



计算机化学化工丛书

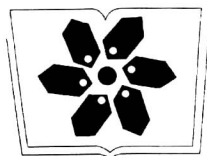
Computer Chemistry and Chemical Engineering Series

药物设计中的 分子模型化方法

周家驹 王 亭 著



科学出版社



中国科学院科学出版基金资助出版

计算机化学化工丛书

药物设计中的分子模型化方法

周家驹 王亭 著

科学出版社

2001

内 容 简 介

本书为《计算机化学化工丛书》之一。

本书以小分子的化学结构信息处理为核心,介绍药物设计中的分子模型化方法。全书分为九章:前三章为第一部分,是一些通用的方法、原理和技术,分别介绍分子结构的计算机处理、分子的相互作用、分子结构信息数据库;后六章为第二部分,分别介绍人工智能中的遗传算法,定量结构活性关系(QSAR)研究、结构匹配和数据库搜索技术、虚拟受体模型方法、药物与受体的对接技术和模拟肽学等六个专题。

本书是一本面向应用、介绍计算机辅助药物分子设计方法的入门书,适合从事计算机化学的研究人员及从事新医药、新农药研制与开发的研究人员阅读和参考,也可供化学各专业大学生及研究生阅读。本书有助于读者了解计算机处理化学信息的基本方法和用于分子设计的各种实用技术。

图书在版编目(CIP)数据

药物设计中的分子模型化方法/周家驹,王亭著. —北京:
科学出版社, 2001

(计算机化学化工丛书/许志宏主编)

ISBN-7-03-0098690-X

I. 药… II. ①周… ②王… III. 药物-设计-分子结构-化学模型 IV. TQ460

中国版本图书馆CIP数据核字(2000)第81844号

科学出版社 出版

北京东黄城根北街16号

邮政编码:100717

源海印刷厂 印刷

科学出版社发行 各地新华书店经销

*

2001年4月第一版 开本:850×1168 1/32

2001年4月第一次印刷 印张:77/8

印数:1—3000 字数:198000

定价:21.00元

(如有印装质量问题,我社负责调换〈北燕〉)

《计算机化学化工丛书》

编 委 会

主 编 许志宏

副主编 杨小霞

编 委 (以姓氏笔画为序)

马沛生	王淀佐	王 夔	许 禄
李 科	来鲁华	陈丙珍	陈冀胜
陈凯先	陈念贻	陈敏伯	陈维明
杨友麒	严新建	林少凡	郑崇直
周家驹	胡鑫尧	俞汝勤	郭传杰
郭 力	袁身刚	麻德贤	惠永正
潘忠孝			

《计算机化学化工丛书》序

化学化工是信息量特别大的一门学科。到 2000 年 4 月底,在 CAS 录的化学物质约 2400 万种,它们的各种性质加上多元体学的性质,可以说是一个无边无际的数据海洋。于是,化学数据库的建设就成为 20 世纪后 20 年国际上的一件大事。中国科学院从 1979 年开始建设化学数据库,迄今已经整整 20 年。其学院从 1979 年开始建设化学数据库,迄今已经整整 20 年。其间,多次得到国家和中国科学院的奖励。

长期以来,人们不仅希望能定性地掌握而且希望能定量地了解化学化工学科的规律,而计算机的能力及相关技术高速发展,正在帮助人们一步一步地实现着这个愿望。从理论化学的计算、谱图解析、化学计量学、化工过程模拟、热力学的复杂计算,都在计算机的帮助下得到了很大的发展。

上述基础工作为化学化工领域的工作者增加了很大的自由度:可以用计算机帮助对化合物的谱图解析,帮助选择合成路线,帮助进行药物分子设计,可以进行新过程、新技术的开发,可以进行大型工业装置的设计,可以对工厂的生产过程进行优化,……等等。

在改革开放的 20 年里,我们的计算机化学从无到有,形居了一支生机勃勃的研究队伍,也形成了一个很大的用户和读者群体。到 21 世纪,这个群体更需要有能力利用计算机帮助自己的工作,所以本套丛书中也包含一些计算机化学化工的教材,以利于化学化工本科生和研究生的培养和工程人员自学。

所以,我们希望通过本套丛书介绍一些解决问题的方法,帮助读者在遇到问题时,知道如何去解决问题。为此,要求作者在自己的著作中,要给出软件、数据的出处、网络地址或光盘。时代发展

得很快,仅做到这一点显然还是不够的。我们特别注意到近年来,Internet 网络的高速发展已经给我们的时代带来了巨大变化。到1999年7月,Internet 已经是一个连接5600多万个节点的一个网络系统,它上面的文本信息已经超过600亿字节。这些信息一个最大的长处是时间的滞后最少,易于通过计算机获取。

