



大数据 解读

邓莎莎
著



支持决策研讨的文本
分析方法研究

中西書局

大数据解读

邓莎莎
著

常州大学图书馆
藏书章

支持决策研讨的文本
分析方法研究

中西書局

图书在版编目(CIP)数据

解读大数据：支持决策研讨的文本分析方法研究 /
邓莎莎著. —上海：中西书局，2017.10
ISBN 978 - 7 - 5475 - 1326 - 2

I. ①解… II. ①邓… III. ①数据处理—研究 IV.
①TP274

中国版本图书馆 CIP 数据核字 (2017) 第 215486 号

本书出版得到上海汽车工业教育基金会资助

解读大数据

——支持决策研讨的文本分析方法研究

邓莎莎 著

责任编辑 伍珺涵

装帧设计 梁业礼

出版发行 上海世纪出版集团

中西书局(www.zxpress.com.cn)

地 址 上海市陕西北路 457 号(200040)

印 刷 上海世纪嘉晋数字信息技术有限公司

开 本 700×1000 毫米 1/16

印 张 14

字 数 156 000

版 次 2017 年 10 月第 1 版 2017 年 10 月第 1 次印刷

书 号 ISBN 978 - 7 - 5475 - 1326 - 2 / T · 014

定 价 49.00 元

本书如有质量问题,请与承印厂联系。T: 021-69214196

表格目录

表 1-1	用于文本分析的各种语言特征	011
表 2-1	群体研讨文本连贯性分析研究文献归纳分析	058
表 2-2	Searle 的言语行为概述	063
表 2-3	候选发言对的特征集	074
表 3-1	会话主题拆解实验结果	091
表 3-2	连贯性分析特征对比实验结果	095
表 3-3	连贯性分析方法对比实验结果	097
表 3-4	标注者之间的信度值	104
表 3-5	言语行为分类标注结果	104
表 3-6	言语行为分类对比实验结果	106
表 3-7	SATrees 与 Conversation Tree 的信息质量 实验结果	109
表 3-8	社会网络中心性测量的平均绝对百分比误差	111
表 3-9	用户实验中两组研讨文本的统计指标	116
表 3-10	受试者基本信息	117
表 3-11	所有 8 个问题的用户实验结果	119

表 3-12	Action, Situated Action 及 Symbolic Action 问题回答结果统计分析	119
表 4-1	虚假评论总数信息	139
表 4-2	真实评论总数信息	140
表 4-3	语言特征和度量	146
表 4-4	不同的特征集上分别使用 SVM、NB 和 C4.5 分类器的准确率、精度、召回率和 F 值	151
表 4-5	采用其他表示信息丰富度方法的平均特征 权重	153
表 4-6	采用词性分布的平均特征权重	154
表 4-7	采用语言接近度的平均特征权重	155
表 5-1	两个研讨问题的网络民意建模过程	174

插图目录

图 1-1	相邻语轮对紊乱的例子	010
图 1-2	Toulmin 逻辑结构图	025
图 1-3	认知适应理论示意图	027
图 1-4	本书结构	042
图 2-1	支持在线研讨意义构建的基于 LAP 的文本 分析框架	052
图 2-2	支持在线研讨意义构建的基于 LAP 的文本 分析系统(LTAS)	068
图 2-3	研讨文本切割算法中各种变量的定义	072
图 2-4	在线研讨切割算法	072
图 2-5	TBL-RM 中的剩余匹配算法	078
图 2-6	言语行为分类方法示意	081
图 2-7	包含回复关系和言语行为的 SATrees	083
图 3-1	GSS 系统原型系统用户界面	087
图 3-2	20 个会话 DSA 算法与其他算法 F 值比较	092
图 3-3	20 个会话 TBL-RM 方法和其他方法的比较 结果	098

图 3-4	群体研讨会例子中的社会网络	112
图 4-1	网络推广平台界面	125
图 4-2	在线评论欺骗识别系统构架	147
图 5-1	面向研讨问题的网络民意分析系统框架	164

目 录

表格目录	001
插图目录	001

第1章 绪论

1.1 研究背景	002
1.2 研究对象与意义	006
1.2.1 研究对象	006
1.2.2 研究意义	007
1.3 研究文献综述	008
1.3.1 研讨文本数据分析简述	008
1.3.2 在线研讨中欺骗问题研究简述	015
1.3.3 社会化媒体数据集成研究简述	020
1.3.4 决策支持系统研究简述	024
1.3.5 研究评述	034
1.4 研究问题的提出	036

1.5	研究思路与本书框架	037
1.5.1	研究思路	037
1.5.2	研究方法	037
1.5.3	研究内容	039
1.5.4	技术路线	042
1.6	本书主要创新之处	043

