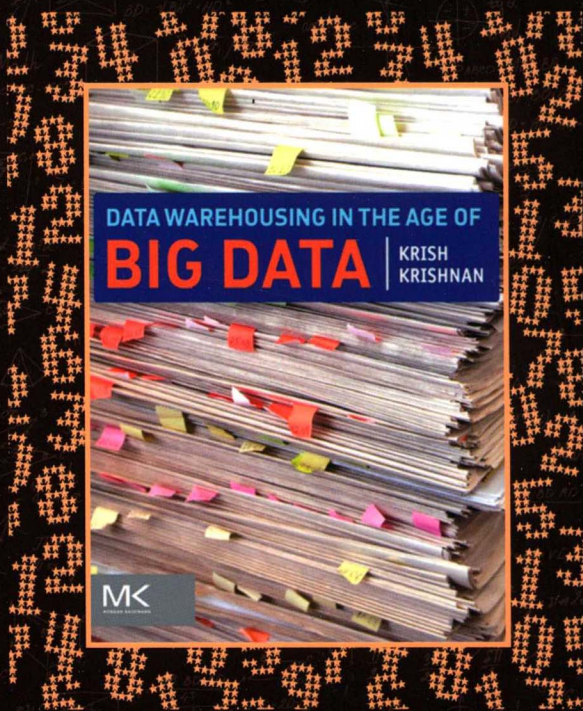


# 大数据与数据仓库

## 集成、架构与管理

[美] 克里什·克里希南 (Krish Krishnan) 著

邢春晓 张勇 张桂刚 译



DATA WAREHOUSING IN THE AGE OF  
**BIG DATA**

数据科学与工程技术丛书

ISBN 978-7-111-59482-6  
I·大·11.01.01  
中国版本图书馆CIP数据核字(2018)第026262号  
本书版权登记：图字：01-2013-1854

DATA WAREHOUSING IN THE AGE OF  
BIG DATA

# 大数据与数据仓库

## 集成、架构与管理

常州大学图书馆

[美] 克里希南 (Krishna) 著

藏书章

张勇 张桂刚 译

大数据与数据仓库：集成、架构与管理

ISBN 978-7-111-59482-6  
I·大·11.01.01  
中国版本图书馆CIP数据核字(2018)第026262号  
本书版权登记：图字：01-2013-1854



机械工业出版社  
China Machine Press

## 图书在版编目 (CIP) 数据

大数据与数据仓库：集成、架构与管理 / (美) 克里什·克里希南 (Krish Krishnan) 著；邢春晓，张勇，张桂刚译. —北京：机械工业出版社，2018.4

(数据科学与工程丛书)

书名原文：Data Warehousing in the Age of Big Data

ISBN 978-7-111-59482-6

I. 大… II. ① 克… ② 邢… ③ 张… ④ 张… III. 数据库系统 IV. TP311.13

中国版本图书馆 CIP 数据核字 (2018) 第 056562 号

本书版权登记号：图字 01-2013-7854

Elsevier

Elsevier (Singapore) Pte Ltd.

3 Killiney Road, #08-01 Winsland House I, Singapore 239519

Tel: (65) 6349-0200; Fax: (65) 6733-1817

Data Warehousing in the Age of Big Data

Krish Krishnan

Copyright © 2013 Elsevier Inc. All rights reserved.

ISBN-13: 978-0-12-405891-0

This translation of Data Warehousing in the Age of Big Data by Krish Krishnan was undertaken by China Machine Press and is published by arrangement with Elsevier (Singapore) Pte Ltd.

Data Warehousing in the Age of Big Data by Krish Krishnan 由机械工业出版社进行翻译，并根据机械工业出版社与爱思唯尔（新加坡）私人有限公司的协议约定出版。

《大数据与数据仓库：集成、架构与管理》(邢春晓 张勇 张桂刚 译)

ISBN: 978-7-111-59482-6

Copyright © 2018 by Elsevier (Singapore) Pte Ltd.

All rights reserved. No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from Elsevier (Singapore) Pte Ltd. Details on how to seek permission, further information about the Elsevier's permissions policies and arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: [www.elsevier.com/permissions](http://www.elsevier.com/permissions).

This book and the individual contributions contained in it are protected under copyright by Elsevier (Singapore) Pte Ltd. and China Machine Press (other than as may be noted herein).

### 注意

本译本由 Elsevier (Singapore) Pte Ltd. 和机械工业出版社完成。相关从业及研究人员必须凭借其自身经验和知识对文中描述的信息数据、方法策略、搭配组合、实验操作进行评估和使用。由于医学科学发展迅速，临床诊断和给药剂量尤其需要经过独立验证。在法律允许的最大范围内，爱思唯尔、译文的原文作者、原文编辑及原文内容提供者均不对译文或因产品责任、疏忽或其他操作造成的人身及 / 或财产伤害及 / 或损失承担责任，亦不对由于使用文中提到的方法、产品、说明或思想而导致的人身及 / 或财产伤害及 / 或损失承担责任。

Printed in China by China Machine Press under special arrangement with Elsevier (Singapore) Pte Ltd. This edition is authorized for sale in the People's Republic of China only, excluding Hong Kong SAR, Macau SAR and Taiwan. Unauthorized export of this edition is a violation of the contract.

本书封底贴有 Elsevier 防伪标签，无标签者不得销售。

## 大数据与数据仓库：集成、架构与管理

出版发行：机械工业出版社（北京市西城区百万庄大街 22 号 邮政编码：100037）

责任编辑：谢晓芳

印刷：北京市荣盛彩色印刷有限公司

开本：185mm×260mm 1/16

书号：ISBN 978-7-111-59482-6

责任校对：李秋荣

版次：2018 年 4 月第 1 版第 1 次印刷

印张：17.75

定价：79.00 元

凡购本书，如有缺页、倒页、脱页，由本社发行部调换

客服热线：(010) 88379426 88361066

购书热线：(010) 68326294 88379649 68995259

投稿热线：(010) 88379604

读者信箱：hzit@hzbook.com

版权所有·侵权必究

封底无防伪标均为盗版

本书法律顾问：北京大成律师事务所 韩光 / 邹晓东

## 译者序

数据仓库最早用于为企业决策提供所有类型数据支撑的大型数据集。随着大数据时代的到来，数据量越来越大，对数据的处理速度越来越快，同时数据的价值密度也变得越来越小。大数据时代下具有 4V 特征的数据环境中，数据仓库又会变得怎样呢？本书为回答这一问题做了详细的分析和呈现。

