

基于话题自适应的中文微博情感分析

任远 巢文涵 周庆 李舟军

(北京航空航天大学计算机学院 北京 100191)

摘要 近年来,随着社会网络的迅速兴起,面向社会网络的情感分析技术逐渐成为数据挖掘领域新的研究热点。中文微博以其语言简短、文法灵活的特点,给情感分析的研究工作带来了新的挑战。对数据预处理、情感词典构造、话题元素引入等中文微博情感分析技术进行了系统的研究,提出了给情感词分级的方法以提升情感分析的准确度;同时提出了面向话题的自适应方法以更准确地识别情感词;最后实验结果验证了以上方法的有效性。

关键词 中文微博,情感分析,情感词,话题自适应

中图分类号 TP391 文献标识码 A

Sentiment Analysis of Chinese Microblog Using Topic Self-adaptation

REN Yuan CHAO Wen-han ZHOU Qing LI Zhou-jun

(School of Computer Science and Engineering, Beihang University, Beijing 100191, China)

Abstract Recently, with the rapid development of social networks, sentiment analysis over social networks has gradually become a new hot research topic, especially in the field of data mining. The typical features of Chinese microblog (such as “short” and “flexible”) bring some new challenges for the researcher to analyze its sentiment. So this paper carried out a systematic study on Chinese microblogging emotional analysis technology, including data preprocessing, sentimental lexicon construction, topic adjunction. In addition, to improve the precision of sentiment analysis, a novel emotional words classification approach was proposed. Meanwhile, we proposed a topic-oriented adaptive method to promote the work of emotional words identification. And the experimental results demonstrate the feasibility and effectiveness of our approach.

Keywords Chinese microblog, Sentiment analysis, Sentimental words, Topic self-adaptation

1 引言

随着 Web2.0 的蓬勃迅猛发展,越来越多的用户习惯于在互联网上发布自己的观点、分享个人生活体验、与朋友交流互动,由此推动了社会网络的快速发展。微博,作为一个综合性的社会网络平台,给用户提供了便捷的网上沟通渠道与丰富而新鲜的信息内容,迅速成为最受网民欢迎的新兴社会网络媒体。以我国的新浪微博为例,据统计截至 2012 年 12 月,新浪微博用户数量超过 5 亿,日活跃用户量达到 4620 万,可见微博已经成为人们生活中不可或缺的一部分,并在一定程度上影响着人们的日常工作和生活。正是由于微博具有如此庞大的活跃用户量,使得其蕴含的数据与信息规模相当庞大,仅靠传统的人工方法已无法对这些海量的数据与信息进行有效的收集和处理,因此迫切需要利用计算机技术来帮助用户对其进行分析与处理。在此背景下,社会网络分析技术应运而生,并迅速引起学术界和互联网企业的广泛关注和高度重

视,而面向社会网络的情感分析技术就是其重要的组成部分和研究热点。

情感分析是对带有情感色彩的主观性文本进行分析、处理、归纳和推理的过程,目的是为了从用户发布的带有主观情感色彩的文本信息中提取用户的观点,并判断其情感极性。情感分析一般分为情感信息抽取和情感倾向分类两大任务。目前针对英文微博(如 twitter)的情感分析技术已经取得了重要的进展,积累了不少研究成果。但面向中文微博的情感分析技术的研究尚处于初级阶段。由于中英文在语法、语义和语用等方面存在着巨大的差异,使得在处理中文微博文本的过程中面临着更多的问题。例如,英文微博限制每条文本最多包含 140 个英文字符,大约相当于 7~15 个英文单词组成的一整句话,表达的话题单一、情感极性容易判别;中文微博规定每条文本最多包含 140 个中文字符,这其中可以有几句话,每句话表达的话题和情感可以各不相同。例如下面一条中文微博:“今天洛杉矶湖人队的比赛真是糟糕透了!加索

到稿日期:2013-02-10 返修日期:2013-04-22 本文受国家自然科学基金项目(61170189,61370126,61202239),高等学校博士学科点专项科研基金(20111102130003),软件开发环境国家重点实验室自选课题(SKLSDE-2013ZX-19),中央高校基本科研业务费专项基金(YWF-13-T-RSC-072)资助。

任远(1987-),男,硕士生,主要研究方向为数据挖掘与社交网络分析,E-mail:renyuan@cse.buaa.edu.cn;巢文涵(1979-),男,博士,讲师,主要研究方向为数据挖掘、自然语言处理、机器翻译;周庆(1991-),男,主要研究方向为数据挖掘;李舟军(1963-),男,博士,教授,博士生导师,CCF 高级会员,主要研究方向为智能信息处理技术、信息安全技术。

