

基于预训练模型的受控文本生成研究综述

周强伟¹, 施水才^{1,2}, 王洪俊²

(1. 北京信息科技大学计算机学院, 北京 100101; 2. 拓尔思信息技术股份有限公司, 北京 100096)

摘要: 自然语言生成(NLG)作为人工智能的一个分支,近年来随着预训练语言模型(PLMs)的发展取得了显著进展。NLG旨在根据多种输入源(如文本、图像、表格和知识库)生成连贯、有意义的文本。研究者通过架构扩展、微调 and 提示学习等方法提升了PLMs的性能。然而,NLG在处理非结构化输入和低资源语言生成方面仍面临挑战,尤其是在缺乏足够训练数据的环境中。为探讨NLG的最新发展、应用前景以及所面临的挑战,通过文献分析,提出PLMs性能改进策略,并展望未来研究方向。研究表明,尽管存在诸多限制,但NLG在内容创作、自动新闻报导、对话系统等领域已展现出潜力。随着技术的不断进步,NLG在自然语言处理和人工智能领域将扮演更重要的角色。

关键词: 人工智能; 自然语言生成; 受控文本生成; 预训练语言模型; 提示学习

DOI: 10.11907/rjdk.231393

开放科学(资源服务)标识码(OSID):

中图分类号: TP391.1

文献标识码: A

文章编号: 1672-7800(2024)004-0199-09



Overview of Controlled Text Generation Based on Pre-trained Models

ZHOU Qiangwei¹, SHI Shuicai^{1,2}, WANG Hongjun²

(1. School of Computer, Beijing Information Science and Technology University, Beijing 100101, China;
2. TRS Information Technology Co Ltd, Beijing 100096, China)

Abstract: Natural language generation (NLG), a branch of artificial intelligence, has seen significant progress in recent years, particularly with the development of Pre-trained language models (PLMs). NLG aims to generate coherent and meaningful text based on various input sources such as texts, images, tables, and knowledge bases. Researchers have enhanced the performance of PLMs through methods like architectural expansion, fine-tuning, and prompt learning. However, NLG still faces challenges in dealing with unstructured inputs and generating text in low-resource languages, especially in environments lacking sufficient training data. This study explores the latest developments in NLG, its application prospects, and the challenges it faces. By analyzing existing literature, we propose strategies to improve the performance of PLMs and anticipate future research directions. Our findings indicate that despite limitations, NLG has shown potential in areas such as content creation, automated news reporting, and conversational systems. The conclusion is that, with technological advancements, NLG will play an increasingly significant role in natural language processing and other related fields of artificial intelligence.

Key Words: artificial intelligence; natural language generation; controlled text generation; pre-trained language models; prompt learning

0 引言

近年来,人工智能领域发展迅速,其中一个重要分支就是自然语言处理(Natural Language Processing, NLP)。自然语言处理是指利用人类交流所使用的自然语言与机器进行交互通讯的技术,要实现人机自然语言交互意味着计算机不仅能理解自然语言的意义,也能够以自然语言传递信息、表达情感。上述提到的两种功能,前者称为自然语言理解(Natural Language Understanding, NLU),后者称为

自然语言生成(Natural Language Generation, NLG),也称为文本生成,旨在从包括文本、图像、表格和知识库等各种形式的输入数据中生成可信和可读的人类语言文本。受控文本生成是指在生成文本的过程中,加入一些额外的信息以约束生成的文本符合特定条件。这些条件可以包括输入的语义信息、领域特定知识、结构化数据等。在过去的几十年中,受控文本生成技术已广泛应用于各种应用领域。例如,它们已用于对话系统中,以生成对话中用户话语的响应^[1],在机器翻译中将一种语言的文本翻译成另一种语言,以及在文本摘要中生成源文本的简短摘要^[2]。

收稿日期: 2023-05-24

作者简介: 周强伟(1998-),男,北京信息科技大学计算机学院硕士研究生,研究方向为自然语言生成;施水才(1966-),男,硕士,北京信息科技大学计算机学院教授、硕士生导师,研究方向为信息检索、大数据分析和挖掘、人工智能。本文通讯作者:施水才。

1 文本生成发展研究现状

1.1 早期文本生成

文本生成的主要目标是自动学习从数据中构建端到端解决方案的输入输出映射,以最小化人类干预。这种映射函数允许生成系统在更广泛的领域中进行泛化,并在给定条件下生成自由文本。早期方法通常采用统计语言模型对给定 n -gram 上下文情况下单词的条件概率进行建模^[3-4]。这种统计方法已知存在数据稀疏问题,因而许多平滑方法被提出以缓解该问题,并更好地估计未观察到的术语出现次数。

然而,在这些方法中使用单词令牌作为基本表示单位会导致相似令牌之间难以轻松映射的问题。随着深度学习技术的出现,神经网络模型已经成为文本生成主流方法并在生成自然语言文本方面取得了显著成效。深度神经网络生成模型通常采用基于编码器—解码器框架的序列到序列架构:编码器首先将输入序列映射到固定大小低维向量(称为输入嵌入),然后解码器根据输入嵌入生成目标文本。嵌入表示与早期统计方法的区别在于,它使得处理输入和输出之间可能存在的关系更加容易。在深度学习领域,各种神经模型被提出,不同的编码器—解码器架构设计被采用,例如用于编码图形输入的图形神经网络(GNN)^[5]和用于解码文本的递归神经网络(RNN)^[6]。此外,注意力机制和复制机制已广泛用于改善文本生成模型性能^[7]。深度神经网络对于文本生成的一个重要优点是它们可以实现从输入数据到输出文本语义映射端到端学习而无需费力进行特征工程。此外,深度神经模型通过低维语义表示捕获语言的基本特征,这有助于缓解数据稀疏性问题。

尽管深度神经模型在文本生成方面取得了成功,但是在大规模标记数据集可用性方面仍然存在主要性能瓶颈,大多数文本生成方法需要大量手动标记平行数据,且其在许多领域中的应用受注释示例匮乏所限制。迄今为止,针对文本生成任务的现有标记数据集通常很小,在这种情况下,深度神经网络很可能会过拟合这些小数据集,并且在

实践中无法很好地进行泛化^[8]。此外,文本生成的早期神经模型仍然相对较浅,只有1~3个神经元,这些模型难以建模上下文和单词含义之间的复杂关系,并推导出更好的上下文词表示以进行更好的生成。

