



普通高等教育“十一五”国家级规划教材

化学信息学

李梦龙 文志宁 熊庆 编



化学工业出版社

普通高等教育“十一五”国家级规划教材

化学信息学

李梦龙 文志宁 熊庆 编



化学工业出版社

· 北京 ·

本书为普通高等教育“十一五”国家级规划教材，是教育部“使用信息技术工具改造课程”项目的研究成果。全书主要分为四大部分，其中第1章概述了化学信息学的产生及特点；第2~4章讲述了化学信息的来源，包括手册、书籍、搜索引擎以及目前广为使用的期刊文献数据库；第5~7章介绍了化学信息的处理工具（即化学软件）、处理方法（相关化学计量学算法）以及定量构效关系（QSAR）的原理及应用；第8章对生物信息学领域的研究进行了概述。

本书可作为高等院校化学化工专业本科“化学信息学”课程的入门教材，另外，书中提供了大量与信息学相关的网址，也可作为研究生的参考书籍。

图书在版编目（CIP）数据

化学信息学 / 李梦龙，文志宁，熊庆编. —北京：化学工业出版社，2011.6

普通高等教育“十一五”国家级规划教材

ISBN 978-7-122-11203-3

I. 化… II. ①李… ②文… ③熊… III. 计算机应用—化学—情报检索—高等学校—教材 IV. G252.7

中国版本图书馆 CIP 数据核字（2011）第 080602 号

责任编辑：杜进祥

文字编辑：昝景岩

责任校对：陈 静

装帧设计：韩 飞

出版发行：化学工业出版社（北京市东城区青年湖南街 13 号 邮政编码 100011）

印 装：大厂聚鑫印刷有限责任公司

787mm×1092mm 1/16 印张 12 1/4 字数 306 千字 2011 年 6 月北京第 1 版第 1 次印刷

购书咨询：010-64518888（传真：010-64519686）售后服务：010-64518899

网 址：<http://www.cip.com.cn>

凡购买本书，如有缺损质量问题，本社销售中心负责调换。

定 价：26.00 元

版权所有 违者必究

前 言

HUAXUE XINXIXUE
化学信息学

化学是一门以实验为基础的古老学科，科学家们致力于探索新物质的各类化学属性，随着实验技术的进步，人们获取数据的能力已有了很大的提高。时至今日，面对呈指数级增长的化学数据，化学家们发现在本学科的探索中，最大的瓶颈已不再是如何获取未知物质的性质，而是如何从已有的实验数据中提取更多的化学信息，总结规律。1995年，美国著名化学家 Brown 曾指出，化学家习惯于将 99% 的精力和资源用在数据的收集上，只余下 1% 用于数据的分析和处理，将其转化为信息。

庆幸的是，现在越来越多的化学工作者们开始注意到这个问题并投身于化学数据的分析当中，化学信息学在这种情况下应运而生，它是汇集化学、数学、信息科学等交叉学科知识的研究领域，其主要是通过对化学信息的检索、整理、分析以及可视化，最终完成将数据转化为信息的过程。2003 年德国 Johann Gasteiger 和 Thomas Engel 出版的 “Chemoinformatics A Textbook”（中文版为《化学信息学教程》）一书亦指出，化学信息学的任务就是运用信息学的方法来解决化学的问题。

目前，化学信息学主要涵盖了化学信息的获取、化学信息的表达以及化学信息的处理三个方面的内容。作为一门新的基础课程，如何尽快地让化学专业的本科生了解并掌握其中涉及的概念、方法以及网络资源是该学科建设亟待解决的主要问题。本书的出版正是基于这一目的，系统列举了各类文献信息、网络化学数据资源；对于常用的化学信息软件，如 ChemBioOffice、MATLAB、Amber 等，亦做了概要的介绍；在化学信息处理方面，我们详细阐述了使用频率较高的模式识别及 QSAR 等方法的基本原理；由于生物信息学近来发展迅速，在药物设计等方面与化学信息学也有交叉，因此在本书的最后，我们对生物信息学进行了简单概述。此书可作为化学专业本科生的普及教材使用。

感谢清华大学图书馆的战玉华老师撰写了本书 3.4 节的部分内容以及在后期修改过程中提出了宝贵的意见。本书在编写过程中还得到了实验室蒲雪梅副教授、王智猛副教授、印家健副教授、郭延芝博士、刁元波博士、方亚平博士的热心帮助；实验室的博士研究生李益洲、杨刚、孙婧、余乐正以及硕士研究生张娟、唐小净、朱丽娟、孟艳艳、田雪、李功兵、胡美、谭颖、尹辉、王翠翠、罗杰斯、吴镝、敬闰宇、刘雯、张丽芳、肖秀婵、杨魏、林娇等同学亦参与了本书资源的收集与整理工作；另外，化学工业出版社在此过程中提供了大力支持，在此一并表示诚挚的感谢。

