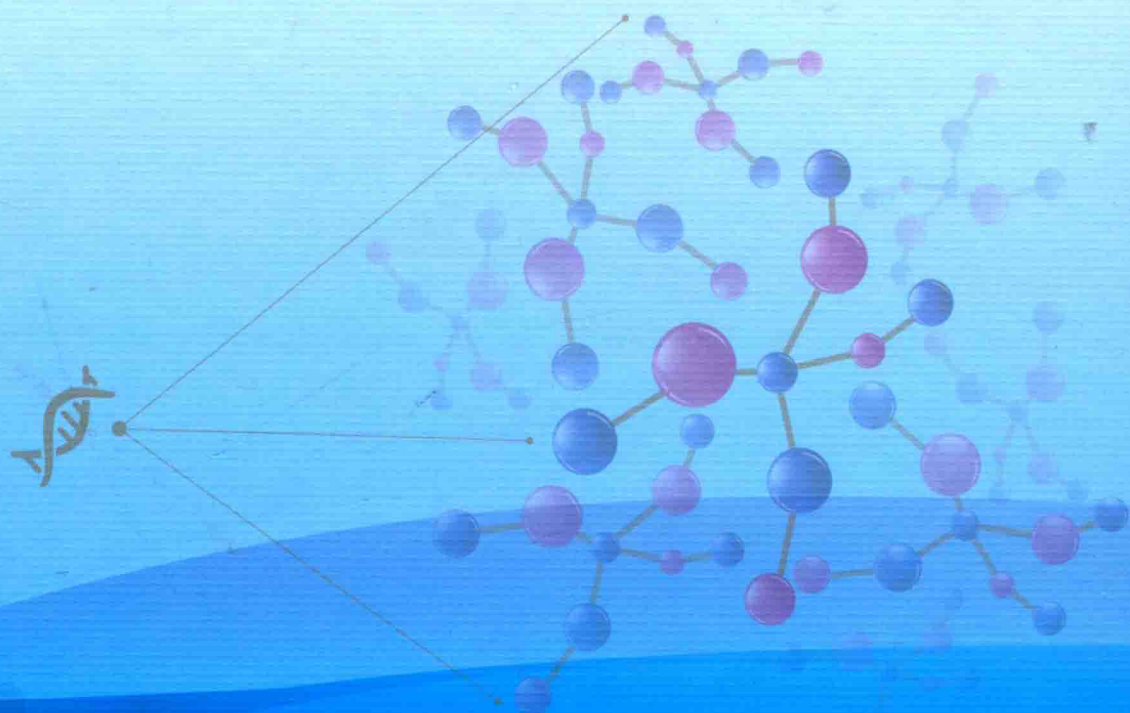


“十二五”国家重点图书出版规划项目

— 中医药信息学丛书 —

中医药信息学

崔 蒙 吴朝晖 乔延江 主编



科学出版社

“十二五”国家重点图书出版规划项目

中医药信息学丛书

中医药信息学

崔 蒙 吴朝晖 乔延江 主编

科学出版社

北京

内 容 简 介

本书主要介绍中医药信息学的理论、研究方法、应用领域及其研究进展。本书分上、下篇：上篇系统论述了中医药信息学作为一门新学科应具备的理论与方法学基础，包括其产生的背景、概念、研究内容、基本原理、与相关学科的关系，以及其基础标准、数据和知识服务的方法等。下篇介绍了中医药信息学的应用领域，包括中药信息学、中医临床信息学、中医药图书馆学、中医药情报学，重点论述了中医药信息学方法在中药信息、中医临床信息、中医药图书管及中医药情报研究中的实际应用进展与所取得的成果。

本书是第一部分系统论述中医药信息学这门新学科的专著，在内容和写作上体现“求实、创新”的宗旨，在阐述理论和信息学发展时，有着实际工作的基础。本书对中医药信息学的研究者具有启发性，也可作为中医药科研人员的参考书。

