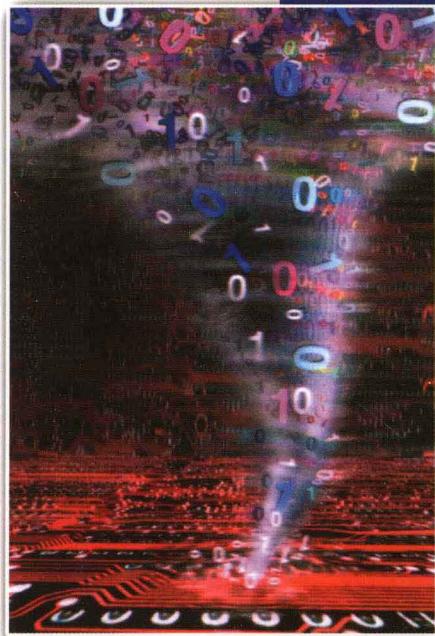


预测性文本挖掘基础

Fundamentals of
Predictive
Text Mining



[美] 绍洛姆·韦斯
[澳] 露廷·因杜尔亚 著

[美] 张潼
赵仲孟 侯迪 译

Sholom M. Weiss Nitin Indurkha Tong Zhang



西安交通大学出版社
XI'AN JIAOTONG UNIVERSITY PRESS

Fundamentals of Predictive Text Mining

预测性文本挖掘基础

〔美〕绍洛姆·韦斯

〔澳〕霓廷·因杜尔亚 著

〔美〕张 潼

Sholom M. Weiss

T. J. Watson Research Center, IBM Corporation

Nitin Indurkha

School of Computer Science & Engg. University of New South Wales

Tong zhang

Dept. Statistics, Hill Center Rutgers University

赵仲孟 侯 迪 译



西安交通大学出版社

Xi'an Jiaotong University Press

Translation from the English language edition:

“Fundamentals of Predictive Text Mining” by S. Weiss & N. Indurkhy & T. Zhang (edition: 2; year of publication: 2010)

Copyright © 2010 Springer Verlag London as a part of Springer Science+Business Media
All Rights Reserved

本书中文简体字版由施普林格科学与商业媒体集团授权西安交通大学出版社独家出版发行。未经出版者预先书面许可,不得以任何方式复制或发行本书的任何部分。

陕西省版权局著作权合同登记号 图字 25-2011-002 号

图书在版编目(CIP)数据

预测性文本挖掘基础/(美)韦斯(Weiss,S.),张潼,
(澳)因杜尔亚(Indurkhy,N.)著;赵仲孟,侯迪译.

—西安:西安交通大学出版社,2012.10

ISBN 978-7-5605-4427-4

I. ①预… II. ①韦… ②张… ③因… ④赵…
⑤侯… III. ①数据采集-研究 IV. ①TP274

中国版本图书馆 CIP 数据核字(2012)第 142928 号

书 名 预测性文本挖掘基础
著 者 (美)绍洛姆·韦斯,(澳)霓廷·因杜尔亚,(美)张潼
译 者 赵仲孟 侯 迪

出版发行 西安交通大学出版社
(西安市兴庆南路 10 号 邮政编码 710049)
网 址 <http://www.xjtupress.com>
电 话 (029)82668357 82667874(发行中心)
(029)82668315 82669096(总编办)
传 真 (029)82668280
印 刷 西安建科印务有限责任公司

开 本 700mm×1000mm 1/16 印张 15.75
印 数 0001~3000 字数 247 千字
版次印次 2012 年 10 月第 1 版 2012 年 10 月第 1 次印刷
书 号 ISBN 978-7-5605-4427-4/TP·569
定 价 43.00 元

读者购书、书店添货、如发现印装质量问题,请与本社发行中心联系、调换。

订购热线:(029)82665248 (029)82665249

投稿热线:(029)82665397

读者信箱:banquan1809@126.com

版权所有 侵权必究

推荐序

预测性文本挖掘基础(Fundamentals of Predictive Text Mining)为一本介绍文本挖掘基础的工具书,引导这个迅速发展的领域。本书结合许多不同之主题,跨越各个学科及研究平台(数据挖掘、机器学习、数据库和计算语言学等),理论与实务兼备。

