

# 基于BERT-PGN模型的中文新闻文本自动摘要生成

谭金源<sup>1</sup>, 刁宇峰<sup>1</sup>, 祁瑞华<sup>2</sup>, 林鸿飞<sup>1\*</sup>

(1. 大连理工大学 计算机科学与技术学院, 辽宁 大连 116024; 2. 大连外国语大学 语言智能研究中心, 辽宁 大连 116024)

(\* 通信作者电子邮箱 hflin@dlut.edu.cn)

**摘 要:**针对文本自动摘要任务中生成式摘要模型对句子的上下文理解不够充分、生成内容重复的问题,基于BERT和指针生成网络(PGN),提出了一种面向中文新闻文本的生成式摘要模型——BERT-指针生成网络(BERT-PGN)。首先,利用BERT预训练语言模型结合多维语义特征获取词向量,从而得到更细粒度的文本上下文表示;然后,通过PGN模型,从词表或原文中抽取单词组成摘要;最后,结合coverage机制来减少重复内容的生成并获取最终的摘要结果。在2017年CCF国际自然语言处理与中文计算会议(NLPCC2017)单文档中文新闻摘要评测数据集上的实验结果表明,与PGN、伴随注意力机制的长短时记忆神经网络(LSTM-attention)等模型相比,结合多维语义特征的BERT-PGN模型对摘要原文的理解更加充分,生成的摘要内容更加丰富,全面且有效地减少重复、冗余内容的生成,Rouge-2和Rouge-4指标分别提升了1.5%和1.2%。

**关键词:**生成式摘要模型;预训练语言模型;多维语义特征;指针生成网络;coverage机制

**中图分类号:**TP391.1 **文献标志码:**A

## Automatic summary generation of Chinese news text based on BERT-PGN model

TAN Jinyuan<sup>1</sup>, DIAO Yufeng<sup>1</sup>, QI Ruihua<sup>2</sup>, LIN Hongfei<sup>1\*</sup>

(1. School of Computer Science and Technology, Dalian University of Technology, Dalian Liaoning 116024, China;

2. Language Intelligence Research Center, Dalian University of Foreign Languages, Dalian Liaoning 116024, China)

**Abstract:** Aiming at the problem that the abstractive summarization model in text automatic summarization task does not fully understand the context of sentence and generates duplicate contents, based on BERT (Bidirectional Encoder Representations from Transformers) and Pointer Generator Network (PGN), an abstractive summarization model for Chinese news text was proposed, namely Bidirectional Encoder Representations from Transformers-Pointer Generator Network (BERT-PGN). Firstly, combining with multi-dimensional semantic features, the BERT pre-trained language model was used to obtain the word vectors, thereby obtaining a more fine-grained text context representation. Then, through PGN model, the words were extracted from the vocabulary or the original text to form a summary. Finally, the coverage mechanism was combined to reduce the generation of duplicate contents and obtain the final summarization result. Experimental results on the single document Chinese news summary evaluation dataset of the 2017 CCF International Conference on Natural Language Processing and Chinese Computing (NLPCC2017) show that, compared with models such as PGN and Long Short-Term Memory with attention mechanism (LSTM-attention), the BERT-PGN model combined with multi-dimensional semantic features has a better understanding of the original text of the summary, has the generated summary content richer and more comprehensive with the generation of duplicate and redundant contents effectively reduced, and has Rouge-2 and Rouge-4 indicators increased by 1.5% and 1.2% respectively.

**Key words:** abstractive summarization model; pre-trained language model; multi-dimensional semantic feature; Pointer Generator Network (PGN); coverage mechanism

## 0 引言

随着近些年互联网产业的飞速发展,大量的新闻网站、新闻手机软件出现在日常生活中,越来越多的用户通过新闻网站、手机软件快速获取最新资讯。根据中国互联网络信息中心(China Internet Network Information Center, CNNIC)第42次

发展统计报告,到2018年6月,中国的移动电话用户规模达到7.88亿,网民接入互联网的比例也在增加,通过手机达到98.3%<sup>[1]</sup>。网友人数增多、新闻媒体网络平台使用率不断提升,网友们使用今日头条等新闻媒体的频率也不断提升。

为了适应当下快节奏的生活,网友需要阅读最少的新闻字数,获取新闻文章的关键内容。网友们可以通过文本自动

收稿日期:2020-05-31;修回日期:2020-07-07;录用日期:2020-07-08。

基金项目:国家重点研发计划项目(2019YFC1200302);国家自然科学基金重点项目(61632011)。

作者简介:谭金源(1997—),男,辽宁大连人,硕士研究生,主要研究方向:自然语言处理;刁宇峰(1987—),女,辽宁沈阳人,博士研究生,主要研究方向:自然语言处理;祁瑞华(1974—),女,湖北襄樊人,教授,博士,主要研究方向:自然语言处理;林鸿飞(1962—),男,辽宁大连人,教授,博士,主要研究方向:自然语言处理。

