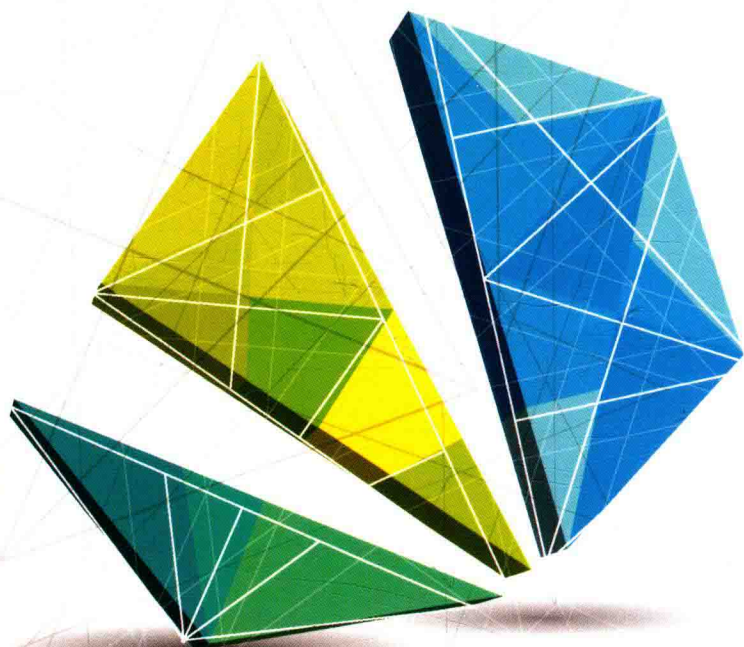


高等学校大数据技术与应用规划教材

# 大数据 及其可视化

DASHUJU JIQI KESHIHUA

周苏 王文 等编著



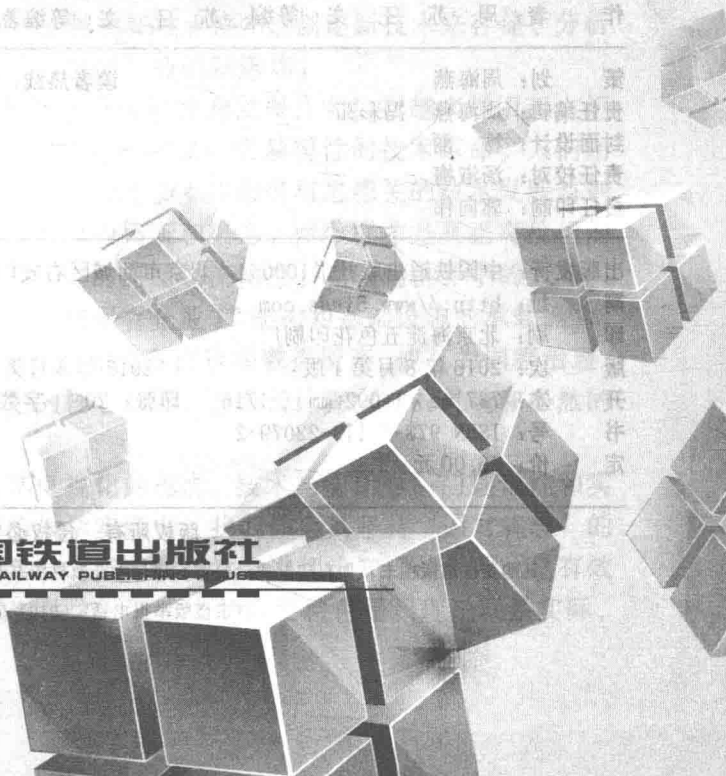
中国铁道出版社  
CHINA RAILWAY PUBLISHING HOUSE

高等学校大数据技术与应用规划教材

# 大数据及其可视化

周 苏 王 文 等编著

中国铁道出版社  
CHINA RAILWAY PUBLISHING HOUSE



## 内 容 简 介

大数据及其可视化是一门理论性和实践性都很强的课程。本书针对计算机、信息管理、经济管理和其他相关专业学生的发展需求,系统、全面地介绍了关于大数据及其可视化技术的基本知识和技能,详细介绍了大数据与大数据时代、数据可视化之美、Excel 数据可视化方法、Excel 数据可视化应用、大数据的商业规则、大数据激发创造力、大数据预测分析、支撑大数据的技术、数据引导可视化、Tableau 可视化初步、Tableau 数据管理与计算、Tableau 可视化设计、Tableau 地图与预测分析和 Tableau 分享与发布等内容,具有较强的系统性、可读性和实用性。

本书适合作为普通高等院校相关专业“大数据基础”“大数据导论”“大数据可视化”等课程的教材,也可供有一定实践经验的软件开发人员、管理人员学习参考。

### 图书在版编目(CIP)数据

大数据及其可视化/周苏,王文等编著. —北京:中国铁道出版社,2016.8  
高等学校大数据技术与应用规划教材  
ISBN 978-7-113-22079-2

I. ①大… II. ①周… ②王… 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(2016)第177191号

书 名: 大数据及其可视化  
作 者: 周 苏 王 文 等编著

---

策 划: 周海燕  
责任编辑: 周海燕 冯彩茹  
封面设计: 穆 丽  
责任校对: 汤淑梅  
责任印制: 郭向伟

读者热线: (010) 63550836

---

出版发行: 中国铁道出版社(100054,北京市西城区右安门西街8号)  
网 址: <http://www.51eds.com>  
印 刷: 北京海淀五色花印刷厂  
版 次: 2016年8月第1版 2016年8月第1次印刷  
开 本: 787 mm×1 092 mm 1/16 印张: 20 字数: 448 千  
书 号: ISBN 978-7-113-22079-2  
定 价: 46.00 元

