



CCF 大数据教材系列丛书
CCF 大数据专家委员会 组编

主编 程学旗

大 数 据

分 析

高等教育出版社

数据教材系列丛书
数据专家委员会 组编

程学旗

大	数	据
分	析	

大	数	据		
Doshuju Fenxi		分	析	

图书在版编目(CIP)数据

大数据分析 / 程学旗主编. -- 北京: 高等教育出版社, 2019.4
ISBN 978-7-04-051632-6

I. ①大… II. ①程… III. ①数据处理-高等学校-教材 IV. ①TP274

中国版本图书馆CIP数据核字(2019)第050472号

郑重声明
高等教育出版社依法对本书享有专有出版权。任何未经许可的复制、销售行为均违反《中华人民共和国著作权法》，其行为人将承担相应的民事责任和行政责任；构成犯罪的，将被依法追究刑事责任。

为了维护市场秩序，保护读者的合法权益，避免读者误用盗版书造成不良后果，我社将配合行政执法部门和司法机关对违法犯罪的单位和个人进行严厉打击。社会各界人士如发现上述侵权行为，希望及时举报，本社将奖励举报有功人员。

反盗版举报电话
(010) 58581999 58582371

58582488
反盗版举报传真

(010) 82086060

反盗版举报邮箱
dd@hep.com.cn

通信地址
北京市西城区德外大街4号
高等教育出版社法律事务
与版权管理部
邮政编码 100120

防伪查询说明

用户购书后刮开封底防伪涂层，利用手机微信等软件扫描二维码，会跳转至防伪查询网页，获得所购图书详细信息。也可将防伪二维码下的20位密码按从左到右、从上到下的顺序发送短信至106695881280，免费查询所购图书真伪。

反盗版短信举报
编辑短信“JB，图书名称，出版社，购买地点”发送至10669588128
防伪客服电话
(010) 58582300

策划编辑 张江漫
责任编辑 张江漫
书籍设计 张申申
插图绘制 于博
责任校对 刁丽丽
责任印制 赵义民

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
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.hepmall.com.cn>
<http://www.hepmall.com>
<http://www.hepmall.cn>
印刷 北京中科印刷有限公司
开本 787mm×1092mm 1/16
印张 22
字数 370千字
版次 2019年4月第1版
印次 2019年4月第1次印刷
定价 48.00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 51632-00

内容简介

本书系统地介绍了大数据分析的理论、算法及应用,全书共10章,包括基本概念和基本知识、大数据统计分析方法、大数据机器学习方法、机器学习理论、大数据算法理论、文本大数据分析、知识计算、网络数据挖掘、社交媒体分析、大数据分析系统架构等内容。除章节内容外,每章还提供“小节及进一步阅读”“参考文献”“习题”等内容帮助读者学习。

本书既可作为普通高等学校大数据相关专业的教材使用,也可供大数据分析领域的专业技术人员参考。

主任 梅 宏

成员 (按姓名拼音排序)

卜佳俊 陈宝权 陈恩红

程学旗 杜小勇 方 粮

胡 斌 黄宜华 金 海

马华东 潘柱廷 王建民

王晓阳 王元卓 袁晓如

周傲英 周 涛 周晓方

随着大数据的蓬勃发展，大数据领域人才的需求越来越大，大数据人才培养受到了各界的广泛关注。2016年，教育部开始批准设立“数据科学与大数据技术”本科专业，越来越多的高校申请开设“数据科学与大数据技术”专业或开设大数据方向的相关课程，截至2018年3月，已有近三百所高校获批建设“数据科学与大数据技术”专业。虽然大数据专业和大数据方向的课程不断开设，但是，当前我国高校的大数据教学尚处在摸索阶段，尤其缺乏成熟的、系统性和规范性的大数据教学体系和教材。

在此背景下，中国计算机学会大数据专家委员会成立了大数据教材系列丛书编委会，着手编著系列化、规范化的大数据教材。自2017年6月，经编委会多次研讨，形成丛书框架，作者们随即开始紧张的编写工作。编委会和作者间也有多轮的初稿审阅和研讨交流。数易其稿，终于付梓。

