

Greenplum

从大数据战略到实现

冯雷 姚延栋 高小明 杨瑜◎著



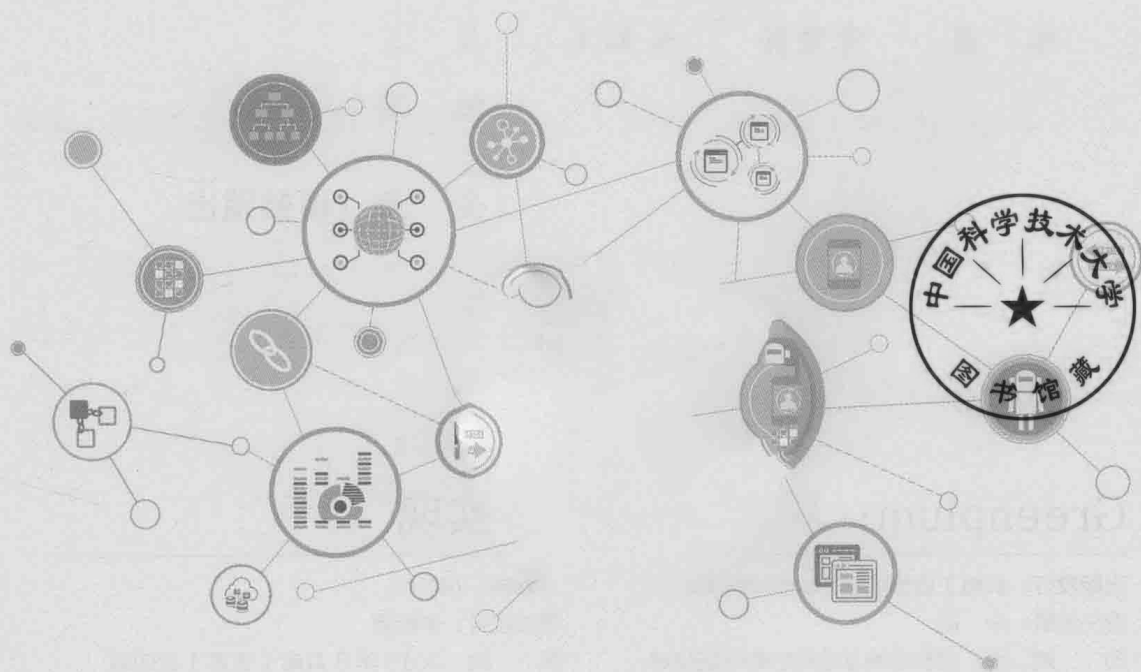


技术丛书

Greenplum

从大数据战略到实现

冯雷 姚延栋 高小明 杨瑜◎著



机械工业出版社
China Machine Press

图书在版编目 (CIP) 数据

Greenplum: 从大数据战略到实现 / 冯雷等著. —北京: 机械工业出版社, 2019.7
(大数据技术丛书)

ISBN 978-7-111-63216-0

I. G… II. 冯… III. 关系数据库系统 IV. TP311.132.3

中国版本图书馆 CIP 数据核字 (2019) 第 143118 号

Greenplum: 从大数据战略到实现

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 劼

责任校对: 李秋荣

印 刷: 大厂回族自治县益利印刷有限公司

版 次: 2019 年 8 月第 1 版第 1 次印刷

开 本: 186mm × 240mm 1/16

印 张: 25.5

书 号: ISBN 978-7-111-63216-0

定 价: 119.00 元

客服电话: (010) 88361066 88379833 68326294

投稿热线: (010) 88379604

华章网站: www.hzbook.com

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

编委会
(按姓氏拼音顺序)