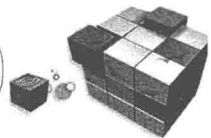
---

版权所有 侵权必究

凡购买铁道版图书,如有印制质量问题,请与本社教材图书营销部联系调换。电话:(010) 63550836

打击盗版举报电话:(010) 51873659

# 前 言



大数据 (Big Data) 的力量正在积极地影响着社会的方方面面, 它冲击着许多主要的行业, 包括零售业、电子商务和金融服务业等, 同时, 也正在彻底地改变人们的教育方式、生活方式、工作方式。如今, 通过简单、易用的移动应用和基于云端的数据服务, 人们能够追踪自己的行为以及饮食习惯, 还能提升个人的健康状况。因此, 有必要真正理解大数据这个极其重要的议题。

中国是大数据最大的潜在市场之一。据估计, 中国有近 6 亿网民, 这就意味着中国的企业拥有绝佳的机会来更好地了解其客户并提供更个性化的体验, 同时, 为企业增加收入并提高利润。阿里巴巴就是一个很好的例子, 其不但在商业模式上具有颠覆性, 而且还掌握了与购买行为、产品需求和库存供应相关的海量数据。除了阿里巴巴高层的领导能力之外, 大数据必然是其成功的一个关键因素。

然而, 仅有数据是不够的。对于身处大数据时代的企业而言, 成功的关键还在于找出大数据所隐含的真知灼见。“以前, 人们总说信息就是力量, 但如今, 对数据进行分析、利用和挖掘才是力量之所在。”

很多年前, 人们就开始对数据进行利用。例如, 航空公司利用数据为机票定价, 银行利用数据搞清楚贷款对象, 信用卡公司则利用数据侦破信用卡诈骗等。但直到最近, 数据才真正成为人们日常生活的一部分。随着谷歌 (Google) 以及 QQ、微信、淘宝等的出现, 大数据游戏被永远改变了。你和我, 或者任何一个享受这些服务的用户都生成了一条数据足迹, 它能够反映出人们的行为。每次进行搜索时, 如查找某个人或者访问某个网站, 都加深了这条足迹。互联网企业开始创建新技术来存储、分析激增的数据——结果就迎来了被称为“大数据”的创新爆炸。

进入 2012 年以来, 由于互联网和信息行业的快速发展, 大数据越来越引起人们的关注, 已经引发云计算、互联网之后 IT 行业的又一大颠覆性的技术革命。人们用大数据来描述和定义信息爆炸时代产生的海量数据, 并命名与之相关的技术发展与创新。云计算主要为数据资产提供保管、访问的场所和渠道, 而数据才是真正有价值的资产。企业内部的经营信息、互联网世界中的商品物流信息, 互联网世界中的人与人交互信息、位置信息等, 其数量将远远超越现有企业 IT 架构和基础设施的承载能力, 实时性要求也将大大超越现有的计算能力。如何盘活这些数据资产, 使其为国家治理、企业决策乃至个人生活服务, 是大数据的核心议题, 也是云计算内在的灵魂和必然的升级方向。

对于在校大学生来说, 大数据及其可视化的理念、技术与应用是一门理论性和实践性都很强的“必修”课程。在长期的教学实践中, 我们体会到坚持“因材施教”的重要原则, 把实践环节与理论教学相融合, 抓实践教学促进理论知识的学习, 是有效改善教学效果和提高教学水平的重要方法之一。本书的主要特色是: 理论联系实际,

结合一系列了解和熟悉大数据理念、技术与应用的学习和实践活动，把大数据及其可视化的相关概念、基础知识和技术技巧融入实践中，使学生保持浓厚的学习热情，加深对大数据技术的兴趣、认识、理解和掌握。

本书系统、全面地介绍了大数据及其可视化的基本知识和应用技能，详细介绍了大数据与大数据时代、数据可视化之美、Excel 数据可视化方法、Excel 数据可视化应用、大数据的商业规则、大数据激发创造力、大数据预测分析、支撑大数据的技术、数据引导可视化、Tableau 可视化初步、Tableau 数据管理与计算、Tableau 可视化设计、Tableau 地图与预测分析，以及 Tableau 分享与发布等内容，具有较强的系统性、可读性和实用性。

本课程的教学评测可以从这样几个方面入手，即：

- (1) 每章课前【案例导读】（14次）。
- (2) 每章课后【实验与思考】（14次）。
- (3) 课程设计（附录）。
- (4) 课程实验总结（附录）。
- (5) 结合平时考勤。
- (6) 任课老师认为必要的其他考核方法。

与本书配套的教学 PPT 课件等文档可从中国铁道出版社教学资源网站（[www.tdpress.com/51eds](http://www.tdpress.com/51eds)）的下载区下载，欢迎教师与作者交流并索取为本书教学配套的相关资料并交流：zhousu@qq.com，QQ：81505050，个人博客：<http://blog.sina.com.cn/zhousu58>。

本书由周苏、王文等编著，并得到浙江大学城市学院、浙江商业职业技术学院、温州安防职业技术学院等多所院校师生的支持，王硕莘、张丽娜、张健、吴林华等参与了本书的部分编写工作，在此一并表示感谢！