文本挖掘(Text Mining)又可称为文字知识发掘(Knowledge Discovery from Text ,KDT)或是文件信息探勘(Document Information Mining),其应用了信息检索、信息萃取、计算语言学、自然语言处理、数据挖掘……等。文本挖掘特别着重于利用这些技术,自非结构或半结构的文字中发掘出先前未知、隐含而有用的信息,Dan Sullivan (2001)定义文本挖掘为“一种编辑、组织及分析大量文件的过程,为了提供特定用户特定的信息,以及发现某些特征及其间的关联”。相较于传统的数据挖掘,文本挖掘需要加上额外的数据选择处理程序,以及复杂的特征萃取步骤。文本挖掘整合了许多传统信息检索技术,包括了关键词萃取、全文检索、档自动分类、自动摘要等等,以提供文字处理更强大的功能。故本书为文本挖掘入门最佳的书。

本书在 Dr.Sholom M.Weiss & Dr.Nitin Indurkhyia 以及 Dr.Tong Zhang 努力合作下,理论与实务兼容并包,不论对初学者还是深入研究者都是很好的入门教材及继续研究文本挖掘的经典参考教材。

谢邦昌 敬上,于 2012 年 9 月
台湾辅仁大学商学研究所博士班 所长
台湾辅仁大学统计资讯学系暨应用统计所 教授
中华资料采矿协会 理事长

译者序

电脑的普及,使得越来越多的文档以数字化形式出现,互联网的发展使得这些数字化文档得以快速传播。如今海量电子化文档的规模仍在急剧增长,这些电子化文档(例如研究报告、学术论文、在线文献库、E-mail、Web 页面、公司内部公告、会议纪要等)包含了大量的信息,也是重要的知识源。但是很多情况下,由于文档的数量非常庞大、缺乏组织整理并且格式多种多样,导致人们无法充分利用这些数据。这就要求开发自动化的大规模文本数据分析技术,来帮助我们处理和利用这些文档。文本挖掘就是研究如何从这些文件中提取信息,也就是关注分析非结构化自然语言文字的过程。作为一个新的正在迅速成为热点的研究领域,它致力于从文本数据中发现新的事实和知识,帮助我们获得有价值的信息。

预测性文本挖掘研究涉及不同学科的论题,如数据挖掘、机器学习、模式识别、信息处理、数据库和计算语言学等,目的是寻找隐藏在文本数据之中具有预测功能的模式。目前对于预测问题的研究比较广泛,并且已经获得很好的算法,但这些算法专门处理结构化的数值数据,而文本文件往往是非结构化数据,这就要求研究新的方法处理文本信息。

呈现在读者面前的这本《预测性文本挖掘基础》,是近年来国际上少见的专门介绍预测性文本挖掘技术的入门性教材和指南,原文作者总结了多年来在这一领域研究的大量优秀成果。作为入门级读物,论述由浅入深、理论结合实践、语言风趣、样例详实,详细地分析和总结了该领域的研究现状及未来发展。深入地讨论了文档分类、信息检索、聚类与组织文档、信息提取、基于 Web 的数据源、预测与评价等方面提出的问题。时值互联网和云计算蓬勃发展时代,为关注海量非结构化信息处理的学习者了解文本挖掘研究前沿打开了一扇窗口。

本书前言、摘要、第 2、4、6、8 章以及封底、附录及索引等由赵仲孟副教授

负责翻译,第1、3、5、7、9章等由侯迪副教授负责翻译。赵仲孟副教授负责对全书内容进行整理、修改和校审。参与本书翻译和修订工作的有王文科、沈辉、曹志、张新峰、辛蓉、岳鹏、张丹、杨丁福、赵静、赖文清等硕士研究生,特别是王文科同学对本书修订做出了巨大贡献。李颖编辑也给予了宝贵的修改意见,她大量专业而细致的工作保证了本书的质量。在此向他们表示衷心感谢!

本书适用于大学高年级本科生或者研究生,既可作为计算机科学、软件工程、信息技术等专业相关课程的教材,也可作为文本挖掘爱好者及工程研究人员的参考书。

限于译者的经验和水平,书中难免有诸多纰误与不足,欢迎读者通过电子邮件提出您的宝贵意见!

赵仲孟、侯迪
zmzhao@mail.xjtu.edu.cn
2012年夏于西安

前 言

五年前,拙作《文本挖掘:非结构化数据预测方法》付梓。那时,我们认为读者是有文本挖掘经验的业内人士,后来学生也在老师的指导下拿它当教材来用。《文本挖掘:非结构化数据预测方法》是最早尝试去收集整理预测性文本挖掘的相关资料之一,它包含了丰富的内容。自出版至今,在众多老师的帮助之下,我们收到了大量学生读者的积极反馈。近年来,互联网为我们提供了越来越多的数据资源,随之而来文本挖掘领域也出现了大量有趣新颖的技术,考虑到这些新的内容,再版实有必要。一年前,出版商希望《文本挖掘》再版,以使它更好地成为一本教科书。所以我们改进了许多章节,添加了新的能够反应当今技术发展的内容,同时也增加了习题和总结部分。

