



计算机科学

COMPUTER SCIENCE

基于GAN和中文词汇网的文本摘要技术

刘晓影, 王淮, 乌吉斯古楞

引用本文

刘晓影, 王淮, 乌吉斯古楞. [基于GAN和中文词汇网的文本摘要技术](#) [J]. 计算机科学, 2022, 49(12): 301-304.

LIU Xiao-ying, WANG Huai, WU Jisiguleng. [GAN and Chinese WordNet Based Text Summarization Technology](#) [J]. Computer Science, 2022, 49(12): 301-304.

相似文章推荐 (请使用火狐或 IE 浏览器查看文章)

Similar articles recommended (Please use Firefox or IE to view the article)

[一体化网络多终端接入智能路由技术](#)

Intelligent Routing Technology for Multi-terminal Access in Integrated Network
计算机科学, 2022, 49(12): 332-339. <https://doi.org/10.11896/jsjcx.210900042>

[事件抽取技术研究综述](#)

Survey on Event Extraction Technology
计算机科学, 2022, 49(12): 264-273. <https://doi.org/10.11896/jsjcx.211100226>

[软件需求工程技术综述](#)

Review on Technologies of Requirement Engineering of Software
计算机科学, 2022, 49(11A): 210900132-14. <https://doi.org/10.11896/jsjcx.210900132>

[基于记忆增强 GAN 的异常检测](#)

Memory-augmented GAN-based Anomaly Detection
计算机科学, 2022, 49(11A): 211000202-9. <https://doi.org/10.11896/jsjcx.211000202>

[基于双重指针网络的车货匹配双重序列决策研究](#)

Study on Dual Sequence Decision-making for Trucks and Cargo Matching Based on Dual Pointer Network
计算机科学, 2022, 49(11A): 210800257-9. <https://doi.org/10.11896/jsjcx.210800257>

基于 GAN 和中文词汇网的文本摘要技术

刘晓影 王淮 乌吉斯古楞

华北计算技术研究所网络安全工作组 北京 100083

(xiaoying81@126.com)

摘要 随着神经网络技术的广泛应用,文本摘要技术吸引了越来越多科研人员的注意。由于生成式对抗网络(GANs)具有提取文本特征或学习整个样本的分布并以此产生相关样本点的能力,因此正逐步取代传统基于序列到序列(Seq2seq)的模型,被用于提取文本摘要。利用生成式对抗网络的特点,将其用于生成式的文本摘要任务。提出的生成式对抗模型由3部分组成:一个生成器,将输入的句子编码为更短的文本表示向量;一个可读性判别器,强制生成器生成高可读性的文本摘要;以及一个相似性判别器,作用于生成器,抑制其输出的文本摘要与输入的摘要之间的不相关性。此外,在相似性判别器中,引用中文的WordNet作为外部知识库来增强判别器的作用。生成器使用策略梯度算法进行优化,将问题转化为强化学习。实验结果表明,所提模型得到了较高的ROUGE评测分数。

关键词: 文本摘要;生成式对抗网络;WordNet;强化学习;自然语言处理

中图分类号 TP391

GAN and Chinese WordNet Based Text Summarization Technology

LIU Xiao-ying, WANG Huai and WU Jisiguleng

Network Security Group, North China Institute of Computing Technology, Beijing 100083, China

Abstract Since the introduction of neural networks, text summarization techniques continue to attract the attention of researchers. Similarly, generative adversarial networks(GANs) can be used for text summarization because they can generate text features or learn the distribution of the entire sample and produce correlated sample points. In this paper, we exploit the features of generative adversarial networks(GANs) and use them for abstractive text summarization tasks. The proposed generative adversarial model has three components: a generator, which encodes the input sentences into shorter representations; a readability discriminator, which forces the generator to create comprehensible summaries; and a similarity discriminator, which acts on the generator to curb the disconnection between the outputted text summarization and the inputted text summarization. In addition, Chinese WordNet is used as an external knowledge base in the similarity discriminator to enhance the discriminator. The generator is optimized using policy gradient algorithm, converting the problem into reinforcement learning. Experimental results show that the proposed model gets high ROUGE evaluation scores.

