



大数据丛书系列之四

总主编◎曾 羽 龙奋杰

大数据 科学

DASHUJU
KEXUE



主 编◎霍雨佳 周若平 钱晖中



电子科技大学出版社

大数据丛书系列之四

总主编◎曾 羽 龙奋杰



主 编◎霍雨佳 周若平 钱晖中



电子科技大学出版社

图书在版编目(CIP)数据

大数据科学 / 霍雨佳, 周若平, 钱晖中主编. -- 成都: 电子科技大学出版社, 2017.7

ISBN 978-7-5647-4820-3

I. ①大… II. ①霍… ②周… ③钱… III. ①数据处理 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2017) 第 176277 号

大数据科学

霍雨佳 周若平 钱晖中 主编

策划编辑 杨仪玮 李燕苓

责任编辑 杨仪玮 李波翔

出版发行 电子科技大学出版社

成都市一环路东一段 159 号电子信息产业大厦 邮编 610051

主 页 www.uestcp.com.cn

服务电话 028-83203399

邮购电话 028-83201495

印 刷 成都市火炬印务有限公司

成品尺寸 165mm × 240mm

印 张 18.5

字 数 352 千字

版 次 2017 年 7 月第一版

印 次 2017 年 7 月第一次印刷

书 号 ISBN 978-7-5647-4820-3

定 价 65.00 元

目 录

第1章 大数据概述	2
1.1 产生背景	2
1.1.1 互联网与大数据	3
1.1.2 信息技术与大数据	6
1.1.3 云计算与大数据	9
1.1.4 物联网与大数据	11
1.1.5 智能终端与大数据	13
1.2 大数据的概念	15
1.2.1 狹义的大数据概念	15
1.2.2 广义的大数据	16
1.3 大数据特征	17
1.3.1 大数据体量巨大	17
1.3.2 大数据类型多样	19
1.3.3 数据处理与流动速度快	20
1.3.4 数据价值密度低	21
1.4 发展大数据的意义	21
1.4.1 大数据创新科学的研究	23
1.4.2 大数据是实现行业融合发展的需要	25
1.4.3 大数据是助推产业转型升级的加速器	26
1.4.4 大数据助力智慧城市建设	28
1.4.5 大数据创新商业模式	29
第2章 起源与发展历程	34
2.1 大数据发展起源	34
2.2 大数据发展历程	35
2.2.1 突破阶段（2000—2006年）	35

2.2.2 成熟阶段（2006—2009年）	36
2.2.3 完善发展阶段（2010年至今）	36
2.3 大数据发展趋势	38
第3章 大数据引发的哲学思考	46
3.1 大数据与世界观	46
3.1.1 数据的本体论主张	46
3.1.2 大数据及其本质	48
3.2 大数据与方法论	49
3.2.1 整体和部分的结合，实现了还原论与整体论的辩证统一	49
3.2.2 承认事物的多样性，地方性知识引起重视	50
3.2.3 突出相关性而不是因果性	52
3.3 大数据与认识论	53
3.3.1 数据挖掘与科学知识	54
3.3.2 数据规律及其真理性	56
第4章 大数据来源	60
4.1 数据的概念和分类	60
4.1.1 什么是数据？	60
4.1.2 我们生活在数据的世界里	60
4.1.3 数据的分类	60
4.1.4 数据科学	63
4.2 常用数据采集方法	64
4.2.1 数据采集的概念	64
4.2.2 传统的数据来源	65
4.2.3 传统的数据采集	66
4.2.4 大数据环境下的数据来源	66
4.2.5 大数据的数据采集方法	68
4.3 常用的数据采集工具	71
4.3.1 传统数据采集的常用工具	71
4.3.2 大数据采集的常用工具	77
第5章 大数据存储技术	82
5.1 数据库系统原理	82

