

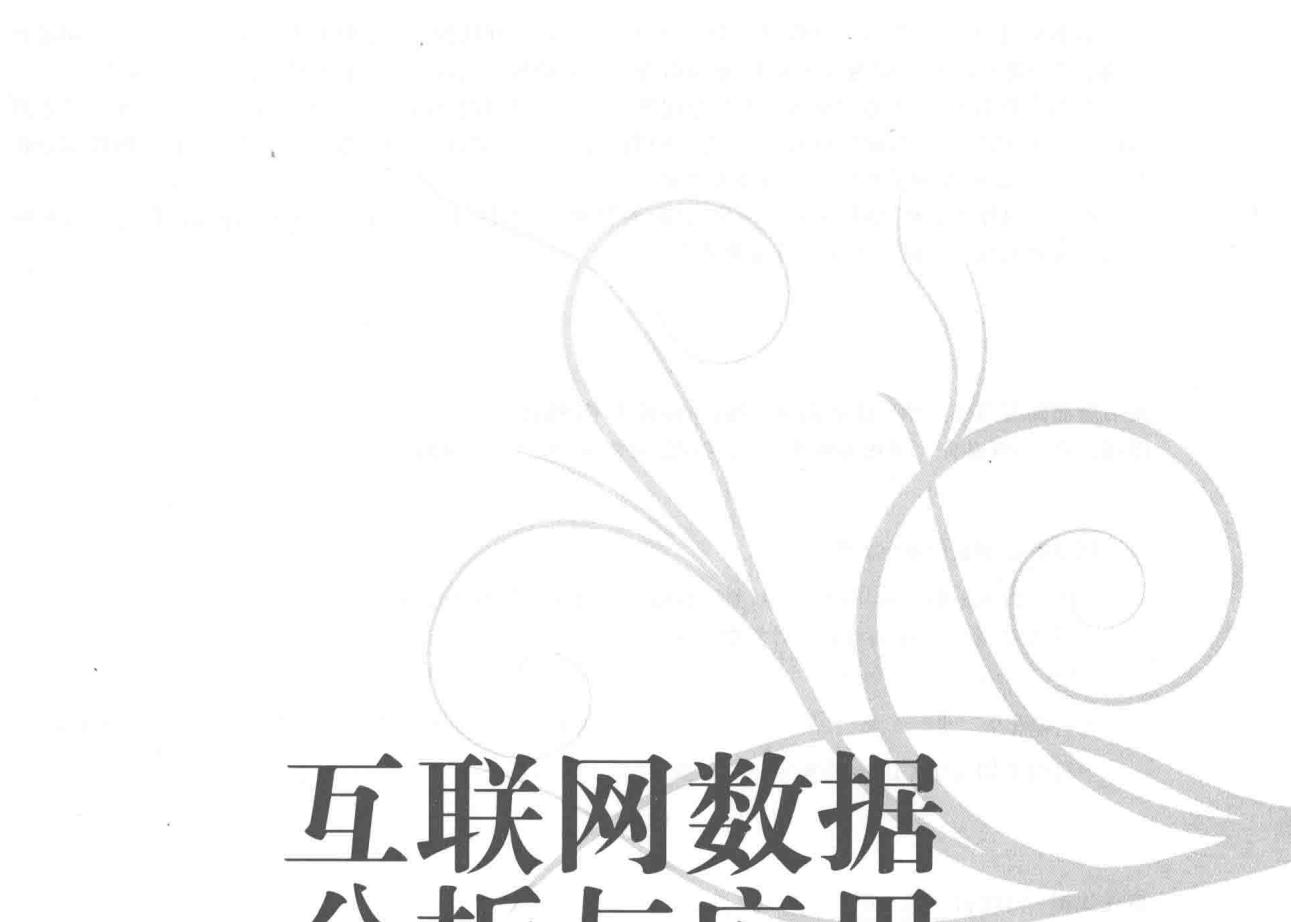


互联网数据分析与应用

赵守香 姜同强 编著

清华大学出版社

21世纪高等学校规划教材 | 计算机应用



互联网数据 分析与应用

赵守番 姜同强 编著

清华大学出版社
北京

内 容 简 介

本书从“挖掘数据的潜在价值和应用”的视角出发,针对互联网企业、政府、传统企业、个人等不同的社会角色对数据挖掘和分析的应用需求,系统地介绍了大数据时代数据分析的作用、技术和具体应用。

