

基于级联模型的中文情感要素抽取

王亚坤, 黄河燕, 冯冲, 刘全超

(北京理工大学计算机学院北京市海量语言信息处理与云计算应用工程技术研究中心, 北京 100081)

摘要: 随着社交媒体的发展及成熟, 每天在互联网环境中都会产生大量的用户评论信息. 抽取评价短语、评价对象和观点持有者等情感要素, 已经成为了中文观点挖掘和情感分析的重要先决任务. 针对中文情感要素抽取任务, 本文提出了一个统计和规则相结合的级联模型, 主要贡献包括: (1) 针对汽车领域评论信息, 构建情感要素标注语料库和相关词典; (2) 对于以往研究较少关注的中文评价短语, 本文详细分析阐述其定义和分类; (3) 结合统计和规则, 分别针对评价短语和情感要素提出级联抽取策略. 实验结果充分证明了该级联模型的有效性, 相比较于其它基于规则的情感要素抽取算法有效提升了召回率, 同时为后续社交媒体情感分析任务提供了有力的支持.

关键词: 信息抽取; 情感要素; 评价短语; 评价对象; 观点持有者

中图分类号: TP391.1 **文献标识码:** A **文章编号:** 0372-2112 (2016)10-2459-07

电子学报 URL: <http://www.ejournal.org.cn>

DOI: 10.3969/j.issn.0372-2112.2016.10.025

Chinese Evaluation Element Extraction Based on Cascaded Model

WANG Ya-shen, HUANG He-yan, FENG-Chong, LIU Quan-chao

(Beijing Engineering Research Center of High Volume Language Information Processing and Cloud Computing Applications,
School of Computer, Beijing Institute of Technology, Beijing 100081, China)

Abstract: With the development of social media, massive reviews are generated by users every day. The extraction of evaluation elements, including evaluation phrase, comment target and opinion holder, is an important pre-task of Chinese opinion mining and sentiment analysis. This paper proposes an efficient method for extracting Chinese evaluation elements based on cascaded model and mainly makes three contributions: (i) to implement and evaluate the method, we construct an original annotated corpus for Chinese evaluation elements of automobile; (ii) we provide specific definition and classification of Chinese evaluation phrase; (iii) combining statistic method and rule-based method, we present cascaded strategy for extraction of evaluation phrase and evaluation elements, respectively. According to the experiment results, the proposed method performs well, and effectively improve the recall compared with other rule-based algorithm. Meanwhile it contributes greatly to our subsequent tasks, such as sentiment analysis of social media.

Key words: information extraction; evaluation element; evaluation phrase; comment target; opinion holder

1 引言

情感分析 (Sentiment Analysis) 和观点挖掘 (Opinion Mining)^[1,2] 旨在根据文本的话题或者情感极性来判断产生该内容的用户的观点和态度; 而情感要素 (Evaluation Element) 抽取则是其一项重要的先决任务, 不仅直接决定了后续任务的效果, 而且在实际应用中有着巨大的需求. 借鉴文献[3]中提出的“评价表达式 (Appraisal Expression)”概念, 本文所研究的中文“情感要

素”包括评价短语 (Evaluation Phrase, EP)、评价对象 (Comment Target, CT) 和观点持有者 (Opinion Holder) 等三部分 (三元组如下所示), 并设计级联模型完成对上述三者的抽取.

情感要素 = 〈评价短语, 评价对象, 观点持有者〉

目前学界尚无对“评价短语”的权威定义, 文献[4]曾将其粗略定义为“连续出现的一组评价词语”, 文献[5]和文献[6]分别提出过类似的概念“Appraisal Groups”和“Evaluative Expression”. 在上述研究基础之

收稿日期: 2015-02-11; 修回日期: 2015-06-26; 责任编辑: 马兰英

基金项目: 国家重点基础研究发展计划 (973 计划) 资助项目 (No. 2013CB329605, No. 2013CB329303); 国家自然科学基金 (No. 61132009, No. 61201351)

上,本文将“评价短语”定义为:针对某特定评价对象,表达一种观点态度和情感倾向的连续的词语组合.评价短语富含情感信息,能够提供有价值的特征以服务于后续的情感分析任务,而且可以作为结果直接提供给用户,为用户展示丰富的“全景式”信息,进而帮助用户全面理解相关产品(服务)或者事件.

