

大数据 智能分析

张华平 商建云 刘兆友◎编著

大数据、人工智能与自然语言处理三者融合贯通之作

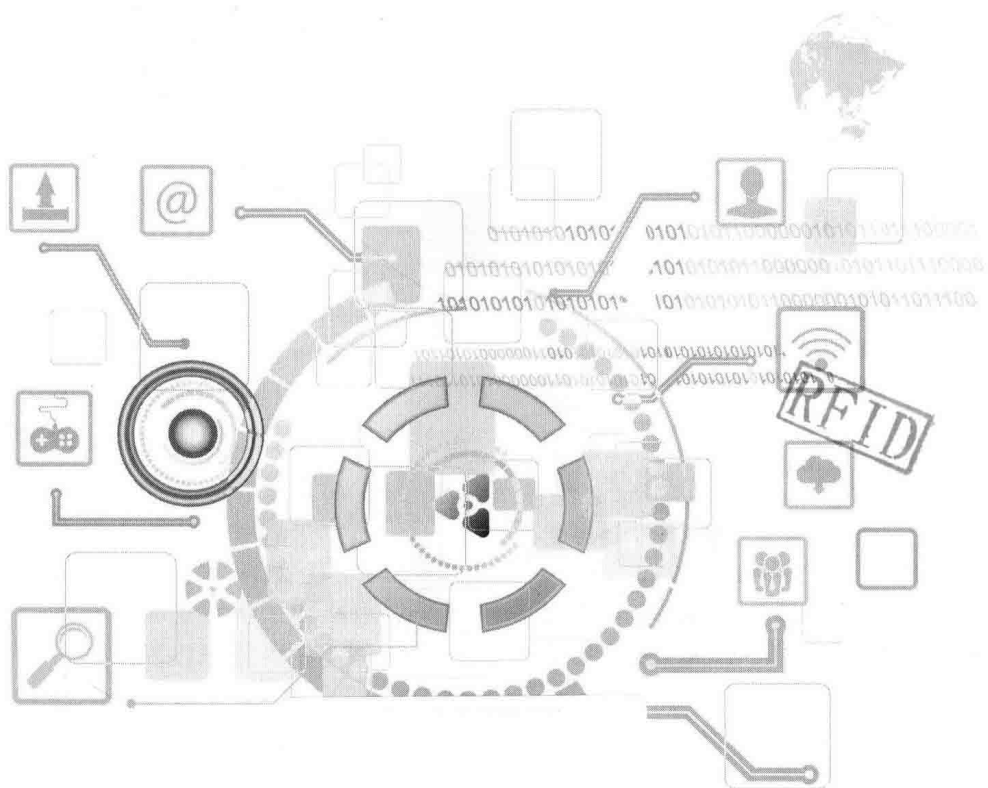
- ◎ 立足于作者二十年大数据智能的前沿研究与工程实践，结合北京理工大学“大数据分析与应用”研究生课程三年教学经验，体系完整，理论与实践并重。
- ◎ 可作为高校研究生与本科生的专业教材，也可作为大数据智能方向的科研人员和工程技术人员的实践参考书。

清华大学出版社



大数据 智能分析

张华平 商建云 刘兆友◎编著



清华大学出版社
北京

内 容 简 介

大数据智能是大数据、人工智能与自然语言处理等学科交叉融合的关键技术。本书主要讲述大数据智能的框架平台、理论算法、关键技术和应用实践：在大数据与人工智能方面主要讲述了大数据智能概述、大数据技术平台与架构、传统机器学习与深度学习算法；在自然语言处理方面详细讲解了大数据精准搜索、汉语分词、新词发现、文本分类聚类、情感分析等当前热门的自然语言处理关键技术；在应用实践方面，本书进一步提供了自主研发的 NLP-IR 大数据智能分析工具平台，具体介绍警情大数据、网络赌博、微博挖掘、看图说话等多个实际的大数据应用项目，也引入《红楼梦》前后作者分析、二手房房价、歌词生成等有意思的课程实践案例。

本书立足于作者近 20 年的前沿研究进展和工程实践，结合北京理工大学“大数据分析与应用”研究生课程讲授经验，体系完整，内容深入浅出，理论与实践并重，吸收了当前的技术前沿成果，同时突出原创的研究成果。本书可作为大数据、人工智能与自然语言处理方向的科研人员、高校研究生与本科生的教材，也可作为大数据智能方向的工程技术人员和爱好者的参考书。

