

● 李悦晨, 钱玲飞, 马 静 (南京航空航天大学经济与管理学院, 江苏 南京 211106)

基于 BERT-RCNN 模型的微博谣言早期检测研究^{*}

摘 要: [目的/意义] 为了解决传统谣言检测算法在实际应用中存在滞后性的问题, 尝试在不使用评论和转发数据的基础上实现微博谣言早期检测。[方法/过程] 针对传统谣言检测模型需要大量特征以及难以实现及时检测的问题, 使用 BERT 模型对微博原文进行向量表示, 然后将获取的语义特征输入到 RCNN 模型中进行谣言检测。[结果/结论] 基于微博谣言数据集进行对比实验, 实验结果显示基于 BERT-RCNN 的微博谣言早期检测模型的准确率为 95.16%, $F1$ 指标为 95.14%; 该模型与其他主流的谣言检测模型相比, 能在较短时间内完成对谣言的检测, 证明了方法的有效性。[局限] 文章提出的检测模型针对微博文本进行检测, BERT 模型需要依赖谷歌发布的预训练模型。

关键词: 微博; 谣言检测; BERT 模型; 预训练; 深度学习

DOI: 10.16353/j.cnki.1000-7490.2021.07.025

引用格式: 李悦晨, 钱玲飞, 马静. 基于 BERT-RCNN 模型的微博谣言早期检测研究 [J]. 情报理论与实践, 2021, 44 (7): 173-177, 151.

Early Detection of Micro Blog Rumors Based on BERT-RCNN Model

Abstract [Purpose/significance] In order to solve the problem of lag in the actual application of traditional rumor detection algorithms, an attempt was made to achieve early detection of microblog rumors without using comments and forwarding data. [Method/process] In view of the problems that traditional rumor detection models require a large number of features and are difficult to achieve timely detection, the BERT model is used to vectorize the original text of Weibo, and then the acquired semantic features are input into the RCNN model for rumor detection. [Result/conclusion] A comparative experiment was carried out on the Weibo rumors dataset. The experimental results proved that the accuracy of the early detection model of Weibo rumors based on BERT-RCNN was 95.16%, and the $F1$ index was 95.14%. [Limitations] The detection model proposed in this article is for microblog text detection, and the BERT model needs to rely on the pre-training model released by Google.

Keywords: micro blog; rumor detection; BERT model; pre-training; deep learning

随着互联网的不断发展, 用户在中文社交媒体上所产生的网络数据呈指数型增长, 同时促进网络谣言的生成与传播。社会学家将谣言定义为“广泛传播但未经官方认证的与当前事实相关的信息, 并且试图让人们相信是真实的。”^[1] 微博平台上用户可以直接发布任何信息, 而且用户并没有足够的能力来鉴别这些信息的真实性。微博是极具代表性的中文社交平台, 以微博数据为实验对象进行谣言检测研究具有广泛借鉴意义。目前, 微博平台主要依靠人工审核和辟谣平台等方式来检测谣言, 不仅效率低下, 同时也需要消耗大量的人力成本和时间, 并不能满足当下谣言实时检测的需要^[2]。因此, 如何实现微博谣言自动检

测是未来研究的重点。

传统谣言检测方法通常会结合大量特征, 需要使用谣言传播产生的大量评论信息和转发信息进行计算, 因此传统模型检测效果是相对滞后的, 往往将谣言检测出来时已经产生了极大的不良影响。

谣言早期检测是在谣言产生之初就进行检测, 更加贴合当下网络环境监管的需要, 能够有效避免谣言传播所带来的社会危害。

本文针对传统谣言检测模型存在一定滞后性的问题, 提出一种基于 BERT-RCNN 的微博谣言早期检测模型, 尝试在不使用评论特征、传播特征等数据的情况下, 仅通过分析微博文本的内容特征来实现谣言早期检测。通过在微博谣言数据集上进行实验, 证明本文所提出的模型可以在微博谣言产生早期就实现对其的检测, 相较于传统模型更早取得检测结果, 从而大大避免谣言广泛传播带来的危害。

