

高等学校信息管理与信息系统专业系列教材

数据仓库与数据挖掘 的原理及应用

李志刚 马 刚 编著



高等教育出版社
HIGHER EDUCATION PRESS

TP311. 13/333

2008

高等学校信息管理与信息系统专业系列教材

数据仓库与数据挖掘
的原理及应用

李志刚 马 刚 编著

高等教育出版社

内容提要

本书详细阐述数据仓库与数据挖掘的基本原理，系统而全面地介绍数据仓库与数据挖掘的概念、作用、算法以及应用领域、相关学科和发展趋势，并着重讨论数据仓库和数据挖掘在企业管理中的应用及构建策略。基于 SQL Server 2005 介绍数据仓库与数据挖掘工具的操作和应用，并结合具体实例，阐述企业数据仓库和数据挖掘的实施过程。最后，以证券行业为对象提供一个数据挖掘的开发实例。本书的指导思想是在系统阐述基本知识和基本理论的基础上，强调实际应用能力的培养，充分体现数据仓库和数据挖掘技术作为解决实际问题的方法和工具的特点。

本书既可以作为信息系统、电子商务、管理科学与工程、计算机应用、软件工程等专业的本科高年级和研究生教材，又可以作为从事竞争情报、信息管理、知识管理、战略管理和软科学的研究人员的参考资料。

图书在版编目 (CIP) 数据

数据仓库与数据挖掘的原理及应用 / 李志刚, 马刚
编著. —北京：高等教育出版社，2008.2

ISBN 978 - 7 - 04 - 023014 - 7

I . 数… II . ①李… ②马… III . ①数据库系统 -
高等学校 - 教材 ②数据采集 - 高等学校 - 教材 IV .
TP311. 13 TP274

中国版本图书馆 CIP 数据核字 (2008) 第 000470 号

策划编辑 刘艳 责任编辑 康兆华 封面设计 于文燕 责任绘图 郝林
版式设计 张岚 责任校对 金辉 责任印制 陈伟光

出版发行	高等教育出版社	购书热线	010 - 58581118
社址	北京市西城区德外大街 4 号	免费咨询	800 - 810 - 0598
邮政编码	100011	网 址	http://www.hep.edu.cn
总机	010 - 58581000	网上订购	http://www.landraco.com
经 销	蓝色畅想图书发行有限公司	畅想教育	http://www.landraco.com.cn
印 刷	北京奥鑫印刷厂		http://www.widedu.com
开 本	787×960 1/16	版 次	2008 年 2 月第 1 版
印 张	22.25	印 次	2008 年 2 月第 1 次印刷
字 数	410 000	定 价	27.80 元

本书如有缺页、倒页、脱页等质量问题，请到所购图书销售部门联系调换。

版权所有 傲权必究

物料号 23014 - 00

前　　言

.....

随着计算方法和信息技术的不断发展,大量数据的产生和收集导致信息爆炸。现代社会的竞争趋势要求对这些数据进行实时的和深层次的分析。虽然现在已经出现更强大的存储系统和检索系统,但是使用者发现在分析所拥有的信息方面变得越来越困难。数据仓库提供了容纳大量信息的场所,但它只有和数据挖掘技术相结合才能最终解决用户的困惑,使用户能够从繁杂的数据中找出真正有价值的信息和知识。数据仓库可以加强企业对信息的管理能力,数据挖掘可以改善企业的经营状况,使企业的决策制定过程更加科学化和快速,为企业带来巨大的收益,增强企业的竞争优势。

数据仓库和数据挖掘是 20 世纪 90 年代中期兴起的决策支持新技术,它们是基于大规模数据库的决策支持系统的核心。数据仓库是区别于数据库的一种新型数据存储形式,它是面向主题的、集成的、不可更新的、随时间不断变化的数据集合,用以支持经营管理中的决策制定。数据挖掘是从数据库中发现知识的核心技术,它从大量的数据中提取隐含的、人所未知的、可信而有效的知识。数据挖掘能够对数据进行再分析,以期获得更加深入的了解。它具有预测功能,可以通过已有数据预测发展趋势。数据仓库与数据挖掘技术相结合,与现代的管理决策方法相结合,就能够使数据仓库在组织机构的经营管理决策中发挥巨大的作用。

