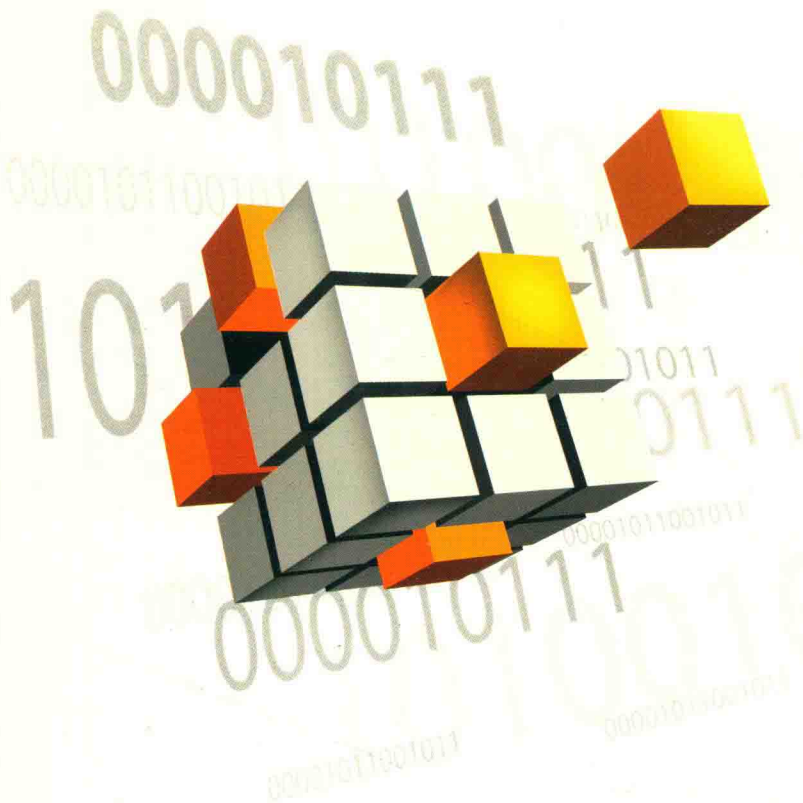




新世纪高等学校规划教材·大数据系列



陈明◎编著

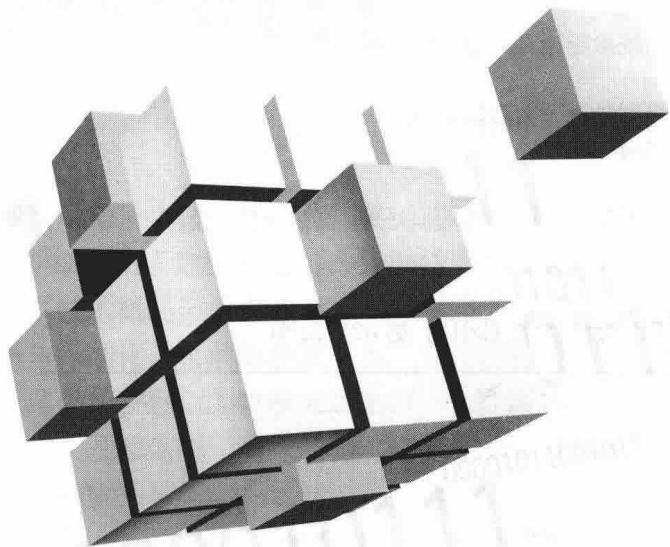
# 数据科学与 大数据技术导论

- ◆ 注重基本内容与基本方法的介绍
- ◆ 走进数据科学与大数据技术大门



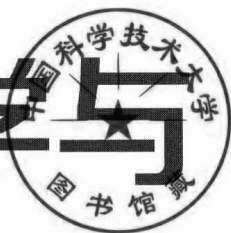
北京师范大学出版集团  
BEIJING NORMAL UNIVERSITY PUBLISHING GROUP  
北京师范大学出版社

新世纪高等学校规划教材·大数据系列



陈明◎编著

# 数据科学与 大数据技术导论



SHUJU KEXUE YU DASHUJU JISHU DAOLUN



北京师范大学出版集团  
BEIJING NORMAL UNIVERSITY PUBLISHING GROUP  
北京师范大学出版社

---

图书在版编目 (CIP) 数据

数据科学与大数据技术导论/陈明编著. —北京: 北京师范大学出版社, 2018. 6

新世纪高等学校规划教材·大数据系列

ISBN 978-7-303-23452-3

I. ①数… II. ①陈… 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 (2018) 第 025112 号

---

营 销 中 心 电 话 010-62978190 62979006

北师大出版社科技与经管分社网 <http://jsws.bnupg.com>

电 子 信 箱 [kjjg@bnupg.com](mailto:kjjg@bnupg.com)

---

出版发行: 北京师范大学出版社 [www.bnupg.com](http://www.bnupg.com)