Keywords Text summarization, Generative adversarial network, WordNet, Reinforcement learning, Natural language processing

1 引言

文本摘要是用较少的语义准确的词来概括给定的整篇文章的技术。通过这样的方式,关键的信息能够得以保留,而文档中不重要的冗余信息则被过滤掉,从而提高人们对信息的接受效率。文本摘要技术分为两种类型:抽取式和生成式。抽取式文本摘要需要从语料库中提取出已经存在的、包含正确含义的词或句子,来维持文档的总体意义。然后,将这些抽取出的词或短语进行组合以形成摘要。而在生成式的文本摘要中,需要在理解了文本的含义之后进行概括,生成新的句子或词,以形成语义相关的摘要。显然,相比抽取式文本摘要,

生成式文本摘要的任务具有更吸引人的应用场景,因为它与人类总结摘要的方式相似。但是生成式文本摘要任务通常是一个难以完成的任务,比抽取式更难执行,因为它需要全面理解文章的内容才能生成一个好的摘要,这对于机器来说是难以实现的目标。进行自动文本摘要时存在的另一个问题是输入语句的长度不一致。如果输入的语料库和其中的句子很长,则很难捕捉到上下文语义关系并找到语义相关的关键词或短语。此外,受语义表示、推理和自然语言生成等技术的牵制,抽取式摘要的结果通常比生成式摘要的结果更好。这些问题比抽取式摘要这样的基于数据驱动的方法更难,因为数据驱动的方法一般只涉及定位和提取目标单词或短语,不

到稿日期:2021-06-21 返修日期:2021-09-14

基金项目:国家重点研发计划(2018YFC0831200)

This work was supported by the National Key R&D Program of China(2018YFC0831200).

通信作者:王淮(498929906@qq.com)

的卷积过滤器,共 256 核,并按照文献[13]中的预训练方法进行了预训练。然后,将编码器输出的 256 个激活特征图连接到一个具有 512 个神经元的全连接层^[10]。相似性判别器用 SGD 优化算法进行优化,学习率为 0.06,训练批次大小为 128。

3.3 WordNet 特征

TF-IDF 的计算基于两种统计方法,即术语频率(Term Frequency, TF)和逆文档频率(Inverse Document Frequency, IDF)。TF 指一个术语或单词在一个文档中出现的总次数, IDF 则表示出现某一特定术语的文档总数。IDF 的计算还依赖于语料库中存在的所有文档数。例如,在文献[12]中,TF-IDF 是 TF 和 IDF 的乘积。IDF 的计算方式是 N 的对数除以 DF 。这里, N 是语料库中所有的文档数, DF 是一个词出现的文档数量。

WordNet 中的每一个词都构成了一个同义词网络。每个同义词组中的一个词可能对应于该词可能出现的不同的语境。模型在编码当前单词或句子的语义时,若能充分考虑到不同语境下单词所包含的不同语义,就能充分丰富语义。因此可以利用 WordNet 识别同义词组,以进行知识语义增强,并根据它们在文本中出现的频率来调整 TF-IDF 的结果。因此,我们可以通过 WordNet 中的同义词来改进 TF-IDF 的语义表示能力。最终,对短文本与长文本分别应用 WordNet,得到两个向量 w_s 与 w_l 。

3.4 可读性判别器

可读性判别器是另一个基于 CNN 的模型,用于判别生成的摘要是否是生成器生成的还是人类生成的。因此,它本质上是一个二元分类器,其结构包括一个卷积层、一个最大池化层,以及一个用于预测两个类别概率分布的 softmax 分类器。同样地,该判别器按照与相似性判别器相同的预训练方式进行预训练。

判别器将人类书写的参考摘要视为正类别,而由生成器生成的摘要则被视为负类别。可读性判别器的二元分类器中使用的损失函数是依照正负样本的输出概率而定义的,表达式如下:

$$L_D = \frac{1}{N} \sum_{i=1}^N -\log(1 - D(x_i^-)) - \log(D(x_i^+))$$

可读性判别器也是一个 CNN 模型,使用的过滤器窗口大小为 3,4,5。每个大小的过滤器有 128 个,因此能产生一个 128 维的激活特征图。在用生成器和相似性判别器对可读性判别器进行训练之前,需要对可读性判别器进行预训练。在预训练阶段,我们只用生成器来对可读性判别器进行预训练。生成器的输出被视为负样本,而参考文本摘要则是正样本。摘要生成模型的整体架构如图 2 所示。

