

基于社交网络大数据的 民众情感监测研究

李爱黎¹, 张子帅¹, 林荫², 王秋菊², 杨建安¹, 孟炜程¹, 张岩峰¹

1. 东北大学计算机科学与工程学院, 辽宁 沈阳 110000;

2. 东北大学外国语学院, 辽宁 沈阳 110000

摘要

近年来, 新浪微博、推特等社交网络平台逐渐成为反映社会舆情的主要载体之一, 为网民发表观点和表达情绪提供了便利。基于社交网络大数据的舆情监控已经成为新的研究热点, 利用各国的社交网络大数据进行民众情感监测, 有助于直接掌握国际关系中的民众情感倾向, 对我国外交、对外贸易等方面都有很重要的作用。基于此, 提出了一种面向中日语料的民众情感监测系统, 该系统能够同时分析新浪微博和推特等社交平台的中日文语料数据中包含的情感倾向, 并以可视化的形式展现给用户。情感分析算法方面, 在BERT模型基础上结合自扩展的中日文情感词典, 提出了一个新的情感分析模型——EmoBERT。实验结果表明, 相比于原始BERT模型, EmoBERT模型在中文情感分类任务和日文情感分类任务上都取得了很好的表现。其中中文模型EmoBERT-C将中文BERT模型准确率从89.68%提升至92.15%, 日文模型EmoBERT-J将日文BERT模型准确率从74.73%提升至78.26%。

关键词

情感分析; 舆情监测; 情感词典; 中日关系; 微博; 推特

中图分类号: TP311.13

文献标志码: A

doi: 10.11959/j.issn.2096-0271.2022054

Research on emotion monitoring of public based on social network big data

LI Aili¹, ZHANG Zishuai¹, LIN Yin², WANG Qiuju², YANG Jianan¹, MENG Weicheng¹, ZHANG Yanfeng¹

1. School of Computer Science and Engineering, Northeastern University, Shenyang 110000, China

2. Foreign Studies College, Northeastern University, Shenyang 110000, China

Abstract

In recent years, social networking platforms such as Sina Weibo and Twitter have gradually become one of the main carriers for reflecting social public opinion, providing a convenient platform for netizens to express their opinions and emotions. Public opinion monitoring based on social network big data has become a new research hotspot. People's emotions monitoring using social network big data in various countries is helpful to directly grasp people's emotional tendencies in international relations,

and has a great impact on the diplomacy, foreign trade, and other aspects. Based on this, a public sentiment monitoring system for Chinese and Japanese data was proposed, which could analyze the emotional tendencies contained in Chinese and Japanese data on social platforms such as Sina Weibo and Twitter simultaneously, and displayed them to users in a visual form. In the aspect of sentiment analysis algorithm, based on the BERT model and combined with the self-expanding Chinese and Japanese sentiment lexicon, a new sentiment analysis model, EmoBERT, was proposed. The experimental results show that, compared with the original BERT model, the EmoBERT has achieved good results on both Chinese sentiment classification tasks and Japanese sentiment classification tasks. Among them, EmoBERT-C increases the accuracy of Chinese BERT from 89.68% to 92.15%, and EmoBERT-J increases the accuracy of Japanese BERT model from 74.73% to 78.26%.

Key words

sentiment analysis, public opinion monitoring, sentiment lexicon, Chinese-Japanese relation, Weibo, Twitter

0 引言

互联网的飞速发展改变了人们传统的交流习惯,人们对网络的利用率越来越高。互联网上相继出现了社区、论坛、微博等形式的社交网络平台,用户在网上通过这些平台表达自己对某一事件的看法和态度,这些信息包含了大量的社会热点及情感倾向^[1]。因此,在大数据技术支撑下,挖掘社交网络中用户的观点、态度和情感,并服务于社会,是一个很有意义的工作。

目前,多数研究主要对主流社交网络平台的热门数据进行情感分析与监测。Zhao J C等人^[2]构建了一个面向中文微博的情感分析系统,对异常或突发事件进行监测。Wang H等人^[3]利用推特(Twitter)数据构建了一个针对2012年美国大选结果的实时预测系统,通过统计美国民众对4位候选人的情感倾向来预测大选结果。Williams J等人^[4]提出了一种预测算法,针对Twitter上发生的某一事件,预测其发生时间。

然而目前的大多数研究仅仅是对微博、Twitter等某单一平台进行舆情数据的情感

分析,并且多数是针对中文语料和英文语料的分析,国内使用日文语料进行情感分析的研究极少。与此同时,针对海量的数据,利用人工浏览、打标签的方式来获取用户情感是一件极其复杂且困难的事情。

因此,本文提出一种面向中日语料的民众情感监测系统,该系统能够同时分析微博和Twitter等社交平台的中日文舆情数据中包含的情感倾向。当某一焦点事件发生时,自动进行中日两国民众的情感对比,供相关舆情部门监测。情感分析算法方面,本文在BERT(bidirectional encoder representations from transformers)模型(以下简称BERT)的基础上提出了一个新的情感分析模型——EmoBERT,利用自扩展的情感词典改进了BERT的预训练任务,并提出了情感词增强的注意力机制,弥补了BERT在预训练阶段情感特征提取不充分的缺陷。实验证明,相比于BERT,EmoBERT在中文和日文情感分类任务上都取得了更好的效果。

此外,为了更好地对算法结果进行存储和展示,本文采用Flask框架搭建网站,设计并实现了一个自动化情感监测系统,使用ECharts库实现情感分析结果的可视化展示,可交互的动态曲线可以让

用户实时监测到中日民众对某一事件的情感态度变化。针对情感突变点和情感差值较大的区间,通过词频-逆向文件频率(term frequency-inverse document frequency, TF-IDF)算法自动生成该区间的关键词,同时给出热度排名前5的博文。经验证,生成的关键词和给出的博文可以很好地对应情感突变点处发生的事件。

1 相关工作

舆情监测系统是一个网络应用系统,用于监测由热门事件或突发事件引发的有影响力且倾向性强的观点和言论^[5]。同时,情感分析技术也被广泛应用于舆情监测的研究,成为当前舆情监测研究中的主流方法之一。进入21世纪后,情感分析在自然语言处理的各个领域都被广泛研究,如数据挖掘、文本挖掘、舆情监测和信息检索等。情感分析最早是在2003年由Yi J等人^[6]提出的,对文本中包含的情感进行计算,进而分析用户的情感倾向和观点。

对文本数据进行情感分析的方法主要有3种,分别是基于情感词典的情感分析方法、基于机器学习的情感分析方法以及基于深度学习的情感分析方法。3种情感分析方法优点与不足见表1。

1.1 基于情感词典的情感分析方法

基于情感词典的情感分析方法是在标

注极性或极性分数单词的基础上,比对情感文本中包含的极性情感词,然后采用简单统计的方法或权值算法进行情感分类。此方法不需要训练数据,因此被广泛地应用于传统的文本情感分析。20世纪90年代末期,国内外开始了有关文本情感分析的研究。Riloff E M等人^[7]基于语料数据构建了语义词典。熊德兰等人^[8]基于知网(HowNet)常识知识库研究了句子的褒贬性。潘明慧等人^[9]提出了基于词典的方法识别微博表达的6种情绪。

1.2 基于机器学习的情感分析方法

基于机器学习的情感分析方法也被广泛应用于情感分析领域。该方法首先建立一个训练集,并根据用户情感标记数据;然后从训练集中提取特征,构建分类模型,进而预测没有标签的数据;最后通过分类器对未标记的数据进行情感倾向判定。在国外,Pang B等人^[10]使用了3种机器学习的方法进行对比试验,分别是朴素贝叶斯、最大熵和支持向量机,对电影评论进行了情感极性分类,将情感极性分为积极和消极,并比较3种方法的实验结果,其中支持向量机方法的分类效果最好。国内也有很多学者比较不同的分类算法,杨艳霞^[11]使用两种机器学习方法对微博数据集进行了情感分析,分别是贝叶斯和支持向量机,同时比较了两种方法在分类性能上的优劣,其中贝叶斯方法的准确率更高。

表1 情感分析方法优点与不足

方法	优点	不足
基于情感词典的情感分析方法	不需要人工标注,可以用简单的统计方法进行情感分类	分类的结果依赖词典的质量和规模
基于机器学习的情感分析方法	适用于较小的数据集,泛化能力强	需要投入大量的人工成本标注数据集
基于深度学习的情感分析方法	学习能力强,覆盖范围广,适应力强,可移植性好	计算量大,硬件成本较高,模型设计复杂性好

1.3 基于深度学习的情感分析方法

2011年, Collobert R等人^[12]在解决词性标注等问题时最早将深度学习应用到了自然语言处理领域。深度学习最大的特点是可以自动学习批量数据,从而挖掘数据中的潜在特征,并通过注意力机制实现对目标内容的增强关注,在训练过程中进行参数的调整^[13]。Schuller B等人^[14]引入了一种基于长短期记忆(long short-term memory, LSTM)循环神经网络语言模型的新方法,不需要进行任何特殊的预处理或特征选择。宋婷等人^[15]为了提取方面级的情感,提出了分层的LSTM模型。徐志栋等人^[16]提出一种基于胶囊网络的方面级情感分类模型——SCACaps,解决了方面级情感分析中多重情感造成的特征重叠问题。张宝华等人^[17]提出了一种多输入模型,该模型结合了多通道卷积神经网络(multi-channel convolutional neural network, MCNN)、LSTM和全连接神经网络。而深度学习中的迁移学习(transfer learning)也常常被应用于舆情分析领域,如以美团外卖的评论数据为原始数据集,抽取其特征,建立美团外卖评论的情感分析模型,再将其应用到相应的目标域(如电影评论的情感分析),以此实现模型的大规模迁移。基于此,迁移学习逐渐成为舆情分析领域的研究热点。Radford A等人^[18]提出了名为OpenAI GPT的预训练模型,该模型可以经过少量的微调后用于各种下游任务。近年来,BERT作为一个强大的预训练模型,首先在大规模的语料库上进行预训练,获取通用的语言模型,然后进行一系列的微调以吸收下游具体任务的相关知识^[19-20]。但是在情感分类任务上,BERT还有一定的提升空间,这是因为BERT在预训练阶段并没有考虑任何情感

