基于 ChineseBERT 的双通道隐式情感分类

朱士成¹ 钱 钢¹ ZHU Shicheng QIAN Gang

摘要

隐含情感分析是当前自然语言处理领域的热门研究方向之一。相比于传统的文本情感分析方法,隐含情感分析面临着一些挑战,例如表达方式较为隐晦、缺乏明确的情感词汇等。针对这些问题,提出一种基于 ChineseBERT 的双通道中文隐式情感分类模型。首先,采用嵌入汉字音形向量的 ChineseBERT 预训练模型来提取文本词的动态向量表征。然后,并行联合使用 CNN 与 BiLSTM 混合神经网络模型,通过嵌入自注意力机制的多尺度 CNN 网络捕捉文本局部特征,同时引入结合自注意力机制的 BiLSTM 提取文本深层次上下文信息特征,将改进后的 CNN 与 BiLSTM 进行特征拼接。最后,输入全连接层获得情感分类结果。经过实验,所设计的模型在 SMP2019 "拓尔思杯"数据集上 Acc 值达到了 82.5%,分类效果显著提升,验证了模型具有可行性和有效性。

关键词

隐式情感分析:注意力机制:ChineseBERT:双向长短期记忆神经网络:卷积神经网络

doi: 10.3969/j.issn.1672-9528.2024.06.005

0 引言

文本情感分析,又称为意见挖掘,是对主观性文本进行分析的过程,旨在揭示其中蕴含的情感倾向并对情感态度进行分类^[1]。与显式情感分析不同,隐式情感分析专注于解读隐含主观情感的文本,即使缺乏直接的情感词汇,也能捕捉隐式情感句的情感倾向^[2]。

隐式情感分析因缺少显性情感词,导致情感识别难题。 在隐式情感分析初期研究阶段,主要依赖上下文对隐式情感 句进行语义挖掘。潘东行等人[3]聚焦于挖掘隐式情感句中重 要的上下文语境特征,提出了一种新颖的中文隐式情感分类 模型。该模型有效融合了上下文特征,进一步提升了情感分 析的精准度。赵容梅等人^[4]则创新性地结合 CNN、BiLSTM 和 Attention 机制,提出了一种用于精准分析文本中隐式情 感的混合神经网络模型。Yuan 等人 [5] 综合考虑了上下文信 息和目标情感句的时序信息,提出联合 BiGRU+Attention 与 GRU 的隐式情感分析模型。近年来, 多项自然语言处理任 务通过应用预训练语言模型,均实现了性能上的显著提升。 黄山成等人^[6] 结合隐式情感极性与文本要素,提出基于 ER-NIE2.0-BiLSTM-Attention的情感分析方法,有效提升了隐 式情感的识别准确率。张军等人[7]在隐式情感分析领域,针 对句子语义中隐藏情感的捕捉难题进行了探索,提出了基于 RoBERTa 融合双向长短期记忆网络及注意力机制的 RBLA 模 型。陆靓倩等人[8]提出了一种结合文本、词性与依存关系的 图神经网络模型来进行隐式情感分类。模型首先抽取文本的

1. 南京审计大学计算机学院 江苏南京 211815

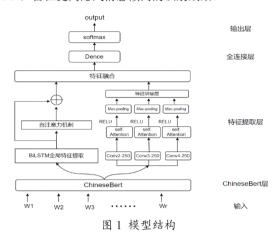
词性和依存特征,然后使用 BERT 预训练语言模型提取文本向量特征,从而构建了一个基于多种语言学特征的图注意力神经网络。李嘉伟等人^[9] 针对隐式情感句会受到上下文影响以及缺乏情感词汇的难题,提出一种基于图神经网络的文本图表征的方法,丰富隐式情感句的词级别的语义。

尽管上述研究已经取得了一定进展,但在隐式情感特征的提取方面仍存在一定的不准确性问题。本文选择融入汉字的音形特征进行隐式情感分类。现有的中文隐式情感分类模型中,没有考虑汉字读音、字形等特征信息,因此不能充分表示中文语义,从而在一定程度上降低了模型的准确性。汉字的字形特征包含了字的结构、笔画特征信息。隐式情感句虽然不含显式情感词,但仍然存在某一个字或词能够暗示句子的情感倾向,因此字形的特征信息可以辅助隐式情感的表达。针对目前主流的隐式情感分析方法缺少对汉字拼音、字形特征的考虑,本文提出基于 ChineseBERT 的双通道中文隐式情感分类模型,主要贡献点如下。