本书共分为8章,内容包括大数据与数据分析、互联网数据存储、互联网数据分析工具、商务网站数据分析与应用、政府网站数据分析及应用、物联网数据分析与应用、移动商务数据分析与应用、微博数据分析与应用等,全面地介绍了数据分析的相关内容。

本书可以作为高校“数据分析与应用”、“商务智能”等本科、研究生课程的教材,也可以作为从事互联网数据管理和数据分析的专业人员的参考书。

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

互联网数据分析与应用/赵守香,姜同强编著. --北京: 清华大学出版社, 2015

21世纪高等学校规划教材·计算机应用

ISBN 978-7-302-37974-4

I. ①互… II. ①赵… ②姜… III. ④互联网—数据管理—高等学校—教材 IV. ①TP393.4

中国版本图书馆 CIP 数据核字(2014)第 209523 号



责任编辑: 吴红梅 薛 阳

封面设计: 傅瑞学

责任校对: 李建庄

责任印制: 李红英

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 刷 者: 北京富博印刷有限公司

装 订 者: 北京市密云县京文制本装订厂

经 销: 全国新华书店

开 本: 185mm×260mm 印 张: 19.5 字 数: 484 千字

版 次: 2015 年 9 月第 1 版 印 次: 2015 年 9 月第 1 次印刷

印 数: 1~2000

定 价: 34.50 元

出版说明

随着我国改革开放的进一步深化,高等教育也得到了快速发展,各地高校紧密结合地方经济建设发展需要,科学运用市场调节机制,加大了使用信息科学等现代科学技术提升、改造传统学科专业的投入力度,通过教育改革合理调整和配置了教育资源,优化了传统学科专业,积极为地方经济建设输送人才,为我国经济社会的快速、健康和可持续发展以及高等教育自身的改革发展做出了巨大贡献。但是,高等教育质量还需要进一步提高以适应经济社会发展的需要,不少高校的专业设置和结构不尽合理,教师队伍整体素质亟待提高,人才培养模式、教学内容和方法需要进一步转变,学生的实践能力和创新精神亟待加强。

教育部一直十分重视高等教育质量工作。2007年1月,教育部下发了《关于实施高等学校本科教学质量与教学改革工程的意见》,计划实施“高等学校本科教学质量与教学改革工程(简称‘质量工程’)\”,通过专业结构调整、课程教材建设、实践教学改革、教学团队建设等多项内容,进一步深化高等学校教学改革,提高人才培养的能力和水平,更好地满足经济社会发展对高素质人才的需要。在贯彻和落实教育部“质量工程”的过程中,各地高校发挥师资力量强、办学经验丰富、教学资源充裕等优势,对其特色专业及特色课程(群)加以规划、整理和总结,更新教学内容、改革课程体系,建设了一大批内容新、体系新、方法新、手段新的特色课程。在此基础上,经教育部相关教学指导委员会专家的指导和建议,清华大学出版社在多个领域精选各高校的特色课程,分别规划出版系列教材,以配合“质量工程”的实施,满足各高校教学质量和教学改革的需要。

为了深入贯彻落实教育部《关于加强高等学校本科教学工作,提高教学质量的若干意见》精神,紧密配合教育部已经启动的“高等学校教学质量与教学改革工程精品课程建设工作”,在有关专家、教授的倡议和有关部门的大力支持下,我们组织并成立了“清华大学出版社教材编审委员会”(以下简称“编委会”),旨在配合教育部制定精品课程教材的出版规划,讨论并实施精品课程教材的编写与出版工作。“编委会”成员皆来自全国各类高等学校教学与科研第一线的骨干教师,其中许多教师为各校相关院、系主管教学的院长或系主任。