1.2 预训练语言模型文本生成

近年来,预训练语言模型(PLMs)的范式在自然语言处理领域蓬勃发展^[9]。其基本思想是:首先在大规模无监督语料库上对模型进行预训练,然后在下游有监督任务中微调这些模型。这种预训练—微调框架实现了最先进的性能。随着Transformer和更强大计算能力的出现,PLMs的架构已经从浅层演变为更深层次结构,例如BERT^[10]和OpenAI GPT^[11]。大量工作表明,PLMs可以将来自于预训练语料库的大量语言知识编码到它们的大规模参数中,并学习具有特殊设计目标(如掩码令牌预测)的通用和上下文表示形式。因此,PLMs通常对下游任务有益,并且可以避免从头开始训练新模型。随着PLMs在其他NLP任务中取得成功,可将PLM应用于具有多个步骤的文本生成任务中。通过对大规模语料库进行预训练,PLMs可以准确理解自然语言并进一步流畅地表达人类语言,这两种能力对于完成文本生成任务至关重要。基于PLMs的文本生成被视为学术界和工业界颇具前景的研究方向,这极大提高了该领域的技术水平。

1.3 研究现状

目前,已有许多关于文本生成和PLM的调查论文。例如,Qiu等^[9]总结了整个NLP领域中的两代PLMs,并介绍了各种扩展和适应方法。Kalyan等^[12]简要概述了基于Transformer的PLMs中自监督学习的进展。Han等^[13]深入研究了预训练的历史,特别是它与迁移学习和自监督学习之间的特殊关系。

此外,El-Kassas等^[14]专注于将PLMs应用到文本摘要领域。Zaib等^[15]讨论了将PLMs应用于对话系统,特别强调了问答系统。这些调查重点关注具体应用程序(例如摘要和对话系统),但没有深入探讨核心技术即文本生成方面。由于文本生成是各种应用程序中的关键组成部分,因而提供基于PLMs的受控文本生成调查非常有用。基于PLM的文本生成流程如图1所示。

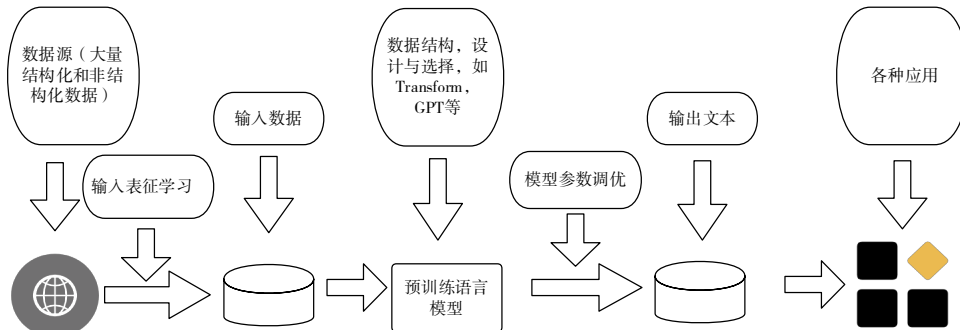


Fig.1 Text generation process based on PLM

图1 基于PLM的文本生成流程

2 基于预训练模型的受控文本生成

2.1 文本生成

通常,文本可以被建模为一个令牌序列 $y = \langle y_1, \dots, y_j, \dots, y_n \rangle$, 其中每个令牌 y_i 都来自词汇表 V 。文本生成任务旨在生成一种可信且易读的人类语言文本。在大多数情况下,文本生成是基于某些输入数据(例如文本、图像、表格和知识库)进行条件化,这些输入数据被表示为 x 。特别地,期望所产生的文本满足一些所需语言属性,如流畅性、自然度和连贯性。将输出文本的所需属性称为属性集 P 。根据上述符号表示法,可用如式(1)所示形式描述文本生成任务。

$$y = f_M(x, P) \quad (1)$$

文本生成模型会根据 PLM 进行特殊设计以产生输出文字 y , 并给定输入数据 x 和满足属性集合 P 的某些特殊要求。具体而言,在不同类型的输入数据 x 和属性集合 P 下可能实例化出不同类型的任务:①当未提供或者提供随机向量作为输入时,则会退化成语言建模或无条件文字生成,此时需要输出文字满足一些普遍的语言属性,如流畅性和自然度;②当输入数据是结构化数据,例如知识库或表格,则被视为是数据到文字生成的任务,该任务旨在生成关于结构化数据的描述性文字,因而输出文本应客观准确;③输入数据最常见的形式是文本序列,这种形式涵盖了许多应用,例如机器翻译、文本摘要和对话系统;对于特定任务,期望输出的文本满足所需属性,例如在文本摘要中,摘要不应与输入文本中描述的事实相矛盾,而在对话中,回复应与输入对话历史和上下文相关。

2.2 预训练语言模型

预训练语言模型 (PLMs)^[9]是在大规模未标记的语料库上预先训练的深度神经网络,可以进一步在各种下游任务中进行微调。研究已经证明,PLMs 可以将大量的语言知识编码到它们庞大的参数中。因此,用 PLMs 增强对语言的理解并提高生成质量具有很大潜力。由于 Transformer^[7]取得了巨大成功,几乎所有 PLMs 都采用它作为骨干。作为两个典型的 PLM, GPT 和 BERT 分别建立在 Transformer 解码器和编码器之上。继 GPT 和 BERT 之后,研究者在文献中又提出 XLNet、RoBERTa、ERNIE、T5 和 BART 等其他 PLMs。其中, XLNet、RoBERTa 和 ERNIE 基于 BERT 模型而开发,而 T5 和 BART 是基于编码器—解码器结构设计。最近研究表明,PLM 性能可以通过增加模型参数规模加以提升,这促使了像 GPT-3(175B)^[14]、PANGU(200B)、GShard(600B)、Switch-Transformers(1.6T)等数十亿或数万亿参数的大规模 PLMs 的发展。此外,PLMs 还被设计用于其他任务,如命名实体识别、编程和网络。根据预训练目标,文本生成的 PLMs 可以分为掩码 LMs、因果 LMs、前缀 LMs 和编码器—解码器 LMs。

2.3 受控文本生成

受控文本生成方法是指在文本生成任务中,通过一定的控制手段实现生成文本的特定要求,如风格、主题、长度、语言等。常见的受控文本生成方法包括以下5种:

(1)控制标记法。在输入文本中加入特定的控制标记,以指示生成文本的特定要求。例如,在生成情感文本时,可以在输入文本中加入"positive"或"negative"等情感标记,以指示生成正面或负面情感的文本。