图书在版编目(CIP)数据

中医药信息学 / 崔蒙, 吴朝晖, 乔延江主编. —北京: 科学出版社, 2014. 12

(中医药信息学丛书)

“十二五”国家重点图书出版规划项目

ISBN 978-7-03-042775-5

I. 中… II. ①崔…②吴…③乔… III. 中国医药学—信息学 IV. R2-05

中国版本图书馆 CIP 数据核字 (2014) 第 294404 号

责任编辑: 刘 亚 曹丽英 / 责任校对: 朱光兰 刘亚琦

责任印制: 肖 兴 / 封面设计: 黄华斌 陈 敬

科学出版社 出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

中国科学院印刷厂 印刷

科学出版社发行 各地新华书店经销

*

2015 年 1 月第 一 版 开本: 787×1092 1/16

2015 年 1 月第一次印刷 印张: 35

字数: 815 000

定价: 128.00 元

(如有印装质量问题, 我社负责调换)

《中医药信息学丛书》编委会

主 编 崔 蒙 吴朝晖 乔延江

编 委 王映辉 李海燕 张华敏 赵英凯

李园白 王 耘 姜晓红

《中医药信息学》编委会

主 编 崔 蒙 吴朝晖 乔延江

副主编 王映辉 李海燕 张华敏 赵英凯 李园白

编写人员 (按姓氏汉语拼音排序)

白 岩	陈广坤	陈华均	董 燕	段 青	高 博
顾珮菝	郭玉峰	何 巍	胡艳敏	贾李蓉	姜晓红
姜又琳	焦宏官	亢 力	李凤玲	李鸿涛	李敬华
李 萌	李彦文	连超杰	廖利平	林兆洲	刘 辉
刘 静	刘堃靖	刘丽红	马兆辉	孟凡红	潘艳丽
史新元	宋观礼	田 野	佟 琳	童元元	王 静
王俊文	王 琳	王 星	王映辉	王 耘	徐丽丽
薛清录	杨坤杰	杨 阳	于 彤	张百霞	张 红
张 晶	张润顺	张伟娜	张燕玲	郑金生	周建伟
周霞继	周雪忠	朱 玲			

序

21 世纪是世界科学技术迅猛发展的时期，学科之间的交叉融合成为科技发展的重要趋势之一。其中，信息科学技术产生了广泛而深远的影响，对于医学领域也不例外。医学信息学是医学、计算机科学、人工智能、决策学、统计学和信息管理学的新兴交叉学科，在电子病历、医院信息系统、临床决策支持系统、远程医疗及数据交换标准等方面取得了丰硕的成果，已经在医院管理、教学和科研，疾病的预防、诊断和治疗等方面发挥了不可替代的作用。不言而喻，中医药信息学的发展历程更为年轻，富有潜力。中医中药流传数千年，至今仍然保持旺盛的生命力，在维护生命健康中发挥着独特而重要的作用。纵观中医药发展历程，总是与时代紧密相连，唯其如此，方能历久弥新。当今，现代科技背景下，中医药学术繁荣复兴，与现代医学乃至其他学科的汇聚、交流、融合、互补，逐渐成为中医药时代发展的显著态势。

中医药文献典籍浩如烟海，学术经验传承异彩纷呈，蕴藏着极为宝贵的学术资源，有待深入发掘。信息科学技术方法为此提供了崭新的机遇，对中医药学术的当代传承与发展发挥了重要的作用，中医药信息学这门新兴的学科也由此应运而生。同时，也应当看到，缘于学科性质、理论钩沉、社会文化背景、语言表述、思维模式、时代变迁等差异，中医药学术内容本身与信息科学技术的融合过程中必然存在重大挑战，中医药信息的获取、转化与共享等面临许多困难。这一点是医学信息学、地理信息学等其他与信息学交叉的学科发展过程中较少遇到的。所以尽管呈现出蓬勃的生机与巨大的潜力，但至今尚少有学者，也无专著对其内涵、外延进行详细论述。虽然已经成为国家中医药管理局重点建设学科，但其具体的学科建设仍是筚路蓝缕，充满艰辛，亟需奠基性著作充实其理论内核，支撑后备学术人才的教育培养。幸而，以崔蒙研究员等为首的学术团队，多年来致力于中医药信息学原理与方法学的研究、中医药信息数据库及中医药信息国际标准的研制，其进行了大量基础性的研究工作，积累了较丰富的经验和学识，很多工作与研究都充实了学科领域，为中医药信息学学科的设置、建设与发展提供了极其坚实的基础和有益的借鉴。

