

文章编号: 1003-0077(2019)01-0133-10

融合卷积神经网络与层次化注意力网络的中文 文本情感倾向性分析

程 艳, 叶子铭, 王明文, 张 强, 张光河

(江西师范大学 计算机信息工程学院, 江西 南昌 330022)

摘 要: 文本情感倾向性分析是自然语言处理研究领域的一个基础问题。基于深度学习的模型是处理此问题的常用模型。而当前的多数深度学习模型在中文文本情感倾向性分析方面的应用存在两个问题: 一是未能充分考虑到文本的层次化结构对情感倾向性判定的重要作用, 二是传统的分词技术在处理文本时会产生歧义。该文针对这些问题基于卷积神经网络与层次化注意力网络的优点提出了一种深度学习模型 C-HAN(Convolutional Neural Network-based and Hierarchical Attention Network-based Chinese Sentiment Classification Model), 先用并行化卷积层学习词向量间的联系与组合形式, 再将其结果输入到基本单元为双向循环神经网络的层次化注意力网络中判定情感倾向。实验表明: 模型在中文评论数据集上倾向性分类准确率达到 92.34%, 和现有多个情感分析模型相比有所提升; 此外, 对于中文文本, 选择使用字级别词向量作为原始特征会优于词级别词向量作为原始特征。

关键词: 卷积神经网络; 层次化注意力网络; 情感倾向性分析; 词向量

中图分类号: TP391 **文献标识码:** A

Chinese Text Sentiment Orientation Analysis Based on Convolution Neural Network and Hierarchical Attention Network

CHENG Yan, YE Ziming, WANG Mingwen, ZHANG Qiang, ZHANG Guanghe

(School of Computer Information Engineering, Jiangxi Normal University, Nanchang, Jiangxi 330022, China)

Abstract: Text sentiment orientation analysis is a fundamental problem in natural language processing. To further improve the deep learning based models used in this issue, this paper proposes a new model named C-HAN, i. e. Convolutional Neural Network-based and Hierarchical Attention Network-based Chinese Sentiment Classification Model. It utilizes a convolution layer to extract a sequence of higher-level phrase representations, which are then fed into a hierarchical attention network to obtain the final representations. On the Chinese sentiment analysis corpus, the character level C-HAN achieves a sentiment prediction accuracy of 92.34%, slightly better than the word level C-HAN yielding 91.96% accuracy.

Keywords: convolutional neural network; hierarchical attention network; sentiment orientation analysis; word embedding

0 引言

情感分析, 也称观点挖掘, 是指人们对服务、产品、组织、个人、问题、事件、话题及其属性的情感、观点、评价、态度和情绪^[1]。文本情感倾向性分析是情感分析的一个分支, 其目的在于从原始文本中判断

说话者对事物的情感倾向性。基于知识的方法是最初被广泛应用于此领域的技术, 但其需要编写十分复杂的规则, 才能让计算机较为准确地理解人类语言, 难度较大。故该类方法仅能在小规模的数据上取得一定的成果^[2]。随着文本数据量的增多, 使用基于知识的方法处理文本已是捉襟见肘。自 20 世纪 90 年代以来, 机器学习方法开始在文本情感分析

收稿日期: 2018-08-02 定稿日期: 2018-11-05

基金项目: 国家自然科学基金(61262080, 61562043); 江西省科技重点项目(20151BBE50121, 20161BBE50086); 江西省教育厅科技重点项目(GJJ150299)

领域崭露头角^[3-4]。但这些方法都属于浅层学习范畴,函数模型和计算方法相对简单,导致它们在有限样本和计算单元下无法表达一些复杂的函数,泛化能力较弱,同时也需要人工选择大量数据特征。这些缺陷导致机器学习方法在此任务上遇到了瓶颈。深度学习能够从原始数据中自动学习重要的特征与特征表达方式来处理各种复杂任务,在建模、解释、表达能力以及优化等方面优势明显。卷积神经网络(Convolutional Neural Network, CNN)和循环神经网络(Recurrent Neural Network, RNN)是目前深度学习领域比较热门的两种模型。卷积神经网络能够提取数据中的局部化结构信息,循环神经网络则能够处理序列化结构信息。近年来,也出现了结合两者模型结构的复合模型,在文本情感分析领域取得了优异效果。注意力机制是当前深度学习领域的最新成果,它能够捕捉文本中最具代表性的特征,优化模型结构。使用深度学习模型分析文本情感是当前热门的研究方向^[5]。

本文认为文本的情感倾向是由句子层面和词语或字符层面两个层级共同决定的。首先,文本是由句子组成的,不同的句子对于情感倾向性分析结果而言拥有不同程度的重要性。例如,若文本整体情感倾向为正,其中一些情感色彩较为负面的句子就不是重要的,并不是所有的句子都会影响最终结果的判定。同理,句子又是由词语或字符构成,不同词语或字符对于句子的情感倾向判定又有不同程度的影响。现有模型很少从这个角度出发探索文本的情感,未能很好地体现文本结构的层次化和文本内容的上下文关联对倾向性分析结果的影响。故本文建立了层次化的情感倾向性分析模型,并引入注意力机制,从两个层面筛选出对倾向性分析结果影响最高的文本信息。另一方面,词向量的表示对于文本分类任务非常重要,近年来对词向量粒度的研究越来越细,出现了一些基于字符级别的工作,但这些工作大都基于英文文本数据^[6-7],对于中文数据的研究较少。中文文本与英文文本的差别尤其体现在字符级别上。英文单词由 26 个字母组成,单个字母往往不代表特殊含义。中文则不同,很多单个汉字就能表示明确的含义,组合起来能够表达的语义更是多种多样,故对中文文本进行字符级别的分析是很有意义的。由于中文的特殊性,大部分中文文本分类任务都会使用分词操作。但分词操作执行的同时,固定了汉字间的组合形式,有时易导致歧义,无法切分出正确的汉字组合形式。为了解决此问题,本文

使用卷积神经网络的并行卷积层学习中文文本字级别特征,不依靠解析树等句法分析方法,同时也避免了语言知识层面的分析与复杂的数据预处理过程。本文实验结果表明:对于中文语料,使用训练过后的字级别词向量作为原始特征会好于使用训练过后的词级别词向量作为原始特征。