按照教育部的要求,“编委会”一致认为,精品课程的建设工作从开始就要坚持高标准、严要求,处于一个比较高的起点上;精品课程教材应该能够反映各高校教学改革与课程建设的需要,要有特色风格、有创新性(新体系、新内容、新手段、新思路,教材的内容体系有较高的科学创新、技术创新和理念创新的含量)、先进性(对原有的学科体系有实质性的改革和发展,顺应并符合21世纪教学发展的规律,代表并引领课程发展的趋势和方向)、示范性(教材所体现的课程体系具有较广泛的辐射性和示范性)和一定的前瞻性。教材由个人申报或各校推荐(通过所在高校的“编委会”成员推荐),经“编委会”认真评审,最后由清华大学出版

社审定出版。

目前,针对计算机类和电子信息类相关专业成立了两个“编委会”,即“清华大学出版社计算机教材编审委员会”和“清华大学出版社电子信息教材编审委员会”。推出的特色精品教材包括:

- (1) 21世纪高等学校规划教材·计算机应用——高等学校各类专业,特别是非计算机专业的计算机应用类教材。
- (2) 21世纪高等学校规划教材·计算机科学与技术——高等学校计算机相关专业的教材。
- (3) 21世纪高等学校规划教材·电子信息——高等学校电子信息相关专业的教材。
- (4) 21世纪高等学校规划教材·软件工程——高等学校软件工程相关专业的教材。
- (5) 21世纪高等学校规划教材·信息管理与信息系统。
- (6) 21世纪高等学校规划教材·财经管理与应用。
- (7) 21世纪高等学校规划教材·电子商务。
- (8) 21世纪高等学校规划教材·物联网。

清华大学出版社经过三十多年的努力,在教材尤其是计算机和电子信息类专业教材出版方面树立了权威品牌,为我国的高等教育事业做出了重要贡献。清华版教材形成了技术准确、内容严谨的独特风格,这种风格将延续并反映在特色精品教材的建设中。

清华大学出版社教材编审委员会

联系人:魏江江

E-mail:weijj@tup.tsinghua.edu.cn

前言

移动互联时代，数以百亿计的机器、企业、个人随时随地都会获取和产生新的数据。1分钟之内，新浪微博发送数万条微博，苹果应用商店下载次数以万计，淘宝卖出了几万件商品，百度产生了百万次搜索查询……所有这些行为都由海量的数据来呈现。

据CNNIC第35次调查报告显示：截止到2014年12月，我国的网民数量达到6.49亿，手机网民数达到5.57亿，互联网普及率为47.9%。这些网民每天产生海量的浏览数据、交易数据和原创数据，这些数据都存储在经营者的数据库里，占用了大量的存储资源。如何发现这些数据背后所隐藏的有价值的信息？如何利用这些数据来发现机会？

目前，很多高校在本科高年级和研究生阶段都开设了“数据挖掘”、“商务智能”、“客户关系管理”、“企业信息管理”、“物联网技术与应用”等课程，这些课程从不同侧面介绍了数据分析的技术、数学模型、行业应用等。随着大数据时代的到来，越来越多的高校已经意识到数据分析的重要性，培养学生，尤其是高年级本科生和研究生数据分析和应用的能力已经被列入了新的教学计划。

本书从“挖掘数据的潜在价值和应用”的视角出发，从互联网企业、政府、传统企业、个人等不同的社会角色对数据挖掘和分析的应用需求，系统地介绍了大数据时代数据分析的作用、技术和具体应用。

本书共分为8章：第1章 大数据与数据分析，第2章 互联网数据存储，第3章 互联网数据分析工具，第4章 商务网站数据分析与应用，第5章 政府网站数据分析及应用，第6章 物联网数据分析与应用，第7章 移动商务数据分析与应用，第8章 微博数据分析与应用，分别从不同侧面介绍了数据分析的相关内容。

在本书的写作过程中，参考了很多互联网上的相关资料，由于数量众多，不可能在参考文献里一一列出，在此一并向无私奉献自己的智慧和经验的同行们表达衷心的感谢！王晨、杨致远、邵庆月、张学伟为本书搜集了一手资料，对本书的体系架构提出了很好的建议，在此表示感谢！

