

国外计算机科学经典教材

Data Mining
with SQL Server 2005

数据挖掘原理与应用

—— SQL Server 2005 数据库

(美) ZhaoHui Tang 著
Jamie MacLennan
邝祝芳 焦贤龙 高升 译
杨大川 审校



清华大学出版社

国外计算机科学经典教材

数据挖掘原理与应用

—— SQL Server 2005 数据库

(美) ZhaoHui Tang 著
Jamie MacLennan
邝祝芳 焦贤龙 高升 译
杨大川 审校

清华大学出版社

北 京

ZhaoHui Tang and Jamie MacLennan

Data Mining with SQL Server 2005

EISBN: 0-471-46261-6

Copyright © 2005 by John Wiley & Sons, Inc.

All Rights Reserved. This translation published under license.

本书中文简体字版由 John Wiley & Sons, Inc. 授权清华大学出版社出版。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。

北京市版权局著作权合同登记号 图字: 01-2005-6695

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 13501256678 13801310933

图书在版编目(CIP)数据

数据挖掘原理与应用——SQL Server 2005 数据库/(美)唐(Tang, Z. H.), (美)麦克雷南(MacLennan, J.)著; 邝祝芳, 焦贤龙, 高升译.—北京: 清华大学出版社, 2007.1

书名原文: Data Mining with SQL Server 2005

ISBN 978-7-302-14000-9

I. 数… II. ①唐… ②麦… ③邝… ④焦… III. 关系数据库—数据库管理系统, SQL Server 2005
IV. TP311.138

中国版本图书馆 CIP 数据核字(2006)第 121563 号

责任编辑: 王 军 郑雪梅

装帧设计: 孔祥丰

责任校对: 成凤进

责任印制: 杜 波

出版发行: 清华大学出版社 地 址: 北京清华大学学研大厦 A 座

<http://www.tup.com.cn> 邮 编: 100084

c-service@tup.tsinghua.edu.cn

社 总 机: 010-62770175 邮购热线: 010-62786544

投稿咨询: 010-62772015 客户服务: 010-62776969

印 刷 者: 北京鑫丰华彩印有限公司

装 订 者: 三河市新茂装订有限公司

经 销: 全国新华书店

开 本: 185×260 印 张: 24.5 字 数: 536 千字

版 次: 2007 年 1 月第 1 版 印 次: 2007 年 1 月第 1 次印刷

印 数: 1~4000

定 价: 46.00 元

本书如存在文字不清、漏印、缺页、倒页、脱页等印装质量问题, 请与清华大学出版社出版部联系调换。联系电话: (010)62770177 转 3103 产品编号: 020414-01

出版说明

近年来，我国的高等教育特别是计算机学科教育，进行了一系列大的调整和改革，亟需一批门类齐全、具有国际先进水平的计算机经典教材，以适应我国当前计算机科学的教學需要。通过使用国外优秀的计算机科学经典教材，可以了解并吸收国际先进的教学思想和教学方法，使我国的计算机科学教育能够跟上国际计算机教育发展的步伐，从而培养出更多具有国际水准的计算机专业人才，增强我国计算机产业的核心竞争力。为此，我们从国外多家知名的出版机构 Pearson、McGraw-Hill、John Wiley & Sons、Springer、Thomson 等精选、引进了这套“国外计算机科学经典教材”。

作为世界级的图书出版机构，Pearson、McGraw-Hill、John Wiley & Sons、Springer、Thomson 通过与世界级的计算机教育大师携手，每年都为全球的计算机高等教育奉献大量的优秀教材。清华大学出版社和这些世界知名的出版机构长期保持着紧密友好的合作关系，这次引进的“国外计算机科学经典教材”便全是出自上述这些出版机构。同时，为了组织该套教材的出版，我们在国内聘请了一批知名的专家和教授，成立了专门的教材编审委员会。

教材编审委员会的运作从教材的选题阶段即开始启动，各位委员根据国内外高等院校计算机科学及相关专业的现有课程体系，并结合各个专业的培养方向，从上述这些出版机构出版的计算机系列教材中精心挑选针对性强的题材，以保证该套教材的优秀性和领先性，避免出现“低质重复引进”或“高质消化不良”的现象。

