

Big
Data
Analytics

and
Computation

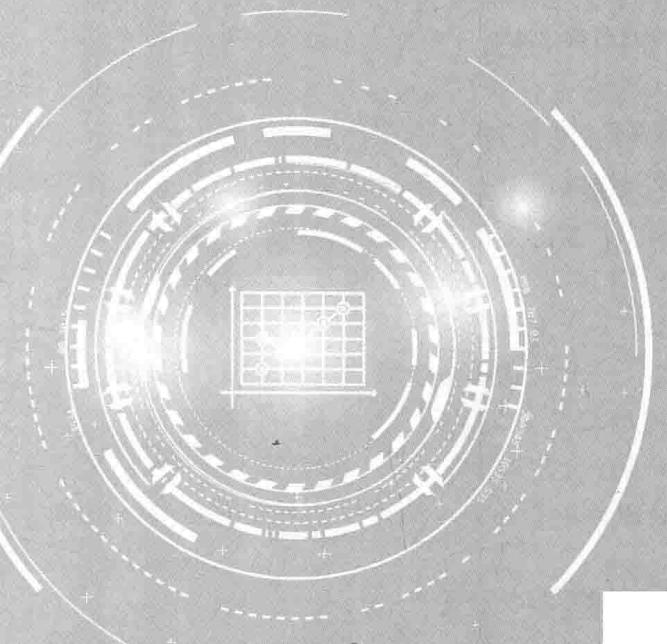
大数据分析 与计算

◎ 汤羽 林迪 范爱华 吴薇薇 编著



清华大学出版社

Big
Data
Analytics



and
Computation

大数据分析 与计算

◎ 汤羽 林邁 范華 吴薇薇 编著



清华大学出版社
北京

内 容 简 介

大数据应用已成为行业热点和产业发展新增长点,数据科学与计算技术也是最新的前沿领域,其中,大数据计算分析提供了核心的技术支撑。本书从大数据计算系统的三个层次对数据模型、处理算法、计算模型与架构、开发技术标准等内容进行了综合性的介绍,重点阐述了各类数据分析算法和MapReduce,图并行计算,交互式处理,流计算,内存计算等计算架构。本书适合作为数据科学与大数据计算技术、计算机科学与技术、互联网应用系统、物联网工程等专业相关课程的教材。

本书封面贴有清华大学出版社防伪标签,无标签者不得销售。

版权所有,侵权必究。侵权举报电话: 010-62782989 13701121933

图书在版编目(CIP)数据

大数据分析与计算/汤羽等编著. —北京: 清华大学出版社, 2018

ISBN 978-7-302-48586-5

I . ①大… II . ①汤… 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(2017)第 316141 号

责任编辑: 贾斌 薛阳

封面设计: 刘键

责任校对: 梁毅

责任印制: 王静怡

出版发行: 清华大学出版社

网 址: <http://www.tup.com.cn>, <http://www.wqbook.com>

地 址: 北京清华大学学研大厦 A 座 邮 编: 100084

社 总 机: 010-62770175 邮 购: 010-62786544

投稿与读者服务: 010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈: 010-62772015, zhiliang@tup.tsinghua.edu.cn

课件下载: <http://www.tup.com.cn>, 010-62795954

印 装 者: 北京鑫海金澳胶印有限公司

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 31.25

字 数: 759 千字

版 次: 2018 年 3 月第 1 版

印 次: 2018 年 3 月第 1 次印刷

印 数: 1~2000

定 价: 89.00 元

产品编号: 070307-01

大数据(Big Data)已被视为硬件、软件、网络之外的第四种计算资源,随着各类大数据应用的兴起,大数据的采集、存储、建模及计算处理已成为分布式计算领域的热门研究课题,也引起产业界极大的兴趣和关注。大数据的计算处理不仅涉及各类数据分析挖掘算法,其计算系统的性能更多依赖于计算模型与计算架构。目前,比较一致的看法是大数据计算系统大致可分为三个层次:数据存储层、数据处理层和数据应用层。数据存储层提供海量数据存储架构与数据访问界面;数据处理层提供对数据分析算法和计算模型的支持;数据应用层则包含各种基于大数据计算分析的应用软件系统。这三个层面都涉及不同的数据模型、计算架构及开发技术标准,目前主流的有两个主线:以 Google 为代表的商业产品和以 Hadoop 为代表的开源技术。在学习和研究大数据计算技术时,需要对上述计算架构、技术和标准有一个总体的了解,这样才能做到不限于一点而把握全局。