赵守香
北京工商大学计算机与信息工程学院
2015年5月

目 录

第1章 大数据与数据分析	1
1.1 概述	1
1.1.1 大数据的含义	2
1.1.2 大数据的定义	2
1.1.3 大数据的特征	3
1.1.4 大数据与云计算	5
1.1.5 大数据与商业模式变革	5
1.1.6 大数据带来的问题	6
1.2 大数据与云计算	8
1.3 大数据与电子商务	14
1.3.1 电子商务催生大数据	15
1.3.2 数据分析给电子商务带来更多机会	15
1.3.3 网站分析与应用	17
1.4 大数据与物联网	18
1.4.1 物联网的含义	18
1.4.2 物联网与大数据的关系	19
1.4.3 美国物联网应用	19
1.5 移动互联网与智能终端	20
1.6 大数据应用的机会与挑战	23
1.6.1 挖出“潜伏”的数据价值	24
1.6.2 大数据面临的挑战	27
1.6.3 大数据思维	30
1.7 银行业大数据应用	36
第2章 互联网数据存储	37
2.1 大数据对数据存储的要求	37
2.1.1 数据存储面临的问题	38
2.1.2 与大数据存储基础设施相关的属性	39
2.1.3 数据存储技术面临的挑战	41
2.1.4 存储技术趋势预测与分析	42
2.2 存储技术	44
2.2.1 DAS 存储	46

2.2.2 RAID 存储	46
2.2.3 NAS	48
2.2.4 SAN	50
2.2.5 IP 网络存储	51
2.2.6 iSCSI	52
2.2.7 存储技术比较	54
2.3 云存储技术	56
2.3.1 云存储技术与传统存储技术	56
2.3.2 云存储的优点	57
2.3.3 云存储的分类	58
2.3.4 云存储的技术基础	60
2.3.5 云存储系统的结构模型	61
2.3.6 云存储的用途	62
2.4 大数据存储解决方案	63
2.4.1 戴尔的流动文件系统	64
2.4.2 华为的集群存储系统	65
2.4.3 戴尔的自动分层存储	66
2.4.4 EMC 的闪存存储技术	68
第3章 互联网数据分析工具	70
3.1 数据分析概述	70
3.1.1 数据分析过程	72
3.1.2 数据分析框架的主要事件	73
3.2 数据分析与数据挖掘	74
3.2.1 数据挖掘的任务	75
3.2.2 数据挖掘的过程	77
3.2.3 数据挖掘的主要算法	78
3.2.4 数据挖掘的应用领域	81
3.2.5 数据挖掘和 OLAP	83
3.3 关联分析	84
3.3.1 关联规则挖掘过程	84
3.3.2 关联规则分类	85
3.3.3 关联规则算法	86
3.3.4 关联规则应用	87
3.4 聚类分析	88
3.5 分类分析	93
3.5.1 决策树	93
3.5.2 其他分类算法	95
3.6 时间序列分析	99

3.6.1	时间序列的构成要素	100
3.6.2	时间序列的分类	101
3.6.3	预测方法	102
3.6.4	确定性时间序列分析	104
3.6.5	随机性时间序列分析	107
第4章	商务网站数据分析与应用	109
4.1	概述	111
4.1.1	商业活动与商业数据	111
4.1.2	电子商务数据的特点	112
4.1.3	商务数据的挖掘利用	113
4.2	网站数据分析	115
4.2.1	为什么需要数据分析	115
4.2.2	网站数据分析的内容	115
4.2.3	怎么做数据分析	115
4.3	网站数据分析的指标体系	117
4.3.1	相关术语介绍	117
4.3.2	网站数据分析的指标分类	122
4.3.3	数据分析的内容指标体系	126
4.3.4	网站分析的商业指标	130
4.3.5	网站数据分析的应用价值	132
4.4	网站流量数据的获取	133
4.4.1	监听网络数据包	133
4.4.2	分析服务器日志	134
4.4.3	添加页面脚本	134
4.4.4	三种方法的比较	134
4.5	网站数据分析技术	135
4.5.1	数据收集系统	137
4.5.2	数据转发系统	139
4.5.3	实时数据分析系统	140
4.5.4	离线数据平台系统	143
4.6	数据分析应用	145
4.6.1	网站优化	146
4.6.2	个性化推荐	146
4.6.3	网页设计优化	150
4.6.4	服务提升与优化	152
4.6.5	网络营销	153
4.7	案例分析	155