为了保证出版质量，我们为这套教材配备了一批经验丰富的编辑、排版、校对人员，制定了更加严格的出版流程。本套教材的译者，全部由对应专业的高校教师或拥有相关经验的 IT 专家担任。每本教材的责编在翻译伊始，就定期不间断地与该书的译者进行交流与反馈。为了尽可能地保留与发扬教材原著的精华，在经过翻译、排版和传统的三审三校之后，我们还请编审委员或相关的专家教授对文稿进行审读，以最大程度地弥补和修正在前面一系列加工过程中对教材造成的误差和瑕疵。

由于时间紧迫和受全体制作人员自身能力所限，该套教材在出版过程中很可能还存在一些遗憾，欢迎广大师生来电来信批评指正。同时，也欢迎读者朋友积极向我们推荐各类优秀的国外计算机教材，共同为我国高等院校计算机教育事业贡献力量。

清华大学出版社

国外计算机科学经典教材

编审委员会

主任委员：

孙家广 清华大学教授

副主任委员：

周立柱 清华大学教授

委员（按姓氏笔画排序）：

王成山	天津大学教授
王 珊	中国人民大学教授
冯少荣	厦门大学教授
冯全源	西南交通大学教授
刘乐善	华中科技大学教授
刘腾红	中南财经政法大学教授
吉根林	南京师范大学教授
孙吉贵	吉林大学教授
阮秋琦	北京交通大学教授
何 晨	上海交通大学教授
吴百锋	复旦大学教授
李 彤	云南大学教授
沈钧毅	西安交通大学教授
邵志清	华东理工大学教授
陈 纯	浙江大学教授
陈 钟	北京大学教授
陈道蓄	南京大学教授
周伯生	北京航空航天大学教授
孟祥旭	山东大学教授
姚淑珍	北京航空航天大学教授
徐佩霞	中国科学技术大学教授
徐晓飞	哈尔滨工业大学教授
秦小麟	南京航空航天大学教授
钱培德	苏州大学教授
曹元大	北京理工大学教授
龚声蓉	苏州大学教授
谢希仁	中国人民解放军理工大学教授

审校者序

接受这本书的审校任务时，我正在美国进行商务访问。期间，遇到了我的好友，本书的作者之一：ZhaoHui Tang。谈起这本书已被译为中文，并很快在国内出版，大家都感到非常的高兴和欣慰。借此机会，我不妨谈谈自己的感想。

数据挖掘，作为商业智能(Business Intelligence)实现的最深层次，在商业智能解决方案的体系中占据着重要的位置。

数据库中存在的只是数据，对于业务人员来说，只是一些无法看懂的天书，没有人会去拿放大镜分析数据库服务器硬盘上的磁轨。他们需要的是信息。那么，我们以前如何解决这个矛盾的呢？一般的答案是报表系统。简单说，业务人员看到的是美观的界面，便捷的操作，鼠标点击后，报表系统生成 SQL 语句，数据库服务器收到以后，返回所需要的信息。不错，报表系统已经可以称作是 BI 了，它是 BI 的低端实现。

现在国外的企业，大部分已经进入了更深一些层次商业智能，叫做数据分析，即基于多维数据库的在线分析系统(OLAP)。还有一些企业已经开始进入更深一个层次商业智能，叫做数据挖掘。从广义上说，任何从数据库中挖掘信息的过程都叫作数据挖掘。从这点看来，报表、多维分析和深度的挖掘都是挖掘数据的手段。但是，从技术术语上说，数据挖掘(Data Mining)特指的是：源数据经过清洗和转换等成为适合于挖掘的数据集，然后建立特定的挖掘模型，利用这些数据集训练模型，最后利用模型找出的知识模式进行预测，从而辅助决策工作。

过去，谈起数据挖掘，大家想到的往往是那些专业数学家、统计学家，一般的技术人员和业务人员望而却步。现在，随着 IT 技术的发展和工业化，SQL Server 2005 提炼了上述的各种复杂知识，加工成友好的视窗工具，嵌入到分析服务(Analysis Services)中，使得数据挖掘的用户扩展到了大量的开发者人群，甚至是经过培训的业务人员。它使我们的员工和程序变得更聪明了。