目前对数据仓库进行论述的著作都没有明确指出大数据这一特殊时代的特点以及面临的问题和相应的解决方法，而本书将数据仓库与大数据两者进行了有机融合与统一，从而使得数据仓库技术符合大数据这一特定环境的需求。同时，也为新型企业基于大数据 - 数据仓库的管理决策提供了有效的途径。

本书逻辑清晰，内容较为全面，具有很强的适用性。全书由三大部分组成，分别为大数据、数据仓库以及构建大数据 - 数据仓库。第一部分主要包含：大数据简介，使用大数据，大数据处理架构，大数据技术简介，以及大数据驱动的商业价值。第二部分主要包括：再论数据仓库，数据仓库的再造，数据仓库中的工作负载管理，以及应用到数据仓库的新技术。第三部分主要包括：大数据和数据仓库的集成，大数据的数据驱动架构，大数据的信息管理和生命周期，大数据分析、可视化和数据科学家，以及实施大数据 - 数据仓库的现实情况。值得一提的是，本书还对客户案例研究与建设医疗保健信息工厂两个案例做了分析，有助于读者更好地理解本书。

本书的翻译工作主要由如下人员完成。清华大学信息技术研究院邢春晓研究员负责译稿的审校工作，清华大学信息技术研究院张勇副研究员和中国科学院自动化研究所张桂刚副教授负责本书的翻译工作。

# 前言

Web 2.0 改变了我们的生活和工作方式，比如开展业务、与客户沟通、与朋友和家人共享信息、用业务收入和客户花销份额来衡量成功，以及定义品牌管理。最重要的是，它创造了一种独一无二的生财之道。无论是安排度假地点、购买最新型的电视、更换移动服务供应商，还是想要为郊游买最好的食材，你都可以通过互联网查看顾客的评论和读者的推荐。同样，在个人生活中，你可以使用 Facebook、YouTube、iTunes、Instagram 和 Flickr 分享你喜欢的音乐、电影、照片和视频。

当今，企业所提供的产品和服务的个性化为消费者创造了许多机会，同时也大大促进了数据量增大、数据格式（品种）增多和数据生产速度加快。数据的关键价值是，当我们使用地理和人口学数据建模来创建关于相似人群的个性、行为和影响的聚类时，能够找到在数据中隐含的智慧。

向服务的个性化和以客户为中心的商业模式进行转变形成了三个不同的趋势。

- **众包。**这是 Jeff Howe 于 2006 年在《连线》杂志上提出的术语。众包是在当今世界使用协同智能研究人类行为的过程。信息管理和个人层次上的推荐共享共同形成了业界的趋势。
  - 众包已演变成一个强有力的工具。它现在在商业上有很多用途，例如寻找有竞争力的研究、客户情感分析和因果分析等。同时还部署了其他的分析模型，例如协作过滤、推荐和机器学习算法。
  - 众包的最佳案例之一是当时身为参议员的奥巴马在 2008 年的总统候选人提名竞选中筹款。通过使用互联网和社交媒体作为一种个性化的联系渠道，他在筹措资金方面明显超过了其他候选人，从而能够进行有效的竞选。
- **社交媒体分析。**今天的消费者依靠的数据和信息是通过社交媒体渠道获得的，而这些数据和信息又依赖于将这个平台作为其“个人决策支持平台”的广大用户所做出的个人决策。这使得更多的人利用社交媒体作为与客户、合作伙伴和供应商直接和间接的沟通渠道。今天，如果你没有使用社交媒体，那么你是过时的，尤其是与 90 后和新千年的客户群相比。
  - 如何度量你的社交媒体渠道和沟通策略的有效性？这表明你从哪里开始实施一项社交媒体分析战略。该战略应从两个角度进行度量，包括从内向外和从外向内。在这一领域一个企业的成熟和演变往往需要经过多个阶段。在现在的新闻和互联网上，你会发现一些使用该策略成功进行业务变革的例子。
- **游戏化。**今天的另一个热门趋势是在企业内外使用游戏化策略来吸引员工、客户、潜在客户和任何对你的业务和服务感兴趣的人。

○ 游戏化基于博弈理论和统计模型的组合，在对“长尾”得到的结果进行建模时，这已成为一种非常有效的工具。它也是在 Web 2.0 时代由 Chris Andersson 定义的术语，他还专门就这个主题写了一本书。

○ 这一策略最好的例子是在 2012 年美国总统大选中，竞选策略专家使用博弈论和统计模型寻找目标选民，而且非常有效。奥巴马总统的竞选专门使用这种组合作为一种有效的和颠覆性的策略，从个人层次在候选人和选民之间创建了很多需要的连接。

从 Web 2.0 的观点来看，所有趋势、理论和成果的共同主线可以归结为两点。

- 使用 Web 2.0 平台所需要的数据量远远大于现今企业所用到的。
- 在计算的历史上，使用统计模型和分析的需求比以往更加强烈。

这两个事实已经被 Facebook、Groupon、Google、Yahoo、Apple 和其他财富 500 强公司证明是成功的。

数据带来了如下问题：如何计算海量和多样的数据，以及如何应对数据体量。这是 Google、Facebook 和 Yahoo 清楚展示的方式；前者创造了一种新的计算模型，该模型基于文件系统和一种叫作 MapReduce 的编程语言。MapReduce 扩展了搜索引擎的能力，能够同时处理多个查询。2002 年，架构师 Doug Cutting 和 Mike Carafella 正在做开源搜索引擎项目 Nutch，这促使他们基于 Google 模型来对底层架构进行建模。这也使得 Nutch 成为一个开源的顶级 Apache 项目。该项目于 2006 年被 Yahoo 所采用，称为 Hadoop。在过去的几年中 Hadoop 成就了大量的公司，这些公司有商业化的解决方案，同时将相应功能回馈给基础的开源项目，这是一种真正基于协作的软件和框架开发。

