

一种基于情感句模的文本情感分类方法

陈涛^{1,2}, 徐睿峰^{*1}, 吴明芬², 刘滨¹

(1.哈尔滨工业大学深圳研究生院,广东省深圳市 518055; 2.五邑大学计算机学院,广东省江门市 529020)

摘要: 考虑到同类型的情感句往往具有相同或者相似的句法和语义表达模式, 本文提出了一种基于情感句模的文本情感自动分类方法。首先, 将情感表达相关句模人工分为3大类105个二级分类; 然后, 设计了一种利用依存特征、句法特征和同义词特征的句模获取方法, 从标注情感句中半自动地获取情感句模。最后, 通过对输入句进行情感句模分类实现文本情感分类。在NLP&CC2013中文微博情绪分类评测语料及ReLabEAC1.0文学语料的实验结果显示本文提出的分类方法准确率显著高于基于词特征支持向量机分类器。

关键词: 情感句模; 情感分类; 句法特征; 依存特征

中图分类号: TP391

文献标识码: A

A sentiment classification approach based on sentiment sentence framework

Chen tao^{1,2}, Xu ruifeng^{*1}, Wu mingfen², Liu bin¹

(1. Harbin Institute of Technology Shenzhen Graduate School, Shenzhen, 518055, China;

2. School of Computer Science Wuyi University, Guangdong, 529020, China)

Abstract: Considering that opinionated sentences always have the same or similar syntax and semantic expression frameworks, this paper proposes a sentiment analysis approach based on sentiment sentence framework. Firstly, we divided sentiment sentence frameworks into three categories and 105 subcategories. A sentence framework extraction method is designed to semi-automatically extract sentiment sentence frameworks from annotated sentiment sentences using dependency features, syntactic features and synonym features. The polarity of input sentence is determined through the classification of its sentiment sentence frameworks. The evaluations on NLP&CC 2013 micro-blog emotion analysis corpus and RelabEAC 1.0 literature corpus show that our proposed sentiment classification approach achieves better precision performance compared to word-based support vector machine classifiers.

Key words: sentence framework; sentiment classification; syntactic feature; dependency feature

1 引言

情感(Sentiment)是一种复杂的生理和心理现象, 包括情绪(Emotion)、感觉(Feeling)和心情(Mood)等。它是人类智能的重要特征, 是人类生活的重要内容。随着 Web 2.0、社交网络的兴起, 越来越多的人通过网络分享自己的观点、体验和心情, 包含有情感的文本也越来越多。对文本中蕴含的情感和情绪进行自动分析正在成为新的研究热点。这些分析技术的研究对于大数据行业挖掘文本潜在的情感表达, 发现用户兴趣与需求, 提高服务质量等应用领域都具有重要的意义。

目前文本情感分类的典型方法包括: (1) 基于关键词的方法。如 Turney^[1]等使用词之间的点式互信息 (Point-wise mutual information)和 SO (Semantic orientation)对评论进行非监督的分类; Kamps^[2]等利用 WordNet 记录的信息来分析形容词的极性; 朱嫣岚^[3]等基于 HowNet 分析词汇的倾向性进行句子倾向性分析。这类方法对分析显式的、含有情感词的文本比较有

基金项目: 高等院校博士学科点专项基金(20122302120070); 国家自然科学基金(11271040); 深圳市基础研究计划(JCYJ20120613152557576; 深圳市国际合作计划(GJHZ20120613110641217); 广东省自然科学基金(S2011010003681); 模式识别国家重点实验室开放课题基金。

作者简介: 陈涛(1981-), 男, 博士研究生, 主要研究方向为自然语言处理; 通讯作者: 徐睿峰(1973-), 男, 副教授, 主要研究方向为自然语言处理; 吴明芬(1964-), 女, 教授, 主要研究方向为模糊集、模糊集理论及其在智能信息处理中的应用; 刘滨(1981-), 男, 助理教授, 主要研究方向为生物信息学, 自然语言处理;

