

## << Clementine 数据挖掘方法及应用 >>

### 图书基本信息

书名 : << Clementine 数据挖掘方法及应用 >>

13位ISBN编号 : 9787121117787

10位ISBN编号 : 7121117789

出版时间 : 2010-9

出版时间 : 电子工业

作者 : 薛薇//陈欢歌

页数 : 303

版权说明 : 本站所提供下载的PDF图书仅提供预览和简介 , 请支持正版图书。

更多资源请访问 : <http://www.tushu007.com>

## << Clementine 数据挖掘方法及应用 >>

### 前言

数据挖掘是当前数据分析领域中最活跃最前沿的地带。

任何事物都有定性和定量两个方面，定量则产生数据。

从数据分析入手是我们认识事物本质的基本手段。

任何事物都是互相关联着的，从数据分析入手是我们把握事物之间联系的基本方法。

任何事物都在永恒地变化发展着，从数据分析入手是我们探索事物发展规律的基本思路。

所以我们进行数据分析，既是一种世界观，也是一种方法论。

我们在研究着丰富多彩的客观世界的同时，也体现着分析者主观的智慧和自身的价值。

随着中国社会经济的蓬勃发展，在错综复杂的宏观、中观和微观的共同作用下，战略决策和战术选择都显得敏感而关键，越来越多的人们加入到数据分析的行列中来。

这是一个非常富有挑战性的工作，不但有意思而且有意义。

IBM公司于2009年1月公布了其“智慧地球”战略。

该战略的主要思想是，将传感设备或智能仪表嵌入到建筑、电力、交通、管道等各种物体中，进行数据自动采集，之后基于互联网形成物物相联的物联网，然后通过超级计算机和云计算将数据整合，进行智能化分析和建模，从而实现社会与物理世界的融合。

这是一个未来理想化的信息世界图景。

在这个智慧系统中，其核心是数据处理。

为此，IBM公司于2009年7月斥资12亿美元收购了著名的SPSS统计分析软件公司，将其应用广泛的SPSS统计分析软件和Clementine数据挖掘软件纳入麾下。

同时对软件产品进行了整合，将Clementine更新命名为PASW ( Predictive Analytics Software ) Modeler，并快速推向市场。

目前，SPSS Clementine软件已经连续若干年蝉联数据挖掘应用的王者，而业界对于PASW Modeler的认知则刚刚开始。

所以本书继续沿用为广大读者所熟悉的Clementine这个名字。

Clementine软件不但将计算机科学中许多机器学习的优秀算法带入到数据分析中来，同时也综合了一些行之有效数据挖掘方法，成为内容最为全面、功能最为强大的数据挖掘产品。

Clementine软件充分利用计算机系统的运算处理能力和图形展现能力，将方法、应用与工具有机地融合为一体，是解决数据挖掘问题的最理想工具。

Clementine软件继续保持了SPSS产品的一贯风格：界面友好且容易使用。

复杂的数学算法和冗余的输出结果被软件隐藏在程序系统内部。

## 内容概要

数据挖掘是当前数据分析领域中最活跃最前沿的地带。

本书以数据挖掘的实践过程为主线，通过生动的应用案例，从数据挖掘实施角度，系统介绍了经典的数据挖掘方法和利用Clementine实现数据挖掘的全部过程，讲解方法从易到难，说明问题从浅至深。本书力求以最通俗的方式阐述数据挖掘方法的核心思想与基本原理，同时配合Clementine软件操作的说明，希望读者能够直观了解方法本质，尽快掌握Clementine软件使用，并应用到数据挖掘实践中。为方便读者学习，书中所有数据和案例与所附光盘内容一致。

本书适合于从事数据分析各应用领域的读者，尤其适合于商业管理、财政经济、金融保险、社会研究、人文教育等行业的相关人员。

同时，也能够作为高等院校计算机类、财经类、管理类专业本科生和研究生的数据挖掘教材。

## &lt;&lt; Clementine数据挖掘方法及应用 &gt;&gt;

## 书籍目录

第1章 数据挖掘和Clementine概述 1.1 数据挖掘的产生背景 1.1.1 海量数据的分析需求催生数据挖掘  
 1.1.2 应用对理论的挑战催生数据挖掘 1.2 什么是数据挖掘 1.2.1 数据挖掘的概念 1.2.2 数据挖掘能做什么 1.2.3 数据挖掘得到的知识形式 1.2.4 数据挖掘的算法分类 1.3 Clementine软件概述 1.3.1 Clementine的窗口 1.3.2 数据流的基本管理和执行 1.3.3 数据流的其他管理 1.3.4 从一个示例看Clementine的使用