信息。为了解决这个问题,本文将融合情感词典和BERT,将情感特征引入预训练过程。

2 数据来源与数据预处理

2.1 数据来源

利用Python语言进行编程,使用微博应用程序接口(application programming interface, API)、日本Twitter API、网络爬虫等技术完成舆情数据的获取。

(1) 中文舆情数据

在中国,新浪微博(以下简称微博)具有用户多、消息数量大、更新快等特点,成为人们获取信息、发表舆论的主要平台之一,越来越多的民众习惯在微博这一社交网络平台上交流观点、分享信息。这些信息包含了大量的社会热点及情感,能很好地反映民众对话题的关注和态度。本研究以“日本”为关键词,通过网络爬虫技术,爬取了2013—2021年的舆情数据。数据来自科技、体育、娱乐、经济、疫情5个类别,数据主要包括微博标题、统一资源定位符(uniform resource locator, URL)、时间、内容、点赞数、评论数及转发数等。

(2) 日文舆情数据

在国外, Twitter无疑是拥有巨大访问量的社交网络平台之一。在日本, Twitter作为互联网Web2.0时代的最新应用,逐渐影响和改变世界的交流和沟通方式。因此, Twitter数据十分适合进行国外舆情分析。本研究以“中国”为关键词,通过网络爬虫技术,同样爬取了2013—2021年的舆情数据。数据来自科技、体育、娱乐、经济、疫情5个类别,数据主要包括推文标题、URL、时间、内容、点赞数、评论数及转发数等。

2.2 数据预处理

微博和Twitter数据文本含有很多标签、注释等特殊符号,使用Python自然语言处理工具包NLTK和正则表达式等工具对数据进行清洗。由于情感分析的质量依赖情感词典,因此必须对清洗后的数据做分词处理,同时移除停用词。本文以jieba为中文分词工具,以MeCab为日文分词工具。主要处理过程包括中文分词、日文分词、提取词元(token)、词根化(stemming)、移除停用词等。

其中微博的中文数据集采用“百度停用词表”进行过滤。由于现有的日文停用词表中停用词较少且不够全面,因此,笔者在现有的停用词表基础上扩充了新的日文停用词表,并将其用于Twitter日文数据集。

3 基于自扩展情感词典的情感分析模型

微博和Twitter中的文本具有领域广、

更新速度快等特点,而通用情感词典存在领域差异、知识覆盖率较低、情感词权值过于固定等问题,因此,本文利用自扩展的中日文情感词典提出一种情感极性量化算法来计算文本的情感强度值,通过计算情感词权值来量化该文本的情感强度。具体方案如下。

首先,在通用情感词典的基础上,构建适合本研究领域的中日文情感词典,并对预处理后的数据进行情感倾向性分析;然后,分别构建中文和日文的程度词表和否定词表,之后对特殊标点符号进行量化加权;最后,考虑点赞数,对情感分值进行加权计算,得到最终的情感分值,并进行情感分类。情感分析框架如图1所示。

3.1 情感词典构建

情感倾向是用户对某一事物主观的内心喜恶及主观评价的一种倾向。不同的情感词或情感语气可以表达不同程度的情感倾向。通常给每个情感词赋予不同的权值。例如,“楽しい”和“嬉しい”,都表达开心,

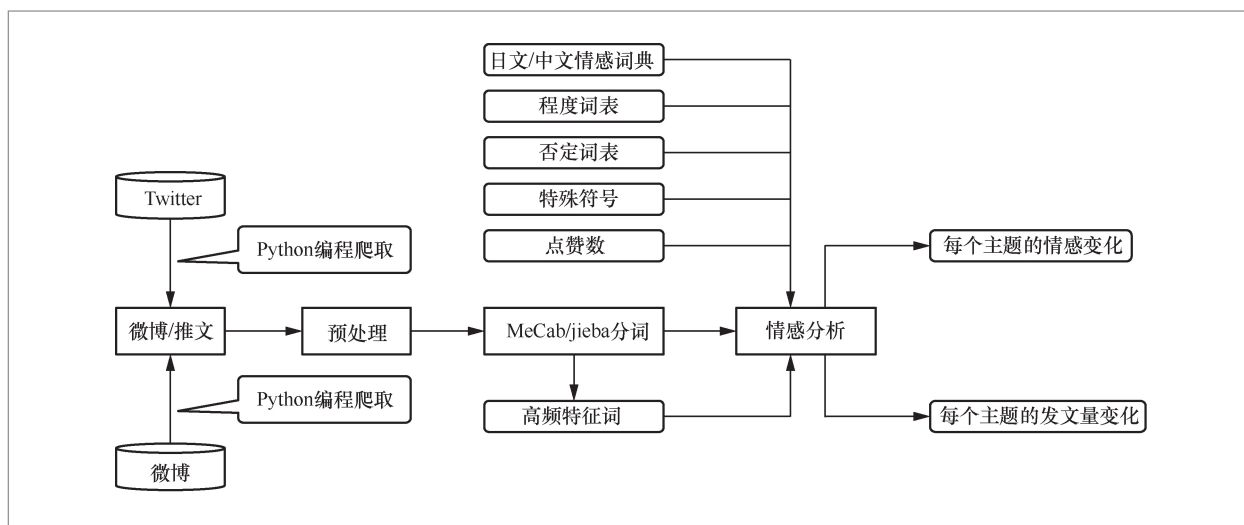


图1 情感分析框架

但是在表达情感程度上,“嬉しい”要比“楽しい”更强烈。中文的“讨厌”与“厌恶”都表达消极情感,但是“厌恶”的情感程度会更强烈。因此,情感词典能否覆盖全面在一定程度上影响着情感分类结果,情感词典的构建是情感分析研究的基础,本文尽可能构建一个足够大、覆盖面足够广的情感词典,并将其应用于中日舆情研究领域。本文构建情感词典组成如图2所示

- 基础情感词典。基础情感词典总结并整理当前已有的情感词典资源。对于中文情感词典,本文将HowNet情感词典作为中文基础情感词典,其组成见表2。对于日文情感词典,本文集成了多个开源的日文情感词典用于Twitter数据集。

- 网络情感词典。随着互联网的高速发展,网络用语应运而生。网络用语的形式和传统词语有着很大区别,它们往往具有强烈的感情色彩。这些词语是不包含在基础情感词典当中的,但在判别情感倾向中起着重要的作用。

- 领域情感词典。从微博和Twitter

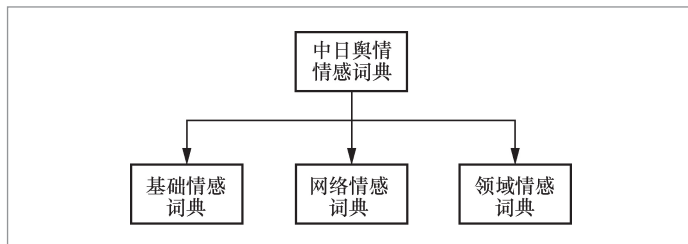


图2 情感词典组成

表2 基础情感词典组成

类别	个数/个
积极评价词语	3 730
消极评价词语	3 116
积极情感词语	837
消极情感词语	1 255

获取的中日舆情数据集中选取适合本研究领域、情感鲜明的词,并将其作为基准词,通过基于扩展的情感倾向点间互信息 (semantic orientation pointwise mutual information, SO-PMI) 算法计算候选词与基准词的相似度,以此判断候选词的情感倾向,将领域情感词自动加入基础情感词典,构建适用于中日舆情领域的中日文情感词典。

情感词典中积极词语分值为1,消极词语分值为-1。假设一个句子中包含pos_num个积极词语和neg_num个消极词语,那么该句子的情感分值score计算过程如下:

$$\text{pos_score} = \text{pos_num} \quad (1)$$

$$\text{neg_score} = \text{neg_num} \quad (2)$$

$$\text{score} = \text{pos_score} + \text{neg_score} \quad (3)$$

3.2 程度词表构建

笔者发现,在网民发布的微博和推文中,情感词语前大有含副词修饰,如在“非常喜欢”“很喜欢”中,“非常”修饰“喜欢”,“很”也修饰“喜欢”,但是,“非常”所表达的情感强度显然要多于“很”。为了更加准确地计算文本的情感倾向,本文构建了中文程度词表,将程度副词分为5个等级(极强级、中强级、中级、中弱级、微弱级),并将程度副词的强度取值范围限定在[0,3]。人工标注这些程度副词语气的强弱,并用一个二元组level(adv, intensity)来表示,其中adv表示词语名称, intensity表示该词的语气强度,一个副词的语气强度取值范围为[0,3],越接近0说明该词表达的情感强度越弱,越接近3说明该词表达的情感强度越强烈。例如,“出头”的强度设置为0.5,“更”的强度设置为2,“极