效。(2) 基于规则或常识知识库的方法,如 Hu^[4]等使用关联规则挖掘客户的主观评论;姚天昉^[5]等使用句法规则的方法对汽车评论中的情感倾向进行挖掘;刘鸿宇^[6]等基于句法树中的路径对评价对象进行抽取;任巨伟^[7]等在陈健美^[8]等人的情感常识表示框架基础上构建了二元结构的情感常识库,进行文本情感分析和倾向性分析。这类方法具有一定的分析隐含情感和领域相关情感文本的能力。(3) 基于统计和机器学习的方法。如 Pang Bo^[9]应用朴素贝叶斯、最大熵、支持向量机 SVM(Support Vector Machine)等分类器对电影评论进行分类;谷学静^[10]等利用隐马尔科夫模型 HMM 对情感进行建模;王根^[11]等采用多重冗余标记的条件随机场分类器,通过求联合解码最优,减少了单分类的错误传递;Li^[12]等利用多分类器融合的方法改进单一分类器的效果;李寿山^[13]等采用了基于 Stacking 组合分类方法对分类器进行情感倾向分析。这类方法得到了较多的应用。

考虑到同类别的情感句往往有相同或者相似的句法和语义表达模式,本文提出一种基于情感句模的文本情感分类方法。首先,从《现代汉语基本句模》^[14]中选取与情感表达相关的三大类句模,并进行人工补充获得 105 个二级分类句模。而后,利用情感标注语料,对基础情感句模无法覆盖的情感句进行分词、句法分析和依存关系分析,从中找出句子的核心谓词和与其直接关联的句子主干词以及对句子情感有直接影响的其他词,通过半自动的方法获取情感表达句模,从而建立一个情感句模库。在情感分类任务中,将情感句分类问题转换为最相似句模分类和排序,从而实现情感分类。在 NLP&CC 2013 中文微博情绪分类评测数据集及 ReLabEAC1.0 文学语料进行的评估实验显示,本文提出的方法对语料多标签情感分类准确率分别达到 43%和 60%,明显优于基于词语特征的 SVM 分类器方法,显示本文提出的基于情感句模的方法可以有效地提高文本情感分类性能。

2 基于情感句模的文本情感分类方法

通过对大量情感句的表达方式进行分析和总结,可以发现句子的主要语义往往通过句子的主干来表达,很多时候具有相同或者类似主干的句子所表达的情感也相同。例如:表示喜爱情感的两个句子“我喜欢你”与“我爱自然语言处理”具有共同的句子主干“情感的持有者+表示喜爱的词语+情感的对象”。为此,本文引入朱晓亚^[15]提出的汉语句模的概念进行描述。这里,句模定义为动核结构生成句子时与句型结合在一起的语义成分的配置模式,是根据句子语义平面的特征分出来的类别。上述例子中的句子主干可以用句模“<感事><喜爱词类><向事>”来描述,其中“感事”表示情感的主体,“向事”表示情感施加的对象。每一类句模包含对应的词类。在利用句模对情感句表达方式进行分析和总结的基础上,本文提出以下假设:

假设 1: 情感句模能够表达句子的主要语义。

假设 2: 如果句子 S 能用情感句模 M 表示(即与该句模匹配),则 S 与 M 表达的情感分类相同。

基于情感句模的文本情感分类方法的基本设计思想是:将待分类句子与情感句模进行匹配,找出匹配程度最高的句模,句模所属的情感类别即为此句子的情感分类。

2.1 情感分类及句模库构建

考虑到情感表达的灵活多样,因此需要对情感表达句和对应的情感句模进行相对精细的区分。鲁川^[14]等人在论文《现代汉语基本句模》中将常见的汉语句模分为 26 个大类,122 个二级分类。本文首先从中选择出与情感表达有着密切关系的包括“态度、感受、思想”3 个大类和 14 个二级分类,其中“态度”大类分为“热情类、细心类”等,“感受”大类分为“感知类、感觉类”等,“思想”大类分为“希望类、愿意类”等。从这些分类对应的句模库中抽取了 14 个句模和 14 个对应的词类构成基本情感句模库。考虑到^[14]存在对情感表达句式覆盖率有所不足的问题,结合对大规模情感语料库的观察和分析,本文对上述二级分类进行了扩展,最终得到对应于“态度”大类的 41 个二级分类、对应于“感受”大类的 48