如果能将科学数据库在网上的功能,由数据的存取扩大到运算、绘图、模拟等多方面,必将极大地推动科学数据库工作的发展和广泛使用。在21世纪,将逐步可以做到,人们在用户端将数据从库中取出,在服务器端程序系统上计算,结果以图形或多媒体方式输出到用户端。据了解,我们有些作者在自己的工作中已经能够在网上实现数据查询、计算、绘图、三维图形显示等。

进入21世纪,Internet 网络系统的应用将更加普及,Internet 网络的客户/服务器的应用将进入千家万户,进入教室和办公室的各个角落。所以,如果能将科学数据库和计算程序库在网络服务器上实现,那么它的普及应用,将会随着计算机网络的推广而推广。

如果有的作者,目前仅能给出单机版本的软件,也欢迎他们能再作一点工作,很快能达到上网服务的目的。相信进入21世纪不久,在用户上人们就有可能逐步享受到多种媒体的全方位的科学信息的服务。

这套丛书是我国多位化学化工学科的专家、教授、学者多年辛勤劳动的成果,也是科学出版社、国家自然科学基金委员会优秀研究成果专著出版基金和中国科学院科学出版社基金大力支持的产物,希望它的出版能促进我国21世纪计算机化学化工学科的发展,并有助于相关学科发展。

《计算机化学化工丛书》编委会

2000年5月

前 言

分子的结构信息在计算机内部的表达、存储、转换技术,以及研究结构-性质关系的分子模型化方法是计算机化学的基本问题之一。本书以小分子结构信息处理为核心,围绕作者实验室的研究工作,介绍药物设计中的化学结构信息处理和分子模型化方法,是一本面向应用、介绍计算机辅助分子设计方法的入门书。

全书分为两部分。第一部分(第一至三章)是计算机化学中一些通用的方法、原理和技术:第一章介绍分子结构信息的计算机处理,包括结构信息在计算机中的表达、存储和显示,以及分子结构的几何优化、构象分析、构象搜索等重要概念;第二章介绍分子的相互作用,包括对生物过程至关重要的分子间的各种弱相互作用(非键相互作用),即范德华作用、静电作用、氢键作用和疏水作用;第三章介绍分子结构信息数据库,它们是从事分子模型化研究必不可少的工具。通过对中药化学成分及生物活性数据库的描述,引入了交叉学科数据库的概念。

第二部分(第四至九章)是分子模型化方法中的几个专题:第四章介绍遗传算法,包括了一种基于线杂交和面变异思想的新算法;第五章介绍定量结构活性关系(QSAR)研究,包括传统的2D-QSAR和近年的3D-QSAR。其中强调了结构信息自动数值化的思想,因为只有实现了结构信息的数值化,才有可能顺理成章地使用各种数学工具,而只有实现了自动数值化, QSAR方法才能实现大面积的普及应用;第六章介绍结构匹配和数据库搜索技术,其中二维的GMA算法和三维的3DFS软件都是有特色的成果,而基于三维药效团的大型数据库搜索法可用来发现新的先导化合物;第七章介绍的虚拟受体模型方法是一种很实用的技术,特别是对于受体结构未知的情形;第八章介绍药物与受体的对接技术。对接是基

于受体结构的药物设计中的核心技术，对药物设计感兴趣的读者不可以不知道对接的概念；第九章介绍近年来药物设计中的一个热门话题——模拟肽学。有理由认为，在 21 世纪模拟肽学将在药物设计中占据重要地位，特别是各种非肽的模拟物。

本书论及一个典型的交叉学科研究领域，涉及基础知识比较多，如量子化学基础、分子力学、有机立体化学、生物化学、药物化学等。这些内容不可能在本书中都做介绍，有需要的读者可阅读有关书籍。

本书第一、二、五、六、八章由王亨执笔，第三、七、九各章由周家驹执笔，第四章由何险峰执笔，全书由周家驹定稿。

书中涉及到的 LCC 实验室的一些研究工作大都是集体完成的，对此作出贡献的有：谢前、陈红明、孙红梅、石乐明、张跃敏、何敏、曹凌霄、李松梅、蔡萍、张晓晨、肖贞、谢桂荣、严新建、唐武成等。多年来，各级领导、科学界前辈和同仁对 LCC 的分子模型化研究给予了支持、指导和帮助，我们表示由衷的感谢。他们是：原化工部凌秋萍司长，郑秀兰、赵忠华、余昌申高工，科技部邵立勤副司长，中国科学院郁小民副局长，南开大学李正名先生，北京大学唐有祺先生、徐筱杰教授，北京大学药学院王夔先生、李仁利教授，中国农业大学陈万义先生，中国科学院生态研究中心白酒彬研究员，中国科学院化工冶金研究所郭慕荪先生和计算机化学开放实验室许志宏先生、杨章远、温浩研究员。