第2章 Clementine数据的读入 2.1 变量的类型 2.1.1 从数据挖掘角度看变量类型 2.1.2 从数据存储角度看变量类型 2.2 读入数据 2.2.1 读自由格式的文本文件 2.2.2 读Excel电子表格数据 2.2.3 读SPSS格式文件 2.2.4 读数据库文件 2.3 生成实验方案数据 2.4 合并数据 2.4.1 数据的纵向合并 2.4.2 数据的横向合并

第3章 Clementine变量的管理 3.1 变量说明 3.1.1 取值范围和缺失值的说明 3.1.2 变量取值有效性检查和修正 3.1.3 变量角色的说明 3.2 变量值的重新计算 3.2.1 CLEM表达式 3.2.2 变量值重新计算示例 3.3 变量类别值的调整 3.4 生成新变量 3.5 变量值的离散化处理 3.5.1 常用的分箱方法 3.5.2 变量值的离散化处理示例 3.6 生成样本集分割变量 3.6.1 样本集分割的意义和常见方法 3.6.2 生成样本集分割变量的示例

第4章 Clementine样本的管理 4.1 样本的排序 4.2 样本的条件筛选 4.3 样本的随机抽样 4.4 样本的浓缩处理 4.5 样本的分类汇总 4.6 样本的平衡处理 4.7 样本的其他管理 4.7.1 数据转置 4.7.2 数据的重新组织

第5章 Clementine数据的基本分析 5.1 数据质量的探索 5.1.1 数据的基本描述与质量探索 5.1.2 离群点和极端值的修正 5.1.3 缺失值的替补 5.1.4 数据质量管理的其他功能 5.2 基本描述分析 5.2.1 计算基本描述统计量 5.2.2 绘制散点图 5.3 变量分布的探索 5.4 两分类变量相关性的研究 5.4.1 两分类变量相关性的图形分析 5.4.2 两分类变量相关性的数值分析 5.5 两总体的均值比较 5.5.1 两总体均值比较的图形分析 5.5.2 独立样本的均值检验 5.5.3 配对样本的均值检验 5.6 变量重要性的分析 5.6.1 变量重要性分析的一般方法 5.6.2 变量重要性分析的应用示例

第6章 分类预测：Clementine的决策树 6.1 决策树算法概述 6.1.1 什么是决策树 6.1.2 决策树的几何理解 6.1.3 决策树的核心问题 6.2 Clementine的C5.0算法及应用 6.2.1 信息熵和信息增益 6.2.2 C5.0的决策树生长算法 6.2.3 C5.0的剪枝算法 6.2.4 C5.0的推理规则集 6.2.5 C5.0的基本应用示例 6.2.6 C5.0的损失矩阵和Boosting技术 6.2.7 C5.0的模型评价 6.2.8 C5.0的其他话题：推理规则、交叉验证和未剪枝的决策树 6.3 Clementine的分类回归树及应用 6.3.1 分类回归树的生长过程 6.3.2 分类回归树的剪枝过程 6.3.3 损失矩阵对分类树的影响 6.3.4 分类回归树的基本应用示例 6.3.5 分类回归树的交互建模 6.3.6 分类回归树的模型评价 6.4 Clementine的CHAID算法及应用 6.4.1 CHAID分组变量的预处理和选择策略 6.4.2 Exhaustive CHAID算法 6.4.3 CHAID的剪枝 6.4.4 CHAID的应用示例 6.5 Clementine的QUEST算法及应用 6.5.1 QUEST算法确定最佳分组变量和分割点的方法 6.5.2 QUEST算法的应用示例 6.6 决策树算法评估的图形比较 6.6.1 不同模型的误差对比 6.6.2 不同模型收益的对比

第7章 分类预测：Clementine的人工神经网络 7.1 人工神经网络算法概述 7.1.1 人工神经网络的概念和种类 7.1.2 人工神经网络中的节点和意义 7.1.3 人工神经网络建立的一般步骤 7.2 Clementine的B-P反向传播网络 7.2.1 感知机模型 7.2.2 B-P反向传播网络的特点 7.2.3 B-P反向传播算法 7.2.4 B-P反向传播网络的其他问题 7.3 Clementine的B-P反向传播网络的应用 7.3.1 基本操作说明 7.3.2 计算结果说明 7.3.3 提高模型预测精度 7.4 Clementine的径向基函数网络及应用 7.4.1 径向基函数网络中的隐节点和输出节点 7.4.2 径向基函数网络的学习过程 7.4.3 径向基函数网络的应用示例