针对国家“互联网+”的战略发展需求,近期国内不少高校新开设了数据科学与大数据计算技术专业,大数据分析与计算成为其主干专业课程,其他如计算机科学与技术、互联网应用系统、物联网工程等专业都需要开设大数据计算课程,因此迫切需要一本对大数据处理与计算有一个较全面的论述、适合高年级本科生或研究生学习的教材,正是基于这种需求,本书作者编著了此书,希望对大数据计算系统的各类分析算法、计算模型、计算架构与开发技术做出一个综合性的介绍与阐述,为大家进一步学习大数据技术及应用开发打下基础。

全书共计 20 章,第 1~3 章介绍大数据计算的概念、计算体系总体架构、技术标准等,让读者建立大数据计算的基本概念;第 4~6 章介绍数据采集方法、数据建模及各类分析算法;第 7~10 章介绍文本数据读取、数据处理与分析、数据可视化技术;第 11 章和第 12 章详细介绍 Hadoop 计算平台,包括 HDFS 分布式文件系统与 MapReduce 计算模型;第 13~16 章具体介绍各类大数据计算模型与架构,包括图并行计算、交互式计算、流计算、内存计算等,其中重点阐述了 Pregel、Hama、Storm、Spark 等计算架构;第 17~20 章则介绍了大数据计算技术在医疗保险系统、互联网电子商务、金融信贷系统等领域的应用。本书包含内容较多、篇幅较长,教师在讲授时可根据自己的需要对章节进行选取裁剪。

汤羽教授负责本书的总体结构及第 1~3 章、第 11 章和第 12 章的撰写,林迪副教授负责第 4~10 章,范爱华副教授负责第 13~16 章,吴薇薇硕士负责第 17~20 章。本书部分图片取自互联网,部分文字也参考了网页内容,作者尽可能将引用链接在参考文献罗列中给出,少部分无法给出引用的,作者在此一并致谢。

大数据计算是一个新兴技术领域且仍在高速发展中,新的概念、方法和技术不断涌现。作者因学识有限,本书必然会产生不足,希望得到学界同仁的批评指正,以利我们改进完善。“业精于勤荒于嬉、行成于思毁于随”,作者愿与科学界同行一起努力在这个领域耕耘。

汤 羽

2017年7月于蓉城

目 录

CONTENTS

第 1 章 绪论 ······	1
1.1 数据与数据科学 ······	2
1.2 大数据概念 ······	6
1.3 大数据技术特征 ······	13
参考文献 ······	16
习题 ······	17
第 2 章 大数据计算体系 ······	18
2.1 大数据计算架构 ······	18
2.2 数据存储系统 ······	19
2.2.1 数据清洗与建模 ······	19
2.2.2 分布式文件系统 ······	21
2.2.3 NoSQL 数据库 ······	23
2.2.4 统一数据访问接口 ······	27
2.3 数据处理平台 ······	31
2.3.1 数据分析算法 ······	31
2.3.2 计算处理模型 ······	35
2.3.3 计算平台与引擎 ······	38
2.4 数据应用系统 ······	39
2.4.1 大数据应用领域 ······	39
2.4.2 大数据解决方案 ······	41
参考文献 ······	50
习题 ······	52
第 3 章 大数据标准与模式 ······	53
3.1 大数据标准体系 ······	53
3.2 大数据计算模式 ······	64
参考文献 ······	68

习题	69
第4章 数据采集方法	70
4.1 系统日志采集	70
4.1.1 日志采集的目的	71
4.1.2 日志采集过程	71
4.2 网络数据采集	72
4.2.1 网络爬虫工作原理	72
4.2.2 网页搜索策略	73
4.2.3 网页分析算法	73
4.2.4 网络爬虫框架	74
4.3 数据采集接口	75
参考文献	76
习题	76
第5章 数据清洗与规约方法	77
5.1 数据预处理研究现状	78
5.1.1 数据清洗的研究现状	78
5.1.2 数据规约的研究现状	78
5.2 数据质量问题分类	79
5.2.1 单数据源的问题	79
5.2.2 多数据源的问题	81
5.3 数据清洗技术	82
5.3.1 重复记录清洗	82
5.3.2 消除噪声数据	83
5.3.3 缺失值清洗	83
5.4 数据归约	84
5.4.1 维归约	85
5.4.2 属性选择	85
5.4.3 离散化方法	85
5.5 数据清洗工具	86
参考文献	86
习题	87
第6章 数据分析算法	88
6.1 C4.5 算法	88
6.1.1 算法描述	89
6.1.2 属性选择度量	89
6.1.3 其他特征	91

