

<<数据挖掘导论>>

图书基本信息

书名：<<数据挖掘导论>>

13位ISBN编号：9787111316701

10位ISBN编号：7111316703

出版时间：2010.9

出版时间：机械工业出版社

作者：(美)Pang-Ning Tan,Michael Steinbach,Vipin Kumar

页数：769

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>

前言

Advances in data generation and collection are producing data sets of massive size in commerce and a variety of scientific disciplines . Data warehouses store details of the sales and operations of businesses , Earth-orbiting satellites beam high-resolution images and sensor data back to Earth . and genomics experiments generate sequence , structural , and functional data for an increasing number of organisms . The ease with which data can now be gathered and stored has created a new attitude toward data analysis : Gather whatever data you can whenever and wherever possible . It has become an article of faith that the gathered data will have value . either for the purpose that initially motivated its collection or for purposes not yet envisioned . The field of data mining grew out of the limitations of current data analysis techniques in handling the challenges posed by these new types of data sets . Data mining does not replace other areas of data analysis , but rather takes them as the foundation for much of its work . While some areas of data mining , such as association analysis , are unique to the field , other areas , such as clustering , classification, and anomaly detection , build upon a long history of work on these topics in other fields . Indeed . the willingness of data mining researchers to draw upon existing techniques has contributed to the strength and breadth of the field , as well as to its rapid growth .

<<数据挖掘导论>>

内容概要

本书全面介绍了数据挖掘的理论和方法，着重介绍如何用数据挖掘知识解决各种实际问题，涉及学科领域众多，适用面广。

书中涵盖5个主题：数据、分类、关联分析、聚类和异常检测。

除异常检测外，每个主题都包含两章：前面一章讲述基本概念、代表性算法和评估技术，后面一章较深入地讨论高级概念和算法。

目的是使读者在透彻地理解数据挖掘基础的同时，还能了解更多重要的高级主题。

· 包含大量的图表、综合示例和丰富的习题。

· 不需要数据库背景。

只需要很少的统计学或数学背景知识。

· 网上配套教辅资源丰富，包括PPT、习题解答、数据集等。

作者简介

作者：（美国）谭（Pang-Ning Tan）（美国）斯坦巴克（Michael Steinbach）（美国）库马尔（Vipin Kumar）Pang.Ning Tan现为密歇根州立大学计算机与工程系助理教授，主要教授数据挖掘、数据库系统等课程。

