

基于多输入模型及句法结构的中文评论情感分析方法

张宝华¹, 张华平¹, 厉铁帅², 商建云¹

1. 北京理工大学计算机学院, 北京 100081;
2. 中央军事委员会政法委员会, 北京 100120

摘要

海量的网络文本给情感分析任务带来了巨大的机遇和挑战, 传统基于规则的方法已经很难胜任这类文本的分析工作, 现有的深度学习方法存在一些不足, 一方面模型的输入只包括文本嵌入矩阵, 缺乏其他特征的使用; 另一方面, 词嵌入算法会导致文本结构信息缺失, 进而影响分析效果。在对基于规则的情感分析方法中的句法规则进行研究的基础上, 提出了一种结合MCNN、LSTM和全连接神经网络的多输入模型。同时在深度学习模型中构建了句法特征提取器来提取句法特征。在3个公开数据集上进行了实验, 结果表明, 构建的模型较其他模型拥有更好的分类性能, 且句法规则特征的引入对模型分类效果有一定的提升。

关键词

情感分析; 句法规则; 多输入模型

中图分类号: TP183

文献标识码: A

doi: 10.11959/j.issn.2096-0271.2021059

Chinese comment sentiment analysis method based on multi-input model and syntactic structure

ZHANG Baohua¹, ZHANG Huaping¹, LI Tieshuai², SHANG Jianyun¹

1. School of Computer Science & Technology, Beijing Institute of Technology, Beijing 100081, China
2. Politics and Law Commission of Central Military Commission of the People's Republic of China, Beijing 100120, China

Abstract

Massive network texts have brought huge opportunities and challenges to sentiment analysis tasks. Traditional rule-based methods have been difficult to analyze such texts. Existing deep learning methods have some shortcomings. On the one hand, the inputs of the model only include the text embedding matrix, lack the use of other features. On the other hand, the algorithm of word embedding will lead to the lack of text structure information, then impact the result. Based on the research of syntactic rule in the rule-based sentiment analysis methods, a multi-input model combined with MCNN, LSTM and fully connected neural network was proposed. Meanwhile, a syntactic feature extractor to combine the syntactic features was constructed in the deep learning model. Experiments on three public data sets were conducted. The results

show that the model constructed in this article has better classification performance than other models, and the introduction of syntactic rule features has a little improvement in the classification effect of the model.

Key words

sentiment analysis, syntactic rule, multi-input model

1 引言

随着智能电子设备的普及和网络的发展,大量的社交媒体及电商平台开始走入人们的生活,人们在日常使用过程中会产生海量的评论数据。合理利用观点挖掘技术可以从这些数据中获取巨大的价值,如对电商评论数据进行观点挖掘,可以分析得出商品的优缺点,商家可以对其进行修改;对影评数据进行挖掘,可以看到当前电影的缺点和优点,方便用户进行选择,也方便出版方的宣传工作;对新闻评论区等数据进行挖掘,可以掌握当前群众的态度。情感分析作为观点挖掘的主要技术之一,面临着巨大的挑战。

新词频出、长短不一、结构不定是网络评论数据的主要特点。严重依赖情感词典的传统规则情感分析方法一方面由于情感词典中缺少网络新词,无法得到新词的正确情感权重,在计算时只能忽略这部分词;另一方面又因为网络文本的结构不定,使用现有的规则对网络文本结构进行分析会出现一定的误差,导致该方法对这类数据的分析效果很差。而基于深度学习的分析方法虽然在分析效果上较好,但也存在一些问题。首先,单一的神经网络模型在处理文本数据时会因为模型本身存在的结构缺陷,造成部分情感特征的损失,从而导致分析准确率较低;其次,现有神经网络模型需要将文本数据映射到向量空间,构建文本向量矩阵之后进行模型运算,但是在这个过程中丢失了在传统方法中对句

子情感有很大影响的结构信息;最后,现有的情感分析模型已经开始研究如何将规则方法中使用的一些特征加入深度学习方法中,但是其仍然以输入全部文本数据为主,缺少对其他特征的提取。

虽然基于规则的方法的分析效果要弱于深度学习方法,但是基于规则的方法中的句法规则特征在深度学习方法中仍然具有很重要的作用。因此,需要构建一种新的模型,将这部分句法规则特征融入神经网络模型。为了结合不同模型的优点,并将句法结构规则引入深度学习模型中,本文构建了基于多输入模型及句法结构的神经网络模型。该模型同时将文本向量、情感词向量和语法规则向量输入独立的神经网络模型中,并对模型提出的特征进行拼接,从而得到更加全面的文本特征。实验证明,本文提出的神经网络模型较其他模型的效果更好。

本文主要有以下贡献:

- 本文提出了一种结合多通道卷积神经网络(multi-channel convolutional neural network, MCNN)、长短期记忆(long short-term memory, LSTM)网络和全连接神经网络的合并模型MCNN_S_LSTM_NN,该模型可以结合单个模型的优势,从文本评论中获取更全面的情感特征,从而提高情感分析的准确率;
- 本文针对每部分模型的特点,设计了不同的模型输入,可以从不同的角度对文本进行特征提取;
- 本文首先将句法结构、标点符号等在基于情感词典的情感分析方法中会用到

