

doi: 10.3969/j.issn.0490-6756.2020.02.010

基于混合神经网络的中文隐式情感分析

赵容梅¹, 熊熙¹, 琚生根², 李中志¹, 谢川¹

(1. 成都信息工程大学网络空间安全学院, 成都 610225; 2. 四川大学计算机学院, 成都 610065)

摘要: 隐式情感分析是情感计算的重要组成部分, 尤其是基于深度学习的情感分析近年来成为了研究热点. 本文利用卷积神经网络对文本进行特征提取, 结合长短期记忆网络(LSTM)结构提取上下文信息, 并且在网络中加入注意力机制, 构建一种新型混合神经网络模型, 实现对文本隐式情感的分析. 混合神经网络模型分别从单词级和句子级的层次结构中提取更有意义的句子语义和结构等隐藏特征, 通过注意力机制关注情绪贡献率较大的特征. 该模型在公开的隐式情感数据集上分类准确率达到77%. 隐式情感分析的研究可以更全面地提高文本情感分析效果, 进一步推动文本情感分析在知识嵌入、文本表示学习、用户建模和自然语言等领域的应用.

关键词: 情感分析; 深度学习; 卷积神经网络; 注意力机制; 长短期记忆网络

中图分类号: TP391 **文献标识码:** A **文章编号:** 0490-6756(2020)02-0264-07

Implicit sentiment analysis for Chinese texts based on a hybrid neural network

ZHAO Rong-Mei¹, XIONG Xi¹, JU Sheng-Gen², LI Zhong-Zhi¹, XIE Chuan¹

(1. College of Cybersecurity, Chengdu University of Information Technology, Chengdu 610225, China;

2. College of Computer Science, Sichuan University, Chengdu 610065, China)

Abstract: Implicit sentiment analysis is an important part of emotional computing, especially sentiment analysis based on deep learning has become a research hotspot in recent years. This paper uses convolutional neural network to extract features from text, combines long-short-term memory network (LSTM) structure to extract context information, and adds attention mechanism to the network to construct a new hybrid neural network model to realize implicit emotions analysis on text. The hybrid neural network model extracts more meaningful semantic features such as sentence semantics and structure from the hierarchical structure of word level and sentence level respectively, attention is paid to the characteristics of large emotional contribution rate through attention mechanism. The proposed model has a classification accuracy of 77% on the public implicit sentiment data set and can improve the effect of text emotion analysis more comprehensively.

Keywords: Emotional analysis; Deep learning; Convolutional neural network; Attention mechanism; Long-term and short-term memory network

收稿日期: 2019-09-28

基金项目: 国家自然科学基金(81901389); 中国博士后科学基金(2019M653400); 四川省科技计划项目(2018GZ0253, 2019YFS0236, 2018GZ0182, 2018GZ0093, 2018GZDZX0039)

作者简介: 赵容梅(1996-), 女, 四川南江人, 硕士研究生, 研究方向为自然语言处理与情感分析.

通讯作者: 熊熙. E-mail: flyxiongxi@gmail.com

1 引言

情感分析是自然语言处理领域的一个热点研究问题, 可以理解为通过识别文本观点以分析情感的倾向性, 文本情感分析不同于文本分类和文本挖掘之处在于: 情绪具有抽象性, 仅以字面意思判断文本所表达的情绪是片面的, 这种方式只适用于带有明显的情感词或者是能够充分表达情绪的表情符号的文本。

现今越多越多的人愿意在微博、微信等社交平台上通过发表纯文字或者是带有图片的动态信息表达自己的生活状态和情绪状态, 少部分的动态信息可以直接表达自己当前情绪, 而大量动态语言则较含蓄, 没有出现明显的情感词, 无法直观判断所表达的隐式情感。如表 1 的三句话, 分别表达了褒义、贬义和客观(中性)的隐式情感。

表 1 隐式情感句
Tab. 1 Implicit emotional sentences

序号	隐式情感句
1	这种自然景观, 不输给国外的什么大峡谷的!!!
2	多年前骑行来的感觉全没了完全商业化了!!!
3	2009 年 5 月新馆落成后易名“四川博物馆”

显式情感分析的研究成果目前已经比较丰富, 隐式情感分类的研究仍处于起步的阶段。隐式情感分析在知识嵌入、文本表示学习、用户建模和自然语言理解等方面的研究起着重要作用, 可以更全面地提高文本情感分析效果, 推动文本情感分析在相关领域的应用。