1 相关工作

情感倾向性分析一直是情感分析领域的研究热点。早前使用的方法主要包括基于情感词典的方法和基于机器学习的方法。基于情感词典的方法,通常是将词典中已经记录了情感倾向性的词条对句子中的词语进行匹配,然后通过对词语的情感倾向进行聚合(如求平均或求和)得到最终的情感倾向性。Kamps 和 Marx 使用 WordNet 判断词语的情感倾向性^[8]。Budnitsky 和 Hirst 通过在 WordNet 中计算词语间的路径距离从而得到情感相似度,以此计算词语的情感倾向性^[9]。规范的中文情感词典相对缺乏,最早也是最普遍传播的是知网(HowNet)提供的情感分析用词语集^[10]。其实,真正的情感判断并不是一些简单规则的堆砌,而是一个复杂、系统的工程,且情感词典中的词语需要人工进行选择。因此该方法的性能很大程度上取决于先验知识与人工设计。基于机器学习的情感倾向性分析问题常被看成一个有监督的学习问题。Pang 等^[11]早在 2004 年便利用朴素贝叶斯、最大熵和支持向量机等机器学习方法来尝试解决情感分析问题。但这些方法需要复杂的特征选择过程,此过程同样依赖于人工设计,导致推广能力差。

深度学习方法能够对特征进行自动选择,逐渐发展成为近年来情感分析领域的主流方法。Collobert 等^[12]于 2011 年首先提出使用 CNN 解决词性标注等 NLP 领域的问题。2014 年, Kim^[13]提出将 CNN 应用于情感分析任务, Kalchbrenner 等^[14]在此基础上提出宽卷积和 K-max pooling 方法。Conneau 等^[15]提出 VDCNN 模型,采用了深度卷积神经网络方法。但 CNN 模型有其缺陷,即只能挖掘文本的局部信息。与 CNN 相比, RNN 更能捕捉到文本间的长距离依赖。Tang 等^[16]为了对句子之间的关系进行建模,提出采用层次化 RNN 模型来对篇章级文本进行建模。Wang 等^[17]提出了 DRNN 模型,固定了信息流动的步长。结合 CNN 与 RNN 各自的优点, Siwei Lai 等^[18]提出了 RCNN 模型,先使

用双向循环神经网络得到上下文表示,再经过卷积、池化操作后输出分类结果。Chunting Zhou 等^[19]提出了 C-LSTM 模型,先利用卷积神经网络提取文本特征,再输入循环神经网络得到分类结果。注意力机制能够捕捉到特征的重要性,在文本情感分析任务中亦有应用,例如, Yang 等^[20]提出层次化注意力模型进行情感分析任务。总之,深度学习方法应用在文本情感倾向性分析中,免去了传统方法繁琐的特征工程步骤,具有一定的优势。

本文模型结构是在 Chunting Zhou 等工作^[19]基础上进行的改进。这种传统的 CNN-RNN 模型架构没有充分考虑到文本不同成分对情感倾向判定的重要程度。本文模型在此架构基础上加入层次化注意力机制,有利于模型学习到对情感倾向结果判定最重要的信息。此外,在卷积部分,本文参考了 Kim 等的工作^[13],不同的是将位置向量引入模型中,赋予每个词以实际位置编码从而构造新型词向量编码,使得模型能够学习到更丰富的词向量编码信息。同时,本文也借鉴了 Yang 等的工作^[20],与之不同的是分别采用词级别向量与字级别向量进行实验。

本文模型的工作流程是:首先利用卷积神经网络将词向量编码到新的向量空间中,学习到词的位

置信息与上下文信息,然后通过层次化注意力学习句子和文本的序列化信息及其对文本倾向判定的重要性。总而言之,该模型综合利用了几种深度学习方法的优点,考虑了文本的局部信息和全局信息,既避免了信息的丢失,又能够筛选出对结果影响最大的信息。

2 模型

2.1 卷积神经网络

CNN 最早应用于计算机视觉中。近年来,其在文本分类中也有着优越表现^[12-13,21]。传统的 CNN 模型在处理文本任务时,常将词语转换为向量形式,将不同数目、不同大小的卷积核与向量进行按元素相乘操作,经过卷积、池化、dropout 正则化等一系列操作后得到最终输出。CNN 能够捕捉到文本任务中字或词之间的局部关系,是文本情感分析任务中较为常用的处理手段。由于本文中采用 CNN 的目的在于抽取单个句子中词语的 N-gram 特征输入到模型下一层结构中,故仅使用了卷积操作。其结构如图 1 所示。

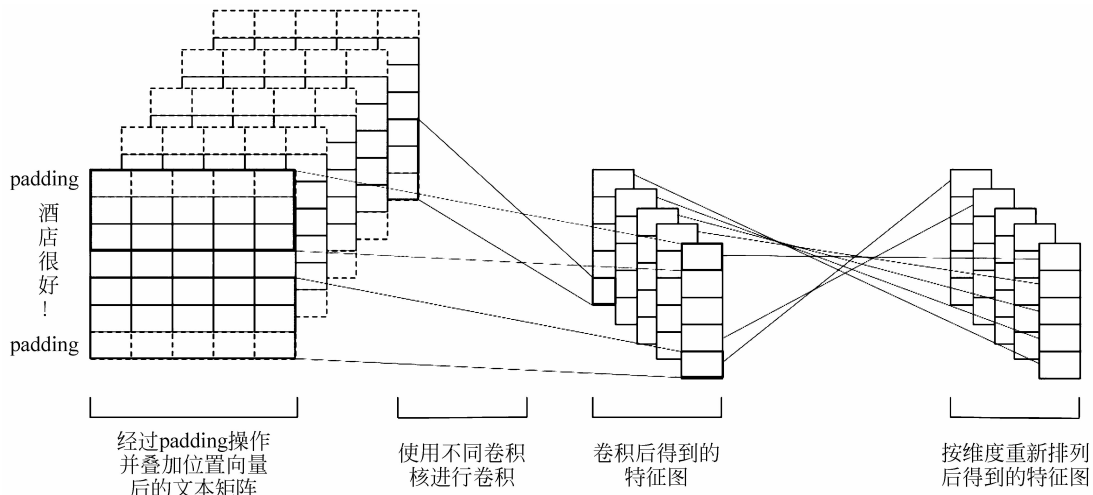


图 1 CNN 卷积层结构图

图 1 中 padding 表示补零操作,目的是为了保证转换后的句子表示矩阵长度和词向量矩阵长度一致。假设当前输入为第 i 个句子中的第 j 个词 $x_{ij} \in R^d$, d 表示词向量维度。在本文中,每一个词均被赋予一个位置编码 $l_{ij} \in R^d$,该编码与词向量语义无关,具体数值通过模型训练学习得到。如此,每一个词便拥有了一个新的编码 a_{ij} ,如式(1)所示。