主 编：冯 雷 姚延栋 高小明 杨 瑜

编写组：郭 罡 李 阳 林 文 任振中

王 昊 王淦舟 翁岩青 吴 疆

张 桓

出版统筹：段 旻

序 *Foreword*

“大数据”一词最早出现于 20 世纪 90 年代，作为一个技术术语流行起来则始于 2012 年。时至今日，该词仍没有统一、明确的定义。人们通常从 Volume、Velocity、Variety 等角度定义大数据，而最吸引大众是 Volume 这一特点。根据维基百科的介绍，自 20 世纪 80 年代起，人均存储信息的能力每 40 个月增加一倍；截至 2012 年，全世界每天产生 2.5 艾字节（ 10^{18} 字节）的数据。IDC 报告预测，全球数据将从 2018 年的 33 泽字节（ 10^{21} 字节）增长到 2025 年的 175 泽字节，其中近 30% 数据需要实时处理。世界正在以前所未有的速度数字化和创造数据。数字化时代到来了，数据时代到来了！

随着数据时代的到来，越来越多的企业和政府开始重视大数据及相关技术。2012 年，美国政府宣布投资 2 亿美元拉动大数据相关产业发展，将“大数据战略”上升为国家意志。美国政府将数据定义为“未来的新石油”，并表示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陆权、海权、空权之外的一种国家核心资产。随后多国政府和很多组织提出了相应的大数据战略。

然而，任何行业的升级、发展都不是一蹴而就的。就目前来看，不同企业和组织处于四种不同的数字化和大数据阶段：传统阶段、数字阶段、数据阶段和数学阶段。传统阶段指企业仍然以传统的方式使用软件技术，其主要特点是用软件支撑企业内部流程，通常是由企业内部的 IT 部门主导；数字阶段指企业开始用全新的视角看待和使用软件，软件成为公司主营业务的重要组成部分或者主要组成部分；数据阶段指企业通过全业务的数字化，积累大量数据，再通过数据分析，从数据中获取洞见，反过来促进业务健康发展；数学阶段是指自动化、智能化达到了高阶阶段，通过算法和模型的自动优化为公司提供动力，数学算法和模型成为公司发展的核心引擎。目前来看，虽然大家已经对“大数据”一词耳熟能详，但大多数企业和组织仍然处于传统阶段或者数字阶段早期。造成这一现象的主要原因之一是人才匮乏。根据中国商业联合会数据分析专业委员会统计，未来我国基础性数据分析人才缺口将达到 1400 万，

而大数据专业技术人才缺口将达数百万。优秀的专业书籍对培育人才、缓解这一缺口大有裨益。

目前，市面上的大多数关于大数据的书籍要么侧重于大数据思维，要么侧重于某种或者某几种具体的大数据技术。与这些大数据书籍不同，本书立意新颖，涵盖范围很广，从多个角度对大数据战略到技术进行系统性介绍。本书横向从商业角度介绍了大数据、云计算和人工智能的关系，站在高阶数字化战略的高度解读大数据；纵向从数据处理背后的技术推动力的角度，阐述了大数据发展的历程及未来趋势；从技术实战角度则详细介绍了如何使用 Greenplum 大数据和机器学习平台实现大数据战略。

Greenplum 是先进的开源分布式数据库之一，创建于 2003 年，2010 年被 EMC 公司收购。它因出色的技术能力、易用性和丰富的企业级特性受到大量用户的欢迎，被广泛应用于金融、保险、证券、通信、航空、物流、零售、媒体、医疗、制造、能源等行业，在国内外有一大批拥趸者。2015 年开源后更是发展迅速，目前在全球拥有大量的开源用户。腾讯云等主流的云厂商都将其列为重要的大数据存储、处理和分析服务之一。

本书作者均为 Greenplum 内核开发团队核心成员，在大数据和机器学习行业具有丰富的经验，全球视野和技术前瞻性都毋庸置疑。我也有幸和作者团队多次深度合作，相信他们精心打造的这本书可以给读者全新的启发，帮助大家用正确的理念和方法论来迎接大数据和人工智能时代的挑战与机遇。