近年利用机器学习方法进行情感分析的研究提升了实验性能。王进等人^[1]针对行人属性分类受行人属性不均衡影响的问题, 提出了一种基于属性敏感卷积神经网络的行人属性分类方法。Hu 等人^[2]将卷积神经网络应用于语义匹配任务中, 陈波等人^[3]提出一种基于循环结构的神经网络文本分类方法, 该方法对文本进行单次正向及反向扫描, 能够在学习单词表示时尽可能地捕获上下文信息。Gu 等人^[4]将 SVM 嵌入到卷积神经网络中代替传统的 softmax 分类器, 该模型在中文句子分类任务中有比较好的实验性能。

Chen 等人^[5]提出了一种 RNN 模型, 在多任务学习框架下同时进行手写体识别和文本行识别。罗帆等人^[6]提出一种多层网络 H-RNN-CNN, 模型是 CNN 和 RNN 的简单结合, 用于处理中文文

本情感分类任务。彭嘉毅等人提出一种基于字符特性, 双向长短时记忆网络(Bi-LSTM)与条件随机场(CRF)相结合的信息安全领域命名实体识别方法。目前有关文本的研究应用比较广泛的模型是 LSTM 和 GRU 和基于这两个模型的改进算法。Tran 等人^[8]将 GRU 模型应用在对话系统中, 实验表明该算法的实验性能高于原有的算法。

注意力机制逐渐被广泛应用于各类自然语言处理任务中。杜天宝等人^[9]将文本中的每个词语映射成情感词向量, 进而将其作为卷积神经网络的输入, 并加入注意力机制对输出结果进行优化。Huang 等人^[10]提出了用于语音情感识别的深度卷积神经网络, 同时在该网络中引入了注意力机制, 用于学习与任务相关的话语结构。刘广峰等人^[11]针对文档水平情感分析传统模型存在先验知识依赖以及语义理解不足问题, 提出一种基于注意力机制与层次网络特征表示的情感分析模型 TWE-ANN。

一些研究者通过组合各类神经网络解决单一神经网络中存在的不足。Chen 等人^[12]提出了一种基于 CNN、BiLSTM 和条件随机场 CRF 的组合模型, 该模型采用分而治之的句子分类方法。赵勤鲁等人^[13]针对当前文本分类神经网络不能充分提取词语与词语和句子与句子之间的语义结构特征信息的问题, 提出一种基于 LSTM-Attention 的神经网络实现文本特征提取的方法。Li 等人^[14]提出了 BiLSTM-CNN 方法研究了自然语言处理在文本分类中的应用问题, 该模型在新闻文本分类方面具有很大的优势。

目前关于情感分析方面的研究工作主要可以分为两大类: 基于规则与词典的情感分析方法, 以及基于机器学习的情感分析方法。基于规则和词典的情感分析方法的分类灵活性较低、可迁移性差; 利用机器学习方法进行情感分析解决了情感词典适应性差的问题, 传统的机器学习方法的实验性能还有待进一步提高; 基于深度学习方法可以提取文本中更深层次句子级特征, 因而提高文本分类的准确率。

本文在深度模型的基础上加入了注意力机制, 在中文隐式情感分析上表现出了良好性能。本文提出的模型利用双层卷积神经网络对文本级和句子级的特征进行提取, 结合双向长短期记忆网络 BiLSTM 结构注意其上下文信息, BiLSTM 替换了传统卷积神经网络中的池化层, 并且在其中加入注

注意力机制,分析中文文本中包含的隐式情感,将其分为不含情感、褒义隐式情感和贬义隐式情感三类。

2 基本原理

2.1 LSTM

长短期记忆网络 LSTM 是 RNN 的一个变体,用于解决 RNN 网络中存在的梯度爆炸和梯度消失问题. LSTM 中包含遗忘门、输入门和输出门三种类型的门,这三种门也分别对应 LSTM 的内部实现的忘记阶段、选择记忆阶段和输出阶段. LSTM 网络如图 1 所示.