6.2 k-均值算法	92
6.3 支持向量机	92
6.4 Apriori 算法	93
6.5 EM 算法	94
6.5.1 案例：估计 k 个高斯分布的均值	94
6.5.2 EM 算法步骤	95
6.6 PageRank 算法	96
6.6.1 PageRank 的核心思想	96
6.6.2 PageRank 的计算过程	97
6.7 AdaBoost 算法	99
6.7.1 Boosting 算法的发展历史	99
6.7.2 AdaBoost 算法及其分析	99
6.8 k-邻近算法	101
6.9 朴素贝叶斯	102
6.9.1 朴素贝叶斯分类器	102
6.9.2 贝叶斯网络	103
6.10 分类回归树算法	104
6.10.1 建立回归树	105
6.10.2 剪枝过程	105
参考文献	105
习题	106
 第 7 章 文本读写技术	107
7.1 读取文本文件	107
7.1.1 读取 txt 文件	108
7.1.2 读取 csv 文件	109
7.2 写入文本文件	110
7.3 处理二进制数据	111
7.4 数据库的使用	112
7.4.1 数据库的连接	112
7.4.2 执行 SQL 语句	112
7.4.3 选择和打印	113
7.4.4 动态插入	113
7.4.5 update 操作	114
参考文献	114
习题	114
 第 8 章 数据处理技术	115
8.1 合并数据集	115

8.1.1 索引上的合并.....	115
8.1.2 轴向连接.....	119
8.1.3 合并重叠数据.....	122
8.2 数据转换	123
8.2.1 移除重复数据.....	123
8.2.2 利用函数进行数据转换.....	125
8.2.3 替换值.....	126
8.2.4 重命名轴索引.....	128
8.2.5 离散化数据.....	129
8.2.6 检测异常值.....	133
8.2.7 排列和随机采样.....	134
8.2.8 哑变量.....	136
8.3 字符串操作	138
8.3.1 内置字符串方法.....	138
8.3.2 正则表达式.....	139
8.3.3 Pandas 中矢量化的字符串函数	142
参考文献.....	144
习题.....	144
 第 9 章 数据分析技术.....	145
9.1 NumPy 工具包.....	145
9.1.1 创建数组.....	146
9.1.2 打印数组.....	147
9.1.3 基本运算.....	148
9.1.4 索引、切片和迭代	149
9.1.5 形状操作.....	150
9.1.6 复制和视图.....	151
9.1.7 NumPy 实用技巧	152
9.2 Pandas 工具包	153
9.2.1 Series	154
9.2.2 DataFrame	157
9.3 Scikit-Learn 工具包.....	160
9.3.1 逻辑回归.....	160
9.3.2 朴素贝叶斯.....	161
9.3.3 k-最近邻	161
9.3.4 决策树.....	161
9.3.5 支持向量机.....	162
9.3.6 优化算法参数.....	162
参考文献.....	163

习题	163
第 10 章 数据可视化技术	164
10.1 Matplotlib 绘图	164
10.1.1 Matplotlib API 入门	164
10.1.2 Figure 和 Subplot 的画图方法	164
10.1.3 调整 Subplot 周围的间距	167
10.1.4 颜色、标记和线型的设置	167
10.1.5 刻度、标签和图例	169
10.2 Mayavi2 绘图	172
10.2.1 使用 mlab 快速绘图	172
10.2.2 Mayavi 嵌入到界面中	174
10.3 其他图形化工具	176
参考文献	176
习题	177
第 11 章 Hadoop 生态系统	178
11.1 Hadoop 系统架构	178
11.2 HDFS 分布式文件系统	188
11.2.1 HDFS 体系结构	189
11.2.2 HDFS 存储结构	191
11.2.3 数据容错与恢复	196
11.2.4 Hadoop/HDFS 安装	198
11.3 分布式存储架构	208
11.3.1 HBase 系统架构	209
11.3.2 数据模型与存储模式	214
11.3.3 HBase 数据读写	218
11.3.4 数据仓库工具 Hive	220
11.3.5 HBase 安装与配置	224
11.4 HBase 索引与检索	229
11.4.1 二次索引表机制	229
11.4.2 二次索引技术方案	231
11.5 资源管理与作业调度	240
11.5.1 分布式协同管理组件 ZooKeeper	240
11.5.2 作业调度与工作流引擎 Oozie	246
11.5.3 集群资源管理框架 YARN	249
参考文献	257
习题	259

