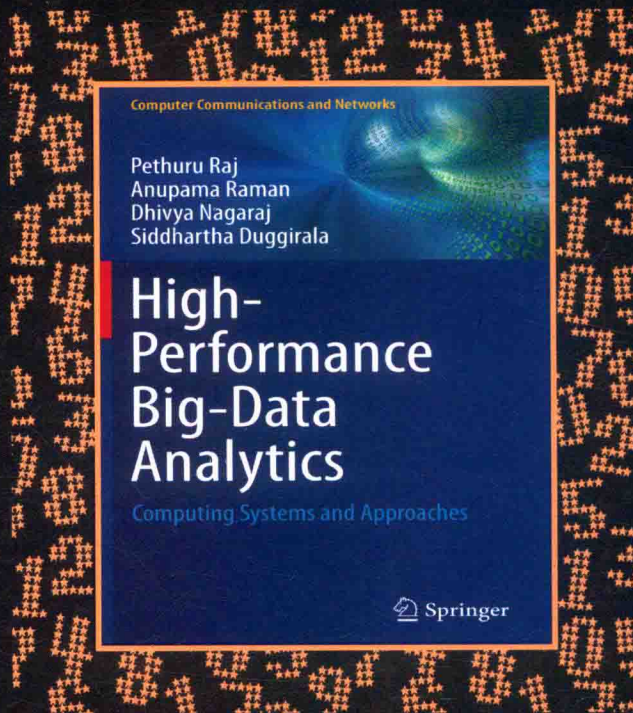


高性能计算系统 与大数据分析

佩瑟鲁·拉吉 (Pethuru Raj)
阿诺帕马·拉曼 (Anupama Raman)
德维亚·纳加拉杰 (Dhivya Nagaraj) 著
悉达多·杜格拉拉 (Siddhartha Duggirala)
齐宁 庞建民 张铮 韩林 译



HIGH-PERFORMANCE
BIG-DATA ANALYTICS
COMPUTING SYSTEMS AND APPROACHES

HIGH-PERFORMANCE
BIG-DATA ANALYTICS
COMPUTING SYSTEMS AND APPROACHES

高性能计算系统 与大数据分析

佩瑟鲁·拉吉 (Pethuru Raj)

[印] 阿诺帕马·拉曼 (Anupama Raman)

德维亚·纳加拉杰 (Dhivya Nagaraj)

悉达多·杜格拉拉 (Siddhartha Duggirala)

齐宁 庞建民 张铮 韩林 译



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

高性能计算系统与大数据分析 / (印) 佩瑟鲁·拉吉 (Pethuru Raj) 等著; 齐宁等译. —北京: 机械工业出版社, 2018.11

(数据科学与工程丛书)

书名原文: High-Performance Big-Data Analytics: Computing Systems and Approaches

ISBN 978-7-111-61175-2

I. 高… II. ①佩… ②齐… III. ①高性能计算机—计算机系统 ②数据处理 IV. ① TP38
② TP274

中国版本图书馆 CIP 数据核字 (2018) 第 234674 号

本书版权登记号: 图字 01-2017-2012

Translation from the English language edition:

High-Performance Big-Data Analytics: Computing Systems and Approaches

by Pethuru Raj, Anupama Raman, Dhivya Nagaraj and Siddhartha Duggirala.

Copyright © Springer International Publishing Switzerland 2015.

This Springer imprint is published by Springer Nature.

The registered company is Springer International Publishing AG.

All Rights Reserved.

本书中文简体字版由 Springer 授权机械工业出版社独家出版。未经出版者书面许可, 不得以任何方式复制或抄袭本书内容。

本书介绍了大数据分析所需的高性能基础设施以及高性能大数据分析领域的新技术和工具。在新兴分析类型方面, 涵盖了传感器分析、机器分析、运营分析、实时分析、高性能分析、社交媒体和网络分析、客户情绪分析、品牌优化分析、金融交易及趋势分析、零售分析、能量分析、药物分析以及效用分析等。在 IT 基础设施方面, 则包含了大型机、并行和超级计算系统、P2P、集群和网格计算系统设备、专业集成和按需定制的系统、实时系统、云基础设施等。

本书适合作为高校大数据、高性能计算相关课程的教材, 也适合业务主管、技术专家、软件工程师、大数据科学家、解决方案架构师等专业人士阅读。

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: 朱秀英

责任校对: 殷 虹

印 刷: 北京瑞德印刷有限公司

版 次: 2019 年 1 月第 1 版第 1 次印刷

开 本: 185mm×260mm 1/16

印 张: 18

书 号: ISBN 978-7-111-61175-2

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) 88378991 88361066

投稿热线: (010) 88379604

购书热线: (010) 68326294 88379649 68995259

读者信箱: hzsj@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

近年来，随着信息技术的发展，特别是互联网和物联网的飞速发展，产生、收集、存储了大量的数据，急需有效的分析方法从数据中挖掘有意义的规律，这使得大数据技术成为当前非常流行的一种技术。