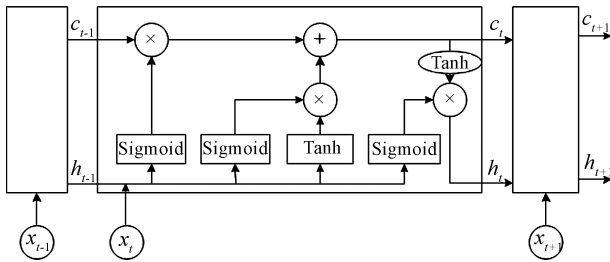


图 1 LSTM 原理图

Fig. 1 Long Short Term Memory network

遗忘门主要是选择性的忘记上一个节点传进来的数据,遗忘门控 f_t 计算方式如下.

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) \quad (1)$$

其中, h_{t-1} 表示上一个 cell 的输出; x_t 表示当前 cell 的输入; σ 表示 sigmoid 函数.

选择阶段是对输入进行选择性的记忆,也就是通过 sigmoid 层决定那些信息需要更新,更新后的信息用 i_t 表示, tanh 层生成的向量表示为 \tilde{C}_t , 最终更新的细胞状态由 C_t 表示, 如下式.

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) \quad (2)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C) \quad (3)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (4)$$

输出门决定模型中最终输出的值, 首先经过 sigmoid 对输入数据进行二次选择, 保留需要输出的部分记为 o_t , 再将最终更新的细胞状态输入 tanh 做一个非线性化处理, 将两个部分相乘则得到了最终的输出 h_t .

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o) \quad (5)$$

$$h_t = o_t * \tanh(C_t) \quad (6)$$

$W_f, b_f, W_i, b_i, W_o, b_o$ 分别表示遗忘门、输入门和输出门对应三个阶段的权重参数和偏置参数.

2.2 注意力机制

在 NLP 中, 注意力机制允许模型根据输入文本和到目前为止生成的文本来学习应该关注什么, 而不是像标准 RNN 和 LSTM 那样将整个源文本编码成固定长度的向量. 注意力机制原理如图 2 所示.

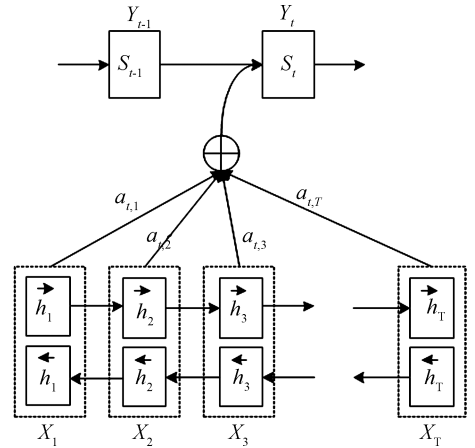


图 2 注意力机制原理图

Fig. 2 Schematic diagram of attention mechanism

在自然语言处理中注意力机制首先用于机器翻译任务中, 在机器翻译过程中注意力机制为每个单词都赋予了一个权重, 而这个权重则是通过计算出来的注意力, 实现在翻译下一个单词时可以只关注局部对自己比较重要的信息, 注意力计算公式如下.

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})} \quad (7)$$

其中, a_{ij} 是一个 softmax 模型输出, 概率值的和为 1; e_{ij} 表示一个对齐模型, 在机器翻译中表示的是在翻译第 i 个词时, 受 encoder 端第 j 个词的影响程度. 常见的对齐计算方式有点乘 (Dot product), 值网络映射 (General) 和 concat 映射三种方式. 注意力机制本质上还是一个 Encoder-Decoder 结构, 早期应用在机器翻译任务的 sequence-to-sequence 结构在翻译单个单词时需要关注整个原句, 不能只关注与之对应的单词和上下文. Attention 机制就解决了上述问题, 这也是在自然语言处理任务中被广泛应用的一个原因.

3 本文模型

3.1 模型简介

本文使用的混合神经网络包括卷积神经网络 (CNN), 双向长短期记忆网络 (BiLSTM) 和注意力

机制,模型结构如图 3 所示. 模型结构的第一层是预处理过程中的词嵌入层,第二层是一个包括三个卷积核的卷积层,使用混合卷积核的主要目的是尽可能多的提取更有意义的抽象化特征,用 BiLSTM 代替卷积神经网络中的池化层,利用 Droupout 防止过拟合.