大数据教材系列丛书采用“1+3+X”的体系，即以1本《大数据导论》为基础，设置《大数据管理》《大数据处理》《大数据分析》3本关键技术教材，以及针对行业领域的X本应用教材。本套教材系列丛书既适合高校大数据专业的专科生、本科生以及研究生系统地学习大数据相关知识与技术，也适合从事大数据相关技术的企事业单位研究人员、工程师作为参考用书。

《大数据导论》是一本全面介绍大数据相关知识的专业通识教材，其系统地介绍大数据涵盖的内容，包括数据与大数据、大数据获取与感知、大数据存储与管理、大数据分析、大数据处理、大数据治理、大数据安全与隐私等，同时还介绍了部分行业中大数据的典型应用案例，反映了大数据在社会经济生活中的重要价值。

《大数据管理》首先综述数据管理系统的发展，指出发展大数据管理系统是历史的必然，并沿着数据模型和系统构件两个维度上展开。在数据模型的维度上，主要介绍关系、键值对、图和文档数据模型及其语言；在系统维度上，介绍系统结构、存储与组织、查询处理、事务管理、故障恢复等话题。

《大数据处理》包括大数据处理基础技术、大数据处理编程与典型应用处理、大数据处理系统与优化三个方面。本教材以大数据处理编程为核心，从基础、编程到优化等多个方面对大数据处理技术进行系统介绍，

使得读者能够快速入门，同时体会大数据处理系统的设计理念与优化方法本质。

《大数据分析》包括大数据分析方法和理论、典型大数据分析任务以及大数据分析系统与应用。本教材特色是理论联系实际，本书从基础理论、典型任务以及系统应用多个方面对大数据分析相关知识进行了系统而详细的介绍，使得读者能够快速入门，体会大数据分析技术的本质特征，领略大数据技术带来的创新理念。

“X”系列教材包含面向各行各业的大数据应用的知识与技术，既可面向工程实践，又可面向职业培训，且将随着产业界大数据应用的发展进行更新迭代。

大数据已成为学术界、产业界和政府共同关注的热点，正在开启信息化的新阶段。大数据人才培养也刚刚起步，还需要付出更多努力去探索。通过汇集中国计算机学会大数据专家委员会的智力资源，丛书编委会希望本系列教材能够为我国大数据人才培养尽到绵薄之力，助力我国大数据事业的蓬勃发展。尽管编委会和作者花费了很大精力规划和编写本系列教材，但是囿于对大数据的认识局限和自身能力限制，难免存在疏漏和错误，欢迎读者批评指正，以待再版时修正完善。

大数据教材系列丛书编委会

2018年7月

大数据开启了信息化发展的第三阶段，大数据分析技术是推动信息化发展新阶段社会进步与经济发展的关键技术。大数据分析涉及统计、机器学习、计算机算法、计算架构、引擎系统等多个方面的技术，是一个综合性技术栈。同时，应用的多样性带来的数据来源、数据类型、数据形态等方面的多样化，使得大数据分析和具体应用紧密关联。撰写大数据分析方面的书籍需要综合考虑大数据分析的理论、算法及应用等多个层面的进展，本书作者从向量化、序列化、网络化三类典型数据出发，结合大数据分析的基础理论和方法，撰写本书。

《大数据分析》一书是大数据系列教材的成员之一，致力于从大数据分析的理论、算法及应用方面为大数据相关专业的本科生、研究生及科研人员提供一本较为全面介绍大数据分析相关技术的专业教材。本书系统地介绍大数据分析方法，包括大数据统计分析方法、大数据机器学习方法、机器学习理论、大数据算法理论、文本大数据分析、知识计算、网络数据挖掘、社交媒体分析、大数据分析系统架构等。除了介绍大数据分析方法之外，本书还介绍了部分行业中大数据的应用案例。本书既可作为普通高等学校大数据相关专业的教材使用，也可供专业技术人员参考。

本书共分10章。第1章简要介绍大数据分析的一些基本概念和基本知识，为读者提供一个使用本书的背景知识，并对本书的架构进行简明扼要的介绍。第2章介绍大数据统计分析方法，包括相关性分析、因果推断和大数据采样分析方法。第3章介绍大数据机器学习方法，由浅入深分别介绍了描述性分析方法、预测性分析方法、深度学习方法和强化学习方法。第4章介绍机器学习理论，从机器学习基础、模型选择、偏差方差分解、PAC学习理论和非独立同分布学习等几个方面展开。第5章介绍大数据算法理论，包括组合优化算法、在线算法、流式算法、参数算法等。第6章着眼于文本大数据分析，包括文本表达、文本匹配、文本生成。第7章围绕知识图谱介绍知识计算，包括知识抽取、知识融合、知识推理等。第8章着眼于网络数据，从网络排序、网络聚类、网络表示学习三个方面介绍了网络数据挖掘的方法和研究进展，包括中心度度量、网络划分、社区发现、网络嵌入等。第9章以社交媒体计算为例介绍了大数据分析的应用，包括社交网络分析、大图异常检测和其他一些社交媒体分析的新应用。第10章介绍大数据分析