协同利用基于统计抽取策略和基于规则抽取策略的各自优势,并综合考虑语料标注的难易程度,我们将“评价短语”分成“简单结构评价短语”和“复杂结构评价短语”两类,并实施“先易后难”的级联抽取策略:先抽取“简单结构评价短语”,在此基础上抽取“复杂结构评价短语”,而最终的“评价短语”由这两部分抽取结果共同组成.其中,“简单结构评价短语”主要是指程度副词和情感词语(主要是形容词和名词)的词语组合(包含使用连词或者顿号连接的情况),该类评价短语一般结构简单而且在文本中出现的位置比较固定(主要集中在定语、状语和补语等位置),例如“非常方便”和“及其无聊”等.

先前大量的相关工作只关注形如“简单结构评价短语”的短语甚至只关注单个情感词语,而很少关注结构复杂的短语^[3,5,7].但是此类结构复杂的评价短语往往富含情感信息(例如介词短语能够表达比较关系).本文主要研究括号短语、介词短语和副词短语等三类“复杂结构评价短语”,并分别制定了抽取规则.

真正对文本情感分析有帮助的不是单独的评价短语,而是评价短语和评价对象的组合^[3,7](即“评价搭配”^[8]);此外,增加观点持有者信息,有助于对进行观点归类 and 摘要^[9].因此,本文以抽取得到的评价短语为核心,采用级联模型抽取情感要素三元组:〈评价短语,评价对象,观点持有者〉.例如下述示例中,抽取的情感要素三元组为〈出色的,外形设计,专家〉.

专家认为,飞思拥有了出色的外形设计。

本研究所提出的级联模型主要关注产品评论信息中的情感要素抽取,但是也同样适用于其他类型文本的分析.该级联模型的流程如图1所示:(1)对输入文本进行分句、分词和词性标注等预处理;(2)基于条件随机场模型抽取简单结构评价短语;(3)基于简单结构评价短语抽取结果,应用规则抽取复杂结构评价短语,进而得到最终的评价短语抽取结果;(4)对于抽取得到的“评价短语”,基于规则定位和抽取其对应的评价对象(评价对象词典和情感词典等资源可以根据应用需求辅助使用^[10]),构成“评价搭配”〈评价短语,评价对象〉;(5)对于抽取得到的“评价搭配”,通过识别观点指示动词,完成对观点持有者的抽取,构成最终抽取结

果:情感要素三元组〈评价短语,评价对象,观点持有者〉.(图1中实线箭头指向为数据流动方向)

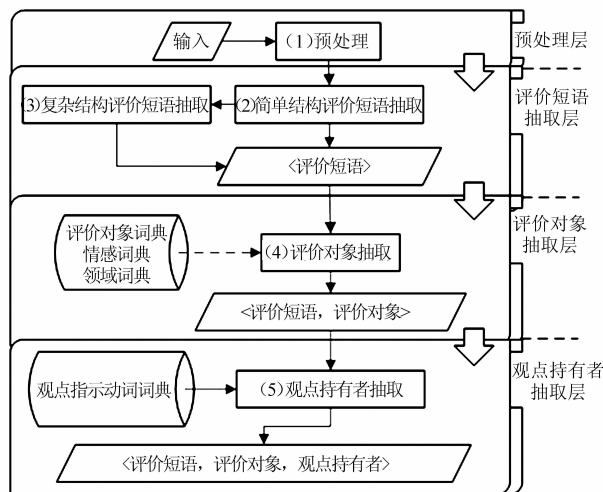


图1 中文情感要素抽取级联模型

2 相关工作

文献[5]认为情感分析的基本单元应该是评价短语而非单个词语;在这种思路的启发下,一系列算法和模型被提出^[6,11,12].但是上述研究所涉及的评价短语只属于本文提及的“简单结构评价短语”范畴而不涉及复杂结构.目前主流的评价对象抽取方法分为非监督学习方法和基于机器学习的有监督抽取方法^[7,12],围绕特征选择问题,条件随机场模型在评价对象抽取中广受青睐^[13].对于“评价搭配”抽取任务^[8],早期研究一般将这项任务分为两个步骤:首先获取情感句中的评价对象,然后评价对象附近窗口为 k 的范围内定位评价词语^[14].随后,部分研究者将对评价对象和评价词语的识别合并为一个独立的任务,提出了基于规则(或模板)的方法来识别评价搭配,其中句法分析结果被广泛用于构造规则^[4,7,15].虽然此类方法使得识别准确率得到提高,但是由于模板或者规则需要手工制定,召回率受限.

