

基于卷积神经网络和注意力模型的文本情感分析*

冯兴杰, 张志伟[†], 史金钊

(中国民航大学 计算机科学与技术学院, 天津 300300)

摘要: 针对社交网络数据的文本情感分析, 目前常用的研究方法主要是基于传统机器学习算法, 根据手工标注好的情感词典, 对文本信息使用朴素贝叶斯、支持向量机、最大熵方法等机器学习算法进行情感分析。为了避免对手工方式建立的情感词典的依赖, 减少机器学习过程中的人工干预, 提出基于卷积神经网络和注意力模型相结合的方法进行文本情感分析。实验表明, 根据准确率、召回率和 F_1 测度等衡量指标, 提出的方法较传统的机器学习方法和单纯的卷积神经网络方法有明显的提高。

关键词: 社交网络; 文本情感分析; 卷积神经网络; 注意力模型

中图分类号: TP393.04

文献标志码: A

文章编号: 1001-3695(2018)05-1434-03

doi:10.3969/j.issn.1001-3695.2018.05.033

Text sentiment analysis based on convolutional neural networks and attention model

Feng Xingjie, Zhang Zhiwei[†], Shi Jinchuan

(College of Computer Science & Technology, Civil Aviation University of China, Tianjin 300300, China)

Abstract: Most of the existing research methods about the text sentiment analysis of social network are based on the hand-marked emotional dictionary to analyze the text information's emotion, such as naive Bayes, support vector machines, maximum entropy method and other machine learning algorithms. However, by manually establishing the emotional dictionary is a type of time-consuming and laborious work, in order to avoid dependence on sentiment dictionary. This paper proposed a method of combining the convolutional neural network and attention model for sentiment analysis. The experimental results show that the proposed method is more effective than the existing machine learning method and the simplex convolution neural network method in terms of accuracy, recall rate and F_1 measure.

Key words: social network; text sentiment analysis; convolutional neural network; attention model

随着微博、Facebook、Twitter等社交网络的兴起,网络不仅成为了人们获取信息的重要来源,同时也成为人们表达自己观点的平台。通过在博客、微博、Twitter等网络社区来评论热点事件、抒写影评观点、描述产品体验等,产生了大量带有情感倾向的文本信息,而通过对这些文本信息进行情感分析,可以更好地理解用户行为,发现用户对产品的倾向性、对热点事件的关注程度等^[1]。随着信息规模的急剧增大,仅仅依靠人工进行处理已经无法完成这一任务,这就促进了自然语言处理领域的一个研究热点,即文本情感分析技术^[2]的发展。目前,文本情感分析的主要研究方法还是基于传统机器学习的算法,通过人工设计特征构造出结构化的文本信息特征,然后用机器学习的方法来进行分析。常用的文本情感分析方法有朴素贝叶斯、支持向量机、最大熵方法等,这些方法都可以被划分为浅层学习方法^[3]。浅层学习方法计算量小且实现容易,但是由于对复杂函数的表达能力的限制,使得对复杂分类问题的泛化能力受到制约。为弥补这一缺陷,人工构造特征被引入这一模型当中,如使用人工标注的情感词典、句法与语法分析等,虽然这些方法可以有效地提高文本情感分析的准确率,但由于需要过多的人工标注数据,费时费力,并且需要一定的先验知识,所以随着互联网规模的不断发展,文本数据规模的不断扩大,从而限制了这些方法的发展。本文采用基于卷积神经网络和注意力模型的结构,避免了依赖人工构造特征的方法,采用相关数据

