



吉林大学哲学社会科学学术文库

# 商务智能 概念、方法及在管理中的应用

**Business Intelligence  
Concept, Method and Application in Management**

刘伟江 / 著



吉林大学哲学社会科学学术文库

# 商务智能 概念、方法及在管理中的应用

Business Intelligence  
Concept, Method and Application in Management

刘伟江 / 著

## 图书在版编目(CIP)数据

商务智能：概念、方法及在管理中的应用 / 刘伟江著 —北京：  
社会科学文献出版社，2012.1  
(吉林大学哲学社会科学学术文库)  
ISBN 978 - 7 - 5097 - 2947 - 2  
I ①商· II ①刘· III ①电子商务 - 研究 IV ①F713.36  
中国版本图书馆 CIP 数据核字 (2011) 第 253825 号

· 吉林大学哲学社会科学学术文库 ·

### 商务智能

——概念、方法及在管理中的应用

著 者 / 刘伟江

出 版 人 / 谢寿光

出 版 者 / 社会科学文献出版社

地 址 / 北京市西城区北三环中路甲 29 号院 3 号楼华龙大厦

邮 政 编 码 / 100029

责 任 部 门 / 财经与管理图书事业部 (010) 59367226

责 任 编 辑 / 陶璇

电 子 信 箱 / caijingbu@ssap.cn

责 任 校 对 / 孙光 远

项 目 统 筹 / 恽薇 陶璇

责 任 印 制 / 岳 阳

总 经 销 / 社会科学文献出版社发行部 (010) 59367081 59367089

读 者 服 务 / 读者服务中心 (010) 59367028

印 装 / 三河市又通印刷包装有限公司

印 张 / 13

开 本 / 787mm × 1092mm 1/16

字 数 / 163 千字

版 次 / 2012 年 1 月第 1 版

印 次 / 2012 年 1 月第 1 次印刷

书 号 / ISBN 978 - 7 - 5097 - 2947 - 2

定 价 / 39.00 元

本书如有破损、缺页、装订错误，请与本社读者服务中心联系更换

【】 版权所有 翻印必究

## | 前 言 |

随着信息技术的不断发展，人们拥有大量的数据，如何运用这些数据并从中挖掘出有意义的信息和知识，进而将之作为人们决策的依据是很多人关注的重点。基于此，本书力图通过介绍与商务智能相关的内容来达到这一目标。

本书注重理论与实践相结合，以商务智能的三个组成部分——数据仓库、联机分析和数据挖掘为主线，分别介绍了各部分所涉及的相关概念、方法和技术。本书的特点是：

第一，面向管理中的问题，注重实际应用。数据仓库、联机分析和数据挖掘都是应用性很强的技术。本书在对这些方法和技术进行介绍的过程中，力图结合管理中的实际问题，以期为企业的决策提供有价值的参考。

第二，可操作性强。本书在介绍相关技术时，根据网上公开的 Foodmart 2000. mdb 数据集中的数据，在 SQL Server 2005 操作环境下，作了丰富的图例和操作讲解，读者可以模拟本书的操作步骤，并把这些技术方法应用到自己需要处理的问题中。

本书主要分为四部分内容。

第一部分为概述，是第 1 章的内容，这章对商务智能的概念、发展等作了相应的介绍。

第二部分为数据仓库和联机分析，包含第 2、3、4 章的内容，

在这部分中主要对数据仓库的建立、数据清理、数据查询等做了说明。

第三部分为数据挖掘，包含第5、6、7章的内容，这部分主要对数据挖掘方法进行了详细的介绍，涉及决策树、神经网络、关联规则、聚类等方法。

第四部分为商务智能在管理中的应用，主要是第8章的内容。在这部分中，根据商业企业在实际中可能遇到的商业问题，通过选择一定的数据挖掘方法来进行处理，并对得到的数据结果进行了分析说明。

由于作者的水平和经验有限，有些方面的研究还有待深入和提高，难免有错讹之处，真诚欢迎广大读者批评指正。

刘伟江

2011年9月于长春

# 目 录

## CONTENTS

前 言 .....	1
<b>第1章 概述 .....</b>	<b>1</b>
1.1 商务智能简介 .....	1
1.1.1 商务智能概念 .....	1
1.1.2 商务智能的发展 .....	3
1.1.3 从数据处理的角度看商务智能的组成 .....	6
1.2 为什么需要商务智能 .....	7
1.3 商务智能工具 .....	8
<b>第2章 数据仓库 .....</b>	<b>12</b>
2.1 数据仓库概述 .....	12
2.1.1 数据仓库的概念及特点 .....	12
2.1.2 数据库与数据仓库的区别 .....	14
2.1.3 数据仓库的技术支持 .....	16
2.2 数据仓库的设计 .....	16
2.3 数据仓库的构建实例——以 Foodmart 2000. mdb 数据集为例 .....	24