摘要技术,概括出新闻的主要内容,节省阅读时间,提升信息使用效率。因此,本文提出的面向新闻的文本自动摘要模型具有重要意义。

国内外学者针对文本自动摘要已经做了大量的研究。文本自动摘要是 20 世纪 50 年代出现的一种用计算机完成的文本摘要技术,帮助人们从信息海洋中解放,提高信息的使用效率<sup>[2]</sup>。自 2001 年美国国家标准技术研究所举办文档理解会议以来,文本自动摘要研究得到了越来越多的关注<sup>[3]</sup>。

本文受文献[4]启发,针对网友阅读理解新闻时需要花费大量时间的问题,基于 BERT (Bidirectional Encoder Representations from Transformers) 和指针生成网络 (Pointer Generator Network, PGN),提出了一种面向中文新闻文本的自动摘要模型——BERT-指针生成网络 (Bidirectional Encoder Representations from Transformers-Pointer Generator Network, BERT-PGN),能够有效节省时间,提高信息使用效率。该模型首先利用 BERT 预训练语言模型获取新闻文本的词向量,结合多维语义特征对新闻中的词所在的句子进行打分,其结果作为输入序列输入到指针生成网络中进行训练,得到新闻摘要的结果。

本文主要贡献如下。

1) 本文提出了一种面向新闻文本进行自动摘要的模型——BERT-PGN,分为两个阶段实现:基于预训练模型及多维语义特征的词向量获取阶段以及基于指针生成网络模型的句子生成阶段。

2) 实验结果表明,该模型在 2017 年 CCF 国际自然语言处理与中文计算会议 (the 2017 CCF International Conference on Natural Language Processing and Chinese Computing, NLPCC2017) 单文档中文新闻摘要评测数据集上取得了很好的效果,Rouge-2 和 Rouge-4 指标分别提升 1.5% 和 1.2%。

## 1 相关工作

自动文本摘要有两种主流方式,即抽取式摘要和生成式摘要<sup>[5]</sup>。在对文本进行语义挖掘的研究中,许多经典的分类、聚类算法被先后提出<sup>[6]</sup>。最早的摘要工作主要是利用基于词频和句子位置的基于统计的技术<sup>[7]</sup>。1958 年,Luhn<sup>[8]</sup>提出了第一个自动文本摘要系统。近十几年来,随着机器学习 (Machine Learning, ML) 以及自然语言处理 (Natural Language Processing, NLP) 的快速发展,许多准确高效的文本摘要算法被提出<sup>[9]</sup>。互联网作为商业媒介快速发展,导致用户吸收了太多信息。为了解决这种信息过载,文本自动摘要起到了关键作用。文本自动摘要可以在屏蔽大量干扰文本的同时,让用户更加快捷地获取关键信息,适应当下快节奏的生活<sup>[10]</sup>。

抽取式摘要方法是为一篇文章分成小单元,然后将其中的一些作为这篇文章的摘要进行提取。Liu 等<sup>[11]</sup>提出了一个抽取式文本摘要的对抗过程,使用生成对抗网络 (Generative Adversarial Network, GAN) 模型获得了具有竞争力的 Rouge 分

数,该方法可以生成更多抽象、可读和多样化的文本摘要;Al-Sabahi 等<sup>[12]</sup>使用分层结构的自注意力机制模型 (Hierarchical Structured Self-Attentive Model, HSSAM),反映文档的层次结构,进而获得更好的特征表示,解决因占用内存过大模型无法充分建模等问题;Slamet 等<sup>[13]</sup>提出了一种向量空间模型 (Vector Space Model, VSM),利用 VSM 进行单词相似性测试,对文本自动摘要的结果进行测评,比较文本摘要实现的效果;Alguliyev 等<sup>[14]</sup>发现,与传统文本自动摘要方法相比,基于聚类、优化和进化算法的文本自动摘要研究最近表现出了良好的效果。但抽取式摘要并未考虑文本的篇章结构信息,缺少对文本中关键字、词的理解,生成的摘要可读性、连续性较差。