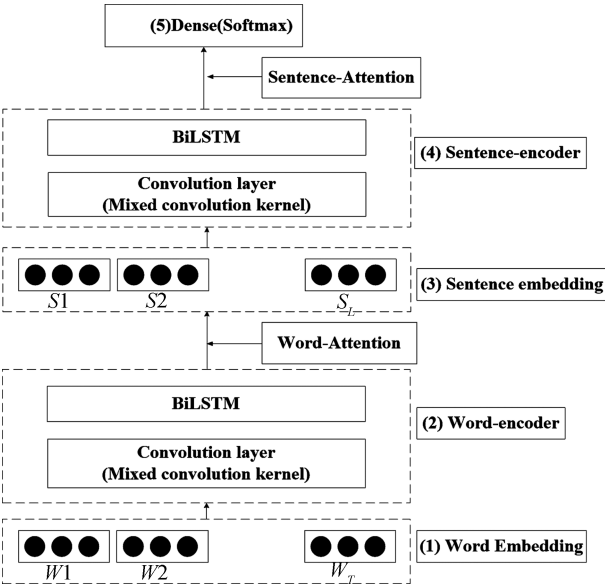


图 3 混合神经注意力网络(CLA)

Fig. 3 The mixed neural attention network(CLA)

BiLSTM 的输出序列结合注意力机制注意贡献率较大的词向量,第三层也就是我们的句子嵌入层,第四层和第二层结构一样,用于提取句子的抽象化特征. 在经过特征提取的句向量中加入注意力后作为最后一层 softmax 分类器的特征进行分类. 词嵌入层的词向量大小为 100 维,每篇文档最大的句子数量设置为 15. 利用卷积神经网络处理文本信息时卷积核的宽度与词向量的大小一致,因此卷积核每次覆盖的上下文的范围就与卷积核的大小一致,因此可以通过设置卷积核的尺寸来决定将几个连续词的特征表示出来. 传统的 CNN 除了有卷积层之外,还通过池化层再次提取较为重要的特征简化网络,本文选择利用双向的 LSTM 来代替池化层,在双向 LSTM 输出的序列基础上加入注意力机制,进一步的将序列中对于分类具有较大贡献率的单词和句子给予更多的注意力.

3.2 BiLSTM 模型编码

本文使用的 BiLSTM 最终对得到的两个方向的 h_i 进行简单的拼接即可, BiLSTM 对句子的编码模型如图 4 所示.

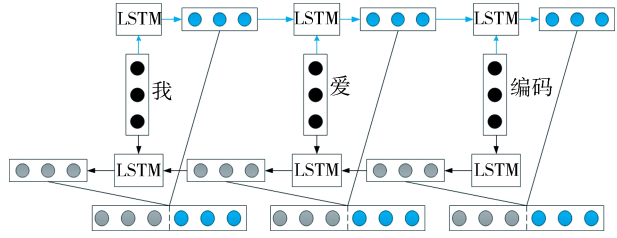


图 4 BiLSTM 对“我爱编码这句话进行双向编码”
Fig. 4 BiLSTM's bidirectional coding of "I love coding"

前向 LSTM 依次输入“我”,“爱”,“编码”得到三个向量,如图 3 蓝色部分表示的向量,后向的 LSTM 依次输入“编码”,“爱”,“我”得到三个向量如图 3 灰色部分表示的向量,将前向和后向 LSTM 得到的隐向量进行拼接就得到了 BiLSTM 编码的向量. 在本文的方法中, BiLSTM 并不是针对原始词经过词嵌入后得到的向量进行编码,而是在经过卷积神经网络对其进行特征提取之后,对含有重要意义的抽象特征的向量进行编码, BiLSTM 编码中含有向量的上下文信息,分类时能充分利用上下文信息,提高分类性能.

3.3 本文注意力机制

BiLSTM 将卷积后的提取的特征从两个方向进行汇总,充分保留特征中包含的上下文信息. 注意力机制使得模型能够将更多的注意力放在对情感分类有较大贡献的单词和句子中,学习局部重要性.