北京新街口外大街 19 号

邮政编码: 100875

印 刷: 北京京师印务有限公司

经 销: 全国新华书店

开 本: 787mm×1092mm 1/16

印 张: 22.75

字 数: 511 千字

版 次: 2018 年 6 月第 1 版

印 次: 2018 年 6 月第 1 次印刷

定 价: 49.80 元

---

策划编辑: 赵洛育

责任编辑: 赵洛育

美术编辑: 刘 超

装帧设计: 刘 超

责任校对: 王 云

责任印制: 赵非非

**版权所有 侵权必究**

反盗版、反侵权举报电话: 010-62978190

北京读者服务部电话: 010-62979006-8021

外埠邮购电话: 010-62978190

本书如有印装质量问题, 请与印制管理部联系调换。

印制管理部电话: 010-62979006-8006

# 内 容 简 介

大数据技术是一个面向实际应用的技术。从大数据中获取有价值信息是大数据技术的精髓。本书详细介绍了数据科学与大数据技术的主要内容，全书分为 13 章，主要包括数据科学与大数据技术概述、大数据获取与存储管理技术、大数据抽取技术、数据清洗技术、大数据去噪与标准化、大数据约简技术、大数据集成技术、大数据挖掘技术、大数据分析、大数据分析结果的解释、大数据机器学习、大数据离线计算技术、大数据流式计算技术等。本书在内容上，注重概念、方法介绍，实例丰富、语言精练、逻辑层次清晰，可作为大学数据科学与大数据技术专业和相关专业的教材，也可以作为科技人员的参考书。

# 前 言

大数据技术与应用展现出锐不可当的强大生命力，科学界与企业界寄予无比的厚望。大数据成为继 20 世纪末、21 世纪初互联网蓬勃发展以来的新一轮 IT 工业革命。

大数据技术是指从数据采集、清洗、集成、挖掘、分析与结果解释，进而从各种各样的巨量数据中快速获得有价值信息的全部技术。大数据技术的精髓是从数据挖掘和分析中获取具有重要价值的信息、产生新见解的能力、识别复杂关系和做出更加精准的预测。

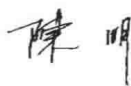
大数据技术是现代科学与技术发展，尤其是计算机科学技术发展的重要成果和结晶，是科学发展史一个新的里程碑。大数据的出现对计算机许多领域提出了挑战与冲击，推动了计算机科学技术的发展。

大数据技术的出现凝集了多学科的研究成果，是一门多学科的交叉融合技术，随着科学技术的进步，大数据技术发展更为迅速，应用更为深入与广泛，并凸显其巨大潜力和应用价值。

本书在体例的设计方面，以大数据技术为核心展开。大数据技术是指在大数据处理周期中所使用的全部技术的集合，也就是说，从数据获取之后，经过存储、抽取、清洗、约简、分析、挖掘等步骤后，获得有价值信息的全过程所需要的技术集合。本书的第 1 章～第 11 章主要介绍了上述过程。除此之外，还在第 12 章～第 13 章中介绍了大数据离线计算技术和大数据流式计算技术等内容。其中，在大数据离线计算技术中，主要包括大数据离线处理架构、MapReduce 的体系结构、基于 Hadoop 框架的分布计算和 MapReduce 程序设计实例分析。在大数据流式计算技术中，主要包括流式数据的计算模式、流式计算的应用、流式计算的系统架构、高可用技术、Storm 处理过程和大数据流式计算的应用案例等。

本书在内容的描述方面，注重大数据技术的主要概念、结构和方法的清晰描述。对主要的算法，如分类算法、聚类算法等给出了形式化描述。

本书在结构上为积木状，各章内容独立地进行概念性与方法性论述。出于篇幅考虑，书中所提及定理没有给出证明，如需要可以查阅相关文献。由于作者水平有限，书中不足之处在所难免，敬请读者批评指正。



2017.8

# 目 录

第 1 章 概述.....	1	1.6.2 数据密集型科学研究第四 范式 .....	22
1.1 数据科学.....	2	本章小结 .....	27
1.1.1 数据科学的产生与发展 .....	2	第 2 章 大数据获取与存储管理技术....	28
1.1.2 数据科学的相关术语 .....	2	2.1 大数据获取 .....	29
1.1.3 数据科学的主要内容 .....	3	2.1.1 大数据获取的挑战 .....	29
1.1.4 数据科学的研究过程与体系 框架 .....	4	2.1.2 传统数据获取与大数据获取 区别 .....	29
1.1.5 数据科学、数据技术与数据 工程 .....	7	2.2 领域数据的获取 .....	30
1.1.6 大数据问题 .....	8	2.2.1 文本数据获取 .....	30
1.2 大数据的生态环境 .....	8	2.2.2 语音视频数据获取 .....	31
1.2.1 互联网世界 .....	9	2.2.3 图片数据获取 .....	31
1.2.2 物理世界 .....	11	2.2.4 摄像头视频数据获取 .....	31
1.3 大数据的概念 .....	12	2.2.5 图像数字化数据获取 .....	31
1.3.1 数据容量 .....	12	2.2.6 图形数字化数据获取 .....	32
1.3.2 数据类型 .....	13	2.2.7 O2O LBS 数据获取 .....	32
1.3.3 价值密度 .....	13	2.2.8 空间数据获取 .....	32
1.3.4 速度 .....	14	2.3 网站数据 .....	34
1.3.5 真实性 .....	14	2.3.1 网站内部数据获取 .....	34
1.4 大数据的性质 .....	14	2.3.2 网站外部数据获取 .....	35
1.4.1 非结构性 .....	14	2.3.3 移动网站数据获取 .....	36
1.4.2 不完备性 .....	16	2.4 大数据存储 .....	37
1.4.3 时效性 .....	16	2.4.1 大数据存储模型 .....	37
1.4.4 安全性 .....	16	2.4.2 大数据存储问题 .....	37
1.4.5 可靠性 .....	16	2.4.3 大数据存储方式 .....	38
1.5 大数据处理周期 .....	16	2.5 大数据的存储管理技术 .....	40
1.5.1 大数据处理的全过程 .....	17	2.5.1 数据容量问题 .....	40
1.5.2 大数据技术的特征 .....	19	2.5.2 大图数据 .....	40
1.5.3 大数据的几个热点问题 .....	20	2.5.3 数据存储管理 .....	42
1.6 科学研究范式 .....	21	2.6 NewSQL 和 NoSQL.....	43
1.6.1 科学研究范式的产生与 发展 .....	22	2.6.1 NoSQL .....	44