由于编者水平有限，加之时间仓促，书中难免存在疏漏和不足之处，恳请读者批评指正。

周 苏

2016年初夏于西子湖畔

# 目 录



## 第 1 章 大数据与大数据时代 ..... 1

### 1.1 大数据概述 ..... 3

#### 1.1.1 数据与信息 ..... 3

#### 1.1.2 天文学——信息爆炸的起源 ..... 3

#### 1.1.3 大数据的定义 ..... 4

#### 1.1.4 用 3V 描述大数据特征 ..... 5

#### 1.1.5 大数据的结构类型 ..... 7

### 1.2 思维变革之一：样本=总体 ..... 8

#### 1.2.1 小数据时代的随机采样 ..... 8

#### 1.2.2 大数据与乔布斯的癌症治疗 ..... 11

#### 1.2.3 全数据模式：样本=总体 ..... 12

### 1.3 思维变革之二：接受数据的混杂性 ..... 12

#### 1.3.1 允许不精确 ..... 12

#### 1.3.2 大数据的简单算法与小数据的复杂算法 ..... 13

#### 1.3.3 纷繁的数据越多越好 ..... 14

#### 1.3.4 5% 的数字数据与 95% 的非结构化数据 ..... 15

### 1.4 思维变革之三：数据的相关关系 ..... 16

#### 1.4.1 关联物，预测的关键 ..... 16

#### 1.4.2 “是什么”，而不是“为什么” ..... 17

#### 1.4.3 通过相关关系了解世界 ..... 18

#### 【实验与思考】深入理解大数据时代 ..... 19

## 第 2 章 数据可视化之美 ..... 21

### 2.1 数据与可视化 ..... 22

#### 2.1.1 数据的可变性 ..... 23

#### 2.1.2 数据的不确定性 ..... 24

#### 2.1.3 数据的背景信息 ..... 25

#### 2.1.4 打造最好的可视化效果 ..... 26

### 2.2 数据与图形 ..... 26

#### 2.2.1 地图传递信息 ..... 26

#### 2.2.2 数据与走势 ..... 27

#### 2.2.3 视觉信息的科学解释 ..... 28

#### 2.2.4 图片和分享的力量 ..... 29

#### 2.2.5 公共数据集 ..... 29

### 2.3 实时可视化 ..... 31

### 2.4 可视化分析工具 ..... 31

#### 2.4.1 Microsoft Excel ..... 32

#### 2.4.2 Google Spreadsheets ..... 32

#### 2.4.3 Tableau ..... 33

#### 2.4.4 可视化编程工具 ..... 33

#### 【实验与思考】熟悉大数据可视化 ..... 35

## 第 3 章 Excel 数据可视化方法 ..... 37

### 3.1 Excel 的函数与图表 ..... 39

#### 3.1.1 Excel 函数 ..... 40

#### 3.1.2 Excel 图表 ..... 41



3.1.3 选择图表类型.....	43	4.3 圆饼图：部分占总体的比例....	68
3.2 整理数据源.....	44	4.3.1 重视圆饼图扇区的 位置排序.....	68
3.2.1 数据提炼.....	44	4.3.2 分离圆饼图扇区 强调特殊数据.....	69
3.2.2 抽样产生随机数据.....	47	4.3.3 用半个圆饼图刻画 半期内的数据.....	70
3.3 数理统计中的常见统计量.....	49	4.3.4 让多个圆饼图对象 重叠展示对比关系.....	71
3.3.1 比平均值更稳定的 中位数和众数.....	49	4.4 散点图：表示分布状态.....	72
3.3.2 正态分布和偏态 分布.....	50	4.4.1 用平滑线联系散点图 增强图形效果.....	72
3.3.3 财务预算中的分析 工具.....	52	4.4.2 将直角坐标改为象限 坐标凸显分布效果.....	73
3.4 改变数据形式引起的 图表变化.....	53	4.5 侧重点不同的特殊图表.....	74
3.4.1 用负数突出数据的 增长情况.....	53	4.5.1 用子弹图显示数据 的优劣.....	74
3.4.2 重排关键字顺序 使图表更合适.....	54	4.5.2 用温度计展示工作 进度.....	75
【实验与思考】体验 Excel 数据 可视化方法.....	55	4.5.3 用漏斗图进行业务 流程的差异分析.....	76
<b>第 4 章 Excel 数据可视化应用.....</b>	<b>57</b>	【实验与思考】大数据如何激发 创造力.....	78
4.1 直方图：对比关系.....	60	<b>第 5 章 大数据的商业规则.....</b>	<b>79</b>
4.1.1 以零基线为起点.....	60	5.1 大数据的跨界年度.....	80
4.1.2 垂直直条的宽度要 大于条间距.....	62	5.2 谷歌的大数据行动.....	81
4.1.3 慎用三维效果的 柱形图.....	63	5.3 亚马逊的大数据行动.....	83
4.1.4 用堆积图表示 百分数.....	64	5.4 将信息变成一种竞争优势.....	84
4.2 折线图：按时间或类别 显示趋势.....	65	5.4.1 数据价格下降， 数据需求上升.....	85
4.2.1 减小 Y 轴刻度单位 增强数据波动情况.....	65	5.4.2 大数据应用程序的 兴起.....	86
4.2.2 突出显示折线图中 的数据点.....	66	5.4.3 实时响应，大数据 用户的新要求.....	87
4.2.3 通过面积图显示 数据总额.....	67	5.4.4 企业构建大数据 战略.....	87