祝各位阅读愉快！

王 龙

腾讯云副总裁

2019 年 3 月

前言 Preface

数字原生

2010年11月，在Greenplum创始人的支持下，我们在北京建立了Greenplum中国研发体系。2013年4月，随着Pivotal公司的建立，我们在Greenplum中国研发的基础上合并了部分VMWare中国研发集团的P层云资产，建立了Pivotal中国办公室。截至本书完稿的时候，我们的中国核心研发团队和全球研发团队一起奋斗了8年，打造的Cloud Foundry产品和Greenplum产品成为Pivotal公司在纽约证券交易所上市荣登PaaS第一股的基础。作为Pivotal中国办公室的创始团队，我们一直在审视和提升Pivotal中国办公室的使命和愿景。高尚的使命和愿景是促使一个机构达到世界一流水平的必要条件，因为使命和愿景比战略更高一层。一个机构在前进的过程中，其战略不可避免地需要调整。在面对战略调整时，如果组织成员缺乏共同的使命和愿景，就很难在变化中存活下来。以PC行业为例，苹果公司由最初的苹果电脑公司（Apple Computers）发展到今天苹果（Apple）公司，业务也从以PC为重心迁移到以移动和云服务为重心。苹果公司的转型一路颠簸但最终成功，这与它们坚持艺术和科技的融合并提供一流的用户体验的使命是分不开的。对于不少没有完成转型的PC企业，仔细观察一下，会发现它们通常不能清楚地表达自己的使命。

那么Pivotal中国办公室的使命是什么？简单地说，是支持全球Pivotal产品和商业战略的成功。但是，这个回答显然不能说服和召集一批学霸把Pivotal中国办公室变成世界一流的创新机构。作者有幸参与Pivotal公司在EMC和VMWare内部的启动倡议（Pivotal Initiative），聆听到董事长Paul Maritz先生对Pivotal宣言（Manifesto）的解读。中国读者可能还不熟悉Maritz先生，根据维基百科的介绍，他是微软Windows平台的主要执行团队成员，负责过Windows 95和Windows NT等关键产品。在创建Pivotal之前，Maritz先生是VMWare公司的CEO，奠定了VMWare在虚拟化和I层云的行业领导地位。鉴于Maritz先生在业内的声望，

作者仔仔细细阅读了他撰写的三页纸篇幅的 Pivotal 宣言，并且思考了 Pivotal 中国办公室如何既能拥抱 Pivotal 宣言又能在自己专注的领域成为国内意见领袖。今天，Pivotal 的使命用一句话描述就是“*The Way The Future Gets Built*”，用中文直接翻译过来就是“构建未来的方式”。这句话显得有些抽象，所以在 Pivotal 中国办公室的日常事务中，我们会针对不同的团队来细化这句话：对于面向数字化转型客户的 Pivotal Lab 团队，这句话被表述为“交付一流的数字化转型体验”；对于云研发团队，这句话被表达为“通过 Cloud Foundry 云平台成为云原生平台的行业标杆”；对于数据库研发团队，这句话被阐述为“通过 Greenplum 成为大数据平台和机器学习的意见领袖”。这些使命背后的共同愿景就是提供“数字原生”世界的新产能，以及企业建立数字化所需要的软件平台和方法论。

数字原生就是从由物理世界为重心向数字世界为中心迁移时思考问题的方式。数字计算机发明之前，我们几乎没有什么数字资产和技术。数字计算机发明至今，我们对于数字资产的积累呈指数级增长，在我国更是呈现出跨越式发展的态势。举个例子，今天，如果我们出门不带手机，就会感觉寸步难行，本质上是因为手机已经成为我们进入数字世界的入口。通过手机，我们可以向数字世界发出各种请求，调度物理世界的资源为我们所用。Pivotal 公司喜欢以“ask+综合部门@pivotal.io”的邮件方式来获得综合部门的支持。早期行政部门的同事刚加入 Pivotal 公司的时候常问我：“为什么不面对面请求，或者打个电话，又或者开个单子？”我的回答是这几种方式看似差别不大，但反映了思考问题方式的差别。Pivotal 公司作为数字化的领导者，把软件和数据平台看作数字世界的入口。我们获取资源的方式是向这个数字世界发出请求。数字世界可能通过它的计算找到最优执行路径。有些工作的执行可能还需要转发给人进行人工处理，例如安装一台打印机。但是，有些请求则可以直接通过软件方式解决，例如申请一台云服务器。对于某些请求，虽然我们今天还无法完全以全数字化、无人干预的方式完成，但是，我们可以先把数字原生的框架奠定起来，为以后的进一步对接和持续改进做好准备。在作者看来，数字原生的持续改进过程分为三个阶段：