5.1.1 数据库系统的基本概念	82
5.1.2 数据库系统的特点	84
5.1.3 数据模型	85
5.1.4 数据库技术的发展趋势	88
5.2 关系型数据库	88
5.2.1 关系模型的基本概念	88
5.2.2 关系操作和关系完整性	90
5.2.3 关系代数	91
5.2.4 SQL	95
5.2.5 关系数据库的其他内容	96
5.3 非关系型数据库	96
5.3.1 非关系型数据库概述	96
5.3.2 非关系型数据库的分类	98
5.3.3 非关系型数据库的发展瓶颈和发展前景	100
5.4 大数据常用存储平台	100
5.4.1 大数据对存储平台的要求	101
5.4.2 大数据常用存储平台介绍	102
5.4.3 大数据存储平台	102
第6章 大数据与统计学	107
6.1 大数据中的统计理论	107
6.1.1 统计学概念和发展历程	107
6.1.2 大数据与统计学的相互作用	109
6.1.3 大数据分析	110
6.2 机器学习和数据挖掘	112
6.2.1 机器学习与数据挖掘概述	112
6.2.2 传统机器学习和基于大数据的机器学习	113
6.2.3 数据挖掘的内涵	114
6.2.4 数据挖掘经典算法	116
6.3 数据学与数据科学	119
6.3.1 数据学与数据科学的基本内容	119
6.3.2 数据学的框架	122
6.3.3 数据学与数据科学的研究对象和内容	123

6.3.4 数据科学家	124
6.3.5 数据学与数据科学和大数据科学的关联	126
6.3.6 小结	128
6.4 常用工具	128
6.4.1 统计产品与服务解决方案软件——SPSS	128
6.4.2 统计分析软件——SAS	130
6.4.3 统计计算和统计制图工具——R	131
6.4.4 小结	132
第7章 大数据常用技术和服务平台	133
7.1 大数据编程模型	133
7.1.1 编程模型的概念以及大数据常用的编程模型	133
7.1.2 常见的大数据编程模型	134
7.1.3 MapReduce 实现单词计数的运算过程	135
7.1.4 MapReduce 的主要特征	138
7.2 大数据处理平台	140
7.2.1 Hadoop 的相关概念和应用场景	140
7.2.2 Storm 的相关概念和应用场景	142
7.2.3 Spark 的相关概念和应用场景	144
7.2.4 小结	145
7.3 数据可视化	146
7.3.1 数据可视化的概念和内涵	146
7.3.2 数据可视化分类	147
7.3.3 大数据时代数据可视化的发展趋势	149
第8章 安全大数据	152
8.1 大数据形式下的反恐新形式	152
8.1.1 什么是恐怖主义	152
8.1.2 我国面临的恐怖主义及反恐工作	152
8.1.3 大数据在反恐工作中的应用	154
8.2 依托大数据的网络信息安全	156
8.2.1 网络信息安全	156
8.2.2 大数据技术在网络信息安全中的应用	158

目 录

8.3 基于大数据的自然灾害预警	162
8.3.1 地震灾害预测	162
8.3.2 总结	166
8.4 大数据下的安全生产	166
第9章 商业大数据	171
9.1 精准营销	171
9.1.1 什么是市场营销	171
9.1.2 什么是精准营销	172
9.1.3 精准营销案例分析	173
9.2 决策支持	174
9.2.1 什么是决策支持系统	174
9.2.2 大数据下的烟草物流决策支持系统	177
9.3 创新模式	178
9.3.1 商业模式创新分析	178
9.3.2 大数据时代商业模式的创新	179
第10章 民生大数据	182
10.1 基于大数据的智慧旅游	182
10.1.1 大数据下智慧旅游的数据特征	182
10.1.2 大数据在智慧旅游中的应用	183
10.1.3 总结与展望	185
10.2 基于大数据的智能交通	185
10.2.1 什么是交通大数据	185
10.2.2 大数据下智能交通的特点	187
10.2.3 大数据在智能交通中的应用	187
10.3 大数据环境下的食品安全	188
10.3.1 我国的食品安全问题	189
10.3.2 大数据在食品安全问题中的应用	190
10.4 大数据环境下的智慧医疗	192
10.4.1 什么是智慧医疗	192
10.4.2 大数据与智慧医疗	193
10.4.3 总结	194