5.5 大数据营销..... 88	6.6.3 以人为本的汽车 设计理念.....111
5.5.1 像媒体公司一样 思考..... 88	6.6.4 寻找最佳音响效果..... 112
5.5.2 营销面对新的机遇 与挑战..... 89	6.6.5 建筑数据取代直觉..... 113
5.5.3 自动化营销..... 90	【实验与思考】大数据如何激发 创造力..... 114
5.5.4 为营销创建高容量 和高价值的内容..... 91	<b>第7章 大数据预测分析..... 116</b>
5.5.5 内容营销..... 91	7.1 预测分析..... 119
5.5.6 内容创作与众包..... 92	7.2 数据情感和情感数据..... 122
5.5.7 用投资回报率评价 营销效果..... 93	7.2.1 从博客观察集体 情感..... 122
【实验与思考】大数据营销的 优势与核心内涵..... 93	7.2.2 预测分析博客中的 情绪..... 122
<b>第6章 大数据激发创造力..... 95</b>	7.2.3 影响情绪的重要 因素——金钱..... 124
6.1 大数据与循证医学..... 97	7.3 数据具有内在预测性..... 125
6.2 大数据带来的医疗新突破..... 98	7.4 情感的因果关系..... 126
6.2.1 量化自我, 关注个人 健康..... 99	7.4.1 焦虑指数与标普 500 指数..... 126
6.2.2 可穿戴的个人健康 设备..... 100	7.4.2 验证情感和被验证 的情感..... 128
6.2.3 大数据时代的医疗 信息..... 101	7.4.3 情绪指标影响 金融市场..... 129
6.3 医疗信息数字化..... 103	【实验与思考】大数据准备度 自我评分表..... 130
6.4 搜索: 超级大数据的 最佳伙伴..... 105	<b>第8章 支撑大数据的技术..... 134</b>
6.5 数据决策的成功崛起..... 106	8.1 大数据在云端..... 135
6.5.1 数据辅助诊断..... 107	8.1.1 云计算概述..... 136
6.5.2 你考虑过……了吗..... 107	8.1.2 云计算的服务形式..... 137
6.5.3 大数据分析使数据 决策崛起..... 108	8.1.3 云计算与大数据..... 137
6.6 大数据帮助改善设计..... 109	8.1.4 云基础设施..... 139
6.6.1 少而精是设计的 核心..... 110	8.2 计算虚拟化..... 140
6.6.2 与玩家共同设计 游戏..... 111	8.3 网络虚拟化..... 140
	8.4 大数据存储..... 141



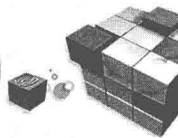


8.4.1 传统存储系统 .....	141	9.3.2 子分类 .....	170
8.4.2 大数据时代的 新挑战 .....	142	9.3.3 数据的结构和模式 .....	171
8.4.3 分布式存储 .....	143	9.4 时序数据的可视化 .....	171
8.4.4 云存储 .....	144	9.4.1 周期 .....	172
8.5 开源技术的商业支援 .....	145	9.4.2 循环 .....	173
8.6 大数据的技术架构 .....	146	9.5 空间数据的可视化 .....	174
8.7 Hadoop 基础 .....	147	9.6 让可视化设计更清晰 .....	175
8.7.1 分布式系统概述 .....	147	9.6.1 建立视觉层次 .....	175
8.7.2 Hadoop 的由来 .....	148	9.6.2 增强图表的可读性 .....	176
8.7.3 Hadoop 的优势 .....	149	9.6.3 允许数据点之间 进行比较 .....	177
8.7.4 Hadoop 的发行 版本 .....	150	9.6.4 描述背景信息 .....	178
8.8 大数据数据处理基础 .....	150	【实验与思考】绘制泰坦尼克 事件镶嵌图 .....	180
8.8.1 Hadoop 与 NoSQL .....	151	<b>第 10 章 Tableau 可视化初步 .....</b>	<b>183</b>
8.8.2 NoSQL 与 RDBMS 的主要区别 .....	151	10.1 Tableau 概述 .....	184
8.8.3 NewSQL .....	153	10.1.1 Tableau 可视化 技术 .....	185
【实验与思考】了解大数据的 基础设施 .....	154	10.1.2 Tableau 主要特性 .....	186
<b>第 9 章 数据引导可视化 .....</b>	<b>156</b>	10.2 Tableau 产品线 .....	187
9.1 可视化对认知的帮助 .....	157	10.2.1 Tableau Desktop .....	187
9.1.1 七个基本任务 .....	157	10.2.2 Tableau Server .....	187
9.1.2 新的数据研究方法 .....	158	10.2.3 Tableau Online .....	188
9.1.3 信息图形和展示 .....	159	10.2.4 Tableau Mobile .....	188
9.1.4 走进数据艺术的 世界 .....	160	10.2.5 Tableau Public .....	188
9.2 可视化设计组件 .....	161	10.2.6 Tableau Reader .....	188
9.2.1 视觉隐喻 .....	162	10.3 下载与安装 .....	189
9.2.2 坐标系 .....	165	10.4 Tableau 工作区 .....	190
9.2.3 标尺 .....	167	10.4.1 工作表工作区 .....	190
9.2.4 背景信息 .....	168	10.4.2 仪表盘工作区 .....	192
9.2.5 整合可视化组件 .....	168	10.4.3 故事工作区 .....	193
9.3 分类数据的可视化 .....	169	10.4.4 菜单栏和工具栏 .....	194
9.3.1 整体中的部分 .....	170	10.5 Tableau 数据 .....	195
		10.5.1 数据角色 .....	196
		10.5.2 字段类型 .....	198