第2章 支持在线研讨意义构建的文本分析方法研究

2.1	引言	048
2.2	在线研讨中的意义构建与语言行为视角	049
2.2.1	意义构建理论	049
2.2.2	语言行为视角	050
2.2.3	基于 LAP 的文本分析框架	052
2.3	研究假设	053
2.3.1	会话主题拆解理论假设	053
2.3.2	连贯性分析理论假设	055
2.3.3	言语行为分类理论假设	063
2.3.4	意义构建理论假设	064
2.4	基于 LAP 的文本分析系统	067
2.4.1	会话主题拆解	069
2.4.2	连贯性分析	073
2.4.3	言语行为分类	078
2.4.4	言语行为树	082
2.5	小结	084

第3章 LTAS 系统的实验与评估

3.1	实验总体设计	086
3.2	实验 1: 会话主题拆解算法的实验与评估	088
3.2.1	实验设计	088
3.2.2	性能指标	089
3.2.3	结果与讨论	090
3.3	实验 2: 连贯性分析方法的实验与评估	093
3.3.1	实验 2a: 连贯性分析特征	093
3.3.2	实验 2b: 连贯性分析方法对比实验	095
3.4	实验 3: 言语行为分类算法的实验与评估	100
3.4.1	言语行为类别定义	101
3.4.2	言语行为人工标注	102
3.4.3	结果与讨论	105
3.5	实验 4: 面向意义构建的信息质量分析	107
3.5.1	实验 4a: SATrees 与 Conversation Tree 的 准确性比较	107
3.5.2	实验 4b: 社交网络中心性测量	109
3.6	实验 5: 面向意义建构的用户实验	113
3.6.1	面向意义建构的问卷设计	113
3.6.2	实验设计	114
3.6.3	测试实验与数据收集	116
3.6.4	结果与讨论	118
3.7	小结	120

第4章 在线研讨过程中欺骗识别研究

4.1	引言	124
4.2	在线评论欺骗行为相关研究	126
4.2.1	欺骗的定义	126
4.2.2	欺骗理论	127
4.2.3	文体分析研究	133
4.2.4	文本分类方法概述	134
4.3	数据集的构建	138
4.3.1	虚假评论的构建	138
4.3.2	真实评论的构建	140
4.4	特征选取	140
4.4.1	词语词频	141
4.4.2	信息丰富度	141
4.4.3	内容信服度	144
4.4.4	特征汇总	145
4.5	在线评论欺骗识别系统设计	146
4.5.1	系统架构	146
4.5.2	预处理	147
4.5.3	特征抽取	148
4.5.4	文本分类	148
4.6	结果与讨论	149
4.6.1	实验设计	149
4.6.2	三种分类算法的实验结果	150
4.6.3	词语词频特征集分析	152

4.6.4	感觉特征集分析	152
4.6.5	词性特征分析	153
4.6.6	语言接近程度特征分析	155
4.6.7	分类技术比较分析	155
4.7	小结	156

第5章 面向研讨问题的网络民意分析研究

5.1	背景介绍	158
5.2	相关研究	160
5.2.1	网络民意与网络舆论	160
5.2.2	文本意见挖掘	162
5.2.3	在线研讨中的利益相关者	162
5.3	面向研讨问题的网络民意建模	164
5.3.1	网络评论网页采集	165
5.3.2	HTML 页面的解析	165
5.3.3	面向研讨问题的主题分析	167
5.3.4	相似度计算	168
5.3.5	利益相关群体提取	169
5.3.6	情感分析	170
5.3.7	“主题—利益相关群体—情感”模型	171
5.4	应用案例与分析	172
5.4.1	网络评论数据	172
5.4.2	决策问题解析	173
5.4.3	网络民意建模	173
5.5	小结	176

第6章 结论

6.1 本书的主要工作与创新点	178
6.2 研究不足及展望	181
附录 1 意义构建实验问卷	183
附录 2 虚假评论问卷	186
附录 3 真实评论问卷	189
附录 4 ICTPOS 词性标注集及含义	192
参考文献	193

绪论

第 1 章

00010000

110000011

00100
00100
0000010001000000100000010000001000000
0001100000100000001100000010001000000010000011000
00000010000100000010000001000000100000010000001000