系统架构，包括数据与计算的演变历程、大数据分布式计算模型、大数据计算系统等。

参与本书撰写的主要人员包括：程学旗、徐君、沈华伟、郭嘉丰、靳小龙、兰艳艳、刘盛华、张家琳、张志斌、庞亮。程学旗牵头确定了全书的知识架构，第1章“基本概念与基本知识”由程学旗起草；第2章“大数据统计分析方法”由徐君起草；第3章“大数据机器学习方法”由郭嘉丰起草；第4章“机器学习理论”由兰艳艳起草；第5章“大数据算法理论”由张家琳起草；第6章“文本大数据分析”由庞亮起草；第7章“知识计算”由靳小龙起草；第8章“网络数据挖掘”由沈华伟起草；第9章“社会媒体分析”由刘盛华、沈华伟起草。第10章“大数据分析系统架构”由张志斌起草。徐君和沈华伟协助我完成了全书统稿。本书由清华大学王建民教授审稿，在此致以衷心的感谢，同时感谢在本书撰写过程中提供素材或建议的老师和同学。

感谢中国计算机学会大数据专家委员会给予的支持和指导，感谢大数据教材系列丛书编委会各位同仁的努力付出。

作者认识和水平所限，参与本书编写人员众多，时间紧，统稿难度大，书中必然存在不少不一致的认知和叙述，欢迎读者批评指正，并积极反馈意见和建议，联系邮箱：bigdata@ccf.org.cn，我们将在再版时吸纳，使本书逐步趋于完善。

程学旗

2019年1月

■ 第1章 基本概念与基本知识001	3.2 预测性分析055
1.1 大数据与大数据分析001	3.2.1 分类分析方法055
1.2 重要的问题和概念004	3.2.2 排序学习064
1.3 大数据分析算法、系统和应用005	3.3 深度学习分析方法072
1.4 大数据分析科学家和工程师007	3.4 强化学习分析方法079
1.5 本书的结构008	3.4.1 代表性方法082
	3.4.2 大数据分析中的 强化学习085
■ 第2章 大数据统计分析方法011	3.5 小结及进一步阅读088
2.1 相关性分析011	习题089
2.1.1 相关性理论的产生011	
2.1.2 相关关系012	
2.1.3 传统的统计相关性 分析方法013	■ 第4章 机器学习理论091
2.1.4 大数据中的统计 相关性分析016	4.1 机器学习基础091
2.2 因果推断020	4.1.1 基本概念092
2.2.1 因果推断简介020	4.1.2 损失函数与风险函数093
2.2.2 相关关系与因果关系020	4.1.3 经验风险最小化与结构风险 最小化095
2.2.3 无模型因果推断021	4.2 过拟合、模型选择以及正则化096
2.2.4 基于模型的因果推断022	4.2.1 训练误差与测试误差096
2.2.5 大数据中的因果推断024	4.2.2 过拟合与模型选择096
2.2.6 Yule-Simpson 悖论025	4.2.3 正则化与交叉验证098
2.3 采样分析026	4.3 偏差方差分解102
2.3.1 采样与随机模拟027	4.4 PAC 学习理论106
2.3.2 蒙特卡罗方法028	4.4.1 一个简单的例子106
2.3.3 马尔可夫链蒙特卡罗 方法031	4.4.2 PAC 学习理论 基本概念107
2.3.4 并行采样方法034	4.4.3 有限假设空间下的推导109
2.4 小结及进一步阅读035	4.4.4 VC 维111
习题037	4.4.5 Rademacher 复杂度113
	4.5 非独立同分布学习115
■ 第3章 大数据机器学习方法039	4.5.1 非独立情形115
3.1 描述性分析039	4.5.2 非同分布情形116
3.1.1 聚类分析039	4.6 小结及进一步阅读118
3.1.2 矩阵分解046	习题118