其”的强度设置为3。同样,本文也构建了日文程度词表,并人工对这些程度词的语气强度进行了标注,与中文的构建方法相同,此处不再赘述。

如果一个积极词语前后出现程度词,那么:

$$\text{pos_score} = \text{pos_score} \times \text{inte} \quad (4)$$

如果一个消极词语前后出现程度词,那么:

$$\text{neg_score} = \text{neg_score} \times \text{inte} \quad (5)$$

其中, inte 为程度词的强度值。

3.3 否定词表构建

在微博和Twitter文本中,否定词也是经常出现的。例如“不公平”“拒绝接受”,其中“不”用来否定“公平”,“拒绝”用来否定“接受”。为了使文本情感倾向性的计算更加准确,本文分别构建了中文否定词表和日文否定词表。对于否定词表,不需要对其进行标注,用一个列表list来表示,当一个积极词语前出现一个否定词,则该词语的情感分值变为原来的相反数。如果是双重否定,则该词语的情感分值不变。消极词语同理。

$$\text{pos_score} = -1 \times \text{pos_score} \quad (6)$$

$$\text{pos_score} = -1 \times (-1) \times \text{pos_score} \quad (7)$$

3.4 感叹句

除此之外,本文认为带有“!”的感叹句往往比陈述句的语气更强烈。因此,本文定义了感叹句的加权计算式:

$$\text{score} = \text{score} \times 1.2 \times n \quad (8)$$

如果一句话是感叹句,那么对该句的情感分值进行加权,其中 n 为感叹号的数量。

3.5 点赞数

对于微博和Twitter这类热门的社交网络平台,笔者认为,一篇微博或推文的点赞数能很好地说明其他网民对该观点的支持度,即点赞数越多的文本,应该被赋予更高的权重。因此,本文将点赞数映射到[0,1],对句子的情感分数进行加权计算,以得到更准确的分析结果。

3.6 TF-IDF关键词抽取

TF-IDF算法是一种常被用来计算一个字或词语对于一篇文档的重要程度的统计方法。如果某个词语在一篇文章中频繁地出现,但在其他文章中很少出现,那么就认为该词或者短语对于该文章具有一定的代表性,适合用来分类。计算式如下:

$$\text{TF} = \frac{\text{某个词在文章中出现的次数}}{\text{文章的总词数}} \quad (9)$$

$$\text{IDF} = \log \left(\frac{\text{语料库的文档总数}}{\text{包含该词的文档数} + 1} \right) \quad (10)$$

$$\text{TF-IDF} = \text{TF} \times \text{IDF} \quad (11)$$

假设一篇微博共有100个词语,其中“会议”出现了5次,那么“会议”一词在该微博中的词频TF就是 $5/100=0.05$ 。而计算逆文档频率IDF的方法是测定有多少篇微博出现过“会议”一词,加1后除数据集中的总微博数,最后求对数。因此,如果“会议”一词在100条微博中出现过,而文件总数是1 000 000份的话,那么,其逆文档频率IDF就是 $\log \left(\frac{1000000}{101} \right)$,约等于4,最后的TF-IDF值约为0.2。

一般来说,一篇文章中某个字或词语的TF-IDF值越大,这个字或词语对于这篇文章越重要,因此通过计算文章中每个字或词语的TF-IDF值,按大小排序,排在最前面的几个字或词语,就是可以代表该文章的关键词。

基于前面的情感分析方法可以得到每个主题的情感变化情况。接下来,通过TF-IDF算法自动生成各个时间区间内词频排名前10的关键词,关键词可以体现出情感变化的原因及引发情感突变的事件。

4 结合中日文情感词典的情感分析模型EmoBERT

本文以BERT预训练模型为基础,利用自采集的中日舆情数据集,提出一种结合中日文情感词典的情感分析模型EmoBERT。

4.1 BERT及其缺陷分析

BERT是一种基于双向Transformer的大规模预训练语言模型,其利用了Transformer的编码层,并在此基础上加入掩码机制,能够自动进行预测训练。“双

向”是指它在处理一个词时,可以考虑到该词前后的单词的信息,进而得到上下文的语义。本质上是在大规模数据集的基础上,使用自监督学习方法对单词学习进行模型训练^[21]。

BERT结构如图3所示,可大致分为输入层、Transformer层以及输出层。其中 E_i 表示BERT输入的编码向量, T_i 表示BERT输出的编码向量, T_m 表示Transformer的编码器结构。

BERT输入层的编码向量包含3种嵌入特征,分别是词嵌入、段嵌入和位置嵌入,如图4所示。为了使BERT适应下游的任务,在输入时,为每个句子附加[CLS]和[SEP],这是两个特殊符号:[CLS]用于下游的分类任务,最终输出时可以用来表示整个句子;[SEP]用来分割两个句子,如[CLS]+句子A+[SEP]+句子B+[SEP]。

- 词嵌入。词嵌入从词汇表学习得到每个特定词的嵌入特征,词嵌入层会将其转换成768维的向量,如图4所示的句子会被转换成一个(10, 768)的矩阵。

- 段嵌入。段嵌入被用来区别两种句子, E_A 表示第一句话, E_B 表示第二句话。文本中的多个句子被拼接在一起后送入BERT,BERT通过段嵌入区分每个句子。

- 位置嵌入。位置嵌入指将单词的位置信息编码成特征向量的形式,将单词位置关系引入BERT,BERT通过学习得到位置向量,实际上,位置嵌入是一个大小为(512,768)的查找表,其中第*i*行是指一个句子中第*i*个位置上的单词的向量。

将3种嵌入向量简单相加,得到模型的输入向量,同时传递给BERT的编码层作为输入表示。**TE**为词嵌入向量,**SE**为段嵌入向量,**PE**为位置嵌入向量。

$$\text{Input Embeddings} = \text{TE} + \text{SE} + \text{PE} \quad (12)$$

然而,尽管输入层中嵌入了词向量、段

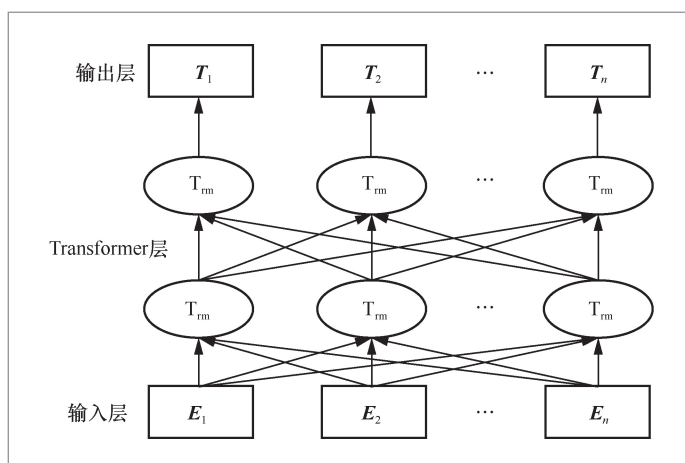


图3 BERT 结构

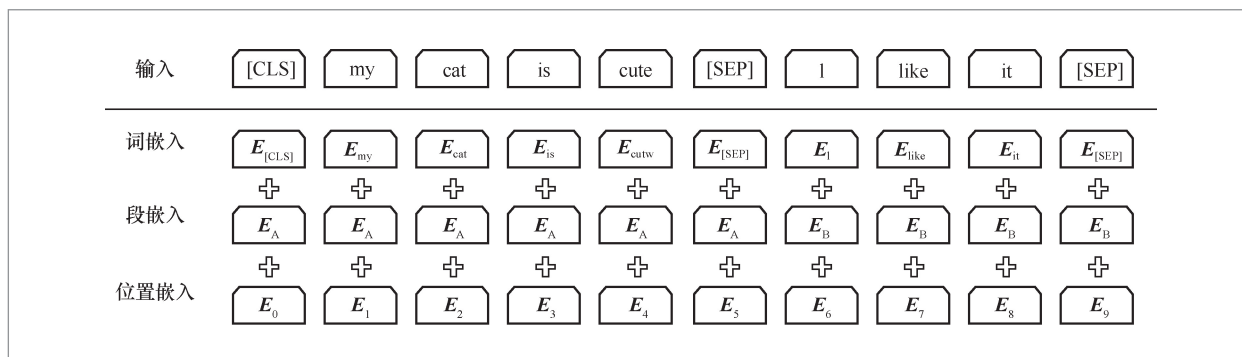


图4 BERT 输入层编码向量

向量和位置向量，能够获取一定的句法和语法信息，可以在Transformer层进行掩码预测，但是在面对情感分类任务时，由于缺乏情感特征，预测效果并不如其他的分类任务。这使在情感分析任务上的预测效果还有待优化和提升，尤其是在一些情感词显著的样本中，情感词的特征无法被充分提取，因而无法发挥价值。

表3展示了BERT对某两个样本预测时的结果，其中第一个文本中出现了“繁杂混乱”，第二个文本中出现了“危害”，这两个词都是具有明显情感倾向的情感词，但BERT在预测时反而给出了相反的结果，这说明BERT的多头自注意力（multi-head attention）机制并没有给予情感词更多的关注，导致包含一个乃至多个情感词的文本无法利用其具备情感标签的优势。这使情感任务的预测准确率并不理想。

事实上，当下游任务为情感任务时，在预训练阶段模型的注意力机制应该把更多的注意力放在情感词上，使模型可以更好地提取整个文本的情感特征。因此，本文提出一种情感词增强的注意力算法。

4.2 情感词增强的注意力机制

Transformer层利用了自注意力

表3 部分预测结果

文本	结果
移动电源已经成为使用智能手机用户常用的配件，但是繁杂混乱的产品却让用户不知如何选择	积极
抽烟的危害性众所周知，但仍无法做到有效的制止	积极

（self-attention）机制来帮助理解上下文语义。实际上，自注意力机制是一种分配机制，由于文本中的每个词都与句子中的其他词进行联系，因此，注意力机制会根据对象的重要程度以及与句子中其他词的关联性，重新分配权重。

