



基于混合神经网络的中文短文本分类模型

陈巧红, 王磊, 孙麒, 贾宇波

(浙江理工大学信息学院, 杭州 310018)

摘要: 针对已有算法中特征表示存在的稀疏问题以及文本高层特征提取效果不佳问题, 提出了一种基于混合神经网络的中文短文本分类模型。该模型首先通过自定义筛选机制将文档以短语层和字符层进行特征词筛选; 然后将卷积神经网络(CNN)和循环神经网络(RNN)相结合, 提取文本高阶特征, 并引入注意力机制优化高阶向量特征; 最后将得到的高阶向量特征输入到全连接层得到分类结果。实验结果表明: 该方法能有效提取出文档的短语层和字符层特征; 与传统 CNN、传统 LSTM 和 CLSTM 模型对比, 二分类数据集上准确率分别提高 10.36%、5.01% 和 2.39%, 多分类数据集上准确率分别提高 12.33%、4.16% 和 2.33%。

关键词: 卷积神经网络; 循环神经网络; 短文本分类; 特征表示; 注意力机制

中图分类号: TP181

文献标志码: A

文章编号: 1673-3851(2019)07-0509-08

Chinese short text classification model based on hybrid neural network

CHEN Qiaohong, WANG Lei, SUN Qi, JIA Yubo

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

Abstract: In the existing algorithms, feature representation has the sparse problem and high-level extractin effect of texts is poor. Aiming at these problems, a Chinese short text classification model based on hybrid neural network was proposed. Firstly, the model screened the feature words at phrase level and character level through a self-defined the filtering mechanism. Then, convolutional neural network (CNN) and recurrent neural network (RNN) were combined to extract high-order features of texts, and the attention mechanism was introduced to optimize high-order vector features. Finally, the obtained high-order vector features were inputted into the full connection layer to obtain classification results. The experimental results showed that the proposed method could extract the features of phrase and character layers. Compared with CNN, LSTM and CLSTM models, the classification accuracy of the proposed model improveby 10.36%, 5.01% and 2.39% on binary dataset respectively, and the classification accuracyimproveby 12.33%, 4.16% and 2.33%on multiclass dataset respectively.

Key words: CNN; RNN; short text classification; feature representation; attention mechanism

0 引言

随着互联网大规模普及和上网人数急剧增加, 网络上产生的各种短文本数量呈指数式增长, 如何使这些海量的短文本数据产生实际价值, 是科研工作者急需解决的问题。提高短文本分类准确率, 在

一定程度上可以解决网络信息杂乱现象, 方便用户分流信息和准确定位所需信息; 帮助企业了解大众喜好, 掌握商机, 改进产品, 增加收益; 协助政府机构了解民情民意, 服务人民, 监管网络内容, 进行正确的舆论引导。互联网短文本分类作为信息处理关键技术之一, 在信息检索和知识挖掘领域已经取得很

收稿日期: 2018-12-13 网络出版日期: 2019-03-31

基金项目: 国家自然科学基金项目(51775513)

作者简介: 陈巧红(1978—), 女, 浙江临海人, 副教授, 博士, 主要从事计算机辅助设计及机器学习技术方面的研究。

大进展^[1]。

短文本分类相关研究主要有文本数据的特征表示和算法模型的选择与改进。王义真等^[2]提出了一种基于 SVM 的短文本分类方法,该方法在提取短文本相关特征时采用高维混合特征模型,但存在大量的特征工程工作,且在对短文本数据进行特征向量表示时会出现维度过高和数据稀疏问题,同时要求使用者对该领域要有足够的专业知识。Kim^[3]采用卷积神经网络(Convolutional neural network, CNN)对电影评论短文本进行分类,仅通过一层卷积和一层最大池化,最后将得到的高阶向量特征输入到全连接层得出分类结果。该方法虽然利用了深度学习模型,但是隐藏层太浅,不足以提取出更高层特征。Kalchbrenner 等^[4]提出了一种全新的模型,该模型采用动态卷积神经网络方法对句子进行自动建模,通过 K-max 池化操作来获取全文的特征向量。黄文明等^[5]将 K 近邻运用在文本加权上,对初始文本通过一定的权重采样,最后运用 K 近邻分类器得出分类结果。但该方法面对海量的数据集运算量过高,训练时间太长,实际效果不好。黄磊等^[6]使用基于循环神经网络(Recurrent neural network, RNN)改进的长短项记忆网络(Long short-term memory, LSTM)和门阈递归单元(Gated recurrent unit, GRU)计算节点的文本特征,最后通过 softmax 函数进行文本分类。该方法对时间序列问题处理效果很好,但是在处理长短不一致的短文本数据时效果不好,虽然该文采用了补齐和剔除的方式,但这种方式会引起信息冗余和信息丢失。