对于一门学科而言，理论与实践工作同等重要。相比中医药信息研究工作的大量开展，学科理论建设工作有所滞后，长期势必会影响与制约学科发展。由此，《中医药信息学》编撰工作的意义与价值显得极为关键。该书从全方位的角度介绍了这门学科的去、现况和未来，对中医药信息的内涵、外延、研究方法、内容及意义等着墨甚多，阐发明晰而深刻，对中医药信息学下中医药信息标准、中医药科学数据、中医药知识服务、中药信息学、中医临床信息学、中医药图书馆学和中医药情报学等七个分支学科均有系统论述。概言之，其研究内容几乎涵盖了一切与中医药活动有关的信息，如临床、科研、教育、管理、文化、生产经营等领域所产生的信息，提高了对中医药信息获取、转化、传播与利用的能力。

尤其值得一提的是，书中认为中医药信息是认识论层次的信息，具有现代整体性、动态时空性、现象理论等特征，其“主客融合的体验”及“包含本质的现象”等导致了辨证诊断和疗效的模糊，以及相对重视客体的整体变化状态，这些特点与大数据的“整体性”、“混杂性”、“相关性”三大特点不谋而合。如果能够借助大数据研究所获得的成果，从理论、方法学上解决体验信息获取、存储及传播的问题，必将对中医药学发展起到至关重要的推动作用。

目前，欧美发达国家对医学信息学的教育与训练非常重视，认为掌握必要的现代信息技术是医务工作者必须具备的一项基础知识和基本技能。这一点在中医药领域同样适用，但纵观国内临床医疗系统尤其是中医药领域，对此认识还尚待深化，这对拓展中医药工作者的视野、提升其临床水平及科研能力显然不利。我希望《中医药信息学》的问世能够在较大程度上引发学界对此问题的关注与重视，推动中医药信息学术的普及与发展，获得更大范围的学界共识。

相比传承千年、博大精深的中医药学，中医药信息学刚刚起步，尚有很多的工作需要一一完成，还有很多的困难需要一一克服，可谓前路漫长且艰、任重而道远。可喜的是，《中医药信息学》的编撰为万里征程开了一个好头，为这门学科的发展奠定了基础，指明了方向，确立了模式。“前人栽树，后人乘凉”，希望广大中医药信息工作者以此为起点，在全面而深刻把握中医药学术特质与发展规律的基础上，有效借鉴、运用信息科学原理、方法、技术，不断丰富中医药信息学的内涵，探寻其内在规律，为中医药学术的传承、发展乃至创新提供更多的助益，充分发挥其独特作用。

传统与现代的交融总是令人充满无限的遐想与期待，处于高概念和大数据时代的中医药信息学更加深化其学科特质，望能引领中医药学科、事业与产业的发展。对于崔蒙、吴朝晖、乔延江主编及编写团队，我比较熟悉他们的工作，感佩学者们孜孜不倦、辛勤耕耘、认真治学的精神，创建一个崭新的二级学科实在不易，此书乃中医药信息学的奠基之作。书濒脱稿邀我作序，是对我的信任和鼓励，谨志数语乐观厥成。