3 数据描述

虽然本文提出的级联模型可以被应用于处理不同种类的观点信息文本,但是本文研究重点关注产品评论信息,并应用该模型抽取汽车领域用户评论信息中的情感要素.目前尚无公开的中文汽车评论信息标注语料库,因此我们从2012年至2013年的腾讯汽车^①、网易汽车^②和凤凰汽车^③等汽车门户网站爬取和标注专家

① <http://auto.qq.com/>

② <http://auto.163.com/>

③ <http://auto.ifeng.com/>

测评文章和用户评论信息来构建中文汽车评论信息语料库. 该语料库目前已部分公开^①, 详情如表 1 所示(其中, #× 表示 × 的数量).

表 1 中文汽车情感要素语料库

#句子	2,872	#分句	14,606
#简单结构评价短语			6,450
#复杂结构评价短语			4,887
#评价短语			8,787
复杂结构评价短语平均长度			9.4

为了方便表述, 本文作如下概念定义:

分句 一个完整的句子(以句号、问号等终止符号作为结尾)被所包含的所有标点符号成多个“分句”. 本文以分句为基本单位抽取评价短语; 而构建情感要素三元组的时候, 则在整个句子中进行扫描和匹配相关要素.

词性序列 分词后, 一个或者连续几个(可包含连词)具有相同词性的词语构成“词性序列”. 例如, 分句“将内饰打造得更加典雅奢华”中下划线部分即为一个“形容词序列”.

窗口-R 分词后, 包含当前词语、当前词语前面 R 个词语和当前词语后面 R 个词语的词语序列. 本文使用中科院计算所汉语词性标注集, 本文章节 4 和章节 5 中规则表达式中的符号说明如表 2.

表 2 规则表达式中相关符号说明

符号	说明	符号	说明
EP	评价短语	CT	评价对象
p	介词	c	连词
n	名词序列	uls	助动词, “来说”“而言”等
a	形容词序列	due3	助动词, “得”
d	副词序列	*	任意词语
v	动词序列	f	方位词

此外, 为了提升抽取准确率以及充分支持本文模型的跨领域应用, 我们构造了评价对象词典和观点指示动词 (Opinion-Bearing Verb) 词典^②.

4 基于级联模型的情感要素抽取

本章节基于级联模型^[16,17], 抽取中文情感要素: 评价短语, 评价对象以及评价短语. “级联模型”的优势在于: (1) 各模块的输出相互作用关联, 模型最终输出结果是各模块输出结果的有机融合, 反映了各模块特征, 因此级联模型比较适用于元组抽取 (Tuple Extraction); (2) 一个模块的变化(如信息更新)会直接影响下一模块, 所以模块之间的数据流是“一体化”动态更新, 因此级联模型能够确保最终输出结果反映局部的信息更新.

4.1 基于条件随机场模型的简单结构评价短语抽取

条件随机场模型^[18]能够高效捕获输入文本的关联特征和识别序列边界, 并最大程度地降低标记偏执问

题, 被广泛应用于序列标注任务; 而简单结构评价短语具有构成简单、位置固定等特点, 因此本文将简单结构评价短语的抽取问题转换为序列标注问题, 并使用条件随机场模型完成该任务. 条件随机场模型所用特征模板主要包含 3 条特征(如表 3 所示), 其中 w_i 和 p_i 分别表示当前词语(第 i 个词语)及其词性.