针对现有方法存在的特征表示稀疏问题以及提取文本高层特征效果不佳问题,本文提出一种基于混合神经网络的中文短文本分类模型。该模型采用自定义特征词筛选机制,通过一个全局字典对每一句短文本进行线性表示,给每一句短文本分词后的短语进行权重赋值;为解决传统分类方法中特征工程工作量大、特征提取不充分问题,采用一种 CNN 和 RNN 相结合的高阶特征提取网络,引入注意力机制对提取到的高阶特征向量进行优化,突出对短文本分类起作用的短语向量的重要性,以保证提取到的高阶特征向量更加符合语义信息。

1 混合神经网络模型

1.1 字符短语相结合的混合双通道

从语义分析的角度看,短语与字符对短文本分类具有重要意义,一句短文本中某几个单词或字

可以确定短文本的类别。例如,“计算机”、“程序”是比较独特的信息技术。在线产品评论中,字符“好”、“坏”和“棒”可以直接显示情感类别信息。Wang 等^[7]提出了基于词嵌入技术和 CNN 的中文分词方法,实验结果表明,该方法解释了字符级分词和短语级分词,可以提升短文本分类准确率;Zhou 等^[8]提出了基于复合递归神经网络的中文短文本分类方法,结果发现,将字符级和短语级相结合能带来更优的分类结果。为了进一步提高中文短文本分类准确率,本文在短语层分词和字符层分词基础上,采用自定义特征词筛选机制进行特征词权重筛选,将 CNN 和 RNN 结合以提取短文本高层特征,引入注意力机制优化高阶特征向量。本文设计的短语字符相结合混合双通道短文本分类模型如图 1 所示。其中: W 表示短语序列, C 表示字符序列, M_w 和 M_c 分别表示经过 CNN 特征提取后的中间短语特征向量和中间字符特征向量, h_w 和 h_c 分别表示将 CNN 提取到的中间特征向量输入给 RNN 进一步高阶特征提取得到的短语特征向量和字符特征向量, S_w 和 S_c 分别表示经过注意力机制优化后的短语向量和字符向量。

本文混合神经网络模型总共分为四层:自定义特征词筛选和词嵌入层、特征提取层、特征向量优化层与全连接层。自定义特征词筛选和词嵌入层同时作用在短语序列 W 和字符序列 C 。特征提取层和特征向量优化层分为两个平行的神经网络模块,分别提取短语向量 S_w 和字符向量 S_c , 然后根据式(1)将 S_w 和 S_c 合并,得到最终的文本特征向量 S 。将得到的 S 输入到全连接层进行分类。全连接层采用 softmax 函数计算分类结果,如式(2)所示。

$$S = [S_w \oplus S_c] \quad (1)$$

$$P = \text{softmax}(W_s S + b_s) \quad (2)$$

1.2 自定义特征词筛选和词嵌入

a)自定义特征词筛选。由于短文本中几个关键词往往就能代表整句的信息,所以本文通过将特定类别下整个训练集经过人工筛选和结合网络信息构造出一个高质量的全局字典,这个字典包含了某类别下经过筛选后的所有高质量短语,筛选的标准是人工判断和结合网络信息提供的与该类别相关度高的短语。最后将该类别下的每一篇文本用这个全局字典进行线性表示,计算如式(3)所示:

$$T(D_i) = \theta_i \sum_{j=1}^n D_j \quad (3)$$

其中: D_i 表示某一特定的文本经过分词后的一

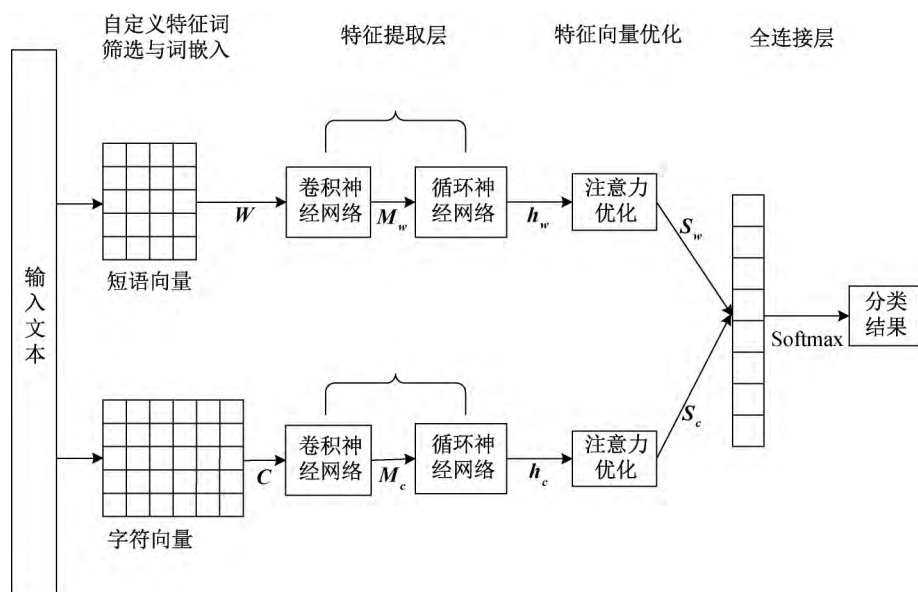


图1 短语字符相结合混合双通道模型

些短语; θ_i 表示一个权重系数向量, T 表示权重系数向量与文本分词后的短语之间的映射关系; 如果该篇文本的短语在全局字典中筛选到, 表示该短语对该类别的分类比较重要, 若没有筛选到, 表示该短语对该类别的分类不是很重要。本文通过实验选择权重系数向量 θ_i 的最优值, 在全局字典 D_i 中筛选到为 0.9, 没有筛选到为 0.1。

b) 词嵌入。经过自定义特征词筛选对每一篇短文本提取到的特征词, 本文采用了 word2vec^[9] 将其训练成一个 100 维的向量, 然后将自定义特征词筛选得到的权重系数与其相对应的短语向量相乘, 得到的最终向量就代表该短语的一个向量表示, 作为卷积神经网络的输入。

1.3 特征提取层

经过自定义特征词筛选和词嵌入技术后得到的短语向量作为卷积神经网络的输入, 进一步提取出该文本的高阶特征, CNN 和 RNN 主要由两个模块组成, 总体结构如图 2 所示。

卷积神经网络文本特征提取主要包括三个部分: 文本卷积、K-max 池化和相同位置向量整合。为避免宽卷积的补零操作, 使用窄卷积作为卷积策略^[10], 卷积核的选取使用长度为 2 的单卷积核。K-max 池化引用 Kalchbrenner 等^[4] 提出的动态卷积神经网络模型学习动态变化的句子, K-max 池化统一长度, 得到长度一致的特征图。相同位置向量整合的具体操作如图 2 中 K-max 池化后的虚线连接所示, 经过将相同位置的特征值合并, 形成中间特征向量 M , 作为循环神经网络的输入。

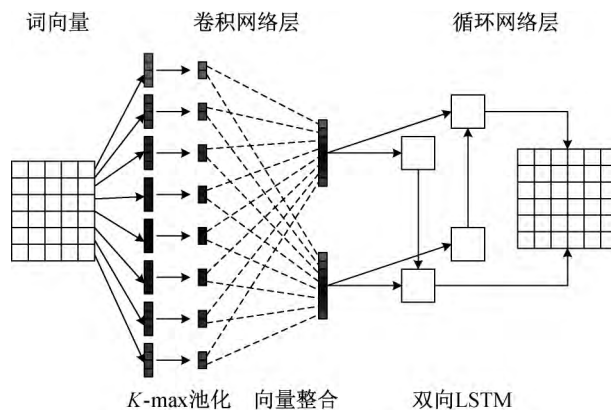


图2 卷积神经网络结构

循环神经网络文本特征提取采用了双向长短期记忆网络 (Bi-directional long short-term memory, Bi-LSTM) 作为循环网络层的实现^[11], Bi-LSTM 网络的整体流程如图 3 所示。

将中间特征向量 M 作为图 3 的输入, 因为经过相同位置向量整合后得到的中间特征向量 M 可以作为一个时间序列的向量表示, 经过 Bi-LSTM 的处理进一步提取到短文本的高阶特征。

1.4 注意力机制优化特征向量

将图 3 中的输出向量矩阵看成是某一篇文章文档经过卷积循环后的高阶向量特征, 注意力机制关注的是一篇文章文档中的某几个重要字词, 此时, 注意力关注的是该输出矩阵中的某几个重要向量, 本文选择 Soft Attention^[12] 实现方式, 其计算流程如图 4 所示, 将图 3 中得到的输出向量矩阵 $H (H = \{h_1, h_2, \dots, h_n\})$ 作为注意力机制的输入, 通过 tanh 函数将 H 的各个子向量进行转换, 如式 (4) 所示:

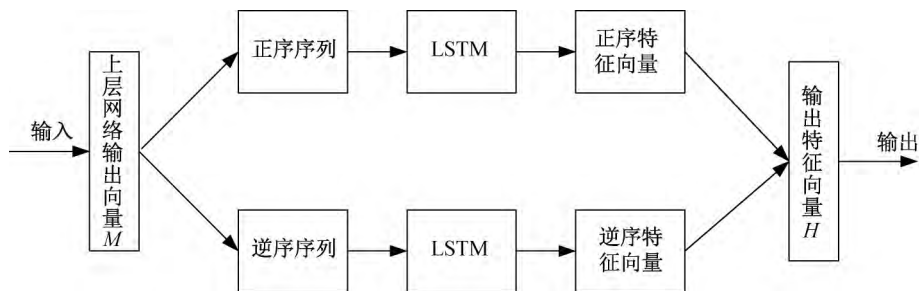


图3 Bi-LSTM网络流程图

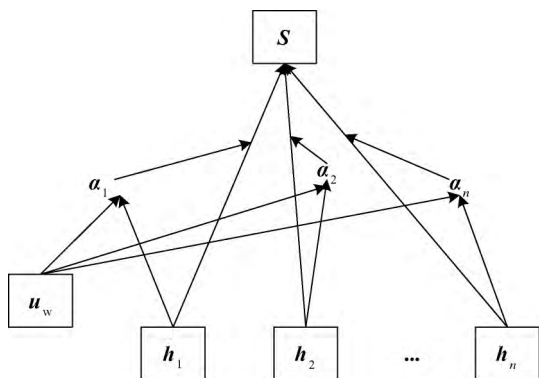


图4 Soft Attention计算流程

$$u_i = \tanh(W_h h_i + b_h) \quad (4)$$

为了获得注意力机制的权重系数,将通过式(4)计算得到的中间向量 u_i 与文本上下文向量 u_w 同时输入给 softmax 函数,计算出两者的相似度 α_i ,如式(5)所示。

$$\alpha_i = \text{softmax}(u_i, u_w) = \frac{\exp(u_i^T u_w)}{\sum_n \exp(u_i^T u_w)} \quad (5)$$

其中文本上下文向量 u_w 也可以称为记忆向量,作为

筛选图3中得到的输出向量矩阵中各个子向量重要性的抽象参数。通过式(5)计算出的相似度 α_i 作为归一化后的注意力权重,将得到的注意力权重和向量矩阵 H 中对应的子向量相乘,两者相乘后得到的特征向量 S_i 就是经过注意力机制优化后的最终文档向量特征,如式(6)所示,将其输入给全连接层得到最终的分类结果。

$$S_i = \sum_n \alpha_i h_i \quad (6)$$

2 实验结果分析

2.1 数据准备

本文使用 Sougou 语料库提供的多类别数据进行多分类实验,使用爬虫库在电影网站爬取一部分电影评论数据进行二分类实验,这两类数据集都没有验证集,本文随机选择了 10% 的数据集作为验证集,分别都在短语层和字符层进行分词操作,统计结果见表1。其中: $|W|$ 为中文文本分词后的短语个数, $|C|$ 为中文文本分词后单个字符个数。

表1 训练集测试集数据

数据集	类别	$ W $	$ C $	训练集	测试集	平均句子长度/个	最大句子长度/个
Sougou 语料库	9	62410	4852	48562	15950	89	112
电影评论数据	2	81204	6014	61248	27141	18	64

2.2 衡量指标

文本分类的评测标准可以采用准确率、精确率、召回率、F1 值作为指标,本文使用准确率作为衡量标准,在计算准确率时用到混淆矩阵进行解释,根据分类结果可建立混淆矩阵如表2所示。

表2 分类结果混淆矩阵

判别结果	真正属于该类别文档	真正不属于该类别文档
判别属于该类别文档	P	Q
判别不属于该类别文档	R	U

准确率在文本分类算法中表示的是分类正确的文档数除以整个训练集的文档总数,如式(7)所示:

$$f/\% = \frac{P+U}{P+Q+R+U} \times 100 \quad (7)$$

其中: f 表示准确率, P 表示模型判定的类别为真且该样本标注的类别也为真, Q 表示模型判定的类别为真且该样本标注的类别为假, R 表示模型判定的类别为假且该样本标注的类别为真, U 表示模型判定的类别为假且该样本标注的类别也为假。