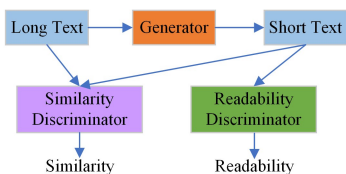


图 2 模型整体结构

Fig. 2 Overall architecture of the proposed model

3.5 策略梯度优化算法

由于无差异性问题的存在,生成器的输出并非直接输入两个判别器中,而是对输出进行采样,将得到的离散序列输入两个判别器中。通过使用 SeqGAN^[11-12] 中的技巧,文本摘要的任务被转换为一个强化学习的问题。当生成器产生语义相关的摘要时,相似性判别器就会给予生成器奖励,鼓励生成器继续生成好的语义相关的摘要。同样地,当生成器产生的摘要是人类可理解的摘要时,可读性判别器也会为生成器产生奖励。生成器的训练用策略梯度优化算法进行优化。策略梯度有助于以强化学习的方式训练生成式对抗网络,并消除监督训练时引入的“暴露偏差”(Exposure Bias)。

4 实验设置与结果分析

4.1 实验数据及评价标准

实验采用的数据集来自国内实际新闻网站的新闻数据。为了充分检验文本摘要模型的有效性,本文使用了多个不同规模大小的数据集。

(1)数据集 1。该数据集包含 350 个新闻网页,有 87 个主题。这些网页均发表于 2007 年 3 月和 4 月,内容都是关于搜索引擎公司的。它们都来自国内的新闻网站,包括新闻报道和评论。最大的主题有 20 篇报道,而最小的只有 1 篇。

(2)数据集 2。该数据集包含 953 个新闻网页,有 108 个主题。这些网页来自国内主要新闻门户网站之一的搜狐网体育频道的新闻。最大的主题有 151 篇新闻报道,而最小的只有 1 篇。

(3)数据集 3。该数据集包含了几十家国内新闻网站体育频道的 24782 个新闻页面以及指定网络日志和网络论坛列表中的 1339 篇文章,发表于 2007 年 9 月和 10 月。

实验采用 ROUGE-N 来评价摘要的效果。ROUGE(Recall-Oriented Understudy for Gisting Evaluation)是一套用于评估自然语言处理中自动文本摘要和机器翻译效果的指标。这些指标将自动产生的摘要或翻译与参考(或一组参考)摘要(如人类书写的摘要)进行比较,得出相应的分值,以衡量自动生成的摘要与参考摘要之间的相似度。ROUGE-N 指摘要系统生成的摘要与参考摘要之间 n -gram 的重叠程度。ROUGE-1 指一元语法(Unigram,或单个单词)的重叠程度,ROUGE-2 指二元语法(Bigram)的重叠程度。其定义如下:

$$ROUGE-N = \frac{\sum_{S \in (\text{ReferenceSummaries})} \sum_{gram_n \in S} \text{Count}_{\text{match}}(gram_n)}{\sum_{S \in (\text{ReferenceSummaries})} \sum_{gram_n \in S} \text{Count}(gram_n)}$$

其中, n 表示 n 元语法的长度,也即 $gram_n$,且 $\text{Count}_{\text{match}}(gram_n)$ 表示候选摘要和一组参考摘要中共有的最大数量。显然,ROUGE-N 是一个基于召回率的衡量指标,因为其分母是参考摘要的 n 元语法数量的总和。此外,为了便于计算,基于 ROUGE-N 还衍生出了 ROUGE-L,其是基于句子 X 与句子 Y 之间的最长公共子序列长度计算得来的。

4.2 实验结果分析

本实验基于 ROUGE-1,ROUGE-2 以及 ROUGE-L 这三个指标来对比本文的摘要生成模型与其他摘要模型之间效果

的优劣。模型在 3 个数据集上的 $R-1$, $R-2$ 以及 $R-L$ 平均值如表 1 所列。

表 1 各个模型 $R-1$, $R-2$ 和 $R-L$ 的结果对比

Table 1 Comparison of results of each model in $R-1$, $R-2$ and $R-L$

方法	$R-1$	$R-2$	$R-L$
Pointer-generator	39.51	17.18	36.28
SummaRuNNer	39.57	16.1	35.3
GAN+policy gradient	37.83	15.7	37.33
Deep-RL	39.85	15.86	36.87
WGAN	35.51	9.36	23.97
GAN	39.91	17.64	36.69
Zhuang 等	40.1	16.22	37.16
本文模型	40.65	17.01	37.53