10.5.3 文件类型 .....	198	11.6.2 公式的自动完成 .....	231
10.6 创建视图 .....	199	11.6.3 临时计算 .....	232
10.6.1 行列功能区 .....	199	11.6.4 创建计算成员 .....	232
10.6.2 标记卡 .....	201	11.6.5 聚合计算 .....	233
10.6.3 筛选器 .....	204	11.7 表计算 .....	235
10.6.4 页面 .....	204	11.8 百分比 .....	237
10.6.5 智能显示 .....	205	【实验与思考】熟悉 Tableau	
10.6.6 度量名称和 度量值 .....	205	数据可视化设计 .....	238
10.7 创建仪表板 .....	206	<b>第 12 章 Tableau 可视化设计 .....</b>	<b>242</b>
【实验与思考】熟悉 Tableau 数据 可视化设计 .....	208	12.1 条形图与直方图 .....	243
<b>第 11 章 Tableau 数据管理与计算 ...</b>	<b>211</b>	12.1.1 条形图 .....	243
11.1 Tableau 数据架构 .....	212	12.1.2 直方图 .....	245
11.2 数据连接 .....	212	12.2 饼图 .....	247
11.2.1 连接数据源 .....	212	12.3 折线图 .....	248
11.2.2 组织数据 .....	214	12.4 压力图与突显表 .....	250
11.2.3 实现多表联结 .....	215	12.4.1 压力图 .....	250
11.3 数据加载 .....	215	12.4.2 突显表 .....	251
11.3.1 创建数据提取 .....	216	12.5 树地图 .....	253
11.3.2 刷新数据提取 .....	217	12.6 气泡图与圆视图 .....	253
11.3.3 向数据提取 添加行 .....	217	12.6.1 气泡图 .....	253
11.3.4 优化数据提取 .....	217	12.6.2 圆视图 .....	254
11.4 数据维护 .....	218	12.7 标靶图 .....	255
11.5 高级数据操作 .....	219	12.8 甘特图 .....	257
11.5.1 分层结构 .....	219	12.9 盒须图 .....	257
11.5.2 组 .....	220	12.9.1 创建盒须图 .....	258
11.5.3 集 .....	221	12.9.2 图形延伸 .....	259
11.5.4 参数 .....	223	【实验与思考】熟悉 Tableau	
11.5.5 参考线及参考 区间 .....	226	数据可视化分析 .....	260
11.6 计算字段 .....	229	<b>第 13 章 Tableau 地图与预测分析 ...</b>	<b>262</b>
11.6.1 创建和编辑计算 字段 .....	229	13.1 Tableau 地图分析 .....	263
		13.1.1 分配地理角色 .....	263
		13.1.2 创建符号地图 .....	264
		13.1.3 创建填充地图 .....	266
		13.1.4 创建多维度地图 .....	266



13.1.5	创建混合地图 .....	267	14.1.5	仪表盘 Web 视图 安全选项 .....	288
13.1.6	设置地理信息 .....	267	14.2	布局容器 .....	288
13.2	Tableau 预测分析 .....	268	14.3	组织仪表盘 .....	289
13.2.1	指数平滑和趋势 .....	268	14.3.1	平铺和浮动布局 .....	289
13.2.2	季节性 .....	268	14.3.2	显示和隐藏工作表 的组成部分 .....	290
13.2.3	模型类型 .....	269	14.3.3	重新排列仪表盘 视图和对象 .....	290
13.2.4	使用时间进行 预测 .....	269	14.3.4	设置仪表盘大小 .....	291
13.2.5	粒度和修剪 .....	270	14.3.5	了解仪表盘和 工作表 .....	291
13.2.6	获取更多数据 .....	270	14.4	Tableau 故事 .....	291
13.3	建立预测分析 .....	271	14.4.1	故事工作区 .....	292
13.3.1	创建预测 .....	271	14.4.2	创建故事 .....	294
13.3.2	预测字段结果 .....	273	14.4.3	设置故事的格式 .....	295
13.3.3	预测描述 .....	274	14.4.4	更新与演示故事 .....	296
13.4	合计 .....	274	14.5	Tableau 发布 .....	296
13.5	背景图像 .....	276	14.5.1	导出和发布 数据 .....	296
13.5.1	添加背景图像 .....	276	14.5.2	导出图像和 PDF 文件 .....	301
13.5.2	设置视图 .....	278	14.5.3	保存和发布 工作簿 .....	302
13.5.3	管理背景图像 .....	278	【实验与思考】熟悉 Tableau 分享与发布 .....	303	
13.6	趋势线 .....	279	附录 .....	305	
【实验与思考】熟悉 Tableau 预测分析 .....	280	参考文献 .....	310		
<b>第 14 章 Tableau 分享与发布 .....</b>	<b>281</b>				
14.1	Tableau 仪表盘 .....	282			
14.1.1	创建仪表盘 .....	282			
14.1.2	向仪表板中添加 视图 .....	283			
14.1.3	添加仪表盘对象 .....	286			
14.1.4	从仪表板中移除 视图和对象 .....	288			



### 【案例导读】亚马逊推荐系统

虽然亚马逊<sup>①</sup>的故事大多数人都耳熟能详，但只有少数人知道它早期的书评内容最初是由人工完成的。当时，亚马逊公司（见图 1-1）聘请了一个由 20 多名书评家和编辑组成的团队，他们写书评、推荐新书，挑选非常有特色的新书标题放在亚马逊的网页上。这个团队创立了“亚马逊的声音”版块，成为当时公司皇冠上的一颗宝石，是其竞争优势的重要来源。《华尔街日报》的一篇文章中热情地称他们为全美最有影响力的书评家，因为他们使得书籍销量猛增。

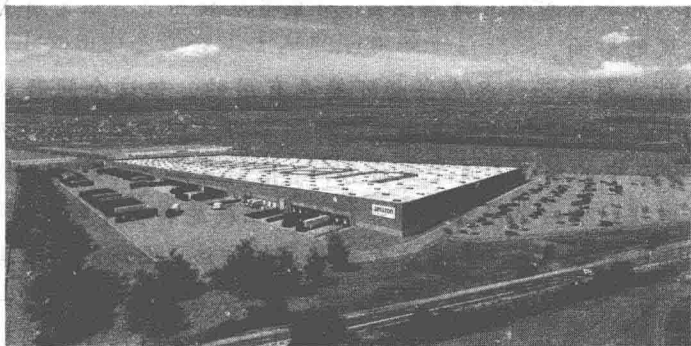


图 1-1 亚马逊公司

亚马逊公司的创始人及总裁杰夫·贝索斯决定尝试一个极富创造力的想法：根据客户个人以前的购物喜好，为其推荐相关的书籍。

从一开始，亚马逊就从每一个客户那里搜集了大量的数据。比如说，他们购买了什么书籍？哪些书他们只浏览却没有购买？他们浏览了多久？哪些书是他们一起购买的？客户的信息数据量非常大，所以亚马逊必须先用传统的方法对其进行处理，通过样本分析找到客户之间的相似性。但这些推荐信息是非常原始的，就如同你在买一件婴儿用品时，会被淹没在一堆差不多的婴儿用品中一样。詹姆斯·马库斯回忆说：“推荐信息往往为你提供与你以前购买物品有微小差异的产品，并且循环往复。”