王永炎
甲午季秋

前 言

21 世纪是信息化的时代，信息技术正在深刻地改变着我们的学习、工作与生活。延用了千余年的中医学，虽然与信息科学起源、哲学基础各异，但均从整体、动态的角度观察、研究事物，两个学科交叉融合形成的中医药信息学则开启了中医药传统经验管理转向新型知识管理的崭新模式，标志着中医药学望、闻、问、切获取及利用信息的传统手段即将开始一次新的革命。

中医药信息学是一门生机勃勃的新兴学科，有关学科的内涵、外延、方法与技术等方面的诸多研究成果一直未得到过系统梳理，尚未形成完整的学科体系。

《中医药信息学》作为“中医药信息学丛书”的第一部，是提纲挈领的一部，分为上、下两篇，共八章。一至四章为上篇，系统介绍了中医药信息学的基本理论、方法与应用。五至八章为下篇，系统介绍了中医药信息学下各分支学科。

第一章为中医药信息学概论。主要介绍了中医药信息与大数据特征的相似性、大数据处理环境下中医药信息学的发展契机，概述了中医药信息学的内涵、外延、研究内容、方法论体系，中医药信息的主要特征、数据特点，中医药信息的形成、获取、转化、反馈、激活、传播的原理，以及在此过程中产生的中医药学独特的意象世界，及其与虚拟世界的沟通。

第二章为中医药信息标准。主要内容包括标准的定义和分级、各级标准的制定程序等，并系统介绍了目前国际、国内现有的中医药信息标准，包括中医药数据标准、中医药语义信息标准、术语系统标准等，以及与健康领域相关的信息标准的基本内容与研究现状。

第三章为中医药科学数据。系统介绍了中医药临床及科研活动所产生的中医药科学数据，以及科学数据的收集、分析及利用情况。内容包括中医药科学数据资源的研制、中医药科学数据常用数据挖掘方法及应用案例，如数据关联分析、数据聚类分析、数据分类分析、数据相似度分析等。

第四章为中医药知识服务。在系统调研、分析了中医药知识资源的基础上，提出了“知识即服务”的理念，对中医药知识服务平台进行了综合性论述，并分别介绍了基于“语义网”技术的中医药信息处理方法、中医药本体服务、中医药百科服务、中医药搜索服务、中医药知识发现服务、中医药决策支持服务，以及面向中医药知识服务的 3D 虚拟社区等具体的知识服务模式。

第五章为中药信息学。从学科内涵、外延、发展阶段、研究内容、基本理论体系等方面系统介绍了中药信息学的基本情况，并就中药药性与方剂配伍、中药有效成分族群辨识、中药生产过程质量控制三大关键科学问题举例说明中药信息学的研究思路与方法。

第六章为中医临床信息学。将中医临床诊疗过程归纳为信息获取、信息传递、信息处

理、信息再生、信息施效、信息组织的信息流程模型，将信息科学理论与方法融合应用于中医临床，介绍了中医临床信息学学科的基本概念、理论基础和发展历程，并从中医临床信息的采集、数据存储与管理、数据整合、分析利用等几个方面系统介绍了中医临床信息学的具体应用。

第七章为中医药图书馆学。系统介绍了中医药图书馆的产生、发展、类型、特性，界定了该学科概念的内涵、外延，重点介绍了中医药文献信息资源的建设与组织、中医药古籍资源积累方法、中医药图书馆核心竞争力等方面的研究成果。

第八章为中医药情报学。内容包括中医药情报学的基本概念、研究方法、研究范式、相关学科等。系统介绍了常用情报研究方法在中医药领域的应用情况，以及中医药战略情报、中医药竞争情报、中医药循证医学、中医药科技查新、中医药统计数据应用、中医药文献计量等中医药情报学各领域的研究进展。

