

大数据分析5个重要阶段：

数据获取和记录  
数据抽取、清洗和注记  
数据集成、聚集和表示  
数据分析和建模  
数据解释

大数据分析3个应用实例：

股票市场预测系统  
海量视频检索系统  
HDFS云文件系统

大数据分析技术：

Hadoop  
HDFS  
HBase  
MapReduce

# 实战大数据

怎么做大数据分析 大数据分析怎么应用在业务系统上

鲍亮 李倩 编著

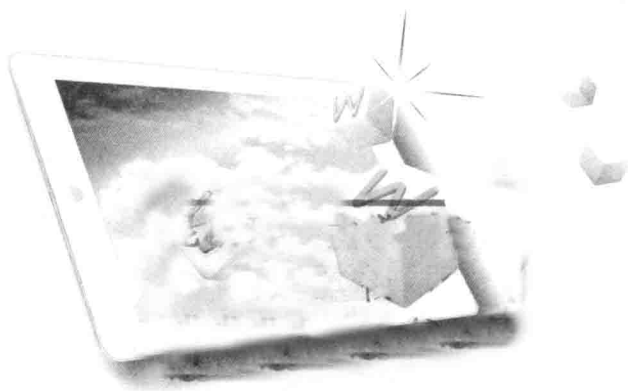


清华大学出版社



# 实战大数据

鲍亮 李倩 编著



清华大学出版社  
北京

## 内 容 简 介

“数据是重要资产”已成为大家的共识，众多公司都在争相分析、挖掘大数据背后的信息资源。本书在此背景下，对目前大数据及其相关技术的发展进行总结，理论联系实际，既不乏理论深度又具有实用价值。

本书共 12 章，内容包括大数据的概念、特点、发展历史，数据获取与存储，数据抽取和清洗，数据集成，数据的查询、分析与建模，异构数据采集，文档的存储与检索，异种数据的统一访问与转换，基于微博的股票市场预测系统实例，海量视频检索系统实例，HDFS 云文件系统实例。

本书适合大数据技术初学者、大数据从业人员和研究人员，也可以作为高等院校相关专业师生的教学参考书。

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

实战大数据 / 鲍亮, 李倩编著. - 北京: 清华大学出版社, 2014

ISBN 978-7-302-34866-5

I. ①实… II. ①鲍… ②李… III. ①数据处理—研究 IV. ①TP274

中国版本图书馆 CIP 数据核字 (2013) 第 310963 号

责任编辑: 夏非彼

封面设计: 王 翔

责任校对: 闫秀华

责任印制: 宋 林

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](mailto:c-service@tup.tsinghua.edu.cn)

质 量 反 馈: 010-62772015, [zhiliang@tup.tsinghua.edu.cn](mailto:zhiliang@tup.tsinghua.edu.cn)

印 装 者: 河北新华第一印刷有限责任公司

经 销: 全国新华书店

开 本: 190mm×260mm

印 张: 33.75

字 数: 860 千字

版 次: 2014 年 3 月第 1 版

印 次: 2014 年 3 月第 1 次印刷

印 数: 1~3500

定 价: 79.00 元



# 前 言

大数据时代已经到来，大数据处理已经成为当今信息处理的热点研究内容。不同于大规模数据，大数据具有自身鲜明的**4V**特征：**Volume**（规模性）、**Variety**（多样性）、**Velocity**（高速性）和**Veracity**（真实性）。大数据不仅规模大，更需要采取新的数据思维来应对，其必然导致理论和技术上的革新。因此，大数据分析也被认为是继实验、理论和计算之后的科学研究第四范式。大数据的出现必将颠覆传统的数据管理方式，在数据来源、数据处理方式和数据思维方面都会对其带来革命性的变化。

2013年初，美国计算机协会数据库专家委员会联合研究界、产业界和政府部门的相关研究人员，发布了大数据研究白皮书，提出了大数据分析的**5**个重要阶段：数据获取和记录，数据抽取、清洗和注记，数据集成、聚集和表示，数据分析和建模，数据解释。在这**5**个阶段中需要考虑数据的异构性、规模、时效性、复杂性和隐私问题。本书以此为提纲进行内容组织，首先介绍了**5**个阶段中相关的科学与技术问题，然后以实际案例的形式详细介绍了数据采集、数据存储与检索、数据处理、数据访问与转换**4**个大数据领域的重要问题，最后以股票市场预测系统、海量视频检索系统和云文件系统**3**个大数据实际应用系统为例详细介绍如何进行问题分析、数据建模以及系统的设计与实现。本书强调理论联系实际，重点在于介绍如何利用现有技术解决实际的大数据问题。