基于自注意力的这一性质，本文利用自扩展的中日文情感词典改进了注意力的计算规则，提出了一种更注重情感词增强的注意力算法，以突出对象的情感特征。

Transformer的编码器结构如图5所示，Transformer由两个子层组成，分别是多头自注意力机制和前馈神经网络，每个子层后连接了一个规范化层及残差单元对输出进行控制，使向量的标准差和均值均为一个固定的数值。输入层的数据和多头自注意力层输出的结果进行残差相加后进行标准化，经过反馈层之后，再进行上述环节，最后输出结果。其中，多头自注意力机制是Transformer层的核心，输入层的3个箭头分别对应多头自注意力的3个输入

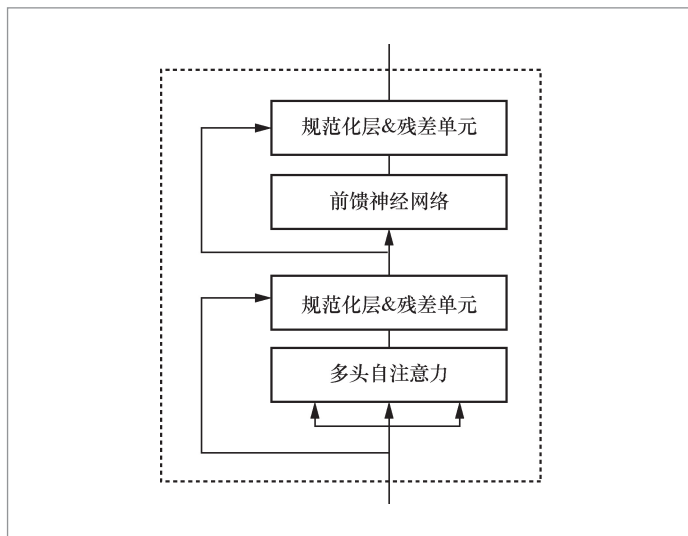


图5 Transformer 编码器结构

向量, 分别是 Q 、 K 、 V , 这3个向量是由输入层的词嵌入 X 和一个矩阵相乘得到的, 注意力的计算式如下:

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK}{\sqrt{d_k}}\right)V \quad (13)$$

其中, d_k 表示 K 的维度。

根据式(13)可以看出, 一个词的注意力在一定程度上会受到 Q 、 K 、 V 这3个向量大小的影响。而 Q 、 K 、 V 的值是通过原始的词嵌入得到的, 它相当于底层的特征信息。因此, 在情感分类任务中, 为了情感词能得到更大的关注, 本文提出一种情感词增强的注意力算法, 希望通过以下两个方面增强情感词的注意力。

(1) 增强情感词的词嵌入 X , 进而增大 Q 、 K 、 V 的值, 增强该词的注意力, 从而提高预测的准确率。将在第3节构建的目标领域的情感词典作为输入层的情感嵌入向量 $S'E$, 向模型中加入外部情感信息, 将4种嵌入特征相加后, 得到模型的输入向量如下:

$$\text{Input Embeddings} = \text{TE} + \text{SE} + \text{PE} + S'E \quad (14)$$

(2) 将情感词典中情感词的分值赋给模型作为额外的权重。一般来说, 情感词表达的强度越强烈, 越能代表其所在句子和文本的情感倾向, 越应该得到更大的注意力权重。例如, 对于“这家餐厅的服务太差劲了, 菜再好吃我也不来了”这一样本, 由于“差劲”的情感强度较强, 因此很大概率是整个文本的情感倾向, 应该给予更多的关注。而情感强度相对平和的情感词, 虽然本身具备情感倾向, 但是相比于情感强度大的词, 上下文转折的可能性会更大。例如, 对于“这家餐厅的服务不怎么样, 考虑到菜品做得太美味了, 下次还来”这一样本, 应该把注意力更多放在“美味”上, 而非“不怎么样”上。因此, 本文重新定义了针对情感任务的自注意力算法, 计算式如下:

$$\text{Attention} = \text{Softmax}\left(\frac{QK}{\sqrt{d_k}}\right) \times (1 + \lambda)V \quad (15)$$

其中, λ 为对情感词的分值归一化处理后的(0,1)之间的数值, 给BERT的注意力机制量化加权。

4.3 结合情感词典的预训练模型

本文利用自扩展的情感词典提出了一个新的情感分析模型——EmoBERT, 模型由3个部分构成, 如图6所示。

由于本文改进了注意力机制, 因此, 在BERT基础上, 进行了进一步的预训练, 预训练任务如下。

(1) 输入层。输入层中每个词额外加入了情感向量, 如图7中情感嵌入所示, 情感向量与其他3个嵌入向量相加, 组成了具备情感特征的词嵌入。

(2) Transformer层。在这一层, 根据本文提出的情感词增强的注意力机制, 可以

求出每个子空间的注意力分值,进而计算出多个子空间的输出结果,计算式如下:

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (16)$$

$$\begin{aligned} & \text{MultiHead}(Q, K, V) = \\ & \text{Linear}(W_i \text{concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n) + b) \end{aligned} \quad (17)$$

其中, b 表示误差项。

(3) 输出层。输入层的词嵌入 X 和 Transformer 层输出的结果残差相加,再进行标准化,经过反馈层之后,重复上述过程,最后输出结果。子层最后得到的输出结果如式 (18) 所示:

$$\begin{aligned} & \text{SubLayer}_{\text{output}} = \\ & \text{LayerNorm}(X + \text{SubLayer}(X)) \end{aligned} \quad (18)$$

由于BERT是在通用数据集上训练的,在本特定领域的任务上,原始BERT无法完全抽取出词元的内在含义。因此,需要用本领域语料对其进行微调。此外,在微调前使用目标领域数据集的数据对模型进一步预训练,相当于在预训练阶段实现将模型从通用领域向特定领域提前迁移,然后再执行目标领域的任务^[22],如图8所示。

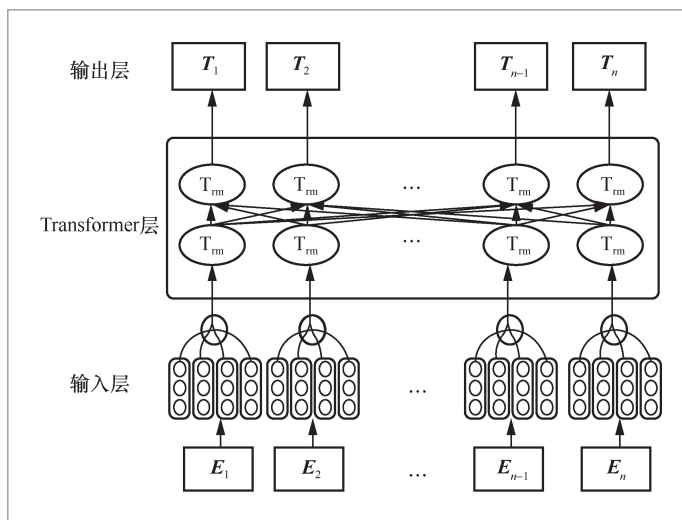


图6 EmoBERT 模型结构

5 实验结果

为了验证本文提出的EmoBERT模型在情感分类任务上的表现,基于自采集的微博和Twitter数据集进行试验。

5.1 数据集

本文将自采集的微博数据集和Twitter数据集作为中日舆情研究的数据

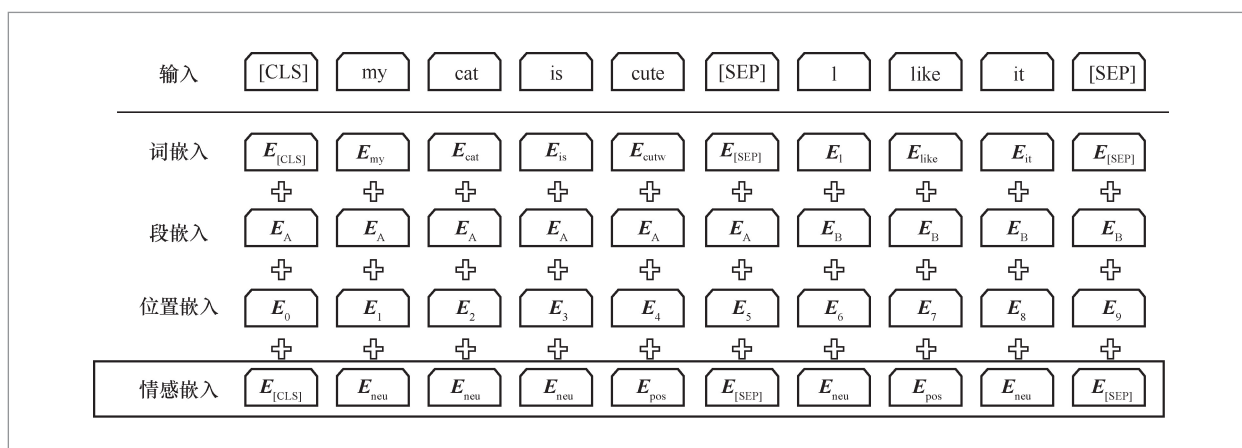


图7 EmoBERT 输入特征向量

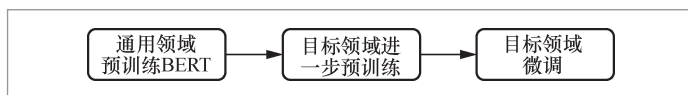


图8 模型流程

集,类别包含娱乐、体育、经济、科技和疫情,具体见表4。

由于自采集的微博数据集和Twitter数据集数据量庞大,并且不包含情感类别标签,因此分别从每个类别提取1 000条数据,并人工标注其情感极性。在加载好实验所需要的数据集之后,分别按照8:2的比例划分这两个数据集的训练集和测试集。表5展示了各个数据集的训练集和测试集包含的情感极性信息。