由于化学信息学涉及面广，编者的水平和时间有限，书中不妥之处在所难免，恳请广大读者批评指正。

编 者
2011 年 4 月

目 录

HUAXUE XINXIXUE
化学信息学

第Ⅰ章 概述 1

1.1 什么是化学信息	1
1.2 化学信息的诞生背景	1
1.3 信息科学在化学领域的应用	2
1.4 化学信息的结构和特点	2
1.5 化学信息的工作方式	3
1.6 信息采集接口	4
1.7 化学信息的应用	4
1.7.1 绘制结构	4
1.7.2 数据库	5
1.7.3 计算机辅助设计反应预测系统	5
1.7.4 预测结构与活性的关系	5
1.8 展望	5

第Ⅱ章 化学信息来源 7

2.1 词典	7
2.2 手册	7
2.3 化学期刊	9
2.3.1 综合类期刊	9
2.3.2 有机化学期刊	10
2.3.3 分析化学期刊	11
2.3.4 无机化学期刊	12
2.3.5 物理化学期刊	12
2.4 图书馆资源	12
2.4.1 生命科学图书馆	13
2.4.2 中国科学院大连化学物理研究所图书馆	13
2.4.3 中国科学院国家科学图书馆	14
2.4.4 国家科技图书文献中心化工分中心	16
2.4.5 清华大学图书馆	17

2.4.6 中国国家图书馆.....	17
2.4.7 哈佛大学图书馆.....	17
2.4.8 斯坦福大学图书馆.....	18
2.5 化学化工信息资源导航系统	19
2.5.1 ChIN	19
2.5.2 Computer Aided Chemistry Tutorial	20
2.5.3 Wilton High School Chemistry	20
2.5.4 化学家链接网站.....	21
第3章 化学信息数据库资源	22
3.1 数据库简介	22
3.1.1 数据	22
3.1.2 数据库	22
3.1.3 数据库管理系统.....	23
3.1.4 数据库系统	23
3.2 数据库历史及分类	24
3.2.1 数据库历史	24
3.2.2 数据库的模型分类.....	25
3.3 三类化学信息数据库	26
3.3.1 文献数据库	26
3.3.2 事实数据库	26
3.3.3 结构数据库	26
3.4 互联网上的化学化工数据库	27
3.4.1 CA.....	27
3.4.2 ISI 数据库	33
3.4.3 OCLC 数据库.....	39
3.4.4 CSA	40
3.4.5 ScienceDirect	40
3.4.6 CNKI	42
3.4.7 万方数据库	48
3.4.8 维普中文科技期刊数据库.....	48
3.4.9 EI	48
3.4.10 出专利数据库	49
3.4.11 Reaxys 数据库.....	51
3.4.12 谱图数据库	52
第4章 信息搜索引擎	54
4.1 概述	54
4.1.1 搜索引擎的原理.....	54
4.1.2 搜索引擎的历史及发展趋势.....	55

4.2 搜索引擎的定义及分类	58
4.2.1 全文搜索引擎	58
4.2.2 目录索引类搜索引擎	58
4.2.3 元搜索引擎	59
4.2.4 垂直搜索引擎	59
4.3 搜索引擎查询方法	59
4.3.1 模糊查询	60
4.3.2 精确查询	60
4.3.3 逻辑查询	60
4.3.4 指定范围查询	61
4.4 常用搜索引擎	61
4.4.1 百度	61
4.4.2 Google 中国	62
4.4.3 维基百科	62
4.4.4 BASE	65
4.4.5 Vascoda	65
4.4.6 Information Bridge	66
4.4.7 Intute	67
4.4.8 Infomine	68
4.5 元搜索引擎	68
4.5.1 Dogpile	68
4.5.2 Excite	69
4.5.3 Ixquick	70
4.5.4 Mamma	70
4.5.5 Metacrawler	71
4.5.6 ProFusion	71
4.5.7 Savvysearch	72
4.6 专业搜索引擎	73
4.6.1 专业搜索引擎的优势	73
4.6.2 著名的专业搜索引擎	73
第5章 化学软件	76
5.1 概述	76
5.2 化学软件的分类	77
5.3 语言软件和依托算法的化学计算软件	78
5.3.1 MATLAB	78
5.3.2 R 语言	91
5.4 绘图软件	101
5.4.1 ACD/ChemSketch5.0	101
5.4.2 Symyx Draw	103
5.4.3 ChemBioDraw	104