生成式摘要方法是一种利用更先进自然语言处理算法的摘要方法,对文章中的句子进行转述、替换等生成文章摘要,而不使用其中任何现有的句子或短语。随着近些年深度学习的快速发展,越来越多的深度学习方法被利用到文本摘要中。Cho 等<sup>[15]</sup>和 Sutskever 等<sup>[16]</sup>最早提出了由编码器和解码器构成的 seq2seq (sequence-to-sequence) 模型;Tan 等<sup>[17]</sup>提出了基于图的注意力机制神经模型,在文本自动摘要的任务中取得了很好的效果;Siddiqui 等<sup>[18]</sup>在谷歌大脑团队提出的序列到序列模型的基础上进行改进,使用局部注意力机制代替全局注意力机制,在解决生成重复的问题上取得了很好的效果;Celikyilmaz 等<sup>[19]</sup>针对生成文档的摘要,提出了一种基于编码器-解码器体系结构的深层通信代理算法;Khan 等<sup>[20]</sup>提出了一种基于语义角色标记的框架,使用深度学习的方法从语义角色理解的角度实现多文档摘要任务;江跃华等<sup>[21]</sup>提出了一种基于 seq2seq 结构和注意力机制并融合了词汇特征的生成式摘要算法,能在摘要生成过程中利用词汇特征识别更多重点词汇内容,进一步提高摘要生成质量。

现阶段大多数的文本自动摘要方法主要是利用机器学习或深度学习模型自动提取特征,利用模型进行摘要句子的选取及压缩。但自动提取的特征和摘要文本会存在不充分、不贴近的情况,不能很好地刻画摘要文本。本文提出的 BERT-PGN 模型基于 BERT 预训练语言模型及多维语义特征,针对中文新闻文本,从更多维度进行特征抽取,深度刻画摘要文本,能够得到更贴近主题的摘要内容。

## 2 BERT-PGN 模型

本文提出的 BERT-PGN 模型主要分成两个阶段实现,即基于预训练模型及多维语义的词向量获取阶段以及基于指针生成网络模型的句子生成阶段,如图 1 所示。该模型第一阶段利用预训练语言模型 BERT 获取新闻文章的词向量,同时利用多维语义特征对新闻中的句子进行打分,将二者进行简单拼接生成输入序列;第二阶段将得到的输入序列输入到指针生成网络模型中,使用 coverage 机制减少生成重复文字,同时保留生成新文字的能力,得到新闻摘要。

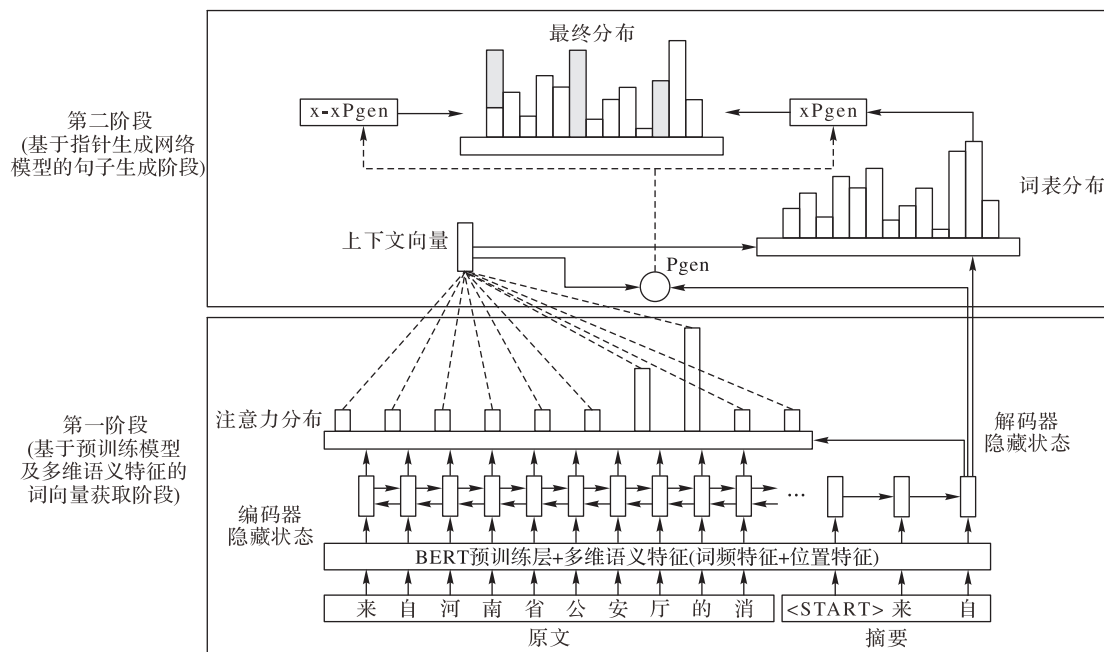


图1 BERT-PGN模型

Fig. 1 BERT-PGN model

## 2.1 基于预训练模型及多维语义特征的词向量获取阶段

### 2.1.1 BERT预训练语言模型

语言模型是自然语言处理领域一个比较重要的概念,利用语言模型对客观事实进行描述后,能够得到可以利用计算机处理的语言表示。语言模型用来计算任意语言序列 $a_1, a_2, \dots, a_n$ 出现的概率 $p(a_1, a_2, \dots, a_n)$ ,即:

$$p(a_1, a_2, \dots, a_n) = \prod_{i=1}^n p(a_i | a_1, a_2, \dots, a_{i-1}) \quad (1)$$