第 12 章 MapReduce 计算模型	260
12.1 分布式并行计算系统	260
12.2 MapReduce 计算架构	266
12.3 键值对与输入格式	269
12.4 映射与化简	274
12.5 应用编程接口	282
参考文献	295
习题	295
第 13 章 图并行计算框架	296
13.1 图基本概念	302
13.2 BSP 模型	304
13.3 Pregel 图计算引擎	307
13.4 Hama 开源框架	316
13.5 应用编程接口	331
参考文献	335
习题	336
第 14 章 交互式计算模式	337
14.1 数据模型	337
14.2 存储结构	339
14.3 并行查询	347
14.4 开源实现	349
参考文献	357
习题	357
第 15 章 流计算系统	358
15.1 流计算模型	359
15.2 Storm 计算架构	369
15.3 工作机制实现	376
15.4 Storm 编程接口	382
参考文献	388
习题	388
第 16 章 内存计算模式	390
16.1 分布式缓存体系	391
16.2 内存数据库	407
16.3 内存云 MemCloud	412

16.4 Spark 内存计算	419
参考文献	439
习题	440
第 17 章 基于医疗数据的临床决策分析应用	442
17.1 国内外研究现状及发展动态分析	443
17.2 技术路线和方案	444
参考文献	451
习题	451
第 18 章 基于医保数据的预测分析应用	452
18.1 数据准备阶段	452
18.2 模型变量选择和转换	452
18.2.1 模型变量的选择	453
18.2.2 模型变量的转换	454
18.2.3 筛选模型变量	455
18.3 建模过程	456
18.4 模型效果	457
参考文献	458
习题	458
第 19 章 互联网电商数据的分析应用	459
19.1 电商流程管理分析	460
19.1.1 行业背景与业务问题	460
19.1.2 分析方法与过程	460
19.2 用户消费行为分析	464
19.2.1 业务问题	464
19.2.2 分析方法与过程	465
19.3 送货速度相关性分析	466
19.3.1 业务问题	466
19.3.2 分析方法与过程	467
19.4 总结	469
参考文献	469
习题	469
第 20 章 金融和经济数据的分析应用	470
20.1 企业对创新经济活动推动的影响分析	470
20.1.1 案例背景	471
20.1.2 分析方法与过程	471

20.2 信贷风险模型评估	477
20.3 中小能源型企业的信用评价分析	480
20.3.1 案例背景	480
20.3.2 分析方法与过程	480
20.3.3 分析结果	483
参考文献	484
习题	485

第 1 章

绪 论

进入 21 世纪,人类发现自己正面临着中国唐代诗人李白描述的“黄河之水天上来”(图 1-1)的大数据场景:互联网搜索引擎 Google 每天完成 10 亿次查询,社交网站 Facebook 每天处理 80 亿条信息;在科学领域,2003 年基本完成的基因组计划完成了四十多种生物全基因组测序以及 3.2×10^9 人类基因组碱基对测序,到 2006 年 DNA 碱基数目已超过 1300 亿,目前全世界每年生物数据产出量估计已达 10^{15} B(1PB),且以每三年翻一番的速度增长;在金融领域,中信银行 2008 年发卡 500 万张,2010 年则翻了一倍,带来了海量数据需要处理;2012 年,国家工业和信息化部宣布,中国移动互联网用户已达 7.5 亿。与此同时,随着智慧城市、物联网等新兴应用模式的发展,各种摄像头、数字标牌、感应装置、检测装置以及嵌入式终端的数量也在急剧增加,有关数据预测显示:作为物联网一个重要组成部分的射频识别(RFID)标签,其销量将从 2011 年的不到 30 亿个发展到 2021 年



图 1-1 黄河之水天上来

的 2090 亿个。在互联网浏览搜索、物联网传感数据、移动终端与 GPS 系统、以及社交网络等领域,全世界的信息量以每两年翻番的速度增长。据国际研究机构 IDC 报告:2011 年,全球数据量为 1.8ZB(1ZB=10⁶PB=10⁹TB=10¹²GB),2015 年达到 8ZB,2020 年将达到 35ZB。