另外一项技术也演化为一个强大的平台，即 NoSQL (Not only SQL) 运动。该平台基于 Eric Brewer 在 2002 年提出的 CAP 定理。根据 CAP 定理，一个数据库不能在任何一个时间点满足 ACID 兼容的所有规则，同时又是可扩展的和灵活的。然而，在一致性、可用性和分区容忍性三个基本性质中，一个数据库可以满足三个性质中的两个，从而创建可扩展的分布式结构，该架构可以演变成满足水平方向上缩放的可扩展性要求并提供更高的吞吐量，因为在这种环境中计算和存储是非常接近的，同时是一个允许多种一致性级别的分布式架构。

Facebook 是 NoSQL 架构的最早提倡者之一，因为他们要解决用户的可扩展性和可用性要求，其用户量仅次于中国和印度的人口。Cassandra 是一个流行的数据库，在 Facebook 经历了很长时间的开发和使用（现在由于更大的可扩展性需求，它已经被 Facebook 抛弃）。许多其他公司把它与 Hadoop 以及其他传统的 RDBMS 解决方案一起使用。它仍然是一个顶级的 Apache 项目，并且正在添加更多的功能。

随着这些新技术和商业模式的出现，也出现了大量噪音，并导致了混乱。这些趋势或噪音之一是“数据仓库的死亡”，这在全球都带来了严重影响，因为企业已不只投入数以百万计美元来搭建这种决策支持平台，而且基于其输出开发了若干下游系统。

作为传统的数据仓库领域和大数据领域中都有经验的数据实践者和咨询师，我开始在数据仓储研究所 (TDWI) 教授课程，在许多国际峰会和其他会议中谈论大数据和数据仓库，以消除数据仓库的“死亡”所带来的恐怖。在过去四年中，在全球关于这个话题展开了大量讨论之后，我决定写这本书并讨论大数据。包括谁使用大数据，它是如何

影响数据仓库世界的，以及数据分析的未来，更重要的是，下一代数据库仓库的概念以及它是如何构建的。

坦白地说，我们将继续构建和使用数据仓库，而且它将仍然是“单一版本的事实”，但我们不再使用 RDBMS 作为数据仓库和分析的平台。在写这本书的时候，我们看到每隔几个月，有时是几周，Hadoop、MapReduce 和 NoSQL 就会发生变化，新功能就会浮出水面。人们正在设计和搭建这些架构，它们可以处理大型和复杂的数据，能够在批处理环境中有效处理数据，但是比起关系数据库管理系统在实时和交互能力方面比较有限。该架构的最终状态将是这些架构的异构组合，以共同创建一个强大和巨大的决策支持架构，这个系统的名称依然是数据仓库。

在读这本书时，你会发现三个不同的部分。第一部分讨论大数据，包括大数据技术及来自早期实践者的用例。第二部分介绍数据仓库、它失败的原因、新的架构选项、工作负载、工作负载驱动的架构，以及大数据和数据仓库的集成技术。第三部分涉及数据治理、数据可视化、信息生命周期管理、数据科学家，以及适合大数据的数据仓库。附录包括来自供应商的实现和一个关于如何建立医疗保健信息工厂的案例研究。

本书的总体目标是帮助你了解大数据和数据仓库的复杂层次，同时为你提供关于如何有效使用所有这些技术和架构来设计下一代数据仓库的信息。

下面描述各章的内容和全书组织结构，为你提供阅读路线图。在逐章阅读时，这些章节结合起来就会为你提供简洁而深入的理解。

## 第一部分：大数据

第 1 章的重点是让你彻底理解大数据。我们避免使用流行词，探讨了新兴的大数据领域和它对企业的重要性。

第 2 章的重点是大数据隐含的复杂之处（即三个 V——体量、速度以及多样性和多义性），如何处理这些特点，以及在哪些主题域有哪些隐藏的陷阱。

第 3 章重点讨论需要或者设计什么架构以进行大数据处理，还讨论了算法级的方法、分类系统、集群和其他内容。

第 4 章重点讨论的是为解决大规模数据处理，核心技术是如何演化的。这些技术包括 Hadoop 及其生态系统、NoSQL 数据库和其他技术。这一章对于这些技术的介绍是极其浓缩的，建议你进一步阅读有关这些主题的核心书目。

第 5 章论述在现实生活中不同公司利用大数据实现价值的各种用例。这些用例涉及 B2B、B2C、C2C 等场景，该章还介绍在每个场景中是如何定义和实现价值的。

## 第二部分：数据仓库

第 6 章重点追溯数据库仓库的起源以及这些年来的演化。该章讨论早期版本的缺陷所导致的数据仓库的失败，以及如何识别和避免这些缺陷。

第 7 章主要介绍如何以及为什么要现代化数据仓库架构。这将为提供概念上的思想以及实现上的一些选项。

第 8 章重点介绍工作负载，及其在数据仓库和大数据领域中的真正含义，理解工作负载的重要性，以及基于工作负载如何创建数据仓库的架构。对于任何数据管理解决方案来说，这都是其未来架构最重要和最关键的方面之一。

第 9 章重点讨论那些已持续应用到企业中的新兴技术，特别是在处理数据库仓库的性能和可扩展性方面。该章还讨论数据仓库一体机、云计算、数据虚拟化和内存计算。

## 第三部分：构建大数据 – 数据仓库

第 10 章重点介绍将数据仓库与大数据集成的方法和相关的技术，这些技术的采用基于公司的数据类型、当前演化状态和现有技术。

第 11 章重点讨论在大数据领域中通过部署有效的 MDM 和元数据策略来创建数据驱动的架构。它强调对数据管理的这两大支柱的需求，特别是在大数据领域。该章还讨论语义层和基于语义网的方法。

第 12 章的重点是管理大数据的生命周期，包括哪些数据是基本的，在处理前和处理后如何以及在哪里保存数据。还将讨论企业大数据中如果不实现一个鲁棒的 ILM 策略会带来哪些问题。

第 13 章涉及使用大数据的最终目标，也就是提供强大的可视化，分析大数据，最重要的是，新兴的数据科学家的角色。这里的目标是为你提供关于这些主题的概念性的想法以及它们如何影响整体的大数据策略。