我在国外学习、工作了多年之后，深感商业智能即将成为未来几年 IT 领域的核心价值，因此从 2003 年开始创建了北京迈思奇科技有限公司，致力于将国外的先进 BI 技术和工具引进国内，帮助国内的企业提高数据分析效率、增强竞争实力。公司成立三年来，与微软密切合作，通过近百次讲座和培训，为企业培养 BI 专业人员；同时，在承担 BI 项目开发的过程中，公司也积累了优秀的团队和丰富的项目案例，创立了国内一流的数据挖掘品牌。

唐朝晖(ZhaoHui Tang)博士是微软公司 SQL Server 2005 数据挖掘产品的产品经理，具有数据挖掘技术的深厚理论根基，同时对该产品有丰富的第一手资料。这本书著成以后，在国外获得了很高的评价。虽然已经多次拜读过英文原著，这次在审校的过程中，

VI 数据挖掘原理与应用

依然感到有新的收获。需要提出的一点是：原书定稿的时候，SQL Server 2005 还是 Beta 版，实际上，我们审校的时候，已经是正式发布版本(加上 SP1 补丁)了。这个细节的变化包括了算法程序的改进，反映在书中，大家会发现不少实际案例的截图与原书是不一样的。这并非翻译和审校的错误。实际上，译书中的所有案例都是我和我公司的同事们基于 SQL Server 2005 SP1 版本全部重新搭建并且重新截图的，付出了相当的精力。希望读者朋友们也能够亲手尝试搭建这些案例，这对于掌握这门技术是很有帮助的。

感谢清华大学出版社及时为我们引进了这部优秀的教材；也感谢本书的译者，准确而清晰地传达了原著的精华。

杨大川

北京迈思奇科技有限公司

译者序

存储技术的迅速发展，特别是硬件价格的下降，使得数据的积累速度不断提高，面对日益庞大的数据资源，我们迫切需要强有力的工具来挖掘其中有用的信息。Microsoft 最新的数据库平台 SQL Server 2005 中的数据挖掘组件是数据挖掘工具的典型代表。

SQL Server 2000 中包括的数据挖掘算法只有决策树算法和聚类算法，与之相比，SQL Server 2005 中引入了多个新的数据挖掘算法，包括贝叶斯算法、时间序列算法、序列聚类算法、关联规则算法和神经网络算法。与传统数据挖掘工具相比，SQL Server 2005 数据挖掘功能具备许多的优势，SQL Server 2005 数据挖掘功能与所有 SQL Server 产品实现了集成，包括 SQL Server、SQL Server Integration Services 和 Analysis Services。SQL Server 2005 数据挖掘功能具有易用性、可伸缩性和可扩展性等特点，同时它包含简单而丰富的 API。

从 SQL Server 2000 到 SQL Server 2005，经历 5 年 SQL Server 数据挖掘的功能实现了一个跨跃式的发展，可谓“五年磨一剑”，也正是本书的 ZhaoHui Tang 和 Jamie MacLennan 两位作者带领他们的团队紧密合作的顶峰。对于本书讲述算法的每一章(包括 4、5、6、7、8、9、10 章)，作者不仅对每一个算法进行了详细的讲述，还在每一章的引言部分给您描述本章算法能应用的一个实际场景，很容易吸引读者往下阅读。

我们作为数据挖掘的研究者，在翻译 SQL Server 2005 数据挖掘时，发觉在 SQL Server 2005 中包含的数据挖掘算法是如此之多，以至于超过了 SAS、SPSS 和 IBM 的 Intelligent Miner 等数据挖掘工具。同时我们作为数据挖掘的工作者，在使用 SQL Server 2005 进行数据挖掘的同时，切身体会到了 SQL Server 2005 数据挖掘功能与所有 SQL Server 产品实现集成给我们的数据挖掘工作带来的便利，以及它的易用性、可伸缩性和可扩展性。

正如本书前言所述，本书对于使用 SQL Server 2005 进行数据挖掘的用户很有用，可谓及时性、实用性、可靠性集于一体。其中，第 1 章对数据挖掘进行了一个基本的介绍，第 2 章讲述了 OLE DB for DM 规范，第 3 章讲述了如何使用 SQL Server 2005 的数据挖掘工具，第 4 章到第 10 章每一章讲述了一个数据挖掘算法，第 11 章讲述了如何对 OLAP 立方体进行数据挖掘，第 12 章讲述了如何使用 SSIS，第 13 章和第 14 章对数据挖掘的体系结构和 API 进行了讲述，第 15 章实现了一个 Web 交叉销售应用程序，第 16 章讲述了如何使用 Microsoft Excel 进行高级的预测，第 17 章讲述了有关扩展 SQL Server 数据挖掘的知识，第 18 章对本书以及 SQL Server 2005 中数据挖掘的功能进行了总结，以及给出了一些附加资源。附录 A 描述了本书使用的 4 个数据集以及如何导入这些数据集；附录 B 描述了 SQL Server 2005 支持的 VBA 函数和 Excel 函数。对于打算采用本书作为