■ 第 5 章 大数据算法理论119	■ 第 7 章 知识计算209
5.1 组合优化算法119	7.1 知识图谱简介209
5.1.1 近似算法119	7.2 知识抽取210
5.1.2 次模优化121	7.2.1 实体抽取210
5.2 在线算法127	7.2.2 关系抽取215
5.2.1 秘书问题127	7.2.3 属性抽取223
5.2.2 在线调度129	7.2.4 实体关系联合抽取228
5.2.3 在线二部图匹配132	7.3 知识融合232
5.2.4 在线学习中的多臂老虎机 问题137	7.3.1 实体对齐232
5.3 流式算法141	7.3.2 实体链接235
5.3.1 流模型和流算法简介141	7.3.3 知识更新238
5.3.2 图上的流模型142	7.4 知识推理239
5.3.3 统计类问题的流模型144	7.4.1 基于逻辑的推理模型239
5.3.4 聚类问题的流模型148	7.4.2 基于图的推理模型241
5.4 参数算法150	7.4.3 基于表示学习的 推理模型245
5.4.1 参数算法设计基本技巧152	7.5 小结及进一步阅读251
5.4.2 参数算法下界156	习题254
5.5 小结及进一步阅读157	
习题158	
	■ 第 8 章 网络数据挖掘255
■ 第 6 章 文本大数据分析159	8.1 网络排序255
6.1 文本表达159	8.1.1 节点中心度255
6.1.1 单词的表示159	8.1.2 边中心度260
6.1.2 句子的表示172	8.2 网络聚类264
6.2 文本匹配181	8.2.1 网络划分264
6.2.1 文本匹配任务182	8.2.2 社区发现268
6.2.2 基于规则的文本匹配184	8.3 网络表示学习272
6.2.3 基于学习的文本匹配187	8.3.1 Laplacian eigenmaps273
6.3 文本生成192	8.3.2 DeepWalk275
6.3.1 文本生成简介192	8.3.3 LINE276
6.3.2 人机对话生成196	8.3.4 SDNE277
6.3.3 图片标题生成201	8.4 小结及进一步阅读279
6.3.4 文本生成的评价203	习题280
6.4 小结及进一步阅读205	
习题207	

■ 第9章 社交媒体分析281	■ 第10章 大数据分析系统架构311
9.1 网络影响力最大化281	10.1 数据与计算的演变历程311
9.2 基于位置的社交网络285	10.1.1 数据规模的演变311
9.3 大图的异常检测288	10.1.2 计算范式的演变313
9.3.1 基于密度子图的 检测方法290	10.2 大数据分布式计算模型317
9.3.2 基于谱图子空间的 检测方法292	10.2.1 大数据分析算法的 挑战317
9.3.3 信念传播294	10.2.2 数据与参数分发策略319
9.3.4 视觉引导的自动检测295	10.2.3 数据更新策略321
9.3.5 基于信号处理的 检测方法296	10.3 大数据计算系统323
9.4 社交媒体分析新应用297	10.3.1 MapReduce 系统324
9.4.1 社交媒体中的广告投放297	10.3.2 Spark 系统326
9.4.2 移动互联网环境的 推荐系统300	10.3.3 参数服务器 Parameter Server330
9.4.3 社交网络中的风险控制303	10.3.4 TensorFlow 系统334
9.5 小结及进一步阅读308	10.4 小结及进一步阅读337
习题308	习题337

1.1 大数据与大数据分析

近几年,大数据迅速发展成为科技界和企业界乃至世界各国政府关注的热点。《自然》和《科学》等刊物相继出版专刊探讨大数据带来的机遇和挑战。著名管理咨询公司麦肯锡称:“数据已经渗透到当今每一个行业和业务职能领域,成为重要的生产因素。人们对于大数据的挖掘和运用,预示着新一波生产力增长和消费盈余浪潮的到来。”美国政府认为大数据是“未来的新石油”,一个国家拥有数据的规模和运用数据的能力将成为综合国力的重要组成部分,对数据的占有和控制将成为国家间和企业间新的争夺焦点。大数据逐渐成为社会各界关注的新焦点,“大数据时代”已然来临。