- (1)预训练语言模型 ChineseBERT^[10] 融入了汉字拼音和字形特征信息,增强了汉字的语义信息,可以解决中文隐式情感分类中多音字消歧、字形信息补充等问题,进而提高隐式情感分类的准确性。
- (2) 使用 BiLSTM 融合 Attention 机制对隐式情感文本的上下文信息进行深层次特征提取,上下文特征信息可以更好地被捕捉。
- (3)使用 CNN 并设置不同大小的卷积核,捕捉到不同 角度的文本局部特征,能够学习更加丰富和抽象的特征表示, 并且内部嵌入自注意力机制,能更好地获取文本的语义信息。

1 C-BLAC 双通道模型概况

如图 1 所示,为本文所提出的隐式情感分析方法,即以融合汉字拼音、字形特征的预训练模型 ChineseBERT 为基础联合使用 BiLSTM -Attention 机制与嵌入自注意力机制的多尺度 CNN,旨在提高隐式情感倾向的识别效果。



本模型主要包括 ChineseBERT 层、特征提取层、全连接层、输出层。

1.1 ChineseBERT 层

经典的中文预训练模型忽略了中文特有的两个重要方面:字形和拼音,它们携带着用于语言理解的重要语法和语义信息。而中文预训练模型 ChineseBERT,将汉字的字形和拼音信息纳入语言模型预训练中。其中,字形嵌入是根据汉字的不同字形得到的,能够从视觉特征中捕捉字符语义,而拼音嵌入则表征了汉字的读音,解决了中文中非常普遍的多音字问题。除此之外,提出全词掩码与字掩码两种训练策略,旨在帮助模型更全面地融合字词语义、字体形状、拼音信息以及序列上下文信息。

ChineseBERT 预训练模型的输入由位置嵌入向量和融合嵌入向量两个部分拼接而成,而融合嵌入向量是由字符向量、字形向量和拼音向量三者融合而成,如图 2 所示。

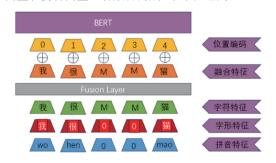


图 2 ChineseBERT 的输入

经过多层 Transformer 编码器的动态学习,文本向量表示为 $K = (k_1, k_2, \dots, k_s)$,作为特征提取层的输入。

1.2 特征提取层

特征提取层为双通道模式,输入的文本向量均为K。其

中,通道一采用 BiLSTM-Attention 进行深层次特征提取;通 道二采用多尺度 Att-CNN 提取更全面的文本局部特征。分别 输出文本上下文语义特征信息向量 Q_1 ,文本局部特征信息向量 Q_2 。

(1) BiLSTM

LSTM 网络通过引入一种被称为"门控单元"的机制来实现对序列中的长期依赖关系的捕捉。门控单元由遗忘门、输入门和输出门组成,每个门控单元都由一个 sigmoid 激活函数和一个点乘操作组成。这些门控单元通过控制信息的流动和存储来影响 LSTM 的状态,其结构如图 3 所示。

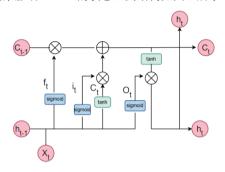


图 3 LSTM 单元

图 3 中, X_t 为 t 时刻输入, h_t 为 t 时刻隐含层输出, C_t 为 t 时刻记忆单元, C_t 为临时记忆单元, \otimes 表示逐元素相乘, \oplus 表示逐元素相加,遗忘门决定哪些信息被添加到当前记忆单元 C_t 中,输出门控制记忆单元 C_t 的输出结果。

单个 LSTM 模型结构的计算过程为:

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

$$\mathbf{i}_{t} = \sigma(\mathbf{W}_{t} \cdot [\mathbf{h}_{t-1}, \mathbf{x}_{t}] + \mathbf{b}_{t}) \tag{2}$$

$$C_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c) \tag{3}$$

$$C_{t} = f_{t} \otimes C_{t-1} \oplus i_{t} \otimes C_{t} \tag{4}$$

$$\mathbf{o}_{t} = \sigma(\mathbf{W}_{a} \cdot [\mathbf{h}_{t-1}, \mathbf{x}_{t}] + \mathbf{b}_{a}) \tag{5}$$

$$\mathbf{h}_{t} = \mathbf{o}_{t} \otimes \tanh(\mathbf{C}_{t}) \tag{6}$$

式中: W_{Γ} , W_{c} , W_{c} , W_{o} , 为各个模块对应的权重矩阵, b_{Γ} , b_{c} , b_{a} , 为偏置项, σ 为激活函数, tanh 为双曲正切激活函数。