我国数据挖掘技术的研究始于 20 世纪 90 年代,经过十几年的发展,这一领域目前正处于蓬勃发展时期。但是由于数据仓库、数据挖掘技术都是数据处理与分析领域出现的新技术,大部分人把目光投向基于这两项技术的基础理论的研究,特别是具体技术和算法的实现,而忽略对数据仓库、数据挖掘理论与实践相结合的应用研究。因此,笔者在结合科研项目的基础上,对数据仓库、数据挖掘技术进行较系统的研究,并将研究成果应用于财经、证券投资等领域;同时笔者在研究过程中不断地学习,既对原有的理论和实践进行总结,又不断地将所学到的知识运用到实践中去,丰富了原有的理论。

本书详细阐述数据仓库与数据挖掘的基本原理,系统而全面地介绍数据仓库与数据挖掘的概念、作用、算法以及应用领域、相关学科和发展趋势,并着重讨论数据仓库和数据挖掘在企业管理中的应用及构建策略。基于 SQL Server 2005 介绍数据仓库与数据挖掘工具的操作和应用,并结合具体实例,阐述企业数据仓库和数据挖掘的实施过程。最后,以证券行业为对象提供一个数据挖掘的开发

实例。本书的指导思想是在系统阐述基本知识和基本理论的基础上,强调实际应用能力的培养,充分体现数据仓库和数据挖掘技术作为解决实际问题的方法和工具的特点。本书兼顾理论性与通俗性,注重理论联系实际,叙述时力求深入浅出,着重阐述理论的基本思路与方法的基本步骤。

本书的目的旨在向读者系统阐述数据仓库与数据挖掘的基本原理、方法和实用工具,介绍国内外的最新研究成果。全书共有 11 章,第 1 章介绍数据仓库的基本概念和知识;第 2 章介绍联机分析处理的基本理论;第 3 章介绍数据仓库的设计思想、方法和技巧;第 4 章介绍数据仓库的规划与开发;第 5 章介绍各种数据仓库工具的基本功能及 SQL Server 2005 数据仓库工具的应用;第 6 章介绍数据挖掘的概念和相关知识;第 7 章介绍数据挖掘的算法;第 8 章介绍文本挖掘、Web 挖掘等数据挖掘新技术;第 9 章介绍数据挖掘的工具及其应用;第 10 章介绍数据仓库与数据挖掘的综合应用;第 11 章介绍基于数据挖掘的上市公司财务危机预警应用实例,使读者能结合具体应用进行上机操作,消化和理解所学的知识。

本书由李志刚负责全书的整体策划和最后统稿。编写任务的分工如下:第 1 章、第 2 章、第 6 章、第 7 章由马刚编写;第 3 章由李志刚、宛石峰编写;第 4 章由李志刚、黄艳编写;第 5 章由马刚、李志刚编写;第 8 章、第 10 章由李志刚编写;第 9 章由李志刚、彭易成编写;第 11 章由彭易成、李志刚编写。郭丰恺、聂运洁参加部分的文字和图形处理工作。

在本书的编写过程中,笔者借鉴国内外一些文献和网上资料,由于编写体例的限制未将其在文中一一注明,只在参考文献中列出,在此谨向各位学者表示由衷的敬意和感谢。由于数据仓库和数据挖掘技术发展迅速,尽管笔者付出艰苦的努力,但由于本人水平所限,疏漏甚至错误之处在所难免,恳请专家与读者批评指正。

李志刚
2007 年 12 月

目 录

.....