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

大数据智能分析/张华平等编著. —北京：清华大学出版社，2019
ISBN 978-7-302-53117-3

I. ①大… II. ①张… III. ①数据处理 IV. ①TP274

中国版本图书馆 CIP 数据核字(2019)第 167582 号

责任编辑：白立军

封面设计：杨玉兰

责任校对：焦丽丽

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：20.5

字 数：499 千字

版 次：2019 年 10 月第 1 版

印 次：2019 年 10 月第 1 次印刷

定 价：59.80 元

产品编号：081294-01

前言

大数据智能是指从客观存在的全量超大规模、多源异构、实时变化的微观数据中,利用自然语言处理、信息检索、机器学习等技术抽取知识,转化而来的决策智慧的方法与过程。大数据智能涉及大数据、人工智能、自然语言处理三个相互支撑的关键技术。其中,大数据为大数据智能提供了关键的数据基础与大数据计算平台,是大数据智能的驱动力;人工智能为大数据提供了算法基础,是大数据智能的核心;自然语言处理直接面对数据中的语义内容,是人工智能“皇冠上的明珠”,直接决定大数据智能的广度与深度。

大数据、人工智能与自然语言处理是当前的研究热点,并被金融投资与商业开发广为追捧。习近平总书记在党的十九大报告中明确指出:“推动互联网、大数据、人工智能和实体经济深度融合。”2017年7月20日,国务院印发《新一代人工智能发展规划》,明确了我国发展人工智能的战略目标,到2030年,人工智能核心产业规模超过1万亿元,带动相关产业规模超过10万亿元。自然语言处理(Natural Language Processing, NLP)是计算机科学领域与人工智能领域中的一个重要方向。它研究能实现人与计算机之间用自然语言进行有效通信的各种理论和方法。NLP的目标是让机器在理解语言上像人类一样智能,最终目标是弥补人类交流(自然语言)和计算机理解(机器语言)之间的差距。人工智能已经成为现代科学“皇冠上的明珠”,而自然语言处理号称是“人工智能皇冠上的明珠”。微软全球副总裁、著名人工智能专家沈向洋在中国计算机大会上明确表示:“懂语言者得天下。”自然语言理解已经成为人工智能研究与投资的重点,孕育着改变世界未来的产业机会。Gartner发布的2017年商业智能和分析平台魔力象限报告中明确提出:到2020年,自然语言生成和人工智能将占现代商业智能平台标准特性的90%,同时50%的分析排队查询来自于搜索、自然语言处理或语音发起,或自动生成。2017年始语音及自然语言处理是创业最火热的领域之一,占据人工智能领域的11%。

大数据、人工智能或者自然语言处理的单一课程已经相对成熟,但很难适应当前科学技术的发展趋势,无法解决当前社交网络为代表的用户需求,将三者融会贯通的教材相对匮乏。笔者从2016年开始在北京理工大学为研究生一年级开设了必修课“大数据分析与应用”,尝试将大数据智能的教学体系进行完善,希望能出一本大数据智能的教材,这个想法得到北京理工大学“十三五”教材的立项支持。

本书共计 14 章。第 1 章对大数据智能进行了概述,并阐述关于大数据、人工智能与自然语言处理的主要体系架构,分享了作者关于大数据与人工智能的观点;第 2 章是大数据技术平台与架构,介绍了 Hadoop、Spark 等常用的云计算平台,介绍了机器学习、大数据可视化与大数据挖掘等常用工具;第 3 章介绍经典的机器学习与数据挖掘方法,包括管理规则挖掘、分类与聚类,最后为读者准备了常用的数据挖掘工具;第 4 章介绍经典的深度学习算法与平台,包括神经网络、循环神经网络(Recurrent Neural Network, RNN)、卷积神经网络(Convolutional Neural Network, CNN)、基于长短记忆的神经网络(Long Short-Term Memory, LSTM)、序列到序列模型(Sequence-to-Sequence)、注意力模型(Attention Model)以及生产对抗网络(Generative Adversarial Network, GAN)等,具体介绍了 TensorFlow 与 PyTorch 两个常用的深度学习平台;第 5 章介绍了作者研制的 JZSearch 大数据精准搜索引擎;第 6 章为汉语分词,主要介绍基于层次隐马尔可夫模型的浅层词法分析,以及基于深度学习的汉语分词算法;第 7 章介绍基于角色标注的命名实体识别;第 8 章介绍新词发现,主要基于上下文邻接熵与互信息的社交网络新词发现;第 9 章是文本分类与聚类;第 10 章重点推荐了基于 Top N 关键词的热点话题发现算法;第 11 章是当前关注比较多的情感分析算法,介绍情感词的自动发现与权重自动计算方法,详细介绍基于树模型的无监督情感分析算法与基于深度神经网络的短文本情感倾向性分析方法;第 12 章介绍文本摘要的关键技术;第 13 章与第 14 章都是大数据智能应用部分,其中第 13 章具体介绍了作者在警情大数据、网络赌博、微博挖掘、看图说话等多个实际的大数据应用项目,第 14 章是《红楼梦》前后作者分析、二手房房价、歌词生成等有意思的课程实践案例。