^{*} 本文为国家自然科学基金项目“大数据环境下学术成果真实价值与影响的实时预测及长期评价研究”(项目编号: 19BTQ062)和研究生开放基金项目“基于 BERT 模型的短文本分类方法与应用研究”(项目编号: kfj20200907)的成果。

全面的词向量表示。

Transformer 编码器^[12]由编码器和解码器两部分组成，编码器将可变长度的输入序列转化为固定长度的向量，解码器则将得到的固定长度的向量解码为可变长度的输出序列。

编码器部分是 Transformer 的主要结构，如图 2 所示。

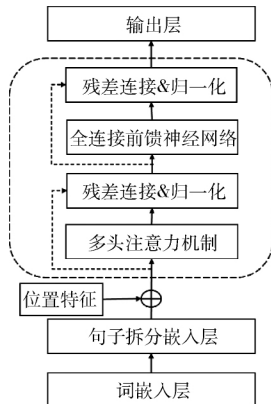


图2 Transformer 编码器结构图

编码器是由 6 个完全相同的层堆叠而成的，每一层有两个子层。第一个子层是多头自注意力机制 (Multi-Head Attention) 层；第二个子层是一个简单的全连接前向神经网络，在两个子层之间通过残差网络结构进行连接，后接一个正则化层。

编码器的主要模块是自注意力层^[13]，其思想是计算一句话中每一词与所有词的相互关系，并利用相互关系来调整每个词的权重来获得新的表达，表示在词本身语义的基础上还包含与其他词的关系，可以实现一词多义的区别。

多头自注意力层的输入是由词向量构成的查询矩阵 Q 、键矩阵 K 、值矩阵 V ，并进行如下计算：

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (1)$$

公式 (1) 根据 Q 矩阵来计算 K 中每个键的权重， K 中每个键的向量维度表示为 d_k ，防止向量维度过高时得到的计算结果过大，对 d_k 进行开方运算对权重进行缩放。得到键的权重矩阵后，与值矩阵相乘得到每个键的最终计算结果。

对矩阵按行进行归一化计算，对行向量元素进行等比例压缩处理，压缩后的向量元素和为 1。

$$\text{softmax}(z_1, z_2, \dots, z_N) = \frac{1}{\sum_i e^{z_i}} (e^{z_1}, e^{z_2}, \dots, e^{z_N}) \quad (2)$$

经过上述一系列计算后得到原始输入语句中每一个字的 Attention 向量，这里的向量融合了其他位置的词信息，将其按行进行排列得到一个矩阵便是最后输出的

Attention 值。

为了更全面地计算注意力，Transformer 引入多头注意力机制。对输入先进行多次不同的线性映射，再计算映射结果的缩放点积注意力，将每次计算结果称为一个头 (head)，再把多次运算得到的 Attention 矩阵横向拼接起来，接着乘以一个权重矩阵压缩成一个矩阵，具体计算公式^[14]如下：

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \text{head}_2, \dots, \text{head}_n) W^O \quad (3)$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

W_i^Q, W_i^K, W_i^V 表示第 i 个头对应的 3 个权重矩阵；Concat 函数将多个头的计算结果进行拼接； W^O 是拼接时使用的权重矩阵。

在 Add & Norm 层对注意力层计算得到的结果进行层次归一化，并采用残差连接的思想来避免出现网络层次过深而导致的退化问题。计算公式如下：

$$\text{LN}(x_i) = \alpha \times \frac{x_i - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}} + \beta \quad (4)$$

Feed Forward 层是一个包括两层全连接计算和一个 RELU 激活函数的全连接网络，计算公式如下：

$$\text{FNN}(x) = \max(0, W_1 \cdot x + b_1) \cdot W_2 + b_2 \quad (5)$$

对全连接网络的输出再进行一次归一化和残差连接处理计算后得到整个模型的输出。BERT 模型的输出有字符级别和句子级别两种向量，本文使用的是句子级别向量加上权重作为语义特征。与传统文本表示方法相比，优势在于不用再去做特征提取和特征向量的拼接。