第 14 章着重介绍在财富 500 强企业的下一代数据仓库的实际实施中的最终架构。目的是当你的企业演化到新的数据领域后，为你提供一些面向未来的想法。

## 附录

附录 A 展示具体的客户案例研究。

附录 B 给出建设医疗保健信息工厂的案例研究。



## 致 谢

本书的出版离不开太多人的支持，我要感谢他们在本书的出版过程中给予的支持和帮助。

首先，我要感谢我的妻子和两个儿子一直以来的支持，他们牺牲了很多周末、假期、看电影、参加学校活动和社会活动的时间（这些我都缺席了），给我提供了考虑如何组织本书内容的黄金时间。没有他们的帮助和支持，我将永远没有机会写作这本书。

接下来，感谢我亲爱的朋友、导师兼商业伙伴 **Bill Inmon**，他的支持和鼓励让我坚持在做一份全职咨询工作的同时写这本书——工作在一家初创公司的同时，又做着无数的事情。谢谢 **Bill** 总是及时提供帮助。

我要特别感谢三个最好的朋友（我很幸运，有这些家伙作为我的朋友）：**Todd Nash**、**Hans Hultgren** 和 **Shankar Radhakrishnan**。他们花了数小时的时间（晚上、周末和在飞机上）审查草稿，并对每一步都提供反馈。没有他们不懈的努力，我们将不能及时完成这本书。他们令我倍感惊喜，他们是这段令人难以置信的旅程中无法用语言描述的一部分。他们是令人难以置信的，很高兴能有这样的朋友，与他们在一起的时光总是很特别。

仅仅有一堆理论是没法写成一本书的——本书获得了很多业内资深人士以供应商支持的方式提供的帮助。**IBM** 的 **Glenn Zimmerman**，他对我的所有案例研究都非常支持，甚至帮我把它们改写成 **Word** 格式。**Teradata** 的 **Kim Dossey**，谢谢他提供的所有支持并快速提供了许多案例研究。**Cloudera** 的 **Alan** 和 **Kate**，他们两个毫不犹豫地让我选择尽可能多的案例进行研究。**Klout** 的案例研究要感谢微软的 **SQL Server** 团队。我想郑重地感谢所有这些供应商的鼎力支持和热心帮助。

我还要感谢在这个征途中以鼓励的言语来支持我的一些朋友：**Paul Kautza**、**Philip Russom**、**Dave Stodder**、**Dave Wells**、**Claudia Imhoff**、**Jill Dyche**、**Mark Madsen**、**Jonathan Seidman**、**Kylie Clement**、**Dave Nielsen**、**Shawn Rogers**、**John Myers**、**John O'Brien**、**William McKinght**、**Robert Eve**、**Tony Shaw**，以及 **John Onder**。

最后一点也很重要，感谢 **Andrea Dierna**、**Heather Scherer** 编辑和整个 **Morgan Kauffman** 团队在这一过程中的所有帮助、指导和支持——如果没有他们，这一切都不可能实现。

谢谢各位！

第1章 大数据技术简介 ..... 31

1.1 大数据 ..... 32

1.2 大数据的5V特性 ..... 34

1.3 大数据的存储 ..... 34

1.4 大数据的传输 ..... 35

1.5 大数据的治理 ..... 36

1.6 大数据的隐私 ..... 36

**Krish Krishnan** 是高性能数据仓库解决方案和非结构化数据的战略、架构和实现等方面全球公认的专家。他是一位很受欢迎且有远见的数据仓库思想领导者和实践者，是世界顶尖的战略和架构咨询师之一。**Krish** 也是全世界大数据相关会议上独立的分析人员和演讲者，且在数据仓储研究所 (TDWI) 讲授这一主题。**Krish** 和其他的专家正在帮助推动下一代数据仓储的行业成熟度，侧重于大数据、语义技术、众包、分析和平台工程。