<b>第3章 数据预处理</b>	38
3.1 为什么需要预处理数据	38
3.2 数据清理	40
3.2.1 空缺值处理	40
3.2.2 异常值检测	41
3.2.3 重复记录检测	42
3.3 数据集成	43
3.4 数据变换	44
3.5 数据归约	47
<b>第4章 联机分析处理</b>	50
4.1 OLAP 的概念与特点	50
4.1.1 OLAP 的概念	50
4.1.2 OLAP 的特点	51
4.1.3 OLTP 和 OLAP 的对比	52
4.2 OLAP 的一些基本概念	53
4.3 OLAP 的分类	55
4.4 OLAP 的基本操作	57
4.5 OLAP——以 Foodmart 2000. mdb 数据集中库存数据表 等相关数据为例	61
<b>第5章 分类</b>	68
5.1 分类的概念	68
5.2 决策树分类	69
5.2.1 基本概念	69
5.2.2 决策树的生成过程	69

5.2.3 决策树停止的条件 .....	74
5.2.4 决策树的修剪 .....	77
5.2.5 决策树的评估 .....	80
5.3 贝叶斯分类 .....	83
5.4 人工神经网络分类 .....	85
5.4.1 人工神经网络概述 .....	85
5.4.2 神经元的数学模型 .....	86
5.4.3 人工神经网络模型 .....	87
5.4.4 神经网络拓扑结构的确定 .....	89
5.5 分类过程中面临的问题——不均衡数据集 .....	90
5.6 其他分类方法 .....	91
5.6.1 k - 最近邻居法 .....	92
5.6.2 粗糙集分类法 .....	94
5.7 Microsoft 分类挖掘模型的操作过程——以基于 决策树的客户分类为例 .....	97
<b>第6章 关联规则 .....</b>	<b>111</b>
6.1 关联规则简介 .....	112
6.2 关联规则的分类 .....	113
6.3 由事务数据库挖掘单维关联规则 .....	115
6.3.1 Aprior 算法 .....	115
6.3.2 频繁模式增长 .....	118
6.4 关联规则的推广 .....	122
6.4.1 多层关联规则 .....	122
6.4.2 多维关联规则 .....	124
6.5 时序关联规则 .....	125
6.6 商品关联关系分析——以 Foodmart 2000. mdb 数据集中	

1997 年销售数据为例 .....	128
<b>第 7 章 聚类 .....</b>	<b>144</b>
7.1 简介 .....	144
7.2 聚类分析算法 .....	146
7.2.1 K - 均值簇算法 .....	146
7.2.2 EM 算法 .....	148
7.3 聚类分析的应用 .....	152
7.4 聚类分析的操作过程——基于客户价值的聚类分析 ..	153
<b>第 8 章 商务智能在管理中的应用 .....</b>	<b>163</b>
8.1 基于决策树的职员职位影响因素研究 .....	163
8.2 基于聚类方法的广告效应差异分析 .....	172
8.3 基于贝叶斯方法和决策树方法的顾客分类效果 比较研究 .....	179
8.4 基于聚类方法的顾客特征分析 .....	186

# 目 录

## C O N T E N T S

<b>Preface</b>	/ 1
<b>Chapter 1 Introduction</b>	/ 1
1. 1 What is Business Intelligence	/ 1
1. 1. 1 Business Intelligence Concept	/ 1
1. 1. 2 The Development of Business Intelligence	/ 3
1. 1. 3 The Form of Business Intelligence	/ 6
1. 2 Why We Need Business Intelligence	/ 7
1. 3 Business Intelligence Tools	/ 8
<b>Chapter 2 Data Warehouse</b>	/ 12
2. 1 Introduction	/ 12
2. 1. 1 Data Warehouse's Concept and Characteristics	/ 12
2. 1. 2 Differences between Database and Data Warehouse	/ 14
2. 1. 3 Data Warehouse Technical Support	/ 16
2. 2 Data Warehouse Design	/ 16
2. 3 Data Warehouse Construction—Taking Foodmart 2000. mdb as an Example	/ 24