表 3 用于识别简单结构评价短语的条件随机场模型的特征模板

特征	特征说明
w_i, p_i	窗口-3 及其词性
$i = -3, -2, -1, 0, 1, 2, 3$	
$w_{i-1}w_i, p_{i-1}p_i$	窗口-3 及其词性的二元组合特征
$i = -2, -1, 0, 1, 2, 3$	
$w_{i-1}w_iw_{i+1}, p_{i-1}p_ip_{i+1}$	窗口-3 及其词性的三元组合特征
$i = -2, -1, 0, 1, 2$	

4.2 基于规则的复杂结构评价短语抽取

本章节基于有限状态机 (Finite State Automaton, FSA) 思想, 设计三种复杂结构评价短语规则: 括号短语 (Parenthesis Phrase) 规则、介词短语 (Preposition Phrase) 规则和副词短语 (Adverb Phrase) 规则. 同时, 本文赋予这三种规则很强的可扩展性, 以便根据应用需求灵活改变规则.

4.2.1 括号短语规则

考虑到括号中的内容一般起到解释说明的作用, 并且往往包含有价值的评论信息, 我们抽取这部分内容作为评价短语. 括号内容紧邻所修饰内容(位于修饰内容的右侧), 所以其所对应的评价对象一般是其左侧的名词序列.

4.2.2 介词短语规则

简单结构介词短语往往表示处所或者状态, 一般不会表达情感倾向; 但是, 如果与其后的补语相结合构成复杂结构介词短语(特别是在评论信息中常见的表示“比较”意义的复杂结构介词短语)之后, 便可以传递一定情感信息.

对于每个分句, 我们从右至左搜索介词: 每搜索到一个介词, 分析其右侧文本是否匹配下述规则, 如果匹配, 则合并该介词连同其右侧符合规则的内容, 并抽取为评价短语; 继续向左进行搜索并重复上述过程, 直至分句搜索完毕. 本文共总结了 8 个基本的介词短语规则:

规则 1 $p + n + EP$

规则描述 如果介词右侧顺序出现名词序列和标注为 EP 的短语, 则合并该介词和这些词语成为一个短

^① <http://hlipca.org/index.php/2014-12-09-02-55-58/2014-12-09-02-56-24/49-chineseevaluationphrase>

^② <http://hlipca.org/index.php/2014-12-09-02-55-58/2014-12-09-02-56-24/49-chineseevaluationphrase>

语,并将词性重新标注为 EP.

规则示例 外观/n 上/f 将/d 会/v 比/p 传祺/nz 轿车/n 硬朗/EP 从该分句中抽取的复杂结构评价短语为“比传祺轿车硬朗”.

其它基本介词短语规则见表 4.

表 4 其它基本介词短语规则

序号	规则	序号	规则
2	p+n+d+v	6	p+n+f+EP
3	p+n+v	7	p+v+n
4	p+n+f+v	8	p+v+n+EP
5	p+n+f+d+v		

上述基本规则可以通过如下方法进行扩展:

(1) 上述规则中的名词序列可以被代词序列所替换(或者附加),规则依然成立.

(2) 上述规则可以通过加入形容词序列或者标注为 EP 的短语来构造更加复杂的规则. 例如,对于规则 3,在名词序列前加入标注为 EP 的短语,规则依然成立: $p+n+v \rightarrow p+EP+n+v$.

(3) 某些中文词语,例如“相比”“对比”等,能够表达“比较(对比)”关系,因此介词短语规则中的“介词(标记为 p)”可以替换成为这些词语,规则依然成立.

4.2.3 副词短语规则

实际应用中,副词可以修饰动词、形容词,甚至整个句子.其中,情态副词(例如“究竟”“简直”等)和程度副词(例如“非常”“很”等),往往引导富含情感信息的短语,所以能够指示观点持有者的态度.

本文主要关注上述副词做谓语和补语的情况.类似于上述介词短语的构造方式,我们通过对右至左搜索分句,判断所出现的每一个副词右侧的文本是否匹配相关副词短语规则.本文共总结了 6 个基本的副词短语规则:

规则 2 d+v+EP

规则描述 如果副词右侧顺序出现动词序列和标注为 EP 的短语,则合并该副词和这些词语成为一个短语,并将词性重新标注为 EP.