亚马逊的格雷格·林登很快就找到了一个解决方案。他意识到，推荐系统实际上

<sup>①</sup> 亚马逊公司（Amazon）：是美国最大的网络电子商务公司之一，也是“财富 500 强”公司，位于华盛顿州的西雅图，成立于 1995 年 7 月，已成为全球商品种类最多的网上零售商。亚马逊致力于成为全球最“以客户为中心”的公司，使客户能在公司网站上找到和发现任何他们想在线购买的商品，并努力为客户提供最低的价格。



并没有必要把顾客与其他顾客进行对比，这样做在技术上也比较烦琐，需要做的是找到产品之间的关联性。1998年，林登和他的同事申请了著名的“item-to-item”协同过滤技术的专利。方法的转变使技术发生了翻天覆地的变化。

因为估算可以提前进行，所以推荐系统不仅快，而且适用于各种各样的产品。因此，当亚马逊跨界销售除书以外的其他商品时，也可以对电影或烤面包机这些产品进行推荐。由于系统中使用了所有的数据，推荐会更理想。林登回忆道：“在组里有句玩笑话，说的是如果系统运作良好，亚马逊应该只推荐你一本书，而这本书就是你将要买的下一本书。”

现在，公司必须决定什么应该出现在网站上，是亚马逊内部书评家写的个人建议和评论，还是由机器生成的个性化推荐和畅销书排行榜？

林登做了一个关于评论家所创造的销售业绩和计算机生成内容所产生的销售业绩的对比测试，结果他发现两者之间相差甚远。他解释说，通过数据推荐产品所增加的销售远远超过书评家的贡献。计算机可能不知道为什么喜欢海明威<sup>①</sup>作品的客户会购买菲茨杰拉德<sup>②</sup>的书。但是这似乎并不重要，重要的是销量。最后，编辑们看到了销售额分析，亚马逊也不得不放弃每次的在线评论，最终，书评组被解散。林登回忆说：“书评团队被打败、被解散，我感到非常难过。但是，数据没有说谎，人工评论的成本是非常高的。”

如今，据说亚马逊销售额的1/3都来自于它的个性化推荐系统。有了它，亚马逊不仅使很多大型书店和音乐唱片商店歇业，而且当地数百个自认为有自己风格的书商也难免受转型之风的影响。

知道人们为什么对这些信息感兴趣可能是有用的，但这个问题目前并不是很重要，而知道“是什么”可以创造点击率，这种洞察力足以重塑很多行业，不仅仅只是电子商务。所有行业中的销售人员早就被告知，他们需要了解是什么让客户做出了选择，要把握客户做决定背后的真正原因，因此专业技能和多年的经验受到高度重视。大数据却显示，还有另外一个在某些方面更有用的方法。亚马逊的推荐系统梳理出了有趣的相关关系，但不知道背后的原因——知道是什么就够了，没必要知道为什么。

(本案例由作者根据相关资料改写)

阅读上文，请思考、分析并简单记录：

(1) 你了解亚马逊等电商网站的推荐系统吗？请列举一个这样的实例（你选择购买什么商品，网站又给你推荐了其他什么商品）。

答：\_\_\_\_\_

(2) 亚马逊书评组和林登推荐系统各自成功的基础是什么？

答：\_\_\_\_\_

① 欧内斯特·米勒·海明威（1899年7月21日—1961年7月2日），美国小说家。被誉为美利坚民族的精神丰碑。出生于美国伊利诺伊州芝加哥市郊区的阿克帕克，代表作有《老人与海》《太阳照常升起》《永别了，武器》《丧钟为谁而鸣》等，凭借《老人与海》获得1953年普利策奖及1954年诺贝尔文学奖。

② 菲茨杰拉德，美国小说家。1920年出版了长篇小说《人间天堂》，从此出名。1925年《了不起的盖茨比》问世，奠定了他在现代美国文学史上的地位，成了20世纪20年代“爵士时代”的发言人和“迷惘的一代”的代表作家之一。



(3) 为什么书评组最终输给了推荐系统? 请阐述你的观点。

答: \_\_\_\_\_

(4) 简单描述你所知道的上一周内发生的国际、国内或者身边的大事。

答: \_\_\_\_\_



## 1.1 大数据概述

信息社会所带来的好处是显而易见的: 每个人口袋里都揣有一部手机, 每台办公桌上都放着一台计算机, 每间办公室内都连接到局域网甚至互联网。半个世纪以来, 随着计算机技术全面和深度地融入社会生活, 信息爆炸已经积累到了一个开始引发变革的程度。信息总量的变化导致了信息形态的变化——量变引起质变。最先经历信息爆炸的学科, 如天文学和基因学, 创造出了“大数据”(Big Data)这个概念。如今, 这个概念几乎应用到所有人类致力于发展的领域中。

### 1.1.1 数据与信息

数据是反映客观事物属性的记录, 是信息的具体表现形式。数据经过加工处理之后, 就成为信息; 而信息需要经过数字化转变成数据才能存储和传输。所以, 数据和信息之间是相互联系的。

数据和信息也是有区别的。从信息论的观点来看, 描述信源的数据是信息和数据冗余之和, 即数据=信息+数据冗余。数据是数据采集时提供的, 信息是从采集的数据中获取的有用信息, 即信息可以简单地理解为数据中包含的有用的内容。