他的研究主要关注于为广泛的应用（包括医学信息学、地球科学、社会网络、Web挖掘和计算机安全）开发适用的数据挖掘算法。

Michael Steinbach拥有明尼苏达大学数学学士学位、统计学硕士学位和计算机科学博士学位，现为明尼苏达大学双城分校计算机科学与工程系助理研究员。

Vipin Kumar现为明尼苏达大学计算机科学与工程系主任和William Norris教授。

1988年至2005年。

他曾担任美国陆军高性能计算研究中心主任。

书籍目录

Preface1 Introduction 1.1 What Is Data Mining? 1.2 Motivating Challenges 1.3 The Origins of Data Mining
1.4 Data Mining Tasks 1.5 Scope and Organization of the Book 1.6 Bibliographic Notes 1.7 Exercises2
Data 2.1 Types of Data 2.1.1 Attributes and Measurement 2.1.2 Types of Data Sets 2.2 Data Quality
2.2.1 Measurement and Data Collection Issues 2.2.2 Issues Related to Applications 2.3 Data Preprocessing
2.3.1 Aggregation 2.3.2 Sampling 2.3.3 Dimensionality Reduction 2.3.4 Feature Subset Selection
2.3.5 Feature Creation 2.3.6 Discretization and Binarization 2.3.7 Variable Transformation 2.4
Measures of Similarity and Dissimilarity 2.4.1 Basics 2.4.2 Similarity and Dissimilarity between Simple
Attributes. 2.4.3 Dissimilarities between Data Objects 2.4.4 Similarities between Data Objects 2.4.5
Examples of Proximity Measures 2.4.6 Issues in Proximity Calculation 2.4.7 Selecting the Right Proximity
Measure 2.5 Bibliographic Notes 2.6 Exercises3 Exploring Data 3.1 The Iris Data Set 3.2 Summary Statistics
3.2.1 Frequencies and the Mode 3.2.2 Percentiles 3.2.3 Measures of Location: Mean and Median
3.2.4 Measures of Spread: Range and Variance 3.2.5 Multivariate Summary Statistics 3.2.6 Other Ways
to Summarize the Data 3.3 Visualization 3.3.1 Motivations for Visualization 3.3.2 General Concepts
3.3.3 Techniques 3.3.4 Visualizing Higher-Dimensional Data 3.3.5 Do's and Don'ts 3.4 OLAP and
Multidimensional Data Analysis 3.4.1 Representing Iris Data as a Multidimensional Array 3.4.2
Multidimensional Data: The General Case 3.4.3 Analyzing Multidimensional Data 3.4.4 Final Comments
on Multidimensional Data Analysis 3.5 Bibliographic Notes 3.6 Exercises Classification:4 Basic Concepts,
Decision Trees, and Model Evaluation 4.1 Preliminaries 4.2 General Approach to Solving a Classification
Problem 4.3 Decision Tree Induction 4.3.1 How a Decision Tree Works 4.3.2 How to Build a Decision
Tree 4.3.3 Methods for Expressing Attribute Test Conditions . 4.3.4 Measures for Selecting the Best Split
4.3.5 Algorithm for Decision Tree Induction 4.3.6 An Example: Web Robot Detection 4.3.7
Characteristics of Decision Tree Induction 4.4 Model Overfitting 4.4.1 Overfitting Due to Presence of Noise
4.4.2 Overfitting Due to Lack of Representative Samples . 4.4.3 Overfitting and the Multiple Comparison
Procedure 4.4.4 Estimation of Generalization Errors 4.4.5 Handling Overfitting in Decision Tree
Induction . . 4.5 Evaluating the Performance of a Classifier 4.5.1 Holdout Method 4.5.2 Random
Subsampling 4.5.3 Cross-Validation 4.5.4 Bootstrap 4.6 Methods for Comparing Classifiers 4.6.1
Estimating a Confidence Interval for Accuracy 4.6.2 Comparing the Performance of Two Models 4.6.3
Comparing the Performance of Two Classifiers 4.7 Bibliographic Notes 4.8 Exercises5 Classification:
Alternative Techniques6 Association Analysis: Basic Concepts and Algorithms

章节摘录

插图：What Is an attribute? We start with a more detailed definition of an attribute. Definition 2.1. An attribute is a property or characteristic of an object that may vary, either from one object to another or from one time to another. For example, eye color varies from person to person, while the temperature of an object varies over time. Note that eye color is a symbolic attribute with a small number of possible values (brown, black, blue, green, hazel, etc.), while temperature is a numerical attribute with a potentially unlimited number of values. At the most basic level, attributes are not about numbers or symbols. However, to discuss and more precisely analyze the characteristics of objects, we assign numbers or symbols to them. To do this in a well-defined way, we need a measurement scale. Definition 2.2. A measurement scale is a rule (function) that associates a numerical or symbolic value with an attribute of an object. Formally, the process of measurement is the application of a measurement scale to associate a value with a particular attribute of a specific object. While this may seem a bit abstract, we engage in the process of measurement all the time. For instance, we step on a bathroom scale to determine our weight, we classify someone as male or female, or we count the number of chairs in a room to see if there will be enough to seat all the people coming to a meeting. In all these cases, the "physical value" of an attribute of an object is mapped to a numerical or symbolic value. With this background, we can now discuss the type of an attribute, a concept that is important in determining if a particular data analysis technique is consistent with a specific type of attribute.

<<数据挖掘导论>>

编辑推荐

《数据挖掘导论(英文版)》是经典原版书库。

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>