1.1 研究背景

时至今日,互联网、物联网和云计算等信息和通信技术(Information and Communication Technologies, ICTs)在社会中的渗透和应用是如此深广。中国互联网络信息中心(CNNIC)于2013年1月15日在京发布的《第31次中国互联网络发展状况统计报告》(以下简称为《CNNIC第31次调查报告》)显示,截至2012年12月底,中国网民数量达到5.64亿,互联网普及率为42.1%,较2011年底,网民增量为5090万,普及率提升3.8个百分点,手机超越台式电脑成为中国网民第一大上网终端。到2012年底,通过台式电脑,包括笔记本上网的网民占上网网民的70.6%,这个数字比2011年下半年的统计数字低了近3个百分点。相比而言,手机上网人数则是增加了74.5%,而这一数字已经超过了使用台式机器的上网人数,截至2012年底,我国利用手机上网的网民规模已经达到了4.2亿。智能手机的普及与移动互联的不断发展为中国广阔的农村地区以及庞大的移动人口提供了更为廉价和简便的接入互联网的方式。

伴随着Web2.0技术的迅速发展,各种提供交流与沟通的社会化媒体(Social Media)平台不断涌现。无论是协同社区(例如:Wikipedia、博客、微博和网络论坛等)和社交网络社区(例如:Facebook、人人网、开心网等),还是各种特定主题的内容社区(例

如：Youtube、豆瓣、点评网等），用户已经从被动地接收信息转换成主动地输出信息，用户的参与度与贡献度已经显著提高。此外，随着移动互联的普及，利用手机即时通信的移动性、碎片性以及随时在线的特点，各种社会化媒体平台纷纷植入了移动互联技术。这种方式使得个人与社区之间的共享、共同创造以及修改用户产生内容变得更加频繁和容易。

社会化媒体的快速增长带来了人们日常交流、决策讨论方式的巨大改变。面对面的会议方式不再是组织内部讨论问题、制定决策的唯一沟通手段。使用社会化媒体技术来支持与商业相关的决策研讨的方式已经被越来越多的组织所采用(Mann 2011a)。根据最近麦肯锡季刊(*McKinsey Quarterly*)的报道，在超过 1 700 个受访的组织中，50% 的组织正在使用社交网络(Social Networking)，41% 的组织在使用博客(Blogs)，25% 的组织正在使用维基(Wikis)，还有 23% 的组织在使用微博(Micro Blogs)。在过去的 4 年中，这些数字正在以超过 2 倍的速度增长。(Bughin et al. 2010)

无论是在组织内部还是组织外部，社会化媒体平台支撑的在线研讨方式从各个方面为组织带来了巨大的收益。在组织内部，社会化媒体平台能够更加灵活地支持组织或者部门内部的决策研讨，提高决策效率，从而加快组织获取知识的速度，降低沟通和操作成本，提高内部专家的识别率并增加员工的满意度。另外，社会化媒体已经逐渐成为与组织外部的客户、商业合作伙伴进行沟通的重要手段。支持组织内部成员与外部顾客的沟通能够有效地提高顾客对组织的知晓度以及忠诚度，降低营销费用和客户服务成本，并且增加新产品或者服务的创新成功率。(Bughin et al. 2010)

尽管社会化媒体技术支持的在线研讨给组织带来的价值巨大,并且已经有很多组织从这种交流沟通方式中得到了实实在在的收益,然而,在互联网环境中研讨文本分析面临着以下四个方面的挑战。

第一,数据量过大且内容复杂,而这超出了人们处理信息的能力。

社会化媒体平台支撑的在线研讨,常常随着事件的发生在短时间内发言数激增,一时间相关信息铺天盖地地出现在互联网中,各种累积的发言数量越来越庞大。因此,仅靠人工阅读所有研讨文本内容,继而综合分析并形成判断进而提供决策建议是一件耗时耗力同时也是非常困难的任务。根据 Hiltz 和 Turoff 的研究,当信息量达到某一信息熵(information entropy)时,群体或者个人就无法有效地组织和理解输入信息。(Hiltz et al. 1985)这时,人们就会抛弃所有信息而完全依靠对前面部分信息产生的主观印象和感受进行判断和决策。因此,仅靠人工处理这些海量信息不仅给处理人员带来极大负担,更重要的是,在很多时候这还是一件不可能完成的工作。

第二,如何快速有效地弄清发言之间的关系,并获得语言文字所要表达的真正意图?

首先,社会化媒体平台支撑的在线研讨产生的数据都是非结构的语言文字,并且这些内容并不是由一个人独立完成,而是由许多人协作完成的。譬如,对某一个问题的讨论常常是由微博博主与成千上万的粉丝互动完成的,论坛中对于热点问题的讨论常常会吸引成千上万的人参与,这时,网络文本已经被大量发言撕扯得支离破碎,上下文之间毫无逻辑关系,难以读懂。其次,语言文字