(2)控制嵌入向量法。在模型中加入控制嵌入向量,以控制生成文本的特定要求。例如,在生成具有特定主题的文本时,可以在模型中加入主题嵌入向量,以指导生成具有特定主题的文本。

(3)控制概率分布法。在生成过程中,根据特定要求的概率分布控制生成文本。例如,在生成长度可变的文本时,可以根据特定长度的概率分布控制生成文本的长度。

(4)基于模板的方法。利用人工设计的模板生成文本。例如,在生成问答对时,可以使用人工设计的问答模板,并通过填充模板中的变量生成具有特定要求的问答对。

(5)强化学习方法。将文本生成任务看作是一个强化学习问题,通过奖励函数指导生成满足特定要求的文本。例如,在生成对话时,可以使用奖励函数指导生成与用户输入相关、流畅、有逻辑的对话。

常见的受控文本生成任务包括机器翻译、文本摘要、对话生成、故事生成等。近年来,随着预训练语言模型的发展,基于预训练模型的受控文本生成方法逐渐成为主流,具有较好的生成质量和泛化能力。此外,也可以结合其他技术,如多模态融合、知识图谱等方法,以进一步提升受控文本生成的效果。

2.4 基于 PLM 的文本生成方法

为了有效地利用 PLMs 完成下游文本生成任务,需要从数据、模型和优化方面重点考虑如下:

(1)输入数据。即如何将输入 x 编码成保留输入语义并且可融合到 PLM 中的表示形式。对于文本生成,包含目标输出关键语义信息的输入数据通常以不同类型出现在不同任务中(例如,顺序文本、结构化表格、多媒体),而大多数 PLMs 通常是在顺序文本数据上进行预训练。因此,开发有效灵活的表示学习方法以捕捉各种类型输入数据中的语义信息是一项重要挑战。

(2)模型架构。即如何设计一个有效的 PLMM 作为生成函数 f_M , 并适应各种文本生成任务非常重要。一些广义架构被开发出来用于一般目的,例如去噪自动编码器或自回归解码器的 PLMs, 尽管这些通用架构无法处理某些特殊情况下的文本生成问题。因此,在底层 PLMs 上进行特定设计以实现良好任务性能非常重要。

(3)优化算法。即如何优化给定参考文字 y 的文字生成函数 f_M , 并确保生成的文本满足特殊的文本属性 P 。为

了生成令人满意的文本,通过开发有效的优化算法学习文本生成函数至关重要,其主要挑战在于一些期望输出文本的所需属性很难被规定或优化。

3 编码输入表示

3.1 非结构化输入

在文本生成中,大多数研究都集中于对非结构化文本(例如句子、段落和文档)进行建模,这需要准确理解输入信息并推导出有意义的文本表示。文本表示学习旨在将输入文本压缩成低维向量,以保留核心语义含义。

(1)段落表征学习。一个段落通常由描述不同主题的多个句子组成,并且每个句子包含一系列单词。为捕捉段落中低级别单词含义和高级别话题语义,在许多研究中提出基于层次或基于图形方法学习段落表征。对于多轮对话这样的多句子段落,典型的方法是将句子连接为整个文本并预测输出文本^[16]。然而,平面连接不能有效地捕捉跨话语的语义动态,这可能导致不准确的生成。为了解决该问题,Gu等^[17]提出分层编码器对输入段落进行建模,利用句子和篇章级Transformer编码器分别编码每个对话话语和话语向量序列。然而,在编码每个单独的话语时,它并没有考虑历史信息,这对于理解对话话语重要。因此,可采用Transformer将每个发言转换成密集向量,在其上设计一个从左到右的流模块以捕获发言级别的动态信息流。

(2)文档表示学习。在许多文本生成任务中,例如文档翻译和文档摘要,输入的文本可能是由多个段落组成的长篇文章。当对这些文件进行编码时,很难建模跨句子(段落)语义并捕获最关键的语义。大多数PLM都是作为掩码语言模型进行训练,它们主要专注于学习标记级别而不是句子级别的表示形式。虽然使用分段嵌入单独表示不同的句子,但它们无法捕获跨句子语义。为了编码跨句子语义,研究者提出以分层方式学习文档表示形式。例如,Liu等^[18]在每个句子开头插入“[CLS]”标记,在较低层次上聚合句子级特征,并将其与更高层次上的自我注意力相结合。此外,Zhang等^[19]提出HIBERT以分层方式学习文档表示形式:使用一个句子编码器将各句子映射到向量,使用一个文档编码器进一步学习给定周围向量并将其作为上下文的情境敏感句子表示,捕获关键语义。在实践中,长篇文档中的句子或段落不可避免地会互补、重叠或冲突。因此,有必要保留最关键内容并在生成的文本中表达出来。为了解决输出文本中缺少关键点的问题,Nguyen等^[20]引入了一个主题模型以捕获全局主题语义,并使用门机制控制提供给文本生成模块的全局语义量。类似地,Liu等^[18]提出了两个基于主题感知对比的学习目标:①一致性检测目标通过检测话题之间的连贯变化识别对话话题;②子摘要生成目标则迫使模型捕获最显著信息并为每个话题生成一个子摘要。

(3)多语言表示学习。现有的PLM主要是在英文文本上进行预训练,而忽略了其他低资源语言。将基于英语的PLM用于解决多语言文本生成任务(例如:多语言机器翻译)颇为困难,于是相应的方法被提出,如跨语言表示法。跨语言表示学习的核心思想是为两种语言学习共享嵌入空间,以提高PLMs在它们之间进行翻译的能力。一种著名的跨语言PLM是XLM^[21],它利用单语言和平行数据学习跨语言表示。然而,在共享字节对编码(BPE)空间上学习到的这些表示是隐式且有限的。因此,Ren等^[41]进一步计算了跨语言 n -gram嵌入,并基于它们推导出一个 n -gram翻译表,以提供显式表示学习信号。再如:多元化代表性。针对超过两种以上不同类型或者领域、语言的情况,多语言PLMs旨在为任何一种语言学习表示形式。基于英文PLMs的BART和T5,Liu等^[42]和Xue等^[43]分别提出mBART和mT5,这些模型可以基于所有语言进行预训练。考虑到不同语言之间的差异(例如语法规则),有研究利用对比学习来学习多元化代表性。特别地,Wang等^[22]提出了两个训练目标:对比句子排名(CSR)和句子对齐替换(SAS)。CSR基于它们的显著性得分创建正面和负面句子对,而SAS则将一个语言中的句子替换成其他语言中的相应内容。通过在共同文本中以对比方式学习这些语言,在跨越各种领域、类型数据集时能够学习到跨越不同领域和类型数据集所支持自然文字表达方式的共享表示空间。