什么是大数据,迄今没有公认的定义。从宏观世界角度来讲,大数据是衔接物理世界、信息空间和人类社会三元世界的纽带。物理世界通过互联网、物联网等信息技术有了在信息空间中的大数据投影,而人类社会则借助人机界面、脑机界面、移动互联等手段在信息空间中产生自己的大数据映像。从信息产业角度来讲,大数据是新一代信息技术产业的强劲推动力。新一代信息技术产业本质上是构建在大数据、云计算、移动互联网等第三代信息技术平台上的信息产业。IDC(International Data Coporation)预测,到2020年第三代信息技术平台的市场规模将达到5.3万亿美元,而从2013—2020年,IT产业90%的增长将由第三代信息技术平台驱动。从社会经济角度来讲,大数据是第二经济(second economy)的核心内涵和关键支撑。第二经济的概念是由美国经济学家Auther在2011年提出的。他指出,由处理器、链接器、传感器、执行器以及运行在其上的经济活动形成了人们熟知的物理经济(第一经济)之外的第二经济(不是虚拟经济)。第二经济的本质是为第一经济附着一个

“神经层”，使人们的经济活动智能化，这是100年前电气化以来最大的变化。Auther还估算了第二经济的规模，他认为到2030年，第二经济的规模将逼近第一经济。而第二经济的主要支撑是大数据，因为大数据是永不枯竭并不断丰富的资源产业。借助于大数据，未来第二经济下的竞争将不再是劳动生产率而是知识生产率的竞争。

从非正式的说法来看，大数据是指无法在一定时间范围内用常规软件工具进行捕捉、管理和处理的数据集合，是需要更新处理模式才能具有更强的决策力、洞察力和流程优化能力的海量、高增长和多样化的信息资产。国际商用机器公司IBM提出用5V特征来描述大数据，得到了业界的广泛认可，这5V特征分别是：

(1) 体量 (volume) 大：2003年，人类第一次破译人体基因密码时，用了10年才完成了30亿对碱基对的排序；而在10年之后，世界范围内的基因仪15min就可以完成同样的工作量。

(2) 种类 (variety) 多：随着传感器、智能终端以及在线社交协作技术的飞速发展，组织中的数据也变得更加复杂，因为它不仅包含传统的关系型数据，还包含来自网页、互联网日志（包括点击流数据）、搜索索引、社交媒体、电子邮件、文档文件、音视频文件各种传感器数据等原始、半结构化和非结构化数据。

(3) 速度 (velocity) 快：所谓“1秒定律”指的是如果不能在秒级时间范围内给出分析结果，数据就会失去价值。而对于大数据来说，这里的快不仅仅指的是实时进行数据处理和结果展示要求数据处理速度快，也包含数据产生速度快的特点。有的数据是爆发式产生，例如，欧洲核子研究中心的大型强子对撞机在工作状态下每秒产生PB级的数据；有的数据是涓涓细流式产生，但是由于用户众多，短时间内产生的数据量依然非常庞大，例如，互联网点击流、电商平台日志、射频识别数据、GPS（全球定位系统）位置信息。

(4) 价值 (value) 大、密度低：数据价值密度相对较低，却如浪里淘沙又弥足珍贵。随着互联网以及物联网的广泛应用，信息感知无处不在，信息海量，但价值密度较低，如何结合业务逻辑并通过强大的机器学习来挖掘数据价值，是大数据时代最需要解决的问题。

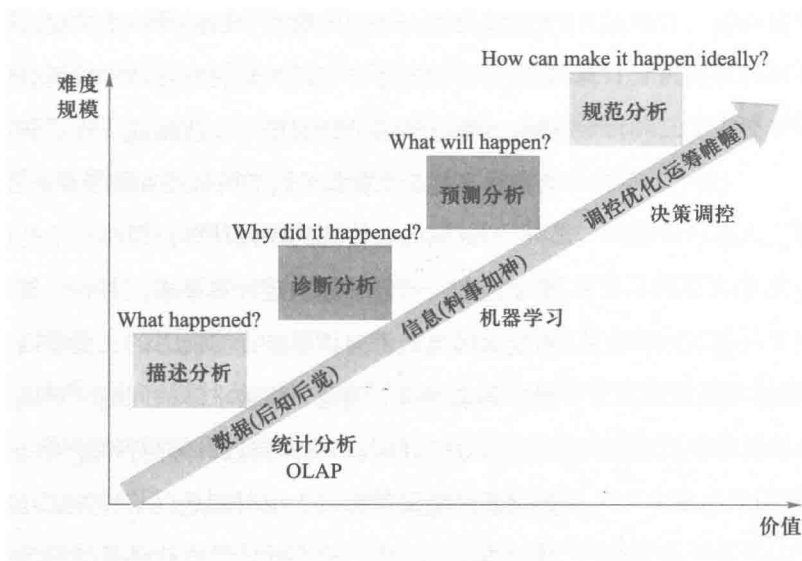
(5) 真实 (veracity) 性难确定：量大导致数据的准确性和可信度难以判定，数据质量良莠不齐。

