



基于用户意图的微博文本生成技术研究

高永兵¹, 黎预璇¹, 高军甜¹, 马占飞²

(1. 内蒙古科技大学 信息工程学院, 内蒙古 包头 014010; 2. 包头师范学院 信息工程系, 内蒙古 包头 014010)

摘要: 微博是个人和组织用户分享或获取简短实时信息的重要社交平台, 微博文本自动生成技术能帮助用户在微博平台上快速实现各种社交意图。为辅助用户发表博文并表达社交意图, 提出一种基于用户意图的微博文本生成技术, 以挖掘提取微博文本特征, 并在给定微博主题的条件下生成与用户意图相一致的微博文本。采用预训练语言模型与微调相结合的方法, 在预训练语言模型 GPT2 上实现联合主题和用户意图的文本控制生成, 以及具备用户对话功能的文本预测生成。实验结果表明, 该技术生成的文本具有较高的可读性且符合微博文本语言风格, 结合主题和 5 类用户意图的生成样本人工评分达 77 分以上。

关键词: 微博文本; 自动生成; 用户意图; 主题; 预训练语言模型; 微调

开放科学(资源服务)标志码(OSID):



中文引用格式: 高永兵, 黎预璇, 高军甜, 等. 基于用户意图的微博文本生成技术研究[J]. 计算机工程, 2022, 48(1): 119-126.

英文引用格式: GAO Y B, LI Y X, GAO J T, et al. Research on Weibo text generation technology based on user intention[J]. Computer Engineering, 2022, 48(1): 119-126.

Research on Weibo Text Generation Technology Based on User Intention

GAO Yongbing¹, LI Yuxuan¹, GAO Juntian¹, MA Zhanfei²

(1. School of Information Engineering, Inner Mongolia University of Science and Technology, Baotou, Inner Mongolia 014010, China;

2. Department of Information Engineering, Baotou Teachers' College, Baotou, Inner Mongolia 014010, China)

[Abstract] Weibo is a mainstream social platform for individuals and organizational users to share or obtain short real-time information. The technique of Weibo text generation can help users quickly realize various social intentions on Weibo. In order to assist users in publishing blog posts and express social intentions, this paper proposes a Weibo text generation technology based on user intention, which mines and extracts Weibo text features, and generates Weibo texts that are consistent with the user intention under a given topic. Using the combination of pre-training language model and fine-tuning, the text control generation of joint topic and user intention and the text prediction generation with user dialogue function are realized on the pre-training language model, GPT2. The experimental results show that the proposed technology can generate texts with high readability and adherence to the language style of Weibo texts. The manual score of the generated samples combined with the theme and five types of user intention is more than 77 points.

[Key words] Weibo text; automatic generation; user intention; topic; pre-training language model; fine-tuning

DOI: 10.19678/j.issn.1000-3428.0060079

0 概述

微博文本自动生成属于文本生成领域的一个研究热点。文本生成是机器根据输入信息, 经过组织规划过程自动生成一段高质量的自然语言文本, 其为自然语言处理研究中一项具有挑战性的任务。文本生成可以形式化为顺序决策过程, 即在每个时间步 t 中, 将先前生成的单词序列作为当前状态, 表示为 $s_t(x_1, x_2, \dots, x_t, \dots, x_n)$, 其中, x_i 为词汇表中所给的

单词, 用当前的单词序列状态 s_t 预测下一时间步单词序列的生成。

近年来, 文本生成技术取得了较大进展。2014年, GOODFELLOW 等^[1]针对图像生成等任务, 提出生成对抗网络(Generation Adversarial Networks, GAN), GAN 由生成器和鉴别器组成, 生成器的作用是模拟真实数据的分布, 鉴别器的作用是判断一个样本是真实样本还是模型模拟生成的样本, GAN 的目标就是训练一个生成器以完美地拟合真实数据分布使得

基金项目: 国家自然科学基金(61762071); 内蒙古自治区自然科学基金(2015MS0621)。

作者简介: 高永兵(1974—), 男, 副教授、硕士, 主研方向为文本挖掘、信息检索; 黎预璇、高军甜, 硕士研究生; 马占飞, 教授、博士。

收稿日期: 2020-11-23 修回日期: 2021-01-22 E-mail: gaoyongbing@163.com

判别器无法区分。自 GAN 被提出以来,基于 GAN 的文本生成引起研究人员的广泛关注,学者们提出多种基于 GAN 的文本生成方法。GOU 等^[2]提出 LeakGAN,其允许鉴别器将自己的高层特征泄漏到生成器中,指导生成器的预测生成,从而解决长文本的生成不连贯问题。ZHU 等^[3]提出 Texygen,其作为一种开放域文本生成的标杆平台,囊括大多数文本生成模型,如 SeqGAN^[4]、MaliGAN^[5]等。

除 GAN 之外,Sequence2Sequence^[6]技术和编解码框架也在自然语言生成任务中得到广泛应用。编码器将输入的序列编码为隐藏状态,经过特征提取作为解码器的输入,解码器预测生成与编码器输入相应的自然语言。

