6	寻找远缘相关的蛋白质: PSI - BLAST	98
3.6.1	基本步骤	98
3.6.2	PSI - BLAST 的结果评估	100
3.6.3	PSI - BLAST 的错误: 破坏的问题	101
3.7	模式识别 BLAST (PHI - BLAST)	101
3.8	用 BLAST 来发现新基因	103
	思考题	104
4	基因组结构注释	105
4.1	引言	105
4.1.1	基因及其结构	105
4.1.2	基因结构预测概述	107
4.2	基于 EST 序列数据识别基因结构	109
4.2.1	判别基因序列的真实 EST 匹配的措施	110
4.2.2	真实 EST 匹配的识别流程	112
4.2.3	确定 EST 对应的基因结构	114
4.3	基因结构预测的统计学建模方法	114
4.3.1	基于多级优化预测基因结构的基本思想	115
4.3.2	基因结构的分级建模	115
4.3.3	基因结构预测的动态规划算法	119

151	4.3.4 基于统计学方法预测基因结构的效果	120
151	4.4 基因组结构的自动注释	121
151	4.4.1 Ensembl 的基因组注释流程	121
151	4.4.2 Ensembl 自动注释结果与人工注释结果比较	122
151	思考题	123
5	分子系统发生分析	124
151	5.1 分子水平的进化介绍	124
151	5.1.1 问题的历史起源	124
151	5.1.2 分子钟	125
151	5.2 基本概念	126
151	5.2.1 系统发生树的基本概念	126
151	5.2.2 直系同源和旁系同源	128
151	5.3 分子系统发生树的构建	128
151	5.3.1 选择可供分析的序列	128
151	5.3.2 多序列比对	129
151	5.3.3 构建系统发生树	131
151	5.3.4 方法的选取	147
151	5.3.5 常用分析软件	147
151	思考题	148
6	蛋白质结构	149
151	6.1 蛋白质结构	149
151	6.2 蛋白质结构数据库和结构可视化	150
151	6.2.1 PDB 数据库	150
151	6.2.2 蛋白质结构家族分类数据库	154
151	6.2.3 蛋白质结构的可视化	156
151	6.3 蛋白质结构分析	157
151	6.3.1 蛋白质结构比对	157
151	6.3.2 结构模型品质的分析	159
151	6.3.3 蛋白质内部相互作用分析	160
151	6.3.4 溶剂可接近表面的计算及分析	161
151	6.3.5 功能位点的分析	162
151	6.4 蛋白质结构预测	163
151	6.4.1 蛋白质结构比较建模	163
151	6.4.2 蛋白质结构从头预测方法	168
151	6.4.3 二级结构预测	168
151	6.4.4 结构预测的策略	169
151	思考题	170

7	蛋白质序列分析与功能预测	171
7.1	引言	171
7.2	功能描述	172
7.2.1	基因本体	173
7.2.2	利用 GO 术语的功能注释	175
7.3	基于序列相似性的功能预测	178
7.3.1	基本预测方法	178
7.3.2	分析与讨论	182
7.3.3	蛋白质家族与序列的相似性聚类	183
7.4	基于蛋白质信号的功能预测	184
7.4.1	蛋白质信号	185
7.4.2	信号的描述	188
7.4.3	蛋白质模体、结构域和家族数据库	193
7.4.4	分析与讨论	198
7.5	基于蛋白质序列特征的功能预测	198
7.5.1	序列的理化性质	198
7.5.2	跨膜与卷曲螺旋分析	200
7.5.3	蛋白质翻译后修饰分析	202
7.5.4	亚细胞定位预测	204
7.5.5	基于序列特征的蛋白质分子功能预测	205
7.6	功能预测的其他思路	205
	思考题	207
8	微阵列数据分析	208
8.1	微阵列	208
8.1.1	微阵列实验过程	208
8.1.2	微阵列制备	209
8.1.3	杂交方式	209
8.1.4	图像分析	210
8.1.5	数据标准化	210
8.1.6	基因表达矩阵	211
8.1.7	基因表达数据分析	211
8.2	数据预处理	211
8.2.1	全局归一化	212
8.2.2	散点分析	212
8.2.3	数据全局归一化中的局部归一化	214
8.3	差异表达基因的检测	215
8.3.1	基本检验方法	215