5.5 化学分析仪器数据处理软件	105
5.5.1 GRAMS	106
5.5.2 MestReNova	109
5.5.3 Origin	110
5.6 分子模拟软件	112
5.6.1 Gaussian 软件	112
5.6.2 Amber 软件	114
第6章 信息处理与数据挖掘	117
6.1 概述	117
6.2 数据的标准化	118
6.3 特征提取与优化	118
6.3.1 主成分分析	118
6.3.2 偏最小二乘法	121
6.3.3 逐步回归分析	122
6.3.4 遗传算法	123
6.4 信号处理方法	125
6.4.1 协方差与相关系数	126
6.4.2 自、互相关分析	126
6.4.3 功率谱密度	127
6.4.4 傅里叶变换	127
6.4.5 小波变换	128
6.5 机器学习方法	132
6.5.1 K 最近邻法	132
6.5.2 概率神经网络	132
6.5.3 分类回归树	133
6.5.4 助推法	134
6.5.5 人工神经网络	135
6.5.6 支持向量机	139
6.6 数据库挖掘技术	141
6.6.1 聚类算法	141
6.6.2 决策树算法	142
6.7 Web 数据挖掘技术	142
6.7.1 web 内容挖掘	142
6.7.2 web 结构挖掘	143
6.7.3 web 日志挖掘	143
第7章 QSAR 及药物设计	144
7.1 概述	144

7.2 QSAR 模型的分类	145
7.2.1 二维定量构效关系	145
7.2.2 三维定量构效关系	147
7.2.3 多维定量构效关系	150
7.2.4 方法评价	150
7.3 定量构效关系研究中常用的回归分析法	151
7.3.1 多元线性回归	151
7.3.2 主成分回归	152
7.3.3 偏最小二乘回归	153
7.3.4 投影寻踪回归	154
7.3.5 非线性方法	155
7.4 药物设计	155
7.5 QSAR 方法的应用	157
第8章 生物信息学	161
8.1 什么是生物信息学	161
8.2 生物信息学的发展历程	162
8.3 生物信息学的研究内容	164
8.3.1 生物信息挖掘	164
8.3.2 药物设计	164
8.3.3 基因组学	165
8.3.4 蛋白质组学	165
8.4 生物信息学的研究方法	167
8.5 生物信息学的应用	168
8.6 生物信息学的研究趋势	169
8.7 蛋白质功能研究	170
8.8 蛋白质数据库简介	171
8.8.1 综合性蛋白质数据库	171
8.8.2 专用性蛋白质数据库	172
8.9 蛋白质序列的特征提取方法	173
8.9.1 基于氨基酸组成和位置的特征提取方法	174
8.9.2 基于氨基酸物理化学特性的特征提取方法	175
8.9.3 基于数据库信息挖掘的特征提取方法	177
8.10 蛋白质相互作用	178
8.11 蛋白质网络	183
参考文献	187

第1章

概 述

材料、能源和信息是构成物质世界的三个基本要素。随着社会发展的需要，人们逐渐认识到信息的重要性，并创立了信息论与信息科学。20世纪90年代初，随着美国“信息高速公路”计划的提出，信息科学和信息产业出现了前所未有的飞速增长，成为这一时代的重要标志。同时信息科学加快了向传统科学渗透，化学中的信息学理论基础不断成熟。正是在这一背景下，结合其使用的计算机和互联网工具，化学工作者在科研实践中促成了化学信息这一新兴化学分支的出现。化学信息学(cheminformatics, chemoinformatics, chemical informatics)是化学领域中近几年发展起来的一个新的分支，是建立在多学科基础上的交叉学科，利用计算机技术和计算机网络技术，对化学信息进行表示、管理、分析、模拟和传播，以实现化学信息的提取、转化与共享，揭示化学信息的实质与内在联系，促进化学学科的知识创新。

1.1 什么是化学信息

1987年诺贝尔化学奖获得者、法国化学家J.M.Lehn在获奖报告中首次提出化学信息的概念，对化学的发展而言具有深远的影响，具有深刻的时代意义。虽然众多的化学工作者没有对化学信息展开实质性的工作，但是传统有机化学、无机化学、生物化学、材料化学以及在受体设计、超分子形成过程的结构化学等方面所积累的大量实验数据，却为构建化学信息提供了基础。在今天，化学信息学处于一种呼之欲出的形势，它将给21世纪的化学带来全新的面貌。