2.3 结果分析

本文设计了4组对比实验,分别是传统文本卷积模型^[3](CNN)、传统 LSTM 模型^[13]、传统卷积循环网络模型^[14](CLSTM)和本文提出的混合神经网络模型。

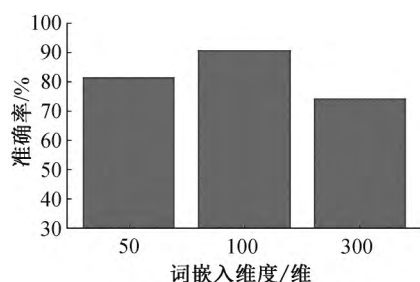
在实验过程中,为了做到实验的可对比性,所有

模型的参数设计都是一致的,显示如表3所示。表3中的模型参数取值这一列是取得最优结果时对应的参数值,模型参数实验范围这一列是实验过程中

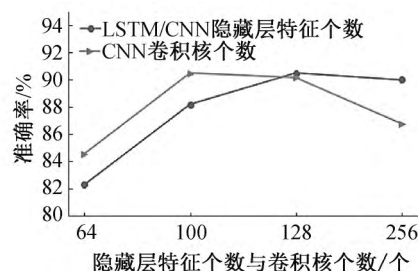
参数的可选值。模型参数取不同的值会得到不同的实验结果,在多分类数据集上模型参数选取不同值对准确率的影响如图5所示。

表3 模型参数设计

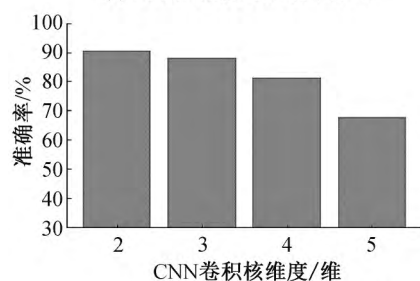
模型参数	模型参数最优取值	模型参数取值范围
词嵌入维度	100	50, 100, 300
LSTM/CNN 隐藏层特征个数	128	64, 100, 128, 256
CNN 卷积核个数	100	64, 100, 128, 256
CNN 卷积核大小	2	2~5
Dropout 随机失活率	0.4	0.1, 0.4, 0.5, 0.6, 0.9
模型迭代轮次	30	10, 30, 50, 100
学习率	0.001	0.100, 0.001, 0.010, 0.050
每批训练集数据大小	16	8, 16, 32, 64
自定义特征词筛选权重值	0.9	0.50, 0.60, 0.70, 0.80, 0.90, 0.95



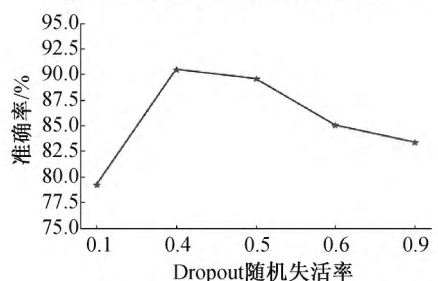
(a) 词嵌入维度对准确率影响



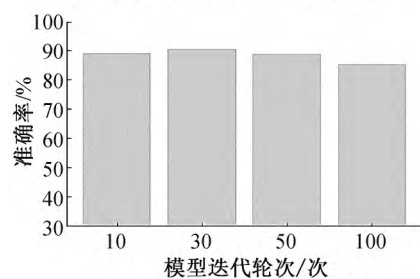
(b) 隐藏层特征与卷积核对准确率影响



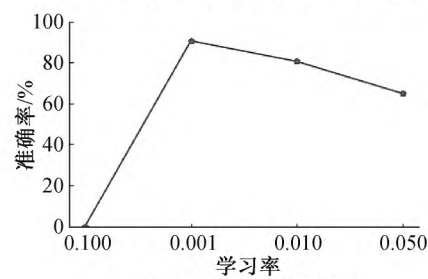
(c) 卷积核维度对准确率影响



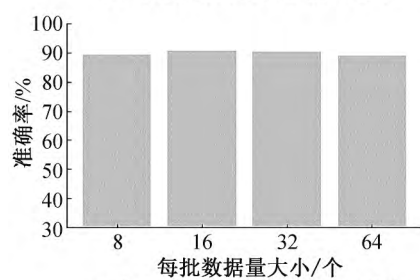
(d) Dropout随机失活率对准确率影响



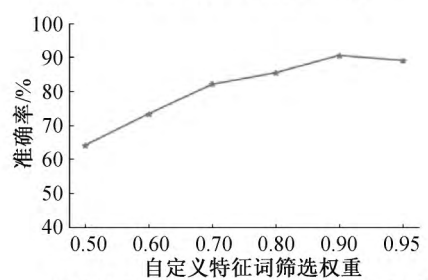
(e) 模型迭代轮次对准确率影响



(f) 学习率对准确率影响



(g) 每批数据量大小对准确率影响



(h) 自定义特征词筛选权重对准确率影响

图5 模型参数选不同值对准确率影响(多分类数据集)