$$a_{ij} = x_{ij} + l_{ij} \quad (1)$$

设卷积核 $B^n \in R^{h \times d}$ ($n \in \{1, 2, \dots, d\}$) 在新的词向量编码矩阵上进行卷积操作。其中, h 表示卷积核长度,卷积核宽度与词向量维度相同, n 表示卷积核的数量。假设每一个 a_{ij} 经过卷积操作后得到对应的向量表示为 $z_{ij} = [z_{ij}^1, z_{ij}^2, \dots, z_{ij}^d] \in R^d$, 则 z_{ij} 中的每一个元素计算过程如式(2)所示。

$$z_{ij}^n = f\left(\sum_{q=-(h-1)/2}^{(h-1)/2} a_{ij+q} \odot B_{(h-1)/2+q}^n\right) \quad (2)$$

其中, \odot 表示向量按元素点乘后求和, B_m^n 表示卷积核矩阵的第 m 个行向量。 f 表示非线性激活函数 ReLU。最终第 i 个句子的序列表示 z_k , $k \in [1, K]$, 其中 K 表示句子中词的总数目。

2.2 GRU 与 Bi-GRU

GRU(Gated Recurrent Unit)由 Cho 等^[22]于 2014 年提出,其结构如图 2 所示。

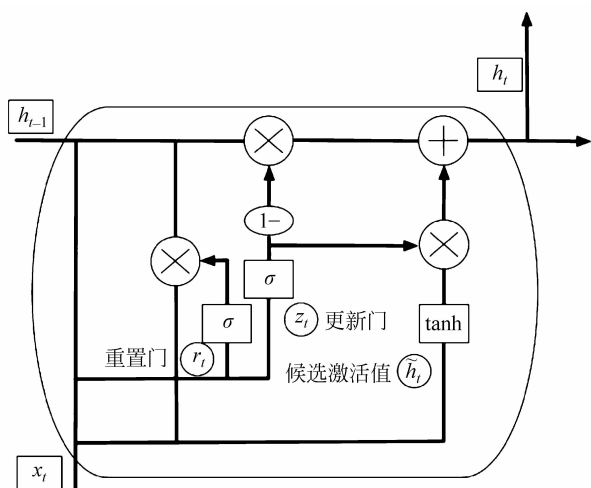


图 2 GRU 结构图

相比 LSTM,GRU 模型更为简单,仅由重置门

r 和更新门 z 组成,用于控制神经元信息的读取、写入。GRU 的计算如式(3)~式(6)所示。

$$h_t^j = (1 - z_t^j)h_{t-1}^j + z_t^j \tilde{h}_t^j \quad (3)$$

$$z_t^j = \sigma(W_z x_t + U_z h_{t-1})^j \quad (4)$$

$$\tilde{h}_t^j = \tanh(W_r x_t + U(r_t \odot h_{t-1}))^j \quad (5)$$

$$r_t^j = \sigma(W_r x_t + U_r h_{t-1})^j \quad (6)$$

其中, h_t^j 表示一个 GRU 单元在 t 时刻的激活值,它同时受到更新门 z , $t-1$ 时刻 GRU 单元激活值 h_{t-1}^j 和候选激活值 \tilde{h}_t^j 的控制。更新门的值 z_t^j 能够决定 GRU 单元激活值的更新程度,由当前状态与上一状态共同决定。候选激活值 \tilde{h}_t^j 的求解过程与标准的 RNN 激活值求解步骤相同。 r_t 为重置门,取值计算过程与更新门类似, \odot 代表按元素相乘操作。例如,若 r_t^j 接近于 0,则 $r_t \odot h_{t-1}$ 也会趋向于 0,表示 \tilde{h}_t^j 只与当前状态有关。总之,GRU 合并了 LSTM 中的遗忘门和输入门,将其统一为更新门,减少了模型的参数与复杂度,是当前较为流行的循环神经网络模型结构之一。

单向 GRU 在使用时是从上文向下文推进的,容易导致后面的词比前面的词更重要。而双向 GRU 是 GRU 的变体,其输出值同时取决于正向计算和后向计算过程,使得输出结果更为精确。其模型结构如图 3 所示。

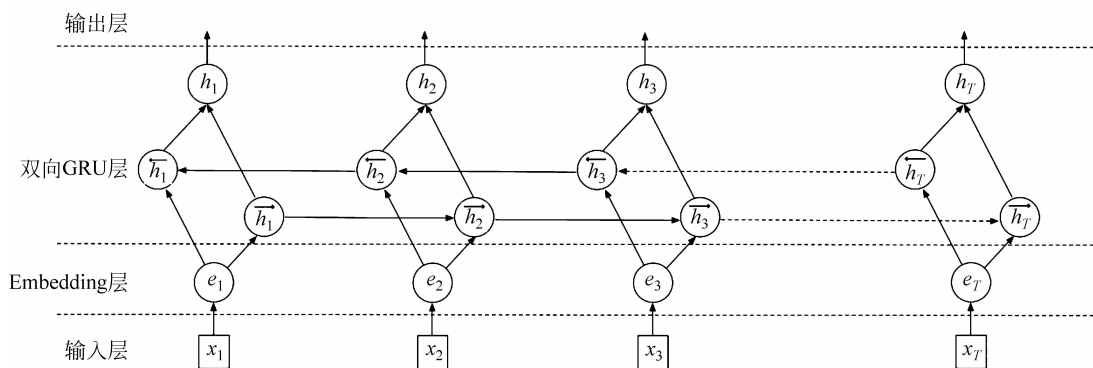


图 3 双向 GRU 结构图

2.3 注意力机制

注意力机制早在 2014 年便在机器翻译任务中得到应用^[23],经过一段时间的发展,亦产生了许多不同形式的变体^[24-25]。

注意力模型可以抽象为由 Module1 和 Module2 两个模块组成。Module1 一般为编码器,对输入数据做一定的变换;Module2 为解码器,同样经过