2.6.2	NewSQL .....	48	3.3.3	非结构化数据组织 .....	80
2.6.3	混合应用模式 .....	48	3.3.4	纯文本抽取通用程序库 .....	82
2.7	分布式文件系统 .....	49		本章小结 .....	83
2.7.1	评价指标 .....	49	<b>第4章 数据清洗技术 .....</b>	<b>84</b>	
2.7.2	HDFS 文件系统 .....	50	4.1	数据质量与数据清洗 .....	84
2.7.3	NFS 文件系统 .....	56	4.1.1	数据质量 .....	85
2.7.4	FastDFS .....	57	4.1.2	数据质量提高技术 .....	87
2.8	虚拟存储技术 .....	59	4.1.3	数据清洗算法的标准 .....	90
2.8.1	虚拟存储特点 .....	60	4.1.4	数据清洗的过程与模型 .....	91
2.8.2	虚拟存储的应用 .....	60	4.2	不完整数据清洗 .....	92
2.9	云存储技术 .....	61	4.2.1	基本方法 .....	92
2.9.1	云存储原理 .....	61	4.2.2	基于 k-NN 近邻缺失数据的 填充算法 .....	94
2.9.2	网络结构 .....	61	4.3	异常数据清洗 .....	96
2.9.3	云的分类 .....	62	4.3.1	异常值产生的原因 .....	96
	本章小结 .....	63	4.3.2	统计方法 .....	97
<b>第3章 大数据抽取技术 .....</b>	<b>64</b>		4.3.3	基于邻近度的离群点检测 .....	98
3.1	数据抽取技术概述 .....	64	4.4	重复数据清洗 .....	99
3.1.1	数据抽取的定义 .....	65	4.4.1	使用字段相似度识别重复值 算法 .....	99
3.1.2	数据映射与数据迁移 .....	66	4.4.2	搜索引擎快速去重算法 .....	100
3.1.3	数据抽取程序 .....	66	4.5	文本清洗 .....	100
3.1.4	Kettle 数据处理工具 .....	67	4.5.1	字符串匹配算法 .....	101
3.1.5	数据抽取方式 .....	71	4.5.2	文本相似度度量 .....	103
3.2	增量数据抽取技术 .....	72	4.6	数据清洗技术的实现 .....	107
3.2.1	增量抽取特点与策略 .....	72	4.6.1	数据清洗的步骤 .....	107
3.2.2	基于触发器的增量抽取 方式 .....	72	4.6.2	数据清洗的工具 .....	108
3.2.3	基于时间戳的增量抽取 方式 .....	74		本章小结 .....	108
3.2.4	全表删除插入方式 .....	75	<b>第5章 大数据去噪与标准化 .....</b>	<b>109</b>	
3.2.5	全表比对抽取方式 .....	75	5.1	基本的数据转换方法 .....	109
3.2.6	日志表方式 .....	75	5.1.1	对数转换 .....	109
3.2.7	系统日志分析方式 .....	76	5.1.2	平方根转换 .....	110
3.2.8	各种数据抽取机制的比较与 分析 .....	76	5.1.3	平方转换 .....	110
3.3	非结构化数据抽取 .....	78	5.1.4	倒数变换 .....	110
3.3.1	非结构化数据类型 .....	78	5.2	数据平滑技术 .....	111
3.3.2	非结构化数据模型 .....	78	5.2.1	移动平均法 .....	111
			5.2.2	指数平滑法 .....	115

5.2.3 分箱平滑法 .....	120	6.8 数值约简 .....	140
5.3 数据规范化 .....	121	6.8.1 有参数值约简 .....	140
5.3.1 最小-最大规范化方法 .....	121	6.8.2 无参数值约简 .....	141
5.3.2 z 分数规范化方法 .....	122	6.9 数值离散化与概念分层 .....	142
5.3.3 小数定标规范化方法 .....	122	6.9.1 基于数值属性的概念 分层 .....	142
5.4 数据泛化处理 .....	123	6.9.2 数值数据的离散化 .....	143
5.4.1 空间数据支配泛化算法 .....	123	本章小结 .....	150
5.4.2 非空间数据支配泛化方法 .....	124	<b>第 7 章 大数据集成技术</b> .....	<b>151</b>
5.4.3 统计信息网格方法 .....	124	7.1 数据集成技术概述 .....	152
本章小结 .....	125	7.1.1 数据集成的概念与相关 问题 .....	152
<b>第 6 章 大数据约简技术</b> .....	<b>126</b>	7.1.2 数据集成的核心问题 .....	155
6.1 数据约简概述 .....	126	7.1.3 数据集成的分类 .....	156
6.1.1 数据约简定义 .....	126	7.2 数据迁移 .....	158
6.1.2 数据约简策略 .....	127	7.2.1 在组织内部移动数据 .....	159
6.2 特征约简 .....	127	7.2.2 非结构化数据集成 .....	160
6.2.1 特征提取 .....	128	7.2.3 将处理移动到数据端 .....	161
6.2.2 特征选择 .....	128	7.3 数据集成模式 .....	161
6.2.3 基于主成分分析的特征约简 方法 .....	129	7.3.1 联邦数据库集成模式 .....	162
6.3 样本约简 .....	130	7.3.2 中间件集成模式 .....	163
6.3.1 随机抽样 .....	130	7.3.3 数据仓库集成模式 .....	164
6.3.2 系统抽样 .....	130	7.4 数据集成系统 .....	165
6.3.3 分层抽样 .....	130	7.4.1 全局模式 .....	166
6.4 数据立方体聚集 .....	131	7.4.2 语义映射 .....	166
6.4.1 多维性 .....	131	7.4.3 查询重写 .....	167
6.4.2 数据聚集 .....	132	7.5 数据集成系统的构建 .....	167
6.5 维约简 .....	133	7.5.1 模式之间映射关系的生成 .....	167
6.5.1 维约简的定义 .....	133	7.5.2 适应性查询 .....	168
6.5.2 维约简的分类 .....	134	7.5.3 XML .....	168
6.6 属性子集选择算法 .....	136	7.5.4 P2P 数据管理 .....	168
6.6.1 逐步向前选择属性 .....	136	7.6 数据聚类集成 .....	169
6.6.2 逐步向后删除属性 .....	136	7.6.1 数据聚类集成概述 .....	169
6.6.3 混合式选择 .....	137	7.6.2 高维数据聚类集成 .....	169
6.6.4 判定树归纳 .....	137	7.7 实时数据集成 .....	172
6.7 数据压缩 .....	138	7.7.1 基于中间件层的实时数据 集成模式 .....	172
6.7.1 离散小波变换方法 .....	138		
6.7.2 主要成分分析压缩方法 .....	139		