通过传统的神经网络语言模型获取的词向量是单一固定的,存在无法表示字的歧义性等问题。预训练语言模型很好地解决了这一问题,能够结合字的上下文内容来表示字。BERT采用双向Transformer作为编码器进行特征抽取,能够获取到更多的上下文信息,极大程度地提升了语言模型抽取特征的能力。Transformer编码单元包含自注意力机制和前馈神经网络两部分。自注意力机制的输入部分是由来自同一个字的三个不同向量构成的,分别Query向量(Q),Key向量(K)和Value向量(V)。通过Query向量和Key向量相乘来表示输入部分字向量之间的相似度,记做 $[QK]^T$ ,并通过 $d_k$ 进行缩放,保证得到的结果大小适中。最后经过softmax进行归一化操作,得到概率分布,进而得到句子中所有词向量的权重和表示。这样得到的词向量结合了上下文信息,表示更准确,计算方法如下:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (2)$$

BERT预训练模型使用了“MultiHead”模式,即使用了多个注意力机制获取句子的上下文语义信息,称为多头注意力机制。BERT预训练语言模型能够使词向量获取更多的上下文信息,更好地表示原文内容。

### 2.1.2 多维语义特征

针对中文新闻重点内容集中在新闻开头、关键词出现频率高等特点,本文引入了传统特征以及主题特征对中文新闻文本中的句子进行细粒度的描述,提升对文本中句子的上下文语义表述性能。

#### 1) 传统特征。

本文所选择的传统特征主要为句子层次的两种特征:句子中的词频以及在文章中的位置。

词频特征是反映新闻文章中最重要信息的一种统计特征,也是最简单、最直接的一种统计特征。新闻文章中出现词的词频可以利用式(3)进行计算:

$$tf(word) = \frac{word}{\sum_{j=1}^n word_j} \quad (3)$$

其中, $word_j$ 代表文章中第j个词出现的次数。

在本文中,选择文章中的句子作为最终的打分基本单位。句子是词的集合,如果句子包含的词语中,有在新闻文章中频繁出现的高频词,则认为这个句子在文章中更加重要。新闻文章中第i个句子的词频特征打分公式如下:

$$TF_i = \sum_{word \in sen_i} tf(word) \quad (4)$$

其中: $TF_i$ 表示第i个句子中包含的词的词频之和, $sen_i$ 代表第i个句子中包含的所有词。

位置特征同样是反映新闻文章中重要信息的一种统计特征。一篇新闻文章是由多个句子组成的,句子所在的位置不同,其代表的重要性也不同,例如文章中的第一个句子大多是新闻文章中最重要的一句话。新闻文章中第i个句子的位置特征打分公式如下:

$$Pos_i = \frac{n - p_i + 1}{n} \quad (5)$$

其中: $Pos_i$ 代表第i个句子的位置得分, $p_i$ 代表第i个句子在新闻文章中的位置, $n$ 代表文章中的句子总个数。

#### 2) 主题特征。

本文选取的主题特征也可表述为标题特征。新闻文章中的标题具有很高的参考价值,很大程度上可以代表文章中的主题。因此,如果文章中的句子与新闻文章的标题有较高的相似度,那么这个句子更容易被选择为文章摘要中的句子。本文使用余弦相似度计算新闻文章中第i个句子的主题特征得分,打分公式如下:



$$Sim_i = \frac{s \cdot t}{\sqrt{\sum_{j=1}^n s_j^2} \cdot \sqrt{\sum_{j=1}^n t_j^2}} \quad (6)$$

其中:  $Sim_i$  表示第  $i$  个句子与新闻文章标题的相似度,  $s$  和  $t$  分别代表标题和新闻文章中句子的向量化表示。

## 2.2 基于指针生成网络模型的句子生成阶段

指针生成网络模型结合了指针网络(Pointer Network, PN)和基于注意力机制的序列到序列模型,允许通过指针直接指向生成的单词,也可以从固定的词汇表中生成单词。文本中的文字  $w_i$  依次传入 BERT-多维语义特征编码器、双向长短时记忆神经网络(Bidirectional Long Short-Term Memory, Bi-LSTM)编码器,生成隐层状态序列  $h_i$ 。在  $t$  时刻,长短时记忆(Long Short-Term Memory, LSTM)神经网络解码器接收上一时刻生成的词向量,得到解码状态序列  $s_t$ 。