第1章 数据仓库概述	1
本章主要内容	1
1.1 从数据库到数据仓库	1
1.1.1 决策支持技术与数据库技术的发展	1
1.1.2 数据仓库与数据库的区别	6
1.2 数据仓库的概念与特点	11
1.2.1 数据仓库概念	11
1.2.2 面向主题	11
1.2.3 数据的集成性	12
1.2.4 数据的非易失性	13
1.2.5 数据因时而变的特点	14
1.3 数据仓库中的关键概念	14
1.3.1 外部数据源	14
1.3.2 数据抽取	15
1.3.3 数据清洗	15
1.3.4 数据转换	16
1.3.5 数据加载	16
1.3.6 元数据	16
1.3.7 数据集市	17
1.3.8 数据粒度	17
1.4 数据仓库的数据组织	18
1.4.1 数据仓库的数据组织结构	18
1.4.2 数据粒度与数据分割	19
1.4.3 数据仓库的数据组织形式	20
1.4.4 数据仓库的数据追加与清理	23
1.5 数据仓库与数据集市的关系	24
1.5.1 数据集市的类型	24
1.5.2 数据集市与数据仓库的区别	26
1.5.3 数据集市的特点	27
1.6 数据仓库体系结构	27
1.6.1 数据仓库系统的层次结构	27
1.6.2 数据仓库的构造模式	30

1.7 操作数据存储 ODS	33
1.7.1 操作数据存储 ODS 的概念	33
1.7.2 操作数据存储 ODS 的应用	33
1.7.3 DB-ODS-DW 三层体系结构	36
1.7.4 ODS/DW、ODS/DB 之比较	38
习题一	39
 第 2 章 联机分析处理	40
本章主要内容	40
2.1 联机分析处理的概念	40
2.1.1 OLAP 的定义	40
2.1.2 OLAP 的相关基本概念	41
2.1.3 OLAP 与 OLTP 的关系及比较	42
2.1.4 OLAP 准则	44
2.2 OLAP 多维数据分析	49
2.2.1 OLAP 基本分析动作	49
2.2.2 广义 OLAP 功能	53
2.2.3 多维数据分析实例	55
2.3 OLAP 数据组织	57
2.3.1 多维数据组织	57
2.3.2 关系数据组织	60
2.3.3 两种数据组织的比较	63
2.3.4 HOLAP	66
2.4 OLAP 的体系结构与展现方式	67
2.4.1 OLAP 体系结构	67
2.4.2 OLAP 前端展现方式	69
2.4.3 OLAP 结果的展现方法	71
2.5 OLAP 工具及评价	73
2.5.1 Oracle OLAP 工具	73
2.5.2 OLAP 服务器和工具的评价指标	74
2.5.3 OLAP 的局限性	77
习题二	78
 第 3 章 数据仓库设计	79
本章主要内容	79
3.1 数据仓库中数据模型概述	79
3.1.1 数据模型的概念	79

3.1.2 数据仓库模型的构建原则	81
3.1.3 企业数据模型	82
3.2 概念模型设计	83
3.2.1 企业模型的建立	83
3.2.2 数据模型的规范化	86
3.2.3 常见的概念模型	89
3.3 逻辑模型设计	90
3.3.1 概念模型到逻辑模型的转换	92
3.3.2 数据表的规范化与分割	95
3.3.3 维度表的设计	96
3.3.4 事实表的设计	96
3.3.5 数据集市的设计	97
3.4 物理模型设计	97
3.4.1 定义数据存储结构	97
3.4.2 索引策略	99
3.4.3 存储分配优化	100
3.4.4 数据加载设计	101
3.4.5 物理模型的设计对数据仓库性能的影响	101
3.5 元数据模型	102
3.5.1 元数据的类型	102
3.5.2 元数据的作用	103
3.5.3 元数据的收集与维护	104
3.5.4 元数据的使用	107
3.6 粒度模型	107
3.6.1 粒度的划分	108
3.6.2 粒度级别的确定	109
习题三	110
第4章 数据仓库的规划与开发	111
本章主要内容	111
4.1 数据仓库的投资分析	111
4.1.1 建立数据仓库的必要性	111
4.1.2 数据仓库的投资回报分析与风险分析	113
4.2 数据仓库的开发方法	115
4.2.1 瀑布式开发	115
4.2.2 螺旋式开发	116
4.3 数据仓库的建立过程	118