化学信息是个广义的概念，它包含对化学相关信息的设计、创造、组织、存储、处理、恢复、分析、再开发、可视化及应用。另一种关于化学信息的定义是从药物研发的角度提出的，认为化学信息是各种信息资源的混合体，目的是将数据转化成信息，再把信息转化成知识，以期更快、更准确地进行药物筛选和设计。

1.2 化学信息的诞生背景

近十年来，由于计算机及网络技术向智能化、网络化方向发展，应用计算机技术能解决的化学问题也愈来愈多，化学工作者不仅获得了大量物质结构的信息，而且这些信息较之从前也更为精确，计算机技术与化学之间的相互渗透已成为化学和计算机科学工作者的研究热点。由于计算机主要是通过数值计算来解决问题，其特点是能快速地进行大量复杂、繁琐的数学运算，而化学是对化学物质进行认识、分析、合成及利用，从而使化学工作者能够对物质化学结构进行解析、表征、模拟与设计，能够处理复杂体系的电子结构、几何结构与其性



能关系，完成微观分子工程设计与化学模拟，开展功能材料的研究，进行生物活性分子和药物分子的相互作用机制及定量构效关系（quantitative structure-activity relationship, QSAR）研究，探讨固态表面结构、固体表面轨道相互作用规律以及实现分子以上层次聚集体（超分子体系、界面体系等）结构和性能的模拟等。

然而，纵观早期这一领域的工作仅仅涉及计算机技术的一些应用层次，要想将计算机技术深入应用到化学中就必须解决化学与计算机的结合问题，从化学工作者的角度应用和设计计算机软、硬件，满足化学工作者处理化学信息的要求。该领域的研究包括计算机与分析仪器的接口、化学类应用软件程序包的开发、化学物质结构数据库的开发和查询。

1.3 信息科学在化学领域的应用

1948 年，Shannon 发表了关于信息论的著名文章，提出了信息熵计算公式

$$H(X) = -\sum_{i=1}^n \frac{1}{p(x_i)} \log_2[p(x_i)] \quad (1-1)$$

式中， $H(X)$ 为事件 X 的信息熵，它可由该事件当中所有可能出现的情况 x_i 的概率 $p(x_i)$ 计算得到。

信息理论开始了它的发展，这一理论最早是与通信技术相关联的，但在其诞生后 10 年左右，即从纯粹数学研究渗透到无线电、电视、雷达、心理学、语义学、经济学、生物学等领域。Wiener 认为信息的实质是负熵，并强调信息这种负熵是在调节过程中相互交换而产生的。

化学科学中的分析化学从其诞生起就具有信息科学的特征，Kateman 等从三个方面阐述分析化学的任务：利用已有的分析方法，提供关于物质化学成分的信息（日常例行分析工作）；研究利用不同学科的原理、方法，取得有关物质系统的相关信息的过程（分析化学的科学研究工作）；研究利用现有分析方法取得关于物质系统的信息的策略（分析实验室的组织工作）。Kowalski 更是明确提出《分析化学作为信息科学》的论文，他认为分析化学不仅在过去是一门信息科学，现在仍然是一门信息科学，在化学的各个分支学科中，分析化学负担的任务与其他分支学科的不同之处在于分析化学的研究对象，它并非某种具体的实物，例如无机或有机材料，而是与这些化学组成或结构相关的信息以及研究获取这些信息的最优方法与策略。

此外，由于化学中熵的概念与 Shannon、Wiener 等提出的信息理论中的熵有着共同的基础，这两门学科之间存在着深刻联系，分析化学从发展分析信息理论作为其基础理论的组成部分，获得了向前发展的动力。随后众多的化学家根据其从事的分析化学工作，发表了多篇用于分析化学的信息理论系列论文，其中捷克学者 Eckschlager 完成了在此领域的第一部专著。

1.4 化学信息的结构和特点

仔细分析可以发现化学信息的主体部分实际上是由三个层次构成的，即信息核心层、信息处理层和信息表示层，如图 1-1 所示，化学信息的这种分层结构本质上是计算机技术分层结构的反映。