注意力分布  $a^t$  用来确定  $t$  时刻输出序列字符时,输入序列中需要关注的字符。计算公式如下:

$$e_i^t = v^T \tanh(W_h h_i + W_s s_t + b_{\text{attn}}) \quad (7)$$

$$a_i^t = \text{softmax}(e_i^t) \quad (8)$$

其中,  $v, W_h, W_s, b_{\text{attn}}$  是通过训练得到的参数。利用注意力分布对编码器隐层状态加权平均,生成上下文向量  $h_i^*$ 。

$$h_i^* = \sum_j a_j^t h_j \quad (9)$$

将上下文向量  $h_i^*$  与解码状态序列  $s_t$  串联,通过两个线性映射,生成当前预测在词典上的分布  $P_{\text{vocab}}$ , 计算公式如下:

$$P_{\text{vocab}} = \text{softmax}(V'(V[s_t, h_i^*] + b) + b') \quad (10)$$

其中,  $V', V, b, b'$  是通过训练得到的参数。

模型利用生成概率  $P_{\text{gen}}$  来确定复制单词还是生成单词, 计算公式如下:

$$P_{\text{gen}} = \sigma(w_h^T h_i^* + w_s^T s_t + w_x^T x_t + b_{\text{ptr}}) \quad (11)$$

其中,  $w_h, w_s, w_x, b_{\text{ptr}}$  是通过训练得到的参数,  $\sigma$  是 sigmoid 函数,  $x_t$  是解码输入序列。将  $a^t$  作为模型输出, 得到生成单词  $w$  的概率分布:

$$P(w) = P_{\text{gen}} P_{\text{vocab}}(w) + (1 - P_{\text{gen}}) \sum_{i: w_i = w} a_i^t \quad (12)$$

为了解决生成词语重复的问题,本文引入了 coverage 机制。通过 coverage 机制对指针生成网络模型进行改进,能够有效减少生成摘要中的重复。引入 coverage 向量  $c^t$  跟踪已经生成的单词,并对已经生成的单词施加一定的惩罚,尽量减少生成重复。coverage 向量  $c^t$  计算方式如下:

$$c^t = \sum_{i'=0}^{t-1} a^{i'} \quad (13)$$

通俗来说,  $c^t$  表示目前为止单词从注意力机制中获得的覆盖程度。使用 coverage 向量  $c^t$  影响注意力分布,重新得到注意力分布  $a^t$ , 计算公式如下:

$$e_i^t = v^T \tanh(W_h h_i + W_s s_t + W_c c_i^t + b_{\text{attn}}) \quad (14)$$

其中  $W_c$  是通过训练得到的参数。

## 3 实验与分析

### 3.1 实验数据

本文的实验部分使用的数据是由 2017 年 CCF 国际自然语言处理与中文计算会议(NLPCC2017)提供,来自于 NLPCC2017 中文单文档新闻摘要评测数据集,包含训练集新闻文本 49 500 篇,测试集新闻文本 500 篇。该任务中要求生成的摘要长度不超过 60 个字符。

### 3.2 评价指标

Rouge 是文本自动摘要领域摘要评价技术的通用指标之一,通过统计模型生成的摘要与人工摘要之间重叠的基本单

元,评判模型生成摘要的质量。本文参考 NLPCC2017 中文单文档新闻摘要评测任务,使用 Rouge-2、Rouge-4 和 Rouge-SU4 作为评价指标,对摘要结果进行评价。

### 3.3 对比实验

本文实验部分选取 8 种基本模型: NLPCC2017 单文档新闻摘要评测任务结果较好团队(ccnuSYS、LEAD、NLP@WUST、NLP\_ONE)提出的模型<sup>[22]</sup>、PGN(without coverage mechanism)<sup>[23]</sup>、PGN<sup>[23]</sup>、主题关键词信息融合模型<sup>[24]</sup>以及 BERT-PGN(without semantic features)。对人工提取的主题特征、传统特征进行特征的有效性验证,验证本文提出方法的有效性。

1) ccnuSYS<sup>[22]</sup>: 使用基于注意力机制的 LSTM 编码器-解码器结构模型生成摘要。

2) LEAD<sup>[22]</sup>: 从原文选取前 60 个字作为文本摘要。

3) NLP@WUST<sup>[22]</sup>: 使用特征工程的方法进行句子抽取,并利用句子压缩算法对抽取的句子进行压缩。

4) NLP\_ONE<sup>[22]</sup>: NLPCC2017 单文档新闻摘要评测任务第一名的算法,包含输入、输出序列的注意力机制。

5) PGN(without coverage mechanism)<sup>[23]</sup>: ACL2017 中提出的一种生成模型,使用指针网络和基于注意力机制的序列到序列模型生成摘要,不使用 coverage 机制。

6) PGN(coverage mechanism)<sup>[23]</sup>: 改进的指针生成网络模型,利用 coverage 机制解决生成重复词和未登录词的问题。

7) 主题关键词融合模型<sup>[24]</sup>: 一种结合主题关键词信息的多注意力机制模型。

8) BERT-PGN(without semantic features): 本文提出的一种基于 BERT 和指针生成网络的模型,利用 coverage 机制减少生成重复内容。