1.1 数据与数据科学

1. 数据定义

数据(Data)被看作现实世界中自然现象和人类活动所留下的轨迹^[1]。在计算机科学中,数据的定义是指所有能输入到计算机并被计算机程序处理的符号的总称,是具有一定意义的数字、字母、符号和模拟量的统称。《韦伯斯特大词典》(Merriam-Webster Dictionary)把数据定义为“用于计算、分析或计划某种事物的事实或信息;由计算机产生或存储的信息(facts or information used usually to calculate, analyze or plan something; information that is produced or stored by a computer)”^[2]。事实上,数据的形式多样化,可以表现为数值、文字、图像、音频、视频或其他计算机可以识别和处理的形式,数据来源也可以是社会数据(商业数据、生产数据、系统数据、媒体数据等)、个人数据(社交网络、个人消费)、政府数据(统计数据、人口普查、经济年报等)。人类四千年历史所产生的所有的文明记录,包括历史、文学、艺术、哲学、考古及一切的科学成就,都可以数据的形式存储和保留下来。

2. 数据简史

在人类文明有记载的四千年历史中,人类活动记录从早期古埃及的结绳记事(图 1-2)到中国殷商时期(公元前 1320—1046 年)的甲骨文(图 1-3),从东汉宦官蔡伦(公元 61—121 年)发明造纸术(图 1-4)到北宋布衣毕昇(公元 61—121 年)发明活字印刷术(图 1-5),文明的记录无不以文字(数据的一种形式)传承下来。

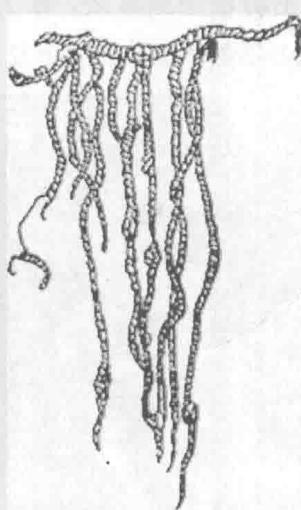


图 1-2 结绳记事



图 1-3 殷商甲骨文



图 1-4 蔡伦造纸术



图 1-5 毕昇活字印刷术

到了近代,面对日益增长的计算量,人工计算的方式已难以应对,人类进入了机械计算时代。这一时期的代表有:1642年,法国数学家布莱士·帕斯卡(Blaise Pascal)发明的机械式加减法计算机;1671年,德国数学家戈特弗里德·莱布尼茨(Gottfried W. Leibniz)制成了第一台能够进行加减乘除运算的机械式计算机;1833年,英国科学家查尔斯·巴贝奇(Charles Babbage)提出了自动化分析机的设计,第一次提出了机器可执行预先记录在穿孔卡上的指令来完成计算的思路。与之合作的艾达·洛夫莱斯伯爵夫人(Augusta Ada Lovelace)更进一步为巴贝奇机器写出了计算注释笔记,其中包含由机器执行的计算伯努利数的算法和步骤,被认为是人类完成的第一个可执行的计算机程序,艾达也因此作为人类历史上的第一个计算机程序员而被纪念。

进入20世纪,1904年,英国工程师约翰·弗莱明(John A. Fleming)发明了真空二极管;1906年,美国工程师李·弗雷斯特(Lee de Forest)发明了真空三极管;1947年,美国贝尔实验室的科学家威廉·肖克利(William B. Shockley)、约翰·巴丁(John Bardeen)和沃尔特·布拉顿(Walter Brattain)组成的研究小组发明了晶体管;1958年,美国物理学家杰克·基尔比(Jack Kilby)发明了集成电路技术。电子管、晶体管、集成电路技术及随后的超大规模集成电路(VLSI)的诞生使得人类真正进入了电子计算机时代。1946年2月,世界上第一台电子管计算机“埃尼阿克”(ENIAC)在美国宾夕法尼亚大学诞生;1956年,美国贝尔实验室研制出第一台晶体管计算机 TRADIC;1964年,美国 IBM 公司推出第一代采用集成电路的电子计算机 IBM360 系列,这以后各种计算能力更为强大、数据处理能力呈爆炸式增长的超级计算机、高端服务器、图形工作站以及计算机集群层出不穷,使得计算机处理数据的能力从早期的 KB、MB 级别达到了今天的 TB 或 PB 量级。2015年,美国著名社交网站 Facebook 每天需处理 100 亿条消息和 3.5 亿张新图片,而谷歌(Google)每天应对的查询请求达到 30 亿次,后台处理的数据量达到 85TB。