基于上述方法,许多文本生成研究取得了显著成果。LI 等^[7]实现了基于对抗条件变分自编码器的中文诗歌生成,其使用 LSTM 训练的微软小冰参加诗词创作比赛,并通过了图灵测试。YAO 等^[8]根据主题、静态和动态情节线实现故事的自动生成,生成的故事连贯、多样且符合主题。ZENG 等^[9]根据微博用户的配置文件、个人描述、历史微博评论,自动生成个性化的微博评论。

文本生成技术在以上任务中取得较好效果,但是上述都是针对特定域的文本生成,在进行其他域的文本生成时需要重新训练模型,相当耗费资源。微博文本复杂多样,所训练模型的健壮性低,难以学习微博文本的语言风格,但是,经过大量规范文本训练的预训练语言模型能直接微调(fine-tuning)处理下游任务,避免从零开始训练模型,从而降低了训练代价。因此,针对微博文本自动生成问题,本文采用预训练语言模型加微调的方法,在给定微博主题的条件下,将微博文本中蕴含的用户意图和用户对话功能作为其个性化特征,微调预训练语言模型生成符合用户预期的微博文本,并整体生成“@用户”,实现用户对话交流功能。具体地,本文从语义和语言风格 2 个角度挖掘微博文本的个性化特征,提出一种基于主题和用户意图的微博文本控制生成技术。微调预训练模型 GPT2-Chinese^[10],实现联合主题和用户意图的微博文本个性化生成。在此基础上,对“@用户”进行特殊处理,自动预测用户交流的对象,在营销推荐意图上实现“@用户”的对话功能。

1 相关工作

近年来,预训练语言模型引起研究人员的广泛关

注,语言模型在高质量、大规模的数据集上预先训练,学习理解自然语言的表示,研究人员可直接使用预训练语言模型微调处理下游任务,省去了繁琐的模型训练过程,从而促进了各种自然语言处理技术的快速发展。

预训练语言模型主要学习词的上下文表示,根据不同的序列预测方式可分为自编码和自回归 2 种语言模型。谷歌提出的 BERT^[11]是典型的自编码语言模型,其使用 Transformer 抽取特征,引入 MLM (Masked Language Model) 和 NSP (Next Sentence Prediction) 预训练目标,能够获取上下文相关的双向特征表示,从而处理句子或段落的匹配任务,但是,该模型预训练过程和生成过程的不一致导致其在生成任务上效果不佳。自回归语言模型的典型代表有 ELMo^[12]、GPT^[13-15]、XLNet^[16]。ELMo 是最早的预训练语言模型,其使用双向长短时记忆 (BiLSTM) 网络串行提取特征,模型按照文本序列顺序拆解的方式在从左至右和从右至左 2 个方向学习词的深度上下文表示,从而获取上下文信息的双向特征,但是,ELMo 本质上是 2 个单向语言模型的拼接,不能同时获取上下文表示,且神经网络 LSTM 不能解决长距离依赖问题,特征提取能力弱。GPT 使用 Transformer 进行特征抽取,能快速捕捉更长范围的信息,目前已经更新到第三代:GPT1 微调阶段引入语言模型辅助目标,解决了微调过程中的灾难性遗忘问题;GPT2 在 GPT1 的基础上进行改进,使用覆盖更广、质量更高的训练数据,认为预训练中已包含很多特定任务所需的信息,其没有针对特定模型的精调流程,在生成任务上取得了很好的效果;GPT3 使用比 GPT2 更多的训练数据和性能更高的计算资源以提高模型性能。XLNet 模型使用 Transformer-XL 抽取特征,该模型针对 BERT 预训练过程与微调不一致的缺点,引入 PLM (Permuted Language Modeling),能学到各种双向上下文表示,分析结果表明,XLNet 在 20 个任务上的性能表现优于 BERT,且性能都有大幅提升。

综上,虽然 XLNet 语言模型学习上下文表示的能力强于其他预训练语言模型,但是单向的自回归语言模型更适合生成任务。微博文本生成任务的训练数据是中文且在生成任务中预训练与微调应具有一致性,因此,本文选择在中文版的预训练模型 GPT2-Chinese 下研究微博文本生成任务。虽然 GPT3 已经问世,但是其模型要求更高的计算资源,实验布置难度较高,因此,本文选择 GPT2 进行研究。

2 用户意图

当今社会快速发展,人们的生活节奏不断加快,信息流动的速度越来越快,碎片化的微博内容比长篇文章更适合阅读。因此,自动生成的微博文本能准确表达用户意图至关重要。一方面,用户意图反映了用户发表博文最直接的社交需求,比如,普通用户发表心情感悟,分享日常生活,企业公司发表博文营销宣传产品和活动,媒体工作者借助微博平台传播新闻事件,呼吁广大群众的关注,文艺工作者借助博文传播知识,实现与其他用户的经验共享、共同进步;另一方面,社交意图蕴含于微博文本中,是更深层次的语义特征,能在一定程度上表现微博文本的个性化特征。因此,具有清晰明确的用户意图的微博文本能促进用户间信息的有效获取和交流,帮助用户在微博平台上实现社交的目的。本文挖掘和提取微博文本中的用户意图类别,在给定微博主题的条件下按照用户意图自动生成微博文本。