本书同市面上常见的介绍大数据技术或工具的书籍有较大的不同，更侧重于介绍大数据分析所需的高性能基础设施以及高性能大数据分析领域的新技术和工具。本书内容非常丰富，在新兴分析类型方面，涵盖了传感器分析、机器分析、运营分析、实时分析、高性能分析、社交媒体和网络分析、客户情绪分析、品牌优化分析、金融交易及趋势分析、零售分析、能量分析、药物分析以及效用分析等。在 IT 基础设施方面，则包含了大型机、并行和超级计算系统、P2P、集群和网格计算系统设备、专业集成和按需定制的系统、实时系统、云基础设施等。

本书由齐宁、庞建民、张铮、韩林完成主要章节的翻译，刘浩、刘镇武也参与了本书的部分翻译工作。在为期近一年的翻译过程中，虽然我们已经对译稿进行了仔细校对，查阅了大量相关资料，使译文尽可能符合中文习惯并保持术语的一致性，但由于本书涉及的范围非常广泛，错误或不当之处仍难以完全避免，敬请各位读者和同行专家谅解，诚挚希望读者将相关意见、建议发送到电子邮箱 qining2005@126.com。

特别感谢机械工业出版社华章公司的朱劼编辑，没有她的信任、耐心与支持，整个翻译工作不可能完成。

译者

2018年9月于郑州

序

近年来，随着大量融合技术（数字化、连接性、集成、感知、小型化、消费化、商品化、知识发现与传播等）的出现，数据的增长幅度令人难以置信。简而言之，在我们的生活和工作环境中，每个普通物品都互相连接，并且支持使得它们能够无缝地、创造性地加入主流计算中的服务。近期，在设备方面也经历了各种各样的创新。各种数字化工件和设备之间深度、极致的互联是产生大量数据的主因。

这种趋势和变化不仅为 IT 专业人员带来了机遇和挑战，同时也为数据科学家带来了机遇和挑战。数据分析学科必然会在多个维度和方向上发展。较新类型的数据分析（通用的和专用的）必然会出现并快速发展，计算、存储和网络方面的挑战也注定变得更加严峻。随着数据规模、结构、范围、速度和价值的不断增加，所有企业和 IT 团队面临的最大挑战就是如何完美捕获和处理数据并及时获取可行的洞见。内部或外部产生的各种类型的数据都可以通过可用的模式、富有成效的联想、及时的警告、机会等形式提供隐藏的洞见。

在本书中，作者揭示了为什么大数据和快速数据分析需要高性能基础设施（服务器、存储设备、网络连接解决方案）来应对下一代数据分析解决方案。作者要阐明高性能大数据分析领域的最新技术和工具，因此有意识地聚焦在各种类型的 IT 基础设施和平台上。

数据分析：处理步骤 众所周知，在所有数据分析任务中，通常有三个主要的阶段：

- 1) 通过数据虚拟化平台进行数据捕捉。
- 2) 数据处理 / 解释平台用于知识发现。
- 3) 通过大量数据可视化平台来完成知识传播。

新兴的分析类型 随着数据的规模性（volume）、高速性（velocity）、多样性（variety）、变化性（variability）、黏性（viscosity）和真实性（veracity）的增长，大量的分析（通用的和专用的）用例正在被挖掘。各种垂直业务和细分行业正在采用不同类型的分析，目的是从数据中获取可行的洞见。通用的分析类型包括：

- 传感器分析
- 机器分析
- 运营分析

- 实时分析
- 高性能分析

特定领域的分析类型包括社交媒体和网络分析、客户情绪分析、品牌优化分析、金融交易及趋势分析、零售分析、能量分析、药物分析以及效用分析。

新兴的 IT 基础设施和平台 IT 基础设施越来越趋于融合、集中、联合、预集成、优化和有组织，并试图演化成未来企业发展的不二之选。分析平台正在以前所未有的力度冲击市场。对各种形式的数据进行仔细分析在商业运营的竞争力中很重要，理解到这一点，全球的企业开始急切寻找高性能的 IT 基础设施和平台，目的是以有效、高效的方式运行大数据和快速数据分析应用。

作者广泛介绍了现有的和新兴的高性能 IT 基础设施和平台，利用它们能够有效地实现灵活的数据分析。本书所讨论的主要 IT 基础设施包括：

- 1) 大型机
- 2) 并行和超级计算系统
- 3) P2P、集群和网格计算系统
- 4) 设备（数据仓库设备及 Hadoop 专用设备）
- 5) 专业集成和按需定制的系统
- 6) 实时系统
- 7) 云基础设施