本文中的注意力机制中采用的对齐模式是点乘模式,将 BiLSTM 输出的向量 h_i 经过 tanh 变换得到隐向量表示 w_i ,再计算每个单词对于情感分类的注意力,最后通过简单的加和得到一个句子中的所有注意力 s_i , u_w 表示单词的上下文向量. 某一句语言所表达的情感总是与上下句密切相关,本文提出的方法能充分利用上下文信息,对语句的情感分析不仅仅对单个语句进行编码,将这个语句的上下文一起编码成为一个简单的文档,将有标记的情感句置于文档中部,前后填充上下文语句. 因此在进行特征提取时句子的上下文语义和结构信息也被包含其中,丰富了特征信息、提高了实验性能.

$$w_i = \tanh(W_w h_i + b_w) \tag{8}$$

$$a_{ij} = \frac{\exp(w_i^T u_w)}{\sum_t \exp(w_i^T u_w)} \tag{9}$$

$$s_i = \sum a_{ij} h_i \tag{10}$$

4 实验结果与分析

4.1 数据集

实验所用数据集是全国社交媒体处理大会(SMP2019)公开的由山西大学提供的数据集^[15],主要是来源于产品论坛、微博、旅游网站等平台的数据. 本文主要的工作是对中文的隐式情感句进行评测,数据集中包含显示情感词的文本已经通过大规模情感词典进行过滤处理. 处理后的数据集中将隐式情感句进行了部分标注,分为褒义隐式情感句(1),贬义隐式情感句(2)和中性情感句(0)三类. 数据以切分句子的文档形式发布,其中包含有句子的完整上下文信息,数据集的详细数据如表 2 所示.

表 2 数据集详细数据

Tab.2 Data set details

数据集	篇章数	褒义句数	贬义句数	中性句数	总数
训练集	12 664	6 992	3 835	3 961	14 788
测试集	4 391	2 554	1 233	1 358	5 145

中性句表示数据集中标注的不含情感句的数量,褒义句表示数据集中标注的含褒义隐式情感句的数量,贬义句表示数据集中标注的含贬义隐式情感句的数量,总数表示含有标注的语句总数.

4.2 评价指标

实验中所用的评价指标为准确率、精确度(P)、召回率(R)、准确率(Acc)、 F_1 值^[16]和汉明损失,汉明损失(Hamming_loss)^[17]用来计算多标签分类模型的精度,可以衡量预测所得标记与样本实际标记之间的不一致程度,其各个评价指标的计算公式如下.

$$P = \frac{TP}{TP + FP} \tag{11}$$

$$R = \frac{TP}{TP + FN} \tag{12}$$

$$F1 = \frac{2 * P * R}{P + R} \tag{13}$$

$$Acc = \frac{TP + TN}{TP + TN + FN + FP} \tag{14}$$

$$HammingLoss = \frac{1}{N} \sum_{i=1}^N \frac{XOR(Y_{i,j}, P_{i,j})}{L} \tag{15}$$

其中, TP 表示将正类预测为正类的数量; FN 表示将正类预测为负类的数量; FP 表示将负类预测为正类的数量; TN 表示将负类预测为负类的数量; N 是样本的总数; L 是标签的个数; $Y_{i,j}$ 是第 i 个预测结果中第 j 个分量的真实值; $P_{i,j}$ 是第 i 个预测结果中第 j 个分量的预测值; XOR 表示异或操作.

