

SQL Server 2005 之数据挖掘

唐朝晖
高级项目经理
SQL Server商务智能
Microsoft公司

提纲

- 数据挖掘概述
 - 什么是数据挖掘?
 - 典型商务问题
- SQL Server 2005 系统中的数据挖掘
 - Key messages
 - DMX
 - 新特征介绍
 - 有关工具
 - 分析引擎
 - BI 集成
- 演示
- Q&A

什么是数据挖掘?



数据探索

模式发现

预测处理

典型的商业问题

- 客户的信用风险是什么？
- 我的客户是哪些？
- 客户比较倾向于同时买哪些商品？
- 下个月我能销售多少Pepsi可乐？
- 我的商店的潜在客户会有多少？
- 在互联网上“我”的主要客户群又是哪些？
- ...

Amazon.com: Books: Jack: Straight from the Gut - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print Mail

Address http://www.amazon.com/exec/obidos/tg/detail/-/10446528382/qid=1042753337/sr=8-1/ref=sr_B_1/104-2680659-3995951?v=glance&s=books&n=50784 Go Links

Books

LOOK INSIDE!

jack

List Price: \$29.95

Sign in to turn on 1-Click ordering.

BOOK INFORMATION

buying info
editorial reviews
customer reviews
look inside

RATE THIS ITEM

I dislike it I love it!

1 2 3 4 5

☐ I own it

Submit

Edit your ratings

Favorite Magazines!

Subscribe to other Business & Finance magazines.

Visit the DVD

Automotive.com | Get a Quote. - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites Media Print Mail

Address http://www.automotive.com/fredir/link.html?src_id=1316&kw_sqld=101860&kw_am=m Go Links

AUTOMOTIVE.com

HOME HOW IT WORKS TESTIMONIALS GET A QUOTE



Get a Quote.

"I got the car I wanted at a price I could afford."

Jenna T., TX

Testimonials >>

Get a price in just **minutes**!

Service Available:  

Make Ford

Model Explorer

Zip* 98007

Get a Price

Our FREE **New Car Quote** allows you to select the exact car you want, get invoice pricing and receive the best possible price from an accredited Dealer.

Questions ? Please contact us at info@automotive.com. All rights reserved © Automotive.com