10.5 大数据与教育	194
10.5.1 什么是教育大数据	194
10.5.2 大数据对教育的影响	195
第11章 政务大数据	199
11.1 基于大数据的网络舆情分析	199
11.1.1 什么是网络舆情	199
11.1.2 网络舆情大数据特征	199
11.1.3 网络舆情分析方法	200
11.1.4 政府网络舆情管理	201
11.2 基于政务大数据的精细化管理和服务	203
11.2.1 以大数据提升政府的科学决策水平	203
11.2.2 以大数据提升政府管理效率降低管理成本	204
11.2.3 以大数据促进政府的精细化、人性化	204
11.2.4 利用政务大数据实现服务精准化	205
11.3 大数据下应急预案处理	207
11.3.1 大数据时代我国应急预案管理体系变革机遇	207
11.3.2 大数据提升政府应急预案管理能力	210
11.3.3 总结	211
第12章 工业大数据	212
12.1 智能装备	212
12.2 智慧工厂	213
12.3 智能服务	216
12.3.1 加速产品创新	217
12.3.2 产品故障诊断与预测	217
12.3.3 工业物联网生产线的大数据应用	218
12.3.4 产品质量管理与分析	218
12.3.5 生产计划与排程	219
第13章 大数据学科发展概述	222
13.1 大数据发展概况	222
13.1.1 大数据概述	222
13.1.2 大数据发展现状	223

13.1.3 大数据应用现状	224
13.1.4 大数据发展前景	225
13.1.5 大数据发展趋势	225
13.2 国家政策环境	226
13.2.1 国外大数据相关政策	226
13.2.2 国外大数据发展举措	228
13.2.3 我国大数据发展启示	230
13.2.4 我国大数据政策支持	231
13.3 大数据学科概述	233
13.3.1 大数据学科发展现状	233
13.3.2 大数据学科发展趋势	235
13.3.3 大数据学科研究热点	236
13.3.4 大数据技术发展趋势	237
13.3.5 大数据学科就业分析	240
第14章 大数据学科构建	242
14.1 大数据学科建设理念	242
14.2 大数据学科建设目标	244
14.3 大数据学科建设方案	245
14.2.1 专业方向设置方案	245
14.2.2 教学内容设置方案	246
第15章 大数据人才培养	250
15.1 大数据学科人才培养概述	250
15.1.1 大数据分析方向	250
15.1.2 大数据平台方向	252
15.1.3 深度计算分析方向	253
15.1.4 国际合作方面	254
15.2 大数据课程体系	255
15.2.1 大数据平台方向	255
15.2.2 大数据分析方向	258
15.3 专科大数据人才培养	260
15.3.1 专科大数据人才培养必要性	260
15.3.2 专科大数据人才培养模式	260

15.3.3 专科大数据专业课程体系	261
15.4 本科大数据人才培养	262
15.4.1 本科大数据人才培养必要性	262
15.4.2 本科大数据人才培养模式	263
15.4.3 本科大数据专业课程体系	264
15.5 研究生大数据人才培养	266
15.5.1 研究生大数据人才培养必要性	266
15.5.2 美国研究生大数据人才培养现状	267
15.5.3 研究生大数据人才培养模式	274
15.5.4 研究生大数据专业课程体系	274
参考文献	278

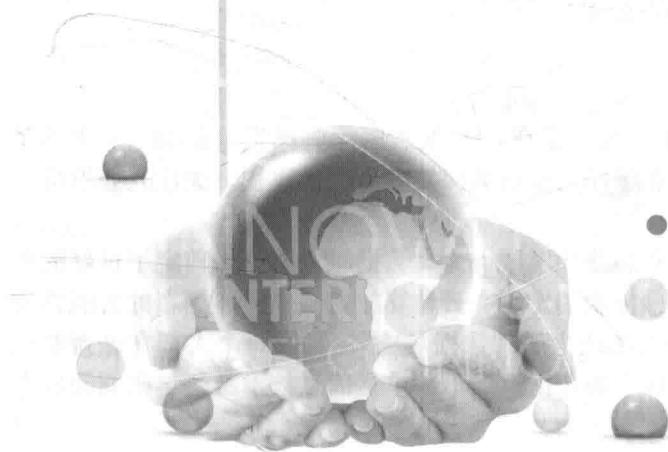
第1篇

基础篇

第1章 大数据概述

第2章 起源与发展历程

第3章 大数据引发的哲学思考



第1章 大数据概述

“大数据”，或称之为海量数据，一般指所含的数据集规模巨大，现在大众的软件工具无法在合理的时间进行采集、存储、分析管理的数据信息。因其在各个行业的广泛应用，使之关注热度历年来居高不下。作为人们获得新的认知、理念和创造价值的源泉，大数据数据来源可以囊括我们从日常生活中可以普遍见到的上传到网页上的图像、视频、录音；高速公路上车辆与收费记录、日常监控录像、医院的治疗病例、高端的基因测序，天文学中通过望远镜收集的信息数据等。