3.2 结构化输入

结构化数据(例如表格、图形和树)是许多实际应用中文本生成的关键输入类型,例如医疗报告和天气预报生成。然而,对于PLMs而言,建模结构化输入并不容易,主要有以下3方面的挑战:①存在结构化数据与PLMs之间的语义鸿沟,因为PLMs通常在自然语言文本上进行预训练;②编码输入数据中的结构信息并不容易;③需要保持所生成文本相对于输入的忠实性。

(1)语义鸿沟。一般而言,PLMs是在非结构化文本上进行预训练,这与结构化数据的形式不同。可利用几种方法弥合该差距,一种方法是结构化数据线性化。为了适应PLMs的结构输入,一个简单的方法是将输入数据线性化成序列。具体而言,Ribeiro等^[23]将知识图(KG)线性化为三元组序列,并连接关系三元组。此外,一些研究采用基于模板的启发式方法对输入数据进行串行化。另一种方法是表示对齐。语义差距使得直接串行化结构化数据时难以有效地注入PLMs,可将结构化数据表示与基于PLM的词嵌入在语义空间中进行对齐。例如,Li等^[52]利用图神经网络(GNN)将KG实体投影到嵌入中,并通过最小Euclidean距离执行表示对齐GNN基础和基于PLM的实体嵌入之间的距离。

(2)捕获结构信息。结构化数据的一个重要特征是以结构化方式表示数据,例如表格中的〈属性,值〉对或知识库中的〈头实体,关系,尾实体〉三元组。这种结构信息可

以更准确地对输入进行建模并帮助生成忠实的文本。为了增强结构信息保留能力,一种典型方法是纳入与结构信息相关的辅助训练目标,如重建输入数据的语义结构;另一种方法是根据结构信息调整输出文本,可将结构信息作为输入。由于PLMs最初是针对顺序输入而开发,因此将额外模块纳入编码结构化输入具有一定意义。

(3)生成高保真文本。在语言学文献中,高保真指生成的文本与结构化数据中的内容一致,生成正确描述结构化输入信息的高保真文本是数据到文本生成算法的关键。为了生成符合输入要求、高忠实度的文本,Gong等^[24]引入基于最优传输匹配损失函数以衡量输入信息和输出文本之间距离的方法。Harkous等^[25]采用语义忠实分类损失函数加以检测并避免产生幻觉等错误。指针生成网络是确保所生成文字对于输入数据具有信仰性的典型方法,可通过将重要单词从输入复制到输出以达成这一目标。例如,Li等^[52]采用该方法将知识库中提取出来的命名体复制到输出文字中,而Suadaa等^[44]将表格值复制到通用占位符中以避免产生不出现在输入表格中但又虚构出现在结果里的短语。为了解决低保真问题,Chen等^[26]认为利用中间意义表示以实现忠实生成非常重要。

4 用于文本生成的PLMs

将输入数据编码为低维表示后,下一步是开发一个有效的PLMM作为文本生成函数 f_M 。基于这样的PLM架构,文本生成目标可以被建模为给定输入数据 x 输出文本 y 的条件概率,它可以令每个token形式化地分解,如式(2)所示。

$$Pr_M(y|x) = \prod_{i=1}^n Pr_M(y_i | y_{<i}, x) \quad (2)$$

其中, y 表示第 i 个输出标记,而 $y_{<i}$ 表示前面的标记 y_1, \dots, y_{i-1} 。为了计算条件概率,传统神经模型主要采用RNN架构及其变体,近年来基于注意力机制的Transformer可以更好地捕捉文本中的长距离依赖关系,这对于建模和生成文本非常有益。由于具有出色的并行化能力,Transformer已成为开发大型PLM(预训练语言模型)的支柱。当在大规模无标注语料库上进行训练时,基于Transformer架构而构建的PLM可以编码丰富的语义或语言知识。此外,PLM可以有效地微调到不同的文本生成任务上,这些都使得PLM成为实现文本生成函数 f_M 的首选。

4.1 标准架构

现有的文本生成预训练语言模型采用单个Transformer或基于Transformer的编码器—解码器,如GPT-3的PLMs使用单个Transformer解码器同时实现输入编码和输出解码过程^[27]。包含3种主要变体:掩蔽语言模型、因果语言模型和前缀语言模型,它们具有不同的注意力掩蔽策略。

(1)掩蔽语言模型。采用Transformer编码器,配备全注意力后,通常会使用掩码式语言建模(MLM)任务进行预训练,即使用双向信息预测被屏蔽的token。最具代表性的模型是BERT,它广泛应用于自然语言理解。但是由于掩码语言模型的预训练任务与下游生成函数之间存在差异,因此很少将其用于文本生成任务。更常见的做法是将掩码语言模型作为文本生成中的编码器部分,以利用其出色的双向编码能力。

(2)因果语言模型。与Transformer解码器类似,因果LM采用对角线掩码矩阵。因果LM旨在进行语言建模,即确定给定单词序列出现在句子中的概率。因果LM对于文本生成很直接,是否可以预测下一个单词取决于所有先前的单词。GPT是文本生成任务的第一个因果LM。GPT-2探索了语言模型零样本生成任务的转移能力,并强调了充足数据量的重要性。此外,GPT-3^[27]表明在一些示例或提示情况下,大规模模型参数可以显著提高下游生成任务性能。CTRL作为有条件的因果语言模型,基于控制代码产生文本以管理风格、内容和特定任务行为。虽然因果LM对于文本生成非常简单直接,但它们具有结构和算法上的限制:因果LM仅从左到右编码标记,而忽略了输入侧双向信息。此外,因果LM并没有专门设计用于序列到序列生成任务,在实践中不会在总结和翻译等任务上达到那般高性能。

(3)前缀语言模型。在单个Transformer上,前缀LM采用输入侧的双向编码方案和输出侧的自然从左到右生成模式。利用混合注意力掩码,输入文本中的标记可以实现相互关注,而目标文本中的标记只能关注所有输入标记和先前生成的标记。UniLM是第一个前缀LM,与因果LM相比,UniLM使用前缀注意掩码解决条件生成任务,类似于编码器—解码器架构。UniLMv2和GLM通过在XLNe中引入排列语言建模改进了普通前缀屏蔽策略^[28]。尽管前缀LM具有自身优点,Raffel等^[45]将单一Transformer前缀LM与基于Transformer的编码器—解码器LM进行了比较,并得出结论:添加显式编码器—解码器注意力可更有效地捕获条件依赖性。