本书由中国中医科学院中医药信息研究所、浙江大学、北京中医药大学、中国中医科学院广安门医院四家单位数十位科研人员通力合作、撰写而成，在成书过程中借鉴了大量国内外信息学相关著述与科研成果，同时也对数十年来中医药信息学学科研究成果进行了沉淀和梳理。中医药信息学是处在高速发展中的学科，其研究内容、研究方法等都在不断进化中，我们很难全面反映所有的研究方法和成果；此外，由于编写人员水平所限，本书还存在很多不足之处，书中舛错、遗漏也在所难免，这都是现阶段无法避免的问题，诚恳欢迎广大同道和读者批评指正，我们将不懈努力、不断完善。

我们正处在一个变革的时代，中医药信息学的发展将给中医药学的传承与创新带来新的思路和方法，希望本书的出版能够为学科再发展提供一个平台和框架，也希望读者能从本书提供的资源及成果中获得有益的启示。

编 者

2014年10月

目 录

序
前言

上 篇

第一章 中医药信息学概论	3
第一节 概述	3
第二节 中医药信息学发展的机遇	5
第三节 中医药信息学的基本概念与原理	13
第四节 发展前景	30
参考文献	31
第二章 中医药信息标准	33
第一节 概述	33
第二节 中医药信息标准的制定程序	35
第三节 中医药信息标准与健康相关领域信息标准	39
第四节 中医药数据标准	45
第五节 中医药语义信息标准与术语系统	61
参考文献	82
第三章 中医药科学数据	86
第一节 概述	86
第二节 中医药科学数据资源研制	88
第三节 中医药科学数据常用数据挖掘方法及案例	103
参考文献	120
第四章 中医药知识服务	124
第一节 概述	124
第二节 基于语义网的中医药信息处理方法	132
第三节 中医药本体服务	143
第四节 中医药百科服务	149
第五节 中医药搜索服务	153
第六节 中医药知识发现服务	160
第七节 中医药决策支持服务	172
第八节 面向中医药知识服务的3D虚拟社区	177
参考文献	182

下 篇

第五章 中药信息学	191
第一节 概述	191
第二节 中药药性与方剂配伍	192
第三节 中药有效成分族群辨识	210
第四节 中药生产过程质量控制	231
第五节 展望	264
参考文献	266
第六章 中医临床信息学	271
第一节 概述	271
第二节 中医临床信息的采集	280
第三节 中医临床信息数据的存储与管理	321
第四节 中医临床信息数据的整合	330
第五节 中医临床信息的分析利用	335
第六节 展望	362
参考文献	365
第七章 中医药图书馆学	370
第一节 概述	370
第二节 中医药文献信息资源建设	382
第三节 中医药文献信息的组织	395
第四节 中医古籍资源积累方法研究	412
第五节 中医药图书馆核心竞争力的研究	438
第六节 展望	449
参考文献	459
第八章 中医药情报学	462
第一节 概述	462
第二节 情报研究方法	467
第三节 中医药战略情报研究	476
第四节 中医药竞争情报	484
第五节 循证医学研究	492
第六节 科技查新	503
第七节 中医药统计数据应用	520
第八节 中医药文献计量学	527
第九节 展望	544
参考文献	545



上 篇



第一章 中医药信息学概论

21 世纪以生命科学、生态科学、信息科学、复杂科学和系统科学为前沿的科学技术迅猛发展, 自然科学与人文科学间的交叉、渗透、融合, 已成为科学技术整体化、综合化发展的重要趋势。这种交叉和融合既是实现科学知识系统整合的重要平台, 也是培育科技重大创新新枝——新兴学科的沃土。

信息科学就是在此环境下催生的一门新兴、综合性的前沿学科, 它是以信息作为主要研究对象, 以信息的运动规律作为主要研究内容, 以信息科学方法论作为主要研究方法, 以扩大人类的信息功能作为主要研究目标的一门科学^[1]。随着信息科学的迅速发展, 并逐步向各个领域的渗透、交叉、融合, 更多新兴交叉学科应运而生, 中医药信息学就是其中之一。

