

<<文本挖掘>>

图书基本信息

书名：<<文本挖掘>>

13位ISBN编号：9787115205353

10位ISBN编号：7115205353

出版时间：2009-8

出版时间：人民邮电出版社

作者：（以）费尔德曼,（美）桑格

页数：410

字数：506000

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

The information age has made it easy to store large amounts of data. The proliferation of documents available on the Web , on corporate intranets , on news wires , and elsewhere is overwhelming. However , although the amount of data available to us is constantly increasing , our ability to absorb and process this information remains constant. Search engines only exacerbate the problem by making more and more documents available in a matter of a few key strokes. Text mining is a new and exciting research area that tries to solve the information overload problem by using techniques from data mining , machine learning , natural language processing ( NLP ) , information retrieval ( IR ) , and knowledge management. Text mining involves the preprocessing of document collections ( text categorization , information extraction , term extraction ) , the storage of the intermediate representations , the techniques to analyze these intermediate representations ( such as distribution analysis , clustering , trend analysis , and association rules ) , and visualization of the results.

This book presents a general theory of text mining along with the main techniques behind it. We offer a generalized architecture for text mining and outline the algorithms and data structures typically used by text mining systems. The book is aimed at the advanced undergraduate students , graduate students , academic researchers , and professional practitioners interested in complete coverage of the text mining field. We have included all the topics critical to people who plan to develop text mining systems or to use them. In particular , we have covered preprocessing techniques such as text categorization , text clustering , and information extraction and analysis techniques such as association rules and link analysis. The book tries to blend together theory and practice; we have attempted to provide many real-life scenarios that show how the different techniques are used in practice. When writing the book we tried to make it as self-contained as possible and have compiled a comprehensive bibliography for each topic so that the reader can expand his or her knowledge accordingly.

## <<文本挖掘>>

### 内容概要

本书是一部文本挖掘领域名著，作者为世界知名的权威学者。书中涵盖了核心文本挖掘操作、文本挖掘预处理技术、分类、聚类、信息提取、信息提取的概率模型、预处理应用、可视化方法、链接分析、文本挖掘应用等内容，很好地结合了文本挖掘的理论和实践。

本书非常适合文本挖掘、信息检索领域的研究人员和实践者阅读，也适合作为高等院校计算机及相关专业研究生的数据挖掘和知识发现等课程的教材。

## <<文本挖掘>>

### 作者简介

Ronen Feldman 机器学习、数据挖掘和非结构化数据管理的先驱人物。以色列Bar-Ilan大学数学与计算机科学系高级讲师、数据挖掘实验室主任，Clearforest公司（主要为企业和政府机构开发下一代文本挖掘应用）合作创始人、董事长，现在还是纽约大学Stern商学院的副教授。

