

NUTE 国家软件与集成电路公共服务平台信息技术紧缺人才培养工程指定教材

大数据技术与应用丛书

有问题，就找黑马程序员问答精灵！

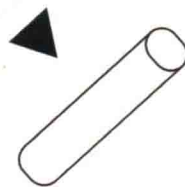
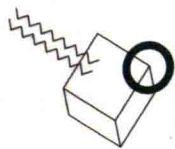
Hadoop

大数据技术原理与应用



黑马程序员 / 编著

清华大学出版社



NITE 国家软件与集成电路公共服务平台信息技术紧缺人才培养工程指定教材

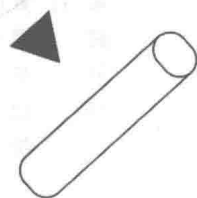
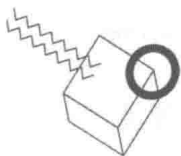
大数据技术与应用丛书

Hadoop



大数据技术原理与应用

黑马程序员 / 编著



清华大学出版社
北京

内 容 简 介

本书围绕 Hadoop 生态圈相关系统介绍大数据处理架构。全书共 11 章,其中,第 1、2 章主要带领大家认识 Hadoop 以及学会搭建 Hadoop 集群;第 3~5 章讲解分布式文件系统(HDFS)、分布式计算框架 MapReduce 以及分布式协调服务;第 6 章讲解 Hadoop 2.0 新特性,包含 YARN 和高可用特性;第 7~10 章主要讲解 Hadoop 生态圈的相关辅助系统,包括 Hive、Flume、Azkaban 和 Sqoop;第 11 章是一个综合项目——网站流量日志数据分析系统,目的是教会大家如何利用 Hadoop 生态圈技术构建大数据系统架构并进行开发,同时加深对 Hadoop 技术的理解。

本书附有配套视频、源代码、习题、教学设计、教学课件等资源。同时,为了帮助初学者更好地学习本书中的内容,还提供了在线答疑,欢迎读者关注。

本书可作为高等院校本、专科计算机相关专业,信息管理等相关专业的大数据课程教材,也可供相关技术人员参考,是一本适合广大计算机编程爱好者的优秀读物。

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

Hadoop 大数据技术原理与应用/黑马程序员编著. —北京:清华大学出版社,2019

(大数据技术与应用丛书)

ISBN 978-7-302-52440-3

I. ①H… II. ①黑… III. ①数据处理软件 IV. ①TP274

中国版本图书馆 CIP 数据核字(2019)第 040842 号

责任编辑:袁勤勇 杨 枫

封面设计:韩 冬

责任校对:李建庄

责任印制:丛怀宇

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

字 数:464 千字

版 次:2019 年 5 月第 1 版

印 次:2019 年 5 月第 1 次印刷

定 价:39.00 元

产品编号:083124-01

序 言



江苏传智播客教育科技有限公司(简称“传智播客”)是一家致力于培养高素质软件开发人才的科技公司。“黑马程序员”是传智播客旗下的高端 IT 教育品牌。

“黑马程序员”的学员多为大学毕业后,想从事 IT 行业,但各方面条件还不成熟的年轻人。“黑马程序员”的学员筛选制度非常严格,包括严格的技术测试、自学能力测试,还包括性格测试、压力测试、品德测试等,以百里挑一的残酷筛选制度确保学员质量,降低企业的用人风险。

自“黑马程序员”成立以来,教学研发团队一直致力于打造精品课程资源,不断在产、学、研三个层面创新自己的执教理念与教学方针,并集中“黑马程序员”的优势力量,有针对性地出版了计算机系列教材 60 多种,制作教学视频数十套,发表各类技术文章数百篇。

“黑马程序员”不仅斥资研发 IT 系列教材,还为高校师生提供以下配套学习资源与服务。

为大学生提供的配套服务

1. 请登录在线平台 <http://yx.boxuegu.com>, 免费获取海量学习资源,还有专业的老师在线为您解答。

2. 针对高校学生在学习过程中存在的压力等问题,我们还面向大学生量身打造了 IT 技术女神——“播妞”,可提供教材配套源代码和习题答案以及更多 IT 学习资源。同学们可以添加“播妞”微信号 208695827 和“播妞”QQ 号 3231342131,获取学习资源。