SQL Server 2005 中的数据挖 掘

微软的解决方案

SQL 2000

SQL Yukon

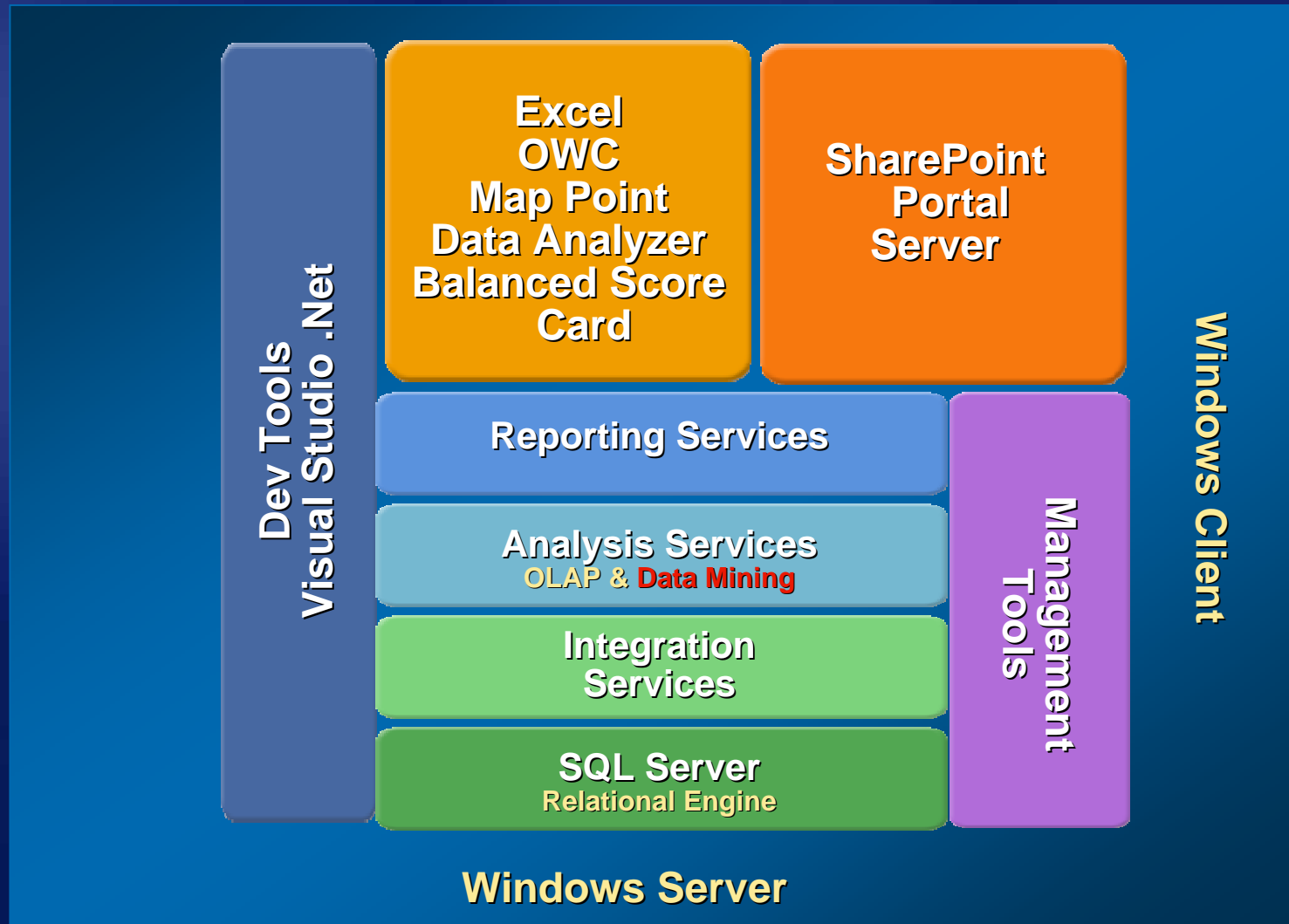
进入DM领域

- 创立工业标准
- 整合开发者建议
- 推出产品、步入市场
- V1.0版本中已包括两个著名算法

成为DM领域领导者

- 面向软件开发人员
- 完善行业标准
- 完备的分析功能
- 向企业级应用渗透

数据挖掘 与 Microsoft 商务智能



Key Messages

■ 嵌入式数据挖掘

- 从DM 到 LOB 都可嵌入应用

■ 集成平台

- DM SQL语言- DMX

- XML/A 支持

■ 高端分析

■ 与ISVs联盟

Key Messages

- 嵌入式数据挖掘

- 集成平台

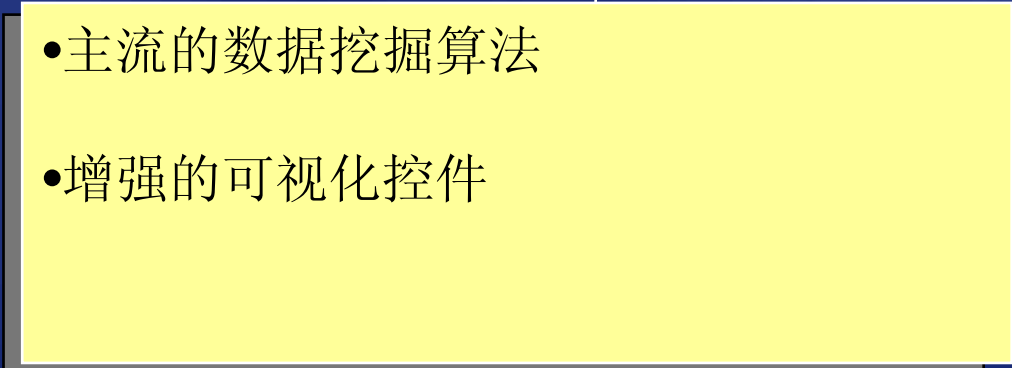
- 高端分析

- 与ISVs联盟

- 与关系数据库(relational), OLAP, DTS和汇报(reporting)技术的紧密集成
- SQL Server: BI的基础平台

Key Messages

- 嵌入式数据挖掘
- 集成平台
- 高端分析
- 与ISVs联盟

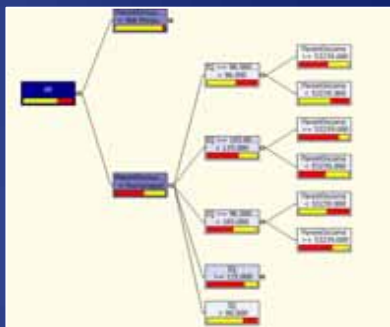
- 
- 主流的数据挖掘算法
 - 增强的可视化控件

Key Messages

- 嵌入式数据挖掘
- 集成平台
- 高端分析
- 与ISVs联盟

- 倾力关注于市场的拓展
- 基于DMX和XML/A的DM工业标准的最终统一
- 第三方算法Plug-in

完备集算法



决策树



聚类



时间序列



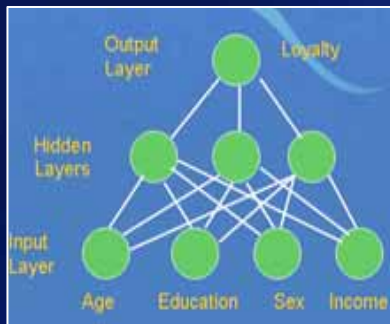
序列聚类



关联

Discrimination scores for Professional/Technical and Service Workers			
Attributes	Values	Favors Professional/Techn...	Favors Service Workers
Education/Years	15-20		
Education/Years	12-13		
Education/Years	7-12		