VIII 数据挖掘原理与应用

教材的老师，则建议您讲述第 1 章到第 10 章，以及第 15 章，如果时间允许，您可以讲述第 11、12、17 章，对其他章感兴趣的学生可以自学。

译者

2006 年 5 月于国防科技大学

前 言

数据库系统在过去的 20 年当中取得了巨大的成功。越来越多的数据被收集并且存储在数据库中。一个数据库拥有海量的数据是很平常的事。在这些数据库中找到有用的信息已经成为许多企业面临的重点问题；数据挖掘作为一个挖掘这些信息的关键组件越来越受到人们的重视。数据挖掘算法和可视化工具适用于挖掘数据中的重要模式，并且提供有价值的预测。这种技术实质上适用于各行各业，包括银行、电信、制造业、营销和电子商务。

在 SQL Server 2000 中引入了数据挖掘算法和可视化工具。从此以后，大多数关系数据库系统包含了数据挖掘的功能。在将数据挖掘技术与数据库技术进行集成方面，SQL Server 2005 中的数据挖掘功能实现了一个跨越式的发展，也正是 SQL Server 产品团队和 Microsoft Research 5 年来紧密合作的一个巅峰。来自这两个组织的项目人员和研究人员为 SQL Server 共同开发了经典的、最新的和前沿的数据挖掘工具。本书作者，ZhaoHui Tang 和 Jamie MacLennan，是这两个组织合作的重要驱动者。

对于使用 SQL Server 2005 进行数据挖掘的用户而言，本书将成为他们非常宝贵的参考手册。作者阐述了每个数据挖掘算法的基本原理和各种可视化工具，并且提供了实用的示例。我确信大多数数据库开发人员、数据库管理人员、IT 专业人员和数据挖掘方面的学生都会从本书中获益。

David Heckerman
Research Manager
Microsoft Research, Redmond

目 录

第 1 章 数据挖掘导论	1
1.1 什么是数据挖掘	1
1.2 数据挖掘解决的商业问题	4
1.3 数据挖掘的任务	5
1.3.1 分类	5
1.3.2 聚类	5
1.3.3 关联	6
1.3.4 回归	6
1.3.5 预测	7
1.3.6 序列分析	7
1.3.7 偏差分析	8
1.4 数据挖掘技术	8
1.5 数据流	9
1.6 数据挖掘项目的生命周期	10
1.6.1 第 1 步: 数据收集	10
1.6.2 第 2 步: 数据清理和转换	10
1.6.3 第 3 步: 模型构建	12
1.6.4 第 4 步: 模型评估	12
1.6.5 第 5 步: 报告	13
1.6.6 第 6 步: 预测(评分)	13
1.6.7 第 7 步: 应用集成	13
1.6.8 第 8 步: 模型管理	13
1.7 数据挖掘当前市场与主要厂商	14
1.7.1 数据挖掘市场的大小	14
1.7.2 主要生产厂商和产品	14
1.8 目前存在的问题及挑战	15
1.9 数据挖掘标准	16
1.10 OLE DB for DM 规范和 XML for Analysis 规范	16
1.10.1 用于数据挖掘的 SQL/Multimedia	17

1.10.2 Java 数据挖掘 API	18
1.10.3 预测模型标记语言	20
1.10.4 Crisp-DM 模型	23
1.10.5 公共仓库元数据	24
1.11 数据挖掘的新趋势	25
1.12 本章小结	26
第 2 章 OLE DB for DM 规范	27
2.1 OLE DB 介绍	27
2.2 为什么使用 OLE DB 进行数据挖掘	29
2.3 OLE DB for DM 规范中的基本概念	31
2.3.1 事例	31
2.3.2 事例键	32
2.3.3 嵌套键	32
2.3.4 事例表和嵌套表	33
2.3.5 标量列和表列	33
2.3.6 数据挖掘模型	33
2.3.7 模型创建	33
2.3.8 模型训练	33
2.3.9 模型预测	34
2.4 DMX	34
2.4.1 数据挖掘的 3 个步骤	34
2.4.2 预测函数	43
2.4.3 单例查询	50
2.4.4 仅仅使用内容进行预测	51
2.4.5 钻取模型的内容	52
2.4.6 内容查询	52
2.5 理解模式行集	52
2.5.1 Mining_Services 模式行集	53