在化学实践中产生出来并被计算机处理过的原始化学信息，比如在一个化学实验中发生的各种实验现象、实验数据以及与化学实验相关的外界条件等，它们组成了化学信息的信息核心层。信息处理层由化学计量学方法、药物分子设计方法、QSAR 方法等能够对信息核心层中的数字化化学信息进行二次开发利用的计算方法组成。表示层处于化学信息学科的最外层，它根据信息核心层的特定要求在计算机信息科学中寻找适合表达化学信息的技术，从多个角度将化学信息以某种直观的形式如基于计算机的图形、音频、视频等多媒体表示手段向化学工作者展示出来。信息处理层和信息表示层统称为化学信息学科的外层。

三个层次中最重要的层次为信息核心层，它在化学信息学科中处于基础核心地位，并决定了其他两个层次的构成。核心层对外层起着决定性作用，外层对核心层也能够产生一定的影响。由于计算机信息科学技术和仪器设备本身存在某些特点的影响，它要求信息核心层中的化学信息必须要按照一定的要求进行组织和编排。例如，连续吸光度曲线在计算机中只能以离散数据点的形式存储，海量的光谱数据、化学化工测量数据或分子结构参数唯有按一定的数据结构规范化并形成数据库甚至是专家系统，才能方便日后使用处理层对这些数据进行开发利用；此外，信息处理层在对来自信息核心层的化学信息进行处理之后，所获得的结果一方面将交由信息表示层处理，另一方面，信息处理层将把某些处理结果和原有数据存储在信息核心层，使该层信息量甚至是一些局部结构发生变化。

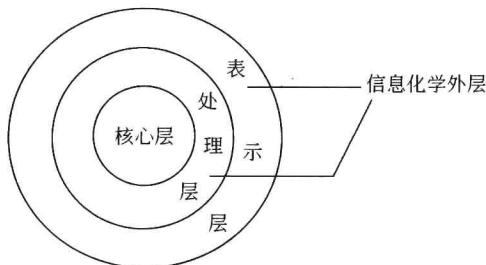


图 1-1 化学信息的结构示意图

1.5 化学信息的工作方式

化学实践与化学信息之间的关系是“母与子”的关系。化学信息通过“信息采集接口”从化学实践这一母体中获取原始信息，原始信息经过接口的处理后以数字化形式存储到信息核心层中，并通过外层将其重现出来给化学工作者。利用这些被数字化技术处理过的化学信息，化学工作者可以进行更深层次的化学研究实践，从而生产出新的数字化的原始化学信息，这就是化学信息的工作方式，如图 1-2 所示。

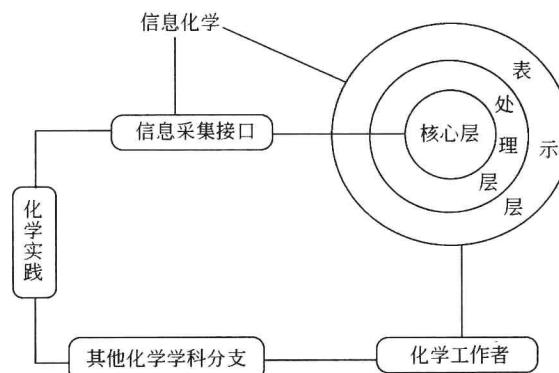


图 1-2 化学信息的工作方式



化学信息的这种工作方式与其他化学学科分支的工作方式基本相同，最大的区别在于化学信息工作者手中的研究工具是以计算机程序表达的化学计量学、QSAR 等可以对大量化学数据进行二次利用的计算化学方法，其研究的直接对象和得到的产品都是数字化了的化学信息。化学信息工作者的任务一是利用现有的计算机软硬件工具研究大量存储在信息核心层的原始化学数据，找出不同化学变量之间的关系，发现有实际意义的化学规律；二是改进现有的研究方法、开发新的研究手段以更新和完善化学信息外层以及不断丰富和修正信息核心层中的化学数据，为今后开展更深层次的研究工作奠定基础。

1.6 信息采集接口

从图 1-2 中可以发现，化学信息和化学实践之间是通过一个信息采集接口相连的，这与其他化学学科分支明显不同。信息采集接口也是化学信息学科一个极为重要的组成部分，它是现实化学世界通往化学信息的桥梁，也是化学信息的生命源泉。信息采集接口和化学信息的三个层次构成了整个化学信息学科。