作者在本书中有意识地提出了下一代 IT 基础设施以及大数据和快速数据分析平台的所有增值信息。除了业务主管、决策者和其他利益相关者之外，本书还对技术专家、顾问以及技术布道者、倡导者等非常有用。非常肯定地说，世界各地的软件工程师、解决方案架构师、云专业人员和大数据科学家都会发现本书非常翔实、有趣，能激励人们深刻理解数据分析将如何作为一种公共服务出现，并在让世界变得更加智慧的过程中发挥不可或缺的作用。

IBM Analytics, Orange County, CA, USA

TBDI, Orange County, CA, USA

Sushil Pramanick, FCA, PMP

前 言

一些行业趋势以及一系列强大的技术和工具无疑将导致大规模的数据爆炸。不经意间，数据已经压倒性地成为各行各业的战略资产。这些前所未有的数据包括以下值得注意的变化：设备生态系统随着人们不断变化的想象而持续扩展；随着智能仪器和互联技术的发展，机器变得智能，并且产生了高达 PB 乃至 EB 级的数据；个人及专业应用都支持服务，从而可以互相操作，进而实现有益的数据共享；社交网站每天产生 TB 级的数据；我们周围的普通物体都被精密地数字化，以不同的速度产生大量的多结构数据。另一方面，ICT 基础设施和平台被高度优化和组织以进行有效的数据存储、处理和分析，具有适应性的 WAN 技术正在形成以加速数据的安全传输，新的架构模式被融入，过程也系统地变得更加灵活，等等，目的是使数据有意义。

仔细分析数据可以提供丰富的信息，这些信息能够彻底改变我们生活的方方面面。这个想法已经在当今 IT 领域演变成游戏规则改变者，被人们称为大数据分析。考虑到数据的规模、速度、范围和结构，计算、存储和网络基础设施需要非常高效。大数据为 IT 带来了三个关键挑战：大数据的存储和管理，大数据分析，产生利用大数据分析的复杂应用。准确地说，大数据分析（BDA）正在迅速成为下一代高性能计算学科，学生、学者和科学家需要挖掘出有效的算法、模式、方法、最佳实践、关键准则、评价指标等。

本书概要介绍这些技术。为了高效率地捕捉、获取、吸收、处理大数据，以便实现知识发现和传播，目前需要对网络和存储基础设施优化进行认真的分析。本书中还包含了大数据分析在各个行业中的应用案例，目的是使读者以简明的方式了解数据分析的重要性。

第 1 章：IT 领域的变革以及未来趋势。本章列出了 IT 领域尤其是大数据和快速数据背景下的新变化。对 ICT 领域有前景的、潜在的技术及工具进行了特别介绍，目的是让读者了解本书中会涵盖哪些内容。

第 2 章：大数据 / 快速数据分析中的高性能技术。本章对高性能大数据及快速数据分析中最具代表性的技术进行了分类。

第 3 章：大数据与快速数据分析对高性能计算的渴望。本章解释了大数据和快速数据分析的本质，目的是强调高性能计算需求的重要性，从而能够从数据堆中获取可行的洞见。

第 4 章：高性能大数据分析的网络基础设施。本章总结了有效地传输大数据的网络基础设施要求。为了能够通过网络有效地进行大数据传输，需要对现有网络基础设施进行一

些改动。可以使用的技术包括网络虚拟化、软件定义网络 (SDN)、两层 Leaf-Spine 架构、网络功能虚拟化, 本章对这些技术进行了详细的讨论。此外, 还需要对现有的广域网基础设施进行优化, 以有效地传输大数据。本章还讨论了一种名为 FASP 的技术, 它能够有效利用 TCP/IP 协议传输大数据。FASP 的一些实现方面的问题也包含在本章中。

第 5 章: 高性能大数据分析的存储基础设施。本章总结了产生大数据的应用程序的存储基础设施需求。目前的存储基础设施没有对存储和处理大数据进行优化, 现有存储技术的主要问题在于缺乏可扩展性, 因此, 设计一种能够有效处理大数据的新存储技术是当务之急。在本章中, 首先介绍了现有存储基础设施以及它们对处理大数据的适合程度。之后, 介绍了一些专门为处理大数据而设计的平台和文件系统, 例如 Panasas 文件系统、Lustre 文件系统、GFS、HDFS。

第 6 章: 使用高性能计算进行实时分析。本章讨论了实时环境中的分析问题, 涵盖了新近的实时分析解决方案, 例如机器数据分析和运营分析。本章可以让读者了解数据是如何进行实时处理的, 以及实时处理对我们更美好的未来生活的价值。