2.1 用户意图示例

用户意图蕴含于微博文本中,能从词和句子中挖掘。华为终端官方微博用户发的一条博文如下:
#华为 Mate40#系列新品发布盛典将于10月30日14:30正式开启!锁定@华为终端官方微博,带你直通发布会现场。
结合“华为”“新品”“锁定”“发布会”等词和整个句子的语义,可以推断出这是华为官博为推出“华为Mate40”发表的博文,旨在营销宣传。用同样的方法对大量的微博文本进行统计分析,可将微博文本中的用户意图大致分为营销推荐、新闻评论、知识传播、心情感悟、日常分享这5个类别。从上述例子来看,基于用户意图的微博文本生成过程如表1所示。

表1 微博文本生成过程示例		
Table 1 Example of the Weibo text generation process		
用户意图	主题	生成博文
营销推荐	华为 Mate40	系列新品发布盛典将于10月30日14:30正式开启!锁定@华为终端官方微博,带你直通发布会现场。

在表1中,给定主题“华为 Mate40”和用户意图“营销推荐”,经过训练的语言模型自动预测与主题、用户意图类别相符合的词,最终构成一个完整的文本段落。

2.2 用户意图识别与分类

为了便于数据处理和辅助 GPT2-Chinese 生成更加符合用户预期意图的微博文本,本文额外训练一个用户意图识别与分类模型。

目前,研究人员已经在消费意图、查询意图、人机对话意图识别中取得了较多成果,但多数是基于传统机器学习的方法^[17-19]。传统机器学习模型在特征工程中需要人为对数据进行提炼清洗,与深度学习模型相比有很大不足。为此,研究人员更加专注在神经网络上进行意图识别^[20-22]并取得了较好的效果。

本文为了学习句子、段落以及关键词的语义编码,从而准确识别微博文本中的用户意图,训练一个用户意图识别与分类模型,该模型采用编解码器框架,融合 BiLSTM、自注意和词句联合训练等要素进行用户意图的识别与分类,模型框架^[23]如图1所示。

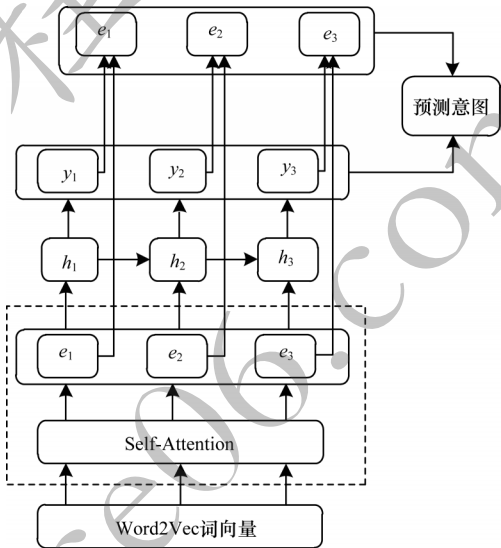


图1 用户意图识别与分类模型框架
Fig.1 User intention recognition and classification model framework

3 本文方法

本文方法的主要目的是使自动生成的样本在内容上与微博主题保持一致且从用户意图的角度展示出微博文本的个性化特征。针对微博文本生成问题,本文作出如下定义:将用户意图定义为特征向量 U_i ,微博主题定义为 T ,模型预测生成的微博文本定义为 X ,生成模型可视作条件概率模型 $P(X|U_i, T)$,最后将微博文本作为训练数据,结合微博主题和用户意图微调 GPT2-Chinese 语言模型,预测生成微博文本。

由于存在多个控制文本生成的用户意图,统一训练可能导致模型学习能力降低,为此,本文使用 MAO 等^[24]提出的两阶段微调方式:第一阶段,微调预训练语言模型 GPT2-Chinese 学习微博文本的语言风格;第二阶段联合主题和用户意图进行微调,使生成的微博文本流畅且符合用户发博的目的。图2所示为微博文本生成任务的 GPT2 两阶段微调流程。