如图 1-3 所示，信息采集接口不是一套单一而又具体的软件或硬件，它实际上是一个综合了其他化学学科分支的某些方法以及原始化学信息数字化方法的方法集合。

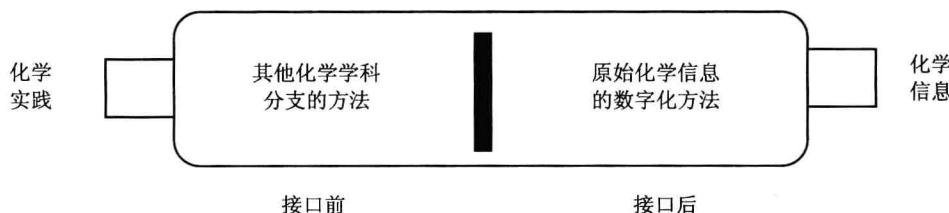


图 1-3 信息采集接口的内部结构

可以用作信息采集接口后端的方法范围非常广泛，一般各种化学图形处理软件、计算化学应用软件、文字处理软件、数码照相机、具有 OCR (optical character recognition) 光学文字识别功能的扫描仪的录入系统等都可以作为信息采集接口的后端，因为它们都具有一个共同的特点，即能够把现实世界中的化学信息以数字化形式存储起来。对于接口的前端，只要是能从化学实践中获取化学信息的研究方法或仪器设备如化学分析测量仪器，量子化学计算方法，分析化学、物理化学实验方法等，都可以用作信息采集接口的前端。

1.7 化学信息的应用

1.7.1 绘制结构

化学信息普遍存在于化学和计算机的结合之中，几乎每一个化学家都是一个绘制结构者，都会使用到 IsisDraw、ChemDraw、JchemPaint 等相关软件去绘制一个分子二维或三维结构，然而，更深入的问题则需要化学信息学的方法来解决，例如如何将化学的结构有效地储存在计算机里，采用何种格式可以用来在不同类型的化学软件中交换数据等。

1.7.2 数据库

化学数据库的开发、维护和更新是化学信息的重要领域。化学结构和一些相关的信息主要被存放在化学数据库中，如 Beilstein 数据库，自 1771 年起，数据库中存放了超过九百多万个有机化合物的信息可以由 CAS 登记号或分子式查询物质的物理和化学性质，包括光谱数据以及热力学参数。另外在构建这些数据库时，基于物质化学数据信息的数据库查询功能同样至关重要。

1.7.3 计算机辅助设计反应预测系统

人们已经做了很多的努力去实现用计算机预测化学反应的进行，模拟化学反应的发生，来合成一个设计好的目标化合物。目前，德国 Gasteiger 的化学信息研究小组，已经有这样一个系统，名叫 ERDS (Elaboration of Reactions for Organic Synthesis)，能够进行包括两种反应物之间的反应结果的预测，或提供采用何种反应物的建议。

1.7.4 预测结构与活性的关系

QSAR 定量构效关系方法尝试通过对一系列结构相似的药物分子进行分析，找出分子性质参数和生物活性之间的关系，并以此为依据去预测具有药效的新型分子的结构与性质。目前发展到三维的 3D-QSAR 实际上是 QSAR 与计算机分子图形学相结合的研究方法，是研究药物与受体间的相互作用，推测受体的图像及进行药物设计的有力工具。3D-QSAR 研究可分为受体结构已知及受体结构未知两种情况。受体结构已知（目前仅限于酶作为受体），可以根据 QSAR 的结果及计算机图形显示受体的三维结构，并随之进行有如“量衣裁衣”式的设计。在受体结构未知的情况下（这是绝大多数情况），则可以根据激动剂或（和）拮抗剂的构效关系及计算机图形显示的化合物优势构象，推测受体的结构，然后进行药物设计，也可以起到“量体裁衣”的作用。

1.8 展望

随着信息时代的到来，各类化学信息的相关数据不断涌现，这些数据在使我们获得更多信息的同时也在信息的分类、分析以及有效应用方面带来了巨大的挑战。化学信息学伴随着计算机技术的发展应运而生，它针对海量的化学信息进行管理、分析，以实现对信息的提取、共享以及应用。

目前一些大的化学公司已经注意到这种挑战，它不仅是对化学信息的存储，而且还需要建立一套体系对化学信息进行分析与处理，例如实验室信息管理系统（LIMS）等。因此，从事化学信息学的研究人员需要掌握多学科、跨专业的知识，它的发展不是光靠某一个人或某一个研究团体就能够做到的，必须依赖于各学科各专业科研人员的通力合作。

在化学研究中，常常需要确定化合物的结构或者设计反应的过程。可以利用化学信息学的原理，建立相应的数据库，基于计算机技术对已有数据进行统计分析，进而预测化学结构与功能、完成反应流程设计、建立相应的模型以及开发专用的软件。