一定的变换后输出数据。每个输出值 m_i 计算过程如式(7)所示。

$$m_i = F(C_i, m_1, m_2, \dots, m_{i-1}) \quad (7)$$

其中, C_i 为每一个输出数据相对应的语义编码,该编码由输入数据的分布生成,如式(8)所示。

$$C_i = \sum_{j=1}^T a_{ij} S(n_j) \quad (8)$$

其中, $S(n_j)$ 表示经过 Module1 处理后得到的

输入数据的隐层状态, T 表示输入数据的个数。 a_{ij} 表示输入 j 对输出 m_i 的注意力分配概率, a_{ij} 计算过程如式(9)、式(10)所示。

$$a_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^T \exp(e_{ik})} \quad (9)$$

$$e_{ij} = V \tanh(W h_j + U s_{i-1} + b) \quad (10)$$

其中, e_{ij} 指第 j 个输入对第 i 个输出的影响力评价分数, h_j 为 module1 中第 j 个输入的隐层状态, s_{i-1} 为上一步过程中 module2 的输出, W 、 U 和 V 为权重矩阵, b 为偏置值, 均由训练过程中学习得

到。Attention 语义编码会作为 module2 的输入, 生成最终的深层特征, 获取最关键的语义信息。

2.4 具体模型结构

本文在上述基础上提出一个融合卷积神经网络与层次化注意力网络的文本情感分析模型。该模型由经卷积神经网络操作的字/词级别初始化向量模块、双向循环神经网络及字/词级别注意力模块、双向循环神经网络及句子级别注意力模块组成。模型结构如图 4 所示。

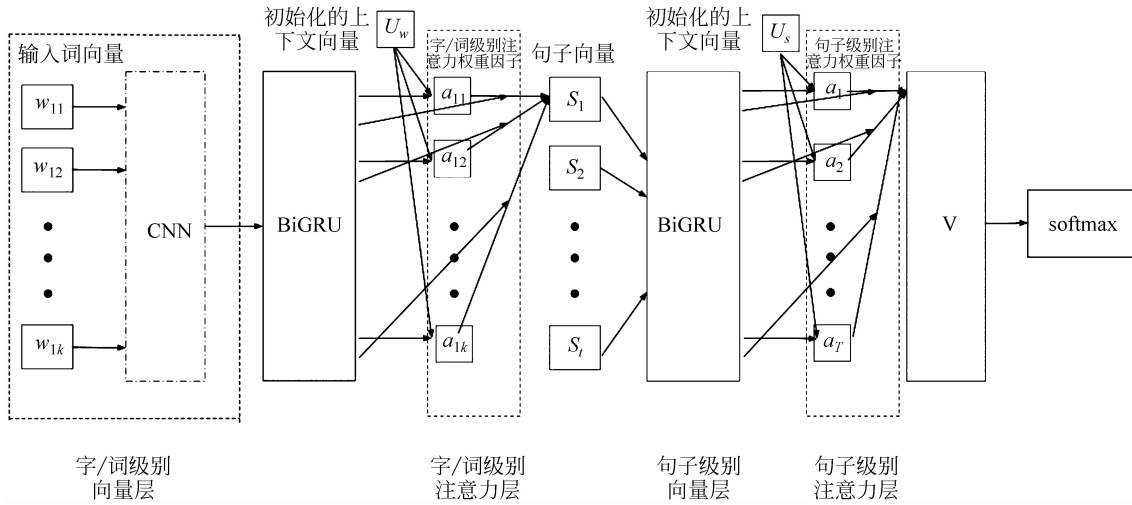


图4 本文模型结构图

模型首先将中文字符或词语通过 CNN 层的操作转化为相应的向量表达形式。假设一段文本有 L 个句子, 表示为 $S_i, i \in [1, L]$ 。句子中又包含 K 个字符或单词。由式(1)已经得到了第 i 个句子经过卷积操作后的向量表示 $z_{ik}, k \in [1, K]$, 故在完成第一步操作后, 将经卷积神经网络操作后输出的结果通过 Bi-GRU 网络将其上下文相关信息结合起来可以获得隐藏层的输出, 具体计算过程如式(11)~式(13)所示。

$$\vec{g}_{ik} = \overrightarrow{\text{GRU}}(z_{ik}), k \in [1, K] \quad (11)$$

$$\overleftarrow{g}_{ik} = \overleftarrow{\text{GRU}}(z_{ik}), k \in [K, 1] \quad (12)$$

$$g_{ik} = [\vec{g}_{ik}, \overleftarrow{g}_{ik}] \quad (13)$$

其中, g_{ik} 即为经过双向 GRU 后得到的向量化表示形式。

这一步后加入 Attention 机制的目的是要把一个句子中对句子含义贡献最大的字或词语找出来。首先, 将 g_{ik} 输入到一个单层的感知机中得到的结果 u_{ik} 作为 g_{ik} 的隐含表示。单词的重要性采用 u_{ik} 和一个随机初始化的上下文向量 U_w 的相似度来决定。