- 1) 软件公司：通过数字应用实现数字世界和物理世界的无缝交互。
- 2) 数据公司：通过大数据平台实现数据积累和数学模型运行支撑。
- 3) 数学公司：通过数学模型的持续改进来最优化数字世界和物理世界资源。

因此，作者和团队希望能够以三部对应的著作（下面简称为“数字化三部曲”）在数字原生的征程上为读者提供战略参考和对应的软件平台及工具指导。

- 第一三部曲：《Cloud Foundry：从数字化战略到实现》——这本书的主要目标是阐述企业如何实现数字原生第一阶段：实现数字化应用。该书讨论了云计算作为第三代技术平台带来的商业模式变更。在云计算的技术栈中，P 层云带动了企业数字化浪潮。传

统企业通过 P 层云可以迅速获得顶级互联网公司的软件迭代和发布速度，把与客户的交互通过消费级的应用数字化。书中例举福特公司通过 FordPass 建立了以汽车实体产品为核心的一系列用户数字化体验：汽车金融、远程监控车辆、停车位预留、旅途产品和服务推荐等。这个阶段也是一个持续改进的过程。以共享出行为例，今天用户通过手机平台进入数字世界，在打车应用中发送订单。打车平台通过选择最优执行路径，把订单发送给打车平台的司机。然后，司机在物理世界中驱车到达用户起点。随着有辅助的无人驾驶技术的成熟，这个数字世界的运行链条会继续延长，数字平台可以直接把无人车派送到用户起点。在其他的行业，数字应用的链条同样也在持续延长。

- 第二部曲：《Greenplum：从大数据战略到实现》（也就是本书）——我们的主要目的是阐述企业如何实现数字原生的第二阶段：大数据平台。随着数字应用的链条不断延长，企业需要一个大数据平台来积累应用生成的数据。这个工作听上去很容易，因为人们很早以前就使用磁带来存储数据，之后，存储媒介发生了巨大的变化，能够便捷地存储大量数据。那么为何还需要 Greenplum 这样一个大数据和机器学习平台？原因有两个：1）量大；2）快速计算。说到大，当数据量达到 PB 级别（相当于 16000 个 64GB 的 iPhone 中存储的数据）时，企业利用廉价但是可靠的存储来备份和管理是非常困难的。说到快，想象让用户从 16000 个 iPhone 的数据中寻找一张 5 年前的照片就可以感受到大海捞针般的困难；更何况企业的数据库平台要支撑的机器学习和人工智能的数学模型的复杂度要比寻找一张照片的复杂度高几十到几万倍。可见，要想用极快的速度处理如此海量的数据是极其困难的。这也是企业在构建大数据平台时步履维艰的原因。Greenplum 团队的优秀专家用企业积累了 15 年的知识和创新来解决这些难题：如何利用低价的存储设备来实现高可靠的数据存储？数据的存储如何为今天模型的计算做准备？如何给模型提供简单但又标准的接口？数据管理如何在“便于存储”和“便于日后查找”之间取得平衡？如何利用现在的 I 层云计算资源？如何访问文本和地理位置信息等各种数据源？如何访问和计算存储在其他系统（例如 Hadoop）的数据？如何支撑今天主流的人工智能和机器学习模型？我们在创新过程中触碰到了很多计算机科学本身的极限。希望这本著作能给读者呈现一个解决了上述问题并可以实操的大数据平台和战略。
- 我们还在酝酿的第三部著作希望能帮助读者更好地实现数字原生的第三阶段：机器学习和人工智能。企业通过第一阶段和第二阶段的努力捕获和存储了大量的数据。为了更好地理解用户的需求，不少企业进入了更高阶的数字化战略：大数据驱动的机器学习和人工智能。在这个阶段的竞争中，企业会增设一个新的岗位：数据科学家。数据