根据国际数据公司（IDC）《数字世界》研究项目在2012年对全球数据量的统计结果可知：2011年，全球数据量的规模将达到里程碑式的1.8ZB，2010年全球产生的数据量为1.2ZB，2009年为0.8ZB，全球数字数据量每两年翻一番^①。据预测，中国数据量规模在2012年至2020年间将从364EB增至8.6ZB，中国在全球数字世界中所占的份额将在2012年至2020年间从13%增至21%。如果你对这些数字仍然感到难以把控的话，接下来一组名为“互联网上一天”的数据可以清晰地告诉你，一天之内，互联网产生的全部内容可以刻满1.68亿张DVD；发出的邮件有2940亿封之多（相当于美国两年的纸质信件数量）；发出的社区帖子达200万个（相当于《时代》杂志770年的文字量）……^②人类计量数据量的单位也从TB级别上升到PB、EB乃至ZB级别。毫无疑问，大数据时代已经来临。

1.1 产生背景

从2006到2016年，大数据在IT行业、医疗、民生、金融、学术等多个领域中炙手可热，行业领导人也对其保持着高度的重视，关注其能够带来的科学价值和社会价值。

第三次科学技术革命的蓬勃发展，为大数据时代的到来奠定了良好的基础。互联网的普及、信息技术的发展、云计算的发展成熟、遍布的智能终端等，无时无刻不记录着人类产生的“数据足迹”，我们每个人都在毫无意识地成为数据的提供者。如今的数据已经不单单是信息技术及科学研究领域人

① <http://news.makpolo.com/1451080.html>

② http://www.360doc.com/content/15/0311/15/4786202_454311432.shtml

员的专有名词，冠有量词属性的“大数据”早已成为一种引人注目的新思潮，成为人们认识事物、分析事物、探索新发现及追求创新的新范畴。故在此，对大数据的产生背景进行一个梳理。

1.1.1 互联网与大数据

互联网技术的不断普及，其每天所催生的巨量数据使得世间万物不断走向数据化，数据量化的节奏不断加快。在由“万事皆数”过渡到“万物皆数”的过程中，互联网每天所产生的数据，对大数据时代的来临有着关键性作用。

(1) 互联网的产生

互联网的发明最初是根据军事的需要而产生的。1969年，美军在ARPA（阿帕网，美国国防部研究计划署）制定的协定下，首先用于军事连接，后将美国西南部的加利福尼亚大学洛杉矶分校、斯坦福大学研究院、UCSB（加利福尼亚大学）和犹他州大学的四台主要的计算机连接起来^①。这个协定由剑桥大学的BBN和MA执行，在1969年12月开始联机执行。随着计算机在战争中军事上的广泛运用，计算机上保存的各自军事机密就越多，其安全性显得尤为重要，军事家们担心如果计算机上的重要军事机密数据泄露的话，会导致整个战争的失败。因此，如何在原有计算机上除能够进行储存数据外，能够通过某种渠道促使两台计算机，或是多台计算机之间能够进行数据之间的相互传递和备份，成为了军事家们的现实需要并极力促进其实现，这促进了早期互联网的形成。

20世纪七八十年代，我们对各类远程信息的接收，大多是依靠收音机或者电视这样的渠道。这种模式的信息数据传播主要是通过信息源利用信号塔进行信息四面八方传播，接收者利用收音机、电视进行信息接收。需要注意的是，这种“信号塔”与“接收”的模式，在整个过程中，接收者是被动的，而且发送信息方也对所传播的信息质量、受众数量、受众偏好等信息数据没有一个完全的记录。

互联网的诞生，使人类进行信息交流的方式发生了质的变化。经济的高速发展，电子产品更新换代进程的加快，使得互联网与大众的基础规模日益扩大。人类也越来越习惯通过互联网接收和传播信息数据。与传统的信息接收和传播模式相比，互联网模式只有当用户利用互联网进行访问，服务站在接收到请求后，会立即在万千信息中寻找到用户所需求的信息内容进行回馈。在这个过程中，我们的用户客户端早已记录下了用户的访问记录数据量、内容、停留的时间等多种信息数据。尽管用户会删除这些数据，但数据