第 5 章 政府网站数据分析及应用	161
5.1 电子政务概述	161
5.1.1 电子政务的含义	161
5.1.2 电子政务价值	163
5.1.3 电子政务环境下的政府信息资源	164
5.1.4 政务网站信息分类	166
5.2 政务网站信息分类	169
5.3 政府数据仓库与数据挖掘	170
5.3.1 政务元数据标准	170
5.3.2 电子政务数据仓库	171
5.3.3 电子政务数据挖掘	173
5.4 电子政务网站数据分析方法	174
5.5 数据分析与政府执行力	179
5.5.1 服务满意度测评	179
5.5.2 主动服务	181
5.5.3 民生热点分析	184
5.5.4 服务流程优化	184
5.6 案例分析	189
第 6 章 物联网数据分析与应用	193
6.1 物联网概述	193
6.1.1 物联网的概念与实质	193
6.1.2 物联网的兴起与发展状况	194
6.1.3 物联网的应用	195
6.1.4 物联网在我国的应用现状	196
6.1.5 应用模式	197
6.2 物联网技术	197
6.2.1 条码技术	197
6.2.2 RFID 技术	200
6.2.3 全球数据同步	203
6.3 物联网数据分析与处理	206
6.3.1 物联网系统中数据的特点	206
6.3.2 物联网数据处理模型	207
6.3.3 物联网与大数据分析	209
6.4 物联网数据分析应用	212
6.4.1 智能家居	212
6.4.2 远程医疗	213
6.4.3 老人关怀	214

6.4.4 药品安全监控.....	214
6.4.5 零售、物流、供应链管理.....	214
6.4.6 食品追踪.....	215
6.4.7 农业育种.....	215
6.5 物联网数据挖掘	215
6.5.1 物联网数据挖掘的关键问题.....	216
6.5.2 物联网环境数据挖掘存在的挑战.....	216
6.5.3 基于云计算的物联网数据挖掘模型.....	216
6.5.4 功能模块.....	218
6.6 应用案例	218
第 7 章 移动商务数据分析与应用.....	221
7.1 移动商务概述	221
7.1.1 移动商务的概念及分类.....	221
7.1.2 移动商务的特点.....	222
7.1.3 移动电子商务的体系与产业链.....	224
7.2 移动商务的应用	226
7.3 移动商务数据分析技术	230
7.3.1 无线与移动技术.....	230
7.3.2 数据仓库技术.....	230
7.3.3 联机事务处理与联机分析处理.....	231
7.3.4 知识发现技术.....	231
7.3.5 信息聚合技术.....	231
7.3.6 智能技术.....	231
7.4 移动商务中的数据挖掘技术	232
7.4.1 数据挖掘基本流程.....	232
7.4.2 关联规则在移动商务客户价值挖掘中的应用案例.....	232
7.5 位置信息分析与应用	235
7.5.1 位置服务的含义.....	235
7.5.2 定位技术.....	238
7.5.3 基于位置服务的推荐算法.....	241
7.5.4 LBS 与物流优化	244
7.6 移动商务发展中的问题	248
7.6.1 技术应用阻力.....	248
7.6.2 商业模式仍需摸索	249
第 8 章 微博数据分析与应用.....	251
8.1 微博概述	252
8.1.1 微博的定义.....	253

8.1.2	微博的特点	254
8.1.3	中外微博的文化差异	255
8.1.4	微博代表	257
8.2	微博应用	259
8.3	微博数据分析技术	262
8.3.1	文本信息抽取技术	263
8.3.2	微博文本处理	263
8.3.3	微博舆情分析	264
8.3.4	基于语义分析的微博文本挖掘技术	265
8.3.5	用户影响力计算的相关算法	269
8.3.6	适于演化的微博信息的数据表达模型	274
8.3.7	适于微博信息的大规模数据集划分方法	274
8.4	企业微博数据分析及应用	276
8.4.1	微博营销：数据分析的应用	276
8.4.2	微应用	278
8.4.3	企业机构话题营销	278
8.4.4	微博营销数据分析案例	279
8.5	政务微博数据分析及应用	281
8.5.1	政务微博的特点	281
8.5.2	政务微博应用存在的问题	282
8.5.3	政务微博应用面临的挑战	282
8.5.4	政务微博的数据分析和应用	284
8.6	大众微博的舆情分析	288
8.6.1	舆情分析的内容	288
8.6.2	新浪微博对网络舆情生成和传播的影响	290
8.6.3	大众网络舆情的作用	291
8.7	微博营销案例	292
	参考文献	297