8.3.2	分析实例	217
8.3.3	疾病基因表达谱差异分析	221
8.4	微阵列数据的分类分析方法	222
8.4.1	聚类分析	223
8.4.2	分类分析	232
8.5	构建基因调控网络	232
8.5.1	基因调控网络的简单例子	233
8.5.2	微分方程模型	235
8.5.3	布尔网络模型	236
8.5.4	贝叶斯网络模型	237
8.6	微阵列数据与分析软件	238
8.6.1	数据交换标准	238
8.6.2	微阵列数据库	239
8.6.3	微阵列数据分析流程	242
8.6.4	微阵列数据分析工具	244
	思考题	247
9	蛋白质组数据分析	249
9.1	二维凝胶电泳数据分析	249
9.1.1	二维凝胶电泳原理	249
9.1.2	二维凝胶电泳数据及其应用	250
9.2	蛋白质质谱数据分析	253
9.2.1	质谱技术	253
9.2.2	蛋白质的质谱分析	255
9.3	蛋白质互作生物信息学	256
9.3.1	亲和层析和质谱	256
9.3.2	酵母双杂交系统	257
9.3.3	蛋白质-蛋白质互作预测	258
9.3.4	蛋白质相互作用数据库	260
9.4	分析细胞通路的生物信息学方法	261
	思考题	263
10	疾病相关研究	265
10.1	疾病基因相关研究的概述	265
10.2	疾病相关的数据资源	266
10.2.1	人类在线孟德尔遗传数据库	266
10.2.2	遗传关联数据库	270
10.2.3	人类基因突变数据库	272
10.2.4	癌症数据库	273

10.2.5	单核苷酸多态性数据库	274
10.3	疾病基因发现	276
	思考题	280
11	SNP 芯片及深度测序数据分析	281
11.1	SNP 简介	281
11.2	结构变异	282
11.3	SNP 实验简介	283
11.3.1	Illumina 芯片	283
11.3.2	Affymetrix 芯片	285
11.4	深度测序技术	286
11.5	序列数据基本格式	288
11.5.1	FASTQ	289
11.5.2	SAM 和 BAM	290
11.5.3	BED	291
11.5.4	VCF	292
11.6	实例数据分析	292
11.6.1	利用深度测序发现 SNV	293
11.6.2	利用 SNP 芯片检测拷贝数变异	294
	思考题	296
	参考书目	297

1 绪言

生物信息学(bioinformatics)是生命科学、计算机科学、现代信息科学、数学、物理学以及化学等多个学科交叉结合形成的一门新学科,是利用信息技术和数学方法对生命科学研究中的生物信息进行存储、检索和分析的科学。自20世纪80年代至今,生物信息学已得到了快速发展,已对基因组分析、大分子的结构和功能、生化途径、疾病发生以及进化相关等诸多重要的生物学问题展开了研究。其研究成果不仅极大地推动生命科学相关学科的发展,还对农学、医药、食品和环境等领域产生巨大的影响,很有可能引发新的产业革命。

1.1 生物信息学的发展历史

生物信息学作为一门独立的学科只有20多年的历史,但是与生物信息学相关的研究却可以追溯至20世纪中期。为帮助读者更好地了解生物信息学的起源和发展历史,下面简要地列出与生物信息学发展相关的部分重要事件。

(1) 萌芽期(20世纪50—70年代)

1953年,J. D. Watson和F. H. C. Crick根据R. E. Franklin和M. H. F. Wilkins的X-晶体衍射数据提出了DNA双螺旋结构。

1955年,F. Sanger发表了牛胰岛素的蛋白质序列,是生物信息学发展的基础。

1962年,L. Pauling提出了分子进化理论。

1967年,M. O. Dayhoff构建了蛋白质序列数据库。

1970年,用于序列比较的Needleman - Wunsch算法发表。

1971年,蛋白质结构数据库(protein data bank, PDB)在美国纽约Brookhaven国家实验室创建。

1974年,欧洲分子生物学实验室(european molecular biology laboratory, EMBL)建立。

1977年,A. M. Maxam和W. Gilbert发表了化学降解法,F. Sanger和A. R. Coulson发表双脱氧终止方法,用于DNA测序。F. Sanger等完成了第一个基因组序列——噬菌体 ϕ X174。

(2) 形成期(80年代)

1980年,K. Wüthrich发明了利用核磁共振技术测定溶液中生物大分子三维结构的方法;EMBL核酸序列数据库建立。

1981年,用于序列比对的Smith - Waterman算法发表。

1982年,创建了GenBank数据库; λ 噬菌体基因组序列完成。

1984年,日本国立遗传学研究所(national institute of genetics, NIG)开始信息服务。

1986年,SwissProt蛋白序列数据库创立。R. Dulbecco在《科学》杂志上撰文首次提出人类基因组计划的设想。美国能源部正式提出实施测定人类基因组全序列的计划。