最后要感谢中国科学院科学出版基金的支持，以及科学出版社刘俊来编辑为此书的出版所付出的辛勤劳动。

限于作者的水平，书中难免有错误及不妥之处，恳切期望读者批评指正。

作者

2000 年 3 月于中国科学院化工冶金研究所

通讯处：周家驹，北京 353 信箱，100080，E-mail: jizhou@lcc.icm.ac.cn

目 录

《计算机化学化工丛书》序.....	i
前 言.....	v
第一章 分子结构的计算机处理.....	1
1.1 分子结构信息的表达.....	1
1.1.1 化合物的命名.....	1
1.1.2 碎片码、线形码、拓扑码.....	2
1.1.3 连接表.....	3
1.2 分子结构的输入.....	5
1.3 分子结构的几何优化.....	6
1.3.1 分子力学方法.....	6
1.3.2 量子力学方法.....	8
1.3.3 优化算法.....	10
1.4 构象分析.....	12
1.4.1 系统搜索方法.....	13
1.4.2 Monte Carlo 方法.....	14
1.4.3 模拟退火算法.....	14
1.4.4 遗传算法.....	15
1.4.5 分子动力学方法.....	16
参考文献.....	18
第二章 分子的相互作用.....	21
2.1 范德华相互作用.....	22
2.2 静电相互作用.....	23
2.2.1 计算原子点电荷的方法.....	24
2.2.2 分子的静电势.....	27

2.3	氢键相互作用	27
2.4	疏水相互作用	29
2.4.1	疏水效应和疏水作用	29
2.4.2	分配系数	31
2.4.3	疏水场	33
2.5	分子相互作用场的开发	36
	参考文献	36
第三章	分子结构信息数据库	39
3.1	结构数据库的三维信息来源	39
3.1.1	三维结构的实验来源	39
3.1.2	三维结构的计算来源	39
3.2	基于实验数据的结构数据库	40
3.2.1	蛋白质结构数据库 PDB	40
3.2.2	有机小分子结构数据库 CSD	41
3.2.3	无机晶体结构数据库 ICSD	43
3.3	基于计算结果的结构数据库	43
3.4	中药成分结构及生物活性数据库	45
3.4.1	需求分析	45
3.4.2	结构信息的计算机化	46
3.4.3	中药英文信息的规范化	48
3.4.4	研究进展	51
3.5	小结	52
	参考文献	52
第四章	遗传算法	54
4.1	遗传算法	54
4.1.1	遗传算法的表达式	54
4.1.2	遗传算法常用算子	56
4.1.3	遗传算法的其他参数及步骤	57
4.1.4	遗传算法的数学原理	59

4.2	经典遗传算法示例	61
4.2.1	初始试群体的产生	62
4.2.2	复制	62
4.2.3	杂交	63
4.2.4	变异	64
4.3	有指导性的遗传算法 DGA	66
4.3.1	初始群体的产生	67
4.3.2	线杂交	69
4.3.3	面变异	71
4.3.4	线杂交和面变异的联合作用	74
4.3.5	排序	75
4.3.6	群体规模	76
4.3.7	停止判据	76
4.3.8	DGA 的步骤	76
4.3.9	对 DGA 在优化中的测试	78
	参考文献	78
第五章	定量结构活性关系研究	80
5.1	用于 QSAR 的生物活性数据	80
5.2	用于 QSAR 的分子结构参数	83
5.2.1	电性参数	83
5.2.2	立体参数	84
5.2.3	疏水参数	85
5.3	经典 QSAR 模型	86
5.3.1	Hansch 模型	86
5.3.2	Free-Wilson 模型	86
5.3.3	其他模型	87
5.3.4	QSAR 模型的评价	87
5.4	CASAC 软件	89
5.4.1	设计思想	89

