

马律倩. 基于 BERT-Attention-BiLSTM 的多特征主题情感分析[J]. 智能计算机与应用, 2024, 14(5): 205-208. DOI: 10.20169/j.issn.2095-2163.240528

基于 BERT-Attention-BiLSTM 的多特征主题情感分析

马律倩

(浙江理工大学 计算机科学与技术学院, 杭州 310018)

摘要: 关注微博用户对于事件的情感倾向, 有利于平台了解用户心声, 也能为决策者的舆情处理工作提供参考和方向。然而, 当前大部分微博情感分析研究仍是基于文本的, 忽略了表情、图片等要素。针对上述问题, 本文提出了一个多模型融合的情感分析模型, 以 BERT 预训练模型为基础, 融合情感词典, 并采用双向 LSTM 获取文本特征, 有效联系前后文, 并引入注意力机制, 同时提出了一种 emoji 表情特征计算方法, 得到一个情感分类更准确的多特征主题情感分析模型。

关键词: 情感分析; 注意力机制; 预训练模型; 深度学习

中图分类号: TP391.1

文献标志码: A

文章编号: 2095-2163(2024)05-0205-04

Multi-feature topic sentiment analysis based on BERT-Attention-BiLSTM

MA Lüqian

(School of Computer Science and Technology, Zhejiang Sci-Tech University, Hangzhou 310018, China)

Abstract: Currently, paying attention to Weibo users' emotional tendencies towards events will help the platform understand users' voices, and can also provide reference and direction for decision makers in handling public opinion. However, most of the current microblog sentiment analysis research is still based on text, ignoring elements such as expressions and pictures. In response to the above problems, this paper proposes a multi-model fusion sentiment analysis model, which integrates sentiment lexicon based on the BERT model, and uses two-way LSTM to obtain text features, effectively connects context, and introduces an attention mechanism. At the same time, a calculation method of emoji expression features is proposed, and a multi-feature topic sentiment analysis model with more accurate classification is proposed.

Key words: sentiment analysis; attention mechanism; pre-training model; deep learning

0 引言

当前,随着互联网的迅猛发展,人们日常的学习、工作和生活与互联网紧密相连,各类社交平台也成为分享生活日常、表达观点、普及知识的重要渠道,而由平台发布的各类文本构成了互联网信息的一部分。各式各样的信息和多样的传播途径可以让用户更加方便地获取需要的信息,也能得到其他领域如金融、人文、艺术等信息,满足多元化的精神需求^[1]。丰富的互联网商品、服务也扩展了网络在生活中的应用,如消费者在各电商平台购物后,通常会发表对商品的评价;而在众多外卖平台点餐后,会对菜品的口味以及商家的服务进行评论、打分。对于这些呈快速增长态势的文本信息,仅靠人工收集整理已经绝无可能完成,因此就需要借助计算机识别出其中隐藏的商业、及社会价值信息,而据此通过情感分析来挖掘用户在

文本中所表达情感的任务也随即应运而生。

情感分析又称为观点挖掘,是当前自然语言处理领域的一个重要分支。文本情感分析是指用数据挖掘算法,对各种社交平台信息以及带有情感表达的文本自动进行情感分类,通过结合上下文信息分析得出指代的积极/消极态度。通过文本情感分析,可以分析出公众对某一事件或话题的看法,为日常决策提供参考。在其他领域,该研究也具有十分重要的现实意义,商家能够为用户提供更有针对性的服务,从而做出更加个性化的服务推荐等。

当前,文本情感分析主要有基于情感词典^[2-3]、其中,基于机器学习^[4-5]和基于深度学习^[6]三种方法。基于深度学习的方法主要应用于情感分析的2个方面:

(1)用于文本建模。Bengio 等学者^[7]于 2003 年首次提出了神经网络语言模型,先将第一层输出的矩

阵作为文本表示,再通过特征距离计算其相似度。Mikolov 等学者^[8]于 2013 年提出了 Word2Vec 语言模型,该模型能够结合上下文信息,并根据需要的输入输出类型,选用不同模型结构,训练得到更准确的词向量表达。