5.2 超参数设置

预训练阶段,批量大小(batch size)设置为256,学习率(learning rate)为 5×10^{-5} ,持续训练1 000 000步。模型中所有丢失率(dropout)都为0.1,参数Adam β_1 和Adam β_2 的值分别为0.9和0.999,L2权重衰减为0.01。

在微调过程中,对批量大小、学习率和训练周期数量等进行一定调整。其中批量大小的值是每次迭代(epoch)训练的句子数。如果设置过小,会使训练时间延长;如果设置过大,损失函数曲线比较平坦时,将无法得到最优模型。不同的下游任务对应不同的最佳超参数值,为了让模型得到最佳

表4 数据集

类别	微博/条	Twitter/条
娱乐	40 768	8 603
体育	19 843	8 788
经济	16 808	6 266
科技	35 538	8 746
疫情	7 520	3 198

分类效果,本文经多组实验验证最优的参数值,具体见表6。

5.3 实验结果及模型对比

为了评估本文提出的EmoBERT模型在情感二分类任务上的分类效果,针对自采集的中日文舆情数据集分别进行了实验,同时与BERT及领域迁移后的EmoBERT进行对比,实验过程中的准确率变化如图9和图10所示。其中,EmoBERT-C是中文模型,EmoBERT-J是日文模型,BERT*为在目标领域数据集上预训练之后的模型。

从图9和图10中可以清晰地看出准确率的对比情况,与原始BERT相比,EmoBERT-C和EmoBERT-J都有一定的提升。其中EmoBERT-C相比于原始BERT模型,准确率从89.68%提升到92.15%。EmoBERT-J相比于原始BERT模型,准确率从74.73%提升到78.26%。

此外,针对中日文舆情数据集中的积极文本和消极文本也分别进行了实验,同时与BERT进行了对比,实验结果见表7。从表7中可以看出,本文提出的EmoBERT模型在积极文本和消极文本中分类的效果均优于BERT,进一步验证了模型的有效性。表7中BERT-C表示中文模型,BERT-J表示日文模型。

5.4 中日舆情分析——新冠肺炎疫情

新冠肺炎疫情(以下简称疫情)带来的影响是多元复杂的。病毒的攻击具有无差别性、跨国性和极大的不确定性,疫情给全球的经济和金融市场造成了剧烈的冲击,并且在很大程度上催化了国际关系的演变。因此,利用该数据集分析相关的中日舆情是很有意义的。

在实验中,本文以“新冠肺炎”为关键

词,采集了2020年1月至2020年9月的舆情数据。使用本文提出的模型对数据进行情感分析,结果如图11和图12所示。图11中横坐标为时间,左侧纵坐标为发文量,右侧纵坐标为情感分值。从图11和图12中可以看到,两国民众的情感态度普遍是积极的。其中2月和3月两平台发文量非常大,两平台对应时间的关键词见表8和表9。通过分析可知,2月和3月是春节期间,人口流动量巨大且新冠疫苗还未研制成功,正是疫情的暴发期,也是民众讨论最多的时期。

6 系统构建及数据可视化

整体分析中日民众情感检测系统需求,将该系统功能分成三大模块,分别为:用户交互模块、情感分析模块以及可视化模块。

6.1 用户交互模块

(1) 实时采集数据并分析

在系统功能界面提供用户输入数据采集条件的接口,为了实时监测社交平台的数据,本系统服务端连接微博和Twitter等社交平台的API,以响应用户输入的关键词、详细程度、起始时间、结束时间等条件,从平台实时采集数据信息,并以数据流的方式将信息传输到服务器端,解析获得合法可分析的数据集。

(2) 接收用户上传的数据并分析

为了满足相关研究领域的专业用户,本系统支持用户自行上传数据集。考虑服务端情感分析算法的计算方式,数据集应包含每条社交平台微博/推文的具体内容、点赞数(喜欢数)以及发布时间,所上传的数据集最终以表单形式将数据集传给服务器端。

表5 数据集情感极性信息

数据集	类别	积极/条	消极/条	合计/条
微博	训练集	1 894	2 106	4 000
	测试集	473	527	1 000
Twitter	训练集	1 768	2 232	4 000
	测试集	442	558	1 000

表6 最优参数值

参数名	值
batch size	64
learning rate	5×10^{-5}
epoch	3

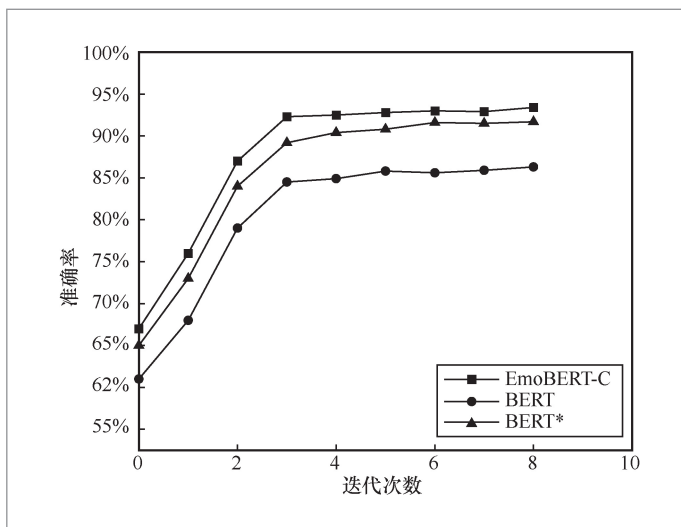


图9 EmoBERT-C及原始模型准确率

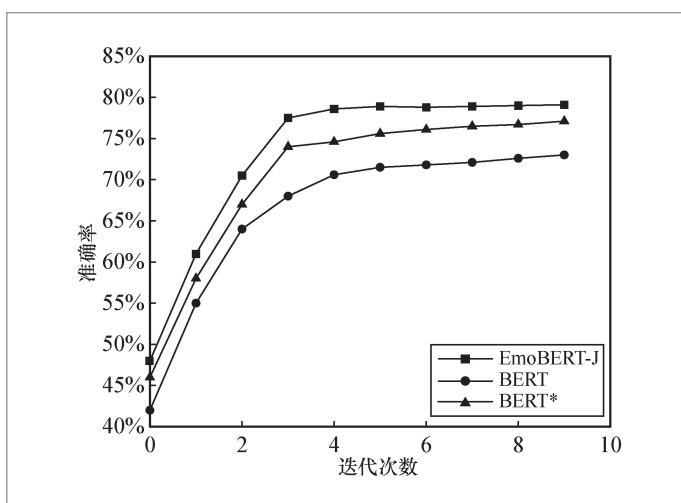


图10 EmoBERT-J及原始模型准确率

表7 针对积极文本和消极文本对比实验结果

模型	情感态度	准确率	F1值
BERT-C	积极	89.81%	83.26%
	消极	89.55%	82.83%
EmoBERT-C	积极	91.83%	83.15%
	消极	92.47%	84.02%
BERT-J	积极	75.92%	70.34%
	消极	77.04%	71.28%
EmoBERT-J	积极	77.92%	71.53%
	消极	78.60%	72.96%

6.2 情感分析模块

服务器端作为自然语言处理的主要执行端,以Python为主要语言,通过封装好的情感分析算法,依赖系统预先保存的中日文情感词典、程度词表、否定词表等情

感资源,分析用户提交或实时采集的数据集,本系统所使用算法将从3个方面进行情感分析。第一,根据用户指定的关键词和数据集计算中日民众的情感分值随时间的变化情况,并实现实时监测;第二,对于每个情感突变点都能分析提取该时间区间热度排名前10的微博和推文,方便分析中日民众情感分歧较大的起因事件;第三,根据算法分析结果,给出中日民众对于该主题的情感极性分布情况。三者最终通过Jinja2模板引擎中的模版函数动态渲染前端HTML模版文件中的JavaScript脚本变量,再通过JSON解析生成可用于展示的格式,最终传递到展示界面,动态生成用户需要的分析结果。