尔状态低迷,但是科比的表现依旧出色!”在此条评论 NBA 比赛的微博中,作者用 3 句话分别针对 3 个不同的话题对象表达了不同的情感。因此,在进行中文微博情感分析时,不能完全依靠普通的情感分析方法,需要针对中文微博自由、灵活、口语化的特点,对微博数据进行更为细致的分析和处理。

本文主要研究面向句子级的中文微博情感分析技术。在情感词处理方面,相比于传统情感词极性的粗粒度划分(正向情感词、中性情感词、负向情感词),本文提出了更为精细的情感词极性程度的分级方法;同时,为提升情感分析的质量,针对中文微博短文本的特点,本文尝试引入话题因素,提出了面向话题的自适应方法,从话题相关文本中抽取出情感词作为话题情感词,辅助微博情感分析任务;最后,通过基于规则的和基于支持向量机(SVM)的情感倾向分类方法的对比实验,验证了本文方法的可行性。

本文第 2 节介绍情感分析相关研究工作;第 3 节详细阐述本文的方法原理和算法设计;第 4 节展示实验结果和相关分析;最后展望下一步工作。

2 相关工作

2.1 情感词典的构建

无论是英文还是中文,目前均没有一个完整的涵盖所有词语的情感词典。因此,情感词典的构建是情感分析中一项非常重要的基础工作。目前,情感词典的构建方法主要有两种:基于语料的方法和基于词典的方法。

基于语料的方法主要通过计算不同词语之间的相似度,并利用词语相似度计算词语语义倾向。Wiebe 等人^[1]运用词聚类的方法,基于大语料库完成了形容词词性的情感词语的获取,Riloff 等人^[2]手工制定一些模板并选取种子情感词语,使用迭代的方法获取了名词词性的情感词语,这两种方法均有一定的局限性;Turney 和 Littman^[3]提出了点互信息(Point Mutual Information)方法,用以实现情感词语的判别,这种方法适用于各种词性的情感词语的识别,但依赖于种子词语的选取。由此可见,基于语料库的方法简单易行,但可利用的语料库有限,且很难归纳情感词语在大语料库中的分布。

基于词典的方法通常采用语义词典来判断词语相似度,所用词典一般为 WordNet、HowNet 等。文献[4-6]利用词典将手工采集的种子情感词语进行扩展来获取大量的情感词语,这种方法简单易行,但是较依赖于种子情感词语的个数和质量,并且容易因一词多义现象引入噪声。为了解决上述问题,文献[7-10]使用词典中词语的注释信息来完成情感词语的识别与极性判断。此外,文献[11]沿用了 Turney 等人的点互信息方法,通过计算 WordNet 中的所有形容词与种子词代表(褒义词“good”和贬义词“bad”)之间的关联度值来识别出情感词语。

2.2 评价对象的抽取

评价对象是指评论文本中情感词所修饰的对象。部分学者^[12-14]使用基于规则(模板)的方法抽取评价对象。规则的制定通常要基于一系列的语言分析与预处理过程,如词性标注、命名实体识别和句法分析等。此类方法可以直接针对待解决的问题制定规则(模板),有较强的针对性;但规则(模板)的可扩展性差,成本较高。

从另一个角度分析,由于评价对象常常蕴涵于情感文本的某些话题中,也有一些学者^[15,16]使用话题模型进行评价对象的识别。文献[17]采用多粒度的话题模型来挖掘产品领域情感文本中的评价对象,并将相似的评价对象进行聚类。

2.3 基于机器学习的情感倾向分类

基于机器学习的情感倾向分类,是采用机器学习方法,通过对标注语料的训练生成倾向分类器,对测试文本进行分类。目前主流的分类方法有支持向量机(support vector machine, SVM)、朴素贝叶斯(naive Bayes, NB)和最大熵(maximum entropy, ME)等。其流程大致如下:先对文本的情感倾向性进行人工标注,提取文本特征表示,并将其作为训练集,利用机器学习的方法构造分类器,通过分类器就可得到待测文本的情感倾向性类别。Pang 等^[18,19]使用朴素贝叶斯、最大熵和 SVM 的分类方法,分别对 Usenet 中的电影评论进行文本的情感倾向性分类,并将它们和手工分类结果进行比较。实验结果显示, SVM 在几种分类方法中效果最好。