4.2 架构扩展

为了得到性能良好的文本生成PLM,改进Transformer骨干网络的方法被提出。下文介绍两种主要的改进技术,即扩展输入嵌入和改进注意力机制。

(1)扩展输入嵌入。除词嵌入外,几乎所有PLM都使用位置嵌入来指示输入单词的索引。与CNN和RNN相比,自注意操作通常是无序的。因此,提供明确的位置信息以捕获文本顺序性至关重要。原始Transformer使用预先确定的正弦函数绝对位置嵌入,而大多数PLM(例如BERT和GPT以绝对信息作为位置编码,例如T5、UniLMv2和ProphetNet采用桶式相对定位方法进行处理。

(2)改进注意力机制。Transformer中存在各种模块

(例如 position-wise FFN、自注意力等),但相关工作主要集中在改进文本生成的自我和交叉注意力机制上。为了适应长格式文本输入并减轻全局注意力计算的二次复杂度,用稀疏注意替换原始的自我关注以处理长格式输入。每个标记只与特定策略下的其他标记进行关注,如窗口关注、全局关注、随机关注等。实际上,许多文本生成任务需要从多个来源处理输入数据。通常利用一个或多个编码器对多个输入进行编码。因此,可采取不同的策略将跨源输入聚合到交叉注意力模块中。Golovanov等^[29]就对话历史信息、当前状态和角色信息采用平均池化方法;Chen等^[46]和Liu等^[47]提出多视图注意和知识感知注意以处理来自多个视图或知识源的嵌入。Ribeiro等^[48]用GNN替换自我注意模块以更好地提取结构信息;Zeng等^[30]在自我注意之后添加门控机制,注入条件感知信息。

5 预训练语言模型优化

为了获得良好的性能,优化PLM以进行文本生成尤为必要,主要有3种类型的优化方法,即微调、提示学习和属性调整。

5.1 文本生成微调

在预训练期间,PLMs能够从大规模语料库中捕获一般的语言知识。然而,执行下游文本生成任务需要特定于任务的知识。为此,基于下游文本生成数据集调整其权重,微调是将任务特定信息纳入PLMs的流行方法。根据PLMs参数如何更新,可将用于文本生成的现有微调方法分为原始微调、中间微调、多任务微调和轻量化微调。与原始微调相比,中间和多任务微调可以在某种程度上缓解小型文本生成数据集上过拟合问题。由于原始微调需要对整个模型进行参数修改,因此像适配器这样的参数高效微调方法可以轻量级方式对PLMs进行精细化处理。

(1)原始微调。原始微调直接使用具有任务特定损失(例如交叉熵损失)的下游文本生成数据集以更新PLMs。Zhang等^[31]通过将多轮对话会话建模为长文本,并使用语言建模目标优化生成模型,基于GPT-2架构训练了DialogPT模型。原始微调的一个主要问题是它通常在小数据集上优化不充足,容易过拟合。

(2)中间微调。中间微调的基本思想是加入包含足够标记实例的中间数据集。中间数据集可以专注于相同目标文本生成任务但来自于不同领域,或者是相似NLP任务且来自于相同目标领域。从中间数据集注入领域或任务特定知识有助于缓解过拟合问题并提高小目标文本生成数据集上的性能。根据中间数据集与目标文本生成数据集之间的相关性,可以将中级精细划分为两类:领域能适应性中间微调(DAIFT)和任务适应性中间微调(TAIFT)。

(3)多任务微调。多任务微调可利用跨任务知识,通过整合辅助任务,以改进主要的文本生成任务。通过从相

关NLP任务中获取知识,多任务微调可以增强PLMs的鲁棒性,并减少在文本生成任务中需要大量标记实例的需求。根据主要文本生成任务和辅助任务之间的相似度,多任务微调(MTFT)可分为两类:纯MTFT和混合MTFT。

(4)轻量化微调。由于上述微调方法需要更新所有PLM参数,在资源有限的情况下执行整个微调非常耗时。可为文本生成任务开发轻量化微调(PEFT),如基于适配器的轻量化微调。适配器是Houlsby等^[51]提出的一种特殊神经层,用于以参数有效的方式对PLMs进行微调。适配器模块将输入向量投影到一个小向量中,然后使用两个前馈层和一个非线性层将其投影回原始维度。具体而言,适配器首先将原始 d 维特征投影到较小的 m 维度中,应用非线性函数,然后再次投影回 d 维度。每层添加的总参数(包括偏差)为 $2md + d + m$ 。通过设置 $m \ll d$,可以限制每个任务附加参数数量。因此,固定原始PLM的参数而仅微调其适配器是非常高效的^[31]。

5.2 文本生成提示学习

大多数生成式PLM都是使用语言建模目标进行预训练,然后使用特定于任务的目标对文本生成任务进行微调。这种预训练和微调之间的差异影响了PLM在文本生成任务上的性能。作为一种新的学习范式,提示学习^[32]将下游任务重新制定为预训练中的语言建模任务。提示学习有离散提示和连续提示两种方法。

(1)离散提示。早期提示研究通过人工设计基于人类自身模板创建提示。作为开创性研究之一,《GPT-2》使用各种手动创建的提示执行文本生成任务。例如,在机器翻译中使用了“translate to french, [input], [output]”这个提示。该提示定义了从输入数据到输出文本的语义映射,在特定的文本生成任务中使用。通过利用多样化的提醒方式,单个PLM能够执行许多不同类型的文本生成任务。这些方法严重依赖于手动创建提示。但是PLM对提示非常敏感,不正确的提示会导致低性能。为避免需要手动提示,Shin等^[35]提出AutoPrompt以自动搜索模板标记。此外,自动生成离散提示的方法还有重述现有提示^[36]、使用PLM生成提示和从语料库中挖掘提示^[37]。

(2)连续提示。也称为软提示,由嵌入向量组成,已广泛应用于文本生成任务。其具有两个主要特点:放松提示模板必须受自然语言词汇限制;消除模板受PLMs参数化限制。相反,连续提示具有自己的参数,可以基于文本生成任务的训练数据进行优化。使用连续提示进行文本生成的最著名方法是前缀调整,它生成PLMs(例如GPT-2、BART),并优化了一系列特定于任务的向量序列(称为前缀)。与需要存储每个文本生成任务调整后模型的全参数微调不同,前缀调整仅针对每个文本生成任务优化前缀。