第1章

大数据与数据分析

【内容提要】

本章从大数据的出现、大数据的影响及大数据对数据处理的要求出发，分析了大数据环境下对数据分析与利用的重要性，主要内容包括：

- (1) 大数据产生的背景和特征；
- (2) 电子商务的快速发展与大数据；
- (3) 物联网的兴起与大数据；
- (4) 移动商务的快速渗透与大数据；
- (5) 数据分析给企业带来的机会与好处；
- (6) 大数据环境下数据分析的需求；
- (7) 数据分析的含义；
- (8) 数据分析的作用；
- (9) 大数据与数据分析。

1.1 概述

互联网、移动互联网、物联网、云计算的快速兴起以及移动智能终端的快速普及，使得当前人类社会的数据增长比以往任何一个时期都要快。数据的爆炸式增长正在出乎人们的想象。据预计，2020年，全球以电子形式存储的数据量将达35ZB，是2009年全球存储量的40倍。而在2010年年底，根据IDC的统计，全球数据量已经达到了1.22ZB。如果将这些数据的50%刻录在DVD上，那么把这些DVD盘片堆起来就可以从地球垒到月球一个来回。

与此同时，伴随着物联网、移动智能终端以及移动互联网的快速发展，移动网络中数据流量的增长速度也非常迅猛。从2011年开始，全球移动数据流量年增长率将保持在50%以上，并处于一个稳定增长的态势。到2016年，全球移动数据流量将达到2011年全球移动数据流量的18倍，达到129.6EB。

数据的疯狂增长，使得适应和应对数据增长成为整个社会关注的焦点。“大数据”的概念正是在这一背景下应运而生。

1.1.1 大数据的含义

我们先来看一组数据：1分钟之内，新浪微博发送数万条微博，苹果应用商店下载次数以万计，淘宝卖出了几万件商品，百度产生了百万次搜索查询……所有这些行为都由海量的数据来呈现。

在2013年12月12日电商的促销期，淘宝网推出“时光机”——一个根据淘宝买家几年来的购买商品记录、浏览点击次数、收货地址等数据编辑制作的“个人网购志”，从而记录和勾勒出让人感怀的生活记忆。背后是基于对4.7亿淘宝注册用户网购数据的分析处理，这正是大数据的典型应用。

随着传统互联网向移动互联发展，全球范围内，除了个人电脑、平板电脑、智能手机、游戏主机等常见的计算终端之外，更广阔的、泛在互联的智能设备，如智能汽车、智能电视、工业设备和手持设备等都连接到网络之中。基于社会化网络的平台和应用，让数以百亿计的机器、企业、个人随时随地都会获取和产生新的数据。

互联网搜索引擎是大数据最为典型的应用之一。百度日处理数据量达到数十拍字节，并呈现高速增长的态势。如果一张光盘容量为1GB，这相当于垒在一起的几千万张光盘。微软Bing(在中国为必应)搜索引擎，一周需要响应100亿次量级的搜索请求。通过和Facebook的合作，每天有超过10亿次的社交网络搜索请求通过Bing来处理。

短短的18个月，中国移动互联网流量增加了10倍。随着社交网络的逐渐成熟、移动带宽迅速提升，更多的传感设备、移动终端接入网络，产生的数据及其增长速度比历史上任何时期都要多，互联网上的数据流量正在迅猛增长。在云计算、物联网等技术的带动下，中国的移动互联网已经步入“大数据”时代了。

而根据市场调研公司IDC的报告，全球信息总量每过两年就会增长一倍，2011年全球产生的数据总量为1.8ZB(1ZB为百万拍字节)，相当于全球历史数据总和。

继云计算后，大数据(Big Data)成为信息技术领域最为热门的概念之一。

简单说，大数据是指那些超过传统数据库系统处理能力的数据。但是，大数据的问题并不仅仅是规模，数据产生的速度以及数据的多样性同样是大数据不可忽略的两个基本特性。根据摩尔定律，计算能力每一年半到两年的时间将增加一倍。可是，现有的网络带宽并没有以同样的速度在增加。因此，如此迅猛的数据洪流的产生，正在给电信运营商的网络运营带来极大的挑战。