**Krish** 是 **Sixth Sense Advisors** 公司的创始人兼 CEO，提供行业分析服务，覆盖数据仓库、分析、云计算、社交媒体和商务智能。**Krish** 还与其合伙人的组织一起在全球各地提供战略和创新咨询服务。

第2章 大数据的架构 ..... 37

2.1 大数据的架构 ..... 37

2.2 大数据的存储 ..... 37

2.3 大数据的传输 ..... 38

2.4 大数据的治理 ..... 38

2.5 大数据的隐私 ..... 38

2.6 大数据的扩展 ..... 39

2.7 大数据的集成 ..... 39

2.8 大数据的集成 ..... 39

2.9 大数据的集成 ..... 39

2.10 大数据的集成 ..... 39

2.11 大数据的集成 ..... 39

2.12 大数据的集成 ..... 39

2.13 大数据的集成 ..... 39

2.14 大数据的集成 ..... 39

2.15 大数据的集成 ..... 39

2.16 大数据的集成 ..... 39

2.17 大数据的集成 ..... 39

2.18 大数据的集成 ..... 39

2.19 大数据的集成 ..... 39

2.20 大数据的集成 ..... 39

2.21 大数据的集成 ..... 39

2.22 大数据的集成 ..... 39

2.23 大数据的集成 ..... 39

2.24 大数据的集成 ..... 39

2.25 大数据的集成 ..... 39

2.26 大数据的集成 ..... 39

2.27 大数据的集成 ..... 39

2.28 大数据的集成 ..... 39

2.29 大数据的集成 ..... 39

2.30 大数据的集成 ..... 39

2.31 大数据的集成 ..... 39

2.32 大数据的集成 ..... 39

2.33 大数据的集成 ..... 39

2.34 大数据的集成 ..... 39

2.35 大数据的集成 ..... 39

2.36 大数据的集成 ..... 39

2.37 大数据的集成 ..... 39

2.38 大数据的集成 ..... 39

2.39 大数据的集成 ..... 39

2.40 大数据的集成 ..... 39

2.41 大数据的集成 ..... 39

2.42 大数据的集成 ..... 39

2.43 大数据的集成 ..... 39

2.44 大数据的集成 ..... 39

2.45 大数据的集成 ..... 39

2.46 大数据的集成 ..... 39

2.47 大数据的集成 ..... 39

2.48 大数据的集成 ..... 39

2.49 大数据的集成 ..... 39

2.50 大数据的集成 ..... 39

2.51 大数据的集成 ..... 39

2.52 大数据的集成 ..... 39

2.53 大数据的集成 ..... 39

2.54 大数据的集成 ..... 39

2.55 大数据的集成 ..... 39

2.56 大数据的集成 ..... 39

2.57 大数据的集成 ..... 39

2.58 大数据的集成 ..... 39

2.59 大数据的集成 ..... 39

2.60 大数据的集成 ..... 39

2.61 大数据的集成 ..... 39

2.62 大数据的集成 ..... 39

2.63 大数据的集成 ..... 39

2.64 大数据的集成 ..... 39

2.65 大数据的集成 ..... 39

2.66 大数据的集成 ..... 39

2.67 大数据的集成 ..... 39

2.68 大数据的集成 ..... 39

2.69 大数据的集成 ..... 39

2.70 大数据的集成 ..... 39

2.71 大数据的集成 ..... 39

2.72 大数据的集成 ..... 39

2.73 大数据的集成 ..... 39

2.74 大数据的集成 ..... 39

2.75 大数据的集成 ..... 39

2.76 大数据的集成 ..... 39

2.77 大数据的集成 ..... 39

2.78 大数据的集成 ..... 39

2.79 大数据的集成 ..... 39

2.80 大数据的集成 ..... 39

2.81 大数据的集成 ..... 39

2.82 大数据的集成 ..... 39

2.83 大数据的集成 ..... 39

2.84 大数据的集成 ..... 39

2.85 大数据的集成 ..... 39

2.86 大数据的集成 ..... 39

2.87 大数据的集成 ..... 39

2.88 大数据的集成 ..... 39

2.89 大数据的集成 ..... 39

2.90 大数据的集成 ..... 39

2.91 大数据的集成 ..... 39

2.92 大数据的集成 ..... 39

2.93 大数据的集成 ..... 39

2.94 大数据的集成 ..... 39

2.95 大数据的集成 ..... 39

2.96 大数据的集成 ..... 39

2.97 大数据的集成 ..... 39

2.98 大数据的集成 ..... 39

2.99 大数据的集成 ..... 39

3.00 大数据的集成 ..... 39

3.1 大数据的集成 ..... 39

3.2 大数据的集成 ..... 39

3.3 大数据的集成 ..... 39

3.4 大数据的集成 ..... 39

3.5 大数据的集成 ..... 39

3.6 大数据的集成 ..... 39

3.7 大数据的集成 ..... 39

3.8 大数据的集成 ..... 39

3.9 大数据的集成 ..... 39

3.10 大数据的集成 ..... 39

3.11 大数据的集成 ..... 39

3.12 大数据的集成 ..... 39

3.13 大数据的集成 ..... 39

3.14 大数据的集成 ..... 39

3.15 大数据的集成 ..... 39

3.16 大数据的集成 ..... 39

3.17 大数据的集成 ..... 39

3.18 大数据的集成 ..... 39

3.19 大数据的集成 ..... 39

3.20 大数据的集成 ..... 39

3.21 大数据的集成 ..... 39

3.22 大数据的集成 ..... 39

3.23 大数据的集成 ..... 39

3.24 大数据的集成 ..... 39

3.25 大数据的集成 ..... 39

3.26 大数据的集成 ..... 39

3.27 大数据的集成 ..... 39

3.28 大数据的集成 ..... 39

3.29 大数据的集成 ..... 39

3.30 大数据的集成 ..... 39

3.31 大数据的集成 ..... 39

3.32 大数据的集成 ..... 39

3.33 大数据的集成 ..... 39

3.34 大数据的集成 ..... 39

3.35 大数据的集成 ..... 39

3.36 大数据的集成 ..... 39

3.37 大数据的集成 ..... 39

3.38 大数据的集成 ..... 39

3.39 大数据的集成 ..... 39

3.40 大数据的集成 ..... 39

3.41 大数据的集成 ..... 39

3.42 大数据的集成 ..... 39

3.43 大数据的集成 ..... 39

3.44 大数据的集成 ..... 39

3.45 大数据的集成 ..... 39

3.46 大数据的集成 ..... 39

3.47 大数据的集成 ..... 39

3.48 大数据的集成 ..... 39

3.49 大数据的集成 ..... 39

3.50 大数据的集成 ..... 39

3.51 大数据的集成 ..... 39

3.52 大数据的集成 ..... 39

3.53 大数据的集成 ..... 39

3.54 大数据的集成 ..... 39

3.55 大数据的集成 ..... 39

3.56 大数据的集成 ..... 39

3.57 大数据的集成 ..... 39

3.58 大数据的集成 ..... 39

3.59 大数据的集成 ..... 39

3.60 大数据的集成 ..... 39

3.61 大数据的集成 ..... 39

3.62 大数据的集成 ..... 39

3.63 大数据的集成 ..... 39

3.64 大数据的集成 ..... 39

3.65 大数据的集成 ..... 39

3.66 大数据的集成 ..... 39