图 5(a)~(h)中每个子图的最高值代表的是本文设计的模型在该参数选取时选到对应值时的最优结果(多分类数据集上)。词嵌入维度选择 100 时最优,当增大到 300,准确率会快速下降,表明当增大词向量维度时,在计算机中存储会过于稀疏,缺失精度。隐藏层和卷积核个数分别选到 128 和 100 时最优,当低于或高于该值时,准确率会下降,因为隐藏层和卷积核个数代表特征提取的大小,过低时对特征的提取就不够,存在欠拟合问题,过高时又会发生过拟合问题。CNN 卷积核大小在选取 2 时最好,对于中文短文本分类问题,文本在经过预处理以及自定义特征词筛选后剩下的词会减少,如果此时卷积核还取得比较大,卷积就不能起到作用。Dropout 失活率选到 0.4 时效果最优,但选择过低和过高时,准确率会大幅下降,选择过低时,模型在训练集上就不能很好的进行学习,会出现欠拟合问题,过高时又起不到正则化的作用,在测试集上泛化能力会很差。模型的迭代次数和每批数据量的大小从各子图发现,差别很小,实验过程中,分别取 30 和 16 时准确率最高。学习率选到 0.001 时结果最好,当选择 0.100 时,发现模型不能进行学习训练,这是因为学习率过大,用随机梯度下降法更新模型参数时,在接近局部最低点时会发散掉,导致模型参数无法更新。自定义特征词筛选权重设置在 0.90 时结果最好,当低于该值时,比如 0.50,起不到筛选的作用,0.60、0.70 等也是如此,当选到 0.95 时结果比选 0.90 差,因为当该值设定得太大,就相当于将某一篇文章特有的词直接删掉,没有考虑某些文档中特定词特有的属性,对结果也会造成一定影响。

二分类数据集上模型参数选取不同值对准确率的影响如图 6(a)~(h)所示。

图 6(a)~(h)中每个子图的最高值代表的是本文设计的模型在该参数选取时选到对应值时的最优结果(二分类数据集上)。通过图 6 的各子图和图 5 中的各子图对比分析,可以得出,不同的模型参数选取得到的结论是一样的,这也进一步验证了本文所选取得到的参数最优值不管在二分类数据集上还是多分类数据集上都是一致的。

在对比实验过程中,本文同时做了短语层和字符层实验,实验结果如表 4—表 5 中的前 4 行所示,最后做了将短语层和字符层混合实验,实验结果如表 4—表 5 中最后一行所示。

表 4 显示了在 Sougou 多分类数据集上该论文

模型与对比模型的详细比较,从第 1 行到第 4 行结果表明,基于短语层分词的结果比基于字符层分词的结果要好,而且从表 4 也可以看出,卷积循环相结合的 CLSTM 模型不管在短语层还是在字符层上都比单模型的传统 CNN 和传统 LSTM 效果更好,这也验证了本文的卷积循环文档特征提取是可以提取到更加合理的高阶向量特征。从第 4 行的结果来看,不管在短语层还是在字符层上又比 CLSTM 的结果好,这进一步验证了本文提出的基于注意力机制优化特征向量方法对卷积循环提取到的高阶特征进行优化是有效的。最后一行的结果显示,将经过注意力机制优化后的文档向量在短语层和字符层进行合并还可以进一步提升文本分类的准确率,从结果上看,虽然只比第 4 行的短语层分类准确率高了 0.60%,但总体来看这种将短语层和字符层相结合的方式可以提高中文短文本分类准确率。从表 4 的各模型分类准确率可以得出,本文模型比传统 CNN、LSTM 和 CLSTM 模型的准确率分别提高了 12.33%、4.16%和 2.33%。