基于机器学习的情感倾向分类方法关键在于特征信息的有效提取,随着语义特征信息的加入和训练语料库的发展,基于机器学习的情感倾向分类将会有广阔的发展前景。

3 算法设计

本文的中文微博情感分析主要聚焦于句子级的情感分析,情感倾向分为正面、负面与中性 3 类。中文微博短文本具有简短灵活、口语化、话题丰富等特点,大大增加了对其进行情感分析的难度。针对此类问题,下面分别详细阐述本文在数据预处理、情感词极性程度分级、话题元素引入等方面所采用的方法。

3.1 数据处理流程

本文的数据处理流程如图 1 所示。

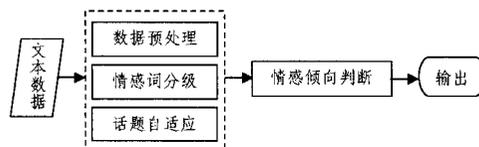


图 1 数据处理流程

其中,数据预处理主要完成规范化、特殊词处理等任务;情感词处理结合原始情感词库,形成分级情感词库;接下来加入话题元素,基于话题自适应的方法,进行情感倾向的判断。

3.2 数据预处理

由于中文微博文本简短,文字口语化特征明显,并存在大量的网络用语,相比于英文微博,其文本往往更不规范,增加了中文微博数据分析与处理的难度;而另一方面,微博上独有的表情标记可转换为相应的情感词语,在一定程度上有助于情感倾向的判断。由此可见,需要针对中文微博的上述特点,进行细致的数据预处理。本节主要讨论规范化处理以及特殊词的处理。

规范化处理,是指识别微博文本中的全角符号,并全部转换为半角符号;对于微博中大量存在的表情符号,文献[20,21]表明表情符和内容的情感具有很强的相关性,因此本文将表情符号按照其所具有的情感倾向(正向情感、负向情感)进行人工划分,结果如表 1 所列。

表1 表情符号分类

情感类别	数目	示例表情符
正向情感	50	
负向情感	45	

对于特殊词语的处理,主要针对汉语中起修饰作用的词,它们在一条语句中或增强、或减弱作者所要表达的情感。本文总结了3类修饰词,如表2所列。

表2 特殊词语汇总

类别	数目	示例词语
反问意义的词	5	难道、怎么会
讽刺意义的词	19	就算、竟然
表示程度的词	131	极度、有点

3.3 情感词极性程度分级

情感词是情感极性判定中较为重要的考量依据。目前公开的情感词典,绝大多数仅仅将情感词按照正向、负向、中性进行了粗粒度划分。而中文情感词的用法比较灵活,在不同的语言环境下可能表达不同的情感。为提高情感分析的精度,本文提出了更为精细的情感词极性程度的分级方法,并结合 Hownet 情感词、台湾大学情感词¹⁾和清华大学情感词²⁾来构建情感词典。

本文将情感词的极性程度(简称为情感词的极值)定义为情感词的“专注性”,规定若情感词修饰的对象越宽泛,则它在句子中作为修饰词语的可能性越低,其在情感词极性程度分级中所获得的等级越低;反之,则等级越高。例如,情感词“好”所能修饰的成分很多,当它表示为一个情感词时,它在相应语句中所表达的情感极性程度并不高;而情感词“任劳任怨”的专注程度高于“好”,在绝大多数情况下,它作为情感词所具有的正向情感极性也高于“好”,因此“任劳任怨”在分级中等级更高。

情感词的极性程度分级(即极值计算)的具体算法可表述如下。

算法1 极值计算

输入:包含情感词的句子

输出:情感词极值

- 第1步 对话料库中包含情感词的句子进行依存分析,找出包含情感词的依存关系对,记为二元组(情感词,关联词);
- 第2步 计算每个情感词对应的依存关系端点上的熵;
- 第3步 选定阈值,根据熵的大小将情感词的极性程度划分为5个等级,熵越小极值越大。

情感词极性程度分级样例,如表3所列。

表3 情感词极性程度分级样例

	等级1	等级2	等级3	等级4	等级5
负向词	冷、沉默	潜逃、诡异	讥讽、残废	气馁、内疚	萎靡不振、 骄奢淫逸
正向词	透明、决心	坚定、准时	诚挚、高雅	审慎、赞誉	实至名归、 栩栩如生

3.4 话题自适应

如本文第1节所述,博主可以在一条微博中对多种话题表达观点。针对此问题,本文融合话题特征以提高微博情感分析的精度。从微博数据中的 hashtag 获得话题元素,对于