rielson hit(YOUNG AND THE RES...	Missing		
rielson hit(YOUNG AND THE RES...	Existing		
rielson hit(S THE WORLD TURN...	Existing		
rielson hit(S THE WORLD TURN...	Missing		

Naïve 贝叶斯



神经网络

Text Mining

简单数据

StudentID	Gender	Parent Income	IQ	Encouragement	College Plans
1	Male	23400	120	Not Encouraged	No
2	Female	79200	90	Encouraged	Yes
3	Male	42000	105	Not Encouraged	Yes

复杂数据

Cust ID	Age	Marital Status	IQ	Favorite Movies	
				Title	Score
1	35	M	2	Star Wars	8
				Toy Story	9
				Terminator	7
2	20	S	3	Star Wars	7
				Braveheart	7
				The Matrix	10
3	57	M	2	Sixth Sense	9
				Casablanca	10

API: DMX

CREATE MINING MODEL CreditRisk

(CustID LONG KEY,
Gender TEXT DISCRETE,
Income LONG CONTINUOUS,
Profession TEXT DISCRETE,
Risk TEXT DISCRETE PREDICT)

USING Microsoft_Decision_Trees

INSERT INTO CreditRisk

(CustId, Gender, Income, Profession,
Risk)

Select

CustomerID, Gender, Income,
Profession,Risk

From Customers

Select NewCustomers.CustomerID, CreditRisk.Risk,
PredictProbability(CreditRisk)