“播妞”微信



“播妞”QQ

为教师提供的配套服务

针对高校教学,“黑马程序员”为 IT 系列教材精心设计了“教案+授课资源+考试系

统+题库+教学辅助案例”的系列教学资源。高校老师请登录在线平台 <http://yx.boxuegu.com> 或关注码大牛老师微信/QQ2011168841, 获取配套资源, 也可以扫描下方二维码, 加入专为 IT 教师打造的师资服务平台——“教学好助手”, 获取最新教师教学辅助资源的相关动态。



前言

我们生活在一个充满“数据”的时代,刷微信、聊 QQ、网购、旅游、看病等一系列行为无时无刻不在产生新的数据,日积月累形成巨大的数据集,迎来了大数据时代。大数据时代的力量,正在积极地影响着人们生活的方方面面,深刻改变着人类的思维、生产、生活、学习方式,深刻展示了世界发展的前景。

大数据时代,数据的存储与挖掘至关重要。企业在追求高可靠性、高扩展性及高容错性的大数据处理平台的同时还希望能够降低成本,而 Hadoop 为实现这些需求提供了解决方案。这里列举 3 条使用 Hadoop 作为大数据业务的基础原因,具体如下。

(1) Hadoop 底层的分布式文件系统具有高拓展性,通过数据冗余保证数据不丢失和提升计算效率,同时可以存储各种格式的数据。它还有多种计算框架,既可以进行离线计算也可以进行在线实时计算。

(2) Hadoop 是架构在廉价的硬件服务器上,且产品是开源的,供开发者免费使用,开发成本和维护成本都降低很多。

(3) Hadoop 具有成熟的生态圈,有许多辅助系统对数据进行处理。

本书作为大数据技术 Hadoop 的入门教程,最重要又最难的一件事情就是将一些复杂、难以理解的问题简单化,让初学者能够轻松理解并快速掌握。本教材对每个知识点都进行了深入分析,并针对每个知识点精心设计了相关案例,然后模拟这些知识点在实际工作中的运用,真正做到了知识的讲解由浅入深、由易到难。

全书共分为 11 章。

第 1 章主要讲解什么是大数据以及 Hadoop 相关概念。通过本章的学习,读者可对大数据有简单的认识,并了解 Hadoop 生态圈工具及各自的用途。

第 2 章主要讲解 Hadoop 集群的构建。通过本章的学习,读者能掌握 Linux 系统网络配置、独立搭建 Hadoop 开发平台,以及简单操作 Hadoop 系统。

第 3 章主要讲解 Hadoop 分布式文件系统(HDFS)。通过本章的学习,读者可以掌握 HDFS 的架构和工作原理,并能够通过 Shell 接口和 Java API 操作 HDFS。

第 4 章主要讲解 MapReduce 的相关知识。通过本章的学习,初学者可以了解 MapReduce 计算框架的思想并且能够使用 MapReduce 解决实际问题。

第 5 章主要讲解 Zookeeper 分布式协调服务。通过本章的学习,读者能够对 Zookeeper 分布式协调服务有基本的认识,掌握 Zookeeper 内部运行原理,并会通过 Shell 和 Java API 操作 Zookeeper。

第 6 章主要讲解 Hadoop 2.0 的新特性,包括 YARN 资源管理框架和 HDFS 的高可用。其中,YARN 作为资源管理框架,读者需要明白它的体系结构和工作流程;HDFS 的高

可用性能够解决集群的单点故障问题,读者要掌握高可用架构的部署方式,并能独立参考文档搭建高可用的 Hadoop 集群。

第 7 章主要讲解 Hive 的相关知识。读者需要了解 Hive 架构、数据模型、Hive 的安装和管理以及 Hive 的数据操作。这里建议初学者在学习 Hive 时多动手操作 Hive,通过丰富的案例练习,掌握 Hive 的使用。