电影评论二分类数据集的结果显示,整体的分类准确率要比在多分类数据集上好。从前 3 行的分类结果看,可以得出和多分类数据集上一样的结论,循环卷积相结合的 CLSTM 模型在短语层和字符层的分类准确率比单模型得到的准确率高,而且也都是短语层上的准确率普遍高于字符层上得到的准确率,从第 4 行的结果看,得到的高阶特征向量经过注意力机制优化后,在短语层和字符层上准确率可以大幅提升。最后从第 5 行的结果可以得出,将短语层和字符层分别得到的向量相结合的方式,输入给分类器得到的结果是最好的,可以达到 94.04%。从表 5 的各模型分类准确率可以得出,本文模型比传统 CNN、LSTM 和 CLSTM 模型的准确率分别提高了 10.36%、5.01%和 2.39%。所以,不管是在多分类数据集上还是二分类数据集上,本文提出的这种短语字符相结合的混合神经网络模型是行之有效的。

3 结 论

本文提出了一种基于中文短文本分类的混合注意力网络模型。该模型基于短语层和字符层将每一个句子切分成一个个单一的短语或字符,运用 word2vec 训练成词向量,将二者分别输入给卷积循环网络进行高阶特征的提取,再将提取到的高阶向量经过注意力机制优化,最后将短语层和字符层优

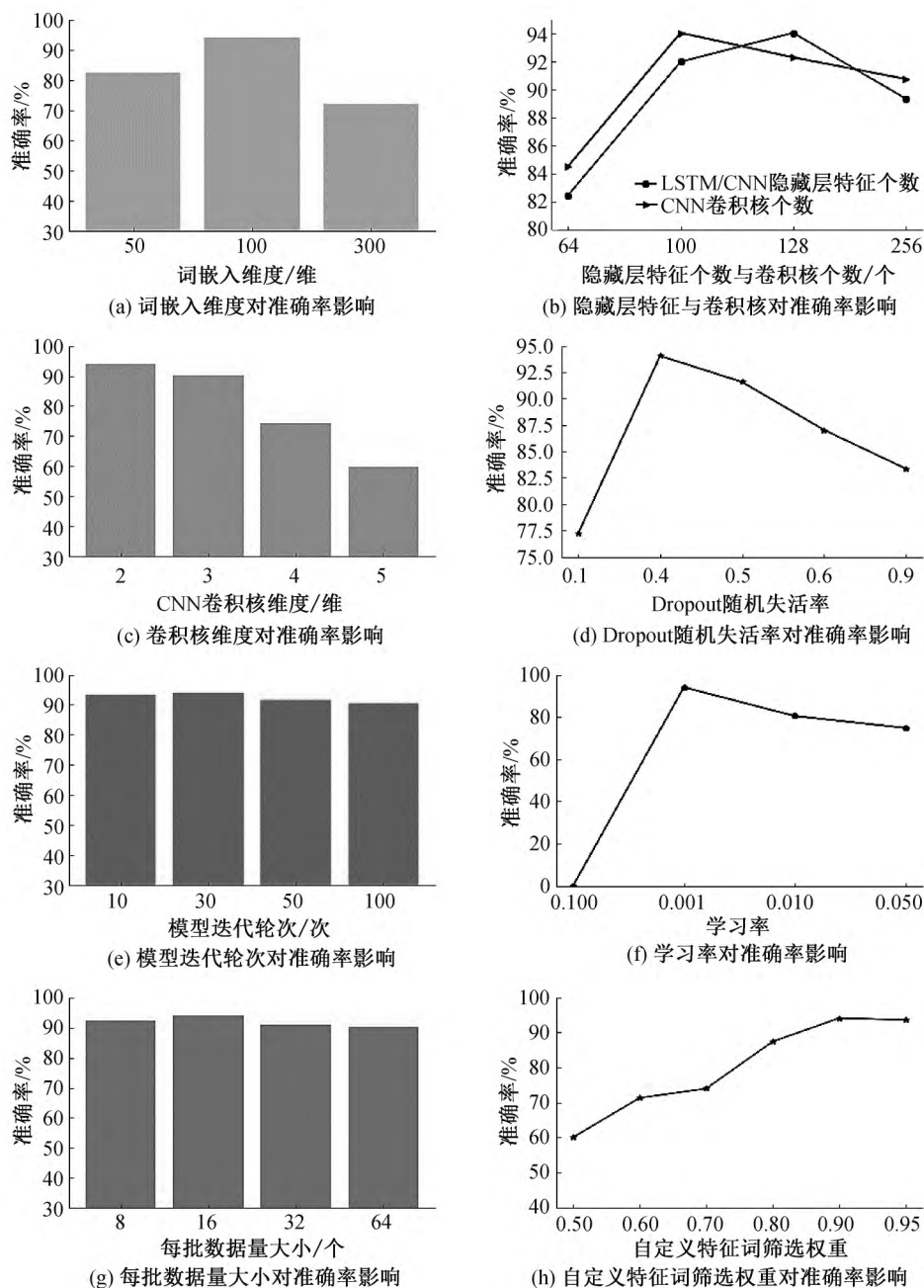


图 6 模型参数选不同值对准确率影响(二分类数据集)