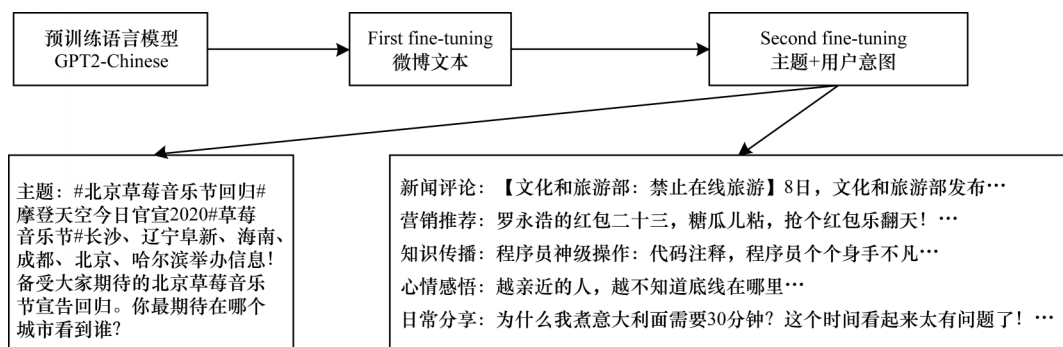


图2 GPT2两阶段微调流程

Fig.2 GPT2 two-stage fine-tuning procedure

3.1 文本生成语言建模

在实际应用中,自然语言模型是学习序列单词的概率分布 $p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n)$,文本生成可以形式化为在给定单词序列的情况下预测下一个单词的条件概率 $p(\mathbf{x}_n | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1})$,其中 \mathbf{x}_i 为单词序列。自回归语言模型按照序列预测的条件概率形式可进行单向或双向训练。单向语言模型按照前向序列预测的训练定义为:

$$p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N) = \prod_{k=1}^N p(\mathbf{x}_k | \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1}) \quad (1)$$

预训练语言模型GPT2采用与上述相同的训练方法,其过程是给定提示(Prompt)序列 $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{k-1})$ 与生成文本序列 $(\mathbf{x}_k, \mathbf{x}_{k+1}, \dots, \mathbf{x}_n)$,共同构成一段连贯的文本,具体形式如下:

$$P(\mathbf{X}) = \prod_{t=1}^n p(\mathbf{x}_{k:n} | \mathbf{x}_{1:k-1}) \quad (2)$$

本文在GPT2下完成微调微博文本生成任务,可将微博主题 T 作为提示条件,式(2)可更新为:

$$P(\mathbf{X}) = \prod_{t=1}^n p(\mathbf{x}_{k:n} | \mathbf{t}_{1:k-1}), \mathbf{t}_i \in T \quad (3)$$

为了在生成的微博文本中表现用户意图(U_i),可将用户意图与GPT2预训练语言模型进行简单地模型融合,融合方式如下:

$$P(\mathbf{X}) = \prod_{t=1}^n p(\mathbf{x}_{k:n} | \{\mathbf{U}_i \oplus \mathbf{t}_{1:k-1}\}) \quad (4)$$

3.2 GPT框架

GPT1使用半监督的方式学习理解自然语言,利用大量的无监督语料库训练语言模型,然后模型以很小的微调迁移到众多特定的有监督学习任务上^[13],其训练的最大似然目标函数如下:

$$L_1(\mathbf{u}) = \sum_i \log_a p(\mathbf{u}_i | \mathbf{u}_{1-k}, \dots, \mathbf{u}_{i-1}; \theta) \quad (5)$$

GPT1语言模型不采用RNN或LSTM,而是采用深层的transformer decoder(Masked),其计算公式如下:

$$\begin{aligned} \mathbf{h}_0 &= \mathbf{U}\mathbf{W}_e + \mathbf{W}_p \\ \mathbf{h}_i &= \text{transformer_block}(\mathbf{h}_{i-1}), \forall i \in [1, n] \\ P(\mathbf{u}) &= \text{softmax}(\mathbf{h}_n \mathbf{W}_e^T) \end{aligned} \quad (6)$$

无监督的语言模型经过式(6)的预训练后,将该模型参数作为初始参数应用到一些特定的有监督任

务上。例如,对于有标签的数据集 C ,利用网络中最后一个transformer_block的输出 \mathbf{h}_i^m 进行如下处理:

$$P(\mathbf{y} | \mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m) = \text{softmax}(\mathbf{h}_i^m \mathbf{W}_y) \quad (7)$$

此时最大化目标函数如下:

$$L_2(C) = \sum_{(\mathbf{x}, \mathbf{y})} \log_a P(\mathbf{y} | \mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^m) \quad (8)$$

将语言建模作为微调的辅助目标,改进监督模型的泛化能力,加速收敛。优化后的目标训练函数如下:

$$L_3(C) = L_2(C) + \lambda L_1(C) \quad (9)$$

GPT2是在GPT1基础上做的改进,是完全无监督学习,其在模型结构上移动了Layer normalization的位置,修改了残差层的初始化方式,增加了模型的部分参数,从而提高了生成文本的质量。

3.3 微博文本生成训练

GPT2预训练语言模型允许对下游微博生成任务采取微调的方法,即不需要从零开始训练模型,只需针对特定的任务调整训练参数。目前已开源了许多版本的GPT2模型,本文选择中文的GPT2-Chinese作为预训练语言模型。GPT2-Chinese采用的训练集是15 GB的Large Scale Chinese Corpus for NLP和THUCNews,能够自动生成诗歌、小说、新闻等内容,且生成的内容连贯。因此,在微博文本生成任务中,本文选择用GPT2-Chinese模型来实现两阶段微调任务:第一阶段,用一个高质量微博训练集进行域适应训练,即微调GPT2-Chinese让模型学会从通用自然语言预测转向符合微博文本语言风格的预测生成,得到中间微调的模型;第二阶段,用微调好的GPT2-Chinese模型在带有主题标签和用户意图标签的微博数据集上进行二次微调,最终实现微博文本的个性化生成。