FROM CreditRisk **PREDICTION JOIN** NewCustomers

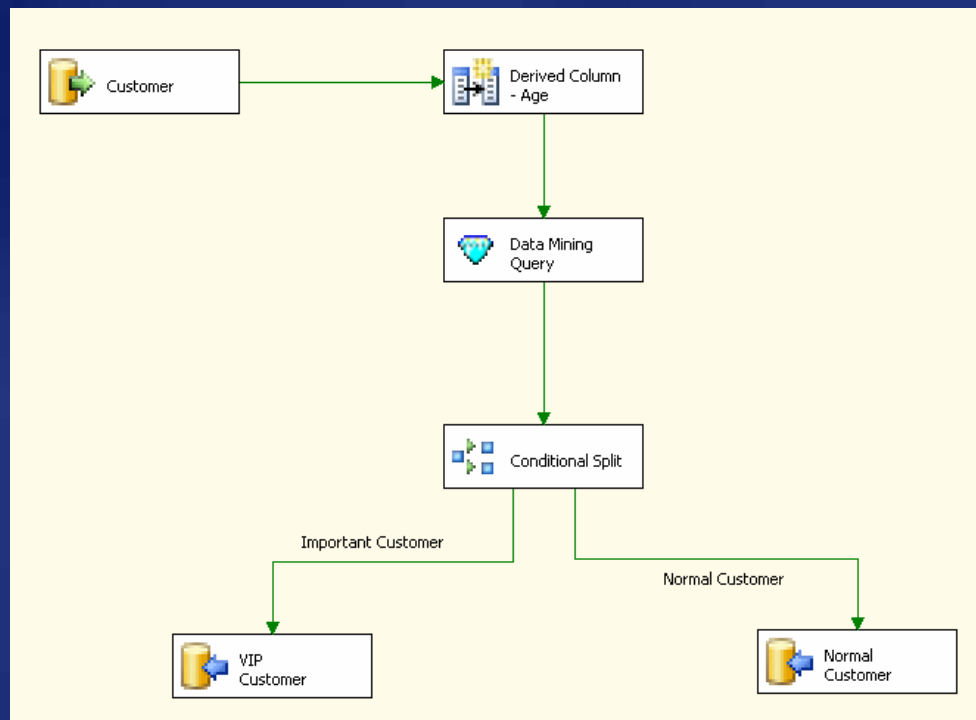
ON CreditRisk.Gender=NewCustomer.Gender

AND CreditRisk.Income=NewCustomer.Income

AND CreditRisk.Profession=NewCustomer.Profession

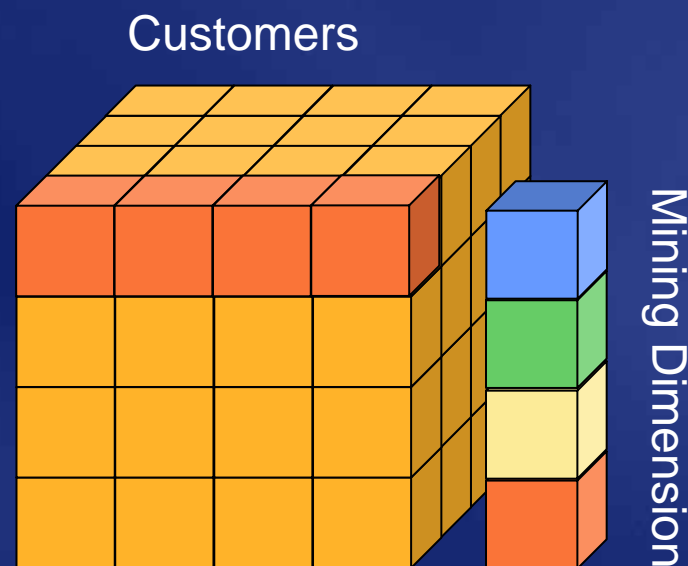
数据挖掘流程

- 与 DTS数据流和任务流的紧密集成
- 数据流
 - 模型训练
 - 预测
 - Text Mining
- 任务流
 - 预测
 - 模型训练



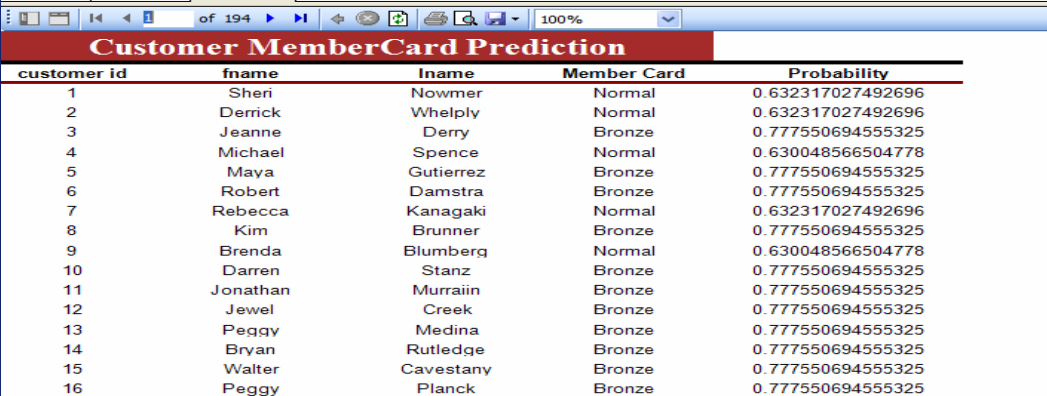
OLAP 挖掘模型增强

- 模型构造上的高灵活性
 - 示例：西北地区上月的市场篮子销量分析
 - 预测下一年的销量
- 数据挖掘维度
- 高性能
- 增强的 OLAP 挖掘型向导和编辑器
- 与关系挖掘模型的致性设计



