

基于风格感知的无监督领域适应算法

宁秋怡 史小静 段湘煜 张民

苏州大学计算机科学与技术学院 江苏 苏州 215006

(qiuyning@stu.suda.edu.cn)



摘要 近年来,神经机器翻译的译文质量取得了显著的进步,但是其在训练过程中严重依赖平行的双语句子对。然而对于电子商务领域来说,平行资源是稀缺的,此外,文化的不同导致产品信息表达存在风格差异。为了解决这两个问题,提出了一种基于风格感知的无监督领域适应算法,该算法在互训练方法中充分利用电子商务单语数据,同时引入拟知识蒸馏的方法处理风格差异。通过获取电商产品数据信息构建非平行双语语料,基于该语料以及中英新闻平行语料进行多组实验,结果表明,相比各种无监督领域适应方法,该算法显著提高了翻译质量,较最强的基线系统提高了约5个BLEU点。此外,将该算法在Ted, Law和Medical OPUS 3类数据上进一步拓展应用,均取得了更佳的翻译效果。

关键词: 机器翻译;无监督;领域适应;风格感知;电子商务

中图法分类号 TP183

Unsupervised Domain Adaptation Based on Style Aware

NING Qiu-yi, SHI Xiao-jing, DUAN Xiang-yu and ZHANG Min

School of Computer Science and Technology, Soochow University, Suzhou, Jiangsu 215006, China

Abstract In recent years, neural machine translation has made significant progress in translation quality, but it relies on parallel bilingual sentence pairs heavily during the training process. However, parallel resources are scarce for the e-commerce domain, in addition, cultural differences lead to stylistic differences in product information expression. To solve these two problems, a style-aware unsupervised domain adaptation algorithm is proposed, which makes full use of e-commerce monolingual data in the mutual training method, while introducing quasi knowledge distillation approach to deal with style differences. We construct non-parallel bilingual corpus by obtaining e-commerce product data information, and then carry out experiments based on the aforementioned corpus and Chinese and English news parallel corpus. The results show that the algorithm significantly improves translation quality compared to various unsupervised domain adaptation methods, improves about 5 BLEU points compared with the strongest baseline system. In addition, the algorithm is further extended to Ted, Law and Medical OPUS data, all of which achieve better translation results.

Keywords Machine translation, Unsupervised, Domain adaptation, Style aware, E-commerce

1 引言

机器翻译(Machine Translation)应用于电子商务产品的信息翻译,可以解决人工翻译低效和成本高的问题。目前机器翻译系统已经达到了很高的水平,特别是神经机器翻译(Neural Machine Translation, NMT)系统。其训练过程依赖于大规模平行语料库,但对于电子商务领域来说此类资源是稀缺的,这极大地阻碍了机器翻译在电子商务产品信息翻译中的应用。

为了避开电子商务领域对平行数据的需求,本文提出了一种新的无监督领域适应方法。在现有的无监督领域适应翻译研究任务中,其主要存在领域内和领域外两者之间不匹配的问题。但在电子商务无监督领域适应机器翻译任务中,不

仅存在领域不匹配问题,还有因区域文化及语言习惯带来的电商产品描述性差异,例如,对于同一类产品(如鞋子),不同的语言给出的相应特性描述示例如下:

中文描述:材料防滑耐磨平衡轻便透气。

英文描述:It's abrasion-resistant material is forged specifically for skating, has great slip resistance, and will keep you balance during sporting.

以上商品描述分别来自英文电子商务平台与中文电子商务平台,可以看出中文的产品描述是以词汇无序组合的形式给出,而英文的产品描述更加流畅自然,这种风格加剧了领域内数据的不平行性,导致电子商务领域翻译更加困难。

为了解决上述问题,即领域不匹配和风格差异问题,本文提出了风格感知的无监督领域适应(Style-Aware Unsuper-

收稿日期:2020-12-09 返修日期:2021-03-21 本文已加入开放科学计划(OSID),请扫描上方二维码获取补充信息。

基金项目:国家自然科学基金(61673289)

This work was supported by the National Natural Science Foundation of China(61673289).

通信作者:段湘煜(xiangyuduan@suda.edu.cn)

vised Domain Adaptation, SA-UDA)算法,该算法应用了一种新的互训练(Mutual Training)过程,通过利用新闻领域的平行数据和电子商务领域的单语数据,相互加强源端到目标端和目标端到源端的翻译。并在此基础上提出拟知识蒸馏(Quasi Knowledge Distillation, QKD)来处理风格差异,旨在翻译过程中保留正确的风格。本文进行了一系列的电子商务产品信息翻译实验,结果表明,本文提出的 SA-UDA 算法明显优于各种无监督领域自适应方法。其中用于对比的无监督领域自适应方法包括拷贝机制(copy)^[1]、反向翻译(back translation)^[2]、领域感知特征嵌入^[3]以及基于词归纳的领域适应^[4]。本文方法较最强的无监督领域自适应方法在翻译性能上提升约5个BLEU点。此外,将SA-UDA的互训练过程以及拟知识蒸馏应用于Ted, Law和Medical OPUS数据集上,其翻译性能较基准系统均有所提升。