BiLSTM 结合了前向和后向两个方向的 LSTM (如图 4 所示),以捕捉序列中的完整上下文信息。

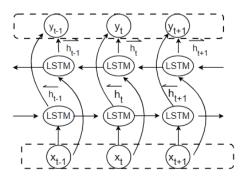


图 4 BiLSTM 结构图

前向 LSTM 生成前向隐藏层向量 \bar{h}_t ,捕捉从序列起始到当前位置的信息;后向 LSTM 则生成后向隐藏层向量 \bar{h}_t ,捕捉从序列末尾到当前位置的信息。 y_t 表示第 t 时刻输出的隐藏层向量, W_y 为权重矩阵, b_y 为偏置项。Y 表示所有时刻隐藏层向量的组合,计算方式为:

$$\overrightarrow{h}_{\bullet} = LSTM(x_{\bullet}, \overrightarrow{h_{\bullet}}) \tag{7}$$

$$\overline{h_{t}} = LSTM(x_{t}, \overline{h_{t-1}}) \tag{8}$$

$$\mathbf{y}_{t} = \sigma(\mathbf{W}_{y} \cdot [\vec{\mathbf{h}}_{t}, \overleftarrow{\mathbf{h}}_{t}] + \mathbf{b}_{y}) \tag{9}$$

$$Y = \{y_t, y_t, \dots, y_t, \dots, y_T\}$$

$$\tag{10}$$

(1) 注意力机制

注意力机制用于处理序列数据和建模序列中不同位置的 重要性或权重。它允许模型根据输入的不同部分,动态地分配注意力或关注度。本文模型通过引入注意力机制,在不同位置上动态分配注意力。这使得模型能够更好地捕捉序列中的重要信息,提高建模的准确性和效果。该层输入为 Y。E 为权重分配后注意力机制层的输出,则其公式为:

$$\mathbf{M} = \tanh(\mathbf{W}_{m} \cdot \mathbf{Y} + \mathbf{b}_{m}) \tag{11}$$

$$\boldsymbol{\alpha} = softmax(\boldsymbol{W}^{\mathsf{T}} \cdot \boldsymbol{M}) \tag{12}$$

$$\mathbf{E} = \mathbf{V} \cdot \mathbf{a}^{\mathrm{T}} \tag{13}$$

式中: W_m 作为注意力网络的权重矩阵, b_m 为偏置向量, W^T 为随机初始化参数矩阵,a 为每个词向量的标准化权重。

最后,将注意力机制层的输出 E 与 BiLSTM 提取的上下文语义特征信息向量 Y进行特征拼接,得到融合特征向量 Q_1 ,其公式为:

$$\mathbf{Q}_{t} = \mathbf{E} \oplus \mathbf{Y} \tag{14}$$

(2) 多尺度 Att-CNN

原始的 CNN 模型由多层结构组成,能有效捕捉文本的局部关键特征。然而,为了更精准地处理长文本数据并捕捉文本的深层语义信息,本文提出了改进方案。该方案在 CNN的基础上增加了自注意力机制层,以建模词语间的关系,并灵活使用不同大小的卷积核来提取多种尺度的局部特征。经过卷积和自注意力处理后,再通过池化层进行特征降维和整合,从而优化模型的性能。

将 ChineseBERT 的输出的文本向量 K 作为卷积通道的输入,利用大小为 2、3、4 三种不同大小的卷积核来提取句子中的 n-gram 信息,从而更好地提取词语之间的相关性和语义信息,其运算过程为:

$$\mathbf{C} = f(\mathbf{K} \otimes \mathbf{W} + \mathbf{b}) \tag{15}$$

式中: \otimes 表示的是卷积运算,W和 b 表示的是权重和偏置, $f(\cdot)$ 表示激活函数。在卷积层输出 C 上应用自注意力机制。自注意力机制将 C 作为输入,计算每个位置与其他位置之间的注意力权重,并将注意力权重与 C 相加,得到自注意力层的输出 A。计算公式为:

$$Q = CW_q \tag{16}$$

$$K = CW_k \tag{17}$$

$$V = CW_{v} \tag{18}$$

$$A = \left(softmax\left(\frac{(QK^{T})}{\sqrt{d_k}}\right)\right)V \tag{19}$$

自注意力机制将词向量分别乘以不同的变换矩阵 W,得到查询矩阵 Q、关键字矩阵 K 和值矩阵 V。

接着,将提取的关键信息 A 进行池化运算,进一步提取主要信息,其运算过程为:

$$\widetilde{A} = \max(A) \tag{20}$$

将使用不同卷积核下的输出特征进行组合拼接作为改进后的 CNN 模型最终输出结果 Q_2 , Q_2 与 融合注意力机制的 BiLSTM 模型得到的上下文语义特征信息向量 Q_1 进行特征融合,得到最终的多尺度特征向量 Z。计算公式为:

$$Z = concat(\mathbf{Q}_1, \mathbf{Q}_2)$$
 (21)

1.3 全连接层

全连接层将高级的特征表示映射到最终的输出或分类结果,即将融合了局部文本特征和全局上下文特征的特征向量**Z**输入到全连接层,输出为**F**,计算公式为:

$$F = \tanh(W_s \cdot Z + b_s) \tag{22}$$

式中: W_s 表示全连接层权重矩阵, b_s 是偏置向量。

1.4 输出层

在输出层,将全连接层的输出 F 通过 softmax 函数进行 归一化处理,得到 F 在不同情感标签类别的近似概率值 n,计算公式为:

$$\mathbf{n} = softmax(\mathbf{A} \cdot \mathbf{F} + \mathbf{b}) \tag{23}$$

式中: A 表示输出层参数矩阵, b 为偏置向量。

2 实验与分析

2.1 数据集

本文所采用的数据集是第二届中文隐式情感分析评测 (SMP-ECISA 2021)数据集,数据源于微博、旅游网站、产 品论坛等社交网络媒体,涵盖国考、旅游等多个领域/主题。

该数据集的情感标签被划分为三类:中性-0、褒义-1和贬义-2。它包含9000条训练数据、5788条验证数据和5145条测试数据,见表2。从数据分布来看,不含情感的数据占据近半,这种不均衡性增加了隐式情感识别的挑战性。

表 2 实验数据统计

数据集	句子	褒义	贬义	中性
训练集	9000	2335	2401	4264
验证集	5788	1493	1557	2738
测试集	5145	1233	1358	2554

2.2 参数设置

训练样本分割成多个 Batch,每批含 Batch_size 个样本。 隐藏单元 1、2 是 BiLSTM 中两层 LSTM 的隐藏层大小。 Epoch 表示数据集完整通过神经网络的正反传播过程。实验 采用 Dropout 减少过拟合,Filter-Size 表示卷积神经网络不同的卷积核大小,具体参数详见表 3。

表 3 参数设置

参数	数值	
Batch_size	6	
隐藏层单元 1	320	
隐藏层单元 2	320	
Epoch	15	
Lr	1E-5	
Dropout	0.5	
Filter- Size	(2,3,4)	

2.3 评价方法

实验采用准确率 Acc、精确率 P、召回率 R 以及 F_1 值共 4 个指标评价模型的性能,其计算公式为:

$$Acc = \frac{TP + TN}{TP + TN + FN + FP}$$
 (24)

$$P = \frac{\text{TP}}{\text{TP} + \text{FP}} \tag{25}$$

$$R = \frac{\mathrm{TP}}{\mathrm{TP} + \mathrm{FN}} \tag{26}$$

$$F_1 = \frac{2PR}{P+R} \tag{27}$$

2.4 对比模型

为验证本文所提出的模型 C-BLAC 的性能,选取现有的 经典隐式情感分类模型作为对比实验模型,下面将简要介绍 这些对比模型。

TextCNN: 卷积神经网络模型,应用于隐式文本情感分类任务。它通过卷积层提取文本中的局部特征,再利用池化层进行特征选择和降维,最后通过全连接层进行分类。

LSTM: 循环神经网络的一种,擅长处理序列数据,特别是长期依赖问题的一类隐式情感分类模型。

BERT: BERT 预训练模型在隐式情感分析中,通过预训练学习文本上下文信息,然后将文本转换为向量表示,捕捉上下文和语义关系,从而识别和理解隐式情感倾向,进行准确分析。