第 7 章: 高性能计算范型。本章详细介绍了多年来高性能计算在大型机上的演变以及背后的原因。几年前, 得出的结论是大型机将随着技术的发展而消失, 但是像 IBM 这样的公司已经证明, 大型机不会消失, 而是通过提供曾经被认为完全不可能的解决方案继续发挥作用。

第 8 章: in-database 处理与 in-memory 分析。本章阐明 in-database 分析技术以及 in-memory 分析技术。当业务系统大规模运行时, 将数据移入或移出数据存储可能是非常令人畏惧且代价昂贵的。当我们将“处理”移动到“数据”的附近时, 数据处理是在数据存储中完成的, 这样做可以减少数据移动成本, 并使用更大的数据集来挖掘数据。随着企业的发展, 速度已经变得至关重要, 此时就需要实时数据库来发挥作用。本章涵盖了 in-database 分析技术及 in-memory 分析技术的方方面面, 并给出了适当的例子。

第 9 章: 大数据 / 快速数据分析中的高性能集成系统、数据库和数据仓库。在即将到来的大数据时代, 对新型数据管理系统有着独特的需求。本章清晰地介绍了新出现的集群 SQL 数据库、NoSQL 数据库和 NewSQL 数据库, 并对专用于大数据的数据仓库进行了解释。

第 10 章: 高性能网格和集群。本章阐明了可用于支持大数据分析 & 数据密集型处理的技术和软件工具。全球的企业都面临着降低分析平台的 TCO (总体拥有成本) 的压力, 同时还要在必要的水平上继续运行。使用这些高性能系统, 企业能够满足所需的性能要求。本章介绍了集群和网格计算系统在大数据分析领域的不同用例。

第 11 章: 高性能 P2P 系统。本章介绍了大数据分析领域中使用的 P2P 技术和工具。由于数据存储或分析系统的大规模性质, 服务器之间通常具有主从关系。这有助于应用程序的并行化, 但是当主节点故障时会产生问题——所有的请求都得不到回复。在这种场景下, 如果软件结构是分散的, 即没有主服务器, 那么就不会发生单点故障, 因此所有的请

求都会得到回复。本章介绍了使用高性能 P2P 系统的不同用例。

第 12 章：高性能大数据分析的可视化维度。本章主要介绍可视化技术和工具。随着数据大小以及数据复杂性的增加，理解数据的含义变得更加困难。如果数据或分析输出以某种可视化形式而不是简单的数字显示，用户可以轻松地获取其含义并据此开展工作。本章介绍了大数据分析领域所使用的信息可视化技术的不同用例。

第 13 章：用于组织增权的社交媒体分析。本章重点介绍社交媒体分析，这是大数据分析的主要技术用例之一。大数据的主要驱动之一就是社交媒体网络所产生的大量非结构化数据。这导致了一种名为社交媒体分析的新分析潮流的出现。本章讨论了社交媒体分析演变的各种驱动因素，详细讨论了描述社交媒体分析用于组织变革的各种用例，此外还详细讨论了跟踪社交媒体对组织的影响时使用的内容指标。用于社交媒体分析的关键预测分析技术是使用文本挖掘的网络分析和情感分析，本章对这两种技术进行了讨论，此外还讨论了一些用于社交媒体分析的工具。

第 14 章：医疗保健的大数据分析。这一章说明了分析在医疗保健领域的重要性。不言而喻，医疗保健的未来是我们所有人的未来。本章涵盖了医疗保健分析的重要驱动因素以及医疗保健中的大数据分析用例。本章提供了一个例子，该例子说明过去未被注意的数据有望以高性价比方式向患者提供优质护理。

目 录

译者序	
序	
前言	
第 1 章 IT 领域的变革以及未来趋势 1	
1.1 引言	1
1.2 新兴的 IT 趋势	1
1.3 数字化实体的实现与发展	4
1.4 物联网 / 万物互联	5
1.5 对社交媒体网站的广泛采用	7
1.6 预测性、规范性、个性化分析时代	7
1.7 用于大数据及分析的 Apache Hadoop	11
1.8 大数据、大洞见、大动作	13
1.9 结论	15
1.10 习题	15
第 2 章 大数据 / 快速数据分析中的高性能技术 16	
2.1 引言	16
2.2 大数据分析学科的出现	17
2.3 大数据的战略意义	18
2.4 大数据分析的挑战	19
2.5 高性能计算范型	19
2.6 通过并行实现高性能的方法	21
2.7 集群计算	22
2.8 网格计算	24
2.9 云计算	27
2.10 异构计算	29
2.11 用于高性能计算的大型机	31
2.12 用于大数据分析的超级计算	32
2.13 用于大数据分析的设备	32
2.13.1 用于大规模数据分析的数据仓库设备	33
2.13.2 in-memory 大数据分析	35
2.13.3 大数据的 in-database 处理	37
2.13.4 基于 Hadoop 的大数据设备	38
2.13.5 高性能大数据存储设备	41
2.14 结论	42
2.15 习题	42
参考文献	43
第 3 章 大数据与快速数据分析对高性能计算的渴望 44	
3.1 引言	44
3.2 重新审视大数据分析范型	45
3.3 大数据和快速数据的含义	47
3.4 用于精确、预测性、规范性洞见的新兴数据源	48
3.5 大数据分析为何不俗	50
3.6 传统的和新一代的数据分析案例研究	51
3.7 为何采用基于云的大数据分析	55