本书由张华平、商建云、郭涛与刘兆友合著。全书成果主要涉及张华平、商建云所在的大数据搜索与挖掘实验室,有些章节内容来自实验室近 20 年发表的学术论文与研究生毕业论文。其中,商建云完成了大数据应用等内容以及部分相关工作综述等的撰写工作;郭涛审阅了全书,并完成最后的统稿工作;刘兆友作为本书的助手,认真完成全部初稿的整理工作。在本书的写作与相关科研课题的研究工作中,得到多方面的支持与帮助,并采用了作者指导的研究生李清敏、张瑞琦、陈晓阳、李然、高莘、杨耀飞、李蕾、孙梦姝、王琦、季欣怡、张红瑾、Asif Khan 等学生的毕业论文及发表的文章;同时采用了北京理工大学“大数据分析与应用”部分同学的课程作业,均在相应的内容部分进行了标识。中国科学院信息工程研究所的张卫博、周艳对深度学习等部分内容进行了修订,在此谨向这些文献的作者以及为本书提供帮助的老师、同仁和课题组成员致以诚挚的谢意和崇高的敬意。本书亦得到 2018 年国家自然科学基金(编号:61772075)、北京理工大学“十三五”教材的资助。

最后,感谢我的太太与家人的支持,家庭的无私支持成就了这部书的写作与出版,本书同时献给望望与杨杨。

由于我们的学识、水平均有限,书中不妥之处在所难免,恳请广大读者批评指正。

张华平
2019 年 3 月

目 录

第 1 章 大数据智能概述	/1
1.1 数据的智能演化过程	/1
1.2 大数据	/2
1.2.1 大数据的概念	/2
1.2.2 大数据的特征	/2
1.2.3 大数据带来的决策方式的革命	/3
1.2.4 大数据面临的挑战及其对应的技术概览	/5
1.2.5 科学的大数据观	/9
1.2.6 大数据架构下的人才需求及产业结构	/10
1.3 人工智能	/12
1.4 自然语言处理	/14
第 2 章 大数据技术平台与架构	/16
2.1 大数据技术概览	/16
2.1.1 大数据技术架构	/16
2.1.2 云计算	/17
2.2 Hadoop、Spark 生态系统	/20
2.2.1 Hadoop 生态系统	/20
2.2.2 Spark 生态系统	/26
2.2.3 Spark 和 Hadoop 的性能对比	/31
2.3 大数据挖掘与可视化工具	/34
第 3 章 传统机器学习与数据挖掘	/40
3.1 机器学习介绍	/40
3.2 关联规则挖掘	/41
3.2.1 Apriori 算法	/43
3.2.2 FP-growth 算法	/43

- 3.3 分类 /45
 - 3.3.1 SVM /45
 - 3.3.2 决策树 /52
 - 3.3.3 朴素贝叶斯 /56
 - 3.3.4 K近邻 /59
- 3.4 聚类 /60
 - 3.4.1 基于划分的聚类方法 /60
 - 3.4.2 基于层次的聚类方法 /65
 - 3.4.3 基于密度的聚类方法 /71
 - 3.4.4 聚类案例：用户细分模型 /74
- 3.5 数据挖掘相关工具 /74
 - 3.5.1 数据获取工具 /75
 - 3.5.2 分词工具 /77
 - 3.5.3 分类聚类工具 /79
 - 3.5.4 Python调用方法 /79

第4章 经典深度学习算法与平台 /81

- 4.1 神经网络基础 /82
 - 4.1.1 神经元 /82
 - 4.1.2 从神经元到神经网络 /82
- 4.2 循环神经网络 /84
 - 4.2.1 RNN基本概念 /84
 - 4.2.2 RNN的长期依赖问题与LSTM /85
 - 4.2.3 深度RNN和双向RNN /88
- 4.3 卷积神经网络 /89
- 4.4 序列到序列模型 /90
- 4.5 注意力模型 /91
- 4.6 生成对抗网络 /93
- 4.7 TensorFlow计算图框架 /95
 - 4.7.1 数据流图 /95
 - 4.7.2 TensorFlow的特征 /95
 - 4.7.3 官方入门教程 /96
- 4.8 PyTorch深度学习框架 /103
 - 4.8.1 PyTorch是什么 /103
 - 4.8.2 自动求导：自动微分 /104
 - 4.8.3 神经网络 /105