(2)用深度学习的方法进行情感分析。常用的深度学习模型主要包括卷积神经网络模型(CNN)、循环神经网络模型(RNN)等。CNN 采用分层特征表示,具有很强的局部上下文特征提取能力^[9]。Hochreiter 等学者^[10]在 RNN 的基础上,提出了长短期记忆神经网络(LSTM),解决了 RNN 天然存在的梯度爆炸问题。Envelope 等学者^[11]采用双向 LSTM,同时考虑前后文对情感分类的影响,并融合分类元数据得到了自定义的分类器。关鹏飞等学者^[12]在 BiLSTM 模型中引入注意力机制,为情感表达相关的词向量分配更多的权重。

在词向量模型的选择上,现有的深度学习方法大多采用 Word2vec、GloVe 等静态词向量模型,BERT、

ELMO 等预训练模型的提出展现出了动态词向量模型对深度学习方法的优化与改进^[13-14]。BERT 预训练模型可以动态地表示上下文信息,解决一词多义的问题,BiLSTM 可以双向获取文本特征^[15-16]。在 2021 年,Li 等学者^[17]基于 BERT 和双向 LSTM 对公开的疫情评论数据进行情感分析,得到了较好表现。从中可以看出预训练语言模型在情感分析任务上的有效性。注意力机制^[18]可以对不同的词向量给予不同的关注度,因此本文提出了基于 BERT-Attention-BiLSTM 的多特征主题情感分析模型 MFTEA,在一定程度上提高了微博主题情感分析的准确率。

1 模型结构

MFTEA 模型框架如图 1 所示,主要由输入层、Senti-BERT 层、BiLSTM 层、语义合成层和情感分类层五个部分构成,各个层之间相互连接,上一层的输出作为下一层的输入。

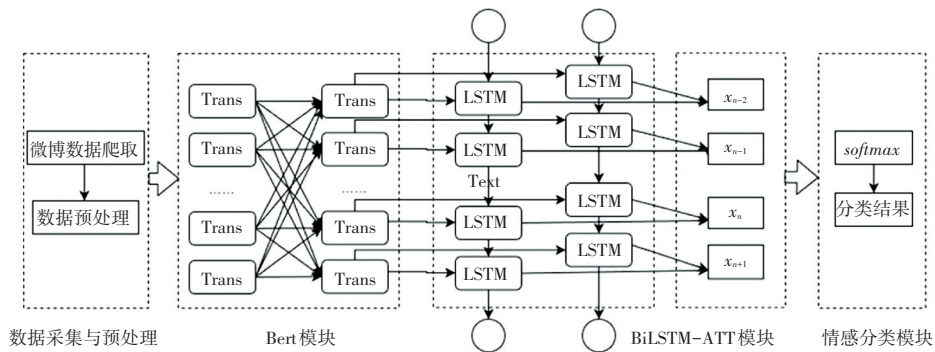


图 1 MFTEA 模型框架图

Fig. 1 Model framework diagram of MFTEA

1.1 输入层

输入层主要用于将采集到的新浪微博原始数据进行预处理后送入 Senti-BERT 模型。预处理操作包括数据清洗、中文分词和去除停用词等。

本文以新浪微博为信息获取平台,通过限定时间范围内的某一话题搜索,获取相关的微博内容作为数据来源,分析新浪微博相关网络舆情。通过设计基于 Scrapy 框架的分布式网络爬虫系统^[19],实现相关原始微博文本的抓取和数据的非结构化存储。数据清洗时,将常用 emoji 表情转化为相应释义,并采用 jieba 分词^[20]工具完成文本分词任务,采用停用词表^[21]去除没有实际意义的字词、数字、标点符号等。

1.2 Senti-BERT 层

本文在 BERT 模型的预训练阶段融入知网情感词典,记作 Senti-BERT 模型。该模型拥有强大的语义表示能力,能够识别出评论过的情感词,获取更丰

富的输入特征,提高模型情感分析的准确率。

Senti-BERT 模型与原始 BERT 模型的主要不同为:BERT 模型采用的方式是将原始句输入模型进行预训练,而 Senti-BERT 模型首先基于情感词典识别出评论句中的情感词、观点词,依照在原句中出现的顺序,将这些词拼接在句首,输入模型进行预训练。