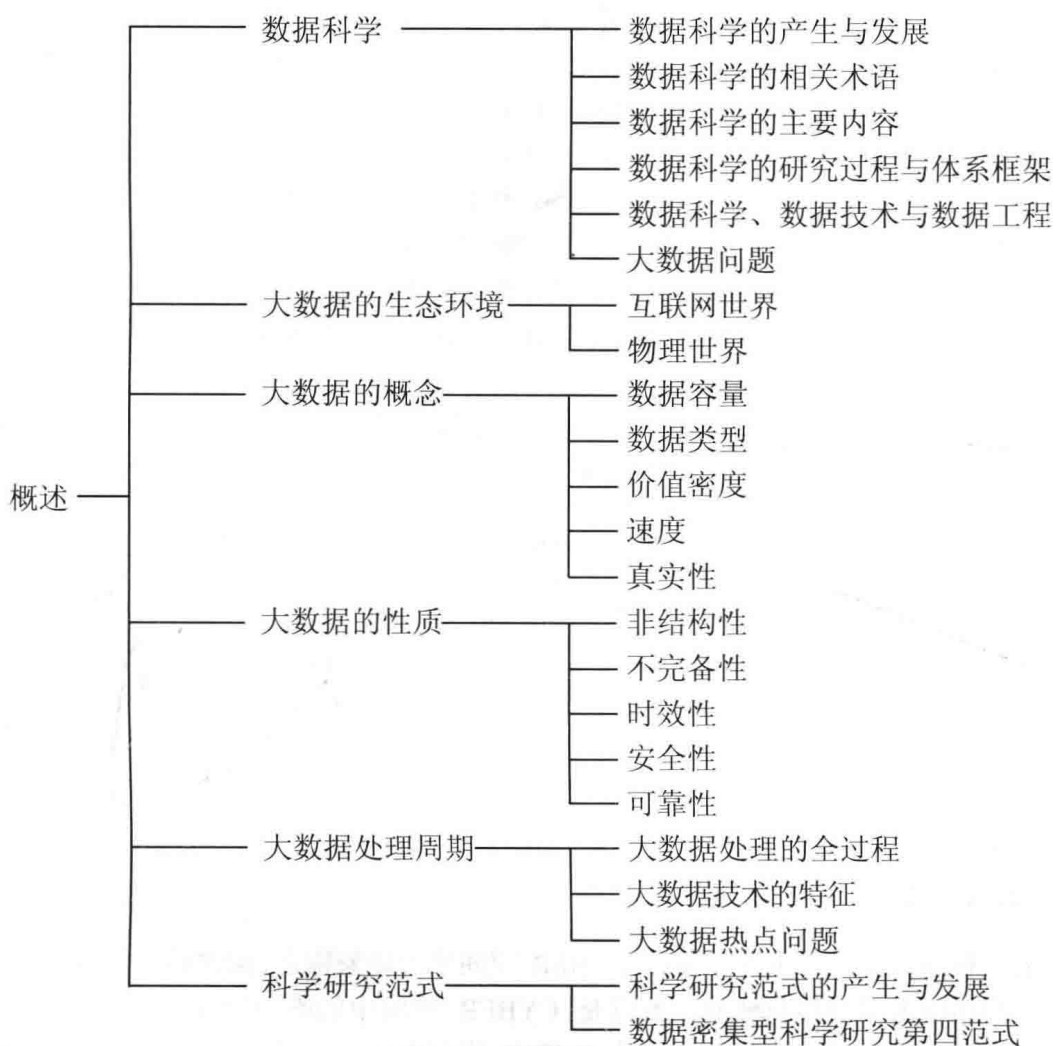
7.7.2 基于数据源层和中间件层的实时数据集成模式 .....	173	8.7.3 文本的自然语言处理 .....	209
7.7.3 基于数据仓库和中间件层的集成模式 .....	174	本章小结 .....	210
7.7.4 基于数据网络的实时数据集成模式 .....	175	<b>第9章 大数据分析</b> .....	<b>211</b>
本章小结 .....	176	9.1 大数据分析定义与方法 .....	211
<b>第8章 大数据挖掘技术</b> .....	<b>177</b>	9.1.1 大数据分析的类型 .....	212
8.1 数据挖掘理论基础 .....	177	9.1.2 数字特征 .....	214
8.1.1 数据挖掘是面向实际应用的 .....	178	9.1.3 统计方法论 .....	217
8.1.2 数据挖掘的理论基础 .....	179	9.1.4 模型与构建 .....	219
8.1.3 基于数据存储方式的数据 .....	180	9.1.5 R语言 .....	221
8.2 关联规则挖掘 .....	182	9.2 统计分析的基本方法 .....	224
8.2.1 经典的频繁项目集生成 .....	183	9.2.1 指标对比分析 .....	224
8.2.2 关联规则挖掘质量 .....	185	9.2.2 分组分析 .....	225
8.3 分类 .....	186	9.2.3 综合评价分析 .....	225
8.3.1 分类定义与分类步骤 .....	186	9.2.4 指数分析 .....	226
8.3.2 基于距离的分类算法 .....	187	9.2.5 平衡分析 .....	226
8.3.3 决策树分类方法 .....	188	9.2.6 趋势分析 .....	227
8.4 聚类方法 .....	191	9.2.7 显著性检验 .....	228
8.4.1 距离与相似性的度量 .....	193	9.2.8 结构分析 .....	231
8.4.2 划分聚类方法 .....	194	9.2.9 因素分析 .....	231
8.4.3 层次聚类方法 .....	196	9.2.10 交叉分析 .....	232
8.5 序列模式挖掘 .....	196	9.3 高级数据分析方法 .....	232
8.5.1 时间序列预测的常用方法 .....	197	9.3.1 动态分析法 .....	232
8.5.2 序列模式挖掘 .....	198	9.3.2 相关分析 .....	233
8.6 Web挖掘技术 .....	200	9.3.3 回归分析 .....	236
8.6.1 Web内容挖掘方法 .....	200	9.3.4 判别分析 .....	240
8.6.2 Web访问信息挖掘方法 .....	202	9.3.5 对应分析 .....	243
8.6.3 Web结构挖掘方法 .....	204	9.3.6 主成分分析 .....	244
8.7 非结构化文本数据挖掘 .....	206	9.3.7 多维尺度分析 .....	245
8.7.1 用户反馈文本 .....	206	9.3.8 方差分析 .....	250
8.7.2 用户反馈文本挖掘的一般 .....	207	本章小结 .....	252
		<b>第10章 分析结果的解释</b> .....	<b>253</b>
		10.1 分析结果的可视化解释 .....	253
		10.1.1 解释的目的与主要内容 .....	254
		10.1.2 检查和验证假设 .....	254
		10.1.3 追踪分析过程 .....	254
		10.2 基本展现方式 .....	255