CA-TRNN^[11]:结合 BiLSTM 与 TRNN 模型提取上下文特征和语义特征的隐式情感分类模型。

文献 [8]: 在 BERT 基础上进行了改进,在融合词性信息的基础上提取语义信息。它是一种结合文本、词性与依存关系的图神经网络模型,能进行隐式情感分类任务。

RoBERTa: 使用 RoBERTa 预训练模型做隐式情感文本分类任务。

2.5 实验结果分析

2.5.1 对比实验结果分析

实验在 SMP-ECISA 2021 数据集上使用本文所提出的

C-BLAC 模型与上述 6 种现有的隐式情感分类模型作对比, 实验结果如表 4。

表 4 各分类模型的各项性能对比

模型	Acc	P	R	F_1
TextCNN	0.691	0.648	0.631	0.636
LSTM	0.779	0.744	0.700	0.721
CA-TRNN	0.801	0.700	0.780	0.738
BERT	0.793	0.785	0.758	0.769
文献 [8]	0.809	0.788	0.804	0.794
RoBERTa	0.816	0.804	0.778	0.789
C-BLAC	0.825	0.809	0.799	0.804

从表 4 可以看出,单一神经网络模型、单一预训练语言模型和混合神经网络模型中各项性能表现最好的模型分别为 LSTM、RoBERTa 预训练模型以及文献 [16] 中所提出的图神经网络模型。

C-BLAC 模型与LSTM 相比, C-BLAC 模型在 Acc、P 值、 R 值和 F_1 值上分别高出了 4.6、6.5、9.9 和 8.3 个百分点。整 体各项性能指标提高显著。相较于 RoBERTa 预训练模型, C-BLAC 模型在性能上取得了显著的提升,具体表现为 Acc 提高了 0.9 个百分点, P 值提升了 0.5 个百分点, R 值增加了 2.1个百分点,同时 F_1 值也上升了1.5个百分点。这一结果 充分证明了本文提出的 C-BLAC 模型在隐式情感分析任务中 相较于预训练模型具备更出色的性能。另外,与当前主流的 文献 [9] 中所提出的图神经网络模型相比, C-BLAC 融合模 型在Acc、P值和F. 值上分别高出了文献中模型1.6个百分点、 2.1 个百分点和 1.0 个百分点。这说明相较于主流的融合模型, 本文的 C-BLAC 模型在融合字形和拼音特征的 BERT 预训练 模型的基础上,通过结合融合注意力机制的双向 LSTM 网 络模型和嵌入自注意力机制的卷积神经网络模型双通道的方 式,实现了对隐式情感文本中深层情感特征的精准提取和高 效分析,从而显著提升了隐式情感分析的整体效果。

2.5.2 消融实验结果分析

分别使用以融合汉字字形与拼音特征的 BERT 模型以及使用传统 BERT 模型为基础并融合自注意力机制的双向长短期记忆网络、多尺度卷积网络的两组消融实验,实验结果如表 5。

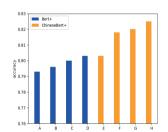
表 5 BERT 与 ChineseBERT 模型下的消融实验结果

模型	Acc	P	R	F_1
BERT	0.793	0.785	0.758	0.769
+ CA	0.796	0.772	0.775	0.773
+LA	0.800	0.780	0.773	0.776
BLAC(+LA+CA)	0.803	0.784	0.773	0.778
ChineseBERT	0.802	0.791	0.774	0.779
+ CA	0.818	0.805	0.795	0.798
+ LA	0.820	0.799	0.801	0.800
C-BLAC(+LA+CA)	0.825	0.809	0.799	0.804

其中,+LA表示在预训练模型的基础上融合了自注意力机制和双向长短期记忆神经网络;+CA表示在预训练模型的基础上,增加了嵌入自注意力机制的多尺度卷积网络;Our model表示在 ChineseBERT 预训练模型上同时增加 LA与 CA的混合模型。

由消融实验结果表可知,相比使用单一的 BERT 预训练模型,融合字形与拼音特征的 BERT 预训练模型的 Acc 值和 F_1 值分别提高了 0.9 和 1.0 个百分点。相比于分别引入 LA、 CA 和同时引入 LA 与 CA 的预训练语言模型,以 BERT 模型 为基础的隐式情感分类方法在 F_1 性能指标值分别依次增加 0.5 和 0.2 个百分点;以 ChinesBERT 为预训练模型的隐式情感分析方法在 F_1 值上提高了 0.6% 和 0.4%。由此可见,相较于仅增加单一特征的模型,同时融入两种特征的模型在 F_1 值上展现出了更为显著的提升。