1987年,克隆容量可达几百至几千kb的酵母人工染色体(yeast artificial chromosomes, YAC)

问世。NIG 发行日本 DNA 数据库 DDBJ 第一版。

1988 年, D. J. Pearson 和 W. R. Lipman 发表 FASTA 算法。美国国立卫生研究院下属的国家生物技术信息中心(National Center for Biotechnology Information, NCBI)成立。国际人类基因组组织(the Human Genome Organisation, HUGO)成立。欧洲分子生物学网络组织(European Molecular Biology Network, EMBnet)创立。GenBank、EMBL 与 DDBJ 共同成立了国际核酸序列联合数据库中心,建立了合作关系。根据协议,这三个数据中心各自搜集世界各国有关实验室和测序机构所发布的序列数据,并通过计算机网络每天都将新发现或更新过的数据进行交换,以保证这三个数据库序列信息的完整性。

1989 年,林华安首先采用“bioinformatics”一词。美国成立国家人类基因组研究中心, J. D. Watson 出任第一任主任。美国 Affymetrix 公司研制出了世界首张基因芯片。

(3) 高速发展期(90 年代至今)

1990 年,美国国会批准正式启动人类基因组计划(Human Genome Project, HGP)研究,计划用 30 亿美元的预算在 15 年的时间内完成人类 30 亿碱基对的测序和基因确定;随后法国、英国、意大利、德国、日本、中国等也相继宣布开始各自的 HGP 研究。S. F. Altschul 发表 BLAST 算法。

1991 年, J. C. Venter 在《科学》杂志上描述表达序列标签(expressed sequence tag, EST)的建立和使用。

1992 年, J. C. Venter 在美国马里兰州成立基因组研究所(the Institute of Genome Research, TIGR),成为细菌基因组测序研究的先驱; M. Simon 和同事宣布细菌人工染色体(bacterial artificial chromosome, BAC)在 DNA 克隆中的应用。

1994 年,欧洲生物信息学研究所(European Bioinformatics Institute, EBI)成立。

1995 年,《科学》杂志首次刊登了 TIGR 采用其创立的全基因组鸟枪法(whole genome shotgun, WGS)完成的流感嗜血杆菌(*Haemophilus influenzae*, *H. inf*)全基因组测序的论文,这是人类完成的第一个单细胞微生物的基因组序列的测定,标志着基因组时代的真正开始。

1996 年,第一个真核生物——酿酒酵母(*Saccharomyces cerevisiae*, *S. cere*)——全基因组测序完成。Affymetrix 开始正式销售商用基因芯片。

1997 年,第一个重要的实验模式生物——大肠杆菌(*Escherichia coli*, *E. coli*)——全基因组测序完成。S. F. Altschul 发表 PSI-BLAST 算法。

1998 年,第一个多细胞真核生物——线虫(*Caenorhabditis elegans*, *C. elegans*)——全基因组测序完成。

2000 年,拟南芥(*Arabidopsis thaliana*)的全基因组测序工作完成,成为植物界第一个被完整测序的物种。黑腹果蝇(*Drosophila melanogaster*)的基因组测序完成。

2001 年,国际人类基因组测序协作组(The International Human Genome Sequencing Consortium, IHGSC)和 J. C. Venter 领导的 Celera 公司分别在《自然》和《科学》杂志同时发表人类基因组草图。

2002 年,小鼠、水稻基因组草图公布。

2003 年,人类基因组测序计划完成。

1.2 本书内容简介

目前的生物信息学研究,已从早期以数据库建立和 DNA 序列分析为主的阶段,转移到后基因组学时代以比较基因组学(comparative genomics)、功能基因组学(functional genomics)和整合基因组学(integrative genomics)为中心的新阶段。本书对以下生物信息学相关内容给予了介绍。

(1) 序列数据资源

序列数据资源贮存了生物信息学研究的原始数据,是生物信息学存在和发展的基础。熟悉并了解这些数据资源将有助于更好地开展生物信息学相关的研究与应用。本书介绍了多个常用的核苷酸和蛋白质序列数据库,以及从这些数据库中获取信息的方法。

(2) 序列比对与比对搜索

序列相似性分析是生物信息学最早涉及的问题之一,常用的分析方法是序列比对(sequence alignment)。本书介绍了多个知名的序列比对算法(包括 Needleman - Wunsch 算法和 Smith - Waterman 算法等),以及目前广泛使用的序列比对工具(例如,NCBI 提供的 BLAST 等)。为了进一步帮助读者理解序列比对的结果,本书还从进化的角度描述了两条序列之间的氨基酸(或者核苷酸)是如何被比对和相互比较的。