2.5.2	Service_Parameters 模式行集	54
2.5.3	Mining_Models 模式行集	54
2.5.4	Mining_Columns 模式行集	55
2.5.5	Mining_Model_Content 模式行集	55
2.5.6	Query_Content 模式行集	58
2.5.7	Mining_Functions 模式行集	59
2.5.8	Model_PMML 模式行集	60
2.6	理解用于挖掘结构的 DMX 扩展	60
2.6.1	挖掘结构	60
2.6.2	挖掘结构的 DMX 扩展	61
2.6.3	Mining_Structure 模式行集	62
2.7	本章小结	63
第 3 章	实践 SQL Server 数据挖掘	65
3.1	BI Dev Studio 介绍	65
3.1.1	理解用户界面	66
3.1.2	脱机模式和即时模式	68
3.2	设置数据源	72
3.2.1	数据源	72
3.2.2	使用数据源视图	74
3.3	创建和编辑模型	83
3.3.1	结构和模型	83
3.3.2	使用数据挖掘向导	83
3.3.3	创建 MovieClick 挖掘 结构和挖掘模型	88
3.3.4	使用数据挖掘设计器	89
3.4	处理	94
3.5	使用模型	96
3.5.1	了解模型查看器	96
3.5.2	使用挖掘准确性图表	98
3.5.3	为 MovieClick 模型 创建一个提升图	101
3.5.4	使用挖掘模型预测	101

3.5.5	针对 MovieClick 模型 执行查询	102
3.5.6	创建数据挖掘报表	103
3.6	使用 SQL Server Management Studio	104
3.6.1	了解 Management Studio 用户界面	105
3.6.2	使用对象资源管理器	106
3.6.3	使用查询编辑器	106
3.7	本章小结	107
第 4 章	Microsoft 贝叶斯算法	109
4.1	贝叶斯算法介绍	109
4.2	理解贝叶斯算法的基本原理	110
4.3	贝叶斯算法的参数	112
4.4	使用贝叶斯算法	113
4.4.1	DMX	114
4.4.2	理解贝叶斯模型的内容	115
4.4.3	浏览贝叶斯模型	117
4.5	本章小结	120
第 5 章	Microsoft 决策树算法	121
5.1	决策树算法介绍	121
5.2	决策树算法的基本原理	122
5.2.1	决策树生成的基本思想	122
5.2.2	处理变量中的多个状态	125
5.2.3	避免过度训练	125
5.2.4	结合先验知识	126
5.2.5	特征选择	126
5.2.6	使用连续的输入属性	127
5.2.7	回归	127
5.2.8	使用 Microsoft 决策树 算法进行关联分析	128
5.3	理解算法参数	129
5.4	使用决策树算法	131
5.4.1	DMX 查询	131
5.4.2	模型内容	135