个二级分类、对应于“思想”大类的16个二级分类。详细的分类列表在附录中给出。

由于基本情感句模库不能够覆盖所有对应类别情感句的表达方式。另一方面，本文新扩展出的二级分类无法从现有的《现代汉语基本句模》资源中获得对应的句模。为此，本文提出了一种情感句模的半自动获取方法，基本过程如下：

(1) 从情感语料库中抽取情感句，利用基本情感句模库进行匹配。对能匹配的句子作为对应句模的实例存储。对不能匹配的句子，人工标注其对应的情感二级分类，并继续处理。

(2) 对这些句子进行分词、句法分析和依存关系分析。从中找出句子的核心谓词和与其直接关联的句子主干词，以及对句子情感有直接影响的其他词（称为：附属词）。

(3) 借助《同义词词林》，查找核心谓词和附属词所在的同义词词类，用同义词词类名代替该核心谓词和附属词。如果这些词语不属于任何同义词词类，则创建新的词类。这里，为区别词和词类，将词类名用尖括号括起来。

(4) 参考《现代汉语基本句模》^[14]中定义的语义角色，将句子主干词抽象成语义角色。这里，为区别词和语义角色，将语义角色用尖括号括起来。

(5) 将制作好的句模存入句模库。

下面以例句1“我爱自然语言处理。”说明“喜爱类”中句模的构建过程：

(1) 生成句子的分词结果“我/爱/自然/语言/处理/。”，以及对应的句法分析和依存关系分析结果如图1和图2所示。结合句法分析和依存分析结果，可知例句1的核心谓词是“爱”，与其直接关联的句子主干词分别是“我”和“处理”，而“自然”和“语言”则不被视为句子主干。



图1 例句1 依存关系树

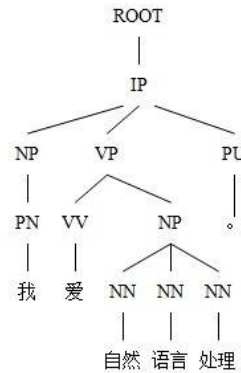


图2 例句1 句法树

表1 典型情感句模及其所属情感分类

句模	一级类	二级类
<施事>站在<向事><站在一起词类>	态度类	支持类
<当事><对_当事介词词类><向事><重视词类>	态度类	重视类
<施事><共事介词词类><共事><恋爱词类>	态度类	恋爱类
<感事><以_方式介词词类><客事>[为]<光荣词类>	感受类	为荣类
<感事><对_当事介词词类><客事><不知道词类>	感受类	不知道类
<感事><知道词类 感知词类 推测估计词类 感到觉得词类><向事><火_词类 吸引力词类 好_词类>	感受类	喜爱类
<施事><终点方向词类><参观词类>	思想类	希望类
<施事><使让词类><致事><否定副词词类><怀疑词类>	思想类	信任类
<施事><在_范围介词词类><范围>上<另眼相看词类><客事>	思想类	另眼相看类

(2) 在句模库中查找“爱”是否属于某个已知词类，如果匹配则使用该词类名代替“爱”；如无法匹配，则在《同义词词林》中查找“爱”的同义词，在句模库中创建“喜爱词类”，并将“爱”和它的同义词添加到该词类中。而后用“<喜爱词类>”替代“爱”，此时例句1的主干为：“我 <

喜爱词类> 处理”。

(3) 参考《现代汉语基本句模》中定义的语义角色，通过将“我”抽象为“感事”，将“处理”抽象为“向事”，则从例句 1 获得句模“<感事><喜爱词类><向事>”。

(4) 将新句模存入情感句模库。

按照上述步骤，我们共从约 3500 个情感句中获得了 413 个情感句模，表 1 列举了一些典型情感句模及其所属的情感分类，其中中括号里面的词是可以省略的词。

2.2 基于情感句模的情感分类算法

2.2.1 分类特征选择

利用情感分类句模库，可以将句子的情感分类转换为对情感句模的分类问题，只要找出与待分类句子匹配程度最高的句模，即可使用句模的情感分类作为此句子的情感分类。本文设计和选择了 4 种句模匹配特征用于描述句模匹配的相似性：

f_1 : 是否匹配核心词，匹配则 $f_1=1$ ，不匹配则为 0。

f_2 : 匹配附属词个数。

f_3 : 匹配依存关系树中与核心词连接的边的个数。

f_4 : 匹配句法树中从根到叶子节点的路径（或称为句法树分支）的个数。

随后，使用下面的线性分类器模型结合上述 4 种匹配特征进行分类：

$$y = \text{sgn}(\mathbf{w} \cdot \mathbf{f} + \mathbf{b}) = \begin{cases} 1, & \text{if } \mathbf{w} \cdot \mathbf{f} + \mathbf{b} > 0 \\ -1, & \text{else} \end{cases} \quad (1)$$