自动生成微博比自动生成诗歌、自动讲故事等生成任务更具挑战性。诗歌有具体形式,且每种形式都有各自的规律,故事有明确的主题和便于规划的故事线,机器更容易学习这种有明确特征的文本。然而微博文本长短不一、书写风格参差不齐,既不似诗歌字数工整、体裁清晰,也不似故事主题明确、上下文连贯。官方公布的GPT2模型更适用于连续的长文本训练,而GPT2-Chinese在训练时设置文本最低长度为128,模型

在自动预测微博文本时表现不佳,上下文的相关性低,无法降低训练的损失率,简单地修改训练文本长度也无法提高生成效果。因此,在训练时本文用微博链接来增加文本长度从而降低损失,生成结果表明,该方法能够提高模型的文本生成性能。

3.4 个性化微博文本生成训练

微调 GPT2-Chinese 语言模型能生成与所给主题相符的样本,且随着实验数据的增多,语言模型学习微博文本特征的能力增强,但是没有明确的控制条件,生成的博文个性化特征不明确。因此,本文从语义特征和语言风格 2 个角度提取微博文本的个性化特征,用其指导博文的预测生成。

通过对微博文本深层次的语义特征进行分析,发现博文中蕴含了微博用户的各种社交意图,且社交意图大致分为新闻评论、营销推荐、知识传播、心情感悟、日常分享这 5 个类别,本文将作为微博文本语义层次上的个性化特征,指导博文生成,且在该过程中采用 2 个方案进行微调:

1)在第二阶段微调时,用 GPT2-Chinese 联合主题和用户意图进行多任务微调,实验中发现预先训练的 GPT2-Chinese 有很强的基础性知识,再次按照用户意图训练能够产生高质量的个性化微博文本,且多个用户意图微调能帮助实现微博文本的个性化预测生成。

2)将本文提出的用户意图识别与分类模型和 GPT2 相结合,辅助 GPT2 生成博文。该方法能提高预测文本中预期的用户意图的准确度。

从语言风格角度可以发现微博文本中会出现大量用户对话的情况,即用户在发表博文时会@另一个用户,从而实现用户之间的交流。本文将用户对话功能看作微博文本的另一个个性化特征,在微调时自动预测用户对话的对象,从而帮助用户实现自动信息交流的目的。

3.4.1 多意图控制生成

为实现个性化微博生成,本文实验中引入 5 个类别的用户意图以控制微博文本的生成。虽然大型预训练语言模型的诞生提高了处理自然语言任务的能力,但直接在特定条件的数据上微调往往很困难。在 CTRL^[25]模型中,将控制条件作为标签加在文本前面,这样在模型训练过程中 attention^[26]会计算控制条件与序列之间的联系,从而实现控制生成。ZACHARY 等^[27]观察到预先训练的 transformer 模型在微调过程中对模型参数的变化很敏感,因此,提出一种直接向 self-attention 注入任意条件的适应方法 pseudo self attention,使用该方法的预训练语言模型适应于任意条件的输入。为了不改变预训练语言模型的结构,DATHATHRI 等^[28]提出用于可控语言生成的即插即用语言模型(PPLM),其将预先训练的 LM 与一个或多个简单的属性分类器相结合以指导文本生成,而不需要对 LM 进行额外的训练。本文在意图控制生成时学习 PPLM 的思想,将用户意图识别与分类模型作为意图分类器,意图分类器判别生成的文本类别并用梯度回传的方式传递给语言模型,预训练语言模型根据意图判别回传的梯度更新模型的内部参数,重新采样生成下一个 token,使生

成的样本接近用户输入的意图。上述过程可以形式化表示为:

$$P(X|U_i) \propto P(U_i|X)P(X) \tag{10}$$

此时用户意图分类模型可看作 $P(U_i|X)$,语言模型 GPT2 中具有历史信息的矩阵 H_i ,再给定 x_i ,GPT2 利用 H_i 预测 x_{i+1} 和更新 H_i :

$$(x_{i+1}, H_{i+1}) = \text{GPT2}(x_i, H_i) \tag{11}$$

其中 x_{i+1} 是经过 softmax 在词汇表中采样生成的文本。PPLM 的思想是从 2 个方向计算分类模型和语言模型的梯度之和,然后通过 ΔH_i 更新历史矩阵 H_i ,实现 x_{i+1} 的重新预测。在用户意图的控制上,本文将意图分类模型 $P(U_i|X)$ 改写为 $P(U_i|H_i + \Delta H_i)$, ΔH_i 的更新方式如下:

$$\Delta H_i \leftarrow \Delta H_i + \alpha \frac{\nabla \Delta H_i \log_a P(U_i|H_i + \Delta H_i)}{\|\nabla \Delta H_i \log_a P(U_i|H_i + \Delta H_i)\|} \tag{12}$$