<<文本挖掘>>

书籍目录

. Introduction to Text Mining .1 Defining Text Mining .2 General Architecture of Text Mining Systems . Core Text Mining Operations .1 Core Text Mining Operations .2 Using Background Knowledge for Text Mining .3 Text Mining Query Languages . Text Mining Preprocessing Techniques .1 Task-Oriented Approaches .2 Further Reading . Categorization .1 Applications of Text Categorization .2 Definition of the Problem .3 Document Representation .4 Knowledge Engineering Approach to TC .5 Machine Learning Approach to TC .6 Using Unlabeled Data to Improve Classification .7 Evaluation of Text Classifiers .8 Citations and Notes . Clustering .1 Clustering Tasks in Text Analysis .2 The General Clustering Problem .3 Clustering Algorithms .4 Clustering of Textual Data .5 Citations and Notes . Information Extraction .1 Introduction to Information Extraction .2 Historical Evolution of IE: The Message Understanding Conferences and Tipster .3 IE Examples .4 Architecture of IE Systems .5 Anaphora Resolution .6 Inductive Algorithms for IE .7 Structural IE .8 Further Reading . Probabilistic Models for Information Extraction .1 Hidden Markov Models .2 Stochastic Context-Free Grammars .3 Maximal Entropy Modeling .4 Maximal Entropy Markov Models .5 Conditional Random Fields .6 Further Reading . Preprocessing Applications Using Probabilistic and Hybrid Approaches .1 Applications of HMM to Textual Analysis .2 Using MEMM for Information Extraction .3 Applications of CRFs to Textual Analysis .4 TEG: Using SCFG Rules for Hybrid Statistical – Knowledge-Based IE .5 Bootstrapping .6 Further Reading . Presentation-Layer Considerations for Browsing and Query Refinement .1 Browsing .2 Accessing Constraints and Simple Specification Filters at the Presentation Layer .3 Accessing the Underlying Query Language .4 Citations and Notes . Visualization Approaches .1 Introduction .2 Architectural Considerations .3 Common Visualization Approaches for Text Mining .4 Visualization Techniques in Link Analysis .5 Real-World Example: The Document Explorer System . Link Analysis .1 Preliminaries .2 Automatic Layout of Networks .3 Paths and Cycles in Graphs .4 Centrality .5 Partitioning of Networks .6 Pattern Matching in Networks .7 Software Packages for Link Analysis .8 Citations and Notes . Text Mining Applications .1 General Considerations .2 Corporate Finance: Mining Industry Literature for Business Intelligence .3 A “ Horizontal ” Text Mining Application: Patent Analysis Solution Leveraging a Commercial Text Analytics Platform .4 Life Sciences Research: Mining Biological Pathway Information with GeneWays Appendix A: DIAL: A Dedicated Information Extraction Language for Text Mining A.1 What Is the DIAL Language? A.2 Information Extraction in the DIAL Environment A.3 Text Tokenization A.4 Concept and Rule Structure A.5 Pattern Matching A.6 Pattern Elements A.7 Rule Constraints A.8 Concept Guards A.9 Complete DIAL Examples Bibliography Index

## 章节摘录

Similarity Functions for Simple Concept Association Graphs      Similarity functions often form an essential part of working with simple concept association graphs, allowing a user to view relations between concepts according to differing weighting measures. Association rules involving sets ( or concepts ) A and B that have been described in detail in Chapter II are often introduced into a graph format in an undirected way and specified by a support and a confidence threshold. A fixed confidence threshold is often not very reasonable because it is independent of the support from the RHS of the rule. As a result, an association should have a significantly higher confidence than the share of the RHS in the whole context to be considered as interesting. Significance is measured by a statistical test ( e.g., t-test or chi-square ) . With this addition, the relation given by an association rule is undirected. An association between two sets A and B in the direction AB implies also the association B A. This equivalence can be explained by the fact that the construct of a statistically significant association is different from implication ( which might be suggested by the notation AB ) . It can easily be derived that if B is overproportionally represented in A, then A is also overproportionally represented in B. As an example of differences of similarity functions, one can compare the undirected connection graphs given by statistically significant association rules with the graphs based on the cosine function. The latter relies on the cosine of two vectors and is efficiently applied for continuous, ordinal, and also binary attributes. In case of documents and concept sets, a binary vector is associated to a concept set with the vector elements corresponding to documents. An element holds the value 1 if all the concepts of the set appear in the document. Table X.1 ( Feldman, Kloesgen, and Zilberstein 1997b ) , which offers a quick summary of some common similarity functions, shows that the cosine similarity function in this binary case reduces to the fraction built by the support of the union of the two concept sets and the geometrical mean of the support of the two sets. A connection between two sets of concepts is related to a threshold for the cosine similarity ( e.g., 10% ) . This means that the two concept sets are connected if the support of the document subset that holds all the concepts of both sets is larger than 10 percent of the geometrical mean of the support values of the two concept sets.

## <<文本挖掘>>

### 媒体关注与评论

“ .....我购买了这本书。  
这本书绝对是非常值得拥有的参考书。

” ——L.Venkata Subramaniam, IBM印度研究实验室 “ 一本由该领域最重要专家编写的文本挖掘导论。

这本书写得非常好。

完美地结合了文本挖掘的理论和实践,既适合研究人员又适合实践者.....极力推荐那些没有任何计算语言学背景而想钻研文本挖掘领域的人阅读本书。

” ——Rada Mihalcea, 北得克萨斯大学 文本挖掘已经成为令人兴奋的新兴研究领域。

本书由世界知名的权威学者编写,除了讲解核心文本挖掘和链路检测算法及技术之外,还介绍了高级预处理技术。

并考虑了知识表示方面的因素以及可视化方法。

此外。

书中还探讨了有关技术在实践中的应用,很好地兼顾了文本挖掘的理论和实践

#### 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>