一个消息越不可预测, 它所含的信息量就越大。事实上, 信息的基本作用是消除人们对事物了解的不确定性。信息量是指从  $N$  个相等的可能事件中选出一个事件所需要的信息度和含量。从这个定义看, 信息量与概率是密切相关的。

### 1.1.2 天文学——信息爆炸的起源

综合观察社会各个方面的变化趋势, 我们能真正意识到信息爆炸或者说大数据的时代已经到来。以天文学为例, 2000 年斯隆数字巡天<sup>①</sup>项目(见图 1-2)启动时, 位于新墨西哥州的望远镜在短短几周内搜集到的数据, 就比世界天文学历史上总共搜集的数据还要多。截至 2010 年, 信息档案已经高达  $1.4 \times 2^{42}$  B。不过, 预计 2016 年底, 在智利投入使用的大型视场全景巡天望远镜在 5 天之内即可获得同样多的信息。

<sup>①</sup> 斯隆数字巡天: 位于新墨西哥州阿帕奇山顶天文台的 2.5 m 口径望远镜红移巡天项目。计划观测 25% 的天空, 获取超过一百万个天体的多色测光资料和光谱数据。2006 年, 斯隆数字巡天进入了名为 SDSS-II 的新阶段, 进一步探索银河系的结构和组成, 而斯隆超新星巡天计划搜寻 Ia 型超新星爆发, 以测量宇宙学尺度上的距离。



天文学领域发生的变化在社会各个领域都在发生。2003年，人类第一次破译人体基因密码时，辛苦工作了十年才完成了三十亿对碱基对的排序。大约十年之后，世界范围内的基因仪每15 min就可以完成同样的工作。在金融领域，美国股市每天的成交量高达70亿股，而其中2/3的交易都是由建立在数学模型和算法之上的计算机程序自动完成的，这些程序运用海量数据来预测利益和降低风险。

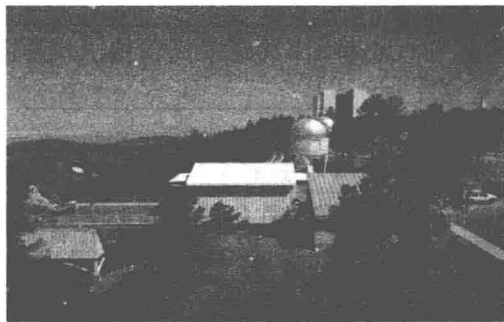


图 1-2 美国斯隆数字巡天望远镜

互联网公司更是要被数据淹没。谷歌公司每天要处理超过24拍字节(PB,  $2^{50}$ B)的数据，这意味着其每天的数据处理量是美国国家图书馆所有纸质出版物所含数据量的上千倍。

从科学研究到医疗保险，从银行业到互联网，各个不同的领域都在讲述着一个类似的故事，那就是爆发式增长的数据量。这种增长超过了人们创造机器的速度，甚至超过了人们的想象。人类存储信息量的增长速度比世界经济的增长速度快4倍，而计算机数据处理能力的增长速度则比世界经济的增长速度快9倍，每个人都受到了这种极速发展的冲击。

以纳米技术为例。纳米技术专注于把东西变小而不是变大。其原理就是当事物到达分子级别时，它的物理性质就会发生改变。一旦知道这些新的性质，就可以用同样的原料做以前无法做的事情。铜本来是用来导电的物质，但它一旦到达纳米级别就不能在磁场中导电了。银离子具有抗菌性，但当它以分子形式存在时，这种性质就会消失。一旦到达纳米级别，金属可以变得柔软，陶土可以具有弹性。同样，当人们增加所利用的数据量时，也就可以做很多在小数据量的基础上无法完成的事情。

大数据的科学价值和社会价值正是体现在这里。一方面，对大数据的掌握程度可以转化为经济价值的来源。另一方面，大数据已经撼动了世界的方方面面，从商业科技到医疗、政府、教育、经济、人文以及社会的其他各个领域。尽管人们还处在大数据时代的初期，但人们的日常生活已经离不开它。

### 1.1.3 大数据的定义

所谓大数据，狭义上可以定义为：用现有的一般技术难以管理的大量数据的集合。对大量数据进行分析，并从中获得有用观点，这种做法在一部分研究机构和大企业中早已存在。现在的大数据和过去相比，主要有3点区别：第一，随着社交媒体和传感器网络等的发展，人们身边正产生出大量且多样的数据；第二，随着硬件和软件技术的发展，数据的存储、处理成本大幅下降；第三，随着云计算的兴起，大数据的存储、处理环境已经没有必要自行搭建。

所谓“用现有的一般技术难以管理”，是指用目前在企业数据库占据主流地位的关系型数据库无法进行管理的、具有复杂结构的数据。或者也可以说，是指由于数据量的增大，导致对数据的查询(Query)响应时间超出允许范围的庞大数据。

研究机构Gartner给出了这样的定义：“大数据”是需要新处理模式才能具有更强的决策力、洞察发现力和流程优化能力的海量、高增长率和多样化的信息资产。

麦肯锡<sup>①</sup>说：“大数据指的是所涉及的数据集规模已经超过了传统数据库软件获取、存储、处理和分析的能力。这是一个被故意设计成主观性的定义，并且是一个关于多大的数据集才能被认为是大数据的可变定义，即并不定义大于一个特定数字的TB才叫大数据。因为随着技术的不断发展，符合大数据标准的数据集容量也会增长；并且定义随不同的行业也有变化，这依赖于在一个特定行业通常使用何种软件和数据集有多大。因此，大数据在今天不同行业中的范围可以从几十TB到几PB。”

随着“大数据”的出现，数据仓库、数据安全、数据分析、数据挖掘等围绕大数据商业价值的利用正逐渐成为行业人士争相追捧的利润焦点，在全球引领了新一轮数据技术革新的浪潮。