的分析规则应用到深度学习中。同时,本文构建了句法规则提取器,可以直接对文本规则进行提取,并映射到向量空间,作为深度学习模型的输入。

2 相关工作

深度学习方法最早由Collobert R等人^[1]在2011年应用到自然语言处理领域,用于解决词性标注等问题。2014年, Kim Y^[2]首先在文本分类方面使用卷积神经网络(convolutional neural network, CNN),并且取得了很好的分类效果。之后Kalchbrenner N等人^[3]提出了一种宽卷积模型,并选择用 k -max池化代替传统CNN的最大池化来保留更多的特征。Zhang Y L等人^[4]通过多次重复实验,比较了不同超参数对CNN模型结构在性能和稳定性方面的影响。Gao J等人^[5]和Shen Y L等人^[6]介绍了如何将句子表示成包含语义的结构。Zhang R等人^[7]提出了可有效获取句子依赖信息的CNN模型,通过处理预训练的词嵌入来构建分层的文本表示。CNN常被用于捕获局部特征,而循环神经网络(recurrent neural network, RNN)由于自身存在反馈环结构,可以保留记忆信息,在时间序列模型中得到了很好的应用^[8]。但是RNN自身存在一定的缺陷,当文本长度增加时,梯度消失和梯度爆炸情况的出现会导致分析效果不理想。LSTM和门控循环单元在传统RNN的基础上引入了门机制,较好地解决了RNN的问题。Socher R等人^[9]通过构建Tree-LSTM获取到更多的文本特征。Tran K等人^[10]为了提升模型对历史信息处理能力,在LSTM的基础上引入了外部记忆单元,但是由于增加了大量的参数,模型准确度提升不大。Chen P等人^[11]使用具有注意力机制的双向长

短期记忆(bidirection long short term memory, BiLSTM)网络获得了较好的分类效果。宋婷等人^[12]建立了分层的LSTM模型,用于提取方面级的情感。Wang Y Q等人^[13]通过对LSTM建模、对上下文建模,结合文本隐藏状态和方面级情感分析中的方面信息生成注意力向量,并建立了AE-LSTM和ATAE-LSTM神经网络模型,最后得到方面级的情感分类模型。LSTM也存在缺点,其虽然能获得文本的上下文语义信息,但是缺少对文本局部信息的获取。

联合多个简单模型的神经网络模型逐渐成为情感分析方法的主流, Cheng Y等人^[14]建立了典型的并行双层网络结构,输入的数据会先经过注意力机制的计算,计算得到的结果分别被输入多通道卷积神经网络(multi-channel CNN, MCNN)和双向门控循环单元(bidirection gated recurrent unit, BiGRU),接着将两个模型得到的特征合并,最后对特征进行训练。Yang L等人^[15]使用串行结构将CNN和BiGRU进行组合,并在最后一层加入了注意力机制。Li W等人^[16]使用并行结构将BiLSTM和CNN组合起来。不同的是,其提出了使用情感词进行填充的方法,可以缓解梯度消失问题。Li W J等人^[17]虽然没有使用CNN,但是并行了3个BiLSTM,将文本的输入向量分别与词性向量、位置向量和依存语法向量拼接,并分别作为3个模型的输入;然后利用注意力机制将3个特征进行拼接。通过拼接这3个向量可以将句子的部分语法规则特征输入模型中,从而提高模型的准确率。Usama M等人^[18]使用了串行的方法,将RNN和CNN合并起来,并在两个模型中间加入了注意力机制。Basiri M E等人^[19]提出了新的既包含并行结构又包含串行结构的模型。该模型首先将BiLSTM和BiGRU进行并行处理,

之后利用CNN进行卷积操作,提取局部特征并降低特征维度。可以看到,并行网络由于可以综合单一网络模型的优点,已经成为当前的分析主流。Jin N等人^[20]提出了MLT-MSCNN-LSTM,该模型同样将MCNN和LSTM网络作为基础模型。此外,该模型还提出了融合网络(fusion net),用于合并MCNN不同卷积核的输出。该模型首先训练word-embedding矩阵,同时作为MCNN和LSTM两个模型的输入;然后分别用LSTM提取全局特征、用MCNN提取局部特征,再使用融合网络合并;最后将两个模型的输出进行拼接得到句子的最终表示特征。在融合网络中,Jin N等人^[20]首先将3个经过最大值池化后的特征拼接,然后使用全连接神经网络获取特征,同时在全连接神经网络中使用dropout机制来提高模型的收敛速度。Li W J等人^[17]提出了SAMF-BiLSTM模型,该模型在LSTM上加入了自注意意义力(self-attention)机制和层正则化(layer normalization),该模型一共包含5层,且通过3个BiLSTM从不同方面提取文本的情感特征。该模型在词向量层做了改进,在训练好的词向量上拼接了语音特征向量、位置值向量和依赖解析向量,并将拼接向量作为模型的输入。Basiri M E等人^[21]提出了ABCDM(attention-based bidirectional CNN-RNN deep model)。该模型主要采用注意力机制,首先将文本数据转换为向量矩阵;然后通过BiLSTM和BiGRU两个模型同时提取文本特征,提取完成后引入注意力机制;接着在最后一层加入卷积神经网络进行卷积,提取相关的局部特征并降低特征维度;最后将输出合并,得到最终的特征向量,利用全连接神经网络进行最后的预测。Usama M等人^[22]提出了基于注意力机制的卷积神经网络和循环神经网络情感分析(ATTConv RNN-