所谓预测问题,即寻找隐藏在数据之中的具有预测功能的数据,目前对于预测问题的研究已经广泛展开,并且已经得到了很多的算法。这些算法专门处理结构化的数值数据,用统一的衡量方法分析被采样的数据。然而,文本往往是非结构化数据,显而易见,文本和数值数据是不同的,这就要求用新的方法处理文本信息。然而我们认为,预测问题有通用的解决方法,无论是对于结构化数据还是非结构化数据而言。文本可以转化成数值格式,例如对于一个词而言,我们可以用存在(1)或者不存在(0)来表示它。用于数据挖掘的预测方法已经被证明同样可以很好地应用于文本预测。当然,二者之间也有一些关键性的区别。评价技术也必须适应数据发布的时序,并且改变评估错误的方法。因为我们要处理的数据是文档,所以要选择一些更特殊的适应分析文本的方法。更重要的是,还需要考虑到我们要处理的数据很可能是高维的、成千上万的单词和文档,但是原理不变。

我们将文本挖掘看成是综合了不同领域的任务,而不再仅仅是依赖于语言学和计算机科学的“自然语言处理过程”,也不再是“区别于其他机器学习的搜索引擎。”我们的视野更加开阔,我们支持你个性化的分析、学习数据,无论

是文本还是数值。因特网上有大量的大文本集合,其中蕴藏了海量等待挖掘的信息,可惜我们今天的语言处理技术还不足以胜任。有一些学者在试图通过发现语义理解的本质来解决文本挖掘问题,我们现在仅仅试图在大容量的文档数据集中寻找可复用的模式。

我们讨论的大多数内容都经过文本挖掘方法的严格验证和一些著名分析方法的应用。本书提供了丰富的实例和软件环境,并且介绍了大量的有实际学术研究价值的内容。本书追求实践,同时也包含了广泛的对文本挖掘有帮助的内容,既涵盖了预测学习方法,又引出了信息检索、搜索引擎和聚类等技术。如果读者可以跟随这本书,使用我们提供的软件动手分析每一实例,相信一定可以大有收获。

随着分析方法的高度发展,文本预测挖掘的应用领域会日新月异。本书总结了笔者多年的经验,并且提供最新的工具和技术,让读者更好地学习和实践。

读者:

本书适用于 IT 开发人员,管理人员,同时也适用于计算机专业的研究生。了解数据挖掘方面的背景知识有助于更好地阅读本书,但这并不是必须的。一些讨论高级概念的章节所涉及的数学知识都附有解释和说明,绝大多数章节对于具有研究分析精神的读者而言都易读易懂。如果你致力于此领域的研究,那么本书将开拓你的视野;如果你志在文本挖掘开发,那么你能从我们推荐的算法和我们举出的实例中有所斩获。我们提供的软件环境要求用户熟悉命令行编程和编辑配置文件。

说明:

这本书所包含的内容已经成功地用于教学工作,从为期一周的短期课程到十二周的全学期课程,都可以使用本书。如果是短期课程,本书所涉及的数学分析内容可以跳过。书中的练习题已经经过了几年的测试。每一章的后面都有与之配套的章节总结、习题。

另外,本书所举的大量例子和图表都可以在 data-miner.com 上免费下载。

其他一些网络资源:

Data-Miner 有限公司为本书的读者提供了一份免费软件,该软件实现了本书介绍的很多算法,可以从 data-miner.com 上下载。大量实例的 Linux scripts 也可以从这里下载。详情请参考附录 A。

致谢：

感谢我的同事 Fred Damerau, 他是我的良师益友, 他也是本书的共同作者。他不仅和我们共同完成了本书的写作, 还为我们的项目做出了杰出贡献, 特别是他在语言学方面的造诣, 给我们带来了难以言表的巨大帮助。本书中第 7 章中的一些例子基于我们之前出版的《文本挖掘》一书, 我们也要感谢那本书的共同作者: Chidanand Apté, Radu Florian, Abraham Ittycheriah, Vijay Iyengar, Hongyan Jing, David Johnson, Frank Oles, Naval Verma, Brian White. Arindam Banerjee 为本书提供了大量建议。书中的习题设计参考了 Statistics. com 中我们定期开设的文本挖掘课程里的内容。最后感谢我们的编辑, Wayne Wheeler, Ann Kostant 和 Wayne Yuhasz, 感谢他们的支持和帮助。