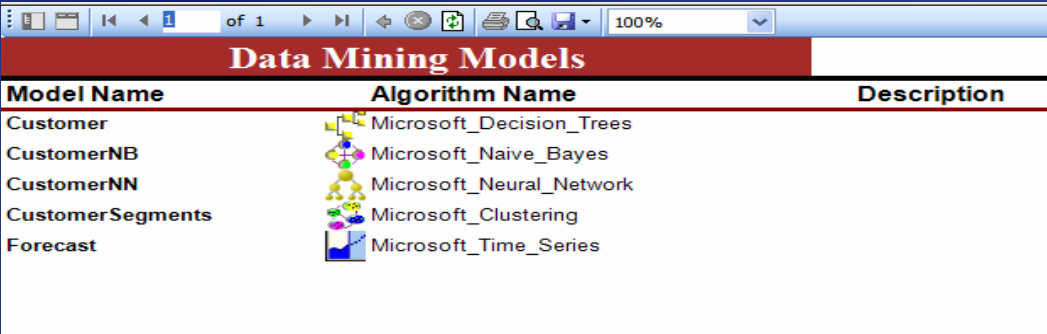
DM报表

- DMX 查询作为 SQL Server 报表数据源
- 报表编译器中嵌入 DMX 查询语言生成器
- Push/Pull mode for DM report delivery



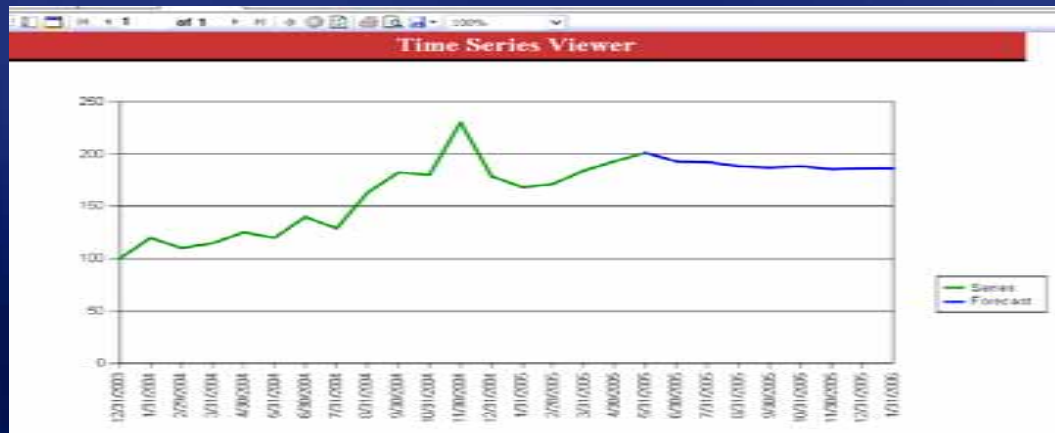
of 194

customer id	fname	lname	Member Card	Probability
1	Sheri	Nowmer	Normal	0.632317027492696
2	Derrick	Whelply	Normal	0.632317027492696
3	Jeanne	Derry	Bronze	0.777550694555325
4	Michael	Spence	Normal	0.630048566504778
5	Maya	Gutierrez	Bronze	0.777550694555325
6	Robert	Damstra	Bronze	0.777550694555325
7	Rebecca	Kanagaki	Normal	0.632317027492696
8	Kim	Brunner	Bronze	0.777550694555325
9	Brenda	Blumberg	Normal	0.630048566504778
10	Darren	Stanz	Bronze	0.777550694555325
11	Jonathan	Murraiin	Bronze	0.777550694555325
12	Jewel	Creek	Bronze	0.777550694555325
13	Peggy	Medina	Bronze	0.777550694555325
14	Bryan	Rutledge	Bronze	0.777550694555325
15	Walter	Cavestany	Bronze	0.777550694555325
16	Peggy	Planck	Bronze	0.777550694555325



of 1

Model Name	Algorithm Name	Description
Customer	Microsoft_Decision_Trees	
CustomerNB	Microsoft_Naive_Bayes	
CustomerNN	Microsoft_Neural_Network	
CustomerSegments	Microsoft_Clustering	
Forecast	Microsoft_Time_Series	



演示

更多信息...

- DM URL
 - www.sqlserverdatamining.com
- 新闻组:
 - Microsoft.Public.SQLserver.Datamining
- 更多:
 - KDNugget.com
 - www.MineSage.com - 迈思奇

Thank You!