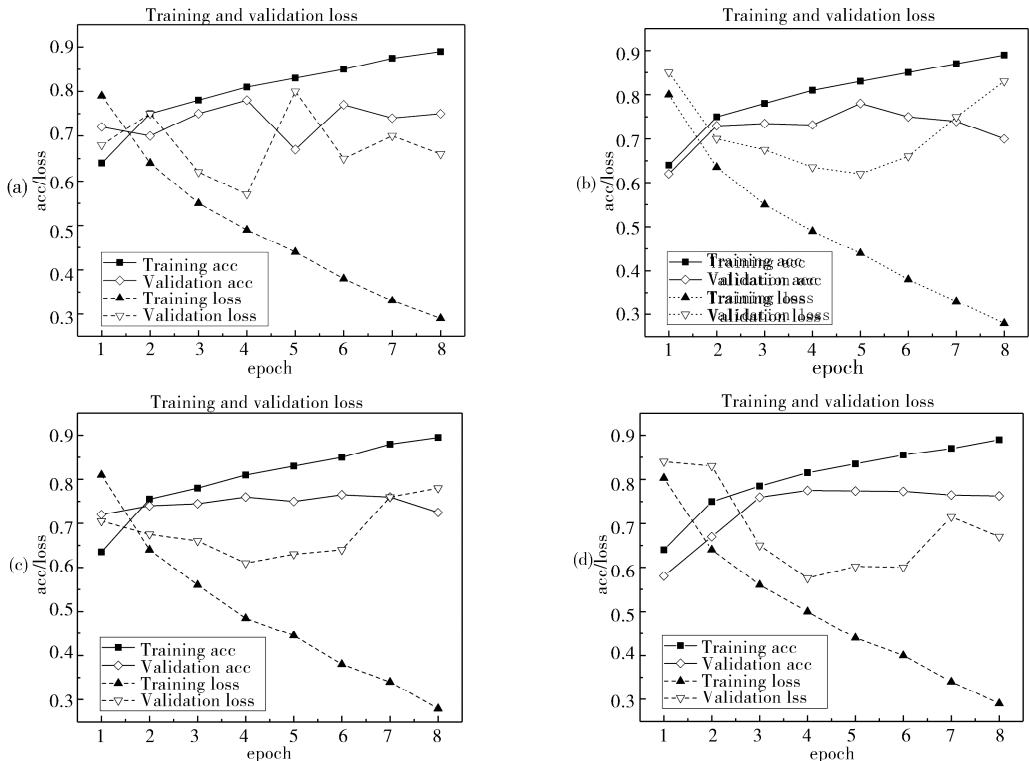


图 5 卷积核对比实验

Fig.5 Convolution comparison experiment

4.3 实验结果分析

卷积神经网络的卷积层将单词和句子的表示中有意义的抽象特征抽取出来,在卷积核大小分别为 3,4,5 以及三个卷积核的混合四种情况下进行对比实验,此时实验的 epoch 设置为 8,实验结果如图 5 所示,图 5 的(a)~(d)图就分别代表卷积核大小为 3、4、5 和混合卷积核四种情况.从图 5 中可以发现,在训练了 8 个 epoch 的情况下,结果开始收敛,训练得到的精度比较高能达到 0.9,测试的精度就在 0.7 上下浮动,可以明显看出的是在 epoch 为 3 或者 4 的时候,训练和测试的拟合情况达到最好,因此在后面的实验中的所有 epoch 都设置为 3.

通过比较图 5(a)~(d)的 4 组图可以看出,图 5(a)中精度和损失变化波动比较大,没有收敛的趋势,这是因为当卷积核为 3 时,每次提取到只是相邻三个单词的特征,此时的句子中的语义和结构信息没有全部包含进去,所以实验效果相对较差,卷积核为 4 和 5 的时候曲线明显趋于缓和,当使用混合卷积核的时候,实验效果达到最好.卷积核为 3 时能够充分提取相邻单词之间的特征信息,卷积核为 5 时,提取特征的前后跨度增大,可以将语句中的语义及结构特征提取出来,将混合卷积核提取的特征结合起来就可以确保卷积神经网络在卷积过程中充分提取单词或语句中有意义的特征信息,因而提高了分类效果.

运用混合卷积核,epoch 设置为 3 的实验结果如表 3 所示.

表 3 混合卷积核下的实验性能

Tab. 3 Experimental performance under the mixed convolution kernel

情感句	精确度	召回率	F_1 值	标记样本数	汉明损失
中性	0.77	0.97	0.86	2554	—
褒义	0.76	0.81	0.78	1233	—
贬义	0.77	0.75	0.76	1358	—
整体	0.77	0.84	0.80	5145	0.23

从表 3 可知,对于中性句(不含情感)的分类 F_1 值达到了 0.86,明显优于含情感的语句分类效果.这是由于中性情感语句的语句结构和语言表达方式相对简单,而含情感的语句表达方式则比较含蓄.实验的整体准确率达到 0.77,汉明损失为 0.23,在隐式情感分类任务中达到了良好的分类效果.

5 结 论

在本文中,我们将传统的卷积神经网络 CNN 中的池化层用 BiLSTM 代替,作用是能够更加完整的保存提取特征中的上下文信息,并且在模型中加入了注意力机制,注意力机制使得模型能够将更多的注意力放在对情感分类有较大贡献的单词和句子中,学习局部重要性.实验结果表明该方法可以显著提高隐式情感句的分类效果.隐式情感分析的难点在于语义表达含蓄,隐式情感句的语义特征、上下文结构特征的信息提取直接影响到最终分类结果,在接下来的隐式情感分析研究中,最主要的工作就是从更细粒度的角度提取隐式情感句的特征信息.现在情感分析研究工作存在的数据稀缺性、类别不平衡、领域依赖性和语言不平衡等难点,因此基于多媒体融合、领域自适应、深层语义和社交网络的情绪分析将是之后研究工作的一个重点和热点.