第 8 章主要讲解 Flume 日志采集系统的基本知识。通过本章的学习,读者应该掌握 Flume 的基本概念、运行机制并且能够掌握 Flume 的安装配置和基本使用。

第 9 章主要讲解 Azkaban 工作流管理器的基本知识。通过本章的学习,读者应该对 Azkaban 有一定的了解,掌握 Azkaban 的部署和使用,并能够使用 Azkaban 进行任务调度管理。

第 10 章主要讲解 Sqoop 数据迁移工具的相关知识。通过本章的学习,读者可以掌握 Sqoop 工作原理,会独立搭建 Sqoop 工具并且能够使用 Sqoop 工具完成常用的数据迁移操作。

第 11 章主要通过开发网站流量日志分析系统来讲解利用 Hadoop 生态体系的技术解决实际问题。通过本章的学习,读者可以了解大数据系统的架构、数据采集、数据预处理、数据仓库的设计、数据分析、数据导出以及最后可视化处理。读者应该熟练掌握系统架构以及业务流程,熟练使用 Hadoop 生态体系相关技术。

致谢

本书的编写和整理工作由传智播客教育科技有限公司完成,主要参与人员有吕春林、高美云、石荣新、翟振方、文燕等,全体参编人员在这近一年的编写过程中付出了许多辛勤的汗水,在此表示衷心的感谢。

意见反馈

尽管我们尽了最大的努力,但书中难免会有欠妥之处,欢迎各界专家和读者朋友们来信提出宝贵意见,我们将不胜感激。您在阅读本书时,如果发现任何问题或有不认同之处可以通过电子邮件与我们取得联系。

请发送电子邮件至 itcast_book@vip.sina.com。

黑马程序员

2019年3月于北京

目 录

第 1 章 初识 Hadoop	1
1.1 大数据概述	1
1.1.1 什么是大数据	1
1.1.2 大数据的特征	2
1.1.3 研究大数据的意义	3
1.2 大数据的应用场景	4
1.2.1 医疗行业的应用	4
1.2.2 金融行业的应用	4
1.2.3 零售行业的应用	5
1.3 Hadoop 概述	6
1.3.1 Hadoop 的前世今生	6
1.3.2 Hadoop 的优势	7
1.3.3 Hadoop 的生态体系	7
1.3.4 Hadoop 的版本	9
1.4 本章小结	11
1.5 课后习题	11
第 2 章 搭建 Hadoop 集群	13
2.1 安装准备	13
2.1.1 虚拟机安装	13
2.1.2 虚拟机克隆	22
2.1.3 Linux 系统网络配置	24
2.1.4 SSH 服务配置	28
2.2 Hadoop 集群搭建	31
2.2.1 Hadoop 集群部署模式	31
2.2.2 JDK 安装	32
2.2.3 Hadoop 安装	33
2.2.4 Hadoop 集群配置	35
2.3 Hadoop 集群测试	38
2.3.1 格式化文件系统	38