特定话题,从互联网上爬取相关的文本,然后从中抽取出话题相关的情感词和评价对象,并计算话题整体极性,辅助微博情感分析任务。具体算法表述如下。

算法2 话题自适应算法

输入:包含情感词的句子及其话题

输出:句子中情感词的极值和

- 第1步 建立一个通用的情感词种子词集,本文选用的种子词有:公正、舒服、赞美、收获、团结、违法、冷漠、故障、邪恶、郁闷等。
- 第2步 根据具体的话题,以微博中的 hashtag 作为话题文本,从互联网搜索并爬取相关文本,以此作为话题语料库;
- 第3步 从话题语料库中迭代抽取话题相关的情感词和评价对象;
- 第4步 抽取过程中,同时统计句子个数、正向和负向的情感词个数及其极值和;
- 第5步 抽取完毕后,计算文本句子中的正向情感词的平均极值和、负向情感词的平均极值和以及句子的全部情感词的极值和;
- 第6步 以前5步得出的话题元素作为分类特征,加入情感分类的任务中。

其中,第3步中使用的抽取模板可形式化描述如下(依存关系工具为 Stanford Parser):

定义1 SentiWord 为所有通用情感词的集合。

定义2 Exp 为已抽取的情感词集合。

定义3 Target 为已抽取的评价对象集合。

定义4 x 为当前分析的词, N 为名词集合。

规则1 若 $e \rightarrow \text{relation} \rightarrow x, e \in \text{Exp}, \text{relation} \in \{\text{'conj'}, \text{'appos'}\}, x \in \text{SentiWord}$, 则 $x \in \text{Exp}$ 。

规则2 若 $e \rightarrow \text{relation1} \rightarrow H \leftarrow \text{relation2} \leftarrow x, \text{relation1} = \text{relation2}, e \in \text{Exp}, x \in \text{SentiWord}$, 则 $x \in \text{Exp}$ 。

规则3 若 $x \rightarrow \text{relation} \rightarrow t, t \in \text{Target}, \text{relation} \in \{\text{'subj'}, \text{'obj'}, \text{'amod'}, \text{'rmod'}, \text{'predet'}, \text{'acomp'}\}, x \in \text{SentiWord}$, 则 $x \in \text{Exp}$ 。

规则4 若 $x \rightarrow \text{relation} \rightarrow H \leftarrow \text{relation} \leftarrow t, t \in \text{Target}, \text{relation} \in \{\text{'subj'}, \text{'obj'}, \text{'amod'}, \text{'rmod'}, \text{'predet'}, \text{'acomp'}\}, x \in \text{SentiWord}$, 则 $x \in \text{Exp}$ 。

规则5 若 $s \rightarrow \text{relation} \rightarrow x, e \in \text{Exp}, \text{relation} \in \{\text{'subj'}, \text{'obj'}, \text{'amod'}, \text{'rmod'}, \text{'predet'}, \text{'acomp'}\}, x \in N$, 则 $x \in \text{Target}$ 。

规则6 若 $e \rightarrow \text{relation} \rightarrow H \leftarrow \text{relation} \leftarrow x, e \in \text{Exp}, \text{relation} \in \{\text{'subj'}, \text{'obj'}, \text{'amod'}, \text{'rmod'}, \text{'predet'}, \text{'acomp'}\}, x \in N$, 则 $x \in \text{Target}$ 。

规则7 若 $t \rightarrow \text{relation} \rightarrow x, t \in \text{Target}, \text{relation} \in \{\text{'conj'}, \text{'appos'}\}, x \in N$, 则 $x \in \text{Target}$ 。

规则8 若 $t \rightarrow \text{relation1} \rightarrow H \leftarrow \text{relation2} \leftarrow x, t \in \text{Target}, \text{relation1} = \text{relation2}, x \in N$, 则 $x \in \text{Target}$ 。

3.5 情感倾向判断

情感倾向判断是在已识别为观点句的基础上,判断句子的情感倾向是正面、负面还是中性。本文采用两种方法进行情感倾向的判断:1)基于规则的方法;2)基于 SVM 的方法。