集对网络模型进行训练后,再进行文本情感分析工作。

1 相关工作

文本情感分析主要通过分析文本内容来判断文本所表达的情感倾向,发现用户对某一事件的关注程度,自2002年由Pang等人^[4]提出关于情感分析的工作后,经过许多学者的研究,获得了很大的发展。情感分析技术大致可以分为基于规则的方法和基于统计的方法,其中基于情感词典的机器学习方法是目前的主要方法。Pang等人根据传统自然语言处理中的文本分类技术,使用了朴素贝叶斯、支持向量机和最大熵等模型,在电影评论上取得了不错的效果;Turney^[5]提出了采用依赖种子情感词集合,使用互信息的方法来判别某个短语是否是评价词语;Ding等人^[6]提出针对特定领域情感词配对的方法来判断情感极性;罗毅等人^[7]通过构建二级情感词典,使用N-gram模型获取文本特征进行情感分析;任远等人^[8]采用支持向量机和TF-IDF计算特征项权值来进行情感分析。自2006年由Hinton等人提出深度学习之后,随着深度学习技术在计算机视觉和语音识别领域的成功应用,越来越多的深度学习技术也被应用于自然语言处理方向。Bengio等人^[9]提出了利用神经网络构建语言模型的方法,把词向量映射到低维空间,通过距离来度量词与词之间的相似性;Mnih等人^[10]提出层析Log-Bilin

收稿日期: 2016-12-23; 修回日期: 2017-02-21 基金项目: 国家自然科学基金资助项目(U1633110, U1233113)

作者简介: 冯兴杰(1969-),男,河北邢台人,博士,主要研究方向为智能算法、智能信息处理理论与技术;张志伟(1992-),男(通信作者),硕士研究生,主要研究方向为人工智能、自然语言处理(zhangzwt@163.com);史金钊(1992-),女,硕士研究生,主要研究方向为人工智能、自然语言处理。

ear模型来训练语言模型;Mikolov等人^[11]借鉴Log-Bilinear模型的思想,提出word2vec模型,实现了CBOW和Skip-gram两种框架,并随着Google开源了其代码,词嵌入被应用到自然语言处理多个领域;随着Kim^[12]利用卷积神经网络来进行句子分类后;梁军等人^[13]利用深度学习来探讨了中文微博情感分析工作的可行性;Bahdanau等人^[14]提出使用注意力模型在机器翻译中取得了不错的效果,随后被应用于Google神经网络翻译系统中。

目前,由于深度学习方法在许多领域都已经取得了不错的效果,本文提出基于卷积神经网络和注意力模型的结构来进行文本情感分析的工作并进行实验,经过实验验证,本文提出的方法是有效的。

2 基于卷积神经网络和注意力模型的结构

为了实现文本情感分析的任务,本文提出了基于卷积神经网络和注意力模型的网络结构,模型结构如图1所示。该模型主要由两部分组成,左边部分为典型的卷积神经网络结构,右边部分为注意力模型的结构,本模型的整体流程为:首先对输入的文本句子利用词向量模型来进行编码,转换为词向量表示后,经过卷积神经网络以后得到该句子的相关特征,然后结合由注意力机制得到的特征进行拼接以后,通过全连接后利用分类器来完成文本情感分析的工作。

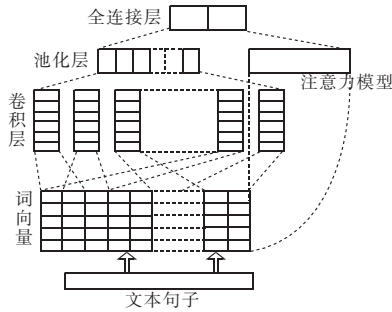


图1 基于卷积神经网络和注意力模型的结构

2.1 模型表示

给定文本句子数据集 D ,其中包含有文本 $X\{x_1, x_2, \dots, x_m\}$ 及每个句子所对应的情感标签 $Y\{y_1, y_2, \dots, y_m\}$,其中每个文本句子 x_i 有 n 个词组成,表示为 $\{x_{i1}, x_{i2}, \dots, x_{in}\}$,将最终的目标函数表示为