5.3 文本属性调整

对于不同的生成任务,需要考虑特定的语言属性以调整PLMs。

(1)相关性。根据语言学文献,在文本生成中,相关性意味着输出文本传达的主题语义与输入文本高度相关。作为代表性例子,在对话系统中,生成的响应应该与历史话语和其他条件(如说话人角色和话题)相关。PLMs利用强大的多层交叉注意机制建模输入和输出之间的语义关联,这可以增强所产生文本与输入数据之间的相关性。此外,Zeng等^[33]利用掩码语言建模目标解决基于各种类型对话上下文生成响应的问题。

(2)忠实度。忠实度是考虑文本生成的重要语言属性,这意味着生成的内容应该遵循输入文本的语义。为了生成忠实的文字,PLMs应准确理解输入的核心语义并获取足够的世界知识以解决下游任务。已有实验表明,PLMs在从纯文本中捕获核心语义方面具有出色的自然语言理解能力,它们确实编码了大量世界知识,这可能通过将背景知识注入到文字中产生益处以生成忠实摘要。

(3)顺序保留。在NLP领域,顺序保留是一个特殊的属性,它指输入和输出文本中语义单元(单词、短语等)的顺序一致。这种属性对于重要的文本生成任务非常关键,例如文本改写和机器翻译。在机器翻译中,在对源语言到目标语言进行翻译时,通常需要保持源文本和目标文本中某些短语的顺序以确保翻译结果的准确性。在机器翻译中,单词对齐是实现顺序保存属性的常用方法。代表性学习是Code-Switching Pre-training(CSP),CSP首先从源和目标单语料库自动提取单词对齐信息,然后为了增强转换期间的有序性质,其通过预测给定目标侧上已经排好队列的片段并不断地预训练PLMs的句子片段以实现有序性质。为了放松离散化对齐约束,其一条线路旨在进行连续表示对齐以提高有序性质。Wada等^[49]专注于将每种语言的单词嵌入映射到公共潜在空间以对齐每种语言的单词表示。

6 文本生成的挑战

文本生成面临的挑战及可能的解决方案如表1所示。

Table 1 Challenge and solutions

表1 挑战与解决方案

方面	挑战	基于PLM的解决方案
数据	训练数据不足	知识转移;数据增强;多任务学习
	语料库数据偏见	减少词嵌入偏见;识别和屏蔽对偏差敏感令牌
模型	模型压缩	量化减少PLM权重;修剪较少的关键权重;知识蒸馏
	模型增强	大规模PLM;知识丰富PLM;领域PLM

6.1 数据

(1)缺乏足够的训练数据。在一些文本生成任务中,很难获得足够的训练数据。知识迁移提供了一个有效解决方案,即将数据丰富的源任务的知识转移到缺乏数据的目标文本生成任务中。此外,还可以使用数据增强和多任

务学习解决该问题。数据增强是通过添加已经存在的稍加修改的副本或从现有数据中新创建的合成数据;多任务学习是利用其他数据丰富的任务和数据集克服数据稀缺问题。大多数研究通常会为增强主要文本生成任务而纳入人类的辅助生成任务。然而,这些方法通常采用每个任务独立的解码器,打破高资源和低资源文本生成任务之间的语义联系。为了弥合这一差距,Bai等^[50]采用了一个统一的解码器,在机器翻译中学习多种语言之间的对齐和模式。

(2)预训练语料库中的数据偏见。通常使用真实世界中的数据进行训练,以便它们能够建模训练数据的统计属性。因此,它们继承了在该数据中普遍存在的偏见和刻板印象,如性别歧视、种族歧视等。这些偏见和刻板印象可能会给下游文本生成任务带来重大挑战。一个简单方法是在训练数据中“交换”有性别特征的术语以生成词嵌入。此外,仅屏蔽名称和代词也可以减少偏差并提高某些语言任务的性能。然而,迄今为止,还没有一种通用、统一的方法可减少PLMs对文本生成产生的数据偏见。

6.2 模型

(1)模型压缩。虽然PLMs在文本生成方面取得了巨大成功,但骨干Transformer仍然笨重且资源消耗大,导致内存消耗高、计算开销大和电力能源成本高。为解决这些问题,越来越多的方法被提出用以压缩PLMs,例如量化、剪枝和知识蒸馏。量化意味着减少用于表示PLMs权重的向量数量,从而可以使用更少的位数表示它们的剪枝;剪枝是指识别和删除冗余和,或不重要的权重;知识蒸馏是指使用PLMs输出训练较小模型,通过复制输入数据和输出文本之间的注意力分布,小模型也可以学习输入和输出之间的上下文依赖关系。

(2)模型增强。尽管PLM在当今取得了巨大成功,但仍未达到人们预期。由此,研究界对加强现有PLM以提高文本生成性能产生了浓厚兴趣。通过扩展PLMs参数数量可以提升其性能,这研究结果引发了大规模文本生成中的PLMs开发。最具代表性的是GPT-3,其中包含1750亿个参数,比任何先前非稀疏PLMs多10倍。由于拥有众多参数,GPT-3可以在各种文本生成任务中实现强劲的性能而无需进行梯度更新或微调。最近研究表明,通过精心构建模型架构,可以使用较少的预训练数据或较低的预训练成本获得等效或更高的文本生成性能。

7 文本生成评估

文本风格转换生成质量评价方法通常分为人工评价和自动评价两种。自动评估较人工评估更便捷,且可重复。这两种方法都从以下3个方面对生成文本进行评估:风格转换准确率、内容保留度和文本流利度。常见自动评估指标如下:

(1) BLEU (Bilingual Evaluation Understudy)^[38]。它是用于机器翻译等任务的评估指标,通过比较生成文本和参考文本之间的 n -gram 重叠度评估生成文本的准确性。优点是简单易实现,但它并不能完全准确地衡量译文的质量,例如可能无法判断机器翻译译文是否确实传达了源文本的意思。其计算公式如式(3)所示。

$$BLEU = BP \times \exp\left(\sum_{n=1}^N W_n \log(P_n)\right) \quad (3)$$

其中, BP 为惩罚因子,可以防止短句子得分过高; W_n 为 n -gram 的权重,通常取 $1/N$; P_n 为生成文本中 n -gram 出现的概率。

(2) ROUGE (Recall-Oriented Understudy for Gisting Evaluation)^[39]。它是用于文本摘要等任务的评估指标,通过比较生成文本和参考摘要之间的 n -gram 重叠度评估生成文本的准确性,优点是可以在短时间内快速评估文本摘要系统的性能,但缺点是忽略了语法和上下文之间的关系。其计算公式如式(4)所示。

$$ROUGE = (1 - \beta) \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \beta \cdot \text{recall}} \quad (4)$$