专属于老师及学生的在线教育平台
<http://yx.boxuegu.com/>

让 IT 教学更简单

教师获取教材配套资源



添加微信/QQ
2011168841

让 IT 学习更有效

学生获取课后作业习题答案及配套源码

添加播妞QQ: 3231342131

添加播妞微信: 208695827



专属大学生的圈子

2.3.2	启动和关闭 Hadoop 集群	39
2.3.3	通过 UI 查看 Hadoop 运行状态	41
2.4	Hadoop 集群初体验	43
2.5	本章小结	46
2.6	课后习题	46
第 3 章	HDFS 分布式文件系统	48
3.1	HDFS 的简介	48
3.1.1	HDFS 的演变	48
3.1.2	HDFS 的基本概念	50
3.1.3	HDFS 的特点	51
3.2	HDFS 的架构和原理	52
3.2.1	HDFS 存储架构	52
3.2.2	HDFS 文件读写原理	53
3.3	HDFS 的 Shell 操作	55
3.3.1	HDFS Shell 介绍	55
3.3.2	案例——Shell 定时采集数据到 HDFS	58
3.4	HDFS 的 Java API 操作	62
3.4.1	HDFS Java API 介绍	62
3.4.2	案例——使用 Java API 操作 HDFS	63
3.5	本章小结	68
3.6	课后习题	69
第 4 章	MapReduce 分布式计算框架	70
4.1	MapReduce 概述	70
4.1.1	MapReduce 核心思想	70
4.1.2	MapReduce 编程模型	71
4.1.3	MapReduce 编程实例——词频统计	72
4.2	MapReduce 工作原理	73
4.2.1	MapReduce 工作过程	73
4.2.2	MapTask 工作原理	74
4.2.3	ReduceTask 工作原理	75
4.2.4	Shuffle 工作原理	76
4.3	MapReduce 编程组件	77
4.3.1	InputFormat 组件	77
4.3.2	Mapper 组件	78
4.3.3	Reducer 组件	78
4.3.4	Partitioner 组件	80
4.3.5	Combiner 组件	80

4.3.6	OutputFormat 组件	81
4.4	MapReduce 运行模式	82
4.5	MapReduce 性能优化策略	84
4.6	MapReduce 经典案例——倒排索引	86
4.6.1	案例分析	86
4.6.2	案例实现	89
4.7	MapReduce 经典案例——数据去重	93
4.7.1	案例分析	93
4.7.2	案例实现	93
4.8	MapReduce 经典案例——TopN	96
4.8.1	案例分析	96
4.8.2	案例实现	97
4.9	本章小结	100
4.10	课后习题	100
第5章	Zookeeper 分布式协调服务	102
5.1	初识 Zookeeper	102
5.1.1	Zookeeper 简介	102
5.1.2	Zookeeper 的特性	103
5.1.3	Zookeeper 集群角色	103
5.2	数据模型	104
5.2.1	数据存储结构	104
5.2.2	Znode 的类型	105
5.2.3	Znode 的属性	105
5.3	Zookeeper 的 Watch 机制	106
5.3.1	Watch 机制的简介	106
5.3.2	Watch 机制的特点	106
5.3.3	Watch 机制的通知状态和事件类型	107
5.4	Zookeeper 的选举机制	107
5.4.1	选举机制的简介	107
5.4.2	选举机制的类型	108
5.5	Zookeeper 分布式集群部署	109
5.5.1	Zookeeper 安装包的下载安装	109
5.5.2	Zookeeper 相关配置	109
5.5.3	Zookeeper 服务的启动和关闭	112
5.6	Zookeeper 的 Shell 操作	113
5.6.1	Zookeeper Shell 介绍	113
5.6.2	通过 Shell 命令操作 Zookeeper	113
5.7	Zookeeper 的 Java API 操作	119

5.7.1	Zookeeper Java API 介绍	119
5.7.2	通过 Java API 操作 Zookeeper	120
5.8	Zookeeper 典型应用场景	122
5.8.1	数据发布与订阅	122
5.8.2	统一命名服务	123
5.8.3	分布式锁	123
5.9	本章小结	123
5.10	课后习题	124
第 6 章	Hadoop 2.0 新特性	125
6.1	Hadoop 2.0 改进与提升	125
6.2	YARN 资源管理框架	125
6.2.1	YARN 体系结构	125
6.2.2	YARN 工作流程	127
6.3	HDFS 的高可用	128
6.3.1	HDFS 的高可用架构	128
6.3.2	搭建 Hadoop 高可用集群	129
6.4	本章小结	134
6.5	课后习题	135
第 7 章	Hive 数据仓库	136
7.1	数据仓库简介	136
7.1.1	什么是数据仓库	136
7.1.2	数据仓库的结构	137
7.1.3	数据仓库的数据模型	138
7.2	Hive 简介	140
7.2.1	什么是 Hive	140
7.2.2	Hive 系统架构	141
7.2.3	Hive 工作原理	141
7.2.4	Hive 数据模型	142
7.3	Hive 的安装	143
7.3.1	Hive 安装模式简介	143
7.3.2	嵌入模式	144
7.3.3	本地模式和远程模式	145
7.4	Hive 的管理	147
7.4.1	CLI 方式	147
7.4.2	远程服务	148
7.5	Hive 内置数据类型	150
7.6	Hive 数据模型操作	151