10.2.1 基于时间变化的可视化 展现 .....	256	11.3.1 大数据的空气质量推断.....	298
10.2.2 由大及小的可视化展现 .....	256	11.3.2 人与建筑的关系分析 .....	299
10.2.3 由小及大的可视化展现 .....	256	11.3.3 针对全球问题的预测模型 ..	299
10.2.4 突出对比的可视化展现 .....	256	11.3.4 地表可视化与数据分析.....	299
10.2.5 地域空间可视化展现 .....	258	本章小结 .....	300
10.2.6 概念可视化展现 .....	260	<b>第 12 章 大数据离线计算技术 .....</b>	<b>301</b>
10.2.7 气泡图可视化展现 .....	261	12.1 数据离线计算概述 .....	301
10.2.8 注重交叉点的数据可视化 展现 .....	262	12.1.1 大数据离线处理特点 .....	301
10.2.9 剖析原因的数据可视化 展现 .....	262	12.1.2 批量计算 .....	302
10.2.10 异常值数据可视化展现 ..	262	12.2 MapReduce 的体系结构 .....	302
<b>10.3 大数据中的常用可视化展现 ..</b>	<b>262</b>	12.2.1 MapReduce 计算描述 .....	302
10.3.1 文本可视化 .....	263	12.2.2 MapReduce 适用的场景 .....	304
10.3.2 网络(图)可视化 .....	265	12.3 Hadoop 分布式计算平台 .....	304
10.3.3 时空数据可视化 .....	268	12.3.1 Hadoop 结构与特点 .....	305
10.3.4 多维数据可视化 .....	269	12.3.2 分布式系统与 Hadoop .....	309
10.3.5 基于 ECharts.js 可视化 工具 .....	271	12.3.3 SQL 数据库系统与 Hadoop .....	309
<b>10.4 大数据可视分析 .....</b>	<b>273</b>	12.3.4 基于 Hadoop 框架的分布 计算 .....	311
10.4.1 可视分析的理论基础 .....	274	12.4 MapReduce 程序设计实例 分析 .....	316
10.4.2 大数据可视分析技术 .....	279	12.4.1 单词计数 .....	316
本章小结 .....	282	12.4.2 MapReduce 的应用 .....	319
<b>第 11 章 大数据机器学习 .....</b>	<b>283</b>	本章小结 .....	321
11.1 机器学习概述 .....	283	<b>第 13 章 大数据流式计算技术 .....</b>	<b>322</b>
11.1.1 机器学习的产生与发展 .....	283	13.1 流式数据的概述 .....	323
11.1.2 机器学习类型 .....	286	13.1.1 流式数据的概念 .....	323
11.1.3 知识表示形式 .....	289	13.1.2 流式数据源 .....	324
11.1.4 机器学习的典型算法 .....	291	13.1.3 流式数据的特征 .....	325
11.2 大数据机器学习的特点与 算法 .....	293	13.2 大数据的计算模式 .....	326
11.2.1 大数据机器学习的特点 .....	294	13.2.1 大数据流式计算模型 .....	327
11.2.2 大数据机器学习的评测 指标 .....	295	13.2.2 流式计算与批量计算的 比较 .....	327
11.2.3 大数据机器学习算法 .....	296	13.2.3 流式计算与实时计算的 比较 .....	329
11.3 大数据机器学习的应用 .....	298	13.3 流式计算技术的应用 .....	329