3.8	大数据分析：主要处理步骤	57	5.6.1	以太网光纤通道	88
3.9	实时分析	58	5.7	网络附属存储	89
3.10	流分析	62	5.8	用于高性能大数据分析的流行文件系统	89
3.11	传感器分析	63	5.8.1	Google 文件系统	89
3.11.1	大数据分析 with 高性能计算的同步：附加价值	63	5.8.2	Hadoop 分布式文件系统	91
3.12	结论	64	5.8.3	Panasas	92
3.13	习题	64	5.8.4	Luster 文件系统	94
第 4 章 高性能大数据分析的网络基础设施			5.9	云存储简介	96
4.1	引言	65	5.9.1	云存储系统的架构模型	96
4.2	当前网络基础设施的局限	66	5.9.2	存储虚拟化	98
4.3	高性能大数据分析网络基础设施的设计方法	68	5.9.3	云存储中使用的存储优化技术	100
4.3.1	网络虚拟化	68	5.9.4	云存储的优点	101
4.3.2	软件定义网络	76	5.10	结论	101
4.3.3	网络功能虚拟化	78	5.11	习题	101
4.4	用于传输大数据的广域网优化	79	参考文献	102	
4.5	结论	81	进一步阅读	102	
4.6	习题	81	第 6 章 使用高性能计算进行实时分析		
参考文献	81	6.1	引言	103	
第 5 章 高性能大数据分析的存储基础设施			6.2	支持实时分析的技术	103
5.1	引言	82	6.2.1	in-memory 处理	103
5.2	直连式存储	83	6.2.2	in-database 分析	105
5.2.1	DAS 的缺点	84	6.3	大规模在线分析	106
5.3	存储区域网络	85	6.4	通用并行文件系统	107
5.3.1	块级访问	85	6.4.1	GPFS 用例	107
5.3.2	文件级访问	85	6.5	GPFS 客户案例研究	111
5.3.3	对象级访问	85	6.5.1	广播公司：VRT	111
5.4	保存大数据的存储基础设施需求	86	6.5.2	石油公司从 Lustre 迁移到 GPFS	113
5.5	光纤通道存储区域网络	87	6.6	GPFS：关键的区别	113
5.6	互联网协议存储区域网络	88	6.6.1	基于 GPFS 的解决方案	114
			6.7	机器数据分析	114

6.7.1 Splunk	114	7.12 Windows 高性能计算	129
6.8 运营分析	115	7.13 结论	130
6.8.1 运营分析中的技术	115	7.14 习题	131
6.8.2 用例以及运营分析产品	116		
6.8.3 其他 IBM 运营分析产品	117	第 8 章 in-database 处理与 in-memory 分析	132
6.9 结论	117	8.1 引言	132
6.10 习题	118	8.1.1 分析工作负载与事务工作负载的对比	132
第 7 章 高性能计算范型	119	8.1.2 分析工作负载的演化	133
7.1 引言	119	8.1.3 传统分析平台	135
7.2 为何还需要大型机	119	8.2 in-database 分析	135
7.3 大型机中 HPC 是如何演化的	120	8.2.1 架构	137
7.3.1 成本: HPC 的一个重要因素	120	8.2.2 优点和局限	138
7.3.2 云计算中的集中式 HPC	120	8.2.3 代表性的系统	138
7.3.3 集中式 HPC 的要求	121	8.3 in-memory 分析	140
7.4 HPC 远程模拟	121	8.3.1 架构	141
7.5 使用 HPC 的大型机解决方案	121	8.3.2 优点和局限	142
7.5.1 智能大型机网格	121	8.3.3 代表性的系统	142
7.5.2 IMG 的工作原理	122	8.4 分析设备	145
7.5.3 IMG 架构	122	8.4.1 Oracle Exalytics	145
7.6 架构模型	125	8.4.2 IBM Netezza	145
7.6.1 具有共享磁盘的存储服务器	125	8.5 结论	147
7.6.2 没有共享磁盘的存储服务器	125	8.6 习题	147
7.6.3 无存储服务器的通信网络	125	参考文献	148
7.7 对称多处理	126	进一步阅读	148
7.7.1 什么是 SMP	126		
7.7.2 SMP 与集群方法	126	第 9 章 大数据 / 快速数据分析中的高性能集成系统、数据库和数据仓库	149
7.7.3 SMP 是否真的重要	126	9.1 引言	149
7.7.4 线程模型	127	9.2 下一代 IT 基础设施和平台的关键特征	150
7.7.5 NumaConnect 技术	127	9.3 用于大数据 / 快速数据分析的集成系统	150
7.8 用于 HPC 的虚拟化	127		
7.9 大型机方面的创新	127		
7.10 FICON 大型机接口	128		
7.11 大型机对手机的支持	129		