3.67 大数据的集成 ..... 39

3.68 大数据的集成 ..... 39

3.69 大数据的集成 ..... 39

3.70 大数据的集成 ..... 39

3.71 大数据的集成 ..... 39

3.72 大数据的集成 ..... 39

3.73 大数据的集成 ..... 39

3.74 大数据的集成 ..... 39

3.75 大数据的集成 ..... 39

3.76 大数据的集成 ..... 39

3.77 大数据的集成 ..... 39

3.78 大数据的集成 ..... 39

3.79 大数据的集成 ..... 39

3.80 大数据的集成 ..... 39

3.81 大数据的集成 ..... 39

3.82 大数据的集成 ..... 39

3.83 大数据的集成 ..... 39

3.84 大数据的集成 ..... 39

3.85 大数据的集成 ..... 39

3.86 大数据的集成 ..... 39

3.87 大数据的集成 ..... 39

3.88 大数据的集成 ..... 39

3.89 大数据的集成 ..... 39

3.90 大数据的集成 ..... 39

3.91 大数据的集成 ..... 39

3.92 大数据的集成 ..... 39

3.93 大数据的集成 ..... 39

3.94 大数据的集成 ..... 39

3.95 大数据的集成 ..... 39

3.96 大数据的集成 ..... 39

3.97 大数据的集成 ..... 39

3.98 大数据的集成 ..... 39

3.99 大数据的集成 ..... 39

4.00 大数据的集成 ..... 39

## 作者简介

第4章 大数据的集成 ..... 39

4.1 大数据的集成 ..... 39

4.2 大数据的集成 ..... 39

4.3 大数据的集成 ..... 39

4.4 大数据的集成 ..... 39

4.5 大数据的集成 ..... 39

4.6 大数据的集成 ..... 39

4.7 大数据的集成 ..... 39

4.8 大数据的集成 ..... 39

4.9 大数据的集成 ..... 39

4.10 大数据的集成 ..... 39

4.11 大数据的集成 ..... 39

4.12 大数据的集成 ..... 39

4.13 大数据的集成 ..... 39

4.14 大数据的集成 ..... 39

4.15 大数据的集成 ..... 39

4.16 大数据的集成 ..... 39

4.17 大数据的集成 ..... 39

4.18 大数据的集成 ..... 39

4.19 大数据的集成 ..... 39

4.20 大数据的集成 ..... 39

4.21 大数据的集成 ..... 39

4.22 大数据的集成 ..... 39

4.23 大数据的集成 ..... 39

4.24 大数据的集成 ..... 39

4.25 大数据的集成 ..... 39

4.26 大数据的集成 ..... 39

4.27 大数据的集成 ..... 39

4.28 大数据的集成 ..... 39

4.29 大数据的集成 ..... 39

4.30 大数据的集成 ..... 39

4.31 大数据的集成 ..... 39

4.32 大数据的集成 ..... 39

4.33 大数据的集成 ..... 39

4.34 大数据的集成 ..... 39

4.35 大数据的集成 ..... 39

4.36 大数据的集成 ..... 39

4.37 大数据的集成 ..... 39

4.38 大数据的集成 ..... 39

4.39 大数据的集成 ..... 39

4.40 大数据的集成 ..... 39

4.41 大数据的集成 ..... 39

4.42 大数据的集成 ..... 39

4.43 大数据的集成 ..... 39

4.44 大数据的集成 ..... 39

4.45 大数据的集成 ..... 39

4.46 大数据的集成 ..... 39

4.47 大数据的集成 ..... 39

4.48 大数据的集成 ..... 39

4.49 大数据的集成 ..... 39

4.50 大数据的集成 ..... 39

4.51 大数据的集成 ..... 39

4.52 大数据的集成 ..... 39

4.53 大数据的集成 ..... 39

4.54 大数据的集成 ..... 39

4.55 大数据的集成 ..... 39

4.56 大数据的集成 ..... 39

4.57 大数据的集成 ..... 39

4.58 大数据的集成 ..... 39

4.59 大数据的集成 ..... 39

4.60 大数据的集成 ..... 39

4.61 大数据的集成 ..... 39

4.62 大数据的集成 ..... 39

4.63 大数据的集成 ..... 39

4.64 大数据的集成 ..... 39

4.65 大数据的集成 ..... 39

4.66 大数据的集成 ..... 39

4.67 大数据的集成 ..... 39

4.68 大数据的集成 ..... 39

4.69 大数据的集成 ..... 39

4.70 大数据的集成 ..... 39

4.71 大数据的集成 ..... 39

4.72 大数据的集成 ..... 39

4.73 大数据的集成 ..... 39

4.74 大数据的集成 ..... 39

4.75 大数据的集成 ..... 39

4.76 大数据的集成 ..... 39

4.77 大数据的集成 ..... 39

4.78 大数据的集成 ..... 39

4.79 大数据的集成 ..... 39

4.80 大数据的集成 ..... 39

4.81 大数据的集成 ..... 39

4.82 大数据的集成 ..... 39

4.83 大数据的集成 ..... 39

4.84 大数据的集成 ..... 39

4.85 大数据的集成 ..... 39

4.86 大数据的集成 ..... 39

4.87 大数据的集成 ..... 39

4.88 大数据的集成 ..... 39

4.89 大数据的集成 ..... 39

4.90 大数据的集成 ..... 39

4.91 大数据的集成 ..... 39

4.92 大数据的集成 ..... 39

4.93 大数据的集成 ..... 39

4.94 大数据的集成 ..... 39

4.95 大数据的集成 ..... 39

4.96 大数据的集成 ..... 39

4.97 大数据的集成 ..... 39

4.98 大数据的集成 ..... 39

4.99 大数据的集成 ..... 39

5.00 大数据的集成 ..... 39

# 目 录

译者序	2.3.4 外部或第三方数据	15
前言	2.3.5 电子邮件	15
致谢	2.3.6 合同	15
作者简介	2.3.7 地理信息系统和地理空间数据	16
	2.3.8 示例: Funshots 公司	17
	2.4 数据速度	19
	2.4.1 Amazon、Facebook、Yahoo 和 Google	19
	2.4.2 传感器数据	19
	2.4.3 移动网络	20
	2.4.4 社交媒体	20
	2.5 数据多样性	21
	2.6 总结	22
	第 3 章 大数据处理架构	23
	3.1 引言	23
	3.2 再论数据处理	23
	3.3 数据处理技术	24
	3.4 数据处理基础设施的挑战	25
	3.4.1 存储	25
	3.4.2 传输	25
	3.4.3 处理	26
	3.4.4 速度或吞吐量	26
	3.5 全共享架构与无共享架构的比较	26
	3.5.1 全共享架构	27
	3.5.2 无共享架构	27
	3.5.3 OLTP 与数据仓库	28
	3.6 大数据处理	28
	3.6.1 基础设施方面	31
	3.6.2 数据处理方面	32
第 1 章 大数据简介		2
1.1 引言		2
1.2 大数据		2
1.3 大数据的定义		4
1.4 为什么需要大数据? 为什么是现在		4
1.5 大数据示例		5
1.5.1 社交媒体的文章		5
1.5.2 调查数据分析		6
1.5.3 调查数据		7
1.5.4 气象数据		8
1.5.5 Twitter 数据		8
1.5.6 集成和分析		8
1.5.7 附加数据的类型		10
1.6 总结		11
延伸阅读		11
第 2 章 使用大数据		12
2.1 引言		12
2.2 数据爆炸		12
2.3 数据体量		13
2.3.1 机器数据		14
2.3.2 应用日志		14
2.3.3 点击流日志		14