4.3.1 数据进入数据仓库的过程与建立数据仓库的步骤	118
4.3.2 需求分析	119
4.3.3 数据路线	120
4.3.4 技术路线	121
4.3.5 应用路线	121
4.3.6 数据仓库部署	126
4.3.7 运行维护	126
4.4 数据仓库的维护	127
4.4.1 数据周期	127
4.4.2 参照完整性	127
4.4.3 数据环境信息	128
4.4.4 数据备份与恢复	129
4.5 提高数据仓库性能	130
4.5.1 提高 I/O 性能	130
4.5.2 缩小查询范围	131
4.5.3 采取并行优化技术	131
4.5.4 选择适当的初始化参数	132
4.6 数据仓库的安全性	132
4.6.1 安全类型	132
4.6.2 安全方法	133
4.7 分布式数据仓库	134
4.7.1 分布式数据仓库的优点	134
4.7.2 分布式数据仓库的模型建立与数据划分	135
4.7.3 分布式数据仓库的建设策略	137
4.7.4 分布式数据仓库的技术缺陷	140
习题四	140
第 5 章 数据仓库的工具	141
本章主要内容	141
5.1 数据仓库工具的选择	141
5.1.1 数据仓库工具的组成	141
5.1.2 数据仓库工具应具备的主要功能	142
5.1.3 数据仓库的发展趋势	143
5.1.4 选取数据仓库工具的方法	144
5.2 常用数据仓库产品简介	145
5.2.1 Oracle 9i	145
5.2.2 NCR TeraData	148

5.2.3 IBM DB2	149
5.2.4 SAS	149
5.2.5 Microsoft SQL Server 2005 的数据仓库工具	150
5.3 SQL Server 数据仓库的操作应用	150
5.3.1 SQL Server 数据仓库的框架	150
5.3.2 连接数据源	150
5.3.3 创建数据源视图	153
5.3.4 创建和部署多维数据集	156
习题五	170
第6章 数据挖掘概述	171
本章主要内容	171
6.1 数据挖掘的定义和对象	171
6.1.1 数据挖掘的由来	171
6.1.2 数据挖掘的技术定义	172
6.1.3 数据挖掘的商业定义	174
6.1.4 数据挖掘的对象	175
6.2 数据挖掘的分类	176
6.2.1 概述	176
6.2.2 描述型数据挖掘	176
6.2.3 预测型数据挖掘	177
6.3 数据挖掘系统	177
6.3.1 数据挖掘系统的结构	177
6.3.2 数据挖掘系统的设计	178
6.3.3 数据挖掘系统的发展	179
6.4 数据预处理	180
6.4.1 概述	180
6.4.2 数据清洗	181
6.4.3 数据集成	182
6.4.4 数据转换	182
6.4.5 数据归约	183
6.4.6 属性概念分层的自动生成	185
第7章 数据挖掘的算法	187
本章主要内容	187
7.1 分类规则挖掘	187
7.1.1 分类过程与方法	187

7.1.2 决策树分类	188
7.1.3 贝叶斯分类	192
7.2 预测分析与趋势分析规则	194
7.2.1 预测的基本概念	194
7.2.2 预测的基本方法	194
7.2.3 趋势分析挖掘	195
7.3 数据挖掘的关联算法	196
7.3.1 关联规则的概念及分类	196
7.3.2 简单形式的关联规则算法	197
7.3.3 多层与多维关联规则的挖掘	199
7.3.4 关联分析的其他算法	200
7.4 聚类分析	201
7.4.1 聚类分析的概念	201
7.4.2 聚类分析中的数据类型	202
7.4.3 主要的聚类分析方法	205
7.4.4 聚类分析算法	207
7.5 神经网络算法	209
7.5.1 神经网络的概念	209
7.5.2 定义神经网络拓扑结构	211
7.5.3 基于神经网络的算法	213
第8章 数据挖掘的新技术	215
本章主要内容	215
8.1 文本挖掘技术	215
8.1.1 文本挖掘概述	215
8.1.2 文本挖掘的过程	216
8.1.3 文本挖掘的方法	217
8.1.4 文本挖掘与信息检索	218
8.1.5 文本挖掘的应用	219
8.2 Web 数据挖掘技术	221
8.2.1 Web 挖掘概述	221
8.2.2 Web 的特点	222
8.2.3 Web 挖掘分类	222
8.2.4 Web 挖掘流程	226
8.2.5 Web 数据挖掘的常用工具	228
8.2.6 Web 挖掘的应用	228
8.3 可视化数据挖掘技术	230