本文的贡献包含以下几点:

(1)通过获取的中英电子商务产品数据信息构建非平行双语语料,搭建电子商务机器翻译模型。在电子商务数据上进行多方面探索和尝试,通过分析对比,挖掘中英电子商务语料的风格特性。

(2)提出了一种新的互训练框架,充分利用领域外的平行数据以及领域内的单语数据,增强双向的翻译模型,解决领域不匹配问题。

(3)提出拟知识蒸馏方法,通过老师模型学习到的风格信息,以蒸馏的方式教给学生模型,使得学生模型在翻译过程中保留正确的风格信息,解决风格差异问题。

(4)在多个领域进行无监督领域适应实验,结果表明,SA-UDA在无监督领域适应中能够有效提升翻译性能,在无监督领域适应中具有通用性。

本文第2节介绍了近年来在领域适应机器翻译方面的相关研究以及本文工作;第3节详细描述了SA-UDA的框架和实现过程;第4节介绍实验设置和实验结果;第5节对相关实验进行详细分析;第6节通过拓展实验进一步论证本文方法的有效性以及通用性;最后总结全文。

2 相关工作

Shen等^[5]利用小型电子商务平行数据集改进词级别的翻译性能,但构建平行数据是一项劳动密集型任务。无监督的领域适应是避开平行语料构建任务的一个很好的选择。

目前在NMT无监督领域适应中主要有两类方法来优化

翻译模型:一类是构造领域内的伪平行数据,另一类是设计新型网络模型。Sennrich等^[2]以及Zhang等^[6]通过将领域内目标语言反向翻译来构建伪平行数据,并在训练过程中始终保持目标端的真实性,从而提升翻译性能。Currey等^[1]采用复制领域内目标语言文本的方法,通过不断更新扩充领域内的数据来提升翻译性能。Hu等^[4]通过词典归纳构建出领域内伪平行语料,极大地提升了翻译质量。此外,He等^[7]在利用伪平行数据的基础上,提出对偶学习(Dual Learning)的方法,提升了模型的性能。Zhen等^[8]在翻译模型中引入对抗分类器的方法实现无监督领域适应。Dou等^[3]通过将隐藏的状态分解成不同的部分来学习特定领域的特征,并采用多任务学习方式训练网络,提升翻译性能。

本文所提方法与文献^[9]提出的镜像生成式机器翻译不同,镜像生成式机器翻译能够在双向训练中借助语言模型从而提升翻译性能,而本文所提方法应用于含有风格差异的翻译训练过程,其目标语言与真实待翻译语句存在差异。SA-UDA不仅使用NMT,而且还利用了统计机器翻译(Statistical Machine Translation, SMT)来完成翻译任务。此外,我们的工作不同于风格转换和Niu等^[10]的敏感机器翻译。风格转换是将一个句子从一种风格转换为另一种风格,同时保持其内容含义不变。Niu等^[10]的敏感机器翻译(sensitive machine translation)是通过两种不同风格表达相同意思的单语句子对(用于单语风格转换的平行数据)和具有不同风格的双语平行句子进行训练,最终产生所需的翻译。而本文工作在处理双语风格差异时维持各自的风格。

3 基于SA-UDA的电商产品翻译

3.1 基本框架

基于风格感知的无监督领域适应算法的整体框架如图1所示。其中, src_{out} 和 tgt_{out} 分别表示领域外平行的真实源端语句与目标端语句, $src_{in}^{src_style}$ 和 $tgt_{in}^{tgt_style}$ 分别为领域内表现出源端风格的源句和目标端风格的目标句,“...”表示用来解码的翻译模型,“ \rightarrow ”表示通过使用相应的翻译模型进行解码,“ \rightarrow ”表示训练相应的翻译模型,“ \Rightarrow ”表示用正确的风格进行翻译的拟知识蒸馏。带^的单词,如 \widehat{src} 表示一个伪的源端句子,它是“ \rightarrow ”的输出结果。框架的左侧顶部用 $M_{src \rightarrow tgt}$ 表示训练源端到目标端的翻译模型,底部用 $M_{tgt \rightarrow src}$ 表示训练目标端到源端的翻译模型,两个模型通过 N 轮互训练,互相加强。