9.3.1 用于大数据分析的 Urika-GD 设备	151	10.2.4 先进集群计算系统	189
9.3.2 IBM PureData System for Analytics	152	10.2.5 网格与集群间的差异	189
9.3.3 Oracle Exadata Database Machine	153	10.3 网格计算	190
9.3.4 Teradata 数据仓库和大数据设备	153	10.3.1 网格计算的动机	191
9.4 大数据分析的融合式基础设施	154	10.3.2 网格计算的演进	192
9.5 高性能分析: 大型机 +Hadoop	155	10.3.3 网格系统的设计原则和目标	192
9.6 快速数据分析的 in-memory 平台	158	10.3.4 网格系统架构	193
9.7 大数据分析的 in-database 平台	160	10.3.5 网格计算系统的优点和局限	196
9.8 用于高性能大数据 / 快速数据分析的云基础设施	161	10.3.6 网格系统和应用	196
9.9 用于大数据的大文件系统	164	10.3.7 网格计算的未来	201
9.10 用于大数据 / 快速数据分析的数据库和数据仓库	166	10.4 结论	202
9.10.1 用于大数据分析的 NoSQL 数据库	167	10.5 习题	202
9.10.2 用于大数据 / 快速数据分析的 NewSQL 数据库	169	参考文献	203
9.10.3 用于大数据分析的高性能数据仓库	170	进一步阅读	204
9.11 流分析	173	第 11 章 高性能 P2P 系统	205
9.12 结论	176	11.1 引言	205
9.13 习题	176	11.2 设计原则与特点	206
第 10 章 高性能网格和集群	177	11.3 P2P 系统架构	207
10.1 引言	177	11.3.1 集中式 P2P 系统	207
10.2 集群计算	179	11.3.2 分散式 P2P 系统	208
10.2.1 集群计算的动机	179	11.3.3 混合 P2P 系统	210
10.2.2 集群计算架构	180	11.3.4 高级 P2P 架构通信协议和框架	211
10.2.3 软件库和编程模型	182	11.4 高性能 P2P 应用	212
		11.4.1 Cassandra	212
		11.4.2 SETI @ Home	214
		11.4.3 比特币: 基于 P2P 的数字货币	215
		11.5 结论	216
		11.6 习题	217
		参考文献	217
		进一步阅读	219

第 12 章 高性能大数据分析的 可视化维度	220
12.1 引言.....	220
12.2 常用技术.....	224
12.2.1 图表.....	224
12.2.2 散点图.....	225
12.2.3 树状图.....	226
12.2.4 箱形图.....	226
12.2.5 信息图.....	227
12.2.6 热图.....	227
12.2.7 网络和图的可视化.....	228
12.2.8 词云与标签云.....	228
12.3 数据可视化工具与系统.....	229
12.3.1 Tableau.....	229
12.3.2 Birst.....	231
12.3.3 Roambi.....	232
12.3.4 Qlikview.....	233
12.3.5 IBM Cognos.....	234
12.3.6 Google Charts 和融合表.....	234
12.3.7 Data-Driven Documents (D3.js).....	235
12.3.8 Sisense.....	236
12.4 结论.....	237
12.5 习题.....	237
参考文献.....	238
进一步阅读.....	238
第 13 章 用于组织增权的社交媒体 分析	239
13.1 引言.....	239
13.1.1 社交数据收集.....	239
13.1.2 社交数据分析.....	240
13.1.3 移动设备的发展.....	240
13.1.4 强大的可视化机制.....	240
13.1.5 数据本身的快速变化.....	240
13.2 社交媒体分析入门.....	241
13.3 建立一个用于企业社交媒体 分析的框架.....	242
13.4 社交媒体内容指标.....	243
13.5 社交媒体分析的预测分析技术.....	244
13.6 使用文本挖掘的情感分析架构.....	245
13.7 社交媒体数据的网络分析.....	246
13.7.1 社交媒体数据的网络 分析入门.....	246
13.7.2 使用 Twitter 的网络分析.....	247
13.7.3 极化网络图.....	247
13.7.4 In-Group 图.....	248
13.7.5 Twitter 品牌图.....	248
13.7.6 Bazaar 网络.....	248
13.7.7 广播图.....	248
13.7.8 支持网络图.....	248
13.8 组织的社交媒体分析的不同 方面.....	249
13.8.1 收入及销售的潜在客户 开发.....	250
13.8.2 客户关系和客户体验管理.....	251
13.8.3 创新.....	251
13.9 社交媒体工具.....	251
13.9.1 社交媒体监控工具.....	251
13.9.2 社交媒体分析工具.....	252
13.10 结论.....	252
13.11 习题.....	252
参考文献.....	252
第 14 章 医疗保健的大数据分析	253
14.1 引言.....	253
14.2 影响医疗保健的市场因素.....	254
14.3 不同的相关方设想不同的目标.....	255
14.4 大数据对医疗保健的好处.....	255
14.4.1 医疗保健效率和质量.....	256