科学家会在大数据平台上创造和优化数学模型，以期待改进数字世界和物理世界的运作来更好地为人服务。前两部曲提供了软件工具和方法论以帮助企业成为基于大数据的人工智能和机器学习战略的数学公司，不少企业在实践过程中希望作者能够分享实践案例并就企业领导力转变提供咨询。考虑到这样一本著作的出版需要两年以上的时间，碰巧出版社和作者看到了顶级大数据咨询公司 Booz Allen Hamilton 的两位高管收集了大量实际案例的著作《The Mathematical Corporation: Where Machine Intelligence and Human Ingenuity Achieve the Impossible》，其中关于“数学公司”的提法和作者的观点不谋而合。通过出版社的努力，作者和团队把这部著作翻译成中文著作，可以作为第二部曲的伴侣著作来阅读。

虽然数字原生第三阶段的探讨还在创新者和早期用户者群体中进行，但是第二阶段大数据平台的建设已经在中国如火如荼地展开。大数据平台在数字原生三部曲中扮演了承上启下的关键角色，中大型的公司已经将大数据纳入信息平台的建设方案中。Greenplum 因为开源生态和杰出的创新能力被列为方案的候选技术选项，这也使 Pivotal 中国办公室的同事们倍感欣慰。伴随 Greenplum 生态的持续发展壮大，希望这部著作能给企业高层制定战略提供建议和参考，既帮助工程团队开发应用，又能指导运营团队运维和保障。

本书内容组织方式

Greenplum 经过 15 年的精心打磨，成为出色的开源 MPP 数据库和数据处理基础平台，已应用于银行、保险、证券、电信、物流、安保、零售、能源和广告等行业。我们希望本书能给已经建立或者准备建立大数据平台的企业决策者、架构师、开发人员、数据工程师、数据科学家和数据库管理员带来帮助，也希望从事大数据科研工作的教育工作者和学生能从中受益。

本书分为四个部分。

- 第一部分介绍大数据战略。其中，第 1 章将分享作者对于 ABC（人工智能、大数据和云计算）之间关系的理解以及对人和人工智能的思考。第 2 章将介绍进取型企业为什么需要大数据战略以及如何建立大数据战略。
- 第二部分介绍大数据平台。其中，第 3 章将以数据平台演进历史和未来趋势为主题，描述三次整合的背景及影响，介绍选择大数据平台需要考虑的因素，以及为什么 Greenplum 是理想的大数据平台。第 4 章为 Greenplum 数据库快速入门指南。第 5 章将介绍 Greenplum 架构的主要特点和核心引擎。第 6 章将介绍数据加载、数据联邦和数据虚拟化。第 7 章将介绍 Greenplum 的资源管理以及对混合负载的支持。

- 第三部分介绍机器学习与数据分析。其中，第 8 章介绍 Greenplum 的各种过程化编程语言（用户自定义函数），用户可以使用 Python、R、Java 等语言实现用户自定义函数，还可以通过容器化技术实现自定义函数的安全性和隔离性。第 9 章将介绍 Greenplum 内建的机器学习库 MADlib，数据科学家可以使用内建的 50 多种机器学习算法基于 SQL 对数据进行高级分析，并介绍如何扩展 MADlib 以实现新算法。第 10 章和第 11 章将分别介绍 Greenplum 如何对文本数据和时空数据（GIS）进行存储、计算和分析。第 12 章将介绍 Greenplum 丰富的图计算能力。
- 第四部分介绍运维管理和数据迁移。其中，第 13 章将介绍各种监控和管理工具及相关企业级产品。第 14 章介绍数据库备份、恢复和迁移。第 15 章和第 16 章将分别介绍如何从 Oracle 和 Teradata 迁移到 Greenplum。