5.4.2	结构框架模型.....	90
5.4.3	结构信息数值化子系统 NumSF.....	92
5.4.4	数据分析及预测子系统 AnaQS.....	93
5.4.5	研究实例.....	94
5.5	用遗传算法挑选变量.....	99
5.6	经典 QSAR 的成功实例.....	104
5.6.1	治偏头痛新药 lomerizine.....	104
5.6.2	杀真菌农药 metconazole 和 ipconazole.....	105
5.6.3	治风湿性关节炎新药 flobufen.....	107
5.7	场分析方法 CoMFA	107
5.8	本征值方法 EVA	114
5.9	偏最小二乘回归.....	116
5.10	小结.....	119
	参考文献.....	120
第六章	结构匹配和数据库搜索技术.....	122
6.1	子结构匹配算法.....	122
6.1.1	回溯法.....	123
6.1.2	划分-松弛法.....	134
6.1.3	筛分法.....	134
6.1.4	其他方法.....	135
6.2	最大共同子结构查找算法.....	135
6.2.1	Brown 算法.....	136
6.2.2	EMCSS 算法.....	138
6.3	药效团识别.....	145
6.4	三维结构数据库搜索技术——3DFS 系统.....	148
6.4.1	提问结构定义.....	150
6.4.2	搜索算法.....	157
6.4.3	搜索例子.....	160
6.5	小结.....	163

参考文献.....	163
第七章 受体模型.....	165
7.1 受体模型概述.....	165
7.2 基于形状的受体模型.....	166
7.2.1 建模思想.....	166
7.2.2 建模步骤.....	166
7.2.3 应用实例.....	167
7.3 基于虚拟原子的受体模型.....	167
7.3.1 建模思想.....	167
7.3.2 方法原理.....	170
7.3.3 应用实例.....	175
7.4 基于虚拟点表面的受体模型.....	176
7.4.1 概述.....	176
7.4.2 虚拟受体表面的产生.....	177
7.4.3 表面性质的计算.....	180
7.4.4 配体分子与受体表面模型的相互作用.....	182
7.4.5 配体分子内能量的计算.....	185
7.4.6 应用实例.....	186
7.5 小结.....	193
参考文献.....	193
第八章 药物与受体的对接技术.....	195
8.1 基于受体结构的药物设计.....	195
8.2 已知配体对接.....	198
8.2.1 已知结合位点的实时图形对接.....	198
8.2.2 自动对接.....	199
8.3 全新配体对接.....	205
参考文献.....	208
第九章 模拟肽学.....	210
9.1 模拟肽学概述.....	210

9.1.1	生物活性肽作为药物的局限性.....	210
9.1.2	模拟肽学中的两类方法.....	211
9.2	以肽主链为基础的模拟肽学.....	211
9.2.1	环化技术.....	211
9.2.2	约束单元.....	215
9.2.3	环连接部件.....	216
9.3	约束氨基酸.....	219
9.3.1	α -甲基化氨基酸.....	219
9.3.2	α, α' -二烷基甘氨酸和 α -氨基环烷羧酸.....	222
9.3.3	N^α - C^α 环化氨基酸.....	224
9.3.4	N^α -甲基化氨基酸.....	226
9.3.5	β 和 γ -氨基环烷羧酸.....	227
9.3.6	α, β -不饱和氨基酸.....	227
9.3.7	β, β -二甲基和 β -甲基氨基酸.....	228
9.3.8	β -置换-2,3-亚甲基氨基酸.....	228
9.3.9	$N-C^\delta$ 和 $C^\alpha-C^\delta$ 环化芳香氨基酸.....	229
9.3.10	置换脯氨酸.....	229
9.3.11	D-氨基酸.....	230
9.4	肽的二级结构的分子模拟.....	230
9.4.1	β -转折.....	231
9.4.2	γ -转折.....	231
9.5	非肽配体的设计.....	231
9.6	小结.....	233
	参考文献.....	234

第一章 分子结构的计算机处理

1.1 分子结构信息的表达

分子在计算机中是如何表达的？这是关心分子模型化方法的人需要了解的首要问题。当我们描述一个分子时最习惯的做法是画出它的分子结构图，以顶点代表原子，边代表原子之间的连接键。然而这样的图形方式却不适合于计算机处理，了解计算机基础知识的人都知道，计算机最擅长接受和处理的是符号、数字以及由符号和数字组成的文件。因此，为了使计算机能方便地处理分子结构信息，需要采取计算机化的方式。

1.1.1 化合物的命名

化合物的命名可以在一定程度上反映分子的结构。早期通常使用简洁的习惯命名法，后来随着化合物种类的急剧增加，这种缺乏规律的命名法逐渐被遵循一套大家所公认的规则的命名法所代替。我们所普遍使用的命名法是 IUPAC 命名法，该命名规则由国际纯粹与应用化学联合会(IUPAC, International Union of Pure and Applied Chemistry)经过多次修订后于 1979 年正式颁布。例如习惯称为三甲胺(trimethylamine)的化合物的系统命名则是 *N,N*-dimethylmethanamine。但是该命名法仍然存在一个缺点，就是对许多化合物无法得到惟一的命名，因而限制了它在化学信息检索系统中的应用。针对这个问题，美国化学文摘社(CAS, Chemical Abstracts Service)提出了一套自己的命名规则，称为 CA 索引名(CA index name)。CA 索引名可以做到对每个化合物只给出一个命名，从而保证了化合物与命名一一对应，便于实现化合物命名的自动化和根据命名对化合物信息的检索。