然后, 经过 softmax 操作获得了一个归一化的 Attention 权重矩阵, 代表句子 i 中第 k 个字或词的权重。最后在得到 Attention 权重矩阵后, 将句子向量看作组成这些字或词向量的加权求和。计算过程如式(14)~式(16)所示。

$$u_{ik} = \tanh(W_w g_{ik} + b_w) \quad (14)$$

$$a_{ik} = \frac{\exp(u_{ik}^T U_w)}{\sum_k \exp(u_{ik}^T U_w)} \quad (15)$$

$$S_i = \sum_k a_{ik} g_{ik} \quad (16)$$

其中, W_w 与 b_w 分别为权重矩阵和偏置矩阵。 a_{ik} 为衡量句子 i 中第 k 个字或词重要性的注意力权重因子。

在求得了 S_i 的表示后, 我们用相似的方法可以对句子进行处理, 获得经过双向 GRU 后得到对应的隐层句子向量 G_i , 如式(17)~式(19)所示。

$$\vec{G}_i = \overrightarrow{\text{GRU}}(S_i), i \in [1, L] \quad (17)$$

$$\overleftarrow{G}_i = \overleftarrow{\text{GRU}}(S_i), i \in [L, 1] \quad (18)$$

$$G_i = [\vec{G}_i, \overleftarrow{G}_i] \quad (19)$$

随后通过引入一个句子级别的上下文向量 U_s 用以衡量句子在整个文本中的重要性程度,得到文本总向量 V ,最后可通过 softmax 层进行情感分析操作。计算过程如式(20)~式(23)所示。

$$u_i = \tanh(W_s G_i + b_s) \quad (20)$$

$$a_i = \frac{\exp(u_i^T U_s)}{\sum_i \exp(u_i^T U_s)} \quad (21)$$

$$V = \sum_i a_i G_i \quad (22)$$

$$p = \text{soft max}(W_2 V + b_2) \quad (23)$$

同上, W_s 、 W_2 与 b_s 、 b_2 分别为权重矩阵和偏置矩阵。 a_i 为衡量句子 i 重要性的注意力权重因子。

除此之外,本文训练的最终目标为最小化损失函数(负对数似然函数),如式(24)所示。

$$L = - \sum_d \log p_{d_j} \quad (24)$$

其中 j 为文本 d 相对应的情感类别标签。

3 实验分析

3.1 情感倾向性分析数据集

本文数据集采用国内学者谭松波整理的酒店评论数据集,部分评论数据格式如表 1 所示,对原始数据集进行整理、欠采样、合并等操作后,得到正面类别情感评论数据与负面类别情感评论数据各 3 000 条共 6 000 条评论数据,在此基础上进行十折交叉验证。

表 1 酒店评论数据集示例

正面	负面
自助早餐非常好,服务很周到,下次还会去住。	非常糟糕的一个酒店,所有的东西都言过其实。我是看了网上的评论才会考虑入住这个酒店的。
相当于五星的服务,好于四星的早餐,差不多三星的房价,就是要借东西。服务员中有个印度人,很有派头,服务热情。	这个酒店以前讲来还是经济实惠,我三月份来这住,服务员服务还不错。今天就不好,来前台登记说酒店系统升级没有之前的资料了,要重新登记,空调嗡嗡响,修了一次又一次,我去找几次,前台小姐说:不是修好了吗?态度不好。我也算老顾客了,太差了。房间装修的白灰袋还在房角堆放,明天就换,下次我不会再住了。

3.2 数据预处理与模型参数

词向量在深度学习模型中具有十分重要的作用。词向量的预训练有助于提高模型准确率^[26]。在词向量的训练过程中,一些句法与语义方面的信息也能够被学习到,这在情感分析的过程中十分重要。本文运用 word2vec 工具^[27],计算词语的向量形式表示,从而进行基于无监督方法的词向量学习。为了预先训练好中文词向量,本文使用大规模中文维基百科数据训练 skip-gram 模型。中文词向量的维度设置为 300 维。以中文字符作为初始化词向量进行训练时,针对句子中的每一个字,为其训练一个词向量放入字典中。以单个字符作为句子层面的基本单位。以词语作为初始化词向量训练时,操作过程同上,但需要使用 Jieba 分词工具对文本进行分词,以分词之后的词语做为句子层面的基本单位。在 word2vec 模型的训练过程中,指定训练字符的最小出现次数为 5,将出现次数超过 5 次的字加入字典中,对于没有在字典中出现的字符,随机初始化其向量形式表示。在句子层面,本文选取逗号、句号、感叹号和问号作为句子间的分隔符进行句子切割。设置最大句子长度为 50,小于该值时,对句子进行补零操作。大于该值时,进行截断操作。设置文本中最大句子数目为 20,预处理时同样进行补零和截断。

本文实验基于 Keras 深度学习框架^[28]。从整体模型架构看,我们使用了一个卷积层。字/词级别注意力层使用一个双向 GRU 层,在句子级别注意力层同样使用一个双向 GRU 层。对于卷积层,本文尝试分别使用卷积窗口大小为 2、3、4 及其组合的卷积方式。根据模型表现最终选取单一窗口大小为 3 的卷积核,卷积核单元数量设置为 300,采用“same”卷积模式;对于双向 GRU 层,将其维度设置为 300,对上下文向量进行随机初始化。

在模型训练方面,对乱序的微批次样本采用随机梯度下降,批量大小设置为 32。训练过程中采用 Adam^[29]更新规则,初始学习率为 0.001,防止过拟合的 dropout 参数设置为 0.2,采用准确率指标对模型表现进行评估。训练过程中对词向量进行微调。

3.3 对比实验与结果分析

本文运用多个模型在此数据集上进行了对比实验并分析实验结果。因超参数选取和具体任务密切相关,故本文参考原论文设置对比试验的参数,以使

模型准确率达到最高。所有对比试验的深度学习模型中词向量均进行微调。实验结果如表 2 所示。对每个实验的具体说明如下。

表 2 模型准确率对比

模型	Accuracy/%
Fasttext	83.67
SVM-word	80.68
SVM-character	81.36
CNN-word	89.14
CNN-character	90.98
RCNN-word	89.67
RCNN-character	90.88
HAN-word	91.32
HAN-character	91.93
C-HAN-word(本文模型)	91.96
C-HAN-character(本文模型)	92.34