第5章 信息检索与大数据搜索 /110

- 5.1 概述 /110

- 5.2 JZSearch 大数据搜索引擎系统架构 /110
- 5.3 大数据精准搜索的基本技术 /112
 - 5.3.1 索引字段类型 /112
 - 5.3.2 索引词项的设计 /113
 - 5.3.3 索引压缩技术 /113
 - 5.3.4 内存交换 /115
 - 5.3.5 增量索引 /116
 - 5.3.6 数据库检索 /117
- 5.4 大数据精准搜索语法 /118
 - 5.4.1 JZSearch 排序算法 /118
 - 5.4.2 JZSearch 结果格式 /119
 - 5.4.3 JZSearch 检索语法说明 /119
- 5.5 JZSearch 大数据精准搜索应用案例 /123
 - 5.5.1 中国邮政集团邮址垂直搜索 /124
 - 5.5.2 标准文档搜索引擎 /124
 - 5.5.3 内网文档的知识搜索门户 /125
 - 5.5.4 商品比价搜索 /125
 - 5.5.5 维吾尔文搜索 /125

第6章 汉语分词 /127

- 6.1 概述 /127
- 6.2 汉语分词的困难性 /129
- 6.3 基于机械匹配的汉语分词算法 /132
 - 6.3.1 词典匹配法 /132
 - 6.3.2 N -最短路径法 /136
- 6.4 基于统计语言模型的汉语分词算法 /137
 - 6.4.1 N 元语言模型 /138
 - 6.4.2 互信息模型 /138
 - 6.4.3 最大熵模型 /140
- 6.5 NLP-ICTCLAS: 基于层叠隐马尔可夫模型的汉语分词算法 /141
 - 6.5.1 层次隐马尔可夫模型 /141
 - 6.5.2 基于类的隐马尔可夫分词算法 /143
 - 6.5.3 N -最短路径的切分排歧策略 /145
- 6.6 基于双向循环神经网络与条件随机场的词法分析 /146
 - 6.6.1 概述 /146
 - 6.6.2 基于双向循环神经网络的序列标注 /146
 - 6.6.3 融合条件随机场的深度神经网络模型 /148
- 6.7 实验与分析 /149
 - 6.7.1 评估方法 /149

6.7.2 实验分析 1 /149

6.7.3 实验分析 2 /153

第 7 章 命名实体识别 /157

7.1 命名实体识别定义 /157

7.2 命名实体识别的研究主体 /158

7.3 命名实体识别的特点及难点 /158

7.4 命名实体识别的研究技术路径 /159

7.5 基于角色标注的命名实体识别 /159

7.6 实验与分析 /162

第 8 章 新词发现 /163

8.1 基于规则的研究方法 /164

8.1.1 规则抽取方法 /165

8.1.2 规则过滤方法 /165

8.2 基于统计模型的研究方法 /166

8.2.1 凝固度 /166

8.2.2 信息熵 /166

8.2.3 新词 IDF /167

8.3 面向社会媒体的开放领域新词发现 /167

8.3.1 引言 /167

8.3.2 新词发现 /168

8.3.3 实验 /171

第 9 章 文本分类与聚类 /175

9.1 文本预处理 /175

9.2 文本表示模型 /176

9.2.1 传统布尔检索与扩展布尔检索模型 /177

9.2.2 向量空间模型 /177

9.2.3 概率检索模型 /180

9.2.4 语言模型 /181

9.3 文本特征选择方法 /182

9.3.1 信息增量 /183

9.3.2 卡方统计 /183

9.3.3 交叉熵 /183

9.4 文本分类概述 /184

9.5 文本聚类概述 /187

9.5.1 聚类算法体系 /187

9.5.2 半监督聚类 /188

第 10 章 话题发现算法 /191

- 10.1 多语语义串自动发现 /195
- 10.2 多语语义关键特征挖掘 /197
 - 10.2.1 关键特征抽取 /197
 - 10.2.2 单个文档 Top N 关键特征挖掘 /198
- 10.3 Top N 热点话题发现和关联归并 /198
 - 10.3.1 Top N 热点话题发现 /198
 - 10.3.2 话题归并 /200
- 10.4 多语文本话题发现与关联归类实验验证 /201