目前市场上以大数据为主题的书籍较多，但经过作者调研，未见以“利用现有技术解决大数据问题”为主题的大数据实战类书籍。本书编写团队核心成员自**2010**年起陆续承担了一些与大数据采集、存储、处理、分析、挖掘和检索方面的研究与应用开发工作，具有丰富的项目实践经验。这些实际项目经验形成了本书最为核心的第**6~12**章的内容。通过项目实战，我们积累了一些解决大数据问题的宝贵经验，对大数据的核心技术有了较为深刻的理解，认为有必要将自己的经验和认识整理出来，以满足广大读者利用现有技术解决大数据实际问题的迫切需求与心情，这也是书名的由来。

本书适合不同层次的读者阅读，建议读者根据自己的兴趣和目的有选择性地阅读：希望了解大数据相关的基础理论与技术的读者，可以重点阅读第**1~5**章；对于大数据领域的初学者，可以重点阅读第**1~9**章；对于已经掌握大数据基础理论，具有一定的技术基础，想解决实际

大数据问题的读者，可以重点阅读第 10~12 章。

除封面署名的作者之外，参与编写的还有李江、张翔、杨阳、王贺、刘凯、王学良、张静、周文琳、刘晓静、张艳华、王炎楠、黄鹏、高小青。还需要感谢阚传奇、蒋帆的大力帮助，感谢我的导师陈平教授在大数据科学研究方面对我的启发与悉心指导。

由于大数据涉及的学科面很广，研究问题纷繁复杂，相关资料目前还比较少，加之作者水平有限，时间紧迫，书中难免存在错误与不当，恳请读者批评指正。建议和意见请发至作者邮箱 [baoliang@mail.xidian.edu.cn](mailto:baoliang@mail.xidian.edu.cn)。

编者

2013 年 12 月

# 目 录

## 第一篇 大数据基础篇

第 1 章 大数据介绍 .....	2
1.1 大数据相关概念 .....	2
1.1.1 大数据的历史 .....	2
1.1.2 大数据的定义 .....	3
1.2 大数据研究内容 .....	6
1.3 大数据研究现状 .....	10
1.3.1 学术界现状 .....	10
1.3.2 产业界现状 .....	12
1.3.3 政府机构现状 .....	15
1.4 大数据的应用领域 .....	18
1.4.1 大数据在制造业的应用 .....	19
1.4.2 大数据在服务业的应用 .....	20
1.4.3 大数据在交通行业的应用 .....	20
1.4.4 大数据在医疗行业的应用 .....	20
1.5 本章小结 .....	21
第 2 章 数据存储技术 .....	22
2.1 数据存储技术介绍 .....	23
2.2 数据采集与存储技术研究现状 .....	25
2.2.1 传统关系型数据库 .....	25
2.2.2 新兴数据存储系统 .....	26

2.3	海量数据存储的关键技术分析 .....	27
2.3.1	数据划分 .....	27
2.3.2	数据一致性与可用性 .....	28
2.3.3	负载均衡 .....	29
2.3.4	容错机制 .....	29
2.3.5	海量数据存储的硬件支持 .....	30
2.4	数据存储技术的实现与工具 .....	36
2.4.1	集中式数据存储管理系统 Bigtable .....	36
2.4.2	非集中式的大规模数据管理系统 Dynamo .....	44
2.4.3	BigTable 的开源实现 HBase .....	50
2.4.4	MongoDB .....	52
2.4.5	CouchDB .....	55
2.4.6	Redis .....	56
2.4.7	Hypertable .....	60
2.4.8	其他开源 NoSQL 数据库 .....	62
2.5	本章小结 .....	69
<b>第 3 章</b>	<b>数据抽取和清洗 .....</b>	<b>70</b>
3.1	数据抽取和清洗技术介绍 .....	71
3.1.1	数据抽取简介 .....	71
3.1.2	数据清洗简介 .....	73
3.2	数据抽取和清洗研究现状 .....	76
3.3	数据抽取技术的实现 .....	78
3.3.1	Web 数据抽取 .....	78
3.3.2	非结构化数据抽取 .....	93
3.3.3	基于云计算的海量数据分析 .....	100
3.4	数据清洗技术的实现 .....	103
3.4.1	数据清洗流程 .....	103
3.4.2	数据清洗框架 .....	105
3.4.3	数据清洗相关技术 .....	109
3.4.4	基于 Hadoop 的数据清洗方案 .....	115