得到微博文本对应的语义特征后，本文引入 RCNN 模型进行谣言分类检测。RCNN^[15]模型将传统卷积神经网络的卷积层替换为内部具有递归结构的循环卷积层，并按前馈连接建立深度网络结构，循环卷积层结构如图 3 所示。

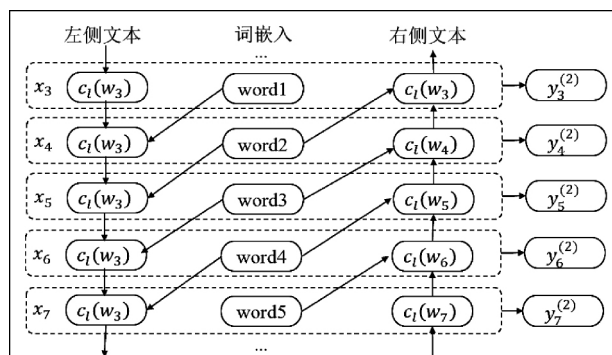


图3 循环卷积层结构图

RCNN 的优势在于可以较为均匀地利用单词的上下文信息，既能够解决在循环神经网络^[16] (Recurrent Neural Network, RNN) 中后面的单词比前面的单词影响力更大的缺点，也不需要像卷积神经网络^[17] (Convolutional Neu-

ral Networks, CNN) 一样需要通过窗口大小来设定对上下文的依赖长度。

将提取到的特征向量放入单层神经网络中,得到潜在语义向量,然后将结果输入池化层,通过 max-pooling 来取得文本中最具有代表性的关键特征,再用全连接层输出,最后通过 softmax 函数来获得分类的结果。

3 实验与分析

3.1 实验数据

为了增加数据量,本实验使用两个公开的中文微博数据集,第一个来源于2016年 Ma^[7]文献中公开的微博数据集,第二个来源于2018年 Song^[18]文献中公开的数据集(CED_Dataset),共包含谣言3889条和非谣言4162条。

原始的微博数据中包含数据较多,为了证明模型可以有效地进行微博谣言早期检测,本文只选取其中的微博原始文本,不使用评论文本。随机选取数据集90%作为训练集,10%作为测试集,在训练集上进行模型训练并调参,最后使用测试集进行测试。

3.2 实验过程

由于BERT模型参数较多,需要使用GPU设备进行计算,本文使用谷歌提供的云平台Google Colab进行模型训练,并使用Pytorch深度学习框架。使用BERT-BASE作为基础模型,模型大小110MB。使用谷歌官方提供的中文预训练模型chinese_L-12_H-768_A-12.zip。

BERT模型其他具体参数如表1所示。

表1 BERT模型训练参数

BERT模型参数	参数取值
学习率	$5e-5$
每批训练数据大小	32
隐层神经元数	768
Transformer 编码器层数	12
自注意力头数	12

RCNN模型训练参数如表2所示。

表2 RCNN模型训练参数

RCNN模型参数	参数取值
词嵌入维度	768
隐藏层节点数	256
学习率	0.001
每批训练数据大小	128
随机失活率	0.1

3.3 对比实验

为了验证本文提出的BERT-RCNN模型在谣言检测任务上的效果,我们将谣言早期检测效果与常用谣言检测模型进行对比实验。实验不仅需要验证基于BERT进行词向量表示的有效性,还需要验证RCNN谣言检测模型的有效