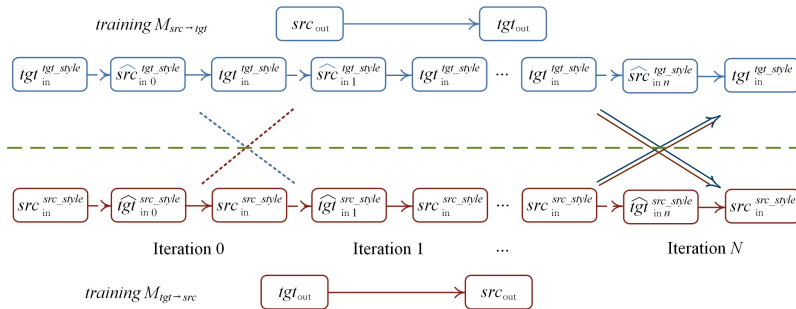


图1 风格感知无监督领域适应的整体框架

Fig. 1 Overall framework of style-aware unsupervised domain adaptation

本文算法通过利用领域外真实平行句和领域内伪平行句

来训练翻译模型。例如,图1的上半部分中, src_{out} 和 tgt_{out}

构成真正的领域外句子对, $\widehat{src}_{in\ n}^{tgt_style}$ 和 $tgt_{in\ n}^{tgt_style}$ 构成每轮领域内句子对。这两部分句子对联合一起训练翻译模型 $M_{src \rightarrow tgt}$ 。其中领域内伪句子对由反向翻译生成,它使用翻译模型 $M_{tgt \rightarrow src}$ 将 $tgt_{in\ n}^{tgt_style}$ 翻译成 $\widehat{src}_{in\ n}^{tgt_style}$, 然后将 $\widehat{src}_{in\ n}^{tgt_style}$ 与 $tgt_{in\ n}^{tgt_style}$ 配对。图 1 的下半部分与上半部分相对。这两个部分通过互训练过程相互加强,详细过程将在 3.2 节进行介绍。

对于领域内数据,在给定输入文本风格的情况下,翻译模型通常以相同的风格输出翻译,但是训练数据和解码数据之间的风格是不一致。例如,对目标风格 $\widehat{src}_{in\ n}^{tgt_style} \rightarrow tgt_{in\ n}^{tgt_style}$ 进行训练,以获得翻译模型 $M_{src \rightarrow tgt}$,但是需要解码的句子 $src_{in\ n}^{src_style}$ 是含有源端句子风格的。在具有特定风格数据上“ \rightarrow ”训练的翻译模型被应用到由“ \leftarrow ”产生的具有另一种风格的输入文本,本文在训练好的翻译模型上进行风格感知的拟知识蒸馏。

下面详细介绍相互训练的过程、拟知识蒸馏以及它们整合的方法。

3.2 互训练过程

互训练的过程如图 1 所示,训练 $M_{src \rightarrow tgt}$ 的伪平行句是通过 $M_{tgt \rightarrow src}$ 解码生成的,反方向的模型训练也是如此,解码结果越好,得到的训练模型的质量越佳,其翻译效果也更好。不断重复互训练过程,直到两个翻译模型收敛。需要注意的是:在每个训练步骤中,来自领域外的真实平行句子与领域内生成的伪平行句子拼接,合并形成领域训练数据。翻译模型 $M_{src \rightarrow tgt}$ 和 $M_{tgt \rightarrow src}$ 可以是统计机器翻译(SMT)模型,也可以是神经机器翻译(NMT)模型。

3.2.1 互训练 SMT 模型

统计机器翻译模型是基于对数线性模型,其概率分布公式如式(1)所示:

$$p(y|x, \lambda) = \frac{\exp(\sum_{m=1}^M \lambda_m f_m(x, y))}{Z} \quad (1)$$

其中, x 和 y 分别是源句子和目标句子, Z 是归一化因子, f 是特征向量, λ 是其权重,总共有 M 个特征和权重用于训练。

本文采用基于短语的统计模型(Phrase-Based Statistical Machine Translation, PBSMT)^[11] 作为无监督领域适应机器翻译的 SMT 模型。

在互训练 SMT 模型的初始化步骤中,仅对域内数据使用无监督的 SMT^[12] 来获得领域内的短语对。使用 phrase2vec¹⁾ 训练 n-gram 嵌入,然后通过向量映射²⁾ 在嵌入空间的近邻短语对中提取其短语表。

利用初始短语表和其他特征,解码每个单词句子以建立一个领域内初始的伪平行语料。然后将领域外真实平行句子和领域内伪平行句子拼接,训练新的 $M_{src \rightarrow tgt}$ 和 $M_{tgt \rightarrow src}$ 。这两个模型都用于生成一个新的领域内伪平行语料,该语料再次与领域外真实数据拼接,训练新的 $M_{src \rightarrow tgt}$ 和 $M_{tgt \rightarrow src}$ 。在这过程中,始终使用无监督调参^[13] 来优化参数 λ ,不断重复直到翻译模型收敛。

3.2.2 互训练 NMT 模型

NMT 机器翻译由编码器和解码器^[14-15] 构成,将源句子 x 编码成隐藏向量,根据已经解码生成的词的隐藏层和对源端隐藏向量的注意力得到目标句子 y ,其表达式如下:

$$p(y|x; \theta) = \prod_i \text{softmax} = \prod_i \text{softmax}(g(\mathbf{h}_{y_i}, \mathbf{h}_{y_{<i}}, \mathbf{c}_i; \theta)) \quad (2)$$

$$\mathcal{L}_{\text{NLL}}(\theta) = - \sum_i \log p(y_i | x, y_{<i}; \theta) \quad (3)$$

其中, θ 表示神经网络的训练参数,每个目标词的概率是模型的输出通过一个线性变换再经过 softmax 得到的。函数 g 将目标生成的隐藏向量 \mathbf{h}_{y_i} , $\mathbf{h}_{y_{<i}}$ 和当前的 \mathbf{c}_i 注意力作为输入,训练最小化损失函数 \mathcal{L}_{NLL} 。