3.5 ETL 现状与发展.....	122
3.5.1 数据 ETL 简介 .....	122
3.5.2 基于 MapReduce 的 ETL 框架 .....	123
3.5.3 ETL 工具 .....	130
3.5.4 ETL 展望 .....	137
3.6 本章小结.....	138
<b>第 4 章 数据集成.....</b>	<b>139</b>
4.1 数据集成技术介绍 .....	139
4.2 数据集成技术研究现状 .....	141
4.2.1 Information Manifold: 具有统一的查询接口 .....	141
4.2.2 数据集成系统的发展建设 .....	144
4.2.3 企业信息集成 .....	147
4.2.4 未来的挑战 .....	148
4.3 数据集成技术的实现与工具 .....	150
4.3.1 Oracle Data Integrator (ODI) 简介 .....	150
4.3.2 ODI 的特点 .....	156
4.3.3 Microsoft SQL Server Integration Services (SSIS) 简介 .....	156
4.3.4 SSIS 的特点 .....	160
4.3.5 IBM InfoSphere Information Server 简介 .....	162
4.3.6 Sybase Data Integrator Suite 简介 .....	168
4.4 本章小结 .....	174
<b>第 5 章 数据查询、分析与建模技术 .....</b>	<b>175</b>
5.1 数据查询、分析与建模技术介绍 .....	175
5.1.1 数据查询 .....	175
5.1.2 数据分析 .....	177
5.1.3 数据建模 .....	177
5.2 数据查询、分析与建模技术研究现状.....	178
5.2.1 并行处理 .....	178
5.2.2 海量数据查询与搜索 .....	180
5.2.3 数据分析中的 OLAP 与数据挖掘技术.....	183



5.2.4	数据模型与数据建模方法 .....	191
5.3	数据查询、分析与建模技术的实现与工具 .....	194
5.3.1	数据查询相关技术实现与工具 .....	194
5.3.2	数据分析相关技术实现与工具 .....	200
5.3.3	数据建模相关技术实现与工具 .....	211
5.4	本章小结 .....	215

## 第二篇 大数据深入篇

第 6 章	采用 OSGi 框架构建可伸缩的异构数据采集平台 .....	217
6.1	应用背景 .....	217
6.2	需求分析与总体设计 .....	219
6.2.1	功能需求 .....	219
6.2.2	非功能需求 .....	220
6.2.3	总体设计 .....	220
6.3	相关技术介绍 .....	222
6.3.1	OSGi 框架介绍 .....	222
6.3.2	多源异构数据的获取 .....	226
6.4	系统设计与实现 .....	232
6.4.1	异构数据采集平台的设计 .....	232
6.4.2	数据采集插件的设计与实现 .....	236
6.4.3	系统服务框架的设计与实现 .....	245
6.5	部署与测试 .....	251
6.5.1	系统部署 .....	251
6.5.2	系统测试 .....	253
6.6	本章小结 .....	257
第 7 章	采用 HBase 实现海量小型 XML 文档的存储与检索 .....	258
7.1	应用背景 .....	258
7.2	需求分析与总体设计 .....	259
7.2.1	需求分析 .....	259
7.2.2	总体设计 .....	265

7.3 相关技术介绍 .....	268
7.3.1 XML 相关技术 .....	268
7.3.2 XQuery 语句 .....	269
7.3.3 XML 检索技术 .....	270
7.3.4 云计算和 HBase .....	272
7.3.5 JavaCC 工具介绍 .....	274
7.4 详细设计与实现 .....	275
7.4.1 数据存储模块的详细设计与实现 .....	276
7.4.2 数据检索模块的详细设计与实现 .....	289
7.4.3 用户模块的详细设计与实现 .....	299
7.5 本章小结 .....	301
<b>第 8 章 采用 Map/Reduce 进行大规模社交网络社团发现 .....</b>	<b>302</b>
8.1 研究背景 .....	302
8.2 相关理论和技术 .....	305
8.2.1 社团结构 .....	305
8.2.2 相关社团发现算法 .....	306
8.2.3 Hadoop 分布计算框架 .....	309
8.3 RMS 算法的并行化实现 .....	312
8.3.1 RMS 算法 .....	312
8.3.2 RMS 算法在 MapReduce 上的实现 .....	314
8.4 AP 聚类算法的并行化实现 .....	317
8.4.1 AP 聚类算法 .....	317
8.4.2 AP 聚类算法在 MapReduce 上的实现 .....	319
8.5 实验与分析 .....	324
8.5.1 实验环境 .....	324
8.5.2 实验与结果分析 .....	325
8.6 本章小结 .....	327
<b>第 9 章 数据统一访问与转换平台 .....</b>	<b>329</b>
9.1 应用背景介绍 .....	329
9.2 数据统一访问需求分析与总体设计 .....	333