7.6.1	Hive 数据库操作	151
7.6.2	Hive 内部表操作	153
7.6.3	Hive 外部表操作	157
7.6.4	Hive 分区表操作	158
7.6.5	Hive 桶表操作	163
7.7	Hive 数据操作	166
7.8	本章小结	170
7.9	课后习题	170
第 8 章	Flume 日志采集系统	172
8.1	Flume 概述	172
8.1.1	Flume 简介	172
8.1.2	Flume 运行机制	172
8.1.3	Flume 日志采集系统结构图	173
8.2	Flume 基本使用	175
8.2.1	Flume 系统要求	175
8.2.2	Flume 安装配置	175
8.2.3	Flume 入门使用	177
8.3	Flume 采集方案配置说明	181
8.3.1	Flume Sources	181
8.3.2	Flume Channels	184
8.3.3	Flume Sinks	186
8.4	Flume 的可靠性保证	189
8.4.1	负载均衡	189
8.4.2	故障转移	195
8.5	Flume 拦截器	196
8.6	案例——日志采集	198
8.6.1	案例分析	198
8.6.2	案例实现	199
8.7	本章小结	204
8.8	课后习题	205
第 9 章	工作流管理器(Azkaban)	206
9.1	工作流管理器概述	206
9.1.1	工作流调度系统背景	206
9.1.2	常用工作流管理器介绍	206
9.2	Azkaban 概述	207
9.2.1	Azkaban 特点	208
9.2.2	Azkaban 组成结构	208

9.2.3	Azkaban 部署模式	209
9.3	Azkaban 部署	210
9.3.1	Azkaban 资源准备	210
9.3.2	Azkaban 安装配置	212
9.3.3	Azkaban 启动测试	220
9.4	Azkaban 使用	224
9.4.1	Azkaban 工作流相关概念	224
9.4.2	案例演示——依赖任务调度管理	226
9.4.3	案例演示——MapReduce 任务调度管理	232
9.4.4	案例演示——HIVE 脚本任务调度管理	235
9.5	本章小结	237
9.6	课后习题	237
第 10 章	Sqoop 数据迁移	239
10.1	Sqoop 概述	239
10.1.1	Sqoop 简介	239
10.1.2	Sqoop 原理	240
10.2	Sqoop 安装配置	241
10.3	Sqoop 指令介绍	242
10.4	Sqoop 数据导入	244
10.4.1	MySQL 表数据导入 HDFS	245
10.4.2	增量导入	247
10.4.3	MySQL 表数据导入 Hive	248
10.4.4	MySQL 表数据子集导入	249
10.5	Sqoop 数据导出	251
10.6	本章小结	253
10.7	课后习题	253
第 11 章	综合项目——网站流量日志数据分析系统	255
11.1	系统概述	255
11.1.1	系统背景介绍	255
11.1.2	系统架构设计	255
11.1.3	系统预览	256
11.2	模块开发——数据采集	257
11.2.1	使用 Flume 搭建日志采集系统	257
11.2.2	日志信息说明	258
11.3	模块开发——数据预处理	258
11.3.1	分析预处理的数据	258
11.3.2	实现数据的预处理	259

11.4	模块开发——数据仓库开发	268
11.4.1	设计数据仓库	268
11.4.2	实现数据仓库	269
11.5	模块开发——数据分析	273
11.5.1	流量分析	273
11.5.2	人均浏览量分析	274
11.6	模块开发——数据导出	275
11.7	模块开发——日志分析系统报表展示	276
11.7.1	搭建日志分析系统	277
11.7.2	实现报表展示功能	285
11.7.3	系统功能模块展示	290
11.8	本章小结	290