rand)模型。该模型使用了双通道的卷积神经网络;然后对经过卷积神经网络计算后的文本特征进行注意力计算,并将其输入循环神经网络中;最后利用全连接神经网络进行学习和训练。Li W J等人^[17]提出了情感词填充的卷积神经网络和长短期记忆网络(CNN_BiLSTM_sentiment_padding)模型,创新性地使用情感词填充代替数字0填充。在对句子长度较短的句子进行填充时,作者使用句子中情感词权重绝对值的大小,确定句子中不同情感词的填充个数,然后使用情感词填充,这可以强调句子中的情感信息。作者构建了并行的BiLSTM和MCNN模型,分别使用两个模型对句子特征进行提取,再将两个模型得到的特征合并,在最后一层采用全连接神经网络进行学习和预测。

但是,上述模型中,只有Yang L等人^[15]使用了部分的语法规则,其他模型都忽略了对语法规则的输入;而Cheng Y^[14]等人的结论中提到,需要考虑将传统方法中的句法结构特征融入深度学习模型中,这也是本文主要研究的内容。

3 模型分析

如图1所示,本文的模型主要包括MCNN、LSTM和全连接神经网络。其中,MCNN的输入为整个句子的词向量矩阵,通过构建多通道CNN提取文本特征;LSTM的输入为句子中的情感词的词向量矩阵。特征提取器主要提取句子中的句法结构信息,并将其映射到向量空间作为全连接神经网络的输入。将这3个网络得到的特征进行融合,然后利用全连接神经网络输出句子的情感类别。

本文模型分为3个部分,MCNN、LSTM和全连接神经网络。其中,MCNN

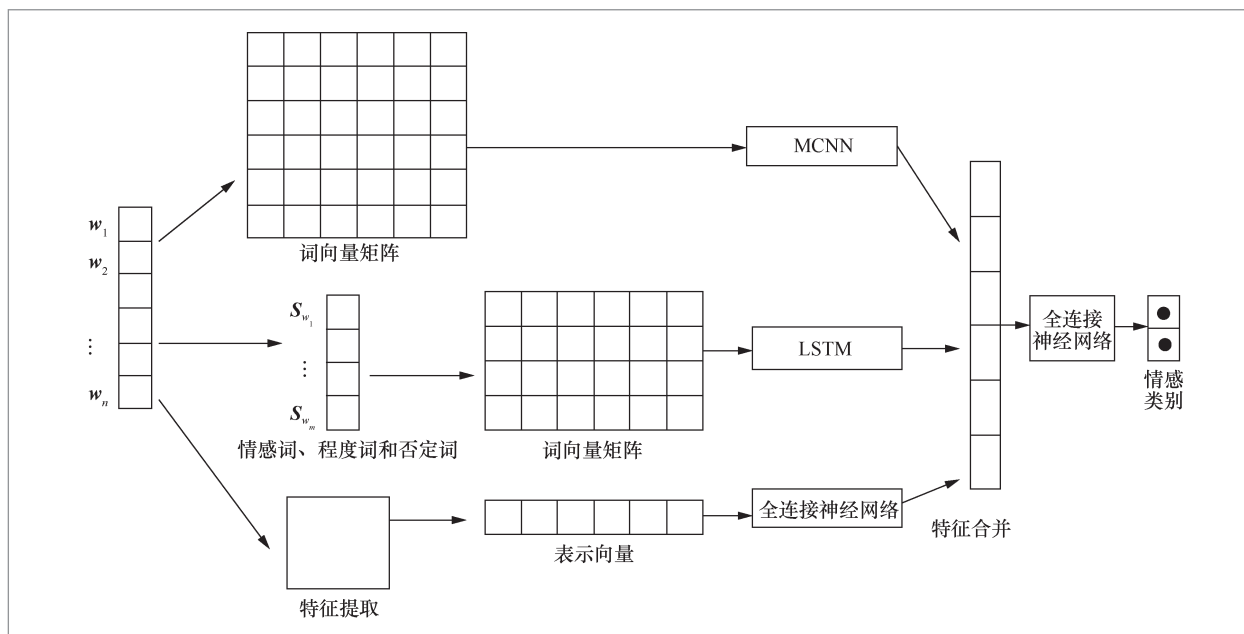


图1 模型的整体架构

的输入为整个句子的词向量矩阵。对于训练集中的数据,首先使用word2vec算法计算其词向量,然后通过MCNN来提取文本特征。LSTM的输入为句子中包含的情感词的词向量矩阵,LSTM主要用来提取文本中的情感词特征。最后一部分的神经网络模型是全连接神经网络,其输入是特征提取器从文本中提取到的句法特征。特征提取器主要包含两个功能:特征提取和特征映射。特征提取是指从文本中提取其包含的句子结构特征、句间关系特征和标点符号特征,详细内容将在第3.1节中介绍。提取到特征之后,需要将这部分特征在特定规则下映射到向量空间,从而输入神经网络模型中。在这三部分特征全部提取结束后,进行拼接处理。最后利用全连接神经网络对整体特征进行学习。