其中: α 是步长; γ 为归一化缩放系数。式(11)就可以更新为:

$$(x_{i+1}, H_{i+1}) = \text{GPT2}(x_i, H_i + \Delta H_i) \tag{13}$$

3.4.2 “@用户”的对话实现

通过对微博文本进行分析可知,“@用户”的功能也是微博文本的特征,因此,本文还研究自动@正确用户名的方法。在生成过程中,本文发现模型能自动生成“@用户名”,但在训练时用户名被当作文本进行了分词处理,因此,模型生成的用户名往往不存在,且由于用户名是一个实体,经过分词后影响了文本的上下文相关性,降低了生成文本的质量。为此,本文将“@用户名”作为一个词并用 UN_i 替换,然后为其建立词表,当生成时出现 UN_i 模型会自动将其转换为“@用户”,其中, i 表示词表中的第 i 个词。最后,本文在数据中发现,在营销推荐意图上使用“@用户”的对话概率更大,因此,将营销推荐数据中的“@用户”经过特殊处理,从而实现能准确“@用户”的对话功能。

4 实验结果与分析

4.1 数据集与实验环境

本次实验数据集是通过新浪微博接口获取的 27 万条微博文本,本文将数据集根据用户意图和微博主题预处理为适合微调的格式。首先标注 7 000 条用户意图数据,进行用户意图识别与分类模型训练;然后使该语言模型自动进行用户意图标注,人工标注和自动标注的标签数据集如表 2 所示;最后在数据集上批量进行微调微博生成实验,观察不同训练样本数量时的微博文本生成效果。

表 2 标签数据集信息

Table 2 Label datasets information

用户意图分类	人工标注数量	自动标注数量
新闻评论	1 548	68 825
营销推荐	1 270	88 753
知识传播	1 441	46 871
心情感悟	1 482	36 333
日常分享	1 259	29 218

实验环境是使用谷歌 Colab 免费的 GPU,虽然谷歌能提供免费的计算环境,但是一天挂载时间只有 12 h,因此,在模型训练过程中本文加入 checkpoint,如果训练中断,还可以接着上次中断位置继续训练。

4.2 用户意图识别与分类实验

用户意图识别与分类实验使用表 2 中的人工标注数据集,用词和句子的意图对文本意图进行打分训练,实验过程中分别用准确率和 F1 值进行评价,实验结果如表 3 所示。其中,词句+BiLSTM 是本文提出的用户意图识别与分类模型,即从词、句 2 个方面关注用户意图,并将 BiLSTM 融入编解码框架中训练用户意图识别与分类模型。从表 3 可以看出,本文模型能获得 93.109% 的 F1 值,优于 BERT 模型。

表 3 用户意图识别与分类实验结果

Table 3 Experimental results of user intention recognition and classification %

模型	准确率	F1 值
BERT	81.673	81.211
词句+BiLSTM	93.176	93.109

4.3 微调微博文本生成任务

在模型训练开始前,本文从 GitHub 上下载了中文的 GPT2 训练模型 GPT2-Chinese,用于微博文本生成任务。首先用无标签的数据训练模型将其转换为

微博文本语言的预测生成,得到中间微调的语言模型;然后用表 2 中的自动标签数据微调生成符合微博主题和用户意图的博文,在训练时,GPT2-Chinese 模型参数量为 81 894 144,为提高 GPU 的利用率,设置 batchsize 为 4,epoch 为 5,每个 epoch 有 10 212 步,设置微博文本的生成温度为 0.8,对下一个单词的预测采用 ANGELA 等^[29]提出的 Top- k 随机抽样方法,其中,下一个单词从 k 个候选词中抽取,设定 $k=10$;接着联合主题标签和用户意图标签,利用中间微调的模型进行二阶段微调学习,从用户意图角度控制微博文本的个性化特征,且在该过程中加入用户意图分类模型,以提高生成样本中指定用户意图的准确性。在此基础上,将“@用户名”进行整体替换,实现用户对话的生成预测。实验结果表明,本文方法生成的微博文本可读性高,且能从用户意图的角度表现微博文本的个性化特征。

4.4 生成样例

本文列出基于主题、基于主题和用户意图 2 种方式预测生成的微博文本样例,分别如图 3、图 4 所示,在此仅展示营销推荐和新闻评论 2 个用户意图生成样例。在营销推荐意图上实现“@用户”的对话功能,生成样例如图 5 所示。

主题: 分享红包
1.分享红包,羊年到,好运先来! 你也来试试手气! <http://t.cn/rzgxntz>
2.分享红包,羊年大吉! 你也来试试手气! <http://t.cn/rzgxntz>
3.分享红包@小年抢个好彩头! 我在小年大家也来试试手气! <http://t.cn/rzdfp5h>