中医学是一门持续发展了数千年的古老传统医学, 它以人的生命活动作为主要研究对象, 以整体观为主导思想, 以中国古代哲学为哲学基础, 以个体化诊疗及辨证论治为临床特点, 是一门以自然科学为主体、多学科知识相交融的医学科学。在几千年的发展过程中, 中医学不仅吸收了古代哲学作为理论体系的主要架构, 还吸纳了古代天文学、气象学、地理学、物候学、农学、生物学、矿物学、植物学、军事学、数学及冶炼技术等诸多学科的知识^[2]。

中医学与信息科学虽然起源不同、哲学基础不同, 但均是建立在整体观基础上的学科。两者都不是从具体结构上对事物加以解剖性分析, 而是重视从整体上、动态中去观察和研究事物, 从而获得关于事物动态现象的运动规律和整体知识, 在信息学中表现为信息方法的功能准则和整体准则, 在中医学中则表现为脏象论和整体观。共同的基础、相似的方法学使得两个学科在交叉融合过程中逐渐形成了中医药信息学, 开启了中医药传统经验管理转向新型知识管理的一种崭新模式; 标志着中医学望、闻、问、切获取信息及利用信息的传统手段即将开始一次新的革命。

中医药信息学隶属于中医学, 是中医学随时代发展而产生的新的分支学科, 也是中医学学科群发展的必然结果。它的诞生关系到中医学的可持续发展, 同时也是标志中医药诊疗手段与经验传承进入飞速发展时代的重要里程碑。

第一节 概 述

古老传统的中医学有其独特的理论体系, 在形成初期, 这个体系相当活跃而包容, 以古代哲学和人体观察作为理论体系的主要架构, 还吸纳了古代自然、人文、社会诸多学科知识。兼收并蓄的发展观、切实有效的疗效, 令中医学保持了旺盛的生命力, 持续了

数千年。在漫长的发展历史中，中医药学理论体系逐渐成熟、稳定下来，同时也渐渐产生了封闭性和排他性，拖慢了中医药学发展的步伐。

近代超过 100 年的时间里，西方文化、西方科学和现代医学对中医药学的冲击，以及自然科学的快速发展，使得中医药学的发展速度相对缓慢。在现代医学刚刚传入中国时，中医药学依然是中国解决医疗卫生问题的主要力量，但随着时间的发展，现代医学逐渐占据了我国医疗卫生的主要市场，中医药学渐渐退居其次。造成这一现象的原因很多，包括西学东渐对中国文化的影响、西方科学对东方科学的冲击，但不可否认，其中一个重要的原因是未能与现代科学技术紧密结合，这导致了中医药学吸收新知识的速度放缓。

还原论是过去三个世纪以来西方科学思想的主要倾向，它主张把高级运动形式还原为低级运动形式，认为复杂的系统、事物、现象可以化解为各部分的组合，通过这种方法，对世界加以理解和描述^[3]。还原论派生出来的方法论手段就是对研究对象不断进行分解，恢复其最原始的状态，化复杂为简单，由此发展出的分析科学指导着现代科学不断取得新的进展。基于现代科学的现代医学同样向着更精、更小的方向不断挺进，这种方法最终是将生命的物质体现或是生命组成部分的加合等同于生命整体本身，将人体视为一架精密仪器，对其不断分割、细化研究，基因的发现是这一研究方向的最新突破，纳米技术的发展则给医疗手段提供了新思路和新方法。而根植于东方哲学，秉持整体观念的中医药学认为人体是一个有机整体，人体与自然、社会是不可分割的，这种整体性、系统性理论体系，使得中医药学无法与基于还原论产生的现代科学技术和成果紧密结合。