1.3 BiLSTM 层

双向长短时记忆网络以合理有效的方式将前向 LSTM 与后向 LSTM 合并,不存在梯度消失问题,且能有效联系文本前后文,获得更准确的句子间依赖关系。在这一层将 Senti-BERT 层输出的文本向量 $(x_1, x_2, x_3, \dots, x_i)$ 送入前后 2 个方向的 LSTM 层进行特征提取。通过前向 LSTM 得到一组特征向量 $(\vec{h}_1, \vec{h}_2, \vec{h}_3, \dots, \vec{h}_i)$,通过后向 LSTM 得到一组特征向量 $(\vec{h}_1, \vec{h}_2, \vec{h}_3, \dots, \vec{h}_i)$,最后将 2 个 LSTM 的输出合并为一个输出向量。计算公式如下:

$$\mathbf{h}_t = \tanh(\mathbf{w}_t \mathbf{x}_t + \mathbf{u}_{h_{t-1}} + \mathbf{b}_t) \quad (1)$$

$$\vec{H} = LSTM(\vec{h}_1, \vec{h}_2, \vec{h}_3, \dots, \vec{h}_t) \quad (2)$$

$$\vec{H} = LSTM(\vec{h}_1, \vec{h}_2, \vec{h}_3, \dots, \vec{h}_t) \quad (3)$$

$$H = \vec{H} + \vec{H} \quad (4)$$

其中, \mathbf{w}_t 是向量 \mathbf{x}_t 的权重矩阵; \mathbf{u} 是在 t 时刻 \mathbf{h}_{t-1} 的权重矩阵; \mathbf{b}_t 是此刻的偏置项; \vec{H} 表示各个时刻正向 LSTM 的输出; \vec{H} 表示各个时刻反向 LSTM 的输出, H 是将 \vec{H} 和 \vec{H} 合并后得到的输出向量。

考虑到加入 emoji 表情对微博文本情感极性的影响, 本文提出了基于 emoji 表情符号的注意力机制。对于句子而言, 每个词对情感极性的贡献度并不相同, 在不同句子中, emoji 表情的权重也是不相等的。emoji 表情特征计算公式如下:

$$\mathbf{S} = \sum_{i=1}^r \alpha_i \mathbf{h}_i \quad (5)$$

$$\text{signi}(e_j) = \mathbf{v}^T \tanh(\mathbf{w}_\theta e_j + \mathbf{b}) \quad (6)$$

$$\alpha_i = \frac{\exp(\text{signi}(e_j))}{\sum_{j=1}^L \exp(\text{signi}(e_j))} \quad (7)$$

$$r = \sum_{i=1}^L \alpha_i e_j \quad (8)$$

其中, \mathbf{S} 表示在句子中计算词语隐藏状态 \mathbf{h}_i 的加权求和; $\{\vec{h}_1, \vec{h}_2, \vec{h}_3, \dots, \vec{h}_t\}$ 表示输入的文本向量; α_i 表示文本中每个隐藏状态的注意力权重; $\{e_1, e_2, e_3, \dots, e_j\}$ 是 emoji 表情向量, $\text{signi}(e_j)$ 表示 emoji 表情向量的重要程度; α_i 表示每个 emoji 表情的权重; \mathbf{w}_θ 和 \mathbf{v} 表示可学习的参数; \mathbf{v}^T 是 \mathbf{v} 的转置矩阵; \mathbf{b} 表示偏置项; r 表示在注意力加权计算后得到的所有 emoji 表情的向量表示。

1.4 语义合成层

在语义合成层, 将 BiLSTM 层得到的文本特征和 emoji 特征放入自注意力机制层进行加权求和, 得到每个分类标签的分数, 并使用线性变换将所有特征表示投影到一个目标空间, 得到文本情感分析的最终特征表示。

1.5 情感分类层

在情感分类层, 使用 *softmax* 函数进行文本情感分类处理。由于 emoji 表情附带了强烈丰富的情感色彩, 本文将情感极性分为积极和消极两种类别。其计算公式和损失函数如下:

$$p_c = \frac{\exp(d_c)}{\sum_{k=1}^c \exp(d_k)} \quad (9)$$

$$L = \sum_{d \in D} \sum_{c=1}^c p_c^e(d) (\log p_c(d)) \quad (10)$$

其中, c 表示情感标签的数量, p_c 表示情感极性的预测公式。假设 d 是爬取的微博数据集合, $p_c^e(d)$ 为某一文本的目标分布, $p_c(d)$ 为其预测的情感分布, 训练的目的是让 $p_c^e(d)$ 与 $p_c(d)$ 之间的交叉熵尽可能损失最小。