8.3.1 数据可视化技术	230
8.3.2 可视化数据挖掘技术的应用	233
8.4 基于 GIS 的空间数据挖掘技术	235
8.4.1 地理信息系统	235
8.4.2 空间数据挖掘	236
8.5 分布式数据挖掘	237
8.5.1 概述	237
8.5.2 分布式数据挖掘系统	239
8.5.3 分布式数据挖掘研究的现状	240
习题八	240
 第 9 章 数据挖掘的工具及其应用	241
本章主要内容	241
9.1 国内外数据挖掘工具及评价	241
9.1.1 数据挖掘软件的特征	241
9.1.2 国外数据挖掘工具	242
9.1.3 国内数据挖掘工具	244
9.1.4 数据挖掘工具的功能分类	246
9.1.5 评价数据挖掘工具优劣的指标	247
9.2 SQL Server 2005 数据挖掘工具应用	249
9.2.1 安装环境要求	249
9.2.2 Analysis Services 功能介绍	250
9.2.3 SQL Server 2005 数据挖掘算法概述	251
9.2.4 SQL Server 2005 数据挖掘算法的选择	252
9.2.5 数据源的准备与创建	258
9.2.6 创建数据挖掘模型	261
9.2.7 处理数据挖掘模型	268
9.2.8 浏览模型	268
9.2.9 测试挖掘模型准确性	272
9.2.10 创建预测查询	275
习题九	277
 第 10 章 数据仓库与数据挖掘的综合应用	278
本章主要内容	278
10.1 数据仓库与数据挖掘的关系	278
10.1.1 数据仓库的观点	278
10.1.2 数据挖掘的观点	281

10.2 数据仓库在企业管理中的应用	281
10.2.1 企业应用数据仓库的意义	282
10.2.2 应用数据仓库弥补 ERP 的不足	284
10.2.3 数据仓库实现分析型 CRM	284
10.2.4 数据仓库提高 SCM 的效率	286
10.3 数据挖掘的社会影响与应用领域	287
10.3.1 数据挖掘的社会影响	287
10.3.2 数据挖掘的应用领域	288
10.3.3 数据挖掘的研究方向	290
10.4 金融业中的数据挖掘应用	291
10.4.1 数据挖掘在银行领域的应用	291
10.4.2 数据挖掘在证券领域的应用	293
10.4.3 数据挖掘在保险领域的应用	296
10.5 数据挖掘与客户关系管理	297
10.5.1 数据挖掘在客户识别和客户保留中的应用	297
10.5.2 客户分类	298
10.5.3 减少信用风险	299
10.5.4 数据挖掘在客户忠诚度分析中的应用	299
10.5.5 个性化营销与销售推荐	300
10.5.6 数据挖掘在客户赢利率分析中的应用	301
10.6 电信业中数据仓库与数据挖掘的应用	302
10.6.1 数据仓库与数据挖掘在电信业中的作用	302
10.6.2 数据挖掘在电信业应用的特点	304
习题十	305
 第 11 章 基于数据挖掘的上市公司财务危机预警应用实例	306
本章主要内容	306
11.1 基本思路和相关知识	306
11.1.1 研究的基本思路	306
11.1.2 财务危机的界定	307
11.1.3 经济预警理论	308
11.2 上市公司财务危机预警模型的指标体系的构建	309
11.2.1 上市公司财务预警指标体系的选取原则	309
11.2.2 上市公司财务预警指标体系的构建	310
11.3 基于数据挖掘建立上市公司财务危机预警模型	313
11.3.1 数据准备	313
11.3.2 建立模型	314