Fasttext^[30]: Fasttext 是 Facebook 开源的文本分类工具。本实验中,将模型学习率设为 0.1,词向量维度选为 300。

SVM-word: 抽取 word2vec 训练出的词级别词向量作为输入,使用 SVM 模型进行词级别情感分类,采用线性核。在数据集上进行十折交叉验证。

CNN-word^[13]: 词级别单层卷积神经网络模型。使用 word2vec 训练出的词级别词向量进行试验。卷积核相关超参数设置与本文相同。

RCNN-word^[18]: 结合双向 LSTM 与 CNN 的模型。使用 word2vec 训练出的词级别词向量进行试验。超参数设置与原论文相同。

HAN-word^[20]: 层次化注意力机制模型。使用 word2vec 训练出的词级别词向量进行试验。超参数设置与原论文相同。

C-HAN-word: 本文模型。使用 word2vec 训练出的词级别词向量进行试验。

SVM-character: 抽取 word2vec 训练出的字级别词向量作为输入,用 SVM 模型分类进行字级别情感分类,使用线性核。在数据集上进行十折交叉验证。

CNN-character^[13]: 字级别单层卷积神经网络模型。利用 word2vec 中文字级别词向量进行实验。卷积核相关超参数设置与本文相同。

RCNN-character^[18]: 结合双向 LSTM 与 CNN 的模型。利用 word2vec 中文字级别词向量进行实验。超参数设置与原论文相同。

HAN-character^[20]: 层次化注意力机制模型。

利用 word2vec 中文字级别词向量进行实验。超参数设置与原论文相同。

C-HAN-character: 本文模型。利用 word2vec 中文字级别词向量进行实验。

3.3.1 模型准确率分析

由表 2 可知, Fasttext 模型的分类准确率(83.67%)高于 SVM 模型(80.68%)。在中文情感分析任务中取得了较好的分类效果,证明了 Fasttext 的优良性能。因其模型简单,拥有极快的训练测试速度,在数据量庞大的情况下,可作为基线模型使用。此外,对比传统的机器学习方法 SVM(80.68%、81.36%),基于神经网络的深度学习方法(CNN 等)对分类结果准确率的提升效果显著。准确率可以达到 90%左右,近乎提升了 10%,证明深度学习方法在中文文本情感分析任务中是更有效的。

对几种深度学习方法进行比较分析,在词级别层面, CNN(89.14%)、RCNN(89.67%)、HAN(91.32%)的准确率依次上升;在字级别层面, RCNN(90.88%)的准确率相较 CNN(90.98%)略有下降但相差无几。而 HAN(91.93%)模型准确率仍能达到将近 92%,这表明注意力机制选择性关注特定目标的优势,在情感倾向性判定中能够得到充分体现。同时,以上结果也表明对于深度学习模型,多种不同类型模型的融合能够带来情感分类准确率的提高。观察准确率,本文模型(C-HAN)在词级别(91.96%)与字级别(92.34%)的准确率相较于之前几种方法均有不同程度的提高,且模型综合表现优于 HAN 模型。原因在于 C-HAN 充分考量到了卷积神经网络、循环神经网络和注意力机制各自的优点。相较于缺乏注意力机制的模型(如 RCNN), C-HAN 引入了层次化注意力机制,使得模型学习过程中只关注有效的信息,降低了噪音影响;而针对缺乏卷积层的模型(如 HAN), C-HAN 引入了卷积层并加入位置向量编码,使得模型学习词向量时能够学习到更加准确且丰富的词向量信息,从而有效提升中文文本情感分类的准确率。

3.3.2 字级别与词级别比较分析

本文在 SVM、CNN、RCNN 与 HAN 模型上分别进行了字级别与词级别的对比实验。从实验结果可看出:使用字级别词向量后,模型表现会优于词级别词向量。原因在于字级别向量特征的粒度更小,在模型训练时,学习到的文本特征更为具体。这

一点可以通过对比分析 word2vec 训练过后的词向量与字向量得证。表 3 给出了词语间的词级别词向量余弦相似度和字级别词向量相加得到的词级别词向量的余弦相似度。

表 3 词向量余弦相似度

词 1	词 2	相似度
酒店	饭店	0.776
酒+店	饭+店	0.803
汉庭	宾馆	0.648
汉+庭	宾+馆	0.724
房间	单间	0.677
房+间	单+间	0.739

由以上三个例子可以看出,经 word2vec 训练过后的字级别词向量相加得到的词级别词向量余弦相似度要高于经 word2vec 训练过后的词级别词向量余弦相似度。当然,在实验过程中我们发现也有少数词语组合不符合上述情况。例如,“君悦”是一个酒店品牌,其与“酒店”间的相似度为 0.647,然而分割开来的相似度为 0.523。这种情况出现的原因是因为训练语料中关于“君悦”的内容少,所以字级别词向量没有较好地学习到相关信息。

3.3.3 卷积层分析

值得注意的是,与 HAN 相比,本文中卷积层的加入使得层次化注意力机制模型的分类效果提升。原因在于:传统分词技术并不是绝对有效的,有时会出现带有歧义的切分,导致切分无法体现句子的正确语义。而在将字级别向量作为输入,通过卷积层操作后,可以学习到相当于 N-gram 的信息与字

符间的抽象联系。

举例来说,“果然是一家高大上酒店”。在这句话中,用传统的中文分词技术会将其切分为“高大/上/酒店”或是“高大/上酒/店”。这些切分都无法让模型学习到句子的正确语义信息,后一种甚至将“酒店”这一关键词错误切分。而以单个汉字字符为单位输入时,通过卷积层设置,譬如设置卷积窗口大小为 3,就可以学习到“高大上”这样的正确语义。此外,语料中出现的一些地名、酒店名的出现亦会导致分词出错,通过以字符为单位增加并行化卷积层操作,可以学习到正确的局部语义信息。

3.3.4 注意力可视化

为了更加直观地展示模型效果,本文在实验时运用 matplotlib 库分别从句子级别和字符级别对注意力权重分配进行可视化展示。分别选取一段短文本与一段长文本进行试验。