2 实例验证

2.1 数据集

MFTEA 是用于新浪微博的主题情感分析模型, 为了保证结果的客观和有效, 经过分析选取公开的 weibo_senti_100k 数据集作为本文实验环节的标准数据集。该数据集是带有情感标签的二分类中文数据集, 将情感极性分为积极(情感标签记为 1)和消极(情感标签记为 0)两类, 其特点是基本每条评论都包含 emoji 表情符号, 适用于本文多特征情感分析的需求。该数据集共有 10 万余条新浪微博评论, 其中正向评论和负向评论各约 5 万条。数据集中部分数据样例见表 1。

表 1 weibo_senti_100k 数据集数据样例

Table 1 Sample data of dataset weibo_senti_100k

微博评论	情感标签
🇨🇷 克罗地亚球迷很爱放烟火! 球又没进, 就硝烟四起。	0
😭😭😭 终于收工啦, 脚丫子快冻掉了	0
😄, 😄 这个太赞了, 生活大爆炸第六季马上就要出啦	1
飘雪的北京 需要双份早餐……	1

2.2 评价指标

实验采用五折交叉验证, 将原始数据随机分成 5 个相等的部分进行实验, 其中 4 等份用于模型训练, 最后一部分用于测试。

为了验证各深度学习模型在文本情感分析领域的表现, 需要制定合理的模型评价方法。通常采用 4 个常用的模型评价指标作为实验的评估标准, 分别是准确率 (*Accuracy*)、精确率 (*Precision*)、召回率 (*Recall*) 和 F_1 值 (*F-score*)。研究给出的计算公式具体如下:

$$Acc = \frac{T_p + T_n}{T_p + T_n + F_p + F_n} \quad (11)$$

$$P = \frac{T_p}{T_p + F_p} \quad (12)$$

$$R = \frac{T_p}{T_p + F_n} \quad (13)$$

$$F_1 = \frac{2PR}{P+R} \quad (14)$$

其中, T_p 、 F_p 、 F_N 、 T_N 分别为真正例、假正例、假反例和真反例的样本数目。

2.3 实验结果分析

为了验证 MFTEA 模型的有效性,将其分别与 BiLSTM-Attention、BiGRU-Attention、CNN-BiLSTM、BERT-TextCNN 进行对比。其中,BiLSTM-Attention 是基于注意力机制的双向 LSTM 模型,通过注意力机制为特征向量分配不同的权重;BiGRU-Attention 是基于注意力机制的双向门控循环单元网络,利用 BiGRU 提取文本特征,并通过注意力机制分配不同权重;CNN-BiLSTM 先利用 CNN 提取文本局部特征,再利用双向 LSTM 提取上下文特征;BERT-TextCNN 先利用 BERT 模型生成词向量,再通过 TextCNN 提取关键特征后送入全连接层和 *softmax* 层进行情感分类。实验结果见表 2。

表 2 情感分类模型实验结果

Table 2 Experimental results of emotional analysis model %

实验序号	实验模型	P	R	F1	Acc
1	BiLSTM-Attention	88.07	84.86	86.43	87.71
2	BiGRU-Attention	92.46	90.58	91.52	90.89
3	CNN-BiLSTM	86.59	81.02	83.71	84.61
4	BERT-TextCNN	80.94	78.13	79.51	80.64
5	MFTEA	93.41	96.27	94.81	95.27

从表 2 中数据可以看出,MFTEA 模型取得的分类效果最好,准确率相比其它几种模型均有所提升。由于 MFTEA 模型在预训练和微调过程中充分学习了上下文信息,能够在不同语义环境下区分情感特征,得到更高质量的文本向量表示,尤其是在一些微博短文本上更高效地获取到了关键特征,文本特征提取能力得到了很大提升,因而能够取得更好的分类效果。

3 结束语

本文介绍了基于 BERT-Attention-BiLSTM 的多特征主题情感分析模型 MFTEA 的构建工作,并通过实验对比证明了该模型在为微博文本情感分析任务中的有效性。

参考文献

[1] GUO Jiang, SHAH D J, BARZILAY R. Multi-source domain adaptation with mixture of experts[J]. arXiv preprint arXiv:1809.02256, 2018.