(3) 基因组结构注释

随着人类基因组测序计划的完成及其他物种基因组测序工作的相继展开,公共数据库中已积累了大量的基因组序列数据。为了充分发挥这些基因组序列的应用价值,一项极其重要的工作即是“翻译”出基因组序列中蕴含的生物学知识。本书主要以“翻译”基因组序列中蛋白质编码基因为例,介绍基因组结构注释相关的生物信息学方法和软件。

(4) 分子系统发生分析

系统发生关系是表示物种进化关系的参考依据。通过分析分子水平的序列数据,可以了解物种系统发生的关系,目前常用树的形式来表示不同物种间的进化关系。本书结合相关实例主要对分子系统发生树的构建方法和过程进行了介绍。

(5) 蛋白质结构

蛋白质的空间结构是其行使功能的基础。在进行蛋白质相关研究时,我们经常会遇到一些难题。比如,跨膜蛋白如何能控制跨膜区域的开合并实现细胞内外分子的转运;又如一些酶被简单的化学修饰后,其活性为何会发生巨大的改变;脲蛋白发生怎么样的构象变化才会导致诸如疯牛病等。要解答上述的问题,都需要深入地了解蛋白质结构。本书介绍了蛋白质结构相关生物信息学的多个重要内容,包括蛋白质结构相关数据库、常用的蛋白质结构分析工具以及利用计算手段预测蛋白质结构的方法。

(6) 蛋白质序列分析与功能预测

蛋白质在生命过程中发挥着巨大的作用,它们执行着大部分生物学功能,包括结构功能、酶功能,以及在细胞内或细胞间转运物质的功能等。本书将会介绍用于描述蛋白质功能的 GO 术语,以及基于不同策略进行蛋白质功能注释的生物信息学方法。

(7) 微阵列数据分析

在考察某个具体生命状态下的基因表达水平时,实验人员最想知道的问题是:在这个生命状态中,哪些基因相对于它们正常状态有较为明显的高或低表达?微阵列(*microarray*)是一种重要的基因表达高通量检测技术。分析微阵列数据时,最大的挑战是如何使用相应的统计学手段判定哪些基因确实存在改变。本书对微阵列的概念、微阵列数据的处理方法和分析软件,以及基于微阵列数据的基因调控网络的构建等内容进行了较详细介绍。

(8) 蛋白质组数据分析

高通量的蛋白质组工程能够大范围地确定蛋白质功能,能确定蛋白质在哪种特殊的生理条件下会出现,还能确定哪些蛋白质之间有相互作用。目前,除了能单独研究一个蛋白质外,同时对数千个蛋白质进行高通量的分析也已成为可能。本书主要介绍了几种对蛋白质进行高通量分析的手段,包括二维凝胶电泳、质谱数据分析、蛋白质-蛋白质互作的生物信息学方法以及分析细胞通路的生物信息学方法等。

(9) 疾病相关研究

寻找疾病相关基因是认识疾病发生机理、研制疾病的基因诊断与防治手段的基础,也是人类基因组研究的重要目标。现在,一些与遗传病有关的重要基因已被分离和测序,另一些常见病,如乳腺癌、结肠癌、高血压、糖尿病和阿尔茨海默病等具有遗传倾向的疾病基因已在染色体遗传图谱上精确定位,寻找致病基因的研究工作也在不断的深入,这些工作为新药的发现和开发提供了可靠的信息和软件工具资源。本书主要介绍了疾病相关的数据资源以及现有的疾病基因预测方法与软件工具。

(10) SNP 芯片及深度测序数据分析

后基因组时代一个典型特征就是生物实验芯片化,深度测序技术的不断发展极大地促进了生物信息学的研究。本章介绍了 SNP 芯片及深度测序技术,介绍了几个重要序列数据格式的规范。在此基础上,通过两个简单的例子,介绍了测序数据分析的基本流程。

1.3 贯穿本书的例子

为了更好地对相关技术与工具资源进行介绍,本书还选用了 J. Pevsner 在其《生物信息学与功能基因组学》中的两个具有代表性的例子,即一个基因和它的对应蛋白质产物——视黄醇结合蛋白(*retinol-binding protein, RBP4*)以及人类免疫缺陷病毒 1(*HIV-1*)的 *pol*(*polymerase, 聚合酶*)基因,进行示例性介绍。