第 1 章

初识Hadoop

学习目标

- 了解大数据的概念及其特征。
- 熟悉大数据的典型应用。
- 了解 Hadoop 的发展历史及其版本。
- 掌握 Hadoop 的生态体系。

随着近几年计算机技术和互联网的发展，“大数据”这个词被提及得越来越频繁。与此同时，大数据的快速发展无时无刻不在影响着我们的生活，例如，医疗方面，大数据能够帮助医生预测疾病；电商方面，大数据能够向顾客个性化推荐商品；交通方面，大数据会帮助人们选择最佳出行方案。

Hadoop 作为一个能够对大量数据进行分布式处理的软件框架，用户可以利用 Hadoop 生态体系开发和处理海量数据。由于 Hadoop 可靠及高效的处理性能，使得它逐渐成为分析大数据的领先平台。接下来，将深入介绍大数据以及 Hadoop 的相关概念，为后面知识的学习建立概念体系。

1.1 大数据概述

1.1.1 什么是大数据

高速发展的信息时代，新一轮科技革命和变革正在加速推进，技术创新日益成为重塑经济发展模式和促进经济增长的重要驱动力量，而“大数据”无疑是核心推动力。

那么，什么是“大数据”呢？如果从字面意思来看，大数据指的是巨量数据。那么可能有人会问，多大量级的数据才叫大数据？不同的机构或学者有不同的理解，难以有一个非常定量的定义，只能说，大数据的计量单位已经越过 TB 级别发展到 PB、EB、ZB、YB 甚至 BB 级别。

最早提出“大数据”这一概念的是全球知名咨询公司麦肯锡，它是这样定义大数据的：一种规模大到在获取、存储、管理、分析方面大大超出了传统数据库软件工具能力范围的数据集合，具有海量的数据规模、快速的数据流转、多样的数据类型以及价值密度低四大特征。

研究机构 Gartner 是这样定义大数据的：“大数据”是需要新处理模式才能具有更强的决策力、洞察发现力和流转优化能力来适应海量、高增长率和多样化的信息资产。

若从技术角度来看，大数据的战略意义不在于掌握庞大的数据，而在于对这些含有意义

的数据进行专业化处理,换言之,如果把大数据比作一种产业,那么这种产业盈利的关键在于提高对数据的“加工能力”,通过“加工”实现数据的“增值”。

1.1.2 大数据的特征

一般认为,大数据主要具有以下 4 个方面的典型特征,即大量(Volume)、多样(Variety)、高速(Velocity)和价值(Value),即所谓的“4V”,接下来,通过一张图 1-1 来具体描述。

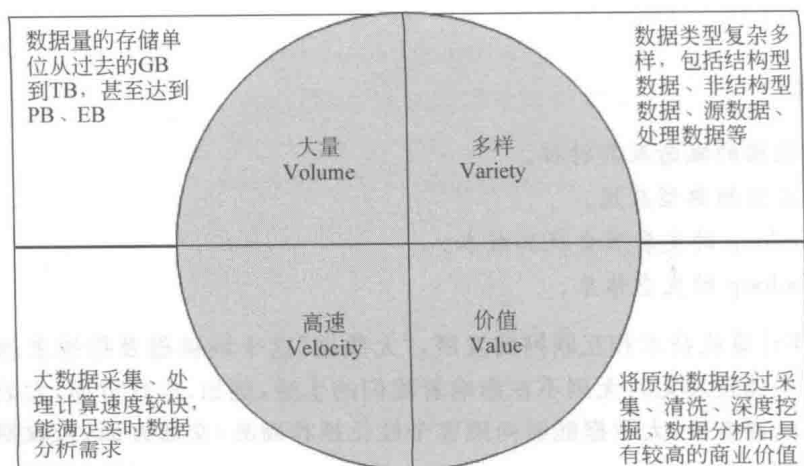


图 1-1 大数据 4V 特征

接下来针对图 1-1 中的 4V 特征进行简要介绍,具体如下。

1. Volume(大量)

大数据的特征首先就是数据规模大。随着互联网、物联网、移动互联技术的发展,人和事物的所有轨迹都可以被记录下来,数据呈现出爆发性增长。数据相关计量单位的换算关系如表 1-1 所示。