9) BERT-PGN(semantic features): 在 BERT-PGN(without semantic features)模型上进行优化得到的模型,结合多维语义特征获取细粒度的文本上下文表示。

### 3.4 实验环境及参数设置

本文实验使用单个 GTX-1080Ti(GPU)进行训练。本实验获取文本词向量使用 BERT-base 预训练模型。BERT-base 模型共 12 层,隐层 768 维。设置最大序列长度为 128,  $\text{train\_batch\_size}$  为 16,  $\text{learning\_rate}$  为  $5E-5$ 。

指针生成网络模型设置  $\text{batch\_size}$  为 8,隐层 256 维,设置字典大小为 50k。训练过程共进行 700k 次迭代,训练总时长约为 7 d5 h(合计 173 h)。

### 3.5 实验结果与分析

#### 3.5.1 总体摘要结果对比实验

本文重新运行了部分 baseline 模型,将获取的结果与本文提出的模型结果做对比,实验结果如表 1。

表 1 总体摘要结果对比

Tab. 1 Results comparison of overall summarization

模型	Rouge-2	Rouge-4	Rouge-SU4
ccnuSYS*	14.98	7.66	13.89
LEAD*	20.88	11.69	19.14
NLP@WUST*	22.53	10.39	20.81
NLP_ONE*	22.87	12.78	21.18
PGN(without coverage mechanism)	23.15	12.33	21.40
PGN(coverage mechanism)	23.57	12.57	21.56
主题关键词融合模型	24.02	12.89	22.35
BERT-PGN(without semantic features)	24.32	13.09	22.68
BERT-PGN(semantic features)	25.53	14.06	23.74

注: 带\*方法的结果引自文献[22]。

从表 1 可以看出,本文提出的模型性能相较于 PGN、NLP\_ONE 等模型有了显著的提升,在 Rouge-2、Rouge-4 以及 Rouge-SU4 的评价指标中有着明显的优势,Rouge 指标提升了 1.2~1.5 个百分点。

由 BERT-PGN(semantic features)模型与 PGN、BERT-PGN(without semantic features)模型进行对比,可以看出使用 BERT 预训练模型并结合有效的多维人工特征,能够显著提升模型效果。使用 BERT 预训练模型并结合人工抽取的特征得到的

句子上下文表示,对文本中句子的语义理解更加深刻、准确,在文本自动摘要任务中能够有效提升性能。

根据表 2 不同模型生成摘要的内容可以发现,本文提出的 BERT-PGN 模型相较于其他模型,在中文新闻文本的自动摘要任务中生成的摘要内容更丰富、更全面、更贴近标准摘要,说明该模型对全文的理解更加充分,能够结合文中句子的上下文充分理解句子、词语的含义,对文中的句子、词语进行更细致的刻画。

表 2 摘要结果示例

Tab. 2 Summarization result examples

模型	摘要内容
标准摘要	河南出台举报暴恐线索奖励办法,最高奖 5 万元;企业商户未及时向公安机关报告,可依法追究责任
LEAD	来自河南省公安厅的消息,为做到预知预警,防患未然,有效维护公共安全,该厅制定出台《举报暴力恐怖违法犯罪线索奖励办法》
NLP@WUST	来自河南省公安厅的消息,该厅制定出台《举报暴力恐怖违法犯罪线索奖励办法》,鼓励公民和组织积极举报暴恐犯罪线索
主题关键词融合模型	河南省公安厅出台《举报暴力恐怖违法犯罪线索奖励办法》,有关企业、商户、从业人员发现违法犯罪可疑人员及时向公安机关报告
BERT-PGN	河南省公安厅制定出台《举报暴力恐怖违法犯罪线索奖励办法》,有关企业、商户、从业人员未及时向公安机关报告,依法追究

### 3.5.2 多维语义特征对比实验

多维特征选取的部分,本文针对新闻文本“主要内容集中在开头部分”的特点,选取传统特征、主题特征中的词频特征、位置特征以及标题特征,分别表示为 TF、Pos 以及 Main。

由表 3 可以看出,同一模型结合人工提取的词频特征和位置特征效果最好,Rouge-2 指标最多提升了 1.2 个百分点,Rouge-4 指标最多提升了 1.0 个百分点。

表 3 特征组合结果对比

单位: %

Tab. 3 Feature combination result comparison

unit: %

特征选择	Rouge-2	Rouge-4	Rouge-SU4
None	24.32	13.09	22.68
Main	24.74	13.47	22.85
Pos	25.04	13.81	23.23
Pos+Main	25.15	13.86	23.30
TF	25.35	13.93	23.46
TF+Main	25.37	13.92	23.51
TF+Pos+Main	25.44	13.96	23.58
TF+Pos	25.53	14.06	23.74