性。对比模型如下所示。

1) baseline: 使用BERT模型对数据集进行预训练并提取深层特征后,连接全连接层将生成的对应向量输入一个Softmax函数进行计算得到分类结果。

2) Word2Vec-CNN: 对微博原文使用Word2Vec来训练词向量表示,并将结果输入到CNN模型中进行分类得到谣言检测结果。

3) Word2Vec-RNN: 对微博原文使用Word2Vec来训练词向量表示,并将结果输入到RNN模型中进行分类得到谣言检测结果。

4) Word2Vec-RCNN: 对微博原文使用Word2Vec来训练词向量表示,并将结果输入到RCNN模型中进行分类得到谣言检测结果。

5) BERT-RNN: 使用BERT模型进行词向量化后输入到RNN模型中完成谣言检测。

6) BERT-CNN: 使用BERT进行词向量化后输入到CNN模型中完成谣言检测。

7) BERT-RCNN: 使用BERT进行词向量化后输入到RCNN模型中完成谣言检测。

具体对比实验结果如表3所示。

表3 对比实验结果

模型	Accuracy	Precision	Recall	F1-score
BERT	0.9268	0.9258	0.9274	0.9266
Word2vec-CNN	0.8139	0.8128	0.8136	0.8132
Word2vec-RNN	0.8313	0.8293	0.8322	0.8308
Word2vec-RCNN	0.8387	0.8377	0.8386	0.8381
BERT-RNN	0.9367	0.9363	0.9369	0.9366
BERT-CNN	0.9442	0.9445	0.9435	0.9439
BERT-RCNN	0.9516	0.9527	0.9506	0.9514

与模型2)到模型4)进行对比实验是为了通过使用文本分类领域常用的词向量文本表示模型Word2Vec的实验结果与使用BERT进行文本表示的模型进行对比,证明BERT提取文本信息的有效性。使用Word2Vec进行文本表示的模型指标表现均不理想,主要原因是Word2Vec仅能表现文本中词语之间的语义关系,不能对不同语境下多义词进行区分,而且忽视了长距离语义关联信息。

与模型5)到模型7)进行对比实验是为了证明使用RCNN的谣言检测模型的有效性,这几组模型均采用BERT模型进行文本特征表示。使用BERT的谣言检测模型与前面几个模型相比,能够取得更好的实验结果。在进行实验的3种使用BERT的模型中,RNN和CNN是文本分类中常用的神经网络模型,RCNN模型各项指标与前两者比较均有不同程度的提升,说明RCNN模型在提取文本的语义特征方面具有更大的优势。

最终实验结果表明,BERT-RCNN模型在原始数据较

少的情况下,与传统模型相比具有更强的语义特征提取能力,在 Accuracy, Precision, Recall 以及 F1-score 几个指标均远远好于其他模型。通过对比实验结果可以看出本文提出的谣言检测模型可以有效提高谣言检测效果,较好地实现微博谣言早期检测,证明了本文提出的谣言检测模型的有效性。

3.4 早期检测效果

在谣言早期检测及时性方面,使用全部评论数据进行训练的谣言检测模型在评论数在 500 条以上才能取得较好的检测效果^[19],而 Ma 等^[7]通过划分时间序列实现早期检测则需要每条微博下有 300 条以上评论时模型才能稳定,从微博发布到产生 300 条评论平均用时 23 小时,也就是说传统研究检测时效在 23 小时以上。而本文提出的 BERT-RCNN 模型不需要结合传播和评论信息,在微博刚发布时就可以对其进行检测,更能满足微博谣言早期检测的实际需要,可以避免谣言大范围传播所带来的负面影响。

使用 4GB 内存 GPU 训练模型,经过 40 轮训练耗时 1 小时后模型已经取得了一个较高的准确率。随着数据量增大,模型训练时间也会相应增加,但是训练好的模型对单条微博的检测时间是不会增长的,在极短时间内便可以得到检测结果,可以满足实际检测中对实时性的需求。

通过以上实验结果可以证明本文提出的微博谣言早期检测模型能够有效地实现谣言早期检测,也能获得比其他检测模型更好的检测效果。

4 结束语

本文针对谣言检测存在的所需特征较多和检测及时性不强问题,尝试在微博谣言传播早期,不依赖微博的评论转发等其他信息的情况下,有效地实现微博谣言早期检测。在微博数据集上展开实验研究,并通过对比实验进一步证明本文提出的 BERT-RCNN 模型的检测效果。结果表明本文方法的检测效果均优于常用谣言检测模型,且通过多个指标证实了本文模型在谣言早期检测任务上能够取得较为理想的效果。