第 11 章 情感分析 /203

- 11.1 概述 /203
- 11.2 情感分类 /205
- 11.3 应用 /208
 - 11.3.1 用户评论分析与决策 /208
 - 11.3.2 舆情监控 /208
 - 11.3.3 信息预测 /209
- 11.4 情感词发现与极性权重自动计算算法 /209
 - 11.4.1 引言 /209
 - 11.4.2 情感词典构建模型 /211
 - 11.4.3 实验 /213
- 11.5 基于树模型的无监督情感分析系统 /216
 - 11.5.1 实现方法 /216
 - 11.5.2 系统架构及流程 /217
 - 11.5.3 实验分析及结论 /219
- 11.6 基于深度神经网络的短文本情感倾向性分析 /221
 - 11.6.1 语料库建设 /221
 - 11.6.2 词袋模型与文本建模 /223
 - 11.6.3 基于 Softmax 和深度神经网络的短文本情感分析算法 /225
 - 11.6.4 实验设计及实验结果 /229

第 12 章 自动摘要 /234

- 12.1 概述 /234
- 12.2 基于关键词提取的自动摘要 /238
- 12.3 面向主题自动摘要 /244
- 12.4 基于主题模型与信息熵的中文文档自动摘要技术研究 /247
 - 12.4.1 主题模型 /248
 - 12.4.2 信息熵 /250
 - 12.4.3 句子信息熵的计算方法 /250



- 12.4.4 算法介绍 /250
- 12.4.5 实验结果 /251
- 12.5 自动摘要应用场景分析及大数据搜索与挖掘软件应用示例 /252

第 13 章 大数据智能应用案例 /254

- 13.1 公安警情大数据挖掘 /254
- 13.2 网络赌博信息文本挖掘 /257
 - 13.2.1 Web 网页信息选择与提取 /257
 - 13.2.2 中文分词及词性标注处理 /258
 - 13.2.3 特征提取 /259
 - 13.2.4 基于网络赌博信息的数据挖掘 /260
 - 13.2.5 网络赌博信息可视化展示 /262
- 13.3 领导人支持信息挖掘 /265
- 13.4 微博博主的特征与行为大数据挖掘 /268
 - 13.4.1 介绍 /268
 - 13.4.2 宏观特征大数据挖掘 /270
 - 13.4.3 实验与分析 /275
 - 13.4.4 微博博主的价值观自动评估方法 /275
- 13.5 看图说话：基于 Mask-RCNN 的图片中文描述生成器 /277
 - 13.5.1 自下而上的注意力机制在图像描述中的应用 /278
 - 13.5.2 Bottom-Up-Attention 和 Top-Down-Attention 图像描述模型 /280
 - 13.5.3 Dense-Attention 图像描述模型 /281
 - 13.5.4 基于语义控制的长短时记忆模型 /281
 - 13.5.5 模型训练相关说明及结果分析 /283
 - 13.5.6 模型测试相关说明及结果分析 /284
 - 13.5.7 测试结果分析 /286

第 14 章 大数据智能课程经典作业汇编 /288

- 14.1 《红楼梦》前后作者同一性分析 /288
- 14.2 党的十九大报告语义智能分析 /293
- 14.3 文章风格对比：方文山与汪峰 /294
- 14.4 智慧旅游大数据应用 /295
- 14.5 某大厦电力数据挖掘 /298
- 14.6 杭州市二手房房价分析 /301
 - 14.6.1 概述 /301
 - 14.6.2 房价分析系统案例介绍 /301
 - 14.6.3 本例设计与实现 /304
- 14.7 数据挖掘在股票分析预测中的应用 /306
 - 14.7.1 概述 /306

- 14.7.2 股票分析预测方法 /307
- 14.7.3 神经网络在股票分析预测应用中的研究现状 /307
- 14.7.4 实验结果 /309
- 14.8 基于 TensorFlow 的歌词自动生成 /310
 - 14.8.1 算法说明 /310
 - 14.8.2 实验结果 /311
- 14.9 基于 LSTM 的购物评论分类 /312
 - 14.9.1 获取语料库比分词 /312
 - 14.9.2 词向量的转换 /313
 - 14.9.3 建立向量和单词列表 /313
 - 14.9.4 将句子转换成序号矩阵 /314
 - 14.9.5 模型训练 /314