3.5.1 基于规则的情感倾向判断

基于规则的方法,对每个观点句提取其所有的情感词,并对情感词的极性作判断。其中,情感词采用极性程度分级的方式,正向情感词取正的极值,负向情感词取负的极值。最后

¹⁾ <http://nlg18.csie.ntu.edu.tw;8080/lwku/index.html>

²⁾ <http://nlp.csai.tsinghua.edu.cn/site2/>

对一个句子的所有情感词进行加权处理,根据所得加权极值所处的区间,判断该观点句的情感倾向。

3.5.2 基于 SVM 的情感倾向判断

本文选用的极性分类特征如表 4 所列。

表 4 极性分类特征

序号	特征内容
1	句子中表情符号个数
2	正向表情符号的个数
3	负向表情符号的个数
4	句子是问句、陈述句、还是感叹句
5	情感词对应的词性标记分别为{"n","v","a","z","d"}的个数
6	情感词对应的词性标记分别为{"n","v","a","z","d"}的情感极值和
7	正向情感词个数
8	负向情感词个数
9	句子的情感词极值和
10	否定词的个数
11	感叹号个数
12	问号个数
13	非情感词的各类词性标记个数 n,t,s,f,v,a,b,z,r,m,q,d,p,c,u,e,y,o,h,x,w
14	话题相关的情感词个数 {"n","v","a","z","d"}等话题相关的情感词对应词性标记的个数
15	句子包含 target 数量
16	话题总体情感倾向
17	话题正向情感倾向
18	话题负向情感倾向

4 实验结果及分析

4.1 实验数据

本文采用 NLP&CC2012 中文微博情感分析的评测数据集对上述方法进行评测。该数据集共包含 2478 条观点句、448 条正向情感句、1994 个负向情感句以及 12 个中性情感句,如表 5 所列。

表 5 实验数据

新闻话题	微博数	正向句	负向句	中性句
90 后当教授	104	110	13	0
奖状植入广告	100	22	53	1
名古屋市长否认南京大屠杀	100	1	56	0
韩寒方舟子之争	106	21	111	0
学雷锋被钓鱼执法	100	6	94	0
疯狂的大葱	100	4	76	1
皮鞋果冻	100	2	110	0
菲军舰恶意撞击	103	5	121	0
食用油涨价	100	7	70	0
90 后暴打老人	99	3	94	0
官员财产公示	100	15	115	1
彭宇承认撞了南京老太	100	5	80	5
苹果封杀 360	100	53	58	3
三亚春节宰客	100	8	110	0
假和尚搂女子	100	8	107	0
国旗下讨伐教育制度	100	57	49	0
官员调研	100	21	106	0
六六叫板小三	105	25	105	0
中国教师收入全球几垫底	100	24	109	0
洗碗工剩菜被开除	106	10	129	0
ipad3	100	41	59	0
官二代求爱不成将少女毁容	105	0	169	1

由表中数据可以看出,每条微博有多个句子,其中分别包含观点句(正向情感句、负向情感句和中性情感句),以及非观点句。

本文面向观点句进行情感倾向判断,采取交叉验证方法,评测标准主要参考准确率(P)、召回率(R)和 F 值(F)。微平均以整个数据集为一个评价单元,计算其整体的评价指标,而不针对具体的话题;宏平均以每个话题为一个评价单元,分别计算各个话题中的评价指标,最后计算所有话题上各指标的平均值。

4.2 整体数据实验

基于表 5 所列的数据,将情感词极性程度分级并融合话题特征,分别采用基于规则的方法^[22]和 SVM 的方法进行情感倾向的判断,实验结果如表 6、图 2 所示。从实验结果可以看出,运用基于 SVM 的方法进行情感倾向判断,效果优于基于规则的方法。

表 6 整体数据实验结果

	基于 SVM 的方法		基于规则的方法	
	宏平均	微平均	宏平均	微平均
正向准确率	0.240759796	0.29787234	0.31869936	0.43651925
正向召回率	0.096273768	0.103194103	0.573608	0.683036
负向准确率	0.814156419	0.809126595	0.875023	0.91413
负向召回率	0.936504572	0.933748584	0.606418	0.641644
正向 F 值	0.10773485	0.153284672	0.38038895	0.53263707
负向 F 值	0.849172089	0.866982124	0.706535	0.754024
总体准确率	0.774006155	0.766198459	0.638463	0.641248

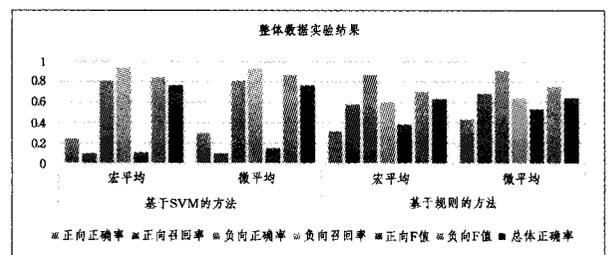


图 2 整体数据实验结果柱状图

4.3 对比实验

由于评测数据偏置性较强,数据包含的正向情感句和负向情感句分布极度不均匀,大部分句子为负向情感的句子,因此本文在整体数据集上进行实验后,选择了一些正向情感句与负向情感句分布相对均匀的话题,分别采用 4 种方法(考虑所有特征、去除分级情感词、去除话题特征、去除话题特征和分级情感词)进行情感分析实验,以验证本文提出的情感词极性程度分级算法和话题自适应算法的有效性。数据选取结果如表 7 所列。