其中，向量 $\mathbf{w}=[w_1, w_2, w_3, w_4]$ ，是对向量 $\mathbf{f}=[f_1, f_2, f_3, f_4]$ 的权重向量。 b 为阈值， $y=1$ 表示匹配成功，该句模所属情感分类记为待分类句子的分类， $y=-1$ 表示不匹配。

例如：例句 2“我 喜爱 信息 检索。”的核心词为“喜爱”，查找句模库得到可能的类别为喜爱类。对应喜爱类中的句模 1“<感事><喜爱词类><向事>”，其中“<喜爱词类>”包括核心词“喜爱”，所以 $f_1=1$ ；句模 1 中没有附属词匹配，所以 $f_2=0$ ；句模 1 的依存关系树与句法树分别如图 3 至 4 所示。依存关系树与核心词连接的边“SBV”和“VOB”，两条边都匹配，所以 $f_3=2$ ；句法树中有 4 个分支匹配：分别是“ROOT→IP→NP→PN”、“ROOT→IP→VP→VV”、“ROOT→IP→VP→NP→NN”和“ROOT→IP→PU”，所以 $f_4=4$ 。因此 $\mathbf{f}=[1,0,2,4]$ 。类似可以生成其他句模对应的 \mathbf{f} 向量。

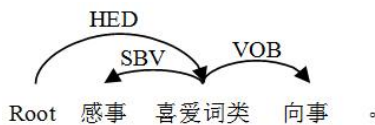


图 3 句模 1 依存树

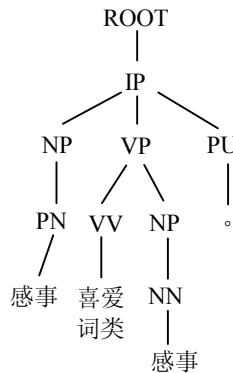


图 4 句模 1 句法树

利用线性分类器分类结果可以到句模 1 为最近似句模分类结果，对应的例句 2 的情感分类结果可以由句模 1 的分类“喜爱”获得。

本文设计的基于情感句模的情感分类算法分类成功的标准是至少找到一个匹配得分超过阈值的句模。每个情感分类中所有句模的最高得分为该分类的得分，按总得分由高到低对情感分类进行排序，分数最高的一个或多个分类为最终分类结果，其它分类作为参考结果。也就意味着本文的分类方法支持多标签分类。

2.2.2 基于感知机的权重参数优化

为提高线性分类器的性能，本文使用感知机学习算法，利用标注训练语料对情感分类算法中的权重参数 w 进行优化。算法伪代码描述如下：

```
Algorithm 1: Perceptron Learning
1: Input: training_set
2:  $w(1) \leftarrow \mathbf{0}$ 
3: do
4:   error_count = 0
5:   for  $[f_i, z_i]$  in training_set
6:      $s = \sum_{i=1}^N w_i \cdot f_i$ 
7:     result =  $(s > \text{threshold}) ? 1 : 0$ 
8:     error =  $z_i - \text{result}$ 
9:     if error  $\neq 0$ 
10:      then error_count++
11:       $w = w + \alpha \cdot (z_i - \text{result}) \cdot f_i$ 
12:    end if
13:  end for
14: until: error_count == 0
15: Output: w
```

其中 f_i 是训练语料中第 i 个句子匹配特征向量，2.2.1 节中例句 1 匹配句模 1 的匹配特征向量 $f_i = [1, 0, 2, 4]$ 。 w 是分类算法的 4 个匹配特征的打分权重向量。 z_i 是期望分类，表示当前句模所在情感类与第 i 个句子标注的情感分类是否相同，相同则 $z_i = 1$ ，不相同则 $z_i = 0$ ，例句 1 标注的情感分类是喜爱类，与句模 1 所在分类相同，所以 $z_i = 1$ 。 training_set 是训练语料对应的匹配特征向量 t 与期望分类 z 的集合。 error_count 记录变量 error 不为 0 的个数，当训练语料中所有的句子对应的 error 变量值都为 0 时程序结束。 α 为学习因子，取值在 0 到 1 之间。