1.1.2 碎片码、线形码、拓扑码

化合物的命名所表达的结构信息是非常笼统的，而且不便于计算机做进一步的结构分析和处理。在计算机处理化学结构的发展过程中，相继出现了碎片码、线形码、拓扑码等一些表达方法。其中碎片码和线形码都是将一些预先定义好的代表分子结构碎片的符号和数字线形排列起来，这两种表达方法的代表分别是 WLN 码^[1]和 SMILES 码^[2,3]，如图 1.1 中 2,3-二甲苯胺的 WLN 码和 SMILES 码分别是 ZR B1 C1 和 Cc1c(cccc1c)N。碎片码和线形码表达的结构信息有紧凑和简明的优点，但同时有一个致命的缺点：难于实现计算机处理化学结构的基本操作——子结构检索。因为用户所提出的提问子结构是不可预测的，当提问子结构与预先定义的结构碎片不相吻合时，就无法得到检索结果。也就是说提问子结构只能限制在预先定义的结构碎片中，而化学结构的千变万化，使得这一限制导致了碎片码和线形码不能被大多数计算机结构处理系统所采用(SMILES 在网络传输中还有应用)。

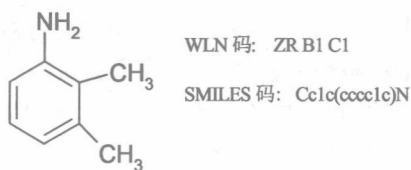


图 1.1 2,3-二甲苯胺的 WLN 码和 SMILES 码

拓扑码则利用了图论的知识，将化学结构图看作数学上的图，原子对应于图中的结点，化学键对应于图中的边，原子的不同类型(碳、氧、氮等)赋予结点不同的颜色，化学键的不同性质(单键、双键、三键、芳香键等)赋予边不同的颜色。最著名的拓扑码是美国化学文摘社采用的 Registry III 和法国学者 Dubois 创立的 DARC 系统。拓扑码虽然没有线形码所面临的子结构检索问题，但面临着一个同样重要的问题：惟一性问题。因为一个含 n 个原

子的结构可以有 $n!$ 种不同的编码方式。为了解决编码的惟一性需要采用一些算法，其中最有名的是 Morgan 算法^[4]。

1.1.3 连接表

无论何种命名法和编码方法都是力图将分子的复杂的结构信息以简略的并且是惟一的方式表达，从这一点看，有些类似于分子的拓扑指数，是对分子结构信息的高度浓缩。而我们在利用计算机处理分子时，希望面对的结构信息越原始越详尽越好，因为计算机不仅需要处理查询等问题，而且要处理结构检索问题，以及在分子设计中经常用到的分子力学、分子动力学、量子力学等计算。因此，以命名或编码的方式表达分子的结构信息是远远不够的。连接表是目前表达分子结构信息的最完善和最计算机化的方式，它除了可以准确地表达分子的二维信息，还能利用原子的空间坐标精确地表达分子的三维结构信息。

连接表本质上是分子中所有原子、键及其空间关系（有时包括立体化学）的一个列表，在计算机上就是一个文本文件，通常包括原子信息部分（原子类型、电荷、坐标等），键信息部分（键的类型，连接原子等）和其他特殊的结构信息部分。连接表不需要考虑惟一性的问题，原子的序号不影响分子的结构信息，也不影响现有的结构处理算法对分子的识别。可以说连接表是一种非常自由的表达方式，随着计算机存储量和计算速度的飞速发展，连接表已经成为绝大多数化学结构处理系统的首选方式。

连接表有各种各样的格式，几乎每个分子模型化软件的推出都伴随着一个新的连接表格式的出现。比较著名的有化学结构协会(Cheical Structure Association)设计的SMD格式^[5]，MDL信息系统公司(MDL Information System)推出的MOL格式^[6]和由Tripos公司推出的MOL2格式^[7]，其中MOL格式和MOL2格式由于相应公司推出的产品在世界各地比较普及，已成为普遍使用的文件格式。表1.1和表1.2分别是苯分子的MOL格式文件和MOL2格式文件。