大数据智能概述

本章将从数据的智能演化过程、大数据、人工智能与自然语言处理 4 部分对大数据智能进行综合介绍。

1.1 数据的智能演化过程

我们先从数据的 4 个层次来洞察数据的智能演化过程,如图 1-1 所示。

数据: 用于表示客观事物的未经加工的原始素材,在计算机系统中,数据以二进制信息单元 0、1 的形式表示。换句话说,只要占据硬盘空间的,都可以认定为数据,但是,数据有可能是杂乱无章的,无法表达实际含义。例如,笔者传给读者 2MB 的数据,存储到计算机后,如果读者完全不了解这个数据的格式,就无从了解其实际含义。

信息: 1948 年,数学家香农在题为《通信的数学理论》的论文中指出:“信息是用来消除随机不定性的东西”。继续上面的例子,如果这 2MB 数据是一个人的 JPEG 格式的图片,采用相应软件打开后,就可以看到其中包含 11 张图片。在这个过程中,这个数据就代表了有价值的信息,读者消除了一个不确定性的东西:数据的格式是什么? 图片中的人究竟长什么样? 这个过程一般可以称为信息管理。信息管理就是人对信息资源和信息活动的管理。信息管理的过程包括信息收集、信息传输、信息加工、信息存储与信息检索。

知识: 人类在实践中认识客观世界的成果,知识是人类从各个途径中获得的经过提升总结与凝练的系统认识。知识本身是信息从量变到质变的产物,也是大数据的本质所在,即从多个表层信息中,利用大数据技术获取表层信息背后进一步的深度知识。依然以刚才的例子来说,11 张普通的“表哥”照片单个来看都是低价值密度的表层信息,但当 11 张普通信息以一定的方式组织起来,并聚焦到其所戴手表时,很容易挖掘出背后的知识:“表哥”拥有多款名贵手表,存在贪腐嫌疑。这个知识的价值远远超出了原始信息,这也是大数据的增值过程。

智能: 指知识进一步归纳总结后的更普世的规律,可演化为更多的知识,用来指导客观实践。例如,我们可以从“表哥”的知识中,举一反三,利用已有的各类报道,自动挖掘更多人

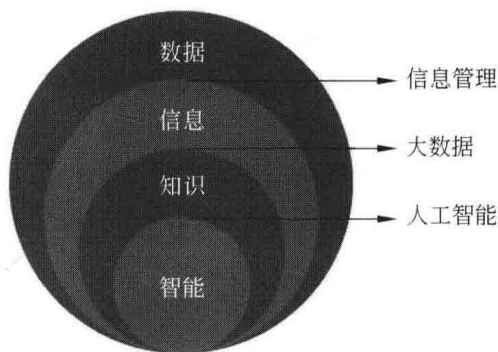


图 1-1 数据的智能演化过程

的奢侈消费,从而获取电子侦察智慧,随时发现更多问题,防患于未然。这个过程,人们常称为人工智能过程。

1.2 大数据

1.2.1 大数据的概念

关于大数据如何定义,不同机构有不同的定义。其中,研究机构 Gartner 的定义:大数据是指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。麦肯锡的定义:大数据是指无法在一定时间内用传统数据库软件工具对其内容进行采集、存储、管理和分析的数据集合。舍恩伯格·维克托的《大数据时代》中的定义:大数据是指不用随机分析法(抽样调查)这样的捷径,而采用所有数据的全域分析方法。

无论哪种定义,都可以看出,大数据并不是一种新的产品,也不是一种新的技术,就如同 21 世纪初提出的“海量数据”的概念一样,大数据只是数字化时代出现的一种现象。那么海量数据与大数据的差别何在?从翻译的角度来看,“大数据”和“海量数据”均来自英文, Big Data 翻译为“大数据”, Large-scale Data 翻译为“大规模数据”, Very Large Data 翻译为“超大规模数据”, Massive Data 则翻译为“海量数据”。从组成的角度来看,海量数据包括结构化和半结构化的数据,大数据除此以外还包括非结构化数据和交互数据。大数据由海量交易数据、海量交互数据和海量数据处理三大主要技术趋势汇聚而成,其规模和复杂程度超出了常用技术,按照合理的成本和时限捕捉、管理及处理这些数据集的能力。

1.2.2 大数据的特征

大数据的特征包含 4 个层面。第一,数据体量巨大。从 TB 级别跃升到 PB 级别。第二,数据类型繁多。例如,网络日志、视频、图片、地理位置信息等。第三,价值密度低。以视频为例,在连续不间断的监控过程中,可能有用的数据仅仅有一两秒。第四,流动速度快。1 秒定律,最后这一点和传统的数据挖掘技术有本质的不同。业界将大数据的特征归纳为 4V,即 Volume、Variety、Value、Velocity。