^① <http://www.jishi60.com/show/20/5455-0.html>

后台已将这些记录保存下来，这些留下来的海量信息数据其背后所蕴含着难以估量的巨大价值。

在淘宝网上营业的一家食品销售店铺，为了扩大自身的竞争力和更准确地知道消费者的需求，需要对消费者购买意愿及购买能力进行掌握。为了解这些信息，需要对进入店铺的消费者进行一个调查研究，所需了解信息包括：长久时间以来顾客进入网上店铺的浏览量？每次浏览的食品有哪些？加入购物车的食品是什么品类？在店铺停留的时间有多长？从确定购买到支付完成的消费周期有多久？如果按照传统的销售策略，进行该项目就需要一系列的工作，或是进行大样本调研问卷进行调查获得数据资料，或是安装摄像头进行监控调查，或是通过选取忠诚的消费者样本，电话咨询；不管怎么样，都需要大量的人力、物力、财力。而如果想节约成本的话，就失去了消费者购买商品后所留下的大量消费信息。互联网大数据的出现会使这一系列的成本降得微乎其微。通过互联网上的用户每次操作访问记录，就可以很清楚地得到这一系列的数据。在这过程中，当消费者点击淘宝上店铺的页面，就相当于进入了店铺，对展示的食品进行点击，就类似于消费者拿起食品进行分析比较，与商家进行谈话交流，可以测算出消费者的关注问题，把食品放在购物车就相当于试穿，通过购物车进行结算就相当于购买，而从最开始进入店铺到支付结算，关闭网页的时间记录就相当于一个消费周期，消费者通过不断的对比所确认购买的服装价格相当于此类消费者的购买能力，同时对于消费者的类型、品味兴趣都有了一个精准的记录。并且获得这些大样本数据的成本为“零”。在实体店中的消费者行为和心理都如实地反映在了访问记录中，而且还可以大幅度降低商家的存储、操作成本！

（2）互联网催生大数据

从1969年互联网诞生到2016年，互联网技术已发展成熟。世界经济技术的高速发展，使互联网作为王榭的堂前燕飞入寻常百姓家的普及面日益扩大。

根据世界银行报告报告，截止到2016年初，七大人口大国之中，美国、日本、俄罗斯的互联网普及率最高，分别为87.4%、90.6%和70.5%，巴西、中国次之，分别为57.6%和49.3%，而印度和印尼最少，分别在18%和33%（见图1-1）。从中国的互联网普及率来看，中国互联网络信息中心（CNNIC）在京发布第38次《中国互联网络发展状况统计报告》显示，截止到2016年6月，中国网民规模达7.10亿，互联网普及率达到51.7%，超过全

球平均水平3.1个百分点（见图1-2）。移动互联网塑造的社会生活形态进一步加强，“互联网+”行动计划推动政企服务多元化、移动化发展。

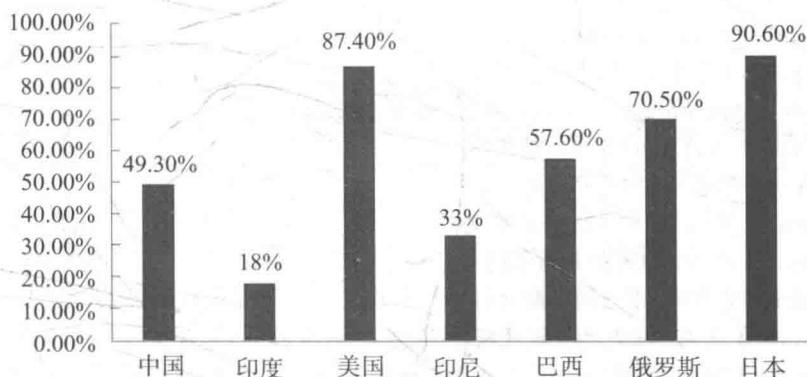
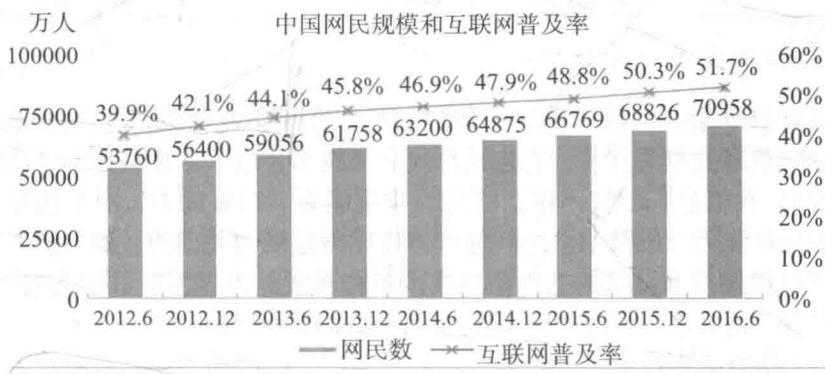