5.4.3 解释模型	136	8.2 Microsoft 序列聚类算法	
5.5 本章小结	139	基本原理	176
第 6 章 Microsoft 时序算法	141	8.2.1 什么是马尔可夫链	176
6.1 Microsoft 时序算法介绍	141	8.2.2 马尔可夫链的阶	176
6.2 Microsoft 时序算法的		8.2.3 状态转移矩阵	177
基本原理	142	8.2.4 使用马尔可夫链来	
6.2.1 自动回归	142	进行聚类	178
6.2.2 使用多个时间序列	144	8.2.5 聚类分解	180
6.2.3 自动回归树	144	8.3 序列聚类算法的参数	180
6.2.4 季节性	145	8.4 使用序列聚类算法	181
6.2.5 预测历史	146	8.4.1 DMX 查询	181
6.2.6 高速缓存预测	146	8.4.2 模型内容	185
6.3 理解时序算法的参数	147	8.4.3 解释模型	185
6.4 使用 Microsoft 时序算法	148	8.5 本章小结	189
6.4.1 DMX 查询	148	第 9 章 Microsoft 关联规则算法	191
6.4.2 模型内容	152	9.1 Microsoft 关联规则算法介绍	191
6.4.3 模型解释	152	9.2 关联规则算法的基本原理	192
6.5 本章小结	155	9.2.1 理解关联规则算法	
第 7 章 Microsoft 聚类算法	157	的基本概念	192
7.1 Microsoft 聚类算法介绍	158	9.2.2 挖掘频繁项集	195
7.2 聚类算法的基本原理	159	9.2.3 生成关联规则	198
7.2.1 硬聚类算法与软聚类算法	160	9.2.4 预测	198
7.2.2 离散聚类	161	9.3 关联算法的参数	199
7.2.3 可伸缩聚类	162	9.4 使用关联算法	200
7.2.4 聚类预测	163	9.4.1 DMX 查询	200
7.3 聚类算法的参数	163	9.4.2 模型内容	202
7.4 使用聚类模型	166	9.4.3 解释模型	203
7.4.1 将聚类作为一个分析步骤	166	9.5 本章小结	205
7.4.2 DMX	167	第 10 章 Microsoft 神经网络算法	207
7.4.3 模型内容	169	10.1 Microsoft 神经网络算法	
7.4.4 理解聚类模型	169	基本原理	207
7.5 本章小结	174	10.1.1 什么是神经网络	208
第 8 章 Microsoft 序列聚类算法	175	10.1.2 组合和激活	209
8.1 Microsoft 序列聚类算法介绍	175	10.1.3 反向传播、误差函数	
		和共轭梯度	211

10.1.4	处理神经网络的 简单示例	212
10.1.5	规范化和映射	213
10.1.6	网络拓扑	214
10.1.7	训练终止条件	215
10.2	神经网络算法的参数	215
10.3	DMX 查询	216
10.4	模型内容	218
10.5	解释模型	219
10.6	本章小结	221
第 11 章	挖掘 OLAP 立方体	223
11.1	OLAP 介绍	224
11.1.1	理解星型模式和 雪花模式	225
11.1.2	理解维和层次	225
11.1.3	理解度量和度量组	226
11.1.4	理解立方体的处理 和存储	227
11.1.5	使用前缀缓存	228
11.1.6	查询立方体	228
11.2	执行计算	229
11.3	浏览立方体	230
11.4	理解统一维度模型	231
11.5	理解 OLAP 和数据挖掘 之间的关系	234
11.5.1	OLAP 在聚集数据方面 给数据挖掘带来的好处	235
11.5.2	OLAP 需要数据挖掘 来发现模式	235
11.5.3	OLAP 挖掘与关系挖掘	236
11.6	使用向导和编辑器来 构建 OLAP 挖掘模型	237
11.6.1	使用数据挖掘向导	237
11.6.2	构建客户细分模型	237
11.6.3	创建购物篮模型	239

11.6.4	创建销售预测模型	242
11.6.5	使用数据挖掘编辑器	245
11.7	理解数据挖掘维	246
11.8	在 DMX 查询内部 使用 MDX	248
11.9	将 AMO 用于 OLAP 挖掘模型	249
11.10	本章小结	253
第 12 章	SQL Server 集成服务 数据挖掘	255
12.1	SSIS 介绍	255
12.1.1	理解 SSIS 包	257
12.1.2	任务流	257
12.1.3	数据流	259
12.2	在 SSIS 环境中进行 数据挖掘	261
12.2.1	数据挖掘任务	262
12.2.2	数据挖掘转换	267
12.3	本章小结	276
第 13 章	SQL Server 数据挖掘 的体系结构	277
13.1	Analysis Services 体系 结构介绍	277
13.2	XML for Analysis	278
13.2.1	XMLA 的 API	279
13.2.2	XMLA 和 Analysis Services	282
13.3	处理体系结构	283
13.4	数据挖掘管理	284
13.4.1	服务器配置	284
13.4.2	数据挖掘安全	285
13.5	本章小结	287
第 14 章	SQL Server 数据挖掘编程	289
14.1	数据挖掘 API	290