应当注意的是,数据(Data)、信息(Information)、知识(Knowledge)与价值(Value)这4个词在信息科学中既相互关联,又具有不同的含义。数据体现的是一种过程、状态或结果的记录,这类记录被数字化后可以被计算机存储和处理。信息则是包含在数据之中的能够为人脑理解和思维推理和结论。例如,“01001000 01100101 01101100 01101100 01101111

00100000 01110111 01101111 01110010 01101100 01100100 00100001”是一串二进制数值，是一组能被计算机识别、存储和处理的数据。经计算机程序识别转换(ASCII 码值字符转换)，我们知道它代表“Hello world!”这样一个字符串，包含向世界问好的特殊信息。更进一步，在计算机编程语言世界，“Hello world!”实际上是一个约定俗成的机器或程序语言启动显示语句，这就上升为知识。最终，如果有人把这一固有的显示方法拿去注册了专利并因此获利，就产生了价值。图 1-6 表征了这一从数据到信息到知识到价值的过程。

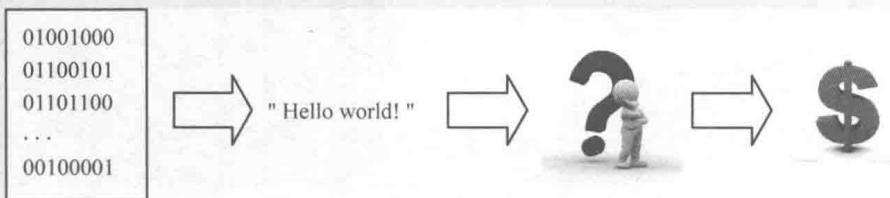


图 1-6 数据-信息-知识-价值的转换过程

3. 数据科学

当信息科学处理的数据发展到 Facebook 和 Google 的数据规模，数据本身(类型、规模、属性、用途等)及相关的规模数据分析计算技术就形成了一门新的学科领域：数据科学(Data Science)或数据工程(Data Engineering)。早在 1974 年，丹麦计算机科学家、2005 年图灵奖得主彼得·诺尔(Peter Naur)即提出了“数据学”(Datalogy)的概念^[3]，但他更多指的是以数据为对象的计算机科学和编程，认为“数据学”是计算机科学的延伸，其研究对象是数值化的数据。事实上，在个人计算机(PC)出现之前的早期的计算机确实是更多地用于处理数据和科学计算，计算机大大加快了数据处理的速度、效率和准确性，但作为计算机运算对象和输出结果的数据本身尚未引起科学家们的特别注意。

1996 年，在日本神户的一个国际会议上第一次正式使用“数据科学”这一名称^[4]。1997 年，密歇根大学教授杰夫·吴(Jeff C. Wu)在演讲中提出“统计学=数据科学?”的观点并建议将统计学改名为数据科学，统计学家改名为数据科学家^[5]。2001 年，贝尔实验室科学家威廉·克里富兰(William S. Cleveland)第一次提出数据科学应作为由统计学延伸出来的一个独立研究领域，认为统计学中与数据分析有关的技术内容(区别于概率论)在下面 6 个方面扩展后形成一个新的独立学科“数据科学”(Data Science)^[6]。

- (1) 多学科研究(Multidisciplinary Investigations)；
- (2) 数据模型与分析方法(Models and Methods for Data)；
- (3) 数据计算(Computing with Data)；
- (4) 数据学教程(Pedagogy)；
- (5) 工具评估(Tool Evaluation)；
- (6) 理论(Theory)。

在 2002 年和 2003 年，国际科学委员会(International Council for Science)和哥伦比亚大学分别创办了数据科学杂志，为这一学科领域的研究工作发表和交流建立了国际学术平台。大规模数据计算的特点和重要性已引起科学界注意，数据科学或数据处理技术被有些科学家认为将成为一个与计算科学并列的新科学领域。已故著名图灵奖获得者 Jim Gray 在 2007 年的一次演讲中提出，数据密集型科学发现(Data-Intensive Scientific Discovery)将