

<<Web知识挖掘>>

图书基本信息

书名：<<Web知识挖掘>>

13位ISBN编号：9787030274991

10位ISBN编号：7030274997

出版时间：2010-6

出版时间：科学出版社

作者：郑庆华

页数：336

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

1989年，欧洲核子研究组织（European Organization for Nuclear Research，CERN）的工程师Tim Berners-Lee针对科学家之间文献交流的需求，首次提出了Web概念与应用架构，其核心是通过超链接实现文本文档的共享。

其后，随着超文本标记语言（HTML）、超文本传输协议（HTTP）等技术标准的逐步成熟，以及Mosaic、Narigatot等浏览器的广泛应用，Web在1995年前后进入了快速发展阶段，表现为Internet上的Web页面数量与服务器数量呈指数级增长。

2004年以后，Internet上的PIW（publicly indexable Web）页面数已达到了10数量级，每天新增页面的数量超过800万，而Web服务器数量的倍增周期仅为23周。

Web已成为一个开放性的、动态的、全球性信息服务中心，以及当前人们获取信息的重要手段。

然而，Web上同样面临着信息社会所共有的“信息爆炸”与“知识贫乏”的矛盾性问题。

如何从这些海量的Web数据中发现有用的知识或者模式，成为人们亟待解决的问题。

传统的数据挖掘技术主要针对结构化的数据对象，还很难适用于具有异构性、半结构化特性以及高度动态性等特点的Web数据。

为此，Etzioni于1996年提出了“Web挖掘”的概念。

Web挖掘是数据挖掘、机器学习、数据库、自然语言处理、Web / Internet等多种信息技术相互渗透与融合的必然结果，旨在研究如何从Web文档与服务中抽取有价值的知识或隐含信息。

近年来，Web挖掘这个研究领域得到了国内外学者越来越多的关注，人们以文本分类、信息抽取、检索结构排序、用户访问模式发现等应用为目标，在Web挖掘的三个子领域——结构挖掘、内容挖掘、日志挖掘方面从事了大量的研究工作，在理论、方法与应用方面取得了一系列研究成果。

## &lt;&lt;Web知识挖掘&gt;&gt;

## 内容概要

本书是一部关于Web知识挖掘的比较系统、完整，且理论和实践相结合的著作，共含7章：第1章与第2章是Web知识挖掘概论，其中，第1章总体上对Web知识挖掘的现状、概念、典型方法、应用领域以及面临的挑战进行综述性说明；第2章介绍了Web知识挖掘的预备知识、分类体系、基本流程等内容。第3~6章是Web知识挖掘的理论与方法，分别论述了Web爬取、Web结构挖掘、内容挖掘、日志挖掘相关理论与方法，并系统总结了我们自己在元数据、概念、知识元等多个层次上的知识获取以及个性化知识服务等方面的工作。

第7章是Web知识挖掘的实践与应用实例，以实例对Web结构挖掘、日志挖掘及内容挖掘的应用进行了说明。

本书不仅系统地介绍了Web知识挖掘领域的基础理论与方法，也阐述了我们在该领域的创新性工作，因而适合不同类型与层次的研究人员及学生。

本书可作为信息领域的科研与工程技术人员的参考书，也可作为计算机与相关专业的研究生和高年级本科生的教材或辅导书目。

## 书籍目录

前言 第1章 Web挖掘概述 1.1 Web发展历史与现状 1.2 Web挖掘的概念 1.3 Web挖掘面临的挑战 1.4 Web挖掘的研究方向 1.5 小结 第2章 Web挖掘的基础知识 2.1 Web挖掘的主要预备知识 2.2 Web挖掘分类 2.3 Web挖掘的主要应用 2.4 Web挖掘的基本流程 2.5 Web挖掘领域的重要文献、国际期刊与会议、标准规范 2.6 小结 第3章 Web爬取与页面组织管理 3.1 Web爬取概述 3.2 Web爬取中的主要技术问题 3.3 隐含Web爬取 3.4 面向主题的Web爬取 3.5 爬取页面的存储与管理 3.6 小结 第4章 Web结构挖掘 4.1 Web结构挖掘概述 4.2 PageRank算法 4.3 HITS算法 4.4 Hilltop算法 4.5 Web宏观结构特性分析 4.6 小结 第5章 Web内容挖掘 5.1 Web页面的特征表示 5.2 Web页面分类 5.3 Web页面聚类 5.4 面向Web的信息抽取 5.5 面向Web的本体学习 5.6 面向Web的知识元及其关联抽取 5.7 多媒体数据挖掘 5.8 Web内容挖掘的未来研究方向 5.9 小结 第6章 Web日志挖掘 6.1 Web日志挖掘概述 6.2 Web日志预处理 6.3 序列模式挖掘 6.4 Web用户行为模式挖掘 6.5 Web用户个性挖掘 6.6 Web用户兴趣感知 6.7 Web日志挖掘的未来研究方向 6.8 小结 第7章 Web挖掘的应用实例 7.1 应用1：面向网络学习的学习者个性挖掘 7.2 应用2：海量Web资源中的知识处理与服务 7.3 小结 参考文献

## 章节摘录

插图：Web挖掘是从数据挖掘发展而来的，但与传统的数据挖掘相比有许多独特之处。数据挖掘，又称为面向数据库的知识发现（knowledge discovery in database, KDD），就是从大量数据中获取新颖的、潜在有用的模式的过程。

数据挖掘的对象是来自关系型数据库或XML数据库中的结构化数据。

而Web挖掘的对象包括网页、图像、声音、视频、网页之间的链接以及网站用户的日志数据。

除了日志数据外，其他类型数据具有海量、异构、非结构化等特性，传统的数据挖掘技术还很难处理这类数据。

因此，必须在Web挖掘领域中，研究专门针对Web数据特点的算法与方法。

在信息检索中，用户以关键词组合表达检索需求，通过关键词匹配的方式从特定文档集中返回与检索需求相关的文档。

信息检索包括文档的建模、分类、索引、结果排序与可视化Web等流程，Web挖掘技术一般用于其中的分类、索引以及结果排序，从这个角度来说，Web挖掘是信息检索过程的重要组成部分（Kosala et al, 2000）。

另一方面，信息检索的结果往往也是Web挖掘的对象，如在HITS算法中，因而信息检索也可作为Web挖掘的组成部分。

信息抽取指从给定的文档中抽取特定类别的信息，例如，从一篇文档中抽取标题、作者等元数据信息。

由于Web站点的异构性，大多数信息抽取都是针对特定网站，一些抽取方法能够自动或半自动地建立抽取模式（Kushmerick, 1999），对于这类信息抽取，Web挖掘可以看做信息抽取的一个过程。

此外，在Web挖掘中，利用信息抽取可以建立文档的压缩版本以提高挖掘效率，从这个角度来说，信息抽取可以作为Web挖掘的预处理过程。

编辑推荐

《Web知识挖掘:理论、方法与应用》由科学出版社出版。

#### 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>