$$p(Y|X) = \arg \max_{\theta} f(Y|X; \theta) \quad (1)$$

其中: θ 表示该模型中涉及到所有的参数; $f(\cdot)$ 表示该模型的形式化表达。

2.2 卷积神经网络

卷积神经网络是一种前馈神经网络,其网络结构主要由输入层、卷积层、池化层(下采样层)、全连接层和输出层组成。其中卷积层为特征提取层,通过滤波器来提取句子的特征;池化层为特征映射层,对经过卷积层后得到的特征进行采样,得到局部最优值。在本模型中将文本句子表示为输入层,其中对每个句子 X 表示为 $n \times k$ 的矩阵,其中 n 表示构成文本句子的词的长度, k 表示词向量 x_i 的维度,文本句子中的词向量采用word2vec模型来进行训练得到。卷积层主要是为了来学习文本句子的局部特征,本层主要对输入层的词向量矩阵进行卷积操作,对每个大小为 k 的连续窗口进行操作,结果表示为

$$c_i = f(w \cdot x_{i:i+h-1} + b) \quad (2)$$

其中: c_i 表示经过卷积操作后对应的第 i 个特征值; $f(\cdot)$ 表示本层卷积核函数的选择; w 表示滤波器里的权重矩阵,其中 $w \in \mathbb{R}^{h \times k}$, $h \times k$ 表示选取的滤波器的大小; b 表示偏置项; $x_{i:i+h-1}$ 表示由文本句子中的第 i 个词到 $i+h-1$ 个词的长度,本文采用多个滤波器来进行学习。经过卷积层后,得到特征矩阵 C 表示为

$$C = [c_1, c_2, \dots, c_{n-h+1}]^T \quad (3)$$

其中: $C \in \mathbb{R}^{n-h+1}$ 。

池化层表示对本文句子经过卷积层后得到的特征矩阵 C 进行下采样,选出其中局部最优特征,本文采用最大池化方式来进行采样,经过池化层以后得到的特征表示为

$$\tilde{C} = \max(c_1, c_2, \dots, c_{n-h+1}) \quad (4)$$

在卷积层本文选用多通道的方式,即选择多个滤波器来进行特征的提取,经过以上操作以后即可得到对原始文本句子的特征。

2.3 注意力模型

注意力模型是用来表征文本句子中的词与输出结果之间的相关性,表示句子 x_i 中的每个词与其相对应标签 y_i 之间的重要程度。在此将采用注意力模型生成的注意力文本用 a_i 表示为

$$s_i = f_{\text{att}}(x_{ij}, y_i) \quad 1 \leq i \leq m, 1 \leq j \leq n \quad (5)$$

$$p_i = \frac{\exp(s_i)}{\sum_{i=1}^m \exp(s_i)} \quad 1 \leq i \leq m \quad (6)$$

$$a_i = \sum_{i=1}^m p_i \times x_i \quad (7)$$

其中: x_i 表示一个文本句子; y_i 表示此句子所对应的标签; f_{att} 表示含一个隐层的前向网络; p_i 与 s_i 表示文本中每个词的重要度信息。经过卷积神经网络与注意力模型得到特征后,将池化层学习到的特征 \tilde{C} 与注意力文本 a_i 连接,作为全连接层的输入,经过全连接层后输出结果,表示为

$$s(x) = f'(w' \cdot (\tilde{C} \otimes a_i) + b') \quad (8)$$

其中: $s(x)$ 表示经过模型后得到的输出值; \otimes 表示向量拼接操作; w' 表示全连接层的权重矩阵; b' 表示偏置项; $f'(\cdot)$ 表示分类器的选择。

2.4 模型训练

文本情感分析本质上是一个分类问题,本文对此作二分类进行处理,将其分为正面和负面两类情感类别。

本文模型通过最小化负对数似然函数来进行训练。对于给定的一个句子 x_i ,通过本文模型经过训练以后,得到给定句子的情感标签 $\tau \in T$ 得分 $s_{\theta}(x)$,其中, T 代表所分的类别,通过选择分类器softmax转换为条件概率:

$$p(\tau|X; \theta) = \frac{\exp(s_{\theta}(x))_{\tau}}{\sum_{\forall i \in T} \exp(s_{\theta}(x))_i} \quad (9)$$

对式(9)取对数得

$$\log p(\tau|X; \theta) = s_{\theta}(x)_{\tau} - \log(\sum_{\forall i \in T} \exp(s_{\theta}(x))_i) \quad (10)$$