文本一: 非常糟糕的一个酒店,所有的东西都言过其实。我是看了网上的评论才会考虑入住这个酒店的。

文本二: 这个酒店以前讲来还是经济实惠,我三月份来这住,服务员服务还不错。今天就不好,来前台登记说酒店系统升级没有之前的资料了,要重新登记,空调嗡嗡响,修了一次又一次,我去找几次,前台小姐说:不是修好了吗?态度不好。我也算老顾客了,太差了。房间装修的白灰袋还在房角堆放,明天就换,下次我不会再住了。

对于文本一,将其以逗号和句号划分,可以分为 3 个子句。绘制句子级别注意力权重热力图如图 5(a)所示。图中灰度值越大表示注意力分配权重越高。

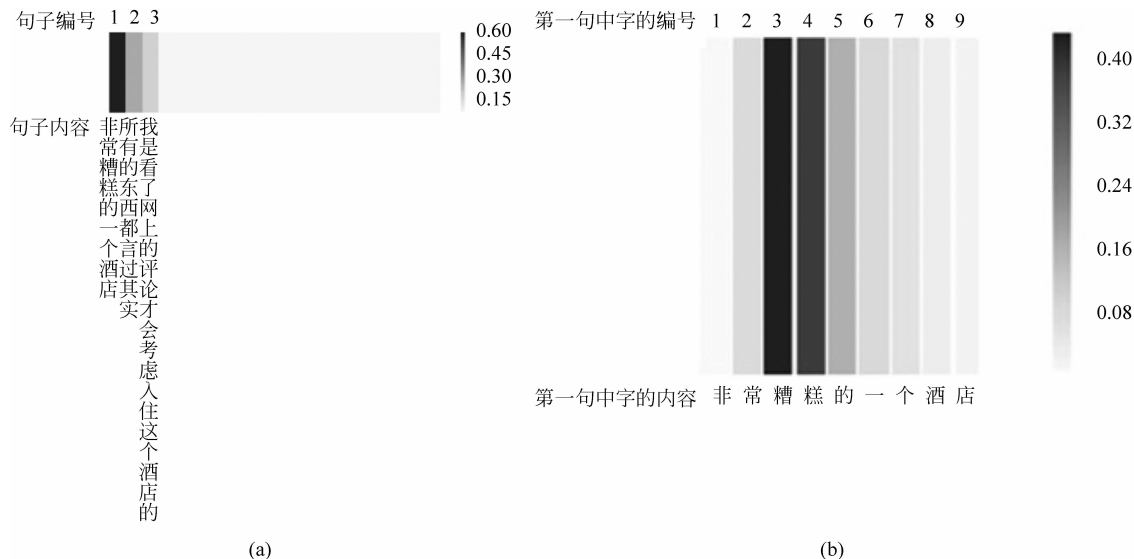


图 5 文本一注意力权重图

不难看出,第一个句子对评论文本的情感倾向影响较大,第二句次之,第三句最小。进一步地,对影响力最大的句子即第一句绘制字符级别注意力权重热力图如图 5(b)所示。

第一个子句中含有 9 个中文汉字,从图中可以看出模型对第三、四、五个字符分配了较高的权重,对应文本中的“糟”“糕”“的”三个汉字。

对于文本二,将其以逗号、句号和问号划分可以分为 16 个子句。绘制句子级别注意力权重热力图如图 6 所示。

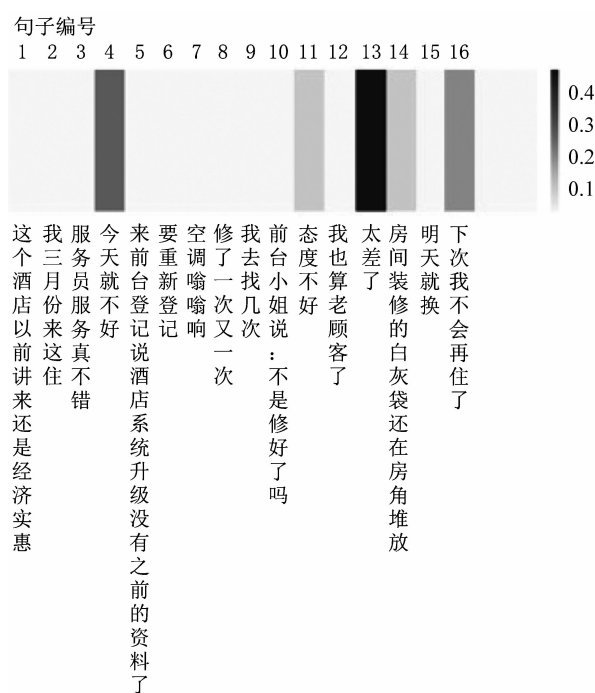


图 6 文本二句子级别注意力权重图

由图 6 可知,本文模型寻找到的对文本情感倾向性影响力较大的句子分别为第四句、第十一句、第十三句和最后一句,对它们分别绘制字符级别注意力权重热力图,如图 7 所示。

由图 7 可知,“就”、“不”、“好”三个汉字在第四句中所占权重较高。同理,对其余三句话分析亦能得出相应分配权重高的汉字。

在两段文本上进行的可视化实验结果表明:本文模型能够找出对情感倾向分析最大的句子,同时亦能在句子中找出对结果影响较大的汉字。特别是在处理文本二这种带有转折的文本时,模型亦有良好表现。具体体现在:文本二中前 3 句话情感倾向是偏积极的,而文本的总情感倾向却是消极的。本文模型能够做到忽略类似于前 3 句话这样偏离文本总情感倾向的句子,找到对文本总体情感倾向结果