2.2.3 特殊词语处理

在算法设计与分析过程中，发现以下几点问题：

(1) 不规则词问题：一些语料尤其是微博语料中经常出现不规则词和短语，分类算法无法识别句模库中没有的词。例如：句子：“刚才的拔河比赛，太鸡冻了”，句模库的激动词类中只有“激动”而没有“鸡冻”。解决办法是搜集不规则词将其添加到词类库中。

(2) 分词错误问题：例如：对“自己是最棒的”的分词结果为“自己/是/最/棒/的”，如果能将“最棒的”作为一个独立单元来处理，更有利于根据关键词选择候选分类。解决办法是将具有明显情感的短语加入到自定义词表中作为一个词处理。

3 实验结果与分析

3.1 实验设置

本文实验使用两个领域的语料：语料 1 为 NLP&CC 2013 中文微博情绪识别评测数据，简称 NLP&CC 语料。共包括 4000 条微博中的 13250 个句子，其中情感句 4949 句，无情感句 8301 句。情感句共分 7 类：Anger 愤怒、Disgust 厌恶、Fear 恐惧、Happiness 高兴、Like 喜好、Sadness 悲伤、Surprise 惊讶。每个句子最多属于两个情感分类。语料 2 为 ReLabEAC1.0 文学语料，包括 1487 篇文学短文共 34954 个句子，其中情感句 32171 句，无情感句 2783 句。情感句共分 8 类：Sorrow, Anger, Anxiety, Surprise, Hate, Love, Joy, Expect。每个句子可属于一个或多个情感分类。

实验使用评估指标为：

$$\text{准确率: Precision} = \frac{\# \text{system_correct}(\text{emotion}=\text{Y})}{\# \text{system_proposed}(\text{emotion}=\text{Y})}$$

$$\text{召回率: Recall} = \frac{\# \text{system_correct}(\text{emotion}=\text{Y})}{\# \text{gold}(\text{emotion}=\text{Y})}$$

$$F\text{-值: } F\text{-Measure} = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

待分类句子的多标签分类结果中任意一个标签分类正确，则视为正确分类句。
 $\#system_correct(emotion=Y)$ 是对语料中情感句分类正确的句子数目， $\#system_proposed(emotion=Y)$ 是语料中的情感句总数， $\#gold(emotion=Y)$ 是语料中被划分为情感句的句子数目。

本文采用的 Baseline 系统是基于词语特征 SVM 分类器的方法。

3.2 实验结果及分析

实验 1: 与 Baseline 系统对比实验。

Baseline 系统使用 SVM 分类器和词语特征分别对 NLP&CC 微博语料和 ReLabEAC1.0 文学语料进行分类，具体方法是：对情感句进行分词标注词性后挑选所有的名词、动词和形容词组成一个词汇表，以待分类句子中的词是否在词汇表中出现以及出现的频率为特征，随机挑选 2/3 的句子为训练语料，1/3 的句子为测试语料，使用 SVM Multi-Class 工具进行训练和测试。Baseline 系统和本文提出的分类器获得的最高准确率统计结果如表 2 所示：

表 2 NLP&CC 微博语料 Baseline 系统对比实验结果

	Precision	Recall	F-Measure
Baseline	0.245	0.567	0.342
基于情感句模的分类算法	0.434	0.323	0.370

表 3 ReLabEAC1.0 文学语料 Baseline 系统对比实验结果

	Precision	Recall	F-Measure
Baseline	0.267	0.584	0.366
基于情感句模的分类算法	0.607	0.353	0.446

由此可见，基于情感句模的分类算法相比 Baseline 方法可以达到较高的准确率。特别是 ReLabEAC1.0 语料上可以达到很高性能，这是由于这个语料中用词比较规范，因此，基于情感句模的分类算法准确率较高。