本文采用随机梯度下降算法来最小化负对数似然函数,得到

$$\theta := \sum_{(x_i, y_i) \in D} -\log p(y_i|x_i; \theta) \quad (11)$$

其中: x_i, y_i 表示训练语料的一条句子及其对应的情感标签; D 表示语料库。

3 实验与分析

3.1 实验数据准备

为了验证模型的有效性,本文选用了有关中文情感挖掘的

酒店评论语料(ChnSentiCorp)。ChnSentiCorp 是中科院谭松波博士收集整理的一个酒店评论的语料,其公布的语料规模为 10 000 篇,被分为四个子集。本文选用 ChnSentiCorp-Htl-ba-6000 数据来进行实验,其为平衡语料,正负类各 3 000 篇。表 1 为数据集样例。

表 1 ChnSentiCorp 数据集样例

积极	消极
商务大床房,房间很大,床有两米宽,整体感觉经济实惠不错!	标准间太差,房间还不如三星的,而且设施非常陈旧。建议酒店把老的标准间重新改善。
非常好的酒店,四星的标准完全超值的享受,服务非常好	肯定我不会再住这里了,太陈旧了,霉味太重,感觉不好!
环境很好,地点很方便,服务也很好,下回还会住	环境一般,住了之后让人感觉价格和服务不成比例!

3.2 实验评价指标

准确率(precision)和召回率(recall)是两个用来评价分类结果好坏的指标,定义如下:

$$\text{precision} = \frac{\text{system. correct}}{\text{system. output}}, \text{recall} = \frac{\text{system. correct}}{\text{human. labeled}}$$

其中:system. correct 表示系统返回的总记录数;system. correct 表示系统为该类返回的正确结果数;human. labeled 表示测试集中该类的总数目。其中,准确率用于衡量分类器准确性,召回率用于衡量分类器是否能找全该类样本,这两个指标应该该兼顾。使用 F_1 测度来均衡这两方面,定义如下:

$$F_1 = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

3.3 词向量训练

本文选用维基百科的中文语料作为训练的语料库,利用 Google 开源的 word2vec tool 来进行词向量的训练,然后对语料库进行预处理,并以此作为本模型的输入层。本文采用的 word2vec tool 的 skip-gram 模型,上下文窗口大小设置为 5,词向量维度大小设为 50,采样值大小设为 $1e-3$,如果有词语不在预训练好的词向量中,则采用随机初始化方式来进行表示。

3.4 超参数设置与训练

对于本文的网络结构,模型中的激活函数选为 relu 函数,采用多组卷积核来进行训练,其中滤波器窗口大小为 2、3、4,每个滤波器的卷积单元个数为 100,隐藏层单元数量为 300,输出层选用 softmax 进行分类,且训练过程中采用 dropout 以防止过拟合,最后采用随机梯度下降算法来进行权重的更新迭代。

3.5 实验结果与分析

本文与传统的基于词袋模型的 SVM 方法、基于 word2vec 的 SVM 方法以及基于 word2vec 的卷积神经网络的方法来进行对比,利用 ChnSentiCorp-Htl-ba-6000 数据集进行 10 倍交叉验证来进行实验分析,实验结果如表 2 所示。

表 2 四种模型的分类结果性能对比(ChnSentiCorp-Htl-ba-6000 数据集)

模型	准确率	召回率	F1 测度
CBOW-SVM	0.775 4	0.772 3	0.773 8
W2V-SVM	0.812 5	0.811 9	0.812 1
W2V-CNN	0.851 6	0.849 6	0.850 1
W2V-Att-CNN	0.872 7	0.871 3	0.871 9