9.2.1	功能性需求分析.....	333
9.2.2	非功能性需求分析.....	338
9.2.3	总体设计.....	339
9.3	数据统一访问与转换关键技术.....	342
9.3.1	SDO 编程技术.....	342
9.3.2	Hadoop MapReduce 框架.....	349
9.3.3	HBase 数据库技术.....	351
9.3.4	模型驱动数据转换技术.....	353
9.4	数据统一访问和灵活转换的详细设计与实现.....	355
9.4.1	数据分析及预处理.....	355
9.4.2	基于 DAS 的数据源统一访问.....	360
9.4.3	映射模式表示与数据存储管理模块.....	369
9.4.4	基于 MapReduce 的数据转换管理模块.....	374
9.5	本章小结.....	378

### 第三篇 大数据应用篇

第 10 章	基于微博的股票市场预测系统.....	380
10.1	应用背景介绍.....	380
10.2	需求分析与总体设计.....	382
10.2.1	需求分析.....	382
10.2.2	总体设计.....	391
10.3	相关技术介绍.....	393
10.3.1	社交网络.....	393
10.3.2	社交网络表示方法.....	395
10.3.3	信息传播模型.....	396
10.4	详细设计与实现.....	398
10.4.1	Twitter 数据采集模块详细设计.....	398
10.4.2	Twitter 数据分析模块详细设计.....	401
10.4.3	用户行为分析模块详细设计.....	407
10.4.4	预测股票价格涨跌模块详细设计.....	413

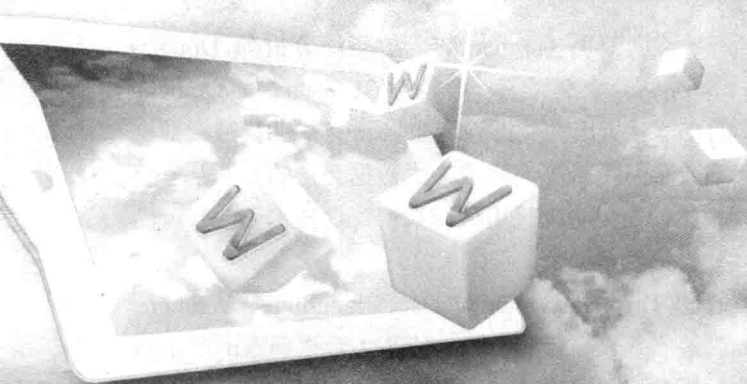
10.4.5 系统实现 .....	419
10.5 本章小结 .....	424
<b>第 11 章 基于内容的大量视频检索系统 .....</b>	<b>426</b>
11.1 应用背景 .....	426
11.2 需求分析与总体设计 .....	427
11.2.1 功能需求 .....	427
11.2.2 非功能需求 .....	431
11.2.3 核心业务处理流程 .....	431
11.2.4 总体设计 .....	435
11.3 相关技术简介 .....	438
11.3.1 MPEG-7 与 OpenCV 简介 .....	438
11.3.2 运动对象提取 .....	440
11.3.3 星形骨架方法 .....	443
11.4 详细设计与实现 .....	449
11.4.1 基于 MapReduce 的视频预处理 .....	449
11.4.2 基于 HBase 的视频数据存储 .....	455
11.4.3 行为识别与运动规则的组合创建 .....	470
11.5 系统运行时截图 .....	475
11.6 本章小结 .....	477
<b>第 12 章 基于 HDFS 的云文件系统 .....</b>	<b>478</b>
12.1 应用背景介绍 .....	478
12.2 需求分析与总体设计 .....	479
12.2.1 需求分析 .....	479
12.2.2 总体设计 .....	488
12.3 相关技术介绍 .....	491
12.3.1 Hadoop HDFS 介绍 .....	491
12.3.2 主控节点和数据节点 .....	493
12.3.3 页面展现技术 .....	494
12.3.4 页面控制技术 .....	494
12.4 详细设计与实现 .....	495