表 1-1 单位换算关系

单 位	换 算 公 式	单 位	换 算 公 式
Byte	1Byte=8bit	TB	1TB=1024GB
KB	1KB=1024Byte	PB	1PB=1024TB
MB	1MB=1024KB	EB	1EB=1024PB
GB	1GB=1024MB	ZB	1ZB=1024EB

2. Variety(多样)

数据来源的广泛性,决定了数据形式的多样性。大数据可以分为三类,一是结构化数据,如财务系统数据、信息管理系统数据、医疗系统数据等,其特点是数据间因果关系强;二是非结构化的数据,如视频、图片、音频等,其特点是数据间没有因果关系;三是半结构化数

据,如 HTML 文档、邮件、网页等,其特点是数据间的因果关系弱。有统计显示,目前结构化数据占据整个互联网数据量的 75% 以上,而产生价值的大数据,往往是这些非结构化数据。

3. Velocity(高速)

数据的增长速度和处理速度是大数据高速性的重要体现。与以往的报纸、书信等传统数据载体生产传播方式不同,在大数据时代,大数据的交换和传播主要是通过互联网和云计算等方式实现的,其生产和传播数据的速度是非常迅速的。另外,大数据还要求处理数据的响应速度要快,例如,上亿条数据的分析必须在几秒内完成。数据的输入、处理与丢弃必须立刻见效,几乎无延迟。

4. Value(价值)

大数据的核心特征是价值,其实价值密度的高低和数据总量的大小是成反比的,即数据价值密度越高数据总量越小,数据价值密度越低数据总量越大。任何有价值的信息的提取依托的就是海量的基础数据。当然目前大数据背景下有个未解决的问题,如何通过强大的机器算法更迅速地在海量数据中完成数据的价值提纯。

1.1.3 研究大数据的意义

现在的社会是一个高速发展的社会,科技发达,信息流通,人们之间的交流也越来越密切,生活也越来越便捷,大数据就是这个高科技时代的产物。阿里巴巴创办人马云曾经说过,未来的时代将不是 IT 时代,而是 DT 的时代,DT 就是 Data Technology,数据科技,这显示出大数据对于阿里巴巴集团来说是举足轻重的。

有人把数据比喻为蕴藏能量的煤矿。煤炭按照性质有焦煤、无烟煤、肥煤、贫煤等分类,而露天煤矿、深山煤矿的挖掘成本又不一样。与此类似,大数据并不在于“大”,而在于“有用”。数据的价值含量、挖掘成本比数量更为重要。对于很多行业而言,如何利用这些大规模数据,发掘其潜在价值,才是赢得核心竞争力的关键。

研究大数据,最重要的意义是预测。因为数据从根本上来讲,是对过去和现在的归纳和总结,其本身不具备趋势和方向性的特征,但是可以应用大数据去了解事物发展的客观规律、了解人类行为,并且能够帮助我们改变过去的思维方式,建立新的数据思维模型,从而对未来进行预测和推测。比如,商业公司对消费者日常的购买行为和使用商品习惯进行汇总和分析,了解到消费者的需求,从而改进已有商品并适时推出新的商品,消费者的购买欲就会提高。知名互联网公司谷歌对其用户每天频繁搜索的词汇进行数据挖掘,从而进行相关的广告推广和商业研究。

大数据的处理技术迫在眉睫,近年来各国政府和全球学术界都掀起了一场大数据技术的革命,众人纷纷积极研究大数据的相关技术。很多国家都把大数据技术研究上升到了国家战略高度,提出了一系列的大数据技术研发计划,从而推动政府机构、学术界、相关行业和各类企业对大数据技术进行探索和研究。

可以说大数据是一种宝贵的战略资源,其潜在价值和增长速度正在改变着人类的工作、生活和思维方式。可以想象,在未来,各行各业都会积极拥抱大数据,积极探索数据挖掘和