实验 2: 权重优化影响实验。

使用 2.2.2 节描述的感知机算法，取 $\alpha=0.1$, $threshold=0.5$ ，对 NLP&CC 语料中情感句进行训练，得到 $w=[0.6, 0.2, 0.6, -0.2]$ ，为方便计算将每个权值放大 10 倍后取整，得到 $w=[6, 2, 6, -2]$ 。分别使用均等权重、经验权重和感知机优化特征权重，对 NLP&CC 语料中情感句和无情感句进行分类，分类结果如表 4 所示：

与均等权重相比，采用感知机学习算法优化特征权重后，分类算法性能提升了约 3%。与经验权重相比，召回率略有上升，准确率有所下降，F 值略微上升。

表 4 评估特征权重影响实验结果

	Precision	Recall	F-Measure
均等权重	0.397	0.306	0.346
经验权重	0.434	0.323	0.370
感知机优化权重	0.413	0.345	0.376

实验 3: 句模数量与分类准确度及句模库对中文情感句的覆盖率的评估。

对 NLP&CC 语料中情感句进行分类，统计与分类正确的情感句匹配频率最高的 10 个句模，如表 5 所示。

进一步，分别统计与分类正确的情感句统计匹配频率最高的 10 个、20 个、50 个、100 个、150 个句模，以及与它们匹配的句子个数，统计结果如图 5 所示。

表 5 与分类正确情感句匹配频率最高的 10 个句模

句模	分类	例句
<感事><希望词类><向事><好_词类 融洽词类>	喜爱类	妈妈希望儿子好
<施事>[<在_范围介词词类><范围>上]<重视词类 1><向事>	重视类	张三在工作上重视李四
<感事><喜爱词类><向事>	喜爱类	我爱你
<施事>[<在_范围介词词类><范围>上]<赏识词类><客事>	赏识类	张三在工作上赏识李四
<受事><有_词类><涉事>[的 之]<嫌疑词类>	怀疑类	张三有盗窃之嫌
<施事><对_当事介词词类><向事><热情词类>	热情类	他对人很热情
<感事><对_对象介词词类 为替对象介词词类><向事><感到词类><愤怒词类 1>	愤怒类	他对她感到愤怒
<感事><感到觉得词类><惊讶词类>	惊讶类	他感到很惊讶
<施事><使让词类><致事><感到词类><悲伤词类>	悲伤类	这件事使我感到悲伤
<感事><感到词类><高兴词类>	高兴类	他感到很高兴
<施事><希望词类><向事><意图词类>	希望类	他希望回家过年
<受事><有_词类><涉事>[的 之]<嫌疑词类>	怀疑类	张三有盗窃之嫌

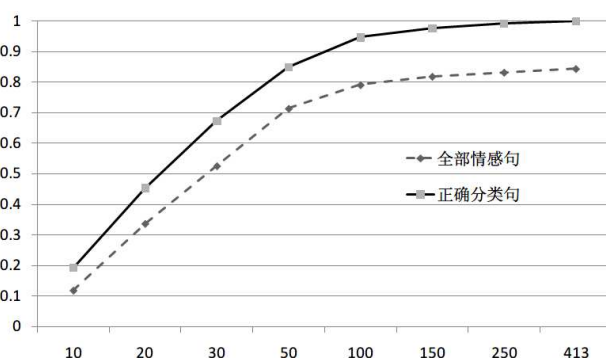


图 5 句模数量与匹配句子百分比统计图

图 5 中横坐标为高频句模数量，实线图的纵坐标为高频句模匹配的句子占全部分类正确的情感句的百分比，虚线图的纵坐标为高频句模匹配的句子占全部情感句的百分比。可以发现前 150 个高频句模匹配了 97.6% 的分类正确的情感句，覆盖了绝大多数分类正确的情感句，对 NLP&CC 语料中全部情感句的覆盖率为 40.7%。