本文选取的主题特征 Main 能够在一定程度上提升模型的 Rouge 指标。从 Pos 和 Pos+Main、TF 和 TF+Main 的特征组合结果对比可以得知,主题特征结合词频特征时提升明显,结合位置特征时基本没有提升。句子在新闻中的位置靠前时,与标题的相似度也更高,说明两种人工特征在衡量句子在新闻中的重要性时起到了相似的作用。通过对比 TF+Main 和 TF+Pos 两种特征组合的结果可以得知,词频信息结合位置信息相较于结合主题信息效果更好,能够充分表达句子在新闻文章中的重要性。因此,本文选择使用词频特征以及位置特征的特征组合作为多维特征。

新闻文章中多次出现的关键词,是反映新闻文章中最重要的信息的一种统计特征,进行词频统计的意义在于找出文章表达的重点;此外,句子出现的位置也是反映句子重要程度的关键,出现的位置越靠前,说明该句子在文章中起到的作用越大。因此,词频、位置特征是自动摘要模型提升的关键。

### 3.5.3 coverage 机制实验分析

本文使用的模型使用了 coverage 机制,试图解决生成重复内容的问题。通过计算生成摘要中 1-gram、2-gram、3-gram 以及 4-gram 所占比例,定量分析引入 coverage 机制解决生成内容重复问题的效果。

由表 4 可以看出,本文提出的 BERT-PGN 模型相较于 NLP\_ONE 能够有效减少生成内容的重复,在解决重复的方面效果明显,在 3-gram、4-gram 的摘要结果定量分析中,接近标准摘要的效果。

表 4 coverage 机制验证

单位: %

Tab. 4 Verification of coverage mechanism

unit: %

数据片段	NLP_ONE	BERT-PGN	参考模型
1-gram	26	19	17
2-gram	18	6	3
3-gram	16	2	1
4-gram	14	1	0

## 4 结语

本文提出了一种面向中文新闻文本的 BERT-PGN 模型,结合 BERT 预处理模型及多维语义特征获取词向量,利用指针生成网络模型结合 coverage 机制减少生成重复内容。经实验表明,BERT-PGN 模型在中文新闻摘要任务中,生成的摘要结果更接近标准摘要,包含更多原文的关键信息,能有效解决生成内容重复的问题。

下一步将尝试挖掘更多要素,例如:面向新闻文本的有效人工特征等,提升摘要结果;简化模型,缩短模型训练时间;提升生成摘要内容的完整性、流畅性;构建新闻领域的外部数据,帮助模型结合句子上下文充分理解句子含义。

### 参考文献 (References)

- [1] LING H. Innovative exploration on empowering students to study and propagandize the new thought under the Internet thinking taking Xi Jinping thought on the socialism with Chinese characteristics in a new era as an example[J]. PEOPLE: International Journal of Social Sciences, 2018, 4(3):1395-1408.
- [2] 刘家益,邹益民. 近 70 年文本自动摘要研究综述[J]. 情报科