(1) 视黄醇结合蛋白

视黄醇结合蛋白是一个相对分子质量小、被大量分泌的蛋白质,能结合血液中的视黄醇(维生素 A)。视黄醇可从胡萝卜中以维生素 A 的形式获得,疏水程度大。RBP4 帮助转运这个配体到眼睛,为视觉系统所用。它有一系列有趣的性质:① 在多个物种中有许多蛋白质和 RBP4 同源,包括人、小鼠和鱼(“直系同源”)中的蛋白质。② 也有许多人类蛋白质和 RBP4 紧密相关(“旁系同源”),它们和 RBP4 的家族称为 *lipocalin* 家族——一群多样的小配体结合蛋白,它们倾向于分泌到细胞外空间。一些 *lipocalin* 蛋白具有的功能与 RBP4 不同,例如结合胆固醇(如 apo-

liprotein D)、与妊娠相关(如 pregnancy-associated lipocalin)、催欲(如仓鼠的催欲蛋白)和气味结合(如黏液中的气味结合蛋白)等。③ 有细菌的 lipocalin 蛋白,它们在对抗生素的抗性中起作用。编码细菌 lipocalin 的基因可能是一古老基因,它通过水平基因转移的过程进入真核生物基因组。④ 一些 lipocalin 蛋白的表达水平受到显著的调控。⑤ lipocalin 蛋白小而丰富,并且是可溶性的,它们的生物化学性质已被详细研究,许多蛋白质的三维结构也以 X 线晶体衍射的方法被解析出来。⑥ 一些 lipocalin 蛋白和人类疾病相关。

(2) 人类免疫缺陷病毒

人类免疫缺陷病毒(HIV)是当今世界上最大的公共卫生挑战。HIV-1 基因组仅编码 9 种蛋白质,包括 *pol*。*pol* 基因的特性、其蛋白质产物以及 HIV-1 基因组具有显著特点:① *pol* 基因编码一种 1 003 个氨基酸的蛋白质。该蛋白质是一多结构域蛋白质:单条肽链但有多个结构和功能不同的结构域。② *pol* 蛋白有反转录酶活性(即 RNA 依赖的 DNA 多聚酶),它也是天冬氨酸蛋白酶,并且还有整合酶(integrase)的活性。有多种活性是多结构域蛋白的典型特征。③ *pol* 蛋白的模块化特点会影响数据库搜索和多序列比对。④ *pol* 基因以相当快的速度发生碱基替换。一个典型的被 HIV 感染的个体可能会有百万种以上的 *pol* 变种。

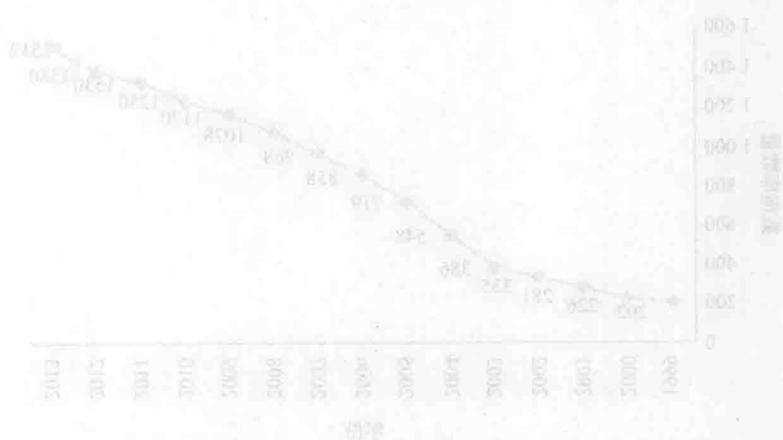


图 1.1 1999 至 2013 年 NCBI 数据库中 HIV-1 *pol* 基因序列的数量

基因组学的应用 (5)

随着测序技术的飞速发展,基因组学在医学、农业、工业和基础科学研究中发挥着越来越重要的作用。基因组学的应用包括:① 疾病诊断和预后:通过分析基因组变异,可以识别与疾病相关的基因,为疾病的诊断和预后提供依据。② 个性化医疗:根据个体的基因组信息,制定个性化的治疗方案,提高治疗效果。③ 药物研发:通过基因组学技术,可以发现新的药物靶点,加速新药的研发进程。④ 农业育种:利用基因组学技术,可以选育优良品种,提高农作物的产量和品质。⑤ 工业生物技术:通过基因组学技术,可以改造微生物,生产高附加值的产品。⑥ 基础科学研究:基因组学为研究生物进化和发育提供了重要的工具,有助于揭示生命的奥秘。