6.3 可视化模块

本系统前端展示界面基于CSS、HTML、JavaScript脚本完成开发。展示模块包括3个

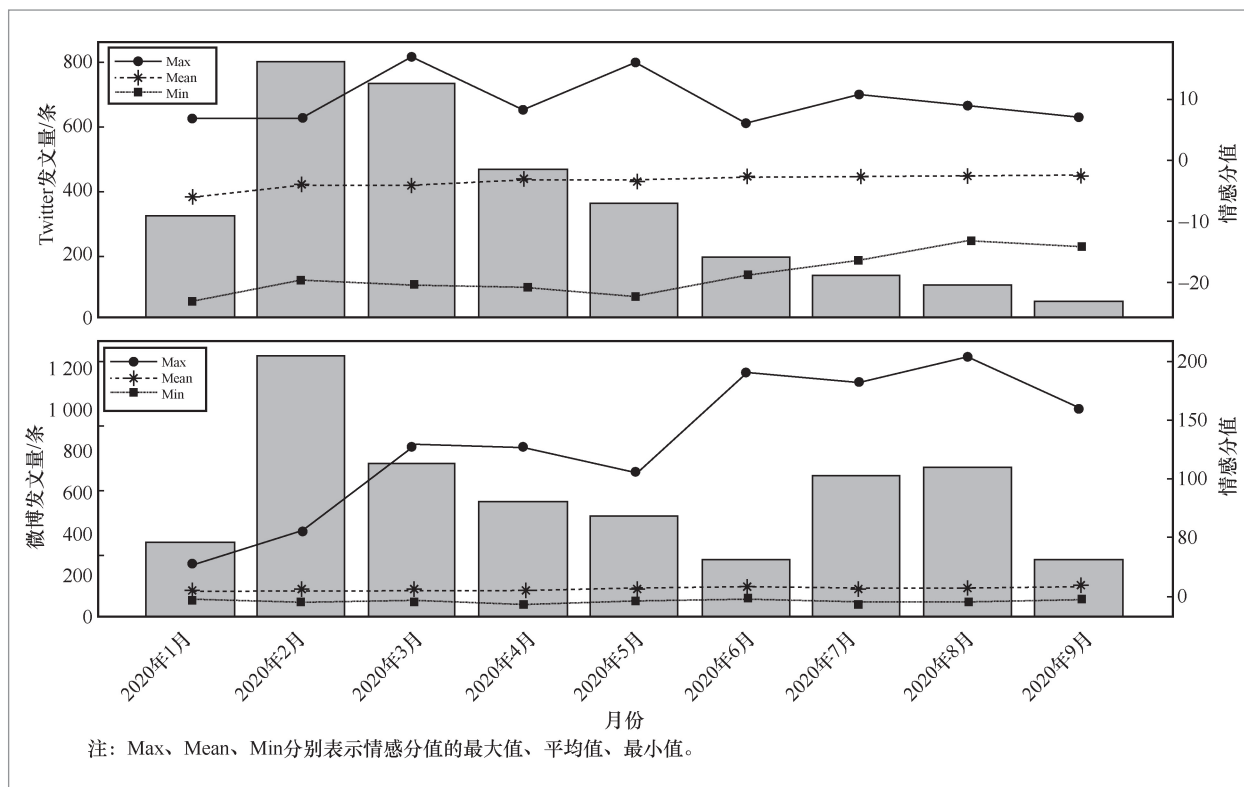


图11 中日民众情感变化

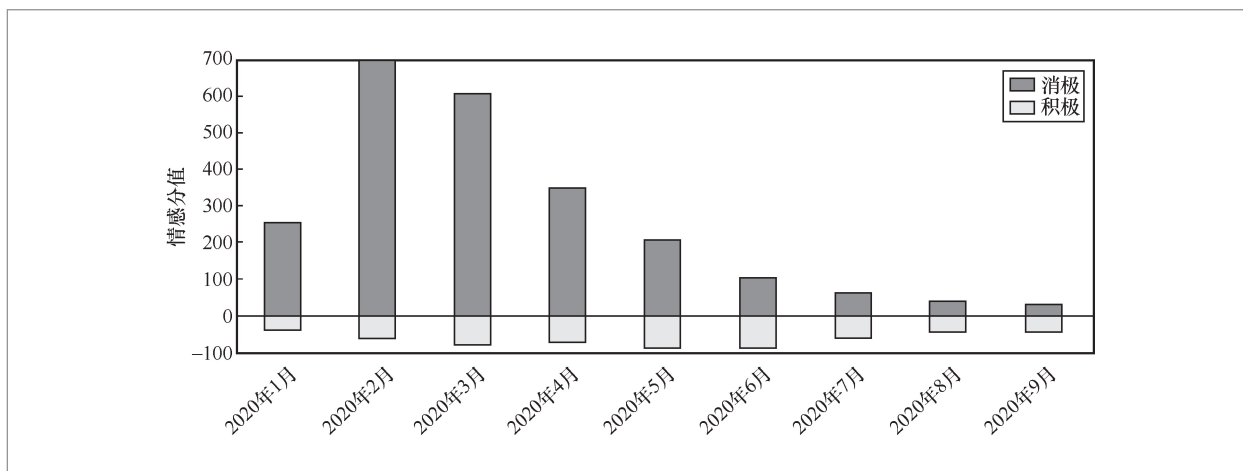


图 12 中日民众情感极性分布

部分，均使用可交互的动态曲线，其中第一部分是微博和Twitter单独的可交互数据展示，包括以日/月为单位的情感分值随时间变化曲线图，以日/月为单位的情感极性随时间变化曲线图，热度前100文章情感散点图，内容来源年份分布扇形图以及词云。第二部分是微博和Twitter的情感随时间变化对比图，每一数值点还支持点击查看影响该点的热度排名前5的微博和推文。第三部分是根据用户输入对微博某一特定关键词进行实时检测，显示实时变化的情感分值随时间变化曲线图。

6.4 系统架构

系统主体基于Web的C/S架构设计，以Flask为开发框架，其负责将HTML文件、JavaScript脚本与后端Python代码进行连接，同时Jinja2模板引擎和Werkzeug WSGI套件负责前后端的数据传递。本系统从规模上来说较为轻量，同时Flask框架高度的灵活性也降低了后续系统功能拓展和维护的成本。

客户端基于Web技术设计，一方面负责处理业务逻辑，另一方面负责返回响应内容。业务逻辑方面，支持用户在页面中输

表 8 微博“新冠肺炎”关键词

时间	关键词
2020年2月	肺炎, 新冠, 疫情, 防控, 患者, 确诊, 工作, 病例, 医院, 武汉, 感染, 新型, 隔离, 冠状病毒
2020年3月	肺炎, 新冠, 疫情, 确诊, 病例, 防控, 视频, 患者, 美国, 累计, 隔离, 医院, 新增, 感染

表 9 Twitter“新冠肺炎”关键词

时间	关键词
2020年2月	新型, ウイルス, コロナ, 感染, 肺炎, 扩大, 武汉, 死者, 对策, 対応
2020年3月	コロナ, 新型, ウイルス, 感染, イタリア, 武汉, 扩大, 肺炎, 死者

入主题关键词、详细程度、起始时间、结束时间等条件实时采集数据集，或接收用户自行上传的本地数据集文件，将主题关键词字符串或合法数据集文件通过表单传递给服务器端进行解析处理。返回响应内容方面，在展示结果页面通过Jinja2模板引擎渲染模版，与服务端进行数据传递，接收服务器端传来的计算结果，客户端展示结果页面使用可视化工具包Chart.js，生成可交互式情感分值随时间变化曲线图和情感极性分布图等，并将分析的结果展示给用户。

服务器端以Flask框架为主体进行开

发,通过不同URL来区分响应前端视图函数的不同表单请求,完成业务逻辑的具体功能实现。本系统主要包含3个不同URL绑定的功能函数。monitor()负责调用社交平台API,即时采集流式数据,并保存到服务端。等待系统进行进一步分析处理。uploader()负责将用户本地数据集上传到服务器端。此外为了防止数据集包含内容缺失或不兼容本算法的运算所需条件,在用户上传数据集后,首先进行合法性检测,要求数据集具有完整性,这保证了系统的稳定性。analyse()负责调用本地情感词典、程度词表,并通过本地情感分析算法计算分析已经上传到服务器端的数据,计算结果通过Jinja2模版引擎渲染动态前端系统客户端的HTML模板文件中的内嵌JavaScript变量,进一步进行JSON解析,并使用eval()进行合法性检测。同时系统使用Flask-Caching扩展的缓存技术提高程序运行速度。

6.5 可视化实现

为了更清晰地展示系统的有效性,本文基于自采集的2018年以“RNG”为主题的中日文数据集对系统的功能进行演示。RNG是当下深受国内外游戏爱好者喜爱的电竞战队。传入数据集后,后台算法开始分析,并将结果反馈给前端可视化界面,图13和图14分别为微博和Twitter单一平台对RNG主题的分析结果,包括情感极性分布、年份分布、热度情况分布、词云、每日情感变化图、每月情感变化图、每日发文量变化图以及每月发文量变化图等,图上的点均可以点击,用于和用户交互,可以展示当前点的详细信息,如情感分值、积极人数、中立人数、消极人数等。

从图13可以看出,中国民众对于“RNG”战队的讨论度很高,数据量达到了38 272条。通过分析情感变化曲线可以看

出,共有5个情感波动较为明显的时间段。在4月28日、5月20日、8月29日以及9月14日这4日前后,中国民众对该主题的情感非常积极,同时通过柱形图也可以看出对应时间的发文量很高,通过点击曲线上对应的点可以看到当天的热门微博。

从图14可以看出,日本民众对“RNG”主题的讨论较少,一年中共爬取到7 437条相关数据。同时,通过分析情感变化曲线可以看出,日本民众对该主题的态度普遍比较积极且比较平稳。图15展示了微博和Twitter两平台中日双语料的对比分析结果,针对中日民众情感分歧较大的点,同样可以通过点击对应的点查看当日两平台关于该主题热度最高的博文,便于用户进一步了解情感分歧的原因事件。

此外,本系统还提供话题情感的实时监测功能,针对用户输入的关键词进行实时监测、实时分析,并实时呈现给用户,如图16所示。

7 结束语

微博、Twitter等社交网络平台的流行使其中蕴含了丰富的情感信息,通过对这些平台上用户发布的内容进行情感分析,可以挖掘其中的社会价值。本文采用了基于自扩展的情感词典结合改进的BERT预训练模型进行了实验,建立的系统可以同时分析中日舆情数据,并自动生成中日民众情感态度对比和情感极性分布,针对情感突变点也能通过分析热点博文和推文来有效地分析出相应的事件,并将分析结果以可视化的形式展现给用户。因此,本文系统是合理的、有效的,且弥补了目前单一平台、单一语料舆情监测的缺陷。未来工作将从以下方面进行。