表 7 对比实验所用数据

新闻话题	微博数	正向句	负向句	中性句
国旗下讨伐教育制度	100	57	49	0
iPad3	100	41	59	0
90 后当教授	104	110	13	0
六六叫板小三	105	25	105	0
奖状植入广告	100	22	53	1
苹果封杀 360	100	53	58	3
中国教师收入全球几垫底	100	24	109	0

对比实验采取表 7 所列的数据,选择基于 SVM 的情感倾向分类方法,对本文提出的情感词极性程度分级方法和话题自适应方法进行了对比实验,结果如表 8—表 11、图 3 所示。

通过对比实验结果可以看出:以微平均为例,当加上所有特征,分类器总体准确率为 0.5575。如果去掉分级情感词组成的特征,则准确率下降至 0.5082。如果去掉话题特征,则准确率下降至 0.5018。如果将分级情感词以及话题特征全部移去,则准确率下降至 0.4664。由此可见,本文提出的情感词极性程度分级和引入话题特征的方法,有助于提升情感倾向识别的准确性。

表 8 4 种总体正确率

方法	总体正确率	
	宏平均	微平均
所有特征	0.561711123	0.557522124
去除分级情感词	0.516413179	0.508217446
去除话题特征	0.511394951	0.501896334
去除话题特征和分级情感词	0.478258981	0.466498104

表 9 4 种方法的类别正确率

方法	类别正确率(宏平均)		类别正确率(微平均)	
	正向情感	负向情感	正向情感	负向情感
所有特征	0.54594812	0.63734847	0.49027237	0.59659090
去除分级情感词	0.51994509	0.62825500	0.41509434	0.56262042
去除话题特征	0.51490943	0.61348080	0.39574468	0.55272727
去除话题特征和分级情感词	0.44947853	0.59466457	0.32599118	0.52678571

表 10 4 种方法的类别召回率

方法	类别召回率(宏平均)		类别召回率(微平均)	
	正向情感	负向情感	正向情感	负向情感
所有特征	0.48784148	0.74738087	0.37951807	0.70627802
去除分级情感词	0.46229869	0.71586504	0.33132530	0.65470852
去除话题特征	0.41407031	0.73847924	0.28012048	0.68161435
去除话题特征和分级情感词	0.34436976	0.73801784	0.22289156	0.66143497

表 11 4 种方法的类别 F 值

方法	类别 F 值(宏平均)		类别 F 值(微平均)	
	正向情感	负向情感	正向情感	负向情感
所有特征	0.43982282	0.62999937	0.42784380	0.64681724
去除分级情感词	0.38165691	0.60139571	0.36850921	0.60518134
去除话题特征	0.35049984	0.60479076	0.32804232	0.61044176
去除话题特征和分级情感词	0.27974784	0.58591923	0.26475849	0.58648111

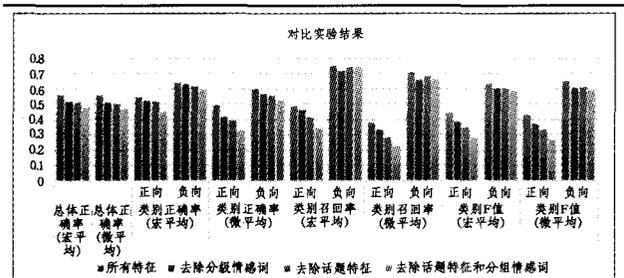


图 3 4 种方法对比实验整体结果

结束语 以微博为代表的网络平台的兴起,改变了传统的人与人之间相互交流的方式,微博也由此迅速成为深受广大网民热爱的社会化网络服务平台。本文针对中文微博文本的特殊性,对中文微博的情感分析技术进行了一定的探索,提出了对情感词极性程度分级和融合话题特征来提高中文微博情感分析精度的方法,并通过对比实验验证了所提方法的有效性。

面向中文微博的情感分析技术目前尚处于初级阶段,本文所提方法仍有很大的提升和扩展空间。今后拟在以下两个方面做进一步的研究工作:

- 1) 引入话题模型,研究面向热点话题的情感分析技术。
- 2) 在情感分析中引入微博的结构信息,以进一步提高微博情感分析的质量。

参考文献

- [1] Wiebe J. Learning Subjective Adjectives from Corpora[C]//Proceedings of AAAI, 2000
- [2] Riloff E, Wiebe J. Learning Extraction Patterns for Subjective Expressions[C]//Proceedings of EMNLP-2003. 2003;105-112
- [3] Turney P, Littman M L. Measuring Praise and Criticism: Inference of Semantic Orientation from Association[J]. ACM Transactions on Information Systems (TOIS), 2003, 21(4): 315-346
- [4] Kim S M, Hovy E. Automatic Detection of Opinion Bearing Words and Sentences[C]//Proceedings of IJCNLP-2005. 2005; 61-66
- [5] Kim S M, Hovy E. Identifying and Analyzing Judgment Opinions [C]// Proceedings of the Joint Human Language Technology/ North American Chapter of the ACL Conference (HLT-NAACL). 2006; 200-207
- [6] Zhu Y L, Min J, Zhou Y Q, et al. Semantic Orientation Computing Based on HowNet[J]. Journal of Chinese information processing, 2006, 20(1): 14-20
- [7] Andreevskaia A, Bergler S. Mining WordNet for a Fuzzy Sentiment; Sentiment Tag Extraction from WordNet Glosses[C]// Proceedings of the European Chapter of the Association for Computational Linguistics (EACL). 2006; 209-216
- [8] Su F, Markert K. Subjectivity Recognition on Word Senses via Semi-supervised Mincuts[C] // Proceedings of NAACL-2009. 2009; 1-9
- [9] Kamps J, Marx M, Mokken R J. Using WordNet to Measure Semantic Orientation of Adjectives [C] // Proceedings of the LREC. 2004
- [10] Esuli A, Sebastiani F. Determining the Semantic Orientation of Terms Through Gloss Analysis[C]// Proceedings of the ACM SIGIR Conference on Information and Knowledge Management (CIKM). 2005
- [11] Esuli A, Sebastiani F. Determining Term Subjectivity and Term Orientation for Opinion Mining[C]// Proceedings of the European Chapter of the Association for Computational Linguistics (EACL). 2006; 193-200
- [12] Yi J, Nasukawa T, Bunescu R. Sentiment Analyzer; Extracting Sentiments about a Given Topic Using Natural Language Processing Techniques[C]// Proceedings of the IEEE International Conference on Data Mining (ICDM). 2003
- [13] Hu M, Liu B. Mining Opinion Features in Customer Reviews [C] // Proceedings of AAAI-2004. 2004; 755-760
- [14] Ni M S, Lin H F. Mining Product Reviews Based on Association Rule and Polar Analysis [C] // Proceedings of the NCIRCS-2007. 2007; 628-634
- [15] Blei D M, Ng A Y, Jordan M I. Latent Dirichlet Allocation[J]. Journal of Machine Learning Research, 2003(3): 993-1022
- [16] Blei D M, Ng A Y, Jordan M I. Correlated Topic Models[C]// Advances in NIPS. 2006; 147-154
- [17] Titov I, McDonald R. Modeling Online Reviews with Multi-grain Topic Models[C]// Proceedings of WWW-2008. 2008; 111-120

维函数来观察 CR_{dn} 的自适应情况,其中 f_6 代表自变量相关 (non-separable) 问题, f_9 代表自变量独立 (separable) 问题。从图中可以看出,针对不同问题, CR_{dn} 在搜索的不同阶段呈现不同的变化趋势,有效地对参数起到了自适应调整的作用。

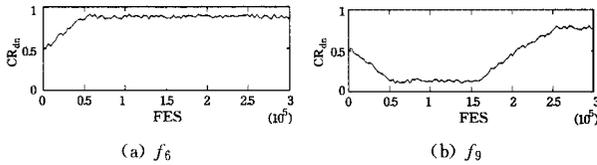


图6 CR_{dn} 的变化趋势图

另外,由 3.2.2 节可知,为了使实际产生的 $F_{i,G}$ 值尽可能落在其限定范围内以减少截断次数,根据柯西分布的特点对 F_{dn} 的取值范围设置了调整量 θ_r ($\theta=1, 2, 3, \dots, r=0.05$)。 θ 取值对 $F_{i,G}$ 值超限概率的影响如表 6 所列。显而易见,随着 θ 值的增大, F_{dn} 的取值范围越来越小,同时 $F_{i,G}$ 值的超限概率越来越低。由于 $F_{i,G}$ 值的超限概率与最大进化代数 G_{max} 无关,仅与 F_{dn} 的初始值和 θ 的取值相关,综合考虑,在设定 F_{dn} 初始值 0.7 的条件下,本文选取 $\theta=2$,以平衡 F_{dn} 的取值范围和超限概率的关系。