限于作者学识，本书难免有疏漏之处，恳请同行和各位读者批判指正，我们将不胜感激。您可以通过数字化三部曲的官网（DigitX.cn）或 Greenplum 中文官方社区（greenplum.cn）给我们留言并了解 Greenplum 的技术信息、获得著作的相关学习资源。

冯 雷

Pivotal 中国常务董事兼研发中心总经理

姚延栋

Pivotal 中国研发中心副总裁

序
前 言

第一部分 大数据战略

第 1 章 ABC: 人工智能、大数据和

云计算..... 2

1.1 再谈云计算..... 2

1.1.1 云计算由南向转为北向..... 2

1.1.2 P 层云的精细化发展..... 3

1.1.3 大数据系统在云中部署不断朝南
上移..... 4

1.2 大数据..... 5

1.2.1 从 CRUD 到 CRAP..... 5

1.2.2 MPP (大规模并行计算)..... 7

1.2.3 大数据系统..... 8

1.2.4 当大数据遇到云计算..... 10

1.3 人工智能..... 11

1.3.1 模型化方法..... 12

1.3.2 AI 的发展史..... 14

1.3.3 对 AI 应用的正确预期..... 15

1.4 ABC 之间的关系..... 16

1.5 AI 和人..... 18

1.5.1 经验与逻辑..... 18

1.5.2 公理化的逻辑系统..... 21

1.5.3 图灵机和可计算数..... 25

1.5.4 认知边界上的考量..... 28

第 2 章 建立基于大数据的高阶数字化

战略..... 32

2.1 基于云原生应用的数字化战略..... 32

2.2 大数据和 AI: 企业未来的终极
竞争点..... 34

2.3 大数据战略的落地..... 36

2.3.1 大数据和 AI 人才..... 36

2.3.2 AI 驱动的开发方法和文化..... 37

2.3.3 大数据基础设施的建设..... 39

2.4 大数据和 AI 的展望..... 41

第二部分 大数据平台

第 3 章 数据处理平台的演进..... 45

3.1 前数据处理时代..... 45

3.2 早期的电子数据处理..... 47

3.2.1	电子计算机的出现	47	4.5	Greenplum 数据库的常用操作	82
3.2.2	软件	47	4.6	Greenplum 数据库的常用命令	83
3.3	数据库	49	4.6.1	gpstart	83
3.3.1	数据模型	50	4.6.2	gpstop	83
3.3.2	数据独立性和高级数据处理语言	54	4.6.3	gpstate	83
3.3.3	数据保护	57	4.6.4	gpactivatestandby	84
3.3.4	数据库早期发展过程中的困境	57	4.6.5	gpconfig	84
3.4	NoSQL 数据库	58	4.6.6	gpdeletesystem	84
3.4.1	NoSQL 出现的背景	58	4.7	小结	85
3.4.2	NoSQL 产品的共性	60	第 5 章 Greenplum 的架构和核心引擎		
3.4.3	NoSQL 的分类	61	5.1	Greenplum 的架构	86
3.5	SQL 数据库的回归	62	5.1.1	Greenplum Master	87
3.5.1	NoSQL 与 SQL 的融合	62	5.1.2	Greenplum Segment	87
3.5.2	Hadoop 不等于大数据	63	5.1.3	Greenplum Interconnect	87
3.5.3	SQL 从未离开	64	5.1.4	Greenplum Standby Master	87
3.6	集成数据处理和分析平台	65	5.1.5	Greenplum Mirror Segment	88
3.6.1	数据类型	65	5.2	Greenplum 查询计划	88
3.6.2	业务场景	66	5.2.1	单机查询计划	89
3.6.3	集中还是分散	67	5.2.2	并行查询计划	90
3.7	数据平台的选型	68	5.3	Greenplum 数据库查询处理的过程	95
3.8	小结	69	5.3.1	Greenplum 数据库的主要功能组件	95
第 4 章 Greenplum 数据库快速入门			5.3.2	Greenplum 数据库查询的执行流程	96
4.1	Greenplum 数据库的发展和现状	72	5.4	小结	97
4.2	Greenplum 数据库的特性	73	第 6 章 从 ETL 到数据联邦和数据虚拟化		
4.3	Greenplum 数据库的组成	75	6.1	Greenplum 中的 ETL	99
4.4	Greenplum 数据库的安装与部署	76			
4.4.1	准备工作	76			
4.4.2	安装 Greenplum	77			
4.4.3	初始化 Greenplum 数据库	80			