14.1.1	ADO	291	15.4.2	设置权限	328
14.1.2	ADO.NET	291	15.4.3	分析 Web 推荐应用 程序的样例代码	329
14.1.3	ADOMD.NET	291	15.5	本章小结	332
14.1.4	Server ADOMD	292	第 16 章 使用 Microsoft Excel 进行高级预测		333
14.1.5	AMO	292	16.1	针对会话模型来配置 Analysis Services	333
14.2	使用 Analysis Services 的 API	292	16.2	使用高级预测工具	334
14.3	使用 Microsoft.AnalysisServices 创建和管理挖掘模型	293	16.3	ExcelTimeSeries 插件的 体系结构	336
14.3.1	AMO 的基本原理	294	16.4	构建输入数据集	336
14.3.2	AMO 应用程序和安全	295	16.5	创建和训练挖掘模型	339
14.3.3	对象的创建	296	16.5.1	连接数据挖掘引擎	339
14.4	浏览和查询挖掘模型	305	16.5.2	创建和训练	340
14.4.1	使用 ADOMD.NET 来预测	306	16.6	预测序列	342
14.4.2	浏览模型	309	16.7	结合所有代码	343
14.4.3	存储过程	311	16.8	本章小结	346
14.4.4	编写存储过程	312	第 17 章 扩展 SQL Server 数据挖掘		347
14.5	本章小结	317	17.1	理解插件算法	347
第 15 章 实现一个 Web 交叉销售 应用程序		319	17.1.1	插件算法的架构	348
15.1	源数据描述	319	17.1.2	插件算法的概念	348
15.2	构建模型	320	17.1.3	模型的创建和处理	350
15.2.1	确定数据挖掘任务	320	17.1.4	预测	351
15.2.2	将决策树算法应用 于关联任务	320	17.1.5	内容导航	352
15.2.3	使用关联规则算法	322	17.1.6	受托管的插件	352
15.2.4	两个模型的比较	324	17.1.7	安装插件算法	353
15.3	执行预测	325	17.2	使用数据挖掘查看器	353
15.3.1	批处理预测查询	325	17.3	本章小结	354
15.3.2	使用单例预测查询	327	第 18 章 总结与其他资源		355
15.4	在 Web 应用程序中 集成预测功能	327	18.1	重新回顾 SQL Server 2005 数据挖掘的亮点	355
15.4.1	理解 Web 应用程序 的体系结构	327	18.1.1	最新的算法	355
			18.1.2	易于使用的工具	356

18.1.3	简单而强大的 API	356
18.1.4	与同类 BI 技术的集成	357
18.2	探讨数据挖掘的新领域 及应用	357
18.3	延伸阅读	358
18.3.1	Microsoft 数据挖掘 的资源	358
18.3.2	数据挖掘的其他资源	358
18.3.3	流行的数据挖掘 Web 站点	359
18.3.4	流行的数据挖掘会议	359

附录 A	导入数据集	361
A.1	数据集	361
A.1.1	MovieClick 数据集	361
A.1.2	Voting Records 数据集	363
A.1.3	FoodMart 2000 数据集	364
A.1.4	College Plans 数据集	364
A.2	导入数据集	364
附录 B	支持的 VBA 函数和 Excel 函数	369
附录 C	学习资源	373

数据挖掘导论

在当今的商业组织中数据挖掘变得越来越受关注。您可能会常听说，“我们应该利用数据挖掘工具对我们的客户进行细分”，“数据挖掘将会增加客户对我们的满意程度”，甚至还有，“我们的竞争对手正在使用数据挖掘工具获得更多的市场份额——我们必须加油！”

那么，什么是数据挖掘？使用数据挖掘将带来什么好处？如何利用这种技术来解决商业问题？数据挖掘使用了什么技术？一个经典的数据挖掘项目的生命周期是怎样的？在本章中，将回答所有这些问题，并且加以延伸，带领您进入数据挖掘世界。

在本章中，将会学习：

- 数据挖掘的定义
- 确定哪些商业问题可以通过数据挖掘来解决
- 数据挖掘的任务
- 使用各种数据挖掘技术
- 数据挖掘流
- 数据挖掘项目的生命周期
- 当前的数据挖掘标准
- 数据挖掘领域的几种新趋势

1.1 什么是数据挖掘

数据挖掘是商务智能(Business Intelligence, BI)产品系列中的关键成员，BI中的关键成员还包括联机分析处理(Online Analytical Processing, OLAP)、企业报表和 ETL。

数据挖掘指的是分析数据，使用自动化或半自动化的工具来挖掘隐含的模式。在过