11.3.3 模型测试	326
11.4 建立财务预警系统——财务预警的自动化	335
11.4.1 财务危机预警过程的自动化	336
11.4.2 预测准确性自动化	337
参考文献	339

第 1 章

数据仓库概述

本章主要内容

本章主要介绍从数据库到数据仓库的发展过程;数据仓库的概念与特点;数据仓库的关键名词;数据仓库的数据组织;数据集市的概念与特点,数据仓库与数据集市之间的区别;数据仓库的体系结构以及操作数据存储系统。

随着计算机应用范围的扩大和网络计算的不断发展,企业和各种类型的组织机构提出新的要求,希望计算机能够越来越多地参与数据分析与决策制定的工作。传统的数据库技术主要面对单一的数据库资源,更适合事务型处理等操作,其分析型处理的能力较弱。数据仓库在这种情形之下应运而生,将事务操作型环境和分析型环境进行分离,划清了数据处理的分析型环境与事务操作型环境之间的界限,从而由原来的以数据库为中心的单一数据环境发展为以数据仓库为中心的一种新型体系化环境。数据仓库技术实现了“数据→信息→知识”的过程,为企业各级管理层提供了有力的决策支持。

1.1 从数据库到数据仓库

1.1.1 决策支持技术与数据库技术的发展

1. 决策支持技术的发展

客观世界是由物质、能量和信息这三大基本要素所组成的,人类社会生活的每分每秒都无法离开信息。从远古时代开始,人类就一直在同信息打交道:原始人通过结绳记事的方式进行信息的存储;古埃及的象形文字主要由当时的账房先生表示庶民欠法老谷子的多少;现今更是一个信息爆炸的时代,人类已经彻底“淹没”在数据的“海洋”之中了。正如 John Naisbitt 在《大趋势》一书中所说:

“我们正在被信息所淹没,但我们却由于缺乏知识而感到饥饿。”

任何技术的演进都是伴随着某种发展过程的,决策支持系统(decision support system,DSS)技术也经历了一个漫长而复杂的演化进程,而且它仍然在继续演化。