3.1 输入层

本文的模型是由3个基本神经网络模型构成的。对于每个基本神经网络模型,

本文构建了不同的输入。对于MCNN,根据CNN可以提取局部特征的能力,选择将整个句子的词向量矩阵作为输入。对于LSTM模型,之前的研究^[23]证明了在LSTM模型中只使用情感词代表句子的情感特征可以获取比使用整个句子作为输入更好的效果。因此,本文选择使用句子中包含的情感词作为输入,在将句子输入之前,先使用情感词典构建方法对数据集构建情感词典,然后提取句子中的情感词,构建情感词的词向量矩阵并作为输入。对于全连接神经网络,本文设计的输入是句子的句法规则特征。本文主要使用了如下3种规则:结构规则、句间关系规则、句型规则。

- 结构规则:根据句子中是否包含多个单句可以将句子分为复句和单句。复句是由多个单句构成的,复句的情感权重由单句的情感权重根据一定的规则累加得到,单句的情感权重就是本身的权重。因此,要先根据句内标点、关联词等将句子标定为单句或复句,若是复句,则根据单句的个数将包含的单词映射到向量空间。

- 句间关系规则：这里的句间关系规则主要有4种，转折关系、递进关系、因果关系和假设关系。在转折关系中根据转折词的不同，可以将转折句分为转折前句和转折后句，转折前后句的情感极性相反，且强调后句的情感；在递进关系中，前后递进句的情感逐渐增强，且更强调后递进句的情感；在因果关系中，更强调原因的情感；在假设关系中，更强调条件，对后假设句的情感有削弱。基于此，本文将每个单句中的关联词映射到向量中作为输入。

- 句型规则：根据结尾标点符号的不同，主要分为陈述句、感叹句和疑问句。其中，以句号结尾的为陈述句，句子情感值不变。以叹号结尾的为感叹句，句子的情感值增强。以问号结尾的为问句，根据是否有反义疑问词，可分为反义疑问句和问句，问句表示无情感，而反义疑问句则强调句子的反向情感。因此根据结尾符号和是否包含反义疑问词，可以将句型映射到向量空间中。具体步骤如下。

(1) 初始化语法规则特征向量 $GV=[0]$ 。

(2) 根据标点符号和句子中是否有连接词判断是否为复句，若是，则将 GV 的第一位置为1，进行下一步；否则根据标点符

号，将 GV 的对应位置置1。

(3) 根据关联词和标点符号对句子进行切分，并根据切分出来的单句数量，在 GV 剩下的位置上置1。

(4) 对切分的每个单句进行关联词匹配和标点符号匹配，并以12位的向量空间标记每条单句的处理结果。其中转折前置词、假设前置词、因果前置词和强调前置词由前4位表示，后关联词由中间4位表示，标点符号类型由后4位表示，每个单句都会得到类似的12位标记，然后按顺序组成 GV 。