在IT业界，有人把大数据产业定义为建立在对互联网、物联网等渠道广泛大量数据资源收集基础上的数据存储、价值提炼、智能处理和分发的信息服务业；也有人概括大数据战略为：致力于让所有用户能够从几乎任何数据中获得可转换为业务执行的洞察力，包括之前隐藏在非结构化数据中的洞察力。

总之，大数据是对大量、动态、能持续的数据，通过运用新系统、新工具、新模型的挖掘，从而获得具有洞察力和新价值的东西。

1.1.2 大数据的定义

根据维基百科的定义，大数据指难以用常用的软件工具在可容忍的时间内抓取、管理以

及处理的数据集。大数据的显著特征有：数据体量巨大；数据类型繁多，包括结构化数据以及非结构化数据如网页、日志、视频、图片等；要求的处理速度快。

大数据技术涵盖了从数据的海量存储、处理到应用多方面的技术，包括海量分布式文件系统、并行计算框架、NoSQL 数据库、实时流数据处理以及智能分析技术如模式识别、自然语言理解、应用知识库等。

1.1.3 大数据的特征

虽然有多种解读，但业界一般认为，大数据有 4 个 V 字开头的特征：Volume(容量)、Variety(种类)、Velocity(速度)和最重要的 Value(价值)。

1. Volume

是指大数据巨大的数据量与数据完整性。IT 业界所指的数据，诞生不过六十多年。而一直到个人电脑普及前，由于存储、计算和分析工具的技术和成本限制，许多自然界和人类社会值得记录的信号，并未形成数据。几十年前，气象、地质、石油勘探、出版业、媒体业和影视业是大量、持续产出信号的行业，但那时 90%以上采用的是存储模拟信号，难以通过计算设备和软件进行直接分析。拥有大量资金和人才的政府和企业，也只能把少量最关键的信号，进行抽取、转换、装载到数据库中。

尽管业界对达到怎样的数量级才算是大数据并无定论，但在很多行业的应用场景里，数据集本身的大小并不是最重要的，是否完整才最重要。

2. Variety

Variety 则意味着要在海量、种类繁多的数据间发现其内在关联。互联网时代，各种设备通过网络连成了一个整体。进入以互动为特征的 Web 2.0 时代，个人计算机用户不仅可以通过网络获取信息，还成为信息的制造者和传播者。这个阶段，不仅数据量开始了爆炸式增长，数据种类也开始变得繁多。

这必然促使我们对海量数据进行分析、处理和集成，找出原本看来毫无关系的那些数据的“关联性”，把似乎没有用的数据变成有用的信息，以支持我们做出的判断。

3. Velocity

Velocity 可以理解为更快地满足实时性需求。数据的实时化需求正越来越清晰。对普通人而言，开车去吃饭，会先用移动终端中的地图查询餐厅的位置，预计行车路线的拥堵情况，了解停车场信息甚至是其他用户对餐厅的评论。吃饭时，会用手机拍摄食物的照片，编辑简短评论发布到微博或者微信上，还可以用 LBS(基于位置的服务)应用查找在同一间餐厅吃饭的人，看有没有好友在附近……

如今，通过各种有线和无线网络，人和人、人和各种机器、机器和机器之间产生无处不在的连接，这些连接不可避免地带来了数据交换。而数据交换的关键是降低延迟，以近乎实时——这意味着小于 250ms——的方式呈献给用户。

4. Value

但比前面三个 V 更重要的,就是 Value,它是大数据的最终意义——获得洞察力和价值。大数据的崛起,正是在人工智能、机器学习和数据挖掘等技术的迅速发展驱动下,呈现了这么一个过程:将信号转化为数据,将数据分析为信息,将信息提炼为知识,以知识促成决策和行动。