另外，药物设计已成为众多科研工作者的研究热点，利用化学信息学的技术，对药效团分子及先导化合物进行预测和筛选仍然是该类研究的重点。

化学信息体现了现代科学的研究由各分支独立发展走向纵横交错共同发展的大趋势，是顺



应时代潮流的一门新型学科，我们坚信，在各领域研究人员的共同努力下，化学信息学一定能够健康地发展下去。



扩展阅读

科学引文索引与影响因子

科学引文索引（Science Citation Index, SCI）的创立基于 Eugene Garfield 的引文思想，于 1961 年由美国科学情报所（Institute for Scientific Information Inc., ISI）出版，逐渐成为国际性检索刊物。被 SCI 收录的论文数量和论文质量是衡量一个国家科研实力和研究水平的重要指标。“越查越旧，越查越新，越查越深”是科学引文索引建立的宗旨，通过登录 Web of Science 数据库可对指定关键词进行溯源或查新。SCI 从来源期刊数量划分为 SCI 和 SCI-E。SCI 指来源刊为 3500 多种的 SCI 印刷版和 SCI 光盘版（SCI Compact Disc Edition, SCI CDE）。SCI-E（SCI Expanded）是 SCI 的扩展库，收录了 6500 多种来源期刊，可通过国际联机或因特网进行检索。

科学引文索引（Science Citation Index, SCI）、工程索引（The Engineering Index, EI）、科技会议录索引（Index to Scientific & Technical Proceedings, ISTP）是世界著名的三大科技文献检索系统。

影响因子（Impact Factor, IF）反映了被 SCI 收录源期刊文章平均被引用次数，是衡量一本期刊在某一研究领域是否处于领先行列的重要指标。相同研究领域中，影响因子越高的期刊，影响力越强。某期刊特定年度的影响因子，等于该期刊前两年中所有论文在当年的总被引用次数除以这两年中发表的文章总数。例如：IF（2011 年）=A/B，其中 A=该期刊 2009~2010 年所有文章在 2011 年中被引用的总次数，B=该期刊 2009~2010 年所有文章数。一般而言，当年只能发布上一年度的期刊影响因子。对于新创立的期刊，从能被检索到的时间算起，两个自然年后会拥有相应的影响因子，在这期间，影响因子计为 0。综述类论文的引用次数明显高于研究类论文，因此综述类期刊或者是出版较多综述类论文期刊的影响因子较高。IF 也是美国科学情报所的一项重要数据，通过登录 Web of Science 数据库可查询期刊的 IF 值。

第2章

化学信息来源

如何获取化学信息是进行化学信息学研究的一个重要方面，传统的化学文献或数据的获取方式主要来源于纸质的资料，如手册等，随着计算机技术的发展以及互联网的普及，信息的存储方式发生了革命性的改变，出现了以磁盘、光盘为存储介质的信息资料，从而使我们不仅可以从图书馆提供的光盘版信息资源中获取相应的标准实验数据和化学文献，还可以通过互联网在办公室或者家中方便地查阅各类化学资源。

2.1 词典

词典是汇集事物词语，解释词义、概念、用法，并且按一定次序编排，以备检索的一类最基本、最常用的工具书。下面列有最常用的化学化工词典。

(1)《英汉化学化工词汇》由科学出版社出版，2003年第四版中收录了与化学化工有关的科技词汇约17万条，除词汇正文外还附有常用缩写词、无机和有机化学命名原则等。

(2)《化工辞典》由化学工业出版社出版，2002年第四版，收集化学化工的专业名词一万多个，解释简明扼要，是中国影响力最大的中型化工专业工具书。

(3)《化学化工大辞典》2003年1月由化学工业出版社出版，是中国规模最大的化学化工类综合性专业辞书，也是目前我国收词量最多、专业覆盖面最广、解释较为详细的化学化工专业词典。

(4)《化工百科全书》由化学工业出版社于1997年出版，全书共19卷，索引1卷，全面介绍了化工领域最新的技术和发展趋势，该书学术性强、覆盖面宽、产业气息浓、实用性高，是一本大型的化学工业及其相关工业技术的百科全书。

(5)《英汉双向精细化工词典》由上海交通大学出版社于2009年9月出版，该词典在包含传统、基本精细化工词汇的基础上收录了最新的有代表性的词汇，比较全面地反映了国内外精细化工领域的最新发展，收词约40000条。