3.7 电信大数据研究	32	5.3.6 先进的光纤网结合实时流数据	85
3.7.1 基础设施	34	5.3.7 解决方案组件	85
3.7.2 数据处理	34	5.3.8 扩展安全边界创建战略优势	85
<b>第4章 大数据技术简介</b>	<b>35</b>	5.3.9 关联传感器数据使得假阳性率为零	86
4.1 引言	35	<b>5.4 案例研究3: 通过大数据分析改善患者预后</b>	<b>86</b>
4.2 分布式数据处理	36	5.4.1 摘要	86
4.3 大数据处理需求	38	5.4.2 业务目标	87
4.4 大数据处理技术	39	5.4.3 挑战	87
4.5 Hadoop	42	5.4.4 概述: 给从业人员新的洞察以指导患者护理	87
4.5.1 Hadoop 核心组件	43	5.4.5 挑战: 将传统数据仓库生态系统与大数据融合	87
4.5.2 Hadoop 总结	69	5.4.6 解决方案: 为大数据分析做好准备	88
4.6 NoSQL	69	5.4.7 结果: 消除“数据陷阱”	88
4.6.1 CAP 定理	69	5.4.8 为什么是 aster	88
4.6.2 键-值对: Voldemort	70	5.4.9 关于 Aurora	89
4.6.3 列簇存储: Cassandra	70	<b>5.5 案例研究4: 安大略大学技术学院——利用关键数据, 提供积极的患者护理</b>	<b>89</b>
4.6.4 文档数据库: Riak	76	5.5.1 摘要	89
4.6.5 图数据库	77	5.5.2 概述	89
4.6.6 NoSQL 小结	78	5.5.3 商业上的收益	90
4.7 文本 ETL 处理	78	5.5.4 更好地利用数据资源	90
延伸阅读	79	5.5.5 智慧医疗保健	91
<b>第5章 大数据驱动的商业价值</b>	<b>80</b>	5.5.6 解决方案组件	91
5.1 引言	80	5.5.7 融合人类知识与技术	92
5.2 案例研究1: 传感器数据	81	5.5.8 扩大 Artemis 的影响	92
5.2.1 摘要	81	<b>5.6 案例研究5: 微软 SQL Server 客户解决方案</b>	<b>93</b>
5.2.2 Vestas	81	5.6.1 客户画像	93
5.2.3 概述	81	5.6.2 解决方案的亮点	93
5.2.4 利用风力发电	81	5.6.3 业务需求	93
5.2.5 把气候变成资本	82	5.6.4 解决方案	94
5.2.6 跟踪大数据的挑战	83	5.6.5 好处	94
5.2.7 维持数据中心的能源效率	83		
5.3 案例研究2: 流数据	84		
5.3.1 摘要	84		
5.3.2 监控和安全: TerraEchos	84		
5.3.3 需求	84		
5.3.4 解决方案	84		
5.3.5 效益	84		

5.7 案例研究 6: 以客户为中心的数据集成 .....	95
5.7.1 概述 .....	95
5.7.2 解决方案设计 .....	98
5.7.3 促成更好的交叉销售和追加销售的机会 .....	99
5.8 总结 .....	100

## 第二部分 数据仓库

第 6 章 再论数据仓库 .....	102
6.1 引言 .....	102
6.2 传统的数据仓库或 DW 1.0 .....	103
6.2.1 数据架构 .....	103
6.2.2 基础设施 .....	104
6.2.3 数据仓库的陷阱 .....	106
6.2.4 建立数据仓库的架构方法 .....	111
6.3 DW 2.0 .....	113
6.3.1 Inmon 的 DW 2.0 概述 .....	114
6.3.2 DSS 2.0 概述 .....	115
6.4 总结 .....	116
延伸阅读 .....	116

## 第 7 章 数据仓库的再造 .....

7.1 引言 .....	118
7.2 企业数据仓库平台 .....	118
7.2.1 事务型系统 .....	119
7.2.2 运营数据存储区 .....	119
7.2.3 分段区 .....	120
7.2.4 数据仓库 .....	120
7.2.5 数据集市 .....	120
7.2.6 分析型数据库 .....	121
7.2.7 数据仓库的问题 .....	121
7.3 再造数据仓库的选择 .....	122
7.3.1 平台再造 .....	122
7.3.2 平台工程 .....	123
7.3.3 数据工程 .....	124
7.4 使数据仓库现代化 .....	125

7.5 使数据仓库现代化的案例研究 .....	127
7.5.1 当前状态分析 .....	127
7.5.2 推荐 .....	127
7.5.3 现代化的业务收益 .....	128
7.5.4 一体机的选择过程 .....	128
7.6 总结 .....	132

## 第 8 章 数据仓库中的工作负载管理 .....

8.1 引言 .....	133
8.2 当前状态 .....	133
8.3 工作负载的定义 .....	134
8.4 了解工作负载 .....	135
8.4.1 数据仓库输出 .....	136
8.4.2 数据仓库输入 .....	137
8.5 查询分类 .....	138
8.5.1 宽 / 宽 .....	138
8.5.2 宽 / 窄 .....	139
8.5.3 窄 / 宽 .....	139
8.5.4 窄 / 窄 .....	139
8.5.5 非结构化 / 半结构化数据 .....	140
8.6 ETL 和 CDC 的工作负载 .....	140
8.7 度量 .....	141
8.8 当前系统设计的局限 .....	142
8.9 新工作负载和大数据 .....	143
8.10 技术选择 .....	144
8.11 总结 .....	144

## 第 9 章 应用到数据仓库的新技术 .....

9.1 引言 .....	145
9.2 重新检查数据仓库挑战 .....	145
9.2.1 数据加载 .....	145
9.2.2 可用性 .....	146
9.2.3 数据体量 .....	146
9.2.4 存储性能 .....	147
9.2.5 查询性能 .....	147
9.2.6 数据传输 .....	147
9.3 数据仓库一体机 .....	147
9.3.1 一体机架构 .....	148

9.3.2	一体机中的数据分布	149
9.3.3	部署数据仓库一体机 最佳实践	150
9.3.4	大数据一体机	152
9.4	云计算	152
9.4.1	基础设施即服务	152
9.4.2	平台即服务	152
9.4.3	软件即服务	153
9.4.4	云基础架构	153
9.4.5	云计算给数据仓库带来 的好处	154
9.4.6	将云计算用于数据仓库 所面临的问题	154
9.5	数据虚拟化	154
9.5.1	数据虚拟化是什么	155
9.5.2	提高商务智能性能	156
9.5.3	工作负载分布	156
9.5.4	实施数据虚拟化项目	156
9.5.5	使用数据虚拟化时应 避免的误区	157
9.5.6	内存技术	157
9.5.7	内存架构的好处	157
9.6	总结	158
	延伸阅读	158