(5) 填充0，将 GV 补足到 N 位，其中 N 表示全连接层的输入长度限制。

由此，针对不同基本模型的特点，构建了不同的输入，并将语法规则特征映射到向量空间作为全连接神经网络的输入，通过全连接神经网络将其补充到最终模型的特征中。

3.2 模型架构

本文提出的基本模型主要有MCNN和LSTM。

3.2.1 卷积神经网络层

MCNN 的架构如图2所示，其主要由

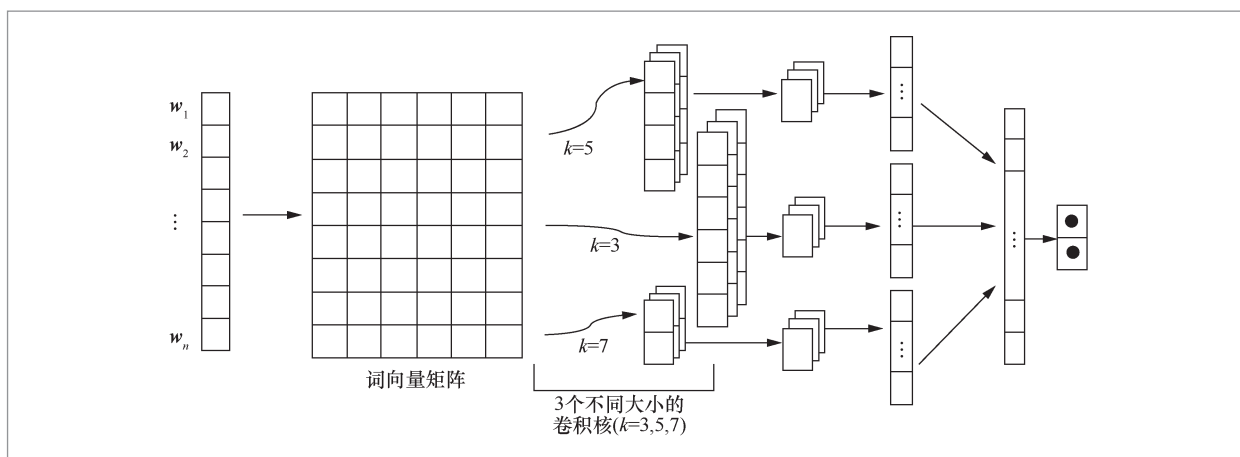


图2 MCNN的架构

卷积层、池化层和全连接层构成。卷积层用来提取输入数据的特征；在卷积层进行特征提取后，输出的特征图会被传递到池化层进行特征选择和信息过滤；全连接层等价于传统前馈神经网络的隐藏层，一般和输出层连接，实现最后的输出。

在本文的模型中，假设句子的最大长度为 N ，其中不足 N 的用 0 补齐；词向量维度为 d ，则输入句子可以用矩阵 $S \in R^{n \times d}$ 表示。假设卷积核 $W \in R^{d \times h}$ ，其中， d 表示卷积核的长度，大小和词向量的维度相同； h 表示卷积核的宽度。本文选择的卷积核大小分别为 3×3 、 5×5 、 7×7 。对于输入 $S \in R^{n \times d}$ ，通过卷积操作可得特征向量 $O = (O_0, O_1, O_2, \dots, O_{n-h}) \in R^{n-h+1}$ ，则 O 中元素的计算式为 $O_i = W \cdot S_{i:i+h-1}$ ，其中 $i=1, 2, 3, \dots, n-h$ ，符号“ \cdot ”表示矩阵的点乘操作。 $S_{i,j}$ 表示矩阵 S 的第 i 行到第 j 行的子矩阵。在卷积神经网络中，卷积完成之后一般进行池化操作，在池化层中使用最大池化操作，在每个过滤器中都可以取得最大值，可以提取到最显著的特征。Zhang Y 等人^[4]的研究工作也表明，在各种句子分类任务中，最大池化操作在性能上始终优于其他

池化策略。因此，这里选用 1-max 池化，其主要思想是通过选择特定特征图的最大值来捕获与特定特征图对应的最重要特征 $v = \max_{0 \leq i \leq s-h} \{O_i\}$ 。如图 2 所示，经过不同大小的卷积核以及同样的池化操作后，提取出不同大小的特征，将这些特征进行合并，然后传递给以 Sigmoid 为激活函数的全连接层，就可以得到不同情感类别的概率。

3.2.2 长短期记忆网络

LSTM 的架构如图 3 所示。RNN 可以处理一定的短期依赖，但是由于序列较长时，序列后部的梯度很难反向传播到前面的序列，因此无法处理长期依赖问题。而在 LSTM 中引入了细胞状态，并使用输入门、遗忘门和输出门来保存和控制信息，可以解决 RNN 的缺点。

LSTM 的某个状态有以下步骤。

(1) 当输入门 i_t 接收到当前输入 x_t 和最后的最终隐藏状态 h_{t-1} 之后， i_t 通过以下计算式进行计算： $i_t = \sigma(w_{ix}x_t + w_{ih}h_{t-1} + b_i)$ 。其中， σ 是一个逻辑 S 型函数， w_{ix} 和 w_{ih} 分别代表两个权重矩阵， b_i 是输入门的偏差向量。