(6)《化合物词典》由上海辞书出版社于2002年6月出版，该词典包括无机化合物和有机化合物两部分，无机化合物部分收词3929条，有机化合物部分收词2652条，为便于查找，书末附词目英汉对照索引。该词典简明扼要、收词面较广、内容较丰富，是简明而实用的化合物专业词典。

(7)《化学辞典》由化学工业出版社2011年出版第二版，全书共选词目近8000条。

2.2 手册

手册是按照某一学科或某一主题汇集需要经常查考的资料，供读者随时翻检的工具书。



对于化学工作者，手册可称之为他们的知识仓库，是他们必不可少的工具书。各国出版的关于化学的手册品种繁多，本书主要给大家介绍一些重要的手册。

(1)《新药化学全合成路线手册》由科学出版社于2008年7月出版，这本手册主要介绍了美国食品与药品管理局(FDA)于1999~2007年批准上市的170余个新分子实体药物的化学合成方法。并对每个药物给出其英文名、中文名、化学结构、化学式、相对分子质量、化学元素分析、药物类别、美国化学会CAS登记号、申报厂商、批准日期、适应症、药物基本信息等。其中，“药物基本信息”部分介绍了对应药物的作用机制、结构信息、合成路线等。这些合成路线大多是目前制药工业中正在使用的生产工艺，有较高的实用性与学术价值。该书共包含了数千个有机合成反应，数百种药物中间体的合成制造方法和数个非常有参考价值的附录。

(2)《分析化学手册》第二版由化学工业出版社出版，包含以下10个分册：基础知识与安全知识、化学分析、光谱分析、电分析化学、气相色谱分析、液相色谱分析、核磁共振波谱分析、热分析、质谱分析和化学计量学。手册涉及的内容包括方法的基本原理、应用技术与重要的应用资料，相关定义、术语及符号等。此外手册还介绍了因特网上的分析化学资源及获取方法。该书为从事分析化学相关工作的技术人员提供了大量丰富翔实的资料，是一部实用性很强的手册。

(3)《兰氏化学手册》[美]J.A.迪安，N.A.兰格(Lange)著。于2003年5月由科学出版社出版第二版译本，这是一部资料齐全、数据翔实、使用方便、供化学及相关科学工作者使用的单卷式化学数据手册，在国际上享有盛誉。自问世以来，一直受到各国化学工作者的重视和欢迎。全书共分11部分，内容包括有机化合物，通用数据，换算表和数学，无机化合物，原子、自由基和键的性质，物理性质，热力学性质，光谱学，电解质、电动势和化学平衡，物理化学关系，聚合物、橡胶、脂肪、油和蜡及实用实验室资料等。书中所列数据和命名原则均取自国际纯粹化学与应用化学联合会最新数据和规定。化合物中文名称按中国化学会1980年命名原则命名。该手册是从事化学方面工作的必备工具书。

(4)《工业聚合物手册》原著是由美国的威尔克斯根据国际权威的《工业化学百科全书(第六版)》编著而成的，于2006年出版中文译本。这是一本关于聚合物的综合性图书，囊括了所有重要的工业聚合物，对其生产基本原理、主要生产过程、表征和应用领域等进行了详尽的介绍，对各种聚合物的发现、发展和现状进行了系统的评述；并对逐步聚合、连锁聚合、开环聚合所得聚合物和其他一些特种聚合物均有详细描述。同时，该手册对天然聚合物及其衍生物也进行了充分介绍。这本手册的技术新颖，实用性强，对从事聚合物科研与生产的工程技术人员是难得的宝贵资料。

(5)《Gmelin无机化学手册》(Gmelin Handbook of Inorganic Chemistry)是世界上最最有威望和最完整的一套无机化合物手册，原名为《理论化学手册》，后来增加了一些内容，又称《Gmelin Handbook of Inorganic and Organometallic Chemistry》。1922年开始出第八版，这套书是西文图书这个家族中的一个大的阵营，该书对各元素及其无机化合物都加以讨论，包括历史、存在、性质、实验室及工业制法，与无机化学相关的许多领域也都包括在内，并引用大量参考文献。1998年，该手册停止出版，全部改为了电子版和网络化检索。

(6)《Beilstein有机化学手册》(Beilsteins Handbuch der Organischen Chemie)是在德国化学会的支持下编著的，是当前国际上最系统、最全面、最权威的有机化合物巨型手册。Beilstein手册包括正编和补编共计566册。收集了各种有机化合物的来源、结构、制备、物理和化学性质、化学反应、化学分析、用途及其衍生物等内容。各种有机物是按结构分类编