就大数据的价值而言,就像沙子淘金,大数据规模越大,真正有价值的数据相对越少。

所以真正好的大数据系统,重要的不是越多越好,其实越少越好。开始数据要多,最后还是要少,把 ZB、PB 最终变成一个比特,也就是最后的决策,这才是最关键的。

正如我们常说的:书刚开始越读越厚,到最后就越读越薄了。

数据挖掘和应用可以多方位创造价值,如表 1-1 和图 1-1 所示。

表 1-1 数据有效性的 5 个方面

数据有效性	内 容
数据质量	准确性: 数据中没有错误的程度 范围: 数据覆盖的深度和广度 合时性: 及时获取数据以采取行动和决策的程度 有效: 相关数据及时更新程度
可用性	数据简明简要呈现的程度 数据易操作和易处理 多种数据来源中数据的一致性
智能性	数据预测准确性 相关数据的趋势性 需求模式分析 精确营销推荐性
远程访问性	数据可被管理者远程获得的程度 管理者可远程应用的程度
支持移动销售	利用便携系统与消费者交换价格信息 利用便携系统与消费者交换订单信息 利用便携系统与消费者交换运输信息

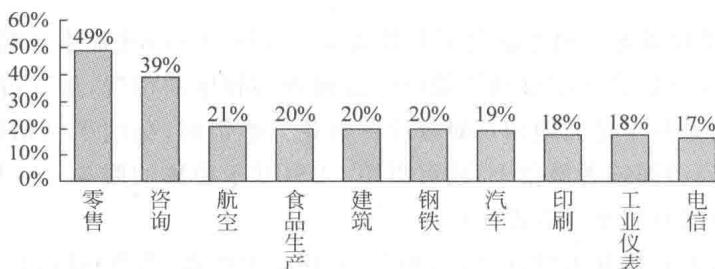


图 1-1 数据可用性提高 10%,各行业员工销售额提高百分比统计

1.1.4 大数据与云计算

云计算和大数据是一个硬币的两面，大数据正在引发全球范围内深刻的技术和商业变革。

云计算是大数据的 IT 基础,而大数据是云计算的一个杀手级应用。云计算是大数据成长的驱动力,而另一方面,由于数据越来越多、越来越复杂、越来越实时,这就更加需要云计算去处理,所以二者之间是相辅相成的。

30 年前,存储 1TB(1000GB)数据的成本大约是 16 亿美元,如今存储到云上只需不到 100 美元;但存储下来的数据,如果不以云计算进行挖掘和分析,就只是僵死的数据,没有太大的价值。

目前,云计算已经普及并成为IT行业的主流技术,其实质是在计算量越来越大、数据越来越多、越来越动态、越来越实时的需求背景下被催生出来的一种基础架构和商业模式。个人用户将文档、照片、视频、游戏存档记录上传至“云”中永久保存,企业客户根据自身需求,可以搭建自己的“私有云”,或托管或租用“公有云”上的IT资源与服务,这些都已不是新鲜事。可以说,云是一棵挂满了大数据的苹果树。

表 1-2 说明了大数据对数据全周期的要求都大大提高了。数据周期为相应产业带来的发展契机见图 1-2。

表 1-2 数据资产化对数据全周期提出更高要求

数据周期	供应商	要 求
数据存储	存储器供应商	需要更大、更快、更准地存储各项数据,同时对所存储数据有 安全性,具有良好的灾备功能
数据处理	服务器供应商	需要更快地计算处理海量数据
数据管理	操作系统与数据库软件商	需要高效地查询和操作海量数据
数据应用	商业智能软件商	对海量数据进行挖掘、分析和优化
数据安全	信息安全提供商	对数据提供安全保护

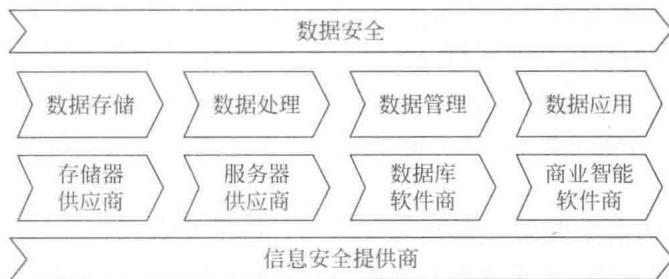


图 1-2 数据周期为相应行业带来发展契机

1.1.5 大数据与商业模式变革

大数据的出现,正在引发全球范围内深刻的技术与商业变革。在技术上,大数据使从数