规则示例 内饰/nz 还/d 算/v 朴素大方/EP

从该分句中抽取的复杂结构评价短语为“还算朴素大方”.

其它基本副词短语规则见表 5,同样可以参照上一章节的扩展方法对基本副词短语规则进行扩展.

表 5 其它基本副词短语规则

序号	规则	序号	规则
2	d+EP	5	d+v
3	d+n	6	d+z+v+n
4	d+v+n		

4.3 基于规则的评价对象抽取

在前述“评价短语”抽取结果的基础上,本章节抽

取其所对应的“评价对象”,构成“评价搭配”(评价短语,评价对象).以每个被抽取的评价短语为中心,我们使用“评价对象构建规则”来定位和构建名词序列作为评价对象候选;然后使用“评价对象抽取规则”从这些候选中挑选出真正与该评价短语配对的对象作为最终的评价对象.此外,也可以根据应用需求引入评价对象词典筛选环节^[10].

评价对象构建规则 文献[3]曾选取距离评价对象最近的形容词作为其对应的评价词语,我们通过考察大规模产品服务类评论信息语料也发现:与某个评价短语配对的评价对象往往是其左侧(或者右侧)最近的名词序列.因此,对于某个评价短语,我们分别向左和向右扫描文本并构建距离其最近的名词序列,作为评价对象候选.例如“专家/d 认为/vo,/wd 飞思/nz 拥有/v 了/ule 出色的/EP 外形/n 设计/vn.”中,以评价短语“出色的”为中心,向左和向右分别定位和构建了名词序列“飞思”和“外形设计”作为评价对象候选:

专家认为, 飞思 拥有了 出色的 外形设计.

评价对象抽取规则 我们使用下述规则(表 6)来从评价对象候选中挑选出最终的评价对象.

表 6 其他评价对象抽取规则

序号	规则描述	规则举例
1	评价对象位于评价短语右侧并且紧邻.	EP + CT
2	评价对象位于评价短语左侧并且紧邻.	CT + EP
3	评价对象位于评价短语右侧,非紧邻.	CT + v + EP; CT + * + uis + EP; CT + * + c + EP; CT + * + f + EP; CT + * + v + * + ude3 + EP;

所以,上述示例中的两个评价对象候选中,只有“外形设计”符合规则,被保留.至此,我们得到了该示例中的“评价搭配”(出色的,外形设计).

专家认为, 飞思 拥有了 出色的 外形设计.

4.4 基于规则的观点持有者抽取

在前述“评价搭配”抽取结果的基础上,本章节抽取其所对应的“观点持有者”,构成最终抽取结果:“情感要素”三元组(评价短语,评价对象,观点持有者).通过对大量评论语料进行调研,我们发现:人名和机构名等命名实体经常出现在观点持有者的位置上,而且观点持有者往往与观点指示动词共现.因此,本文对于观点持有者的抽取策略是基于观点指示动词的位置的,而且我们更多关注的是连续分句之间的观点持有者是否发生变化.

首先,如果分句中出现观点指示动词,则说明观点持

有者可能发生变化,否则认为跟前一分句的观点持有者相同;然后,固定该观点指示动词并向前文进行文本扫描,定位和构建距离最近的命名实体(或名词序列),作为观点持有者候选;最后,我们利用相关规则来判断观点持有者候选是否正确(例如,如果介词“据”出现潜在观点持有者的左侧,则表示抽取成功).上文示例中,“认为”是观点指示动词,其左邻的名词序列“专家”即为观点持有者.至此,我们通过级联模型得到了该示例中的“情感要素”三元组〈出色的,外形设计,专家〉.

专家认为,飞思拥有了出色的外形设计。

5 实验和结果分析

本章节中,我们使用章节 3 所描述的语料来验证本文提出的级联模型的性能.采用 10-折交叉验证的方式分配训练集和测试集,记录 10 次实验结果的均值于相关图表中.实验中,我们使用 NLPiR 汉语分词系统 2014^①完成中文分词和词性标注任务,使用 CRF++ version 0.53^②完成针对简单结构评价短语抽取的条件随机场模型的训练和测试任务.