表 6 θ 的取值与 $F_{i,G}$ 值超限概率的关系

θ 值	超限概率
$\theta=1$	25%
$\theta=2$	15%
$\theta=3$	10%
$\theta=4$	8%

结束语 DE 算法在处理复杂函数最优化时存在后期收敛精度不高、速度较慢以及对参数设置敏感的问题,为此本文提出了一种新的变异策略 DE/current-to-dnbest/1 和新的参数自适应策略,对标准 DE 算法进行了综合改进,构成一种新型的基于动态自适应策略的 dn -DADE 算法。为验证 dn -DADE 算法的整体先进性和高效性,选取多种现有先进 DE 改进算法在 14 个 Benchmark 测试函数上进行对比实验,结果显示该算法具有更好的求解精度、收敛速度和稳定性。同时,通过将本文提出的 dn -DADE 算法与其本身的变种算法进行比较,进一步说明了改进的有效性。

参考文献

[1] Storn R, Price K. Differential Evolution-A simple and efficient heuristic for global optimization over continuous spaces [J]. Journal of Global Optimization, 1997, 11(4): 341-359

[2] Price K, Storn R. Differential Evolution-A practical approach to global optimization [M]. Berlin, Germany: Springer Verlag, 2006: 133-152

[3] Das S, Abraham A, Konar A. Automatic clustering using an improved differential evolution algorithm [J]. IEEE Transaction on Systems, Man and Cybernetics, 2008, 38(1): 218-236

[4] 韩敏, 王明慧, 范剑超. 基于改进差分进化算法的在线轨迹优化 [J]. 控制与决策, 2012, 27(2): 247-251

[5] Das S, Abraham A. Differential evolution using a neighborhood-based mutation operator [J]. IEEE Trans on Evolutionary Computation. 2009, 13(3): 526-553

[6] 袁俊刚, 孙治国, 曲广吉. 差分进化算法的数值模拟研究 [J]. 系统仿真学报, 2007, 19(20): 4646-4648

[7] Qin A K, Huang V L, Suganthan P N. Differential evolution algorithm with strategy adaptation for global numerical optimization [J]. IEEE Transaction on Evolutionary Computation, 2009, 13(2): 398-417

[8] Brest J, Greiner S, Boskovic B. Self-adapting control parameters in differential evolution: A comparative study on numerical benchmark problems [J]. IEEE Transactions on Evolutionary Computation, 2006, 10(6): 646-657

[9] Mallipeddi R, Suganthana P, Pan Q. Differential evolution algorithm with ensemble of parameters and mutation strategies [J]. IEEE Transactions on Evolutionary Computation, 2011, 11: 1679-1696

[10] Salman A, Engelbrecht A P, Omran M G H. Empirical analysis of self-adaptive differential evolution [J]. European Journal of Operational Research, 2007, 183(2): 785-804

[11] Zhang J, Sanderson A C. JADE: Adaptive differential evolution with optional external archive [J]. IEEE Transaction on Evolutionary Computation, 2009, 13(5): 945-958

[12] Mezura-Montes E, Velázquez-Reyes J, Coello C A C. A comparative study of differential evolution variants for global optimization [C] // Proceedings of Genetic Evolutionary Computation Conference (GECCO-2006). 2006: 485-492

[13] Suganthan P N, Hansen N, Liang J J, et al. Problem definitions and evaluation criteria for the CEC 2005 special session on real-parameter optimization [R]. Nanyang Technological University, Singapore, 2005

[14] Noman N, Iba H. Accelerating differential evolution using an adaptive local search [J]. IEEE Transaction on Evolutionary Computation, 2008, 12(1): 107-125

(上接第 235 页)

[18] Pang B, Lee L, Vaithyanathan S. Thumbs Up? Sentiment Classification Using Machine Learning Techniques [C] // Proceedings of EMNLP-2002. 2002: 79-86

[19] Pang Bo, Lee L. Opinion mining and sentiment analysis [J]. Foundations and Trends in Information Retrieval, 2008, 2(1/2): 1-135

[20] Tanaka Y, Takamura H, Okumura M. Extraction and classification of facemarks [C] // Proceedings of the 10th international conference on Intelligent user interfaces (IUI). New York, NY,

USA; ACM, 2005: 28-34

[21] Aoki S, Uchida O. A method for automatically generating the emotional vectors of emoticons using weblog articles [C] // Proceedings of the 10th WSEAS international conference on Applied computer and applied computational science (ACACOS). Venice, Italy: WSEAS Press, 2011: 132-136

[22] Paltoglou G, Thelwall M. Twitter, MySpace, Digg: Unsupervised Sentiment Analysis in Social Media [J]. ACM Transactions on Intelligent Systems and Technology (TIST), 2012, 3(4): 1-19