我们使用一个简单的初始化步骤来相互训练 NMT 模型。首先在外领域的真实平行数据上训练 $M_{src \rightarrow tgt}$ 和 $M_{tgt \rightarrow src}$,然后直接用它们解码领域内的数据,从而得到伪平行数据。接下来互训练过程是将真实平行数据和伪平行数据拼接,以训练新的 NMT 模型,并产生新的伪平行数据。整个过程使用无监督调参,将源语言中的一组句子作为输入,将其翻译成目标语言并返回源端语言,并以原始源句子集为参考计算其 BLEU 分数。重复训练,直到无法获得进一步的改进。

3.2.3 拟知识蒸馏

知识蒸馏^[16] 的基本目的是用一个老师的概率分布最小化交叉熵:

$$\mathcal{L}_{\text{KD}}(\theta, \theta_T) = - \sum_i q(y_i | x, y_{<i}; \theta_T) \times \log p(y_i | x, y_{<i}; \theta) \quad (4)$$

其中, θ_T 表示教师网络的参数; θ 表示学生网络参数; q 和 p 分别是教师和学生的词汇分布; \times 表示点积。通过最小化 \mathcal{L}_{KD} , p 将会类似于 q 。

图 1 所示的翻译模型是在一种风格的数据上进行训练的,而解码是在另一种风格的数据上进行训练的。为了缓解这种训练和解码之间的风格不一致,本文提出了拟知识蒸馏,即使用模拟教师翻译模型以正确的风格来指导学生模型的训练过程。

拟知识蒸馏的主要任务是模拟教师正确的分布,使得学生模型在解码时能够得到正确的风格。正确风格的解码输出与源端输入配对,用于构建训练教师模型的正确风格的平行语料。例如,将目标风格数据 $\{\widehat{src}_{in\ n}^{tgt_style}, tgt_{in\ n}^{tgt_style}\}$ 训练的模型 $M_{src \rightarrow tgt}$ 设置为学生模型,使用上一次迭代的 $M_{src \rightarrow tgt}$ 将 $src_{in\ n}^{src_style}$ 翻译得到 $tgt_{in\ n-1}^{src_style}$,获得源风格平行语料来训练模拟源端风格的教师模型,该模型输出 $src \rightarrow tgt$ 翻译的分布 q 。最后,在 $\widehat{src}_{in\ n}^{tgt_style} \rightarrow tgt_{in\ n}^{tgt_style}$ 的目标风格数据上,教师分布 q 将指导学生模型 $M_{src \rightarrow tgt}$ 的训练,以适应源风格。在 $M_{tgt \rightarrow src}$ 训练中也采用了类似的过程。

此外,拟知识蒸馏过程采取插值损失的方式,其表达式如式(5)所示。通过联合训练学生模型参数 θ 和教师模型参数 θ_T 最小化式(3)和式(4)的损失。

$$\mathcal{L}(\theta, \theta_T) = (1 - \alpha) \mathcal{L}_{\text{NLL}}(\theta) + \alpha \mathcal{L}_{\text{KD}}(\theta, \theta_T) \quad (5)$$

1) <https://github.com/artetxem/phrase2vec>

2) <https://github.com/artetxem/vecmap>

3.2.4 综合训练

本节将上述模块整合到一个训练程序中。首先,使用 SMT 的互训练,因为我们发现在最开始时 SMT 优于 NMT;然后,在互训练 SMT 的基础上继续进行互训练 NMT,该过程进一步提升了翻译性能;最后,对翻译模型风格进行了拟知识蒸馏。

从基于 SMT 的互训练到基于 NMT 的互训练的过渡中,我们采用了一种渐进的方式,即在 NMT 模型的互训练的第一次迭代中,伪平行语料完全由 SMT 系统生成。随着训练的进行,NMT 模型变得越来越好,而伪平行语料将包含越来越多由 NMT 生成的数据。在 NMT 的第 n 次迭代中增加 $n \times a$ 的伪平行句子(其中 a 是控制数据中包含来自 NMT 的句子数量的超参数),并且减少来自 SMT 的相同数量的句子以保持训练数据规模不变。

4 实验

本文进行了中英电子商务产品信息翻译的实验,中文和英文是两种活跃的在线交易语言,有大量的网络文本。

4.1 实验设置

(1) 数据准备

为了构建电子商务产品信息的机器翻译系统,电子商务领域的数据是必不可少的。将电子商务领域作为内领域,对从电商网站获取到的中英文产品文本进行加工处理,将其大致分为 6 类:男士服装、女士服装、家具、电子产品、食品和玩具。电子商务领域数据¹⁾的统计信息汇总于表 1,表中数字为数据集的句子数量,其中训练集是以百万(M)为单位。使用 1.25M 的英中 LDC 平行句²⁾作为外领域数据。英文 LDC2007T07 和中文 LDC2009T27 分别用于训练 SMT 的语言模型。