通过观察图 6 和图 7 的展示,可以发现在各个消融实验模型中,与基于 BERT 预训练模型的方案相比,采用 ChineseBERT 预训练模型作为基础的模型在该数据集上表现出了更高的 Acc 值和 F_1 值。这一结果充分说明了 ChineseBERT 的优势,它不仅继承了 BERT 词向量编码的优秀特性,更在模型预训练的中文语料中巧妙地融入了汉字特有的音形特征信息,从而实现了更为精准的词向量编码。这样的设计使得 ChineseBERT 在中文处理任务中能够发挥出更好的性能。它充分体现了汉字本身的多维特征信息对于提升隐式情感识别能力的重要作用。在同一个预训练模型下可以发现,同时引入 LA 与 CA 的预训练模型的 F_1 与 Acc 的值均最高,由此可知,增加 LA 和 CA 均有助于提高模型对于隐式情感的识别。



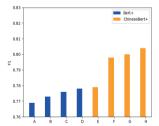


图 6 各模型准确度对比柱形图 图 7 各模型 F, 值对比柱形图

3 结语

针对隐式情感文本中情感特征捕捉不准确的问题,本文提出了 C-BLAC 模型,这是一种结合汉字多维特征的隐式情感分析方法。该模型运用 ChineseBERT 模型精准捕捉隐式情感句中字词的音形和语义特征,同时结合 CNN 模型和自注意力机制,通过设置不同卷积核,充分提取局部特征。此外,通过采用 BiLSTM-Attention 机制提取隐式情感句的上下文特征,增强上下文语意信息,获取更深层次特征,使得模型能够更全面地理解文本内容,更准确地识别情感倾向。经实验验证,与当前隐式情感文本分类模型相比,C-BLAC 模型在整体性能上实现了明显的优化与提升。未来,本文将进一步

深入研究隐式情感分析的中文语言特性,以不断优化和完善模型性能。

参考文献:

- [1] 王婷, 杨文忠. 文本情感分析方法研究综述 [J]. 计算机工程与应用,2021,57(12):11-24.
- [2] 彭俊杰.面向机器智能的情感分析 [J]. 自然杂志, 2024, 46(2): 150-156.
- [3] 潘东行, 袁景凌, 李琳, 等. 一种融合上下文特征的中文隐式情感分类模型[J]. 计算机工程与科学,2020,42(2):341-350.
- [4]ZHAO R, XIONG X, JI S, et al. Chinese implicit affective analysis based on hybrid neural network[J]. Journal of sichuan university(natural science edition),2020,57(2):264-270.
- [5] 袁景凌, 丁远远, 潘东行, 等. 基于时序和上下文特征的中文隐式情感分类模型[J]. 计算机应用,2021,41(10):2820-2828.
- [6] 黄山成,韩东红,乔百友,等.基于 ERNIE2.0-BiLSTM-Attention 的隐式情感分析方法 [J]. 小型微型计算机系统, 2021, 42(12): 2485-2489.
- [7] 张军,张丽,沈凡凡,等.RoBERTa 融合 BiLSTM 及注意力机制的隐式情感分析 [J]. 计算机工程与应用, 2022, 58(23): 142-150.
- [8] 陆靓倩,王中卿,周国栋.结合多种语言学特征的中文隐式情感分类[J]. 计算机科学,2023,50(12):255-261.
- [9] 李嘉伟,张顺香,李书羽,等.基于文本图表征的中文隐式情感分析模型[J/OL].数据分析与知识发现:1-16[2024-03-16]. http://kns.cnki.net/kcms/detail/10.1478.G2.20231225.0956.004.html.
- [10]SUN Z, LI X, SUN X, et al.ChineseBERT:chinese pretraining enhanced by glyph and pinyin information[C]//59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on natural Language Processing,Part 4.Stroudsburg:Association for Computational Linguistics, 2021:2065-2075.
- [11] 陈秋嫦, 赵晖, 左恩光, 等. 上下文感知的树递归神经网络下隐式情感分析 [J]. 计算机工程与应用, 2022,58(7):167-175.

【作者简介】

朱士成(1998—),男,江苏连云港人,硕士,研究方向: 自然语言处理。

钱 钢(1965—), 通 信 作 者 (email:qgmail@vip.sina.com), 男, 江苏常州人, 博士, 教授, 博士生导师, 研究方向: 大数据审计、现代审计技术与方法、自然语言处理。

(收稿日期: 2024-04-07)