美国, 纽约

Sholom Weiss and Tong Zhang

澳大利亚, 悉尼

Nitin Indurkhy

目 录

推荐序

译者序

前言

第1章 文本挖掘概述	(001)
1.1 文本挖掘有什么特别之处?	(001)
1.1.1 结构化或非结构化数据?	(002)
1.1.2 文本数据是否不同于数值数据?	(003)
1.2 文本挖掘可以解决什么类型的问题?	(005)
1.3 文本分类	(006)
1.4 信息检索	(007)
1.5 文档聚类与组织	(008)
1.6 信息提取	(009)
1.7 预测与评估	(010)
1.8 下一章内容	(010)
1.9 小结	(011)
1.10 历史与文献评述	(011)
1.11 问题与练习	(012)
第2章 从文本信息到数值向量	(013)
2.1 文档收集	(013)
2.2 文档标准化	(015)
2.3 标记化	(017)
2.4 词形转化	(020)
2.4.1 词干变形	(020)
2.4.2 化词干为词根	(021)

2.5 预测向量生成	(022)
2.5.1 多词特征	(028)
2.5.2 正确答案的标签	(030)
2.5.3 通过属性分级选择特征	(031)
2.6 语句边界确定	(032)
2.7 词性标签化	(033)
2.8 词义消歧	(035)
2.9 短语识别	(035)
2.10 命名实体识别	(036)
2.11 语法分析	(036)
2.12 特征生成	(038)
2.13 小结	(039)
2.14 历史与文献评述	(040)
2.15 课后练习	(042)
第3章 用文本进行预测	(043)
3.1 识别文档符合模式	(045)
3.2 需要多少文档才可以满足预测需求？	(047)
3.3 文档分类	(048)
3.4 从文本中学习预测	(049)
3.4.1 相似性与最近邻法	(050)
3.4.2 文档相似性	(051)
3.4.3 决策规则	(053)
3.4.4 决策树	(059)
3.4.5 概率估计	(061)
3.4.6 线性评分方法	(063)
3.5 性能评估	(072)
3.5.1 当前与未来的性能估计	(072)
3.5.2 从学习方法中获取最大收益	(074)
3.6 应用	(075)
3.7 小结	(075)
3.8 历史与文献评述	(076)

3.9 问题与练习	(078)
第4章 信息检索和文本挖掘	(079)
4.1 信息检索是文本挖掘的一种形式吗?	(079)
4.2 关键字搜索	(080)
4.3 最近邻法	(081)
4.4 度量相似度	(082)
4.4.1 相同单词计数	(082)
4.4.2 单词计数和奖励	(083)
4.4.3 余弦相似度	(084)
4.5 基于 Web 的文档搜索	(085)
4.5.1 链接分析	(086)
4.6 文档匹配	(090)
4.7 反向列表	(090)
4.8 性能评估	(093)
4.9 小结	(094)
4.10 历史与文献评述	(094)
4.11 问题与练习	(095)
第5章 文档集的结构发现	(096)
5.1 基于相似性的文档聚类	(098)
5.2 复合文档的相似度	(099)
5.2.1 k -means 聚类	(101)
5.2.2 分层聚类	(105)
5.2.3 EM 算法	(107)
5.3 聚类标记有什么含义?	(111)
5.4 应用	(113)
5.5 性能评价	(114)
5.6 小结	(116)
5.7 历史与文献评述	(116)
5.8 问题与练习	(118)
第6章 在文档中查询信息	(119)
6.1 信息提取的目标	(119)

6.2	发现文本模式和实体	(122)
6.2.1	实体提取作为序列标签	(122)
6.2.2	标签预测作为分类	(123)
6.2.3	最大熵方法	(125)
6.2.4	语言特征和编码	(130)
6.2.5	局部序列预测模型	(132)
6.2.6	全局序列预测模型	(135)
6.3	共指和关系提取	(137)
6.3.1	共指消解	(137)
6.3.2	关系提取	(139)
6.4	模板填充和数据库构建	(140)
6.5	应用	(141)
6.5.1	信息检索	(141)
6.5.2	商业化提取系统	(142)
6.5.3	犯罪学	(143)
6.5.4	情报工作	(143)
6.6	总结	(145)
6.7	历史与文献评述	(145)
6.8	问题与练习	(147)
第 7 章	面向预测的数据源：数据库、混杂数据与 Web	(148)
7.1	数据的理想化模型	(148)
7.1.1	预测的理想化数据	(148)
7.1.2	理想的文本数据与非结构化数据	(149)
7.1.3	混杂数据与混合数据	(150)
7.2	实际数据源	(151)
7.3	原型化实例	(153)
7.3.1	基于 Web 的电子表格数据	(153)
7.3.2	基于 Web 的 XML 数据	(154)
7.3.3	观点数据与情绪分析	(157)
7.4	混杂数据实例：独立来源的数值数据与文本数据	(159)
7.5	采用标准表格格式的混合数据	(161)