其中, β 为调整因子,一般取 1, precision 为生成文本中与参考摘要相同的 n -gram 数量占生成文本总数的比例, recall 为生成文本中与参考摘要相同的 n -gram 数量占参考摘要总数的比例。

(3) 困惑度 (Perplexity)^[40]。它用来计算一个句子生成的概率,用于判断一句话是否出自人口,可以衡量语言模型在生成文本时的流畅性和连贯性。其缺点是困惑度只考虑了单个词的生成,忽略了上下文的影响,因此可能无法准确地评估模型性能。其计算过程如式(5)所示。

$$S = W_1, W_2, \dots, W_k \quad (5)$$

则句子的概率可以表示为如式(6)所示。

$$P(S) = P(W_1, W_2, \dots, W_k) \quad (6)$$

困惑度思想是语言模型给测试集的句子赋予较高概率值,原因在于测试集的句子较好。训练之后,语言模型在测试集上的概率越高越好,其公式如式(7)所示。

$$PP(W) = P(W_1, W_2, \dots, W_k)^{-\frac{1}{N}} \quad (7)$$

句子概率大,准确率高,语言模型就越好,困惑度也越小。

8 结语

本文探讨了基于预训练语言模型的受控文本生成技术,这是一种生成特定类型文本的方法,如摘要、翻译、对话等。首先,介绍了应用 PLMs 到文本生成时的输入表示学习、模型架构设计和参数优化;然后,讨论了当前受控文本生成技术存在的问题和挑战,包括对数据集和模型的要求高、受控程度难以平衡、模型可解释性不佳等问题;最后,对各种评估指标进行了回顾总结。未来,期待更多的

研究能够解决这些问题,进一步提高受控文本生成技术的性能和可靠性。总体而言,基于预训练语言模型的受控文本生成技术是一项非常有前途的技术,有着广泛的应用前景。相信,随着技术的不断发展,受控文本生成技术将发挥越来越重要的作用。

参考文献:

- [1] ZHOU L, GAO J, LI D, et al. The design and implementation of xiaoice, an empathetic social chatbot[J]. Computational Linguistics, 2020, 46(1): 53-93.
- [2] LAMPLE G, CONNEAU A. Cross-lingual language model pretraining [DB/OL]. <https://arxiv.org/pdf/1901.07291>, 2019.
- [3] BROWN P F, COCKE J, DELLA P S A, et al. A statistical approach to machine translation[J]. Computational Linguistics, 1990, 16(2): 79-85.
- [4] BROWN R D, FREDERKING R. Applying statistical English language modelling to symbolic machine translation [C]//Proceedings of the Sixth Conference on Theoretical and Methodological Issues in Machine Translation of Natural Languages, 1995: 221-239.
- [5] LI J, LI S, ZHAO W X, et al. Knowledge-enhanced personalized review generation with capsule graph neural network [C]//Proceedings of the 29th ACM International Conference on Information & Knowledge Management, 2020: 735-744.
- [6] LI J, ZHAO W X, WEN J R, et al. Generating long and informative reviews with aspect-aware coarse-to-fine decoding [DB/OL]. <https://arxiv.org/pdf/1906.05667.pdf>, 2019.
- [7] SEE A, LIU P J, MANNING C D. Get to the point: summarization with pointer-generator networks [DB/OL]. <https://arxiv.org/pdf/1704.04368.pdf>, 2017.
- [8] LAQBAL T, QURESHI S. The survey: text generation models in deep learning [J]. Journal of King Saud University-Computer and Information Sciences, 2022, 34(6): 2515-2528.
- [9] QIU X, SUN T, XU Y, et al. Pre-trained models for natural language processing: a survey [J]. Science China Technological Sciences, 2020, 63(10): 1872-1897.
- [10] GARCIA X, FORET P, SELAM T, et al. A multilingual view of unsupervised machine translation [DB/OL]. <https://arxiv.org/pdf/2002.02955.pdf>, 2020.
- [11] RFORS A, WU J, CHILD R, et al. Language models are unsupervised multitask learners [DB/OL]. <https://insightcivic.s3.us-east-1.amazonaws.com/language-models.pdf>, 2019.
- [12] KALYAN K S, RAJASEKHARAN A, SANGEETHA S. Ammus: a survey of transformer-based pretrained models in natural language processing [DB/OL]. <https://arxiv.org/pdf/2108.05542.pdf>, 2021.
- [13] HAN X, ZHANG Z, DING N, et al. Pre-trained models: past, present and future [J]. AI Open, 2021, 2: 225-250.
- [14] EL-KASSAS W S, SALAMA C R, RAFAA A A, et al. Automatic text summarization: a comprehensive survey [J]. Expert Systems with Applications, 2021, 165: 113679.
- [15] ZAIB M, SHENG Q Z, EMMA Z W. A short survey of pre-trained language models for conversational ai—a new age in NLP [C]//Proceedings of the Australasian Computer Science Week Multiconference, 2020: 1-4.
- [16] BAO S, HE H, WANG F, et al. Plato-2: towards building an open-domain chatbot via curriculum learning [DB/OL]. <https://arxiv.org/pdf/2006.16779.pdf>, 2020.
- [17] GU X, YOO K M, HA J W. Dialogbert: discourse-aware response generation via learning to recover and rank utterances [C]//Proceedings of the