图 3 基于主题的生成样例

Fig.3 Generation sample based on topic

主题: 快的打车; 用户意图: 营销推荐
1.快的打车6元红包,快来试试手气,你与幸福之间只有一个红包的距离! <http://t.cn/rzg7ixl>
2.快的打车代金券#不愁天空发脾气,「快的小财神」为您省送“超大代金券礼包”! 福气有你一份,速来认领: <http://t.cn/rwvf4jt>
3.快的打车“寻找miss queen”活动二等您一份,速来认领: <http://t.cn/rw2j5wf>
4.快的打车,绿色同行! 无论风雨,与你同行! 无论论坛也可以与幸福气打车券: <http://t.cn/rwx1ej>
主题: 商丘卫健委通报早产儿死而复生事件; 用户意图: 新闻评论
1.【#商丘卫健委通报早产儿死而复生事件#: 依法依规追责】针对媒体报道河南商丘“一婴儿活着却被医生说是死婴,最后导致重度脑瘫”一事,河南商丘卫健委发布情况说明: 下一步,将组成调查组对相关事件进行认真调查核实。
2.#商丘卫健委通报早产儿死而复生事件#我记得民生频道@小莉帮忙也发过这个视频,怎么就找不到了? 难不成因为什么事情又把原视频删了?
3.#商丘卫健委通报早产儿死而复生事件#对于本已经弱不禁风的医患关系真是雪上加霜//@人民日报

图 4 基于主题与用户意图的生成样例

Fig.4 Generation sample based on topic and user intention

主题: 王俊凯发光同款; 用户意图: 营销推荐
1.王俊凯发光同款的OLAY身体精华乳,买还送小凯的城市精美明信片呢@OLAY身体护理
2.王俊凯发光同款,get帅气代言人@TFBOYS-王俊凯同款一起发光
3.王俊凯发光同款,闪闪发光!!! 就是王俊凯本凯!!! @TFBOYS-王俊凯
4.王俊凯发光同款发光瓶,清醒自律@TFBOYS-王俊凯就是最善良的心。

图 5 实现“@用户”对话功能的生成样例

Fig.5 Generation sample with “@user” dialog function

5 人工评价

与传统文本不同,微博文本生成属于特型文本生成任务,传统的文本评价指标 BLEU 和 ROUGE 不适用于对微博文本生成进行评价。本文通过人工方式评价生成样本的质量。人类评审者需要考虑样本与微博主题、用户意图的匹配度,然后对需要评价的样本进行打分,实验过程中共有 5 位评审,随机采样 600 条生成样本,其中,100 条是只根据主题生成的,500 条是根据主题和用户意图生成的,每类用户意图各 100 条。将每种类型的样本平均分配给 5 个评审,让其对生成的样本进行打分,分数为 0~10 分,最后计算每类样本得分的平均值,结果如表 4 所示,可以看出,基于主题与用户意图生成的微博文本更符合人类预期。

表 4 人工评价得分对比

Table 4 Manual evaluation scores comparison

样本	得分
基于主题的样本	77.0
基于主题+新闻评论的样本	83.1
基于主题+营销推荐的样本	80.0
基于主题+知识传播的样本	78.0
基于主题+心情感悟的样本	79.0
基于主题+日常分享的样本	77.8

6 结束语

本文提出一种基于用户意图的微博文本生成技术,采用预训练语言模型与微调相结合的方法,在给定微博主题的条件下,将微博文本中蕴含的用户意图和用户对话功能作为其个性化特征,微调预训练语言模型生成符合用户预期的微博文本。实验结果验证了该技术的有效性。微博文本生成技术能够为博文编写者提供参考,也可以为后续的文本生成研究提供思路,但是,本文模型仅从统计的角度并利用静态的方式挖掘微博文本的个性化特征,下一步将挖掘兴趣品味、行文风格等更多的用户个性化特征,此外,探索一种性能更好、联合多个控制条件微调预训练语言模型的方法,以提高训练技术同时节省时间和资源成本,也是今后的研究方向。

参考文献

[1] GOODFELLOW I J, JEAN P A, MEHDI M, et al. Generative adversarial nets[C]//Proceedings of the 27th International Conference on Neural Information Processing Systems. Montreal, Canada; NIPS, 2014: 2672-2680.

[2] GUO J X, LU S D, HAN C, et al. Long text generation via adversarial training with leaked information[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1709.08624.pdf>.

[3] ZHU Y, LU S, LEI Z, et al. Taxygen: a benchmarking

platform for text generation models[C]//Proceedings of the 41st International ACM SIGIR Conference on Research & Development in Information Retrieval. New York, USA: ACM Press, 2018: 1097-1100.

[4] YU L, ZHANG W, WANG J, et al. SeqGAN: sequence generative adversarial nets with policy gradient[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1609.05473.pdf>.

[5] CHE T, LI Y, ZHANG R, et al. Maximum-likelihood augmented discrete generative adversarial networks[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1702.07983.pdf>.

[6] SUTSKEVER H, VINYALS ORIOL, QUOC V L, et al. Sequence to sequence learning with neural networks[EB/OL]. [2020-10-02]. <http://cs224d.stanford.edu/papers/seq2seq.pdf>.

[7] LI J, SONG Y, ZHANG H, et al. Generating classical Chinese poems via conditional variational autoencoder and adversarial training[C]//Proceedings of 2018 Conference on Empirical Methods in Natural Language Processing. Washington D. C., USA: IEEE Press, 2018: 3890-3900.

[8] YAO L, PENG N, WEISCHEDEL R, et al. Plan-and-write: towards better automatic storytelling[J]. Proceedings of the AAAI Conference on Artificial Intelligence, 2019, 33: 7378-7385.