图 1.1 表明从 20 世纪 60 年代初期到 1980 年 DSS 技术的演化过程。

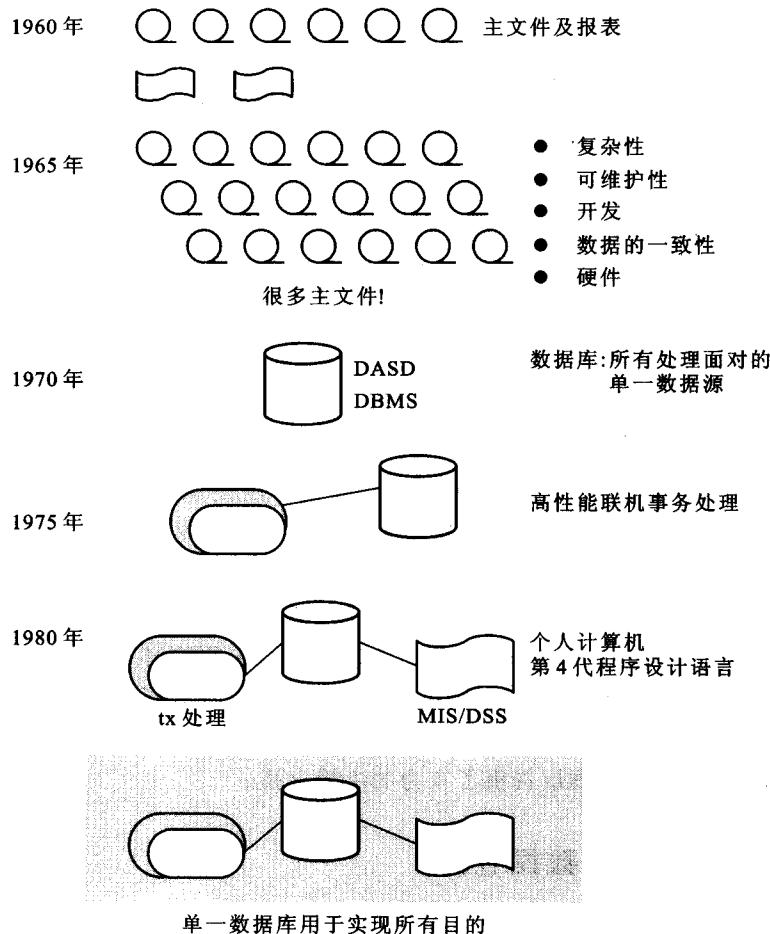


图 1.1 体系化环境的早期演化过程

从时间的角度来看,决策支持技术的发展可以分为以下几个阶段。

(1) 20 世纪 50 年代 ~20 世纪 60 年代: 数据处理阶段

数据处理(data processing)是计算机应用所涉及的最广泛的领域,约占 70%。一个国家的现代化水平越高,数据处理的范围越广,数据量越大,数据处理所占的比例就越高。

(2) 20世纪60年代~20世纪70年代:管理信息系统阶段

随着20世纪五六十年代数据处理应用领域的不断扩大,20世纪六七十年代西方发达国家兴起了管理信息系统(management information system,MIS)的热潮,而我国是在20世纪70年代末到20世纪80年代初,管理信息系统的研究和应用才逐步发展起来。

(3) 20世纪70年代~20世纪80年代:决策支持系统阶段

管理信息系统是在管理科学利用计算机的基础上发展起来的,它促使计算机的应用由数值计算领域拓宽到数据处理(非数值计算)领域,使计算机走向社会活动的各个方面。运筹学和系统工程学利用计算机发展形成了模型辅助决策系统,由于其所采用的模型主要是数学模型,其辅助决策能力主要体现在定量分析上。

20世纪70年代初发展起来的决策支持系统把管理信息系统和模型辅助决策系统有机地结合起来,使得数值计算和数据处理融为一体,提高了辅助决策的能力。

20世纪70年代中期,联机事务处理开始逐步取代数据库,更快速地访问数据成为可能。采用高性能联机事务处理方式,计算机可以完成先前无法完成的工作。

到20世纪80年代,一些更新颖的技术开始涌现,比如个人计算机(personal computer,PC)和第4代程序设计语言(fourth generation language,4GL)。终端用户开始扮演前所未有的角色——直接控制数据和系统,这超出了传统数据处理人员的工作范畴。随着PC与4GL技术的发展,一种新思想诞生了,即除了高性能联机事务处理外,可以对数据进行更多的处理,用来制定管理决策的处理过程可以实现了。在此之前,数据和技术不能一并用来导出详细的操作型决策。一种新的思想体系开始呈现,即单一数据库既能用于操作型的高性能联机事务处理,同时又可用于DSS分析与处理。

(4) 20世纪90年代:智能决策支持系统综合决策支持系统阶段

智能决策支持系统(intelligent decision support system, IDSS)与综合决策支持系统(synthetic decision support system, SDSS)阶段的主要特征是模型技术、专家系统、数据仓库和数据挖掘技术的全方位集成,使得决策支持技术无论是在体系结构还是在信息处理能力方面都产生了较大的变化。

但是, IDSS不能依据数据库中的大量数据进行学习和推理,知识并非由真正的学习而来。

而数据仓库(data warehouse,DW)、联机分析处理(on-line analytical processing,OLAP)、数据挖掘(data mining,DM)相结合构成的新型DSS具有较强的学习能力(从数据中获取新的有用的信息、知识与规则),传统DSS与新型DSS结合(DW+OLAP+DM+MB+DB+ES)形成更高级的SDSS。