参考文献:

- [1] 王进, 黄超, 王科, 等. 基于属性敏感卷积神经网络的行人属性分类[J]. 江苏大学学报: 自然科学版, 2019, 40: 431.
- [2] Hu B, Lu Z, Li H, *et al.* Convolutional neural network architectures for matching natural language sentences[C]// Advances in neural information processing systems, Montreal. San Mateo, CA: Morgan Kaufmann, 2014.
- [3] 陈波. 基于循环结构的卷积神经网络文本分类方法[J]. 重庆邮电大学学报: 自然科学版, 2018, 30: 705.
- [4] Gu C, Wu M, Zhang C. Chinese sentence classification based on convolutional neural network[C]// IOP Conference Series: Materials Science and Engineering, Macau. Bristol: IOP Pub, 2017.
- [5] Chen Z, Wu Y, Yin F, *et al.* Simultaneous script identification and handwriting recognition via multi-task learning of recurrent neural networks[C]// 2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR). Japan: IEEE, 2017.
- [6] 罗帆, 王厚峰. 结合 RNN 和 CNN 层次化网络的中文文本情感分类[J]. 北京大学学报: 自然科学版, 2018, 54: 459.
- [7] 彭嘉毅, 方勇, 黄诚, 等. 基于深度主动学习的信息安全领域命名实体识别研究[J]. 四川大学学

- 报:自然科学版, 2019, 56: 457.
- [8] Tran V K, Nguyen L M. Semantic refinement gr-based neural language generation for spoken dialogue systems[C]// International Conference of the Pacific Association for Computational Linguistics, Burma, Singapore: Springer, 2017.
- [9] 杜天宝, 于纯浩, 温卓, 等. 基于文本情感特征的心理评估模型[J]. 吉林大学学报: 理学版, 2019, 57: 927.
- [10] Huang C W, Narayanan S S. Deep convolutional recurrent neural network with attention mechanism for robust speech emotion recognition[C]// IEEE International Conference on Multimedia & Expo(ICME). Hong Kong: IEEE, 2017.
- [11] 刘广峰, 黄贤英, 刘小洋, 等. 基于主题注意力层次记忆网络的文档情感建模[J]. 四川大学学报: 自然科学版, 2019, 56: 833.
- [12] Chen T, Xu R, He Y, *et al.* Improving sentiment analysis via sentence type classification using BiLSTM-CRF and CNN[J]. *Expert Syst Appl*, 2017, 72: 221.
- [13] 赵勤鲁, 蔡晓东, 李波, 等. 基于 LSTM-Attention 神经网络的文本特征提取方法[J]. *现代电子技术*, 2018, 41: 167.
- [14] Li C, Zhan G, Li Z. News text classification based on improved Bi-LSTM-CNN[C]// 2018 9th International Conference on Information Technology in Medicine and Education (ITME). Hangzhou: IEEE, 2018.
- [15] 中国中文信息学会社交媒体处理专业委员会. 中文隐式情感分析评测[EB/OL]. http://conference.cipsc.org.cn/smp2019/smp_ecisa_SMP.html
- [16] 刘云, 黄荣乘. 最大判别特征选择算法在文本分类的优化研究[J]. *四川大学学报: 自然科学版*, 2019, 56: 65.
- [17] Dembczyński K, Waegeman W, Cheng W, *et al.* On label dependence and loss minimization in multi-label classification [J]. *MACH LEARN*, 2012, 88: 5.

引用本文格式:

中文: 赵容梅, 熊熙, 琚生根, 等. 基于混合神经网络的中文隐式情感分析[J]. *四川大学学报: 自然科学版*, 2020, 57: 264.

英文: Zhao R M, Xiong X, Ju S G, *et al.* Implicit sentiment analysis for Chinese texts based on a hybrid neural network [J]. *J Sichuan Univ: Nat Sci Ed*, 2020, 57: 264.