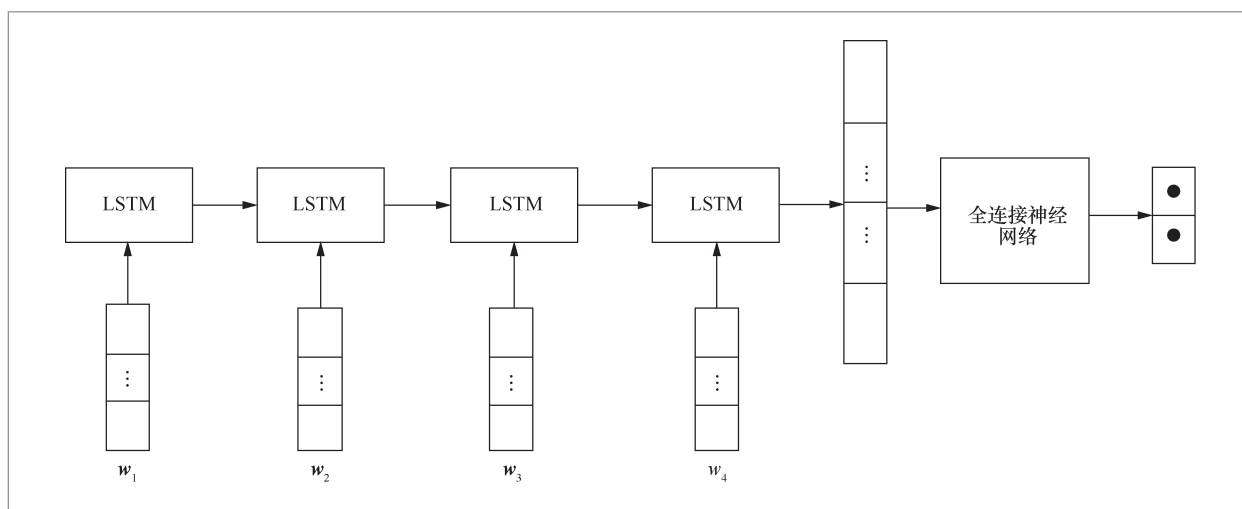


图 3 LSTM 的架构

(2) 若 i_t 的值为 1, 则表示当前输入的信息可以进入细胞状态; 若值为 0, 则表示当前输入的信息不可以进入细胞状态。然后计算候选值 $\tilde{c}_t = \tanh(\mathbf{w}_{cx}x_t + \mathbf{w}_{ch}h_{t-1} + \mathbf{b}_c)$, 其中 \mathbf{w}_{cx} 和 \mathbf{w}_{ch} 分别代表两个权重矩阵, \mathbf{b}_c 是偏差向量。

(3) 之后遗忘门 f_t 将执行操作 $f_t = \sigma(\mathbf{w}_{fx}x_t + \mathbf{w}_{fh}h_{t-1} + \mathbf{b}_f)$ 。 f_t 的值为 0 表示不传递信息 c_{t-1} 给 c_t , f_t 的值为 1 表示将全部信息传递给 c_t 。其中 \mathbf{w}_{fx} 和 \mathbf{w}_{fh} 分别代表两个权重矩阵, \mathbf{b}_f 是遗忘门的偏差向量。

(4) 计算完成后的细胞状态为 $c_t = f_t c_{t-1} + i_t \tilde{c}_t$ 。

(5) LSTM 的最后状态 $h_t = o_t \tanh(c_t)$, 其中 $o_t = \sigma(\mathbf{w}_{ox}x_t + \mathbf{w}_{oh}h_{t-1} + \mathbf{b}_o)$ 。

最终隐藏状态的输出与前一个序列的隐藏状态、当前输入和当前单元状态值有关, 用一个激活函数 \tanh 将当前单元状态的值压缩到 $-1 \sim 1$ 。将先前隐藏状态与当前输入通过 sigmoid 函数转换后的值与当前单元状态经过压缩后的值进行相乘, 保留或舍弃先前的状态信息与此时的输入信息, 从而得到一个新的隐藏状态值。

3.3 输出层

在基本模型都提取到文本的情感特征之后, 将 3 个模型得到的特征进行合并, 这可以取长补短, 综合不同模型的特点, 同时可以将文本的句法规则特征加入文本的最

终特征中, 实现句法规则特征与深度学习模型的融合。最后通过全连接神经网络对最终的文本情感特征进行分类。

4 实验结果及分析

本文在 Linux 环境下使用 Python 2.7 和 Keras 完成了模型的编写, 并在 8 块型号为 1080Ti 的 GPU 上进行了训练。从第 3 节可以看出, 本文的 3 个模型的输入各不相同。其中 LSTM 的输入为情感词, 因此需要先利用情感词典构建方法^[18]构建情感词典, 也可以采用公开的情感词典, 但是效果不如针对数据单独提取的词典。

4.1 数据集

本文的数据集主要使用谭松波酒店评论数据 (Hotel)、NLPCC2014 情感分析 task2 以及笔者从豆瓣上采集的影评数据 (Douban), 将好评认为正面情感, 差评认为负面情感。对 3 个语料库的统计结果见表 1。

4.2 消融实验

为了验证本文模型对最终结果的贡献, 本文设计了消融实验, 实验结果见表 2。

其中, MCNN_S_LSTM_NN 是本文提出的多输入模型, MCNN 表示多通道卷积神经网络, S_LSTM 表示以 LSTM 模型

表 1 对 3 个语料库的统计结果

数据集	总评论数/个	句子最大长度/字符	句子最小长度/字符	平均包含单句数量/个	包含关联词的句子数量/个
Douban	100 000	99	11	4.32	56 998
Hotel	10 000	1 985	4	8.74	8 276
NLPCC2014	10 000	1 004	3	5.48	6 419

为基础且输入为句子中包含的情感词的词向量矩阵, LSTM表示以LSTM模型为基础且输入为整个句子的向量矩阵, NN表示以全连接神经网络为基础且输入为通过句法规则特征提取器提取的句法规则向量。从表2可以看到, 带有NN的模型要比不带NN的模型的准确率高, 说明句法结构特征的引入丰富了情感特征, 从而提高了情感分析的准确率。S_LSTM和LSTM的对比则说明句子中的情感词对句子的情感有极大的影响。