- 学, 2017, 35(7): 154-161. (LIU J Y, ZOU Y M. A review of automatic text summarization research in recent 70 years [J]. Information Science, 2017, 35(7): 154-161.)
- [3] 吴云, 杨长春, 梅佳俊, 等. 词句协同自动摘要提取方法[J]. 计算机工程与设计, 2018, 39(9): 2776-2779, 2810. (WU Y, YANG C C, MEI J J, et al. Method of automatic summarization algorithm based on word-sentence co-ranking [J]. Computer Engineering and Design, 2018, 39(9): 2776-2779, 2810.)
- [4] WEI R, HUANG H, GAO Y. Sharing pre-trained BERT decoder for a hybrid summarization [C]// Proceedings of the 2019 China National Conference on Chinese Computational Linguistics, LNCS 11856. Cham: Springer, 2019: 169-180.
- [5] 石磊, 阮选敏, 魏瑞斌, 等. 基于序列到序列模型的生成式文本摘要研究综述[J]. 情报学报, 2019, 38(10): 1102-1116. (SHI L, RUAN X M, WEI R B, et al. Abstractive summarization based on sequence to sequence models: a review [J]. Journal of the China Society for Scientific and Technical Information, 2019, 38(10): 1102-1116.)
- [6] 戴天, 吴渝, 雷大江. 利用组合模型生成微博热点话题事件摘要[J]. 计算机应用研究, 2016, 33(7): 2026-2029, 2038. (DAI T, WU Y, LEI D J. Hot topic summarization on microblog generated by model combination [J]. Application Research of Computers, 2016, 33(7): 2026-2029, 2038.)
- [7] IBOI H, CHUA S, RANAIVO-MALANÇON B, et al. Performance of opinion summarization towards extractive summarization [J]. Journal of Telecommunication, Electronic and Computer Engineering, 2017, 9(2-10): 57-64.
- [8] LUHN H P. The automatic creation of literature abstracts [J]. IBM Journal of Research and Development, 1958, 2(2): 159-165.
- [9] 侯圣峦, 张书涵, 费超群. 文本摘要常用数据集和方法研究综述[J]. 中文信息学报, 2019, 33(5): 1-16. (HOU S L, ZHANG S H, FEI C Q. A survey to text summarization: popular datasets and methods [J]. Journal of Chinese Information Processing, 2019, 33(5): 1-16.)
- [10] LERMAN K, BLAIR-GOLDENSOHN S, MCDONALD R. Sentiment summarization: evaluating and learning user preferences [C]// Proceedings of the 12th Conference of the European Chapter of the Association for Computational Linguistics. Stroudsburg, PA: Association for Computational Linguistics, 2009: 514-522.
- [11] LIU L, LU Y, YANG M, et al. Generative adversarial network for abstractive text summarization [C]// Proceedings of the 32nd AAAI Conference on Artificial Intelligence. Palo Alto, CA: AAAI Press, 2018: 8109-8110.
- [12] AL-SABAHI K, ZHANG Z, NADHER M. A Hierarchical Structured Self-Attentive model for extractive document Summarization (HSSAS) [J]. IEEE Access, 2018, 6: 24205-24212.
- [13] SLAMET C, ATMADIA A R, MAYLAWATI D S, et al. Automated text summarization for Indonesian article using vector space model [J]. IOP Conference Series: Materials Science and Engineering, 2018, 288: No. 012037.
- [14] ALGULIYEV R M, ALIGULIYEV R M, ISAZADE N R, et al. COSUM: text summarization based on clustering and optimization [J]. Expert Systems, 2019, 36(1): e12340.
- [15] CHO K, VAN MERRIENBOER B, GULCEHRE C, et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation [C]// Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing. Stroudsburg, PA: Association for Computational Linguistics, 2014: 1724-1734.
- [16] SUTSKEVER I, VINYALS O, LE Q V. Sequence to sequence learning with neural networks [C]// Proceedings of the 27th International Conference on Neural Information Processing Systems. Cambridge: MIT Press, 2014: 3104-3112.
- [17] TAN J, WAN X, XIAO J. Abstractive document summarization with a graph-based attentional neural model [C]// Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. Stroudsburg, PA: Association for Computational Linguistics, 2017: 1171-1181.
- [18] SIDDIQUI T, SHAMSI J A. Generating abstractive summaries using sequence to sequence attention model [C]// Proceedings of the 2018 International Conference on Frontiers of Information Technology. Piscataway: IEEE, 2018: 212-217.
- [19] CELIKYILMAZ A, BOSSELTUT A, HE X, et al. Deep communicating agents for abstractive summarization [C]// Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies. Stroudsburg, PA: Association for Computational Linguistics, 2018: 1662-1675.
- [20] KHAN A, SALIM N, JAYA KUMAR Y. A framework for multi-document abstractive summarization based on semantic role labelling [J]. Applied Soft Computing, 2015, 30: 737-747.
- [21] 江跃华, 丁磊, 李娇娥, 等. 融合词汇特征的生成式摘要模型[J]. 河北科技大学学报, 2019, 40(2): 152-158. (JIANG Y H, DING L, LI J E, et al. Abstractive summarization model considering hybrid lexical features [J]. Journal of Hebei University of Science and Technology, 2019, 40(2): 152-158.)
- [22] HUA L, WAN X, LI L. Overview of the NLPCC 2017 shared task: single document summarization [C]// Proceedings of the 2017 National CCF Conference on Natural Language Processing and Chinese Computing, LNCS 10619. Cham: Springer, 2017: 942-947.
- [23] SEE A, LIU P J, MANNING C D. Get to the point: summarization with pointer-generator networks [C]// Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics. Stroudsburg, PA: Association for Computational Linguistics, 2017: 1073-1083.
- [24] 侯丽微, 胡珀, 曹雯琳. 主题关键词信息融合的中文生成式自动摘要研究[J]. 自动化学报, 2019, 45(3): 530-539. (HOU L W, HU P, CAO W L. Automatic Chinese abstractive summarization with topical keywords fusion [J]. Acta Automatica Sinica, 2019, 45(3): 530-539.)

This work is partially supported by the National Key Research and Development Program of China (2019YFC1200302), the Key Project of National Natural Science Foundation of China (61632011).

**TAN Jinyuan**, born in 1997, M. S. candidate. His research interests include natural language processing.

**DIAO Yufeng**, born in 1987, Ph. D. candidate. Her research interests include natural language processing.

**QI Ruihua**, born in 1974, Ph. D., professor. Her research interests include natural language processing.

**LIN Hongfei**, born in 1962, Ph. D., professor. His research interests include natural language processing.