13.3.1 中间计算 .....	329	13.5.1 被动等待策略 .....	338
13.3.2 流式查询 .....	329	13.5.2 主动等待策略 .....	339
13.3.3 流式抽样 .....	330	13.5.3 上游备份策略 .....	339
13.3.4 统计独立元素数 .....	331	13.6 Storm 流处理过程 .....	340
13.3.5 去重计数 .....	332	13.6.1 Storm 特点与架构 .....	340
13.4 流式计算的系统架构 .....	335	13.6.2 topology .....	343
13.4.1 对称式系统架构 .....	335	13.6.3 单词计数 topology .....	346
13.4.2 主从式系统架构 .....	336	13.7 大数据流式计算的应用 .....	347
13.4.3 数据传输方式 .....	337	本章小结 .....	349
13.4.4 编程接口 .....	338	参考文献 .....	350
13.5 高可用技术 .....	338		

# 第1章 概述

## 本章主要内容



世界已经进入了一个模型和假设逐渐清晰的大数据时代，计算机科学是算法与算法变换的科学，数据科学研究范围更为广泛。数据科学研究与进展不仅可以推动数学、计算机科学、统计学、天体信息学、生物信息学、计算社会学等学科的发展，而且还能够大力助推产业发展与技术进步。

## 1.1 数据科学

数据科学是关于数据的科学，数据科学是为研究探索 CYBER 空间中数据界的理论、方法和技术。

### 1.1.1 数据科学的产生与发展

数据科学出现在 20 世纪 60 年代，1974 年彼得·诺尔出版的《计算机方法的简明调查》中将数据科学定义为“处理数据的科学，一旦数据与其代表事物的关系被建立起来，将为其其他领域与科学提供借鉴”。1996 年在日本召开的数据科学、分类和相关方法会议上，将数据科学作为会议的主题词。2001 年美国统计学教授威廉·S.克利夫兰发表了《数据科学：拓展统计学的技术领域的行动计划》，首次将数据科学作为一个单独的学科，并把数据科学定义为统计学领域扩展到以数据作为现金计算对象相结合的部分，奠定了数据科学的理论基础。

### 1.1.2 数据科学的相关术语

#### 1. CYBER 空间

CYBER 空间意译为异次元空间、多维信息空间、计算机空间、网络空间等。其本意是指以计算机技术、现代通信网络技术、虚拟现实技术等信息技术的综合运用为基础，以知识和信息为内容的新型空间，是人类运用知识创造的人工世界并用于知识交流的虚拟空间。信息化是将现实世界中的事物和现象以数据的形式存储到 CYBER 空间中，是一个数据生产的过程。数据是自然和生命的一种表示形式，记录了人类的行为，包括工作、生活和社会的发展。

#### 2. 数据爆炸

将数据快速大量地产生并存储在 CYBER 空间中的现象称之为数据爆炸，数据爆炸在 CYBER 空间中形成数据自然界。数据是 CYBER 空间中的唯一存在，需要研究和探索 CYBER 空间中数据的规律和现象。探索 CYBER 空间中数据的规律和现象就是探索宇宙的规律、探索生命的规律、寻找人类行为的规律、寻找社会发展的规律的一种重要手段。

#### 3. 数据科学的定义

数据科学是关于数据的科学或者研究数据的科学、探索 CYBER 空间中数据界奥秘的理论、方法和技术，研究的对象是数据界中的数据。与自然科学和社会科学不同，数据科学的研究对象是 CYBER 空间的数据，是新的科学。数据科学主要包括两个方面：一个是研究数据本身，研究数据的各种类型、状态、属性及变化形式和变化规律；另一个是为自

然科学和社会科学研究提供一种新的方法，称为科学研究的数据方法，其目的在于揭示自然界和人类行为现象和规律。

#### 4. 数据科学的方法和技术

数据科学采用收集数据的形式，进行开放式分析，不做预先假定。在许多数据科学研究项目中，首先需要浏览原始数据，形成一个假定，然后基于假定进行调查确认。数据科学的关键概念是，数据科学是一个经验科学，直接基于数据进行科学处理。数据科学已经有一些方法和技术，如数据获取、数据存储与管理、数据安全、数据分析、可视化等。其基础理论和新技术主要有数据存在性、数据测度、时间、数据代数、数据相似性与簇论、数据分类与数据百科全书、数据伪装与识别、数据实验、数据感知等。数据学的理论和方法将改进现有的科学研究方法，形成新型的科学研究方法，并且针对各个研究领域开发出专门的理论、技术和方法，从而形成专门领域的数据学，如行为数据学、生命数据学、脑数据学、气象数据学、金融数据学、地理数据学和过程数据学等。

数据科学不仅完成分析，还要涉及整个端到端的生命周期，数据系统本质上是用于研发真实世界理解模型的科学设备。这就表明必须深刻理解数据的来源、数据转换的适用性和准确性、转换算法和过程之间的相互作用以及数据存储机制。进而能够确保所有事物都正确执行，从而探索数据、创建并验证各项科学假设。