4.3 对比实验

本文主要选取CNN、LSTM、SLCABG^[14]、ATTConv RNN-rand^[17]、ABCDM^[18]、SAMF-BiLSTM^[15]、CNN-BiLSTM(sentiment word padding)^[16]、MC-AttCNN-AttBiGRU^[13]作为对比模型, 将精确率(P)、召回率(R)、F1值和准确率(ACC)作为评测指标, 其实验结果见表3。

从表3可以看到, 本文构建的模型MCNN_S_LSTM_NN在3个数据集上的

表2 消融实验的准确率

模型	Hotel	Douban	NLPCC2014
MCNN_S_LSTM_NN	91.72%	86.42%	76.65%
MCNN_LSMT_NN	91.59%	84.55%	75.49%
MCNN_S_LSTM	91.45%	84.61%	75.21%
MCNN_LSTM	91.01%	84.29%	74.20%
MCNN_NN	90.54%	84.52%	72.54%
S_LSTM_NN	90.51%	83.94%	72.90%
MCNN	90.35%	83.55%	71.25%
S_LSTM	89.90%	83.70%	73.50%
LSTM	88.40%	83.64%	71.55%

准确率均最高, 这证明了本文构建的模型的可行性和先进性。另外, 从表3还可以看到, 组合模型的准确率相较于简单模型要更高, 这说明复合模型可以综合简单模型的优点, 可以提取到更全面的文本特征。在模型中使用注意力机制的方法(如MC-AttCNN-AttBiGRU、ABCDM、ATTConv-RNN rand)的准确率也比较高, 这也许是本文可以继续学习和研究的一个方向。

表3 对比实验的实验结果

模型	Hotel				Douban				NLPCC2014			
	P	R	F1	ACC	P	R	F1	ACC	P	R	F1	ACC
CNN	88.18%	89.40%	88.78%	88.61%	82.61%	86.17%	84.36%	83.94%	78.32%	67.67%	72.57%	73.56%
LSTM	87.79%	89.70%	88.74%	88.40%	86.01%	80.54%	83.23%	83.64%	75.92%	68.45%	72.00%	72.65%
MCNNALSTM	90.61%	90.79%	90.70%	90.52%	85.38%	84.16%	84.77%	84.16%	74.36%	75.50%	74.93%	73.80%
SLCABG	91.35%	88.91%	90.12%	90.15%	86.13%	83.89%	85.00%	85.11%	79.08%	71.46%	75.08%	75.40%
ATTConv-RNN rand	92.54%	88.42%	90.43%	90.55%	86.19%	84.85%	85.51%	85.55%	81.70%	67.60%	73.98%	75.35%
ABCDM	89.48%	92.67%	91.05%	90.80%	84.20%	86.91%	85.53%	85.22%	78.59%	73.96%	76.20%	76.05%
SAMF-BiLSTM	93.25%	87.52%	90.30%	90.50%	84.72%	85.49%	85.10%	84.96%	74.82%	78.20%	76.47%	75.05%
CNN-BiLSTM	92.32%	88.12%	90.17%	90.30%	86.13%	83.89%	84.77%	84.79%	78.32%	71.07%	74.52%	74.80%
MC-AttCNN-AttBiGRU	91.35%	90.99%	91.17%	91.10%	85.53%	86.80%	86.16%	85.99%	75.99%	78.11%	77.03%	75.85%
MCNN_S_LSTM_NN	90.74%	93.08%	91.87%	91.70%	88.27%	84.17%	86.17%	86.42%	77.56%	77.34%	77.45%	76.65%

5 结束语

本文构建了基于多输入模型及句法结构的中文评论情感分析方法,通过研究基于情感词典的情感分析方法,将情感分析中对句子情感会产生重要影响的句法规则特征融入深度学习模型中,丰富了模型获取到的文本情感特征,提高了模型的准确率。同时,本文还提出了句法规则特征抽取方法,可以在文本输入时将特征提取出来,构建成向量矩阵并作为模型的一个输入。从消融实验的结果可以看出,句法规则特征的引入可以提高深度学习模型的分析准确率,对比实验则说明单个模型的效果要弱于组合模型,且目前组合模型的情感分析方法已成为当前的研究主流。从3个数据集上的实验结果来看,本文构建的模型可以提取到更多的特征,且情感分析的准确率更高,证明了本方法的有效性。

这是第一次将句法规则特征通过特征提取器提取特征的方式引入深度学习模型中,为之后情感分析方法的研究提供了新的研究思路。当然,本文模型也有一些不足之处,如句法规则的引入过于简单、现有的规则特征提取方法可能会造成一些信息的损失,仍然需要进一步研究解决。

参考文献:

- [1] COLLOBERT R, WESTON J, BOTTOU L, et al. Natural language processing (almost) from scratch[J]. *Journal of Machine Learning Research*, 2011, 12: 2493–2537.
- [2] KIM Y. Convolutional neural networks for sentence classification[J]. *arXiv preprint*, 2014, arXiv:1408.5882.
- [3] KALCHBRENNER N, GREFFENSTETTE E, BLUNSOM P. A convolutional neural network for modelling sentences[J]. *arXiv preprint*, 2014, arXiv:1404.2188.
- [4] ZHANG Y, WALLACE B. A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification[J]. *arXiv preprint*, 2015, arXiv:1510.03820.
- [5] GAO J, PANTEL P, GAMON M, et al. Modeling interestingness with deep neural networks[C]// *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. [S.l.:s.n.], 2014: 2–13.
- [6] SHEN Y L, HE X D, GAO J F, et al. A latent semantic model with convolutional-pooling structure for information retrieval[C]// *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*. New York: ACM Press, 2014: 101–110.
- [7] ZHANG R, LEE H, RADEV D R. Dependency sensitive convolutional neural networks for modeling sentences and documents[J]. *arXiv preprint*, 2016, arXiv:1611.02361.
- [8] CHO K, VAN MERRIËNBOER B, GULCEHRE C, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation[J]. *arXiv preprint*, 2014, arXiv:1406.1078.
- [9] SOCHER R, PERELYGIN A, WU J Y, et al. Recursive deep models for semantic compositionality over a sentiment treebank[C]// *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*. [S.l.:s.n.], 2013: 1631–1642.
- [10] TRAN K, BISAZZA A, MONZ C. Recurrent memory network for language modeling[J]. *arXiv preprint*, 2016, arXiv:1601.01272.
- [11] CHEN P, SUN Z Q, BING L D, et al. Recurrent attention network on memory for aspect sentiment analysis[C]// *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg: Association for Computational Linguistics, 2017: 452–461.

- [12] 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.
- [13] WANG Y Q, HUANG M L, ZHU X Y, et al. Attention-based LSTM for aspect-level sentiment classification[C]// Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. Stroudsburg: Association for Computational Linguistics, 2016: 606-615.
- [14] CHENG Y, YAO L B, XIANG G X, et al. Text sentiment orientation analysis based on multi-channel CNN and bidirectional GRU with attention mechanism[J]. IEEE Access, 2020, 8: 134964-134975.
- [15] YANG L, LI Y, WANG J, et al. Sentiment analysis for e-commerce product reviews in Chinese based on sentiment lexicon and deep learning[J]. IEEE Access, 2020, 8: 23522-23530.
- [16] LI W, ZHU L Y, SHI Y, et al. User reviews: sentiment analysis using lexicon integrated two-channel CNN-LSTM family models[J]. Applied Soft Computing, 2020, 94: 106435.
- [17] LI W J, QI F, TANG M, et al. Bidirectional LSTM with self-attention mechanism and multi-channel features for sentiment classification[J]. Neurocomputing, 2020, 387: 63-77.
- [18] USAMA M, AHMAD B, SONG E M, et al. Attention-based sentiment analysis using convolutional and recurrent neural network[J]. Future Generation Computer Systems, 2020, 113: 571-578.
- [19] BASIRI M E, NEMATIS, ABDAR M, et al. ABCDM: an attention-based bidirectional CNN-RNN deep model for sentiment analysis[J]. Future Generation Computer Systems, 2021, 115: 279-294.
- [20] JIN N, WU J X, MA X, et al. Multi-task learning model based on multi-scale CNN and LSTM for sentiment classification[J]. IEEE Access, 2020, 8: 77060-77072.
- [21] BASIRI M E, NEMATIS, ABDAR M, et al. ABCDM: an attention-based bidirectional CNN-RNN deep model for sentiment analysis[J]. Future Generation Computer Systems, 2021, 115: 279-294.
- [22] USAMA M, AHMAD B, SONG E M, et al. Attention-based sentiment analysis using convolutional and recurrent neural network[J]. Future Generation Computer Systems, 2020, 113: 571-578.
- [23] 张宝华, 李奕林, 张华平, 等. 基于层次结构的情感单元表示方法[J]. 计算机工程与科学, 2021: 已录用.
ZHANG B H, LI E H, ZHANG H P, et al. Representation of sentiment unit based on hierarchical structure[J]. Computer Engineering and Science, 2021: accepted.

作者简介



张宝华(1996-),男,北京理工大学计算机学院硕士生,主要研究方向为自然语言处理、情感分析。



张华平 (1978-), 男, 博士, 北京理工大学计算机学院副研究员, 主要研究方向为大数据搜索与挖掘、自然语言处理、社交网络。



厉铁帅 (1975-), 男, 中央军事委员会政法委员会高级工程师, 主要研究方向为计算机应用。



商建云 (1965-), 女, 博士, 北京理工大学计算机学院高级工程师, 主要研究方向为自然语言处理、数据挖掘。

收稿日期: 2021-04-27

通信作者: 张华平, kevinzhang@bit.edu.cn

基金项目: 国家自然科学基金资助项目 (No.61772075); 北京市自然科学基金资助项目 (No.4212026)

Foundation Items: The National Natural Science Foundation of China (No.61772075), Beijing Municipal Natural Science Foundation (No.4212026)