- AAAI Conference on Artificial Intelligence, 2021: 12911–12919.
- [18] LIU Y, LAPATA M. Text summarization with pretrained encoders [DB/OL]. <https://arxiv.org/pdf/1908.08345v2.pdf>, 2019.
- [19] ZHANG X, WEI F, ZHOU M. HIBERT: document level pre-training of hierarchical bidirectional transformers for document summarization [DB/OL]. <https://arxiv.org/pdf/1905.06566.pdf>, 2019.
- [20] NGUYEN T, LUU A T, LU T, et al. Enriching and controlling global semantics for text summarization [DB/OL]. <https://arxiv.org/pdf/2109.10616.pdf>, 2021.
- [21] CHI Z, DONG L, WEI F, et al. Cross-lingual natural language generation via pre-training [C]//Proceedings of the AAAI Conference on Artificial Intelligence, 2020: 7570–7577.
- [22] WANG D, CHEN J, ZHOU H. Contrastive aligned joint learning for multilingual summarization [C]//Findings of the Association for Computational Linguistics, 2021: 2739–2750.
- [23] RIBEIRO L F R, SCHMITT M, SCHUTZE H, et al. Investigating pretrained language models for graph-to-text generation [DB/OL]. <https://arxiv.org/pdf/2007.08426.pdf>, 2020.
- [24] GONG H, SUN Y, FENG X, et al. Tablegpt: few-shot table-to-text generation with table structure reconstruction and content matching [C]//Proceedings of the 28th International Conference on Computational Linguistics, 2020: 1978–1988.
- [25] HARKOUS H, GROVES I, SAFFARI A. Have your text and use it too! end-to-end neural data-to-text generation with semantic fidelity [DB/OL]. <https://arxiv.org/pdf/2004.06577.pdf>, 2020.
- [26] CHEN Y C, GAN Z, CHENG Y, et al. Distilling knowledge learned in BERT for text generation [DB/OL]. <https://arxiv.org/pdf/1911.03829.pdf>, 2019.
- [27] BROWB T, MANN B, RYDER N, et al. Language models are few-shot learners [J]. *Advances in Neural Information Processing Systems*, 2020, 33: 1877–1901.
- [28] BAO H, DONG L, WEI F, et al. Unilmv2: pseudo-masked language models for unified language model pre-training [C]//International Conference on Machine Learning, 2020: 642–652.
- [29] GOLOVANOV S, KURBANOV R, NIKOLENKO S, et al. Large-scale transfer learning for natural language generation [C]//Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, 2019: 6053–6058.
- [30] ZENG Y, NIE J Y. A simple and efficient multi-task learning approach for conditioned dialogue generation [DB/OL]. <https://arxiv.org/pdf/2010.11140.pdf>, 2020.
- [31] ZHANG Y, SUN S, GALLEY M, et al. Dialogpt: large-scale generative pre-training for conversational response generation [DB/OL]. <https://arxiv.org/pdf/1911.00536.pdf>, 2019.
- [32] CELILYILMZA A, CLARK E, GAO J. Evaluation of text generation: a survey [DB/OL]. <https://arxiv.org/pdf/2006.14799.pdf>, 2020.
- [33] ZENG Y, NIE J Y. Generalized conditioned dialogue generation based on pre-trained language model [DB/OL]. <https://arxiv.org/pdf/2010.11140.pdf>, 2020.
- [34] GARRIDO-MUNOZ I, MONTEGO-RAEZ A, MARTINEZ-SANTIAGO F, et al. A survey on bias in deep nlp [J]. *Applied Sciences*, 2021, 11(7): 3184.
- [35] SHIN T, RAZEGHI Y, LOGAN IV R L, et al. Autoprompt: eliciting knowledge from language models with automatically generated prompts [DB/OL]. <https://arxiv.org/pdf/2010.15980.pdf>, 2020.
- [36] JIANG Z, XU F F, ARAKI J, et al. How can we know what language models know? [J]. *Transactions of the Association for Computational Linguistics*, 2020, 8: 423–438.
- [37] GAO T, FISCH A, CHEN D. Making pre-trained language models better few-shot learners [DB/OL]. <https://arxiv.org/pdf/2012.15723.pdf>, 2020.
- [38] PAPINENI K, ROUKOS S, WARD T, et al. Bleu: a method for automatic evaluation of machine translation [C]//Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, 2002: 311–318.
- [39] LIN C Y. Rouge: a package for automatic evaluation of summaries [C]//Proceedings of the Workshop on Text Summarization Branches Out, 2004: 74–81.
- [40] JELINEK F, MERCER R L, BAHL L R, et al. Perplexity—a measure of the difficulty of speech recognition tasks [J]. *The Journal of the Acoustical Society of America*, 1977, 62(S1): S63.
- [41] REN S, WU Y, LIU S, et al. Explicit cross-lingual pre-training for unsupervised machine translation [DB/OL]. <https://arxiv.org/pdf/1909.00180.pdf>, 2019.
- [42] LIU Y, GU J, GOYAL N, et al. Multilingual denoising pre-training for neural machine translation [J]. *Transactions of the Association for Computational Linguistics*, 2020, 8: 726–742.
- [43] XUE L, CONSTANT N, ROBERTS A, et al. mT5: a massively multilingual pre-trained text-to-text transformer [DB/OL]. <https://arxiv.org/pdf/2010.11934.pdf>, 2020.
- [44] SUADAA L H, KAMIGAITO H, FUNAKOSHI K, et al. Towards table-to-text generation with numerical reasoning [C]//Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, 2021: 1451–1465.
- [45] RAFFEL C, SHAZEER N, ROBERTS A, et al. Exploring the limits of transfer learning with a unified text-to-text transformer [J]. *The Journal of Machine Learning Research*, 2020, 21(1): 5485–5551.
- [46] CHEN J, YANG D. Multi-view sequence-to-sequence models with conversational structure for abstractive dialogue summarization [DB/OL]. <https://arxiv.org/pdf/2010.01672.pdf>, 2020.
- [47] LIU S, ZHAO X, LI B, et al. A three-stage learning framework for low-resource knowledge-grounded dialogue generation [DB/OL]. <https://arxiv.org/pdf/2109.04096.pdf>, 2021.
- [48] RIBEIRO L F R, ZHANG Y, GUREVYCH I. Structural adapters in pretrained language models for amr-to-text generation [DB/OL]. <https://arxiv.org/pdf/2103.09120.pdf>, 2021.
- [49] WADA T, IWATA T. Unsupervised cross-lingual word embedding by multilingual neural language models [DB/OL]. <https://arxiv.org/pdf/1809.02306.pdf>, 2018.
- [50] BAI Y, GAO Y, HUANG H. Cross-lingual abstractive summarization with limited parallel resources [DB/OL]. <https://arxiv.org/pdf/2105.13648.pdf>, 2021.
- [51] HOULSBY N, GIURGIU A, JASTRZEBSKI S, et al. Parameter-efficient transfer learning for NLP [C]//Proceedings of International Conference on Machine Learning, 2019: 2790–2799.
- [52] LI J, TANG T, ZHAO W X, et al. Few-shot knowledge graph-to-text generation with pretrained language models [DB/OL]. <https://arxiv.org/pdf/2106.01623.pdf>, 2021.