14.4.2	早期疾病检测	256	14.11	癌症检测	263
14.4.3	欺诈检测	256	14.12	3D 医学图像分割	263
14.4.4	人口健康管理	257	14.13	新兴医疗方法	264
14.5	大数据技术采纳: 一个新的改进	258	14.14	BDA 在医疗保健方面的用例	264
14.5.1	IBM Watson	258	14.15	人口健康控制	265
14.5.2	IBM Watson 架构	258	14.16	护理流程管理	265
14.6	医疗保健领域中的 Watson	259	14.16.1	核心 IT 功能	265
14.6.1	WellPoint 和 IBM	259	14.17	Hadoop 用例	266
14.7	EHR 技术	259	14.18	大数据分析: 成功案例	268
14.7.1	EHR 数据流	260	14.19	BDA 在医疗保健方面的机会	269
14.7.2	EHR 的优点	261	14.20	Member 360	269
14.8	远程监控和传感	261	14.21	基因组学	269
14.8.1	技术组件	261	14.22	临床监测	271
14.8.2	应用远程监控的医疗 保健领域	261	14.23	BDA 在医疗保健中的经济价值	271
14.8.3	远程监控的局限	262	14.24	医疗保健的大数据挑战	272
14.9	面向医疗保健的高性能计算	262	14.25	医疗保健大数据的未来	273
14.10	人脑网络的实时分析	262	14.26	结论	273
			14.27	习题	273

IT 领域的变革以及未来趋势

1.1 引言

根据大量的报道，IT 领域已经发生了若干可喜的变革以及一些分化。当然，这些变化所带来的后果是多种多样的：灵活的、新一代的特性和功能正在融入现有的以及新兴的 IT 解决方案中；公司和个人正面临着大量的新机会和新可能；新的 IT 产品和解决方案正在以令人难以置信的速度爆发，等等。如同主流市场分析师和研究机构所声称的，有大量颠覆性和革命性的技术正在产生和演化中。例如，Gartner（高德纳，著名市场调研机构）每年都会报告十大技术潮流，这些技术能够给商业组织或大众带来许多微妙的影响。在本章中，为了描述本书的写作背景，将会对 IT 领域的一些相关度最高并最具开创性的趋势进行详细介绍。

有人曾经这样说：IT 领域的第一波浪潮归属于硬件工程。为了满足各种计算、网络、存储的需求，人们细心地设计并集成了各种各样的电子模块（专用的或通用的）。小型化技术带来了大量微米级或纳米级的组件，在硬件的顺利发展中起到了不可或缺的作用。我们即将步入一个计算机无所不在、隐显、用后即弃的时代。IT 领域的第二波浪潮从硬件转移到了软件。从那时起，软件工程开始发挥巨大作用。如今，软件已经变得非常普及且非常有影响力，为人们带来了急需的适应性、可修改性、可扩展性和可持续性，每个有形的事物都通过软件的包裹或嵌入变得智能化。当前 IT 领域处于第三波浪潮之中，这一波浪潮开始于几年前，它是基于对数据（大数据和快速数据）的利用来从硬件和软件的发展中获益。对数据的获取和研究能够产生可行的、及时的洞见（insight），有了这些洞见，就能够实现更聪明的应用程序和设备。

因此，为了通过切实可行的方法实现设想中的智慧地球，数据分析是学习和研究中最令人喜爱和持久的主题。尤其是考虑到异构且分布式的数据源的快速增加，人们对能够满足知识发现和传播目的的数据虚拟化、处理、挖掘、分析和可视化技术情有独钟。数据驱动的洞见使得人们或信息系统能够及时地做出正确的决策。你可以看到席卷 IT 领域的最有前途的趋势就是数据分析，它将给人们带来更好的照顾、选择、便利与舒适。