本文创新点主要在于引入 BERT 和 RCNN 模型,针对微博谣言特征提取及早期检测提出了新的解决方案。BERT 模型的提出使得预训练方法在自然语言处理中得到更多的重视,很多具体下游任务都可以尝试在使用 BERT 模型进行预训练后的基础上再进行微调,为网络舆情管控系统提供了有效方法。取得高质量语义特征是开展多项文本处理任务的基础,该模型也为文献发展脉络分析、技术监测与预警、情报感知与刻画研究等一系列图书情报研究工作提供了参考价值。

同时,本文提出的检测模型也存在一定的局限和不足,一方面是 BERT 模型计算复杂度高,对实验设备要求高,需要使用大量计算资源;另一方面,网络谣言形式多样,大多带有图片、视频等信息,如何将这些信息结合到谣言检测模型中是未来谣言检测研究的重点。在后续的工作中也将针对这些问题展开进一步的研究。□

参考文献

- [1] GIST. Rumor and public opinion [J]. American Journal of Sociology, 1951, 57 (2): 159-167.
- [2] 刘雅辉,靳小龙,沈华伟,等. 社交媒体中的谣言识别研究综述 [J]. 计算机学报, 2018, 41 (7): 108-130.
- [3] CASTILLO C, MENDOZA M, POBLETE B. Information credibility on twitter [C] //Proc of 20th International Conference on World Wide Web Hyderabad, ACM, 2011: 675-684.
- [4] YANG F, LIU Y, YU X, YANG M. Automatic detection of rumor on Sina Weibo [C] //Proceedings of the ACM SIGKDD Workshop on Mining Data Semantics, ACM, 2012.
- [5] KWON S, CHA M, JUNG K, CHEN W, WANG Y. Prominent features of rumor propagation in online social media [C]. 2013 IEEE 13th International Conference on Data Mining. IEEE, 2013.
- [6] 刘政,卫志华,张韧弦. 基于卷积神经网络的谣言检测 [J]. 计算机应用, 2017, 37 (11): 3053-3056.
- [7] MA Jing, GAO Wei, MITRA P, et al. Detecting rumors from microblogs with recurrent neural networks [C] //Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence. IJCAI, 2016: 3818-3824.
- [8] CHEN T, LI X, YIN H, et al. Call attention to rumors: deep attention based recurrent neural networks for early rumor detection [C] //Proceedings of the 2018 Pacific-Asia Conference on Knowledge Discovery and Data Mining, 2018: 40-52.
- [9] 牛海波,赵丹群,郭倩影. 基于 BERT 和引文上下文的文献表征与检索方法研究 [J]. 情报理论与实践, 2020, 43 (9): 125-131.
- [10] 谌志群,鞠婷. 基于 BERT 和双向 LSTM 的微博评论倾向性分析研究 [J]. 情报理论与实践, 2020, 43 (8): 173-177.
- [11] DEVLIN J, CHANG M W, LEE K, et al. Bert: pre-training of deep bidirectional transformers for language understanding [EB/OL]. <http://arXiv.org/pdf/1810.04805.pdf>.
- [12] VASWANI A, SHAZEER N, PARMAR N, et al. Attention is all you need [J]. arXiv preprint arXiv: 1706.03762v5, 2017.
- [13] ADAMS WEI Yu, et al. Qanet: combining local convolution with global self-attention for reading comprehension [J]. arXiv: 1804.09541, 2018.

(下转第 151 页)