### 1.1.4 用3V描述大数据特征

从字面来看，“大数据”这个词可能会让人觉得只是容量非常大的数据集合而已。但容量只不过是大数据特征的一个方面，如果只拘泥于数据量，就无法深入理解当前围绕大数据所进行的讨论。因为“用现有的一般技术难以管理”这样的状况，并不仅仅是由于数据量增大这一个因素所造成的。

IBM说：“可以用3个特征相结合来定义大数据：数量（Volume，或称容量）、种类（Variety，或称多样性）和速度（Velocity），或者就是简单的3V，即庞大容量、极快速度和种类丰富的数据”如图1-3所示。

#### 1. Volume（数量）

用现有技术无法管理的数据量，从现状来看，基本上是指从几十TB到几PB这样的数量级。当然，随着技术的进步，这个数值也会不断变化。

如今，存储的数据数量正在急剧增长中，人们存储所有事物包括：环境数据、财务数据、医疗数据、监控数据等。有关数据量的对话已从TB级别转向PB级别，并且不可避免地会转向ZB级别。但是，随着可供企业使用的数据量不断增长，可处理、理解和分析的数据的比例却不断下降。

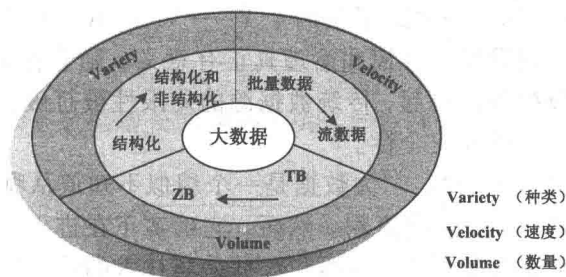


图 1-3 按数量、种类和速度来定义大数据

<sup>①</sup> 麦肯锡公司：是世界级领先的全球管理咨询公司。自 1926 年成立以来，公司的使命就是帮助领先的企业机构实现显著、持久的经营业绩改善，打造能够吸引、培育和激励杰出人才的优秀组织机构。

麦肯锡在全球 52 个国家有 94 个分公司。在过去十年中，麦肯锡在中国地区完成了 800 多个项目，涉及公司整体与业务单元战略、企业金融、营销/销售与渠道、组织架构、制造/采购/供应链、技术、产品研发等领域。

麦肯锡的经验是：关键是找那些企业的领导们，他们能够认识到公司必须不断变革以适应环境变化，并且愿意接受外部的建议，这些建议在帮助他们决定做什么变革和怎样变革方面大有裨益。





## 2. Variety (种类、多样性)

随着传感器、智能设备以及社交协作技术的激增，企业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、互联网日志文件（包括单击流数据）、搜索索引、社交媒体论坛、电子邮件、文档、主动和被动系统的传感器数据等原始、半结构化和非结构化数据。

种类表示所有的数据类型。其中，爆发式增长的一些数据，如互联网上的文本数据、位置信息、传感器数据、视频等，用企业中主流的关系型数据库是很难存储的，它们都属于非结构化数据。

当然，在这些数据中，有一些是过去就一直存在并保存下来的。和过去不同的是，除了存储，还需要对这些大数据进行分析，并从中获得有用的信息，例如监控摄像机中的视频数据。近年来，超市、便利店等零售企业几乎都配备了监控摄像机，最初目的是为了防范盗窃，但现在也出现了使用监控摄像机的视频数据来分析顾客购买行为的案例。

例如，美国高级文具制造商万宝龙（Montblane）过去是凭经验和直觉来决定商品陈列布局的，现在尝试利用监控摄像头对顾客在店内的行为进行分析。通过分析监控摄像机的数据，将最想卖出去的商品移动到最容易吸引顾客目光的位置，使得销售额提高了20%。

## 3. Velocity (速度)

数据产生和更新的频率，也是衡量大数据的一个重要特征。就像搜集和存储的数据量和种类发生了变化一样，生成和需要处理数据的速度也在变化。不要将速度的概念限定为与数据存储相关的增长速率，应动态地将此定义应用到数据，即数据流动的速度。有效处理大数据需要在数据变化的过程中对它的数量和种类进行分析，而不只是在它静止后执行分析。

例如，遍布全国的便利店在24h内产生的POS机数据、电商网站中由用户访问所产生的网站点击流数据、高峰时达到每秒近万条的微信短文、全国公路上安装的交通堵塞探测传感器和路面状况传感器（可检测结冰、积雪等路面状态）等，每天都在产生着庞大的数据。

IBM在3V的基础上又归纳总结了第四个V——Veracity（真实和准确）。只有真实而准确的数据才能让对数据的管控和治理真正有意义。随着社交数据、企业内容、交易与应用数据等新数据源的兴起，传统数据源的局限性被打破，企业愈发需要有效的信息治理以确保其真实性及安全性。

IDC（互联网数据中心）说：“大数据是一个貌似不知道从哪里冒出来的大的动力。但实际上，大数据并不是新生事物。然而，它确实正在进入主流，并得到重大关注，这是有原因的。廉价的存储、传感器和数据采集技术的快速发展、通过云和虚拟化存储设施增加的信息链路，以及创新软件和分析工具，正在驱动着大数据。大数据不是一个‘事物’，而是一个跨多个信息技术领域的动力/活动。大数据技术描述了新一代的技术和架构，其被设计用于：通过使用高速（Velocity）的采集、发现和/或分析，从超大容量（Volume）的多样（Variety）数据中经济地提取价值（Value）。”

这个定义除了揭示大数据传统的3V基本特征，还增添了一个新特征：Value（价值）。总之，大数据是个动态的定义，不同行业根据其应用的不同有着不同的理解，