此外，分别只使用匹配频率最高的 10 个、20 个、30 个、50 个、100 个、150 个句模对 NLP&CC 语料进行分类，统计分类准确率如图 6 所示。

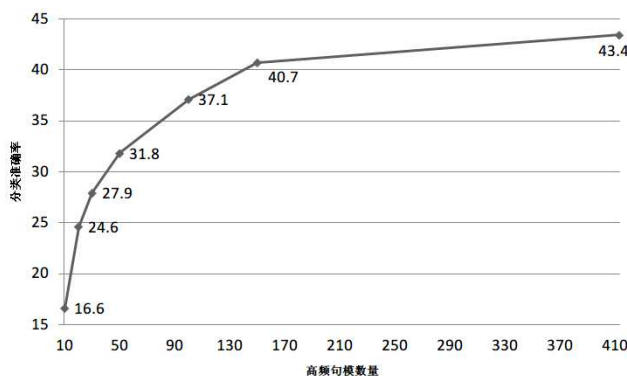


图 6 句模数量与分类准确率统计图

图中横坐标为高频句模的数量，纵坐标为只使用这些句模对 NLP&CC 语料中情感句进行分类的准确率。可以发现只使用前 10 个高频句模对 NLP&CC 语料 7 个分类的分类准确率为 16.6%，只比随机分配 $1/7=14.3\%$ 的概率高 2 个百分点。随着高频句模数量的增加，分类准确率迅速提高。当使用前 150 个句模时，分类准确率为 40.7%，与使用全部 413 个句模的

准确率 43.4% 只相差 3 个百分点。由此可见, 前 150 个高频句模对分类效果起到决定性影响, 继续增加句模数量对分类准确率提高效果不明显。

实验 4: 分类特征、自定义词表与分类效果关系实验。

分别使用以下 3 种方法进行实验:

方法 1: 只使用分类算法 4 个分类特征中词级特征 f_1 和 f_2 , 对 NLP&CC 语料进行分类。

方法 2: 使用全部 4 个特征对 NLP&CC 语料进行分类。

方法 3: 使用全部 4 个特征加自定义词表对 NLP&CC 语料进行分类。

分类统计结果如表 6 所示:

表 6 实验 4 统计结果

	Precision	Recall	F-Measure
方法 1	0.397	0.282	0.330
方法 2	0.427	0.323	0.368
方法 3	0.434	0.323	0.370

由表 6 可知: 方法 2 在方法 1 的基础上使用依存关系特征 f_3 和句法特征 f_4 分类性能提升明显。方法 1 效果较差的原因是句模库中有些同义词的含义并不能与句模匹配。例如:《同义词词林》中“细心”的同义词包括“致密”、“逐字逐句”、“细瞧”、“有心人”等, 与细心类的句模:“<当事><细心词类>”中的“细心词类”并不匹配, 这些词在细心词类中会降低分类系统的召回率。方法 3 比方法 2 多了自定义分词词表, 分类效果略有提高。这说明自定义词表能够提高分类效果, 但自定义词都是针对特定句子的情感表达手工抽取添加, 目前的规模还不够, 覆盖范围有限, 对分类效果提升有限。

3.3 实验结果分析

实验 1 至 5 表明, 继续增加分类系统的句模数量和优化打分权重对分类效果提高影响不大。下一步考虑增加新的匹配特征, 例如: 句模中的语义角色特征等。

实验过程中发现下列问题: (1) 句模库中只有情感类句模, 没有无情感类句模。任何与 4 个匹配特征中任意 1 个匹配的句子都会划分为情感句, 导致无情感句被划为情感句的概率较高, 降低了系统的性能。下一步将考虑优化特征匹配得分的阈值或完善三大情感分类以外的基于谓词的分类, 并构建相应的句模库。(2) 情感分类和句模都是基于谓词和相应规则构建的, 对显式含有情感词或情感搭配的句子比较有效。而类似“如同观看一部真正的大片一样”这样的句子中, 表达情感的要素是名词“大片”和它的修饰语“真正的”, 比较难用句模匹配的方法划分情感类别。类比句或比喻句的情感分类是十分困难的, 下一步考虑增加相应的匹配特征, 尝试结合基于统计的方法, 提高隐含情感句子的分类效果。

4 结语

本文设计和实现了一种半自动获取情感句模的方法, 使用句模分类的方法实现对情感句的分类, 在两个情感测试数据集上的实验结果显示本文提出的方法可以稳定有效地提高文本情感分类性能。目前情感分类的细致划分还在继续进行中, 计划细化同义词词类和尝试添加更多匹配特征。另外还将加入基于统计的情感分类方法, 构筑相应的训练语料和测试语料, 提高对隐含情感句子的分类效果。