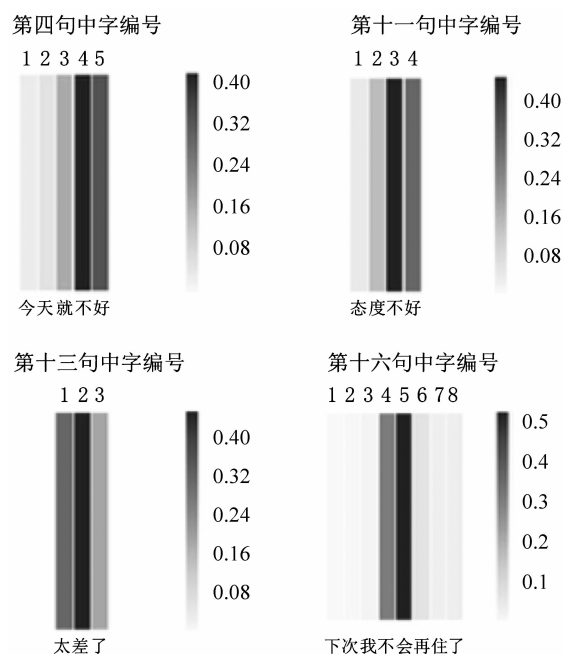


图 7 文本二子句字符级别注意力权重图

影响大的句子。

4 总结与展望

本文结合卷积神经网络、循环神经网络与注意力机制,构建了一种新的层次化中文文本情感倾向性分析模型。其中,卷积神经网络能够抓取到字符或词语间的抽象关系,循环神经网络能够找到语句上下文间的关系,而注意力机制能够有效识别对判定情感倾向有用的隐含信息。实验证明,模型准确率达到 92.34%,优于 SVM 和其他深度学习模型,层次化注意力机制的引入是有效的。此外,本文通过模型在中文文本情感倾向性分析数据集上的表现,证明了对于中文语料,使用字级别词向量作为原始特征会优于使用词级别的词向量作为原始特征。今后,在本文基础上探寻基于树结构等不同形式卷积神经网络对模型的影响,以及尝试更多注意力模型结构的引入,优化文本情感分析模型,会成为进一步研究的方向。

参考文献

- [1] Liu B. Sentiment analysis: Mining opinions, sentiments, and emotions[J]. Computational Linguistics, 2015,42(3): 1-4.
- [2] Ortony A, Clore G L, Collins A. The cognitive structure of emotions[M]. Cambridge University Press,

- 1990.
- [3] Pang B, Lillian L, Vaithyanathan S. Thumbsup?: sentiment classification using machine learning techniques [C]//Proceedings of ACL, 2002: 79-86.
- [4] Hu M, Liu B. Mining and summarizing customer reviews[C]//Proceedings of SIGKDD, 2004: 168-177.
- [5] 梁军, 等. 基于深度学习的微博情感分析[J]. 中文信息学报, 2014, 28(5): 155-161.
- [6] Zhang X, Zhao J, Lecun Y. Character-level convolutional networks for text classification[C]//Proceedings of NIPS, 2015: 645-657.
- [7] Kim Y, et al. Character-aware neural language models [C]//Proceedings of AAAI, 2016: 2741-2749.
- [8] Kamps J, Marx M. Words with attitude[C]//Proceedings of International Conference on Global WordNet, 2002: 332-341.
- [9] Esuli A, Sebastiani F. Pageranking wordnet synsets: An application to opinion mining[C]//Proceedings of Annual Conference of the Association for Computational Linguistics, 2007: 442-431.
- [10] Zhen D D, Qiang D. HowNet and the computation of meaning (With Cd-rom)[M]. World Scientific, 2006.
- [11] Pang B, Lee L. A Sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts[C]//Proceedings of ACL, 2004: 271-278.
- [12] Collobert R, et al. Natural language processing (almost) from scratch[J]. Journal of Machine Learning Research, 2011, 12(8): 2493-2537.
- [13] Kim Y. Convolutional neural networks for sentence classification[C]//Proceedings of EMNLP, 2014.
- [14] Kalchbrenner N, Grefenstette E, Blunsom P. A convolutional neural network for modelling sentences[J]. arXiv preprint arXiv: 1404.2188, 2014.
- [15] Conneau, et al. Very deep convolutional networks for text classification[C]//Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, 2017.
- [16] Tang D, Qin B, Liu T. Document modeling with gated recurrent neural network for sentiment classification [C]//Proceedings of EMNLP, 2015: 1422-1432.
- [17] Wang B. Disconnected recurrent neural networks for text categorization[C]//Proceedings of ACL, 2018: 2311-2320.
- [18] Lai S, et al. Recurrent convolutional neural networks for text classification [C]//Proceedings of AAAI, 2015: 2267-2273.
- [19] Zhou C T, et al. A C-LSTM neural network for text classification [J]. Computer Science, 2015, 1(4): 39-44.
- [20] Yang Z, et al. Hierarchical attention networks for document classification[C]//Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 2016: 1480-1489.
- [21] 刘龙飞, 等. 基于卷积神经网络的微博情感倾向性分析[J]. 中文信息学报, 2015, 29(6): 159-165.
- [22] Cho K, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation[J]. arXiv preprint arXiv: 1406.1078, 2014.
- [23] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate[J]. arXiv preprint arXiv: 1409.0473, 2014.
- [24] Yin W, et al. Abcnn: Attention-based convolutional neural network for modeling sentence pairs[J]. arXiv preprint arXiv: 1512.05193, 2015.
- [25] 栾克鑫, 等. 基于注意力机制的句子排序方法[J]. 中文信息学报, 2018, 32(1): 123-130.
- [26] Zhang Y, Wallace B. A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification[J]. arXiv preprint arXiv: 1510.03820, 2015.
- [27] Mikolov T, et al. Distributed representations of words and phrases and their compositionality[C]//Proceedings of Advances in Neural Information Processing Systems, 2013: 3111-3119.
- [28] Chollet F, Keras[CP/OL]. [2015]. <https://github.com/fchollet/keras>.
- [29] Kingma D, Jimmy B. Adam: A method for stochastic optimization[J]. arXiv preprint arXiv: 1412.6980, 2014.
- [30] Bojanowski P, et al. Enriching word vectors with subword information[J]. arXiv preprint arXiv: 1607.04606, 2016.



程艳(1976—),通信作者,博士后,教授,主要研究领域为智能信息处理、机器学习、教育大数据。
E-mail: chyan88888@jxnu.edu.cn



王明文(1964—),博士,教授,主要研究领域为智能信息处理、信息检索。
E-mail: mwwang@jxnu.edu.cn



叶子铭(1994—),硕士,主要研究领域为自然语言处理、机器学习。
E-mail: 2685487478@qq.com