表 1 电子商务领域的句子数量

Table 1 Number of sentences in e-commerce domain

| | | 男士 服装 | 女士 服装 | 家具 | 电子 产品 | 食品 | 玩具 | 总计 |
|-----|----|----------|----------|-------|----------|-------|-------|-------|
| 训练集 | 中文 | 0.70M | 0.90M | 0.45M | 0.77M | 0.59M | 0.31M | 3.72M |
| | 英文 | 0.62M | 0.67M | 0.54M | 0.53M | 0.70M | 0.64M | 3.73M |
| 验证集 | 中文 | 587 | 731 | 381 | 858 | 780 | 429 | 3766 |
| | 英文 | 591 | 388 | 410 | 387 | 462 | 391 | 2629 |
| 测试集 | 中文 | 889 | 1905 | 887 | 358 | 661 | 911 | 6948 |
| | 英文 | 747 | 653 | 903 | 635 | 1603 | 571 | 3867 |

(2) 训练设置

在 SMT 模型训练过程中,对于使用领域外数据训练的模型和使用领域内、外数据联合训练的模型,其语言模型均设置为 5-gram。在 NMT 模型训练过程中,不仅使用 Transformer^[17]模型框架,还使用 LSTM 作为模型框架。其中 Transformer 的编码器和解码器各为 6 层,词嵌入大小设置为

512,前馈网络的维数设置为 2048。在源语言和目标语言句子上分别设置大小为 64000 的 BPE^[18]。训练 LSTM 模型时,设置编码器和解码器的隐藏层维度和词嵌入维度大小均为 512。

基于对验证集的无监督调整,式(5)中的参数 α 在英文 \rightarrow 中文和中文 \rightarrow 英文两个方向上都设置为 0.9,NMT 相关的实验均在 Nvidia GTX1080TI 上进行。在所有训练设置中使用无监督调参来完全模拟无监督域适应中领域内平行数据的缺乏情景。整个训练耗时大约 5 天,其中基于 SMT 的互训练过程大约需要 50h,基于 NMT 的互训练过程大约需要 32h,拟知识蒸馏大约需要 24h。实验结果采用双语互译评估(Bilingual Evaluation Understudy, BLEU)^[19],目前 BLEU 评测已经被广泛应用于自然语言处理和机器翻译任务,本文使用 Moses 中 multi-bleu.perl 评测工具³⁾来获取实验翻译结果的 BLEU 值。

4.2 实验过程

(1) 基于 Transformer 的实验

在电子商务领域机器翻译中,本文将多种无监督领域适应方法作为基准系统,在此基础上做进一步的研究与分析。下面给出神经网络翻译模型,其均是 Transformer 模型。

SMT_{out} :在领域外平行语料上利用 Moses⁴⁾训练基于短语的 SMT 系统。

NMT_{out} :在领域外平行语料上训练 Transformer 的 NMT 系统⁵⁾。

$NMT_{out}+BT$ ^[2]:将反向翻译应用于 NMT_{out} ,构建伪平行域内数据,然后将领域内伪数据和领域外真实数据结合起来训练新的中英 NMT 模型。

$NMT_{out}+BT+FT$:在 $NMT_{out}+BT$ 的基础上做进一步拓展,其中包含了与反向翻译相反的前向翻译(forward translation),其中前向翻译是由真输入和伪输出组成的伪平行句对。

$NMT_{out}+Copy$ ^[1]:在领域内,直接将目标语言句子拷贝作为源句子,构建领域内的伪平行语料,然后将这份伪数据与领域外真实数据结合起来,以训练新的 NMT 模型。

$DAFE+BT$ ^[3]:在 NMT 编码器中的每一层添加领域感知特征嵌入和任务特定特征嵌入。该网络基于多任务学习框架,通过反向翻译生成的领域内伪数据和领域外真实数据进行训练。

$LEX_{unsupervised}$ ^[4]:使用无监督词典归纳,用所有领域的数据训练词嵌入,通过最近邻搜索获得构建领域内伪平行的词翻译,形成词典,使用该词典对领域内的单语数据逐词翻译,获得的翻译结果与输入配对成伪平行语料,然后将其与领域外真实平行语料结合起来,以训练新模型。

$LEX_{supervised}$ ^[4]:整个训练过程同 $LEX_{unsupervised}$ 一样,但

¹⁾ <https://github.com/nlp-anonymous/MECC>

²⁾ 英中数据集是由语言学数据联盟(Linguistic Data Consortium, LDC)语料库提取,包含 LDC2002E18, LDC2003E14, LDC2004T08, LDC2005T06

³⁾ <http://www.statmt.org/moses/?n=Moses.SupportTools>

⁴⁾ <http://www.statmt.org/moses/>

⁵⁾ <https://github.com/pytorch/fairseq>

LEX_{supervised} 用外领域词典作为种子词典进行监督词归纳。

表 2 列出了电商在中-英、英-中两个方向的各种对比方法以及 SA-UDA 的实验结果。其中 SA-UDA 包含互训练 SMT、互训练 NMT、互训练 SMT+NMT 以及互训练 SMT+

NMT+QKD。互训练 SMT+NMT+QKD 是综合训练,它结合了 SA-UDA 的 3 个组成部分,如第 2 节所述。互训练 SMT+NMT 是先进行互训练 SMT,然后按照综合训练中的方法逐步过渡到互训练 NMT。