表 4 Sougou 多分类数据集分类准确类

模型	分类准确率/%	
	短语层	字符层
传统 CNN	78.32	71.06
传统 LSTM	84.95	80.76
CLSTM	86.14	83.24
本文模型(短语字符分开)	89.87	84.17
本文模型(短语字符混合)	90.47	

表 5 电影评论二分类数据集分类准确率

模型	分类准确率/%	
	短语层	字符层
传统 CNN	82.79	78.43
传统 LSTM	86.36	85.13
CLSTM	89.24	87.91
本文模型(短语字符分开)	92.47	89.46
本文模型(短语字符混合)	94.04	

化后的文档向量相结合,得到文档最终的向量表示,输入给分类器层得到分类结果。实验对比结果显示,本文设计的模型比传统的 CNN 和 LSTM 模型

在准确率上有比较高的提升,表明将卷积循环网络相结合的方式进行高阶特征提取是可行的。与传统的 CLSTM 模型对比,运用注意力机制优化特征向

量还能将准确率进一步提升,最后将短语层和字符层进行合并,得到的分类准确率是最好的。

由于本文的方法虽然在深度学习上使用了 tensorflow 的 gpu 加速功能,但对于文本处理还需很长时间,面对今后海量的数据分类实用性比较低。因此,如何采用分布式平台进行基于深度学习的互联网短文本分类将是笔者的研究重点,该研究不仅能在分类精度上做到显著提高,在分类速度上也可以提高。

参考文献:

- [1] Sebastiani F. Machine learning in automated text categorization[J]. ACM Computing Surveys, 2001, 34(1):1-47.
- [2] 王义真,郑啸,后盾,等. 基于 SVM 的高维混合特征短文本情感分类[J]. 计算机技术与发展, 2018, 28(2), 88-93
- [3] Kim Y. Convolutional neural networks for sentence classification[EB/OL]. (2014-09-03) [2018-12-17]. <https://arxiv.org/abs/1408.5882>.
- [4] Kalchbrenner N, Grefenstette E, Blunsom P. A convolutional neural network for modelling sentences[EB/OL]. (2014-04-08) [2018-12-17]. <https://arxiv.org/abs/1404.2188>.
- [5] 黄文明,莫阳. 基于文本加权 KNN 算法的中文垃圾短信过滤[J]. 计算机工程, 2017, 43(3):193-199.
- [6] 黄磊,杜昌顺. 基于递归神经网络的文本分类研究[J]. 北京化工大学学报(自然科学版), 2017(1):100-106.
- [7] Wang C, Xu B. Convolutional neural network with word embeddings for chinese word segmentation[EB/OL]. (2017-12-13) [2018-12-17]. <https://arxiv.org/abs/1711.04411>.
- [8] Zhou Y, Xu B, Xu J, et al. Compositional recurrent neural networks for Chinese short text classification[C]//Web Intelligence (WI), 2016 IEEE/WIC/ACM International Conference on. IEEE, 2016: 137-144.
- [9] Xu C Z, Liu D. Chinese text summarization algorithm based on Word2vec[C]//Journal of Physics: Conference Series. IOP Publishing, 2018, 976(1): 012006.
- [10] Qu S, Xi Y, Ding S. Visual attention based on long-short term memory model for image caption generation[C]//Control And Decision Conference (CCDC), 2017 29th Chinese. IEEE, 2017: 4789-4794.
- [11] Hochreiter S, Schmidhuber J. Long short-term memory[J]. Neural computation, 1997, 9(8): 1735-1780.
- [12] Xu K, Ba J, Kiros R, et al. Show, attend and tell: Neural image caption generation with visual attention[C]//International Conference on Machine Learning. Lille, France, 2015: 2048-2057.
- [13] Zhou Y, Xu B, Xu J, et al. Compositional recurrent neural networks for Chinese short text classification[C]//Web Intelligence (WI), 2016 IEEE/WIC/ACM International Conference on. IEEE, 2016: 137-144.
- [14] Zhou C, Sun C, Liu Z, et al. A C-LSTM neural network for text classification[EB/OL]. (2015-12-30) [2018-12-17]. <https://arxiv.org/abs/1511.08630>.

(责任编辑:康 锋)