1. 数据体量巨大(Volume)

大数据通常指 10TB(1TB=1024GB)规模以上的数据量。之所以产生如此巨大的数据量,一是由于各种仪器的使用,使人们能够感知到更多的事物,这些事物的部分甚至全部数据可以被存储;二是由于通信工具的使用,使人们能够全时段联系,机器—机器(M2M)方式的出现,使得交流的数据量成倍增长;三是由于集成电路价格降低,使很多东西都向智能化发展。

2. 数据类型繁多(Variety)

随着传感器种类的增多及智能设备、社交网络等的流行,数据类型也变得更加复杂,不仅包括传统的关系数据类型,也包括以网页、视频、音频、E-mail、文档等形式存在的半结构化的和非结构化的数据。

3. 价值密度低(Value)

数据量呈指数级增长的同时,隐藏在海量数据中的有用信息却没有以相应比例增长,反而使人们获取有用信息的难度加大。

4. 流动速度快(Velocity)

通常理解的数据流动速度是指数据获取、存储及挖掘有效信息的速度。由于现在处理的数据是PB级代替了TB级,“超大规模数据”和“海量数据”也有规模大的特点,数据是快速动态变化的。因此,形成流式数据是大数据的重要特征,数据流动的速度快到难以用传统的系统去处理。

大数据的4V特征表明其不仅仅是数据海量,对于大数据的分析将更加复杂,更追求速度,更注重实效。

大数据独立发展形成特有的市场化与规模化,也充分带动了其他行业与大数据的广泛、充分融合,从而推进大数据的全面落地。大数据从产业到行业的成熟将推动更多传统企业向科技智能化转型,将推进政府政务大数据发展,也将鞭策大数据行业在中国平稳落地。

大数据之大,还在于数据结构的有容乃大——它不再需要传统的数据库表格来整齐排列,几乎可以无所不包地记录、存储和计算各种规则的结构化数据和不规则的非结构化数据,于是便有了逐步演变为一个数字化世界的可能。

1.2.3 大数据带来的决策方式的革命

近半个世纪以来,人们经历了计算机时代计算方式的革命、互联网时代信息传播方式的革命、大数据时代决策方式的革命,如表1-1所示。

表 1-1 半个世纪技术的悄然革命

发生时间	时代	带来的技术革命
20世纪70年代	计算机时代	计算方式的革命
20世纪90年代	互联网时代	信息传播方式的革命
21世纪初	大数据时代	决策方式的革命

计算方式大概经历了手工计算、机械计算、电子计算等不同的发展阶段。在运用计算机进行计算之前,所有的计算载体(如算盘等)只有计算功能,没有存储功能。计算机出现后,实现了既可以计算又可以记忆的功能,成为了人类的“外脑”。20世纪60年代,数据一般存储在文件中,由应用程序直接管理;20世纪70年代人们构建了关系数据模型,数据库技术为数据存储提供了新的手段;20世纪80年代中期,数据仓库由于具有面向主题、集成、时变和非易失等特点,成为数据分析和联机分析的重要平台。

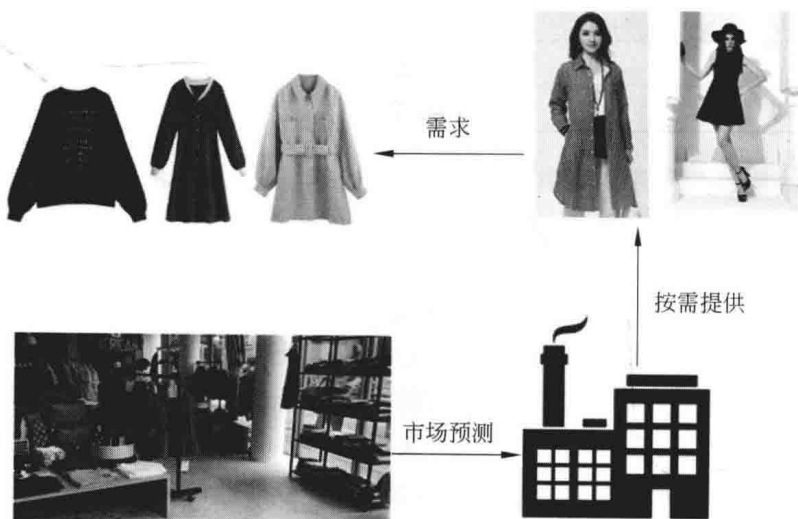
信息传播方式大概经历了口头传播、口头加书面(如报刊)传播、音频/视频和文字传播、新媒体传播等不同发展阶段。在互联网出现之前,信息传播都有一定的滞后性,如书信要邮递,报刊、电视节目都是先编辑或录制,均要经过一定的延时才会到达信息接收者手中。到了20世纪90年代,互联网出现,书信慢慢减少,电子邮件即刻到达。报刊、电视经历了多媒体的融合;拍客、博客的出现也使得人们从被动的信息接收者变成了信息的发布者和接收