<b>Chapter 3 Data Preprocessing</b>	/ 38
3. 1 Why Need Preprocessing	/ 38
3. 2 Data Cleaning	/ 40
3. 2. 1 Missing Values	/ 40
3. 2. 2 Abnormal Values	/ 41
3. 2. 3 Duplicate Values	/ 42
3. 3 Data Integration	/ 43
3. 4 Data Transformation	/ 44
3. 5 Data Reduction	/ 47
<b>Chapter 4 Online Analytical Processing</b>	/ 50
4. 1 OLAP Definition and Characteristics	/ 50
4. 1. 1 OLAP Definition	/ 50
4. 1. 2 OLAP Characteristics	/ 51
4. 1. 3 Differences between OLTP and OLAP	/ 52
4. 2 OLAP Basic Concepts	/ 53
4. 3 OLAP Classification	/ 55
4. 4 OLAP Operations	/ 57
4. 5 OLAP—Taking Foodmart 2000. mdb Inventory Data as an Example	/ 61
<b>Chapter 5 Classification</b>	/ 68
5. 1 What is Classification?	/ 68
5. 2 Decision Tree	/ 69
5. 2. 1 Basic Concepts	/ 69
5. 2. 2 Decision Tree Induction	/ 69
5. 2. 3 The Stopping Condition of Decision Tree	/ 74

5.2.4 Decision Tree Pruning	/ 77
5.2.5 Decision Tree Evaluation	/ 80
5.3 Bayes Classification	/ 83
5.4 Neural Network	/ 85
5.4.1 What is Neural Network?	/ 85
5.4.2 Mathematical Model of Neurons	/ 86
5.4.3 Neural Network Model	/ 87
5.4.4 Defining a Network Topology	/ 89
5.5 Problems in Classification—Unbalanced Data Set	/ 90
5.6 Other Classification Methods	/ 91
5.6.1 k – Nearest Neighbor Classifiers	/ 92
5.6.2 Rough Set Approach	/ 94
5.7 Operating Processing—Taking Classifying Customers by Decision Tree as an Example	/ 97
<b>Chapter 6 Association Rules</b>	/ 111
6.1 What is Association Rules?	/ 112
6.2 Association Rules Classification	/ 113
6.3 Mining Single – Dimensional Association Rules from Transactional Databases	/ 115
6.3.1 Aprior Agrithorim	/ 115
6.3.2 Frequent Pattern Growth	/ 118
6.4 Association Rules Extension	/ 122
6.4.1 Multilevel Association Rules	/ 122
6.4.2 Multidimensional Association Rules	/ 124
6.5 Time Series Association Rules	/ 125
6.6 Association Relationship Analysis among Commodities	

—Taking Foodmart 2000. mdb 1997 Sales Data as an Example	/ 128
<b>Chapter 7 Clustering Analysis</b>	/ 144
7. 1 Introduction	/ 144
7. 2 Clustering Analysis Algorithm	/ 146
7. 2. 1 K – Means Clustering Algorithm	/ 146
7. 2. 2 EM Algorithm	/ 148
7. 3 Clustering Analysis Application	/ 152
7. 4 Operating Processing—Customers Values Clustering Analysis based on Clustering	/ 153
<b>Chapter 8 Business Intelligence Application in Management</b>	/ 163
8. 1 Employee Position Influence Factors Research based on Decision Tree	/ 163
8. 2 Advertise Effect Differences Analysis based on Clustering Analysis	/ 172
8. 3 Customers Classification Results Comparative Analysis by Using Decision Tree and Bayes	/ 179
8. 4 Customers Characteristics Anaylsis based on Clustering Analysis	/ 186

## 1.1 商务智能简介

### 1.1.1 商务智能概念

“啤酒与尿布”的故事最初是美国学者 Agrawal 从沃尔玛商店中大量的顾客消费数据中，利用计算机通过商品关联关系的计算方法——Aprior 算法<sup>[1]</sup>分析得到的一种现象：在某些特定的情况下，啤酒和尿布这两种看起来完全不相干的商品会出现在一个购物篮中，这一现象的发现有助于商家安排商品的合理摆放和选择相应的促销活动。沃尔玛随后对啤酒和尿布进行了捆绑销售，不出意料，两种商品的销售量双双增加。这个有趣的故事带给了人们深深的思考：如何利用信息技术收集数据并采取有效的方法处理这些大量的、有噪声的、不完全的数据，进而从中找出未知的、有用的信息和知识来帮助人们作决策？因此，商务智能（Business Intelligence，简称 BI）作为一种能帮助人们将大量数据转化成有价值的信息和知识的有效过程越来越引起人们的广泛重视，很多公司和个人都根据自己的理解给出了商务智能的定义。

IBM 公司对商务智能的定义是：商务智能是指利用已有的数据资源作出更好的商业决策，它包括数据访问、数据和业务分析及发现新的商业机会。该定义认为商务智能可以从根本上帮助企业把其运营数据转化为高价值的可以获取的知识（或信息），并且在恰当的时候通过恰当的方式把恰当的信息传给恰当的人。Gartner Group 的 Howard Dresner 在 1989 年将商务智能定义为一类由数据仓库（或数据集市）、报表查询、联机分析、数据挖掘等部分组成，以帮助企业决策的技术及应用。该定义的基本层次结构如图 1-1 所示。