表 2 电子商务产品翻译测试的 BLEU 值

Table 2 Experimental results on product translation test sets evaluated by BLEU

| | 男士服装 | 女士服装 | 家具 | 电子产品 | 食品 | 玩具 | 平均 |
|-----|-----------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 英-中 | SMT _{out} | 14.99 | 13.99 | 21.54 | 19.46 | 17.80 | 17.88 |
| | NMT _{out} | 11.16 | 10.77 | 15.76 | 10.10 | 15.76 | 13.30 |
| | NMT _{out} +BT | 12.87 | 11.28 | 15.10 | 17.54 | 17.49 | 14.76 |
| | NMT _{out} +BT+FT | 11.67 | 10.97 | 14.16 | 14.58 | 17.13 | 15.38 |
| | NMT _{out} +Copy | 13.63 | 15.87 | 15.36 | 15.43 | 17.78 | 18.19 |
| | DAFE+BT | 12.74 | 11.23 | 14.42 | 12.01 | 17.14 | 14.92 |
| | LEX _{unsupervised} | 16.16 | 12.79 | 13.50 | 14.57 | 19.88 | 16.47 |
| | LEX _{supervised} | 15.59 | 13.54 | 15.03 | 15.87 | 21.81 | 17.28 |
| | 互训练 SMT | 13.45 | 16.64 | 21.64 | 21.70 | 16.60 | 15.86 |
| | 互训练 NMT | 14.73 | 12.45 | 15.49 | 17.73 | 19.93 | 16.59 |
| | 互训练 SMT+NMT | 18.04 | 18.72 | 22.20 | 22.34 | 22.30 | 19.59 |
| | 互训练 SMT+NMT+QKD | 19.77 | 18.73 | 23.23 | 28.43 | 22.58 | 21.21 |
| 中-英 | SMT _{out} | 7.85 | 11.44 | 9.53 | 18.22 | 12.49 | 10.83 |
| | NMT _{out} | 7.34 | 7.00 | 7.77 | 6.68 | 12.66 | 12.02 |
| | NMT _{out} +BT | 11.94 | 10.87 | 15.14 | 11.61 | 18.92 | 21.27 |
| | NMT _{out} +BT+FT | 8.27 | 7.40 | 9.08 | 7.70 | 13.94 | 13.66 |
| | NMT _{out} +Copy | 12.48 | 16.02 | 14.97 | 15.78 | 18.27 | 19.42 |
| | DAFE+BT | 12.57 | 13.45 | 15.04 | 11.84 | 19.28 | 17.46 |
| | LEX _{unsupervised} | 11.96 | 14.15 | 13.16 | 12.96 | 17.25 | 22.17 |
| | LEX _{supervised} | 14.40 | 15.57 | 12.60 | 15.00 | 18.72 | 17.08 |
| | 互训练 SMT | 9.50 | 13.28 | 13.23 | 20.10 | 14.20 | 20.51 |
| | 互训练 NMT | 13.92 | 13.15 | 15.53 | 13.72 | 20.94 | 19.07 |
| | 互训练 SMT+NMT | 18.75 | 19.59 | 16.75 | 25.27 | 22.58 | 20.82 |
| | 互训练 SMT+NMT+QKD | 18.68 | 20.70 | 17.36 | 26.60 | 24.02 | 22.76 |

(2) 基于 LSTM 的实验

SA-UDA 与对偶学习^[7]相似。我们使用对偶学习¹⁾在电子商务领域构建翻译模型。下面给出用于对比的神经网络翻译模型,其均是 LSTM 模型。

D=NMT_{out}:在领域外平行语料训练基于 LSTM 的 NMT 系统²⁾,直接用于翻译领域内测试集。

D=NMT_{pseudo}:由 D=NMT_{out} 系统对领域内单语数据解码获得伪输出,将领域内的伪平行数据与领域外真实数据相结

合来训练新模型。

Dual+NMT:已训练完成的两个 D=NMT_{pseudo} 模型以及语言模型,通过强化学习过程对偶学习互相教导,迭代更新两个模型,直到收敛。

表 3 列出了电商在中-英、英-中两个方向的基于 LSTM 的各种对比方法以及基于 LSTM 的 SA-UDA 的实验结果,其中 SA-UDA 包含互训练 SMT+NMT、互训练 SMT+NMT+QKD。