CBOW-SVM:向量特征为 uni-gram,利用 SVM 进行训练分类。

W2V-SVM:将词利用 word2vec 训练转换为词向量以后,利用 SVM 进行训练分类。

W2V-CNN:将词利用 word2vec 训练转换为词向量以后,利用卷积神经网络来进行训练。

W2V-Att-CNN:将词利用 word2vec 训练转换为词向量以

后,利用基于卷积神经网络和注意力模型的结构进行训练。

通过对以上对比实验结果进行分析,发现利用 W2V-SVM 训练的结果优于 CBOW-SVM 的训练结果,原因在于利用 word2vec 训练出的词向量相比于 CBOW 模型而言考虑了上下文的语义信息,因此进行情感分析可以得到更好的结果;W2V-CNN 训练的结果优于 W2V-SVM 的训练结果,原因在于利用卷积神经网络包含了多个并行卷积层,从而考虑了不同 N-gram 的信息,得到了更有意义的文本句子中的相关特征,因此得到了更好的结果;W2V-Att-CNN 训练的结果优于 W2V-CNN 的结果,原因在于加上注意力机制后,整个网络模型不仅考虑了不同 N-gram 的信息,而且考虑了文本句子与结果的相关性,将通过卷积神经网络与注意力模型得到的特征相结合以后,从而得到了更好的结果。

4 结束语

本文提出了一种基于卷积神经网络与注意力模型的文本情感分析方法,将文本信息编码为词向量以后,通过卷积神经网络得到文本的相关特征,再结合注意力模型。实验结果表明了本文方法的可行性和有效性,可以更好地发现文本信息的情感倾向性。目前微博等社交网络的信息中图文相关的信息比较多,下一步工作将进一步探讨结合文字与图片信息的情感分析任务,寻找更适合情感分析的深度学习算法。

参考文献:

- [1] Ding Xiao, Liu Ting, Duan Junwen, et al. Mining user consumption intention from social media using domain adaptive convolutional neural network[C]//Proc of the 29th AAAI Conference on Artificial Intelligence. Palo Alto, CA: AAAI Press, 2015: 2389-2395.
- [2] 杨立公,朱俭,汤世平,等. 文本情感分析综述[J]. 计算机应用, 2013, 33(6): 1574-1578.
- [3] 于凯,贾磊,陈雨强,等. 深度学习的昨天、今天和明天[J]. 计算机研究与发展, 2013, 50(9): 1799-1804.
- [4] Pang Bo, Lee L. Opinion mining and sentiment analysis[J]. Journal Foundations and Trends in Information Retrieval, 2008, 2(2): 1-135.
- [5] Turney P D. Thumbs up or thumbs down: semantic orientation applied to unsupervised classification of reviews[C]//Proc of the 40th Annual Meeting on Association for Computational Linguistics. Stroudsburg, PA: Association for Computational Linguistics, 2002: 417-424.
- [6] Ding Xiaowen, Liu Bing, Yu P S. A holistic lexicon-based approach to opinion mining[C]//Proc of Conference on Web Search and Web Data Mining. New York: Association for Computing Machinery, 2008: 231-240.
- [7] 罗毅,李利,谭松波,等. 基于中文微博语料的情感倾向性分析[J]. 山东大学学报:理学版, 2014, 49(11): 1-7.
- [8] 任远,巢文涵,周庆,等. 基于话题自适应的中文微博情感分析[J]. 计算机科学, 2013, 40(11): 231-235.
- [9] Bengio Y, Ducharme R, Vincent P, et al. A neural probabilistic language model[J]. Journal of Machine Learning Research, 2003, 3: 1137-1155.
- [10] Mnih A, Hinton G E. A scalable hierarchical distributed language model[C]//Proc of International Conference on Neural Information Processing Systems. USA: Curran Associates Inc., 2009: 1081-1088.
- [11] Mikolov T, Sutskever I, Chen Kai, et al. Distributed representations of words and phrases and their compositionality [C]//Proc of International Conference on Neural Information Processing Systems. USA: Curran Associates Inc., 2013: 3111-3119.
- [12] Kim Y. Convolutional neural networks for sentence classification[EB/OL]. (2014-08-25). <https://arxiv.org/abs/1408.5882>.
- [13] 梁军,柴玉梅,原慧斌,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.
- [14] Bahdanau D, Cho K, Bengio Y. Neural machine translation by jointly learning to align and translate[C]//Proc of International Conference on Learning Representations. 2015.