### 1.1.3 数据科学的主要内容

数据科学的内容主要包括基础理论研究、数据技术、应用研究及数据科学的学科体系研究，如图 1-1 所示。数据科学学科建立，需要完成知识结构、课程设置、专业设置等学科体系建设，探讨数据科学与自然科学和社会科学之间的关系，数据科学与计算机科学和信息科学之间的关系等。

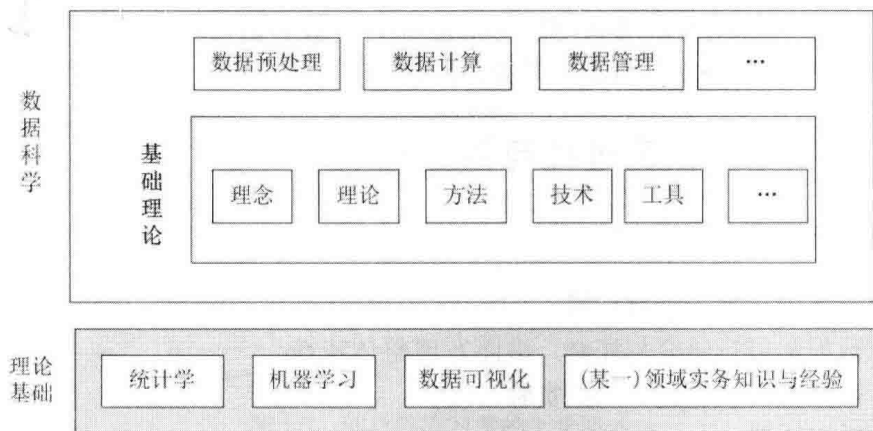


图 1-1 数据科学的内容

数据科学包括用数据的方法研究科学技术和用科学的方法研究数据的两方面。前者包括生物信息学、天体信息学和数字地球等领域，后者包括统计学、机器学习、数据挖掘等领域，这些学科都是数据科学的重要组成部分。

## 1. 用科学的方法研究数据

比较常用的数据有表格、点集、时间序列、图像、视频、网页和网络数据等。基本的数据结构包括度量结构、网络结构和代数结构。

用科学的方法研究数据，主要包括数据处理周期的所有技术，其核心问题是数据分析。数据分析的基本问题是找出模型，由于数据中含有噪声导致了随机模型的出现。为了减少求解的难度，可以找到模型一部分，或将随机模型简化为确定性模型。观察和逻辑推理是科学的基础，数据自然界中主要采用了观察方法与数据推理的理论和方法，包括数据的存在性、数据测度、时间、数据代数、数据分类、数据相似性与簇论等。

## 2. 用数据的方法研究科学

用数据的方法研究科学问题并不是不需要模型，只是模型的出发点不同，不是从基本原理的角度去寻找模型，通常是基于更为简单的数学模型。

需要建立数据科学的实验方法，需要提出科学假说和建立理论体系，并通过这些实验方法和理论体系进行数据自然界的探索研究，从而掌握数据的各种类型、状态、属性、变化形式和变化规律，揭示自然界和人类行为现象和规律。

## 3. 领域数据学

将数据科学的理论和方法广泛应用，开发出专门的理论、技术和方法，从而形成专门领域的数据科学，例如，脑数据学、行为数据学、生物数据学、气象数据学、金融数据学和地理数据学等。

## 4. 数据资源的开发方法和技术

数据资源是重要的现代战略资源，具有巨大的价值，越来越凸显重要性，是继石油、煤炭、矿产等传统资源之后的最重要的资源之一。人类的社会、政治和经济都将依赖于数据资源，而石油、煤炭、矿产等传统资源的勘探、开采、运输、加工、产品销售等都依赖数据资源，离开了数据资源，将无法开展与完成这些工作。

### 1.1.4 数据科学的研究过程与体系框架

#### 1. 数据科学的研究过程

- ① 从自然界中获得一个数据集。
- ② 对该数据集进行研究与探索，进而发现整体特性。
- ③ 进行数据分析或者进行数据实验。
- ④ 发现数据规律。
- ⑤ 将数据进行感知化与可视化等。

#### 2. 数据科学的构成

数据科学的组成要素可以从图 1-2 所示的维恩图中得到线索，主要包括计算机科学、

数学、统计学知识，以及领域的专业知识。

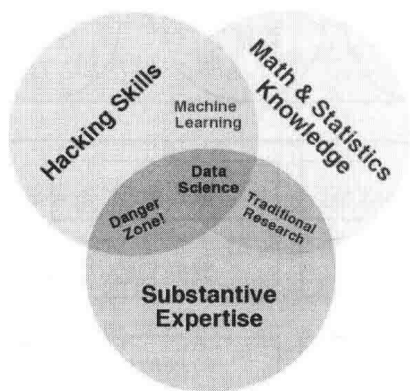


图 1-2 数据科学的组成要素

- Hacking Skills: 编程技巧。
  - Danger Zone: 危险地带。
  - Machine Learning: 机器学习。
  - Data Science: 数据科学。
  - Traditional Research: 传统研究。
  - Substantive Expertise: 行业知识。
  - Math & Statistics Knowledge: 数学与统计学知识。
- 进一步细化数据科学的 12 个主要领域，如图 1-3 所示。