- ternational Conference on Knowledge Discovery And Data Mining, 2012: 1149-1157.
- [25] 商显震, 韩萌, 王少峰, 等. 融合迁移学习和神经网络的皮肤病诊断方法 [J]. 智能系统学报, 2020, 15 (3): 452-459.
- [26] 李渊彤, 罗裕升, 朱珍民. 基于深度学习的舌象特征分析 [J]. 计算机科学, 2020, 47 (11): 148-158.
- [27] WAN S, LIANG Y, ZHANG Y. Deep convolutional neural networks for diabetic retinopathy detection by image classification [J]. Computers & Electrical Engineering, 2018, 72: 274-282.
- [28] SHIE C, CHUANG C, CHOU C, et al. Transfer representation learning for medical image analysis [C]. 2015 37th Annual International Conference of the IEEE Engineering in Medicine and Biology Society (EMBC). IEEE, 2015: 711-714.
- [29] 肖希明. 信息资源建设: 概念、内容与体系 [J]. 中国图书馆学报, 2006 (5): 5-8.
- [30] 洪秋兰, 潘荣. 全媒体环境下高校图书馆信息资源建设研究 [J]. 图书情报研究, 2020, 13 (3): 96-100.
- [31] 朱学芳, 王贵海, 祁彬斌. 5G时代数字信息资源智能服务研究内容及进展 [J]. 情报理论与实践, 2020, 43 (11): 16-21.
- [32] 张莹, 高慧颖, 巴志超. 基于主题关联挖掘的跨类型数字资源分类方法 [J]. 情报理论与实践, 2015, 38 (11): 108-114.
- [33] 李广丽, 朱涛, 刘斌, 等. 面向大数据的数字图书馆多媒体信息检索系统优化研究 [J]. 情报科学, 2019, 37 (2): 115-119.
- [34] 李宇航, 夏绍模, 程华亮. 基于跨域协同的移动图书馆个性化推荐模型研究 [J]. 情报科学, 2017, 35 (3): 82-86.
- [35] ACKHOFF R. From data to wisdom [J]. Journal of Applied Systems Analysis, 1989, 16: 3-9.
- [36] GAO Y, MOSALAM K. Deep transfer learning for image-based structural damage recognition [J]. Computer-Aided Civil and Infrastructure Engineering, 2018, 33 (9): 748-768.
- [37] YANG Q, LIU Y, CHEN T, et al. Federated machine learning: concept and applications [J]. ACM Transactions On Intelligent Systems, 2019, 10 (2): 1-2, 19.
- [38] LIANG H, FU W, YI F. A survey of recent advances in transfer learning [C]. 2019 IEEE 19th International Conference on Communication Technology (ICCT). IEEE, 2019: 1516-1523.
- 作者简介: 王金婉, 女, 1991年生, 博士生。研究方向: 智能信息处理, 机器学习。朱学芳 (通信作者), 男, 1962年生, 教授, 博士生导师。研究方向: 数字信息资源管理及服务, 数字人文, 模式识别, 人工智能。
- 作者贡献声明: 王金婉, 提出研究思路, 设计论文框架, 撰写论文并修改。朱学芳, 讨论研究思路, 提出修改意见。
- 录用日期: 2020-12-28

(上接第177页)

- [14] WANG R, LI Z, CAO J, et al. Convolutional recurrent neural networks for text classification [C]. 2019 International Joint Conference on Neural Networks (IJCNN), 2019.
- [15] ZAREMBA W, SUTSKEVER I, VINYALS O. Recurrent neural network regularization [J]. arXiv preprint arXiv: 1409.2329, 2014.
- [16] KIM Y. Convolutional neural networks for sentence classification [J]. arXiv: 1408.5882v2, 2014.
- [17] CHOI K, FAZEKAS G, SANDLER M, et al. Convolutional recurrent neural networks for music classification [C]. IEEE International Conference on Acoustics. IEEE, 2017: 2392-2396.
- [18] SONG Changhe, TU Cunchao, et al. CED: credible early detection of social media rumors [J]. arXiv preprint arXiv: 1811.04175, 2018.
- [19] LIN D, MA B, CAO D, et al. Chinese microblog rumor detection based on deep sequence context [J]. Concurrency and Computation: Practice and Experience, 2019, 31 (23): E4508.
- 作者简介: 李悦晨, 女, 1997年生, 硕士生。研究方向: 自然语言处理, 谣言检测。钱玲飞 (ORCID: 0000-0001-6406-7069), 女, 1979年生, 博士, 副教授, 硕士生导师。研究方向: 科学评价, 智能信息处理, 大数据分析。马静 (通信作者), 女, 1968年生, 教授, 博士生导师。研究方向: 大数据网络舆情分析, 企业信息化。
- 作者贡献声明: 李悦晨, 资料搜集, 实验验证及初稿撰写。钱玲飞, 文献指导, 修改及补充。马静, 提供研究思路, 论文最终版本修订。
- 录用日期: 2020-12-28