5.1 评价指标

本文使用准确率(Precision, P)、召回率(Recall, R)和 F-值(F-measure, F)作为评价指标对简单结构评价短语、复杂结构评价短语、评价短语和情感要素三元组等测评对象的抽取结果进行测评.

$$Precision = \left(\frac{N_3}{N_2} \right) * 100\% \quad (1)$$

$$Recall = \left(\frac{N_3}{N_1} \right) * 100\% \quad (2)$$

$$F\text{-measure} = \frac{(\beta^2 + 1) * P * R}{R + \beta^2 * P} \quad (3)$$

其中, N_1 表示测试集所含测评对象的个数, N_2 表示本文算法抽取的测评对象的个数, N_3 表示抽取结果中抽取正确的对象的个数, β 取值为 1.此外,考虑到短语覆盖文本范围较大的问题,在评价短语抽取测评中,我们引入下述三种不同的测评粒度^[11](其中后两种统称“松弛匹配”):

精确匹配(又称“严格匹配”) 只有抽取结果严格匹配标准结果,该抽取结果才被认为是正确的.

部分匹配 如果抽取结果包含标准结果,该抽取结果即可被认为是正确的.

范围部分匹配 如果抽取结果和标准结果有重合部分,则将重合部分所占比重加入到 N_3 .

5.2 实验结果

中文评价短语抽取实验的结果如表 7 所示.正如前文分析,简单结构评价短语的构成规律性强、边界

明确,因此即使是在“精确匹配”这种严苛的测评标准下,简单结构评价短语的抽取也拥有很高的准确率和召回率.

表 7 中文评价短语抽取结果

评测粒度	评测指标	简单结构	复杂结构	评价短语
精确匹配	准确率	0.8597	0.4562	0.3921
	召回率	0.8634	0.3657	0.2148
	F-值	0.8615	0.4060	0.2776
部分匹配	准确率	0.8953	0.7823	0.7603
	召回率	0.9570	0.7146	0.6908
	F-值	0.9252	0.7469	0.7238
范围部分匹配	准确率	0.8621	0.7436	0.6874
	召回率	0.8720	0.5943	0.4269
	F-值	0.8726	0.6606	0.5267

表 7 同时也反应出,复杂结构评价短语和评价短语在“精确匹配”测评策略下的性能并不理想,这是因为其复杂且多变的内部结构导致很难准确识别其所辖文本的范围和边界.所以,本文使用“部分匹配”策略作为“评价短语”的主要的测评标准,并且与相关工作文献[6]和文献[11]进行对比(如表 8 所示).

表 8 本文模型与其他算法对于“评价短语”抽取的实验结果对比

算法	测评策略	准确率	召回率	F-值
本文的级联模型	精确匹配	0.3921	0.2148	0.2776
	部分匹配	0.7603	0.6908	0.7238
	范围部分匹配	0.6874	0.4269	0.5267
文献[11]	精确匹配	0.2936	0.1674	0.2208
	部分匹配	0.7918	0.4260	0.5734
	范围部分匹配	0.6328	0.4034	0.5032
文献[6]	精确匹配	0.2174	0.1032	0.1392
	部分匹配	0.7091	0.4427	0.5451
	范围部分匹配	0.6422	0.3987	0.4920

本文提出的级联模型对情感要素抽取结果(采用“部分匹配”策略)如表 9 所示.实验结果显示,本文模型对于观点持有者抽取的准确率较高,而对于情感要素抽取的 F-值也突破了 70%.此外,文献[15]中针对“具有修饰关系的词对”的任务与本文情感要素抽取任务十分相似,因此我们复现了该工作.对比实验显示,作为以往基于规则抽取算法的典型代表,文献[15]虽然取得了较高的准确率,但是召回率存在欠缺;而本文模型将召回率提升了 20.12%,而且在时间消耗方面远优于文献[15].此外,考虑到本文算法所抽取的评价短语更加复杂,而且有助于产品(服务)信息的直观展示,因此本文模型在海量信息处理领域还是具有很大的应用价值.