大数据的5V特征使得大数据分析的主要难点并不仅仅在于数据体量大。实际上，通过对计算机系统的扩展可以在一定程度上缓解数据体量大带来的挑战。而大数据分析真正的挑战来自于数据类型多样、数据质量的不确定性和数据分析的实时性要求。数据类型多样使得一个应用往往既要处理结构化数据，又要处理文本、视频、语音等非结构化数据，这对单一的数据库系统来说是难以应付的；而数据质量的不确定性，使得数据真伪难辨成为大数据应用的一大挑战，追求高质量数据是对大数据处理的一项重要要求，最好的数据清洗方法也难以消除某些数据固有的不可预测性；在数据的实时性分析方面，许多应用中时间就是利益，大数据分析的实时响应、在线更新是关键。

为了应对大数据带来的上述困难和挑战，以Google、Facebook、阿里、腾讯、Amazon等为代表的互联网企业近几年推出了各种不同类型的大数据处理系统。借助于新型的处理系统，深度学习、知识计算、可视化等大数据分析技术也得到迅速发展，并逐渐被推广应用于不同的行业和领域。

从价值角度看，大数据可以提供社会科学的新方法和科学研究的新范式，形成高新科技的新领域，推动社会进步的新引擎，而大数据分析技术是实现从数据到价值的关键。按照分析任务的难度以及其产生的价值划分，大数据分析技术可以分为四个层次：描述分析（descriptive analysis）、诊断分析（diagnostic analysis）、预测分析（predictive analysis）、规范分析（prescriptive analysis），如图1-1-1所示。其依次

图1-1-1 大数据分析的四个层次



进行描述事情发生结果、分析事情发生原因、预测事情未来趋势以及控制事情发生轨迹。

1.2 重要的问题和概念

融合“人、机、物”三元空间的大数据存在规模大、关联复杂、状态演变等显著特征。部分传统的数据分析方法已经不再适用。在原有机器学习的基础上，设计出适用于大数据分析的方法成为关键问题之一。大数据分析算法模型具有训练样本多、输入数据维度高、输入数据结构复杂、模型参数多的特点。研究者们结合机器学习、认知计算和知识工程，研制新型的大数据分析系统算法，将大大提高大数据的价值利用率。

对于大数据基础分析算法，可以利用数据驱动表达学习与强化学习，提升分类、聚类、查询、检索、匹配、关联分析等基础算法的精准性和适用性。对于大数据融合分析技术，可以利用认知计算和深度学习，提升异构表达、跨媒体特征抽取与内容理解、数据融合分析、异常模式识别等高级算法的效能。大数据预测决策与可视化，主要利用机器学习与知识建模，突破大数据预测、知识推演、可视化分析、辅助决策等大数据分析技术。

随着训练数据集规模以及模型参数规模指数级的增长，如何获取强大的计算能力成为大数据分析的另一个关键问题。研究者利用智能芯片和硬件重构技术，可以大幅提升复杂大数据分析计算的速度。预计将来研究弹性泛流式分析引擎，支持批式、流式和在线分析，通过软硬件集合来提高复杂大数据分析的时效性和规模可扩展性。

与此同时，传统的序列计算，即在一个处理器上按照先后顺序进行指令处理，已经不能满足计算要求。并行计算（parallel computing）技术成为提高数据分析系统计算能力的关键手段。并行计算使用多个处理器协同求解同一问题，将被求解的问题分解成若干部分，各个部分均由一个独立的处理器来并行计算。并行计算可以是专门设计的、含有多个处理器的超级计算机，也可以是以某种方式互联的若干台独立计算机构成的集群。其中，分布式计算将软件系统的组成部分多机协作，提高效