表 3 基于 LSTM 的电子商务产品翻译测试 BLEU 值

Table 3 Experimental results on product translation test sets based on LSTM evaluated by BLEU

| | 男士服装 | 女士服装 | 家具 | 电子产品 | 食品 | 玩具 | 平均 |
|-----|-------------------------|--------------|--------------|--------------|--------------|--------------|--------------|
| 英-中 | D-NMT _{out} | 6.18 | 4.80 | 5.21 | 5.07 | 8.32 | 6.51 |
| | D-NMT _{pseudo} | 5.98 | 4.21 | 4.52 | 5.89 | 8.47 | 6.30 |
| | Dual+NMT | 7.74 | 5.62 | 5.83 | 8.71 | 10.62 | 8.09 |
| | 互训练 SMT+NMT | 10.54 | 10.33 | 12.74 | 14.15 | 11.67 | 10.02 |
| | 互训练 SMT+NMT+QKD | 12.24 | 10.06 | 11.10 | 13.61 | 13.98 | 12.23 |
| 中-英 | D-NMT _{out} | 5.21 | 5.14 | 6.81 | 3.06 | 9.28 | 7.95 |
| | D-NMT _{pseudo} | 10.28 | 11.65 | 10.84 | 8.39 | 15.62 | 12.22 |
| | Dual+NMT | 11.78 | 12.35 | 10.51 | 13.68 | 15.92 | 11.52 |
| | 互训练 SMT+NMT | 15.24 | 17.30 | 14.77 | 18.47 | 16.38 | 19.64 |
| | 互训练 SMT+NMT+QKD | 15.68 | 16.53 | 16.61 | 25.41 | 16.53 | 19.61 |

4.3 实验结果

从表 2 和表 3 可以看出,本文方法在电子商务产品翻译的两个方向上均十分有效。

(1) 基准系统

SMT_{out} 始终优于 NMT_{out},这表明直接在领域外数据上训练的 SMT 系统比 NMT 系统更适合电商产品翻译。将拷贝

¹⁾ <https://github.com/yistLin/pytorch-dual-learning>

²⁾ https://github.com/pcyin/pytorch_nmt.git

机制、反向翻译、DAFE 以及词表归纳技术应用于 NMT 时,从表中的实验结果可以观察到,这几种方法在英-中方向翻译性能均有所提升,但这些方法在中-英方向未能击败 SMT_{out} 。所有的 LEX 都比 $NMT_{out} + BT$ 好,这意味着通过统计进行词归纳的方法在电子商务中有一定的优势。由表 2 可见, $NMT_{out} + Copy$ 的性能优于 $NMT_{out} + BT$,这是因为在电子商务领域,双方语言中的品牌名称或产品属性通常都是英文的,如“adidas”,它总是出现在双方语言中。利用拷贝机制能够有效地保留这种电商文本特性,从而获得更好的性能。相较于 NMT_{out} 而言, $NMT_{out} + BT$ 在两个方向上均有所提升,但 $NMT_{out} + BT + FT$ 不能进一步提高性能。这是由于 $NMT_{out} + BT + FT$ 的训练过程中加入了前向翻译数据,其质量较差,会给训练带来负面影响,而 $NMT_{out} + BT$ 目标句始终是真实的,这表明真实目标更有利于训练。

由表 3 可以明显地看出,基于 LSTM 的单层网络翻译模型的效果远低于表 2 中基于 Transformer 网络结构的翻译模型,与 $D-NMT_{out}$ 和 $D-NMT_{pseudo}$ 相比, $Dual + NMT$ 充分利用了单语数据以及语言模型,使训练过程中的双向反馈信息得到增强。

(2) SA-UDA vs. 基准系统

在表 2 中互训练 NMT 在两个方向上均极大地提升了 NMT_{out} 的性能,而 SMT_{out} 仅在中-英方向较 SMT_{out} 有所提高,在英-中方向 SMT_{out} 是较强基线且很难超越。当把互训练 SMT 和互训练 NMT 结合起来时,可以明显地看出其在电子商务领域的两个方向上都显著超过了所有的基线。此外,互训练 SMT+NMT+QKD 进一步提高了翻译性能,在所有方向 and 所有产品类别中达到最佳(除了在中-英方向的男士服装类别略有下降)。具体表现为,互训练 SMT+NMT+QKD 在英-中方向较最强基线系统 SMT_{out} 提高了 4 个 BLEU 点,在中-英方向较最强基线系统 $NMT_{out} + Copy$ 提高了 5 个 BLEU 点以上。

表 3 中,在电子商务翻译中利用领域内数据的 $D-NMT_{pseudo}$, $Dual + NMT$ 和互训练 SMT+NMT 较 $D-NMT_{out}$ 都有显著

提升,此外,互训练 SMT+NMT 比对偶学习 $Dual + NMT$ 表现更佳,这一结果得益于互训练 SMT 在电商领域翻译性能较好。互训练 SMT+NMT+QKD 对大部分测试数据翻译进行了进一步改进,在英-中方向较最强基准系统提高了 4 个 BLEU 点以上,在中-英方向较最强基准系统提高了 5 个 BLEU 点以上。

5 评估与分析

在翻译中,Transformer 模型的表现比 LSTM 模型更好,本节选择基于 Transformer 模型的翻译结果进行评估与分析。首先,分别评估领域内的 SA-UDA 产生的词对和句子翻译的有效性,并分别将英-中、中-英的最强基准系统进行比较,然后分析在 SA-UDA 中各个步骤的增长趋势。

5.1 评估领域内词对

词级别的翻译采用领域内新词对的产生质量进行测评。首先根据人工翻译的测试集建立领域内的词对集合,然后使用词对齐工具 fastalign^[20] 来抽取领域外的对齐词,构建出领域外词对集合,最后利用上述的词对集合排除领域外语料库词对,剩余的对齐词集合构成新词对的参考集合。