① <http://ictclas.nlp.ir.org/newsdownloads? DocId = 389>

② <http://crfpp.googlecode.com/svn/trunk/doc/index.html>

表9 中文情感要素抽取结果及对比

评测对象	准确率	召回率	F-值
评价短语	76.03%	69.08%	72.38%
评价对象	75.91%	68.75%	72.15%
观点持有者	84.47%	76.50%	80.29%
情感要素	67.27%	74.28%	70.59%
文献[15]	78.39%	61.84%	69.14%

5.3 实验结果分析

本文方法的召回率和 F 值相比较于基线算法(包括传统基于规则的方法)有了一定提升,主要原因有如下几点:

(1)本研究针对“评论信息中情感要素抽取”,目的明确、应用性强,而且级联模型中相关规则的设计均基于对中文评论信息扎实的语法、句法分析(特别是各要素之间关系),所以,在平衡规则复杂度和计算复杂度前提下,本研究设计的规则比传统规则更有针对性、更加精确。

(2)传统基于规则的抽取算法,在匹配规则时往往采用“字符连续出现”的匹配模式,导致以往方法的召回率较低;而本研究在匹配规则时,基于有限状态机思想,采取“词性序列顺序出现”的匹配模式,进而有效提升召回率。

(3)借力于级联模型的“联动机制”,本研究所抽取的评价短语、评价对象和观点持有者分别处于级联模型中相连通的不同模块中(图1),因此本文模型能够对三者关系以及句子结构进行更加清晰的刻画。

(4)多策略平衡.通过减少特征种类和数量,实现“质量策略”和“速度策略”的平衡,进而保证系统整体效率;充分发挥“统计策略”和“规则策略”各自优势,实现统计和规则互补。

6 总结

本文着重对中文情感要素中的“评价短语”概念进行了详细的定义和阐述,并且构建了相关的语料库.面向海量中文信息处理需求,通过研究不同情感要素的语法和结构特征,本文设计了统计和规则相结合的级联模型来抽取用户评论语料中的评价短语、评价对象和观点持有者.实验结果充分证明了该级联模型的有效性,相比较于其它基于规则的情感要素抽取算法有效提升了召回率;此外,本研究相关内容已经在部署在实际应用中,并取得了良好的实践效果。

参考文献

[1] Pang B, Lee L. Opinion mining and sentiment analysis Foundations and trends in information retrieval[J]. Foundations & Trends in Information Retrieval, 2008, 2(1-2): 459-526.

[2] Liu B. Sentiment analysis and opinion mining[J]. Synthesis Lectures on Human Language Technologies, 2012, 5(1): 1-167.

[3] Bloom K, Garg N, Argamon S. Extracting appraisal expressions[A]. Proceedings of Human Language Technologies: 2007 Annual Conference of the North American Chapter of the Association for Computational Linguistics[C]. New York: ACL Press, 2007. 308-315.

[4] 赵妍妍, 秦兵, 车万翔, 刘挺. 基于句法路径的情感评价单元识别[J]. 软件学报, 2011, 22(5): 887-898.
ZHAO Yan-yan, QIN Bing, CHE Wan-xiang, LIU Ting. Appraisal expression recognition based on syntactic path[J]. Journal of Software, 2011, 22(5): 887-898. (in Chinese)

[5] Whitelaw C, Garg N, Argamon S. Using appraisal groups for sentiment analysis[A]. Proceedings of 14th ACM International Conference on Information and Knowledge Management[C]. New York, USA: ACM Press, 2005. 625-631.

[6] Nakagawa T, Kawada T, Inui K, Kurohashi S. Extracting subjective and objective evaluative expressions from the Web[A]. Proceedings of 2nd International Symposium on Universal Communication[C]. Osaka, Japan: IEEE Press, 2008. 251-258.

[7] Popescu A M, Etzioni O. Extracting product features and opinions from reviews[A]. Proceedings of Human Language Technology: 2005 Conference on Empirical Methods in Natural Language Processing[C]. Vancouver, Canada: ACL Press, 2005. 339-346.