本文模型在单元语法项的重叠程度(Unigram Overlap-ping)上,即 $R-1$ 和最长公共子序列 $R-L$ 上取得了最优的性能效果。然而,该模型的 $R-2$ 得分,即参考摘要与模型生成的摘要之间的 2 元语法项的重叠程度上稍稍落后。造成这种情况的原因可能是模型失去了完整的 2 元语法项的上下文语义,这可以通过使用在大规模数据集上进行预训练的通用语言模型改进词嵌入表示来进一步改进 2 元语法项的上下文语境表示。

结束语 本文提出了一种基于生成式对抗网络和策略梯度优化算法的生成式文本摘要方法。实验结果表明,所提模型能够对给定长文章或语料生成高可读性的摘要。模型在 3 个测试数据集上的 $R-1$ 和 $R-L$ 得分均达到了最优的效果。所提模型主要包含 3 个部分:生成器、相似性判别器以及可读性判别器。相似性判别器用来提高生成器生成高相关性的文本摘要的能力,而可读性判别器用来提高生成器生成高可读性文本摘要的能力。同时,为了增强相似性判别器的作用,本文提出了一种基于中文 WordNet 的加权 TF-IDF 改进方法,进而提升相似性判别器判断相似程度的能力。模型的 3 个部分均首先进行单独的预训练,然后再用强化学习的策略梯度优化算法进行联合训练。后续的工作可以通过使用文献[14]中介绍的 BERT 词嵌入来进行进一步改善。这些是预先训练好的词向量,可以通过迁移学习使用,也可以自定义训练。这些预训练的通用语言模型有助于进一步提高 ROUGE 指标。

参考文献

- [1] LIN C Y. Rouge: A package for automatic evaluation of summaries[C]//Proceedings of the Workshop on Text Summarization Branches Out(WAS 2004). 2004.
- [2] RUSH A M, CHOPRA S, WESTON J. A neural attention model for abstractive sentence summarization[J]. arXiv:1509.00685, 2015.
- [3] RAMESH N, ZHAI F F, ZHOU B W. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2017.
- [4] ILYA S, VINYALS O, LE Q V. Sequence to sequence learning with neural networks[J]. arXiv:1409.3215, 2014.
- [5] RAMACHANDRAN P, LIU P J, LE Q V. Unsupervised pre-training for sequence to sequence learning [J]. arXiv:1611.02683, 2016.
- [6] SEE A, LIU P J, MANNING C D. Get to the point: Summarization with pointer-generator networks[J]. arXiv:1704.04368, 2017.
- [7] BANAFSHEH R, MOUSAS C, GUPTA B. Generative adversarial network with policy gradient for text summarization [C]//2019 IEEE 13th International Conference on Semantic Computing(ICSC). IEEE, 2019.
- [8] WANG Y S, LEE H Y. Learning to encode text as human-readable summaries using generative adversarial networks[J]. arXiv:1810.02851, 2018.
- [9] ZHUANG H J, ZHANG W B. Generating semantically similar and human-readable summaries with generative adversarial networks[J]. IEEE Access, 2019, 7:169426-169433.
- [10] YU L T, ZHANG W L, WANG J, et al. Seqgan: Sequence generative adversarial nets with policy gradient[C]//Proceedings of the AAAI Conference on Artificial Intelligence. 2017.
- [11] SUTTON R S, MCALLESTER D, SINGH S, et al. Policy gradient methods for reinforcement learning with function approximation[J]. Advances in Neural Information Processing Systems, 2000, 12:1057-1063.
- [12] ZHANG H Y, XU J J, WANG J. Pretraining-based natural language generation for text summarization[J]. arXiv:1902.09243, 2019.
- [13] SARKAR K. Bengali text summarization by sentence extraction [J]. arXiv:1201.2240, 2012.
- [14] DEVLIN J, CHANG M W, LEE K, et al. Bert: Pre-training of deep bidirectional transformers for language understanding[J]. arXiv:1810.04805, 2018.



LIU Xiao-ying, born in 1981, Ph.D., senior engineer, is a member of China Computer Federation. Her main research interests include natural language processing, artificial intelligence and network security.



WANG Huai, born in 1996, master, engineer. His main research interests include network security, threat intelligence analysis and knowledge graph.

(责任编辑:杨雪敏)