者。每个人、每部手机都成为信息源。加上互联网上的人人连接、人物连接及物物连接相互产生了传播放大效应,带来了信息传播方式的革命。以前不会引起重大反响的事件,今天通过互联网的迅速传播,影响巨大。

决策方式经历了从直觉和经验到更加依赖数据和分析的转变。21世纪初,借助计算机软件 and 硬件的不断发展,信息产生方式发生了变化,企事业单位和互联网上都积累并实时产生大量的数据,从而出现了具有4V特征的大数据。用数据分析的方法进行辅助决策是这个时代的特点。图1-2(a)显示了以产品或服务提供者为中心的BSFD(Business, Supply, Feedback, Duration)决策流程;图1-2(b)显示了CDPR(Consumer, Demand, Prediction, Realtime)以消费者为中心的决策流程。



(a) 以产品或服务提供者为中心的BSFD决策流程



(b) 以消费者为中心的决策

图 1-2 大数据分析进行辅助决策的方式

决策方式可以从决策主体、决策依据、决策机制和决策效率4个方面来描述。表1-2显示了大数据带来的决策方式变革。

表 1-2 大数据带来的决策方式变革

名称	描述
决策主体	以商业供给侧为中心→以消费需求侧为中心(Business→Consumer)
决策依据	以商业供给能力为依据→以用户需求为依据(Supply→Demand)
决策机制	反馈机制→预测机制(Feedback→Prediction)
决策效率	期间决策→实时决策(Duration→Realtime)

在大数据的背景下,决策主体从关注商务转移到关注消费者,决策依据不再是生产者的生产能力,而主要是消费者的需求,得益于大数据的支撑,决策机制从反馈机制转变为预测机制,决策效率也从半年、一年决策变为实时决策。

笔者的亲身经历充分证实了这一点。1994年,笔者在日本看到路边有人用不同的砝码记录一天的人流情况,打听后得知他们将此数据带回去分析,用于预测今后的人流状况,进而指导人们的出行或决定道路的改建计划。2013年,笔者在美国,开车在路上行驶,导航仪随时提供道路上的信息。如今,在北京大家可以看到实时的路况信息,实时指导人们的出行道路选择及为道路的改建计划提供参考。

可以看出,决策方式的变化来源于计算方式和信息传播方式的改变。

1.2.4 大数据面临的挑战及其对应的技术概览

鉴于大数据的特征和决策方式的革命要求,在处理大数据时会面临新的挑战,新的数据搜集、存储、传输和处理模式会随之产生,下面结合新的挑战论述大数据技术。

1. 数据量的指数级增长挑战数据存储能力

在大数据的决策中,90%基于数据,10%基于直觉,因此要收集大量的数据,包括实时数据。收集数据量巨大、数据种类繁多、数据潜在商业价值高,这些均要求使用专门的数据库技术和专用的数据存储设备。传统的数据库追求高度的数据一致性和容错性,缺乏较强的扩展性和较好的系统可用性,不能有效存储视频、音频等非结构化和半结构化的数据。目前,数据存储能力的增长远远赶不上数据的增长,设计最合理的分层存储架构成为信息系统的關鍵。

在相关技术中,比较具有代表性的是 Apache 软件基金会开发的 Hadoop。以 MapReduce 和 Hadoop 为代表的非关系数据分析技术,凭借其适合非结构处理、大规模并行处理和简单易用等优势,在互联网搜索和其他大数据分析技术领域取得重大进展,成为主流技术。

MapReduce 是 2004 年谷歌公司提出的用来进行并行处理和生成大数据的模型,是一种线性的、可伸缩的编程模型。其可扩展性得益于 Shared-nothing 结构、各节点间的松耦合性和较强的软件级容错能力。MapReduce 被设计在处理时间内解释数据,所以对非结构化、半结构化的数据处理非常有效。针对 MapReduce 并行编程模型的易用性,产生了多种