- 社交平台用户的情绪比较丰富,应从多个方面分析情感词,不能局限于积极和消极的二分类,应进一步延伸对情感

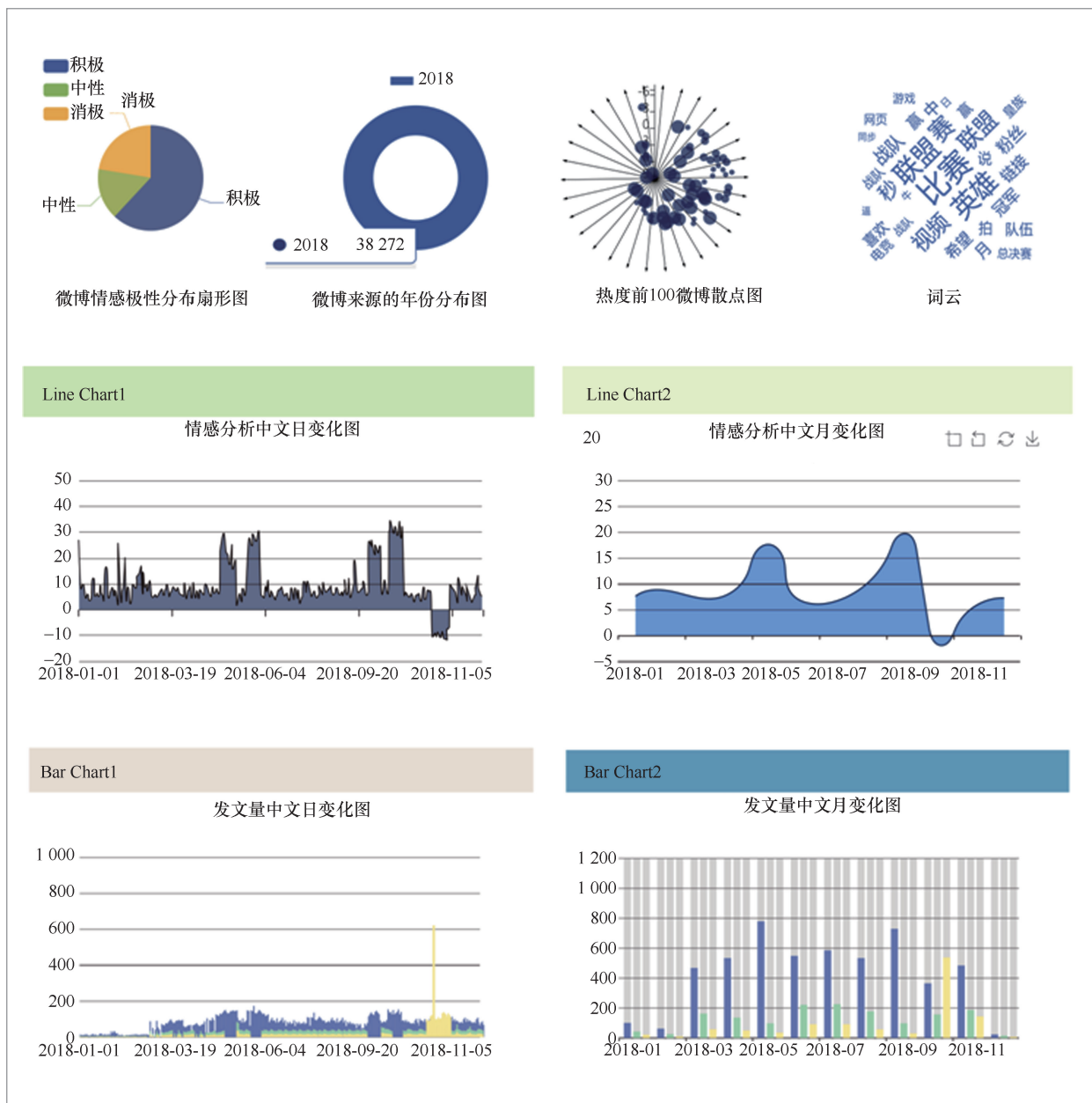


图 13 微博可视化展示界面

情绪的等级判断。

- 舆情监测系统也不应只局限于中文和日文两种语料, 未来应支持更多种语料。

- 只对文本进行了情感分析, 但现实世界中除文本外, 图片、视频、语音等信息也会包含强烈的情感倾向, 系统应实现多模态的情感分析。

参考文献:

- [1] 敦欣卉, 张云秋, 杨铠西. 基于微博的细粒度情感分析[J]. 数据分析与知识发现, 2017, 1(7): 61-72.
DUN X H, ZHANG Y Q, YANG K X. Fine-grained sentiment analysis based on

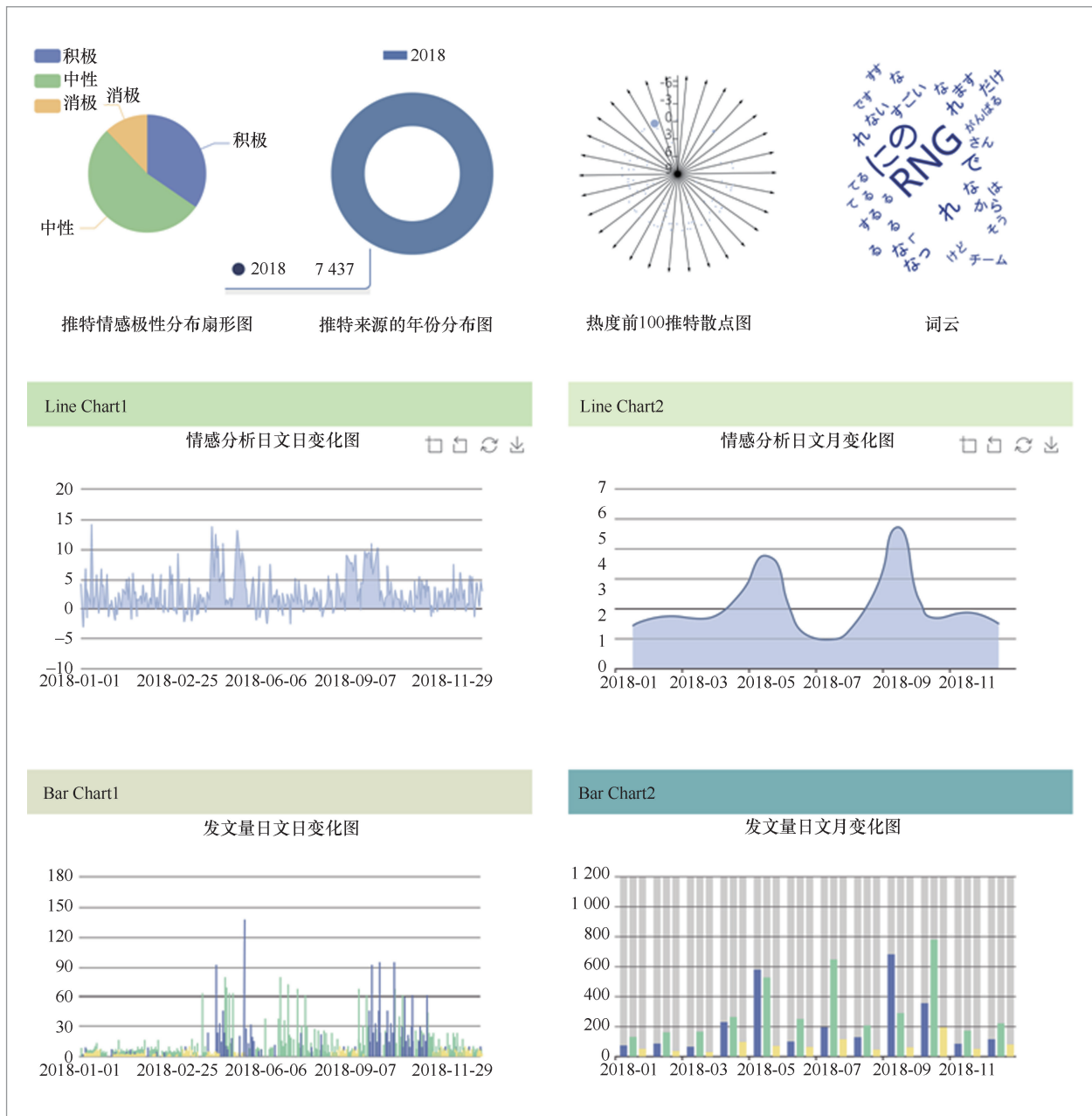


图 14 Twitter 可视化展示界面

weibo[J]. Data Analysis and Knowledge Discovery, 2017, 1(7): 61-72.

- [2] ZHAO J C, DONG L, WU J J, et al. MoodLens: an emoticon-based sentiment analysis system for Chinese tweets[C]// Proceedings of the 18th ACM SIGKDD international conference on Knowledge

discovery and data mining - KDD'12. New York: ACM Press, 2012.

- [3] WANG H, CAN D, KAZEMZADEH A, et al. A system for real-time Twitter sentiment analysis of 2012 U.S. presidential election cycle[C]// Proceedings of the ACL System Demonstrations. [S.l.:s.n.], 2012.