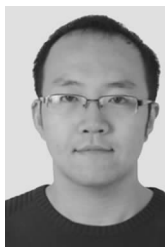
[8] 赵妍妍, 秦兵, 刘挺. 文本情感分析[J]. 软件学报, 2010, 21(8): 1834-1848.
ZHAO Yan-yan, QIN Bing, LIU Ting. Sentiment analysis[J]. Journal of Software, 2010, 21(8): 1834-1848. (in Chinese)

[9] 宋锐, 洪莉, 林鸿飞. 基于 ChunkCRF 的观点持有者识别及其在观点摘要中的应用[J]. 小型微型计算机系统, 2009, 30(7): 1462-1466.
SONG Rui, HONG Li, LIN Hong-fei. Chunk-CRF-based opinion holder identification and application to opinion summarization[J]. Journal of Chinese Computer Systems, 2009, 30(7): 1462-1466. (in Chinese)

[10] Nakagawa T, Inui K, Kurohashi S. Dependency tree-based sentiment classification using CRFs with hidden variables[A]. Proceedings of Human Language Technologies: 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics[C]. Los Angeles, USA: ACL Press, 2010. 786-794.

- [11] Wang Y, Kazama J, Kawada T, Torisawa K. Chinese evaluative information analysis[A]. Proceedings of 24th International Conference on Computational Linguistics [C]. Mumbai, India: ACM Press, 2012. 2773 – 2788.
- [12] 侯敏, 滕永林, 陈毓麒. 评价短语的倾向性分析研究[J]. 中文信息学报, 2013, 27(6): 103 – 109.
HOU Min, TENG Yong-Lin, CHEN Yu-qi. Research on orientation analysis of opinion phrases[J]. Journal of Chinese Information Processing, 2013, 27(6): 103 – 109. (in Chinese)
- [13] 王荣洋, 鞠久朋, 李寿山, 周国栋. 基于 CRFs 的评价对象抽取特征研究[J]. 中文信息学报, 2012, 26(2): 56 – 61.
WANG Rong-yang, JU Jiu-ming, LI Shou-shan, ZHOU Guo-dong. Feature engineering for CRFs based opinion target extraction[J]. Journal of Chinese Information Processing, 2012, 26(2): 56 – 61. (in Chinese)
- [14] Hu M Q, Liu B. Mining and summarizing customer reviews[A]. Proceedings of 2004 ACM SIGKDD International Conference on Knowledge Discovery & Data Mining [C]. New York: ACM Press, 2004. 168 – 177.
- [15] 姚天昉, 等. 一个用于汉语汽车评论的意见挖掘系统[A]. 中国中文信息学会. 中文信息处理前沿进展——中国中文信息学会二十五周年学术会议论文集[C]. 北京: 中国中文信息学会, 2006. 260 – 281.
YAO Tian-fang, et al. An opinion mining system for chinese automobile reviews[A]. Proceedings of the 25th Annual Conference of CIPS [C]. Beijing: Chinese Information Processing Society of China, 2006. 260 – 281. (in Chinese)
- [16] 赵巍, 等. 连续字符识别的级联 HMM 训练算法[J]. 计算机学报, 2007, 30(12): 2142 – 2150.
ZHAO Wei, et al. Cascaded HMM training algorithm for continuous character recognition [J]. Chinese Journal of Computers, 2007, 30(12): 2142 – 2150. (in Chinese)
- [17] 李本阳, 等. 基于单层标注级联模型的篇章情感倾向分析[J]. 中文信息学报, 2012, 26(4): 3 – 8 + 20.
LI Ben-yang, et al. Single-label cascaded model for document sentiment analysis [J]. Journal of Chinese Information Processing, 2012, 26(4): 3 – 8 + 20. (in Chinese)
- [18] Lafferty J, McCallum A, Pereira F. Conditional random fields: Probabilistic models for segmenting and labeling sequence data[A]. Proceedings of 18th International Conference on Machine Learning [C]. Williamstown, MA, USA: ACM Press, 2001. 282 – 289.

作者简介



王亚坤 男, 1989 年出生, 北京理工大学计算机科学与技术专业博士研究生, 主要研究领域为社交网络分析和信息检索.
E-mail: yswang@bit.edu.cn



黄河燕(通讯作者) 女, 1963 年出生, 北京理工大学计算机学院教授、博士生导师, 主要研究领域为语言信息智能处理、社交网络、文本大数据分析处理及云计算.
E-mail: hhy63@bit.edu.cn