第8章 分类预测：Clementine的统计方法 8.1 Clementine的Logistic回归分析及应用 8.1.1 二项Logistic回归方程 8.1.2 二项Logistic回归方程系数的含义 8.1.3 二项Logistic回归方程的检验 8.1.4 二项Logistic回归分析的应用示例 8.1.5 多项Logistic回归分析的应用示例 8.2 Clementine的判别分析及应用 8.2.1 距离判别法 8.2.2 Fisher判别法 8.2.3 贝叶斯判别法 8.2.4 判别分析的应用示例

第9章 探索内部结构：Clementine的关联分析 9.1 简单关联规则及其有效性 9.1.1 简单关联规则的基本概念 9.1.2 简单关联规则的有效性和实用性 9.2 Clementine的Apriori算法及应用 9.2.1 产生频繁项集 9.2.2 依据频繁项集产生简单关联规则 9.2.3 Apriori算法的应用示例 9.3 Clementine的GRI算法及应用 9.3.1 GRI算法基本思路 9.3.2 GRI算法的具体策略 9.3.3 GRI算法的应用示例 9.4 Clementine的序列关联及应用 9.4.1 序列关联中的基本概念 9.4.2 Sequence算法 9.4.3 序列关联的时间约束 9.4.4 序列关联分析的应用示例

第10章 探索内部结构：Clementine的聚类分析 10.1 聚类分析的一般问题 10.1.1 聚类分

## << Clementine数据挖掘方法及应用 >>

析的提出 10.1.2 聚类分析的算法 10.2 Clementine的K-Means聚类及应用 10.2.1 K-Means对“亲疏程度”的测度 10.2.2 K-Means聚类过程 10.2.3 K-Means聚类的应用示例 10.3 Clementine的两步聚类及应用 10.3.1 两步聚类对“亲疏程度”的测度 10.3.2 两步聚类过程 10.3.3 聚类数目的确定 10.3.4 两步聚类的应用示例 10.4 Clementine的Kohonen网络聚类及应用 10.4.1 Kohonen网络的聚类机理 10.4.2 Kohonen网络的聚类过程 10.4.3 Kohonen网络聚类的示例 10.5 基于聚类分析的离群点探索及应用 10.5.1 多维空间基于聚类的诊断方法 10.5.2 多维空间基于聚类的诊断方法应用示例 参考文献

## << Clementine 数据挖掘方法及应用 >>

### 章节摘录

插图：数据挖掘，作为20世纪90年代中后期兴起的，具有鲜明跨学科色彩的应用和研究领域，因其注重减少数据分析方法对数据的限制性和约束性，注重与计算机技术结合以实现数据的可管理性以及分析的易操作性，已成为数据分析应用实践的新生代。

同时，随着数据挖掘方法的不断成熟及其应用的日益普及化，数据挖掘软件的研发也取得了令人可喜的成果。

目前，以Clementine为代表的数据挖掘软件，因其有效地将束之高阁的数据挖掘理论成果解放到数据分析实践中，已普遍应用于商业、社会、经济、教育、金融、医学等领域，并成为数据分析的主流工具，得到数据分析相关领域的极大关注。

1.1 数据挖掘的产生背景数据挖掘的产生和兴起是在计算机数据库技术蓬勃发展，人工智能技术应用领域不断拓展，统计分析方法不断丰富过程中，为有效迎合数据分析的实际需求而逐步形成和发展起来的一门具有鲜明跨学科色彩的应用研究领域。

1.1.1 海量数据的分析需求催生数据挖掘20世纪80年代以来，随着计算机数据库技术和产品的日益成熟以及计算机应用的普及深化，各行业部门的数据采集能力得到了前所未有的提高，组织通过各自内部的业务处理系统、管理信息系统以及外部网络系统，获得并积累了浩如烟海的数据。

以商业领域为例，美国著名的连锁超市Wal-Mart的数据库中已积累了TB级以上的顾客购买行为数据和其他销售数据。

随着互联网和电子商务的普及，各类网上书店、网上银行、网上营业厅和网上商城等积累的Web点击流数据，存储容量也多高达GB级。

另外，国家政府部门所积累的数据量也令人瞠目。

例如，一次全国经济普查或人口普查所采集和处理数据量均在千万级以上。

同时，各经济行业的企业内部也拥有大量的业务数据、财务数据和人事数据。

在严酷的市场竞争压力下，企业为更客观地把握自身和市场状况，提升内部管理和决策水平，管理者们面对如此丰富的海量数据，分析需求越来越强烈。

编辑推荐

《Clementine数据挖掘方法及应用》由电子工业出版社出版。

版权说明

本站所提供下载的PDF图书仅提供预览和简介，请支持正版图书。

更多资源请访问:<http://www.tushu007.com>