[9] ZENG W, ABUDUWEILI A, LI L, et al. Automatic generation of personalized comment based on user profile[C]//Proceedings of the 57th Conference of the Association for Computational Linguistics. Florence, Italy: Association for Computational Linguistics, 2019: 229-235.

[10] DU Z Y. GPT2-Chinese: tools for training GPT2 model in Chinese language[EB/OL]. [2020-10-02]. <https://github.com/Morizeyao/GPT2-Chinese>.

[11] DEVLIN J, CHANG M W, LEE K, et al. BERT: pre-training of deep bidirectional transformers for language understanding[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1810.04805.pdf>.

[12] PETERS M, NEUMANN M, IYYER M, et al. Deep contextualized word representations[C]//Proceedings of 2018 Conference of the North American Chapter of the Association for Computational Linguistics; Human Language Technologies. Florence, Italy: Association for Computational Linguistics, 2018: 2227-2237.

[13] RADFORD A. Improving language understanding by generative pre-training[EB/OL]. [2020-10-02]. https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf.

[14] RADFORD A, JEFFREY W, REWON C, et al. Language models are unsupervised multitask learners[EB/OL]. [2020-10-02]. https://cdn.openai.com/better-language-models/language_models_are_unsupervised_multitask_learners.pdf.

[15] TOM B, BROWN B M, NICK R, et al. Language models are few-shot learners[EB/OL]. [2020-10-02]. <https://papers.nips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf>.

[16] YANG Z L, DAI Z H, YANG Y M, et al. XLNet: generalized autoregressive pretraining for language understanding[EB/OL]. [2020-10-02]. <https://www.cs.cmu.edu/~jgc/publication/>

- XLNET. pdf.
- [17] 付博,陈毅恒,邵艳秋,等. 基于用户自然标注的微博文本的消费意图识别[J]. 中文信息学报,2017,31(4):208-215.
- FU B, CHEN Y H, SHAO Y Q, et al. Consumption intent recognition based on user natural annotation[J]. Journal of Chinese Information Processing, 2017, 31(4): 208-215. (in Chinese)
- [18] 贾云龙,韩东红,林海原,等. 面向微博用户的消费意图识别算法[J]. 北京大学学报(自然科学版),2020,56(1): 68-74.
- JIA Y L, HAN D H, LIN H Y, et al. Consumption intent recognition algorithms for Weibo users[J]. Acta Scientiarum Naturalium Universitatis Pekinensis, 2020, 56(1): 68-74. (in Chinese)
- [19] 桂思思,陆伟,张晓娟. 基于查询表达式特征的时态意图识别研究[J]. 数据分析与知识发现,2019,3(3):66-75.
- GUI S S, LU W, ZHANG X J. Temporal intent classification with query expression feature [J]. Data Analysis and Knowledge Discovery, 2019, 3(3): 66-75. (in Chinese)
- [20] 廖胜兰,吉建民,俞畅,等. 基于BERT模型与知识蒸馏的意图分类方法[J]. 计算机工程,2021,47(5):73-79.
- LIAO S L, JI J M, YU C, et al. Intention classification method based on BERT model and knowledge distillation [J]. Computer Engineering, 2021, 47(5): 73-79. (in Chinese)
- [21] CHEN Q, ZHUO Z, WANG W, et al. BERT for joint intent classification and slot filling[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1902.10909.pdf>.
- [22] QIN L, CHE W, LI Y, et al. A stack-propagation framework with token-level intent detection for spoken language understanding[EB/OL]. [2020-10-02]. <https://aclanthology.org/D19-1214.pdf>.
- [23] 高永兵,李越超. 微博中的社交意图识别与分类技术研究[J]. 内蒙古科技大学学报,2020,39(2):85-89.
- GAO Y B, LI Y C. Social intent recognition and classification in Weibo [J]. Journal of Inner Mongolia University of Science and Technology, 2020, 39(2): 85-89. (in Chinese)
- [24] MAO H H, MAJUMDER B P, MCAULEY J, et al. Improving neural story generation by targeted common sense grounding[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1908.09451v1.pdf>.
- [25] KESKAR N, MCCANN B, VARSHNEY L R, et al. CTRL: a conditional transformer language model for controllable generation[EB/OL]. [2020-10-02]. <https://einstein.ai/presentations/ctrl.pdf>.
- [26] VASWANI A, NOAM S, NIKI P, et al. Attention is all you need[EB/OL]. [2020-10-02]. <https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf>.
- [27] ZACHARY M Z, LUKE M K, SEBASTIAN G, et al. Encoder-agnostic adaptation for conditional language generation[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1908.06938.pdf>.
- [28] DATHATHRI S, ANDREA M, JANICE L, et al. Plug and play language models: a simple approach to controlled text generation[EB/OL]. [2020-10-02]. <https://arxiv.org/pdf/1912.02164.pdf>.
- [29] ANGELA F, MIKE L, YANN D. Hierarchical neural story generation[EB/OL]. [2020-10-02]. <https://aclanthology.org/P18-1082.pdf>.

编辑 吴云芳