12.4.1	云文件系统的操作流程.....	495
12.4.2	云文件系统的模块设计.....	496
12.4.3	云文件系统实现.....	506
12.4.4	云文件系统主要功能截图.....	519
12.5	本章小结.....	525

## 第一篇

---

# 大数据基础篇



# 第 1 章

## 大数据介绍

IT 行业总不乏新鲜的主题，而大数据正当其兴，被业界热情传诵。“数据是重要资产”这一概念已成为大家的共识，众多公司争相分析、挖掘大数据背后的重要资源。为了帮助读者理解大数据的来龙去脉，本章将从大数据的历史与发展、大数据的定义、大数据的研究内容、大数据问题在国内外政府、公司和大学的研究现状等方面进行论述，为这一新兴概念勾勒出一个雏形。

## 1.1 大数据相关概念

### 1.1.1 大数据的历史

大数据 (Big Data) 目前已经成为 IT 领域最为流行的词汇，其实它并不是一个全新的概念。早在 1980 年，著名未来学家阿尔文·托夫勒便在《第三次浪潮》一书中，明确提出“数据就是财富”这一观点，并将大数据热情地赞颂为“第三次浪潮的华彩乐章”。

大数据中的“大”是一个相对概念，数据库、数据仓库、数据集市等信息管理领域的技术，很大程度上也是为了解决大规模数据的问题。被誉为数据仓库之父的 Bill Inmon 早在 20 世纪 90 年代就经常将“大数据”这一概念挂在嘴边了。

目前得到广泛认可的大数据概念首先由知名咨询公司 Gartner 的一位资深分析师 Douglas Laney 提出。他于 2001 年在 *Application Delivery Strategies* 上撰写了一篇名为“3D Data Management: Controlling Data Volume, Velocity, and Variety”的文章，指出大数据管理面临三个 V 的挑战：数据量 (Volume)、数据多样性 (Variety)、高速 (Velocity)。“3V”后来成为大数据公认的三个基本特征。随后，Gartner 发布了大数据的模型，强调大数据需要管理采用传统数据管理技术无法管理的数据，比如微博数据、海量交易数据、多媒体数据，等等。

2008 年 9 月，自然杂志推出《大数据》专刊，通过“The next Google”、“Data wrangling”、“Welcome to the petacentre”、“Distilling meaning from data”等多篇文章，全方位介绍了大数据问题的产生及各个研究领域的影响，首次将“大数据”这一概念引入科学家和研究人员的视野。

2009 年 8 月，Adam Jacobs 在 *ACM Queue* 上发表文章“The Pathologies of Big Data”，文章讨论了大数据问题的起源、发展与现状，指出“大数据”这一概念是相对的，并提出应该考虑为什

么会出现“大数据”这一现象、“大数据”产生的很大一部分原因是数据录入更加容易等观点。

2011年2月11日的《科学》杂志专门推出《数据处理》(*Dealing with Data*)专刊,对大数据现象在科学领域的现状进行了全面分析。该专刊首先联合《科学》杂志的兄弟期刊 *Science Signaling*、*Science Translational Medicine* 和 *Science Careers*, 展开了对各科学领域研究数据规模急剧增大情况下各种问题的调研,问题包括“研究数据的规模”、“研究数据如何存储”,等等。随后,该专刊发表多篇文章,对天文学、气象学、生态学、神经科学、信号处理、社会科学、生物学等多个学科的大数据问题进行了解释和阐述,内容涵盖数据采集、分析、处理、挖掘和可视化等多个方面。

2011年5月,麦肯锡全球研究院发表 *Big data: The next frontier for innovation, competition, and productivity* 白皮书,指出企业正在面临海量的交易数据、顾客信息、供货商信息和运营数据等,需要对这些数据进行管理与挖掘。在物联网环境下,传感器、智能手机、工业设备等都在产生海量数据。互联网中的多媒体数据量也在以指数级上升,如何处理这些数据,为用户提供有用的信息,成为需要考虑的重要问题。

2011年5月26日,经济学人发表“Building with big data”指出在数据极度膨胀的时代,要掌握数据的分析与处理能力,成为数据的主人,而不要成为数据的奴隶。