参考文献:

- [1] Peter D. Turney. Thumbs Up or Thumbs Down? Semantic Orientation Applied to Unsupervised Classification of Reviews [C]// Proceedings of ACL 2002: 417-424
- [2] Kamps J, Marx M, Mokken RJ. Using WordNet to Measure Semantic Orientation of Adjectives. [C]// Proceedings of LREC. 2004: 1115-1118
- [3] 朱嫣岚, 闵锦, 周雅倩, 黄莹菁, 吴立德. 基于 HowNet 的词汇语义倾向计算[J]. 中文信息学报, 2006, 20(1): 14-20

-
- [4] Hu Mingqing, Liu B. Mining Opinion Features in Customer Reviews.[C] Proceedings of AAAI 2004: 755-760
- [5] 姚天畴, 聂青阳, 李建超, 李林琳等. 一个用于汉语汽车评论的意见挖掘系统. [C]//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. 北京: 清华大学出版社, 2006: 260-280
- [6] 刘鸿宇, 赵妍妍, 秦兵, 刘挺. 评价对象抽取及其倾向性分析[J]. 中文信息学报, 2010, 24(1): 84-89
- [7] 任巨伟, 杨源, 王昊等. 二元情感常识库建设及其在文本情感分析中的应用[OL]. 中国科技论文在线, 2013, <http://www.paper.edu.cn/releasepaper/content/201301-158>
- [8] 陈健美, 林鸿飞. 中文情感常识知识库的构建[J] 情报学报, 2009, 28(4): 492-498
- [9] Bo Pang, Lillian Lee, Shivakumar Vaithyanathan. Thumbs up? Sentiment Classification using Machine Learning Techniques [C]//Proceedings of EMNLP 2002: 79-86
- [10] 谷学静, 王志良, 刘冀伟等. 基于 HMM 的人工心理建模方法的研究[C]//第一届中国情感计算及智能交互学术会议, 北京, 2003: 31-36
- [11] 王根, 赵军. 基于多重冗余标记 CRFs 的句子情感分析研究[J]. 中文信息学报, 2007, 21(5) : 51-56
- [12] S. Li and C. Zong, Multi-domain Sentiment Classification[C]//Proceedings of ACL-HLT 2008: 257-260
- [13] 李寿山, 黄居仁, 基于 Stacking 组合分类方法的中文情感分类研究[J] 中文信息学报, 2010, 24(5): 56-61
- [14] 朱晓亚, 范晓. 二价动作动词形成的基干句模[J]. 语言教学与研究, 1999 : 111-122
- [15] 鲁川, 缙瑞隆, 董丽萍. 现代汉语基本句模[J]. 世界汉语教学, 2000, 54(4) : 11-24

附录 1. 情感句模分类

①态度类分为：支持类，反对类，怀疑类，沉默类，耐心类，怨气类，果断类，信心类，冒险类，妥协类，热情类，冷淡类，粗暴类，诚恳类，温柔类，和蔼类，客气类，宽容类，霸道类，谦虚类，细心类，勤奋类，负责类，积极类，谨慎类，粗心类，亲密类，团结类，一见如故类，熟悉类，恋爱类，和睦类，疏远类，友好类，纠纷类，纠缠类，挑逗类，苛刻类，重视类，严格类，轻视类。共 41 个二级分类。

②感受类分为：吸引类，为荣类，自娱类，为耻类，不知所措类，伤感情类，感知类，记得类，生理感觉类，非生理感觉类，听到类，偷听类，看见类，偷看类，知道类，不知道类，发现类，惭愧类，愤怒类，义愤类，幸灾乐祸类，敬佩类，羡慕类，感激类，谴责类，害怕类，喜爱类，溺爱类，讨厌类，仇恨类，宽慰类，失望类，担忧类，高兴类，悲伤类，惊讶类，满意类，不满意类，没耐心类，懊悔类，紧张情绪类，心安类，自豪类，慌张类，眼熟类，耳熟类，眼生类，耳生类。共 48 个二级分类。

③思想类分为：希望类，自愿类，向往类，思考类，想象类，相信类，鉴别类，主张类，接受类，看待类，信任类，宠信类，看得起类，另眼相看类，想念类，着想类。共 16 个二级分类。