但近百年来不止西学东渐，东学也在西渐，东西方文化有了更多的交流，作为文化核心的哲学，东、西方哲学始终是并存的，并在不断的交流中获得进一步的发展^[1]。非线性科学和复杂性科学就是西方哲学的另一个发展方向^[4,5]，其注重从事物的复杂性及系统性去认识和把握事物的整体，关注事物在环境中的整体运动状态的变化，以及事物内部之间、事物与外部之间的复杂关系^[6,7]。这些学说的提出，为东西方文化、科学及医学的协同发展创造了新的条件，促成了中医药学与信息学的交汇融合。

20 世纪 80 年代以来，各地中医药院校与科研院所为了探讨信息学方法在中医药学中的应用，进行了大量工作，包括对中医药信息搜集、加工、处理方法和过程进行相关基础理论和技术方法的研究及探讨，建设了大量中医药学数据库、中医专家咨询系统，对已集成的中医药学数据进行了数据利用和挖掘，还进行了大量的四诊客观化研究，并取得了可喜的成果。但在这一探索过程中，也遇到了许多问题。

首先是怎样保证数据准确性和可用性。一直以来大部分科研工作者采取对数据进行清洗和结构化之后再行数据挖掘，以期发现新知识；在中医临床方面，也采用了随机对照试验（randomized controlled trial, RCT）方法来保证数据的可量化性和可重复性^[8,9]。在这一过程中，耗费了大量人力、物力和财力，而结果并不尽如人意，虽然取得了一些成效，但并不能体现中医药学的精髓。

随着科技的发展、数据的不断积累，当今社会迎来了大数据处理技术高速发展的时代，中医药学也迎来了一个新的发展契机，大数据的三大显著特点，即关注“整体”而非“抽样”、允许数据混杂性、关注事物间的相关关系^[10]，恰好与中医药学的传统观点相吻合，而大数据处理方法与技术的发展则进一步为中医药数据的处理提供了方法与技术，解

决了在小数据处理环境下解决不了的中医药数据处理中所存在的问题。

大数据处理技术的发展,为中医药信息的收集、存储、利用提供了新思路和发展方向,利用好这个契机,可以促进新兴的中医药信息学学科取得长足发展和进步,进而为提高中医药学的经验总结和临床应用提供参考。但目前来说,在中医药信息的收集、利用等方面还缺乏合适的方法与工具。

研究中医药信息固有的特点,以及其生成、获取、转化、激活、传播的原理,数字化后形成的知识密集型数据的处理方法,及其在虚拟世界中部分再现中医意象世界的真实,将对中医药学的发展起到巨大的推动作用。

(崔蒙 高博)

第二节 中医药信息学发展的机遇

在由物质、能量与信息三元素组成的世界中,中医药学主要的研究对象是与人体相关的信息,这是由于中医药学产生的年代使其没有能力深入研究与人体相关的物质和能量。与所有传统医学相似,中医药学是与其产生的文化背景、哲学基础密切相关的。由于受产生时历史环境的影响,中医药学是建立在整体观基础上的,它是研究人体生命运动的科学,因而对产生于还原论基础上的科研成果很难吸收。特别是当它的研究对象是处于开放环境下的复杂巨系统,也就是处于自然和社会环境中的人体整体时,甚至基于信息科学建立起来的小数据处理方法所产生的成果也很难有效地运用于本学科中。直到大数据环境产生后,随着对大数据三大特点——“整体性、混杂性、相关关系性”认识的不断深入,形成了与之相适应的数据处理方法,才为中医药学的发展创造了前所未有的良好机遇,而在这一历史发展机遇中,中医药信息学扮演了重要的角色。