2012年2月11日,纽约时报发表“The Age of Big Data”,向大众宣传大数据时代的到来。

2012年3月22日,奥巴马宣布以2亿美元投资大数据领域,在次日的电话会议上,美国政府将数据定义为“未来的新石油”,美国政府认识到了一个国家拥有数据的规模、活性及解释运用的能力将成为综合国力的重要组成部分,未来对数据的占有和控制甚至将成为继陆权、海权、空权之外的另一种国家核心资产。

2012年7月10日,联合国在纽约总部发布了一份大数据政务白皮书,总结了各国政府如何利用大数据更好地服务和保护人民。

## 1.1.2 大数据的定义

### 1. 维基百科的定义

大数据是指其大小或复杂性无法通过现有常用的软件工具,以合理的成本并在可接受的时限内对其进行捕获、管理和处理的数据集。这些困难包括数据的收入、存储、搜索、共享、分析和可视化。

### 2. Gartner 的定义

Gartner 咨询公司关注大数据的三个量化指标:数据量、数据种类和处理速度。一般企业所面对的数据管理管理的是数据库、结构化数据,以及所能预先安装好的管理软件所带来的数据。大数据管理的往往是我们无法管理的数据,比如来自企业外部,微博、社交网站和多媒体等各种载体的数据。

数据多样性将是大数据的一个重点。它意味着未来数据的产生将更加方便、快捷,无所不在。



数据种类随着物联网等技术的不断兴起而飞速增加，特别是以多媒体数据为代表的非结构化数据迅速增加，为大数据的分析与处理带来了很大难度。处理速度与企业 CIO 关注的系统性能不是等同的关系。这里的速度指的是从数据产生到最终针对数据产生决策的速度，包括存储的过程、计算的过程、系统模型和以什么方式提交出最后的结果。因此，速度不仅是计算能力和存储性能的问题，还要考虑数据管理、数据保护等方面的响应与处理速度。在大数据问题中，速度往往是性命攸关的。比如对于灾难的预测，当灾难发生时，要很快对灾难发生的程度、影响的区域范围、对长远的影响等量化出来。这是大数据很典型的应用，如果短时间内没有计算出来，那么数据就没用了。

另一方面，Gartner 认为在越来越大的数据集上工作能够得到更大的好处，大数据的数据增长挑战（或机遇）是三维立体的：不断增长的数据量、不断增加的速率（数据 I/O 的速度）和不断增加的种类（数据类型、数据源）。而传统的存储技术难以应对大数据处理的三大挑战。

- 挑战一：不断增长的数据量。在大数据背景下，数据通常是不能删除的，这是企业的宝贵的财富，因此数据将不断积累增长。与此同时，增长有加速的趋势，经常会超出人们预计或规划，从而对信息系统带来了极大的挑战。信息中心需要管理 TB 级甚至 PB 级数据。要为这些数据提供存储、保护和使用的方案，IT 系统需要不断地做相应升级或重构，需要投入大量人力物力。
- 挑战二：多格式数据。海量数据包括了越来越多不同格式的数据，这些不同格式的数据也需要不同的处理方法。从简单的电子邮件、数据日志和信用卡记录，再到仪器收集到的科学研究数据、医疗数据、财务数据以及丰富的媒体数据（包括照片、音乐、视频等），都具有这个特点。比如视频文件格式就非常多，有各软件厂商的厂商标准的格式，工业标准组织的工业标准格式。各种格式在当前高清化的趋势下，数据粒度更小，处理更精细，更复杂的格式还不断出现，造成单一文件的体积成倍增加，从而要求处理速度也成倍增加。
- 挑战三：性能。速度是指数据从客户端到处理器和存储的移动速度，涉及终端数据处理能力、数据流访问和交付、服务器计算处理能力以及后端存储的吞吐能力。速度意味着要求数据必须以多快的频率被处理。大数据处理需要不同于交易类应用的速度，通常其对带宽的要求比 IO 操作的速度更重要。

### 3. IBM 对大数据的定义

IBM 专门开辟了大数据专栏，从大数据的定义、大数据处理平台等多个方面对大数据问题及解决方案进行了阐述。

#### **Big data spans three dimensions: Volume, Velocity and Variety.**

Volume: Enterprises are awash with ever-growing data of all types, easily amassing terabytes—even petabytes—of information.

Velocity: Sometimes 2 minutes is too late. For time-sensitive processes such as catching fraud, big data must be used as it streams into your enterprise in order to maximize its value.

Variety: Big data is any type of data - structured and unstructured data such as text, sensor data, audio, video, click streams, log files and more. New insights are found when analyzing these data types together.