[2] KHOO C S, JOHNKHAN S B. Lexicon-based sentiment analysis: Comparative evaluation of six sentiment lexicons[J]. Journal of Information Science, 2018, 44(4): 491-511.

[3] YANG Aimin, LIN Jianghao, ZHOU Yongmei, et al. Research on building a Chinese sentiment lexicon based on SO-PMI[J]. Applied Mechanics and Materials, 2013,263:1688-1693.

[4] AGARWAL B, MITTAL N, AGARWAL B, et al. Machine learning approach for sentiment analysis[M]// AGARWAL B, MITTAL N. Prominent feature extraction for sentiment analysis. Cham;Springer, 2016: 21-45.

[5] MALVIYA S, TIWARI A K, SRIVASTAVA R, et al. Machine learning techniques for sentiment analysis: A review[J]. SAMRIDDHI: A Journal of Physical Sciences, Engineering and Technology, 2020, 12(2): 72-78.

[6] ZHANG Lei, WANG Shuai, LIU Bing. Deep learning for sentiment analysis: A survey[J]. Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery, 2018, 8(4): e1253.

[7] BENGIO Y, DUCHARME R, VINCENT P, et al. A neural probabilistic language model[J]. Journal of Machine Learning Research (JMLR), 2003,3:1137 - 1155.

[8] MIKOLOV T, CHEN Kai, CORRADO G, et al. Efficient estimation of word representations in vector space[J]. arXiv preprint arXiv:1301.3781, 2013.

[9] WANG Libo, FANG Shenghui, ZHANG Ce, et al. Efficient hybrid transformer: Learning global-local context for urban sense segmentation[J]. arXiv preprint arXiv:2109.08937, 2021.

[10] HOCHREITER S, SCHMIDHUBER J. Long short-term memory[J]. Neural Computation, 1997, 9(8): 1735-1780.

[11] ENVELOPE M, MR A, ST B, et al. User sentiment analysis in conversational systems based on augmentation and attention-based BiLSTM[J]. Procedia Computer Science, 2022, 207: 4106-4112.

[12] 关鹏飞, 李宝安, 吕学强, 等. 注意力增强的双向 LSTM 情感分析[J]. 中文信息学报, 2019,33(2): 105-111.

[13] ACHEAMPONG F A, NUNOO-MENSAH H, CHEN Wenyu. Transformer models for text-based emotion detection: A review of BERT-based approaches[J]. Artificial Intelligence Review, 2021,54:85789-85829.

[14] 孔丽雅, 周治平. 基于 ELMo 的混合注意力网络的方面级情感分析研究[J]. 中文信息学报,2023,37(6):147-156.

[15] DONG Yongfeng, FU YU, WANG Liqin, et al. A sentiment analysis method of capsule network based on BiLSTM[J]. IEEE Access, 2020, 8: 37014-37020.

[16] XU Guixian, MENG Yueting, QIU Xiaoyu, et al. Sentiment analysis of comment texts based on BiLSTM[J]. IEEE Access, 2019, 7: 51522-51532.

[17] LI H, ZHANG L. Deep learning based text sentiment analysis during epidemic[J]. International Core Journal of Engineering, 2021, 7(7): 467-472.

[18] WANG Kai, SHEN Weizhou, YANG Yunyi, et al. Relational graph attention network for aspect-based sentiment analysis[J]. arXiv preprint arXiv:2004.12362, 2020.

[19] WANG Xi, CHEN Zhichao, KONG Mingming, et al. A hunger-based scheduling strategy for distributed crawler[J]. Expert Systems with Applications, 2023,222: 119798.

[20] REN Y. The application of case teaching method for "python and application" under the concept of curriculum ideology and politics[J]. Open Access Library Journal, 2022, 9(11): 1-7.

[21] 顾益军, 樊孝忠, 王建华, 等. 中文停用词表的自动选取[J]. 北京理工大学学报, 2005, 25(4): 337-340.