7.6 总结	(163)
7.7 历史与文献评述	(163)
7.8 问题与练习.	(164)
第8章 实例分析	(165)
8.1 互联网市场调研	(165)
8.1.1 问题描述	(165)
8.1.2 解决概览	(166)
8.1.3 方法与过程	(167)
8.1.4 系统部署	(168)
8.2 面向数字图书馆的轻型文档匹配	(169)
8.2.1 问题描述	(169)
8.2.2 解决概览	(170)
8.2.3 方法与过程	(171)
8.2.4 系统部署	(172)
8.3 生成帮助桌面应用的模本范例	(173)
8.3.1 问题描述	(173)
8.3.2 解决概览	(174)
8.3.3 方法与过程	(174)
8.3.4 系统部署	(176)
8.4 新闻文章主题指定	(176)
8.4.1 问题描述	(176)
8.4.2 解决概览	(177)
8.4.3 方法与过程	(178)
8.4.4 系统部署	(181)
8.5 邮件过滤	(181)
8.5.1 问题描述	(181)
8.5.2 解决概览	(182)
8.5.3 方法与过程	(183)
8.5.4 系统部署	(184)
8.6 搜索引擎	(185)
8.6.1 问题描述	(185)

8.6.2	解决概览	(185)
8.6.3	方法与过程	(186)
8.6.4	系统部署	(187)
8.7	文档中命名实体提取	(188)
8.7.1	问题描述	(188)
8.7.2	解决概览	(189)
8.7.3	方法与过程	(189)
8.7.4	系统部署	(191)
8.8	个性化报纸	(192)
8.8.1	问题描述	(192)
8.8.2	解决概览	(193)
8.8.3	方法与过程	(193)
8.8.4	系统部署	(194)
8.9	总结	(195)
8.10	历史与文献评述	(195)
8.11	问题与练习	(196)
第9章	新研究方向	(197)
9.1	摘要	(197)
9.2	主动学习	(200)
9.3	使用未标记的数据学习	(202)
9.4	收集文档样本的不同途径	(202)
9.4.1	文档集合与投票方法	(203)
9.4.2	在线学习	(204)
9.4.3	代价敏感学习	(206)
9.4.4	不稳定样本与罕见事件	(207)
9.5	分布式文本挖掘	(207)
9.6	学习排序	(209)
9.7	问答系统	(210)
9.8	总结	(212)
9.9	历史与文献评述	(212)
9.10	问题与练习	(214)

附录 A 软件说明	(215)
A. 1 软件概要	(215)
A. 2 系统需求	(216)
A. 3 下载说明	(216)
参考文献	(217)
作者索引	(225)
主题索引	(229)

第1章

文本挖掘概述

1.1 文本挖掘有什么特别之处？

你缺乏数据吗？很显然不太可能。由于计算机的普及，目前大部分数据都以电子表格存储。p.1 我们使用计算机系统交易股票、使用文本编辑软件进行写作、在线购买商品，这些行为都离不开计算机。越来越多的纸上交易现在都在朝着无纸化方向发展，因此未来我们将会拥有以电子表格形式存在的“海量数据”等待我们分析。

数据挖掘建立在收集和存储海量数据基础上，它致力于寻找隐藏在数据之中的有用模式。现如今，数据挖掘已不是什么方兴未艾的新兴技术，虽然在业界远未普及，但是，数据挖掘领域理论技术已经高度发达，对于某些问题解决也几近成熟。

我们很想说“只需要提供给我数据，我就能告诉你数据隐藏的模式”。不幸的是，现有数据挖掘方法多数会要求被分析数据格式高度结构化，因此必须对于原始数据进行预处理。最好原始数据本身就是高度结构化的，否则我们必须对原始数据进行格式转换。

数据挖掘依赖于我们的经验。对于致力于做预测数据挖掘的专家而言，他们所处理的数据往往是数值型的，这些人每天和数字打交道。而“文本挖掘者”却从未期望得到一系列有序数字。他们的目标是寻找到一系列文本文件集合，并且这些集合内容是可读的，同时意义明显。

数据挖掘和文本挖掘之间的第一个区别：前者针对数字而后者专注于文