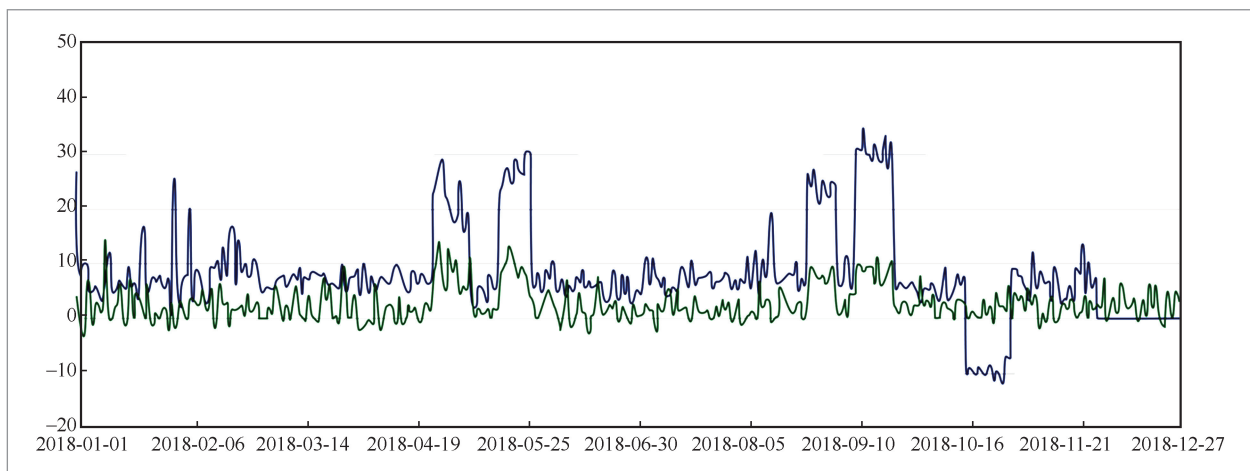


图 15 两平台对比结果可视化展示界面

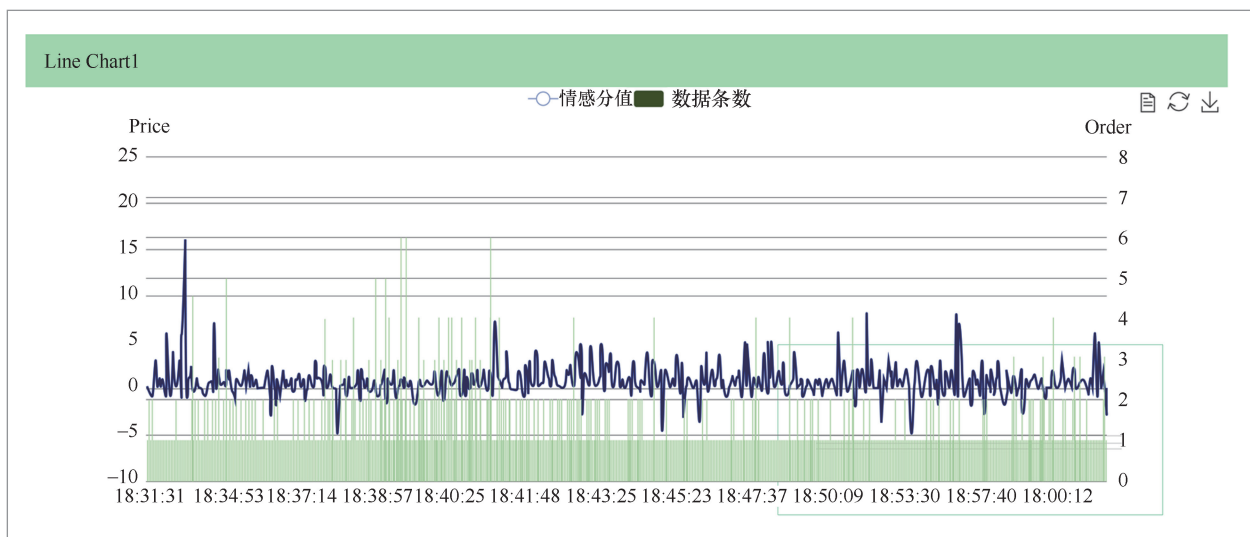


图 16 话题情感实时监测界面

- [4] WILLIAMS J, KATZ G. Extracting and modeling durations for habits and events from Twitter[C]//Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Short Papers – Volume 2. [S.l.:s.n.], 2012: 223–227.
- [5] 李忠俊. 基于话题检测与聚类的内部舆情监测系统[J]. 计算机科学, 2012, 39(12): 237–240.
LI Z J. Internal public opinions monitor system based on topic detection and clustering[J]. Computer Science, 2012, 39(12): 237–240.
- [6] YI J, NASUKAWA T, BUNESCU R, et al. Sentiment analyzer: extracting sentiments about a given topic using natural language processing techniques[C]//Proceedings of 3rd IEEE International Conference on Data Mining. Piscataway: IEEE Press, 2003: 427–434.
- [7] RILOFF E M, SHEPHERD J. A corpus-based approach for building semantic lexicons[J]. arXiv preprint, 1997, arXiv:cmp-lg/9706013.
- [8] 熊德兰, 程菊明, 田胜利. 基于HowNet的句

- 子褒贬倾向性研究[J]. 计算机工程与应用, 2008, 44(22): 143-145.
- XIONG D L, CHENG J M, TIAN S L. Sentence orientation research based on HowNet[J]. Computer Engineering and Applications, 2008, 44(22): 143-145.
- [9] 潘明慧, 牛耘. 基于多线索混合词典的微博情绪识别[J]. 计算机技术与发展, 2014, 24(9): 28-32, 36.
- PAN M H, NIU Y. Emotion recognition of micro-blogs based on a hybrid lexicon[J]. Computer Technology and Development, 2014, 24(9): 28-32, 36.
- [10] PANG B, LEE L, VAITHYANATHAN S. Thumbs up? Sentiment classification using machine learning techniques[C]// Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - EMNLP'02. Morristown: Association for Computational Linguistics, 2002.
- [11] 杨艳霞. 基于分类的微博情感分析算法研究及实现[J]. 计算机与数字工程, 2017, 45(2): 197-200, 396.
- YANG Y X. Microblog sentiment analysis algorithm research and implementation based on classification[J]. Computer & Digital Engineering, 2017, 45(2): 197-200, 396.
- [12] COLLOBERT R, WESTON J, BOTTOU L, et al. Natural language processing (almost) from scratch[J]. The Journal of Machine Learning Research, 2011, 12: 2493-2537.
- [13] AL-RIFAIE M M, BISHOP J. Swarmic sketches and attention mechanism[C]// Proceedings of the International Conference on Evolutionary and Biologically Inspired Music and Art. [S.l.:s.n.], 2013: 85-96.
- [14] SCHULLER B, MOUSA A E D, VRYNIOTIS V. Sentiment analysis and opinion mining: on optimal parameters and performances[J]. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2015, 5(5): 255-263.
- [15] 宋婷, 陈战伟, 杨海峰. 基于分层注意力网络的方面情感分析[J]. 大数据, 2020, 6(5): 82-91.
- SONG T, CHEN Z W, YANG H F. Aspect sentiment analysis based on a hierarchical attention network[J]. Big Data Research, 2020, 6(5): 82-91.
- [16] 徐志栋, 陈炳阳, 王晓, 等. 基于胶囊网络的方面级情感分类研究[J]. 智能科学与技术学报, 2020, 2(3): 284-292.
- XU Z D, CHEN B Y, WANG X, et al. Research on capsule network-based for aspect-level sentiment classification[J]. Chinese Journal of Intelligent Science and Technology, 2020, 2(3): 284-292.
- [17] 张宝华, 张华平, 厉铁帅, 等. 基于多输入模型及句法结构的中文评论情感分析方法[J]. 大数据, 2021, 7(6): 41-52.
- ZHANG B H, ZHANG H P, LI T S, et al. Chinese comment sentiment analysis method based on multi-input model and syntactic structure[J]. Big Data Research, 2021, 7(6): 41-52.
- [18] RADFORD A, NARASIMHAN K. Improving language understanding by generative pre-training[J]. Preprint-Work in Progress, 2018.
- [19] DEVLIN J, CHANG M W, LEE K, et al. BERT: pre-training of deep bidirectional transformers for language understanding[J]. arXiv preprint, 2018, arXiv:1810.04805.
- [20] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need[J]. arXiv preprint, 2017, arXiv:1706.03762
- [21] SUN C, QIU X P, XU Y G, et al. How to fine-tune BERT for text classification?[C]//Proceedings of Chinese Computational Linguistics. [S.l.:s.n.], 2019.
- [22] 杨晨, 宋晓宁, 宋威. SentiBERT: 结合情感

信息的预训练语言模型[J]. 计算机科学与探索, 2020, 14(9): 1563-1570.

YANG C, SONG X N, SONG W.
SentiBERT: pre-training language model

combining sentiment information[J].

Journal of Frontiers of Computer
Science and Technology, 2020, 14(9):
1563-1570.

作者简介



李爱黎(1995-),女,东北大学计算机科学与工程学院硕士生,主要研究方向为情感分析、数据挖掘。



张子帅(2000-),男,东北大学计算机科学与工程学院本科生,主要研究方向为数据挖掘、机器学习。



林荫(1984-),女,东北大学外国语学院讲师,主要研究方向为中日文化比较研究。



王秋菊(1962-),女,东北大学外国语学院教授,主要研究方向为中日文化比较研究、科技与文化研究。



杨建安(2002-),男,东北大学计算机科学与工程学院本科生,主要研究方向为数据挖掘、机器学习。



孟炜程 (2002-), 男, 东北大学计算机科学与工程学院本科生, 主要研究方向为数据挖掘、机器学习。



张岩峰 (1982-), 男, 东北大学计算机科学与工程学院教授, 中国计算机学会高级会员, 主要研究方向为大数据挖掘、大规模机器学习、分布式系统。

收稿日期: 2022-03-08

通信作者: 王秋菊, wangqiuju@fsc.neu.edu.cn

基金项目: 国家自然科学基金资助项目 (No.62072082); 辽宁省重点研发计划(No.2020JH2/10100037); 中央高校基本科研业务费 (No. N2216015)

Foundation Items: The National Natural Science Foundation of China (No.62072082), The National Key Research and Development Program of Liaoning Province (No.2020JH2/10100037), Fundamental Research Funds for the Central Universities (No.N2216015)