1.2 新兴的 IT 趋势

IT 消费化 Gartner 的报告详述了移动设备的多样性，包括智能手机、平板电脑、可穿戴设备等。IT 正在日益接近人类，为了个人目的或工作目的，人们能够在任意时间、任意地点、任意设备、任意网络以及任意媒介访问并使用远程拥有的 IT 资源、业务应用和数

据。大量时尚超薄输入/输出设备的生产,使得终端用户能够直接连接到各种IT领域的新产品,并且从中大大受益。IT消费化的趋势已经发展了一段时间,目前达到了巅峰。也就是说,IT正在直接或间接地成为消费者无法避免的部分,而且随着“自带设备”(Bring Your Own Device, BYOD)成为普遍要求,需要能够提供健壮、灵活的移动设备管理软件解决方案。另一方面是在大量垂直业务市场中下一代移动应用及服务出现。在快速变动的移动空间中,有着大量的移动应用程序、地图、服务开发/交付平台、编程及标记语言、架构与框架、工具、容器、操作系统等。准确来说,IT正在由以企业为中心向面向消费者转换。

IT 商品化 商品化是另一个席卷IT业的潮流。随着云计算和大数据分析被广泛接受和采用,IT的商业价值正在急剧上升。代表性的有嵌入式智能正有意识地从硬件封装及装置中抽离出来,从而使得硬件模块能够被大批量地生产并且可以方便快捷地使用。实现这种精细隔离的另一重要需求是基础设施的可负担性,而且供应商锁定这一长期问题目前正在逐步缓解中,任何产品都可以被来自其他厂商的类似设备替代或更换。随着IT基础设施的巩固、集中化和商品化,对商品化硬件的需求激增。IT行业又重新聚焦于各类IT基础设施(服务器、存储设备、网络解决方案,如路由器、交换机、负载均衡器、防火墙网关等)的商品化。通过虚拟化和集装箱化实现的商品化非常普遍且很有说服力。因此,下一代IT环境肯定是软件定义的,从而可以引入大量可编程以及基于策略的硬件系统。

接踵而至的设备时代 硬件工程的主题就是可以看到许多见所未见的创新产品。毫无疑问,IT市场中最近颇受喜爱的就是各类设备。各大主流厂商正将其资金、时间、人才等投入开发下一代智能集成系统(计算、存储、网络、虚拟化以及管理模块)中,这些系统以即用即用的设备形式存在。IT设备是完全定制化的,而且在工厂内就完成了配置,这样当用户使用它们时,只需要几分钟或几小时就可以发挥它们的作用,而不需要几天的时间。为了尽可能多地自动化,生产预集成、预检测、调试好的融合IT栈成为面向设备的主动战略。例如,在IT融合解决方案的比拼中,FlexPod和VCE处于领先地位。类似地,有很多专业的集成系统,例如IBM的PureFlex系统、PureApplication系统以及PureData系统。此外,Oracle公司的工程系统也逐渐在竞争激烈的市场中赢得份额,例如Oracle Exadata Database Machine以及Exalogic Elastic Cloud。

基础设施优化及弹性 整个IT栈会周期性地发生改造,尤其是在基础设施方面,由于传统基础设施的封闭性、僵化性和整体性,很多人正在致力于将传统基础设施改造成模块化、开发性、可扩展、聚合、可编程的基础设施。另一个让人担忧的方面是昂贵IT基础设施(服务器、存储、网络解决方案)的低利用率。随着IT在不同行业将手动任务自动化,IT无序拓展的问题也随之出现,很多IT基础设施利用率不高,有些甚至长时间都不被使用。理解了IT基础设施的这些问题后,有关方面已经采取了大量措施,目的是增加利用率以及优化基础设施。相关的活动包括基础设施的合理化与简化,也就是说,下一代IT基础设施正在通过整合、集中、联合、聚集、虚拟化、自动化、共享的方式实现。为了带来更多的灵活性,最近规定必须采用软件定义基础设施。

随着大数据分析平台及应用程序的快速普及,商用硬件正在快速、廉价地完成数据密集和处理密集型的大数据分析,也就是说,我们需要具有超级计算能力以及无限存储的廉价基础设施。解决方法是将各类利用率低的服务器收集在一起并构建集群,从而形成动态的、巨大的服务器池,以有效满足对与日俱增的、间歇性的计算能力的需求。准确地说,云是能够优雅且经济地满足上述需求的新一代基础设施。云技术尽管并非全新的技术,但是代表了多个成熟