一、大数据

大数据并非一个确切的观念。最初,这个概念是指需要处理的信息量过大,已经超出了一般电脑在处理数据时所能使用的内存量,因此工程师们必须改进处理数据的工具。

2008年,《自然》杂志推出“Big Data”专辑,从互联网技术、互联网经济学、环境科学、生物医药等多个方面介绍了大数据应用所带来的技术挑战及可以预见的未来发展方向。2009年微软公司发布了《e-Science:科学研究的第四种范式》论文集,全面地描述了大数据时代快速兴起的数据密集型科学研究^[11]。2012年微软亚洲研究院发布了《第四范式:数据密集型科学发现》^[12]论文集中文版,进一步从科学研究模式角度来分析大数据及其深远影响。2011年,《科学》杂志刊登专题——“数据处理”(Dealing with Data),主要针对多个学科相关科研数据的膨胀问题展开讨论,以更好地应对大数据带来的数据组织与访问挑战。2011年,企业界和学术界联合讨论,共同面对“大数据”的机遇和挑战,内容包括大数据的概念、组成、处理的关键技术、服务模式、管理方式等。2012年3月

29日,美国奥巴马政府发布了“Big Data Big Deal”,宣布投资启动“大数据研究和发展计划”,并将其定义为“未来的新石油”,希望增强政府收集、分析和萃取海量数据的能力,通过提高从大型复杂的数字数据集中提取知识和观点的能力,加快科学与工程中的步伐,加强国家安全,改变科学研究。首批共有6个联邦政府部门和机构,2亿美元的投入用于提高大量数据的访问、组织、收集、发现信息的技术水平^[13]。

一些大型公司已经开始赞助大数据相关项目的竞赛,并且为高等院校的大数据研究提供资金。有些大学也已经开始设置大数据相关的新课程,这些课程将培养出新一代的“数据科学家”,扩大大数据技术开发和应用所需的人才供给。一些组织提出,将建立大数据论坛,对公益性的数据进行采集、分析和可视化等。这个项目的推动,标志着“大数据”处理已经在全球获得了高度重视。

二、大数据的特点

(一) 大数据特点之一:关注“整体”而非“抽样”

在小数据处理环境下,因为无法获取全部数据,人们只能采用随机采样的方法,用部分数据推测整体数据,以获得对事物的正确认识,因而,对所采集数据的准确性要求非常严格,样本数据能否代表整体数据成为了数据分析的关键环节。而在大数据处理环境下,由于获取数据的能力大幅度提高,人们不再依赖于随机采样,取而代之的是将所有数据全部纳入分析,保证数据的整体性。这样能更快、更容易地发现问题,从而能够更多地关注到小数据研究所不能发现的细节。

2009年甲型H1N1流感暴发,在暴发前几周,谷歌公司就通过网络检索记录预测了流感的传播^[4]。这一预测并不是依赖于对数据的随机抽样,以及对样本数据的精确分析,而是收集了整个美国几十亿条互联网检索记录,关注的重点是特定检索词条的使用频率与流感在时间和空间传播之间的关联关系。为了处理这些未经选择的全部数据,谷歌公司共建立了4.5亿个不同的数学模型,与2007年、2008年美国疾控中心记录的实际流感病例进行对比,发现了45条检索词的组合,用于特定的数学模型后,他们的预测与官方数据的相关性达到97%,并且非常及时。而疾控中心通常要在流感暴发一两周之后才能做到。谷歌公司对流感的预测,是大数据理论在实践中的一次成功应用,同时也说明了较之随机采样研究法,采用所有数据的研究方法更能够发现事物的发展状态。

在过去,随机采样曾取得了巨大的成功,成为现代社会、现代测量领域的主体研究方法,但对于揭示事物的本来面貌来说,这只是一条捷径。随机采样法是在不可能收集和分析全部数据的情况下,用少量具有整体数据特征的数据尽可能全面地代表全体数据并进行分析,将抽样数据的分析结果视为全部数据的处理所得出的结果。这是在数据收集与处理无法达到大数据水平时的一个折中方法,本身存在许多不可知的风险,它假想的是一种理想状态,即所抽样采集的数据能够代表全部数据。随机采样研究法的成功必须依赖于采样的绝对随机性,且采集的数据确实能够反映整体数据,但这实施起来非常困难,一旦采样过程中存在任何偏见,分析结果就会相去甚远。因此,随机采样法的分析结果能够给出一