### 第三部分 构建大数据 - 数据仓库

#### 第 10 章 大数据和数据仓库的集成 160

10.1	引言	160
10.2	新数据仓库的组件	160
10.2.1	数据层	161
10.2.2	算法	162
10.2.3	技术层	163
10.3	集成策略	164
10.3.1	数据驱动的集成	164
10.3.2	物理组件集成和架构	167
10.3.3	外部数据集成	168
10.4	Hadoop 与 RDBMS	169
10.5	大数据一体机	171

10.6	数据虚拟化	172
10.7	语义框架	173
10.7.1	词法处理	174
10.7.2	聚类	174
10.7.3	语义知识处理	174
10.7.4	信息抽取	175
10.7.5	可视化	175
10.8	总结	175

#### 第 11 章 大数据的数据驱动架构 176

11.1	引言	176
11.2	元数据	177
11.2.1	技术元数据	177
11.2.2	业务元数据	178
11.2.3	上下文元数据	178
11.2.4	过程设计级元数据	178
11.2.5	程序级元数据	178
11.2.6	基础设施元数据	179
11.2.7	核心业务元数据	179
11.2.8	运营元数据	179
11.2.9	商务智能型元数据	180
11.3	主数据管理	180
11.4	处理数据仓库中的数据	181
11.5	处理大数据的复杂性	184
11.5.1	处理能力的限制	184
11.5.2	处理大数据	184
11.6	机器学习	190
11.7	总结	193

#### 第 12 章 大数据的信息管理和 生命周期 195

12.1	引言	195
12.2	信息生命周期管理	195
12.2.1	目标	196
12.2.2	信息管理策略	196
12.2.3	治理	196
12.2.4	信息生命周期管理的 优点	200
12.3	大数据的信息生命周期管理	200

12.3.1 示例：信息生命周期管理和社交媒体数据...200

12.3.2 测量信息生命周期管理的影响...202

12.4 总结...203

**第 13 章 大数据分析、可视化和数据科学家...204**

13.1 引言...204

13.2 大数据分析...204

13.3 数据发现...206

13.4 可视化...206

13.5 数据科学家的角色变迁...207

13.6 总结...208

**第 14 章 实施大数据 – 数据仓库的现实情况...209**

14.1 引言：构建大数据 – 数据仓库...209

14.2 以客户为中心的业务转型...209

14.3 Hadoop 和 MySQL 驱动创新...212

14.4 将大数据集成到数据仓库中...214

14.4.1 增强决策制订...215

14.4.2 成果...216

14.5 总结...216

**附录 A 客户案例研究...217**

**附录 B 建设医疗保健信息工厂...237**

**结束语...269**

可以存储海量的数据，以便进行数学或统计学方面的分析研究。这些理解能够提供更可靠的业务洞察力，帮助公司制定更好的决策并提升质量很有用。

大数据具有 4 个特征：海量的数据体量、多样的数据类型、分布广泛的数据以及快速变化的数据。在数据仓库内外采集的数据包含来自各种高度的结构化、半结构化以及非结构化的数据。这一说法的基本前提是基于传统数据仓库的数据来源，即数据库、数据仓库、数据湖以及数据表上所有者在考虑。全面的大数据包括所有的活动的数据。

大数据非常通用。例如，电影票购买是一个简单的例子。电影的票房取决于电影的受欢迎程度、以及它何时何地（通过设备、平台、地区和语言等）播放。此外，它还取决于天气、经济、情感、记忆和习惯。或者，不要把数据视为决策时使用的信息。大数据可以帮助您了解客户行为、产品结构、营销活动的效果、运营效率、供应链、以及更多其他方面。

大数据可以帮助您了解客户行为、产品结构、营销活动的效果、运营效率、供应链、以及更多其他方面。大数据可以帮助您了解客户行为、产品结构、营销活动的效果、运营效率、供应链、以及更多其他方面。

### 第一部分

# 大数据

- 第 1 章 大数据简介
- 第 2 章 使用大数据
- 第 3 章 大数据处理架构
- 第 4 章 大数据技术简介
- 第 5 章 大数据驱动的商业价值



## 第1章

# 大数据简介

## 1.1 引言

自“因特网”产生以来，吸引现代计算行业注意力的最大的现象是“大数据”。这两个词首次同时出现在麦肯锡公司同主题的论文中，而其基本定义则来自 Gartner 公司的 Doug Laney。

今天的“大数据”受欢迎的根本原因是因为技术平台的同时出现，这些平台有能力处理多种格式和结构的数据，而且不用担心传统的系统和数据库平台的限制。

## 1.2 大数据

数据是表示信息或知识最原始的格式。在计算领域，我们通常认为数据是有组织的值的行和列，这些值代表一个或者多个实体及其属性。在使用电脑帮助处理计算或信息管理的时代之前的很长一段时间，数据是由古希腊人随着计数和贸易的到来而发明的。简单来说，就是把值分配给数字，然后使用这些数字来标记货币价值、人口、日历、税收以及许多和历史有关的实例。这些实例提供了大量证据来证明人类思维在数据和知识的获取与管理方面的魅力。

根据卡内基·梅隆大学的一系列研究，信息或数据管理的必要过程是组织、采集、存储、检索和管理数据。从不同的过程采集的数据使得我们在处理输出结果时，能够理解和满足需求，从而做出决策。这种管理行为是赫伯特·西蒙有限的理性观点<sup>①</sup>或当人类思维应用到数据管理时有限的视野的基本主题。对于决策行为和管理行为的争论的提出是很有意义的，因为我们限制了数据建模过程和应用算法的应用中的数据，而且一直以来都致力数据内的离散关系而不是整体。

然而，在现实中，决策总是超越用于辅助处理的传统系统。例如，患者的治疗和管理不局限于计算机和程序。但是现在在医院里，对于每一个患者，均可通过非结构化的数据集成技术和算法，以电子方式收集和处理由医生、护士、实验室技术人员、急救人员和医

<sup>①</sup> March, J. G., & Simon, H. A. (1958) Organizations (<http://www.amazon.com/Organizations-James-G-March/dp/063118631X>).