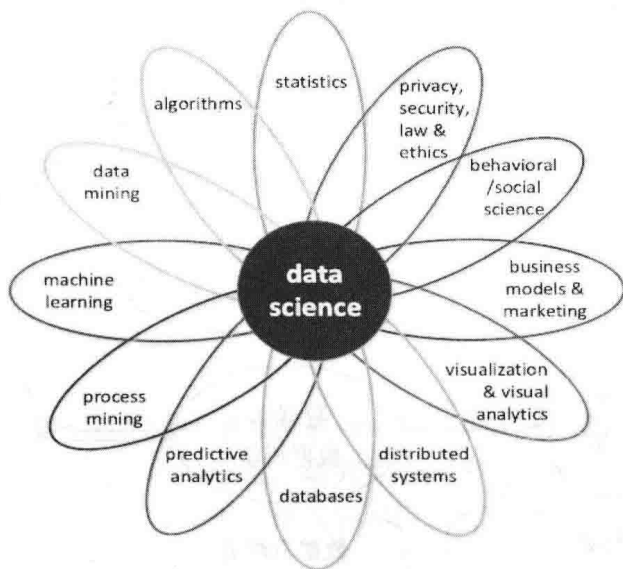


图 1-3 细化数据科学的主要领域

### 3. 数据科学的体系框架

数据科学的体系框架如图 1-4 所示。在图的上部描述了数据的内容，下部分是数据科学的基础描述。



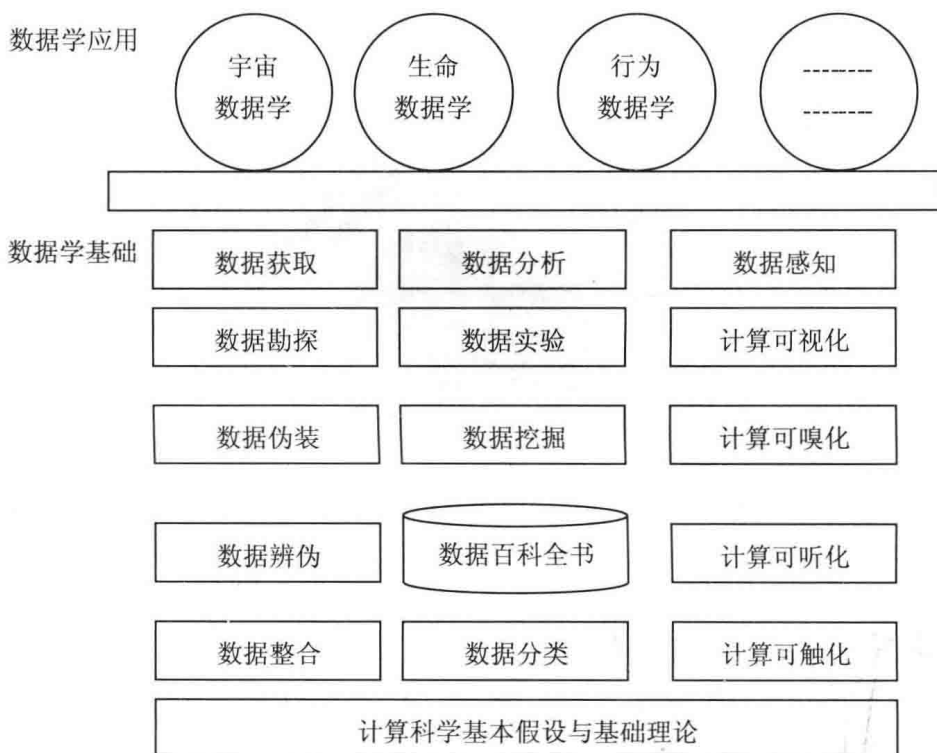


图 1-4 数据科学体系框架

(1) 研究从数据中获取信息与知识

数据科学的研究对象、研究目的与研究方法等与计算科学、信息科学不同。数据存于 CYBER 空间中，信息是自然界、人类社会及人类思维活动中存在和发生的现象，知识是人们在实践中所获得的认识和经验。数据可以作为信息和知识符号的表示或载体，但数据本身并不是信息或知识。数据科学的研究对象是数据，而不是信息，也不是知识。通过研究数据来获取对自然、生命和行为的认识，进而获得重要信息与知识。数据科学与其他学科的关系如图 1-5 所示。

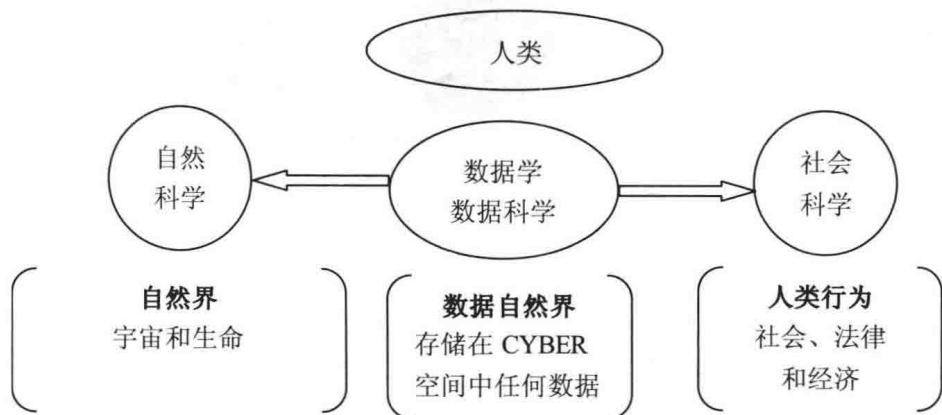


图 1-5 数据科学与其他学科的关系

(2) 通过分析数据来认识自然和行为

自然科学研究自然现象和规律，认识的对象是整个自然界物质的各种类型、状态、属