图1-1 全球人口大国互联网普及率



来源：CNNIC中国互联网络发展状况统计调查

2016.6

图1-2 中国网民规模和互联网普及率

互联网的迅猛发展和快速普及，使得大量的数据信息在采集、存储、传输、处理、管理等方面越来越便捷。同时，互联网的发展也使得其所产生的数据类型变得复杂多样，由最初的结构化数据，发展到非结构化数据、半结构化数据等。就大数据而言，在互联网上一天，都会潜在的拥有众多数据的“产生者”和“发送者”，这些“产生者”和“发送者”每时每刻都贡献出各种各样、难以计量的数据。这些数据可以是结构化数据，如数字、符号，也可以是非结构化数据，如文本、图象、声音、影视、超媒体等。这些接连不断出现的数据，催生着大数据浪潮的来临。根据Excelcom公司2016年发布的一份“互联网一分钟产生数据”的表1-1，可以看出互联网对于大数据的催生作用到底有多大。

表 1-1 互联网一分钟产生的数据量

-
- ① Netflix 共有 69 444 小时长的视频被观看；
 - ② 1 亿 5 千万封电子邮件已发送；
 - ③ Uber 产生了 1389 次驾驶；
 - ④ Snapchat 上分享了 527 760 张照片；
 - ⑤ Apple's AppStore 上已有 51 000 个 app 被下载；
 - ⑥ Amazon 产生了 \$203 596 的销售额；
 - ⑦ LinkedIn 创建了 120 多个新账号；
 - ⑧ Twitter 上发布了 347 222 条新推文；
 - ⑨ Google 上产生 240 万的新搜索请求；
 - ⑩ Tinder 上又有 972 222 的新配对；
 - ⑪ YouTube 上已有 278 万的视频被观看；
 - ⑫ Spotify 上的音乐播放时长已达 38 052 小时；
 - ⑬ Facebook 共产生 701 389 个账号登陆；
 - ⑭ Vine 上的小视频播放了 100 万次。
-

1.1.2 信息技术与大数据

信息技术指人们获取信息、传递信息、存储信息、处理信息、显示信息、分配信息的相关技术，它包括现代通信技术、电子计算机技术、微电子技术等^①。在信息科学技术中，信息的收集能力、存储能力、对于信息数据的处理分析能力、以及信息之间远程传递能力是最为关键的，如果把人类历史上信息科技技术发展的五次信息技术革命看成是世界不断数字化的过程，就会发现一条信息技术进步与大数据时代来临的逻辑线^②。

(1) 信息采集技术

获得信息，是人认识世界的基础，因此也是改造世界和人类本身发展的基础。人类从客观世界获得的信息越多、内容越丰富，则人对客观世界的认识也越深刻，对客观世界的利用、改造及对人类本身发展也越有利。传统数据采集来源单一，且存储、管理和分析数据量也相对较小，大多采用关系型数据库和并行数据中心即可处理。大数据时代下，如何从数据量极大、增长速度极快、数据类型复杂、实用性极高的数据中采集出使用者所需要的信息，为传统的数据采集技术提出了难题。2006 年后出现的数据采集技术使得这一难题得到攻克。这些技术包括 Hadoop 的 Chukwa，Cloudera 的 Flume，

^① 杨强. 信息技术的发展历程及其未来趋势[J]. 美丽中国, 2009 (3): 120-123.

^② 赵国栋, 易欢欢, 麋万军等著. 大数据时代的历史机遇: 产业变革与数据科学[M]. 北京: 清华大学出版社.