6.1.1 PostgreSQL 的 ETL 工具箱	99	8.1.3 安装 Python 包	152
6.1.2 GPLOAD	100	8.1.4 安装 Greenplum 数据计算 Python 包集合	153
6.2 Greenplum 的数据联邦	104	8.1.5 类型转换	153
6.2.1 dblink 简介	104	8.1.6 PL/Python 函数中的数据 共享	154
6.2.2 外部表	107	8.2 PL/R	155
6.2.3 GPFDIST 外部表	109	8.2.1 PL/R 简介	156
6.2.4 可执行外部表	119	8.2.2 安装 R 包	158
6.2.5 Greenplum 的 S3 外部表	120	8.2.3 安装 Greenplum 数据计算 R 包 集合	158
6.2.6 GPHDFS 外部表	127	8.3 PL/Container	158
6.2.7 Spark 连接器	129	8.3.1 PL/Container 简介	159
6.2.8 Gemfire 连接器	129	8.3.2 一个简单的例子	159
6.3 Greenplum 的数据虚拟化框架	130	8.3.3 PL/Container 的基本操作 方法	162
6.3.1 PXF 的架构	130	8.3.4 PL/Container 实践总结	166
6.3.2 PXF 的环境配置	131	8.3.5 关于 PL/Container 的开发	167
6.3.3 GPHDFS 与 PXF 比较	132	8.4 小结	167
6.4 小结	133		
第 7 章 混合负载和资源管理	134		
7.1 混合负载的机遇和挑战	134		
7.2 混合负载的业务和技术要求	136		
7.3 资源管理	139		
7.4 并发管理	145		
7.5 小结	146		
第三部分 机器学习与数据分析			
第 8 章 Greenplum 中的过程化编程 语言	149	第 9 章 MADlib 机器学习库	168
8.1 PL/Python	150	9.1 MADlib 入门	168
8.1.1 PL/Python 简介	150	9.1.1 MADlib 简介	168
8.1.2 受信任的过程化编程语言	151	9.1.2 MADlib 的特点	169
		9.1.3 MADlib 与其他机器学习算法库 的比较	172
		9.1.4 MADlib 的快速安装	173
		9.2 MADlib 的架构	174
		9.2.1 SQL 用户接口	174
		9.2.2 Python 驱动函数	175
		9.2.3 C++ 机器学习算法实现	175

9.2.4 C++ 数据库抽象层	176	10.9 GPText 高级查询	207
9.3 MADlib 应用	177	10.9.1 GPText Facet 查询	207
9.3.1 数据预处理	177	10.9.2 GPText 高亮查询结果	209
9.3.2 监督学习	178	10.10 GPText 分区表查询	210
9.3.3 非监督学习	184	10.11 GPText 对自然语言处理的 支持	211
9.3.4 时间序列	187	10.12 GPText 定制化索引	213
9.3.5 自定义机器学习算法	188	10.13 GPText 管理工具	214
9.4 小结	191	10.14 GPText 用于文本挖掘和分析	215
第 10 章 Greenplum 半结构化文本 数据分析	192	10.15 小结	216
10.1 GPText 文本分析概述	192	第 11 章 地理空间数据分析和 处理	218
10.1.1 GPText 数据提取	192	11.1 概述	218
10.1.2 GPText 的文本处理、索引 流程和高阶分析	193	11.1.1 什么是地理空间数据	218
10.2 GPText 内置的全文检索引擎： Apache SolrCloud	194	11.1.2 地理空间数据应用与分析中 的挑战	220
10.3 GPText 架构：高速并行索引和 查询	195	11.2 Greenplum PostGIS	223
10.4 数据准备	197	11.2.1 Greenplum PostGIS 简介	223
10.5 GPText 的使用：简单的 SQL 和 UDF 函数	198	11.2.2 安装 Greenplum PostGIS 组件	224
10.6 GPText 的安装	200	11.2.3 第一次使用	227
10.7 GPText 索引	201	11.3 Greenplum PostGIS 应用实例	228
10.7.1 创建 GPText 索引	201	11.3.1 GIS 数据准备	228
10.7.2 加载 GPText 索引	204	11.3.2 使用 Greenplum PostGIS 空间数据操作符进行 GIS 数据查询	230
10.7.3 GPText 增减索引列	205	11.3.3 使用 Greenplum PostGIS 的 UDF 进行 GIS 数据分析	233
10.8 GPText 简单查询	205	11.3.4 栅格数据	235
10.8.1 GPText 查询的语法	205	11.4 小结	239
10.8.2 GPText 临近查询	206		
10.8.3 GPText top 查询	206		