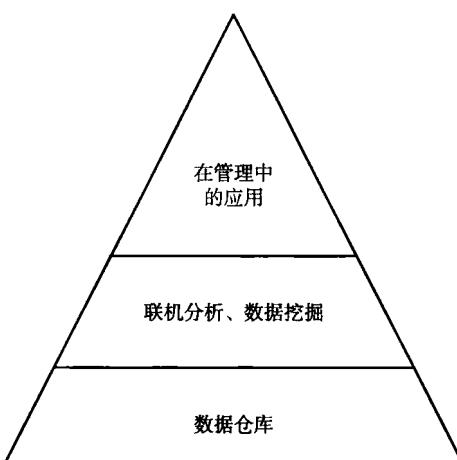


图 1-1 商务智能的基本层次结构

从图 1-1 可以看出，数据仓库中的数据是商务智能的基础，一般来说，数据分为三类：第一类是结构化数据，如数字、符号等，这类数据可以存储在关系型数据库里，用二维逻辑表来表现；第二类是半结构化数据，如文本、网页等；第三类是非结构化数据，如图像、视频、音频等。联机分析、数据挖掘等是商务智能的核心，其中联机分析（Online Analysis Process，简称 OLAP）技术是用来帮助分析人员、管理人员从多种角度把原始数据中转化

出来的、能够真正为用户所理解的，并真实反映数据维特性的信息进行快速、一致、交互的访问，从而让用户可以从各个角度更深入地了解数据的一类软件技术；数据挖掘（Data Mining，简称DM）是与数据仓库密切相关的一种信息技术，近几年来，很多专家和学者从不同的角度给出了数据挖掘的定义，如从技术角度，数据挖掘（Data Mining）就是从大量的、有噪声的、模糊的、随机的实际应用数据中，提取隐含在其中的、人们事先不知道的，但又是潜在有用的信息和知识的过程。该定义包括好几层含义：数据源中的数据是真实的、大量的、含噪声的；发现的内容是用户感兴趣的知识，这些知识是可接受的、可理解的和可运用的；在知识发现过程中，并不要求发现放之四海皆准的知识，这些知识仅支持特定问题的处理就可以。从商业角度来讲，数据挖掘是一种新的商业信息处理技术，其主要特点是对商业数据库中的大量操作型数据进行抽取、转换、分析和其他模型化处理，从中提取辅助商业决策的关键性信息。总之，数据挖掘就是要从大量数据中提取出隐藏在其中的有用信息。

正因为如此，近年来商务智能在管理中有着极其广泛的应用，商务智能的应用问题主要从以下几个方面进行考虑：一是要解决什么样的问题；二是围绕要解决的问题进行相关数据的收集；三是选取合适的数据挖掘方法对数据进行分析；四是分析的结果有个合理的说明并提出切实可行的建议。

### 1.1.2 商务智能的发展

商务智能的发展主要经历了以下几个阶段<sup>[2]</sup>：

#### 1. 萌芽期

1947年卡内基梅隆大学的赫伯特·西蒙（Herbert Simon）教授在《行政组织的决策过程》一书中提出：如果能利用存贮

在计算机里的信息来辅助决策，人类理性的范围将会大大扩大。他的这个观点奠定了决策支持系统的基础，而他对决策支持系统的研究被学术界公认为现代商务智能概念最早的源头和起点。

## 2. 发展期

随着决策支持系统的发展，如何把多个不同数据源的数据有机地整合起来是决策支持系统面临的难题。1988年，为解决企业中的数据集成问题，IBM公司的研究员Barry Devlin和Paul Murphy创造性地提出了一个新的术语：数据仓库（Data Warehouse）。但IBM公司只是把数据仓库这个词当做一个新的概念来宣传，而没有进一步提出实际的架构和设计。1992年，比尔·恩门（Bill Inmon）在《如何构建数据仓库》一书中首次给出了关于数据仓库的清晰定义和操作性极强的建议，从而为数据仓库的广泛应用奠定了基础，也为商务智能的发展提供了支撑。随后的联机分析技术的发展使得数据仓库有了更广泛的用武之地，它使人们可以从多个视角、多个维度观察数据，从而清楚地了解事情的发生过程。

## 3. 成熟期

随着数据仓库、联机分析技术的发展和成熟，商务智能的框架基本形成，但真正给商务智能带来活力的是数据挖掘技术的出现。数据挖掘技术以其能通过对大量数据的分析来揭示数据之间隐藏的、有意义的关系、模式和趋势，为决策者提供新的信息和知识的能力，使得商务智能真正有了“智能”的内涵。

进入21世纪以来，随着技术的发展，信息可视化使得商务智能的应用又有了更进一步的发展。所谓信息可视化（Information Visualization）是指以图形、图像、动画等更为生动、易于理解的方式来展现和诠释数据之间的复杂关系和发展趋势，以