第12章 Greenplum 数据库与图

计算	240
12.1 图的概念	240
12.2 图的应用	241
12.2.1 电子电路设计自动化	241
12.2.2 搜索引擎	242
12.2.3 社交网络	242
12.3 图数据的处理	243
12.4 Greenplum 对图数据的支持	244
12.5 MADlib 中的图结构和算法	245
12.5.1 图的表示	245
12.5.2 MADlib 支持的图算法	245
12.5.3 MADlib 图算法详解	246
12.6 小结	277

第四部分 Greenplum 的运维和迁移

第13章 Greenplum 的监控和

管理

13.1 监控 Greenplum 集群的状态	282
13.1.1 gpstate 命令	282
13.1.2 系统表 gp_segment_	
configuration	283
13.1.3 Segment 的故障恢复和再	
平衡	284
13.1.4 常用的监控命令	287
13.2 管理 Greenplum 集群	289
13.2.1 参数配置	289
13.2.2 访问管理	290
13.2.3 统计信息	292
13.2.4 管理表膨胀	294

13.3 Greenplum 指令中心 (GPCC)	297
13.3.1 GPCC 简介	297
13.3.2 可视化监控	298
13.3.3 查询监控和分析	301
13.3.4 工作负载管理	305
13.3.5 监报告警系统	307
13.4 小结	309

第14章 Greenplum 数据库的备份、 恢复和迁移

14.1 非并行数据库备份	310
14.2 非并行数据库恢复	313
14.3 并行数据库备份	313
14.4 并行数据库恢复	316
14.5 高效的并行数据库备份和恢复	
工具 gpbackup/gprestore	317
14.6 新一代 Greenplum 数据迁移工具	
GPCOPY	322
14.7 小结	324

第15章 从 Oracle 迁移到 Greenplum

15.1 概述	326
15.2 Oracle 与 Greenplum 的架构	
对比	327
15.2.1 Oracle 的主要痛点	329
15.2.2 Greenplum 的优势	330
15.3 从 Oracle 迁移到 Greenplum 的	
流程	331
15.3.1 迁移场景	332
15.3.2 迁移过程	334

15.3.3 特殊场景分析	344	16.3.3 数据操作语句转换	364
15.4 小结	352	16.3.4 函数转换	367
第 16 章 从 Teradata 迁移到 Greenplum	353	16.3.5 ETL 应用工具连接转换	369
16.1 Teradata 产品和用户面临的 问题	353	16.3.6 其他应用接口迁移	372
16.2 从 Teradata 迁移到 Greenplum 的可行性	354	16.4 特殊场景	373
16.3 如何从 Teradata 迁移到 Greenplum	356	16.4.1 事前微批去重	373
16.3.1 迁移流程概述	356	16.4.2 事后批量去重	374
16.3.2 Teradata 数据卸载及 DDL 导出规范	357	16.5 小结	374
		附录 A Greenplum 社区	375
		附录 B 外部表实例	380
		附录 C Greenplum 的 SSL 证书	386
		术语表	390