利用该参考集,可以评估由 SA-UDA 产生的新单词对。表 4 列出了英-中、中-英最强基准系统以及 SA-UDA 各部分的精确率、召回率和 F 值。同构建新词对的参考集合的过程一样,对于 SA-UDA 的每种方法进行新词对集合提取。将领域内的输出与领域外的平行语料连接起来,运用 fastalign 来抽取该方法翻译结果的对齐词,去掉领域外词对来产生词对集合,并根据新词对的参考集合来评估该集合。互训练 SMT 在英-中和中-英两个方向的召回率分别为 64.48% 和 78.53%,远超于基准系统的召回率。而在精确率上,英-中方向的最强基准系统 SMT_{out} 偏高,中-英方向最强基准系统 $NMT_{out} + Copy$ 与 SA-UDA 各步骤之间的差距较小。根据精确率、召回率和综合指标 F 值,SA-UDA 产生的新词对是有效的,其中互训练 SMT+NMT+QKD 在所有方法中表现最好,该方法产生的大部分词对与新词对的参考集合一致。

表 4 领域内单词配对和句子级翻译的评估结果

Table 4 Evaluations of both word pairs and sentence level translations in in-domain

| | | 单词配对 | | | 句子级翻译 | |
|-----|--------------------|--------------|--------------|--------------|-----------------------|-----------------------|
| | | 精确率/% | 召回率/% | F 值/% | PPL _{n-gram} | PPL _{neural} |
| 英-中 | SMT_{out} | 63.16 | 46.27 | 53.41 | 1772.91 | 3528.64 |
| | 互训练 SMT | 42.18 | 64.48 | 51.00 | 3227.14 | 6090.11 |
| | 互训练 NMT | 57.74 | 45.13 | 50.55 | 1055.34 | 849.25 |
| | 互训练 SMT+NMT | 57.46 | 68.47 | 62.49 | 887.89 | 715.91 |
| | 互训练 SMT+NMT+QKD | 57.60 | 92.77 | 71.07 | 838.56 | 671.61 |
| 中-英 | $NMT_{out} + Copy$ | 76.17 | 70.01 | 72.96 | 1003.87 | 830.64 |
| | 互训练 SMT | 74.85 | 78.53 | 76.65 | 1795.06 | 1608.85 |
| | 互训练 NMT | 75.75 | 58.61 | 66.08 | 1334.63 | 622.98 |
| | 互训练 SMT+NMT | 76.92 | 80.82 | 78.82 | 1331.80 | 541.37 |
| | 互训练 SMT+NMT+QKD | 77.00 | 81.27 | 79.08 | 837.86 | 537.33 |

5.2 领域内句子级翻译的困惑度

困惑度(Perplexity, PPL)是在自然语言处理领域中,衡量语言模型好坏的指标。低困惑度的概率分布模型或概率模

型能更好地预测样本。为了更好地判断得到的翻译文本是否保留原始正确数据风格,通过句子级别翻译文本的 PPL 进行衡量。表 4 所列的 PPL 包含了基于 n-gram 语言模型¹⁾及神

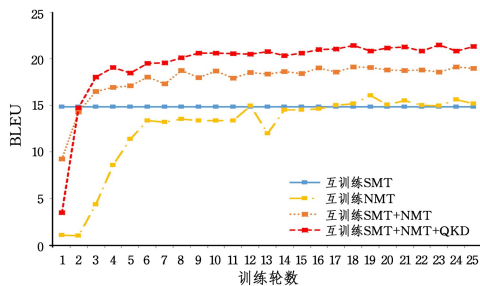
¹⁾ <https://kheffield.com/code/kenlm/>

经语言模型,通过测试集的句子级翻译结果计算得出。这两种语言模型都是基于领域外和领域内数据的真实单语数据。

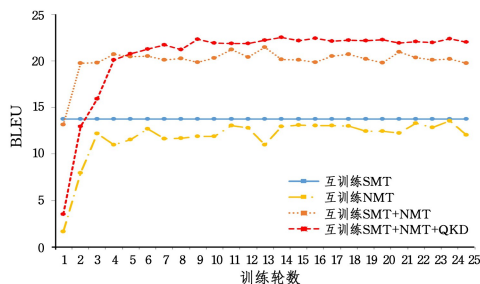
如表4的“句子级翻译”一栏所示,SA-UDA各部分与最强基准系统相比其PPL值显著降低。互训练SMT+NMT+QKD再次成为所有方法中表现最好的。由于互训练SMT使用提取的平行短语进行翻译,其输出不如NMT流畅,导致PPL更高。从表4可以看出,虽然电子商务领域的文本是非正式的,具有较高的PPL值,但SA-UDA仍然可以发现电子商务领域的内在规律,并显著降低句子翻译的PPL。

5.3 SA-UDA 趋势

图2给出了SA-UDA的验证集BLEU趋势曲线,并分别展示了每个方法随着训练的迭代在验证集上的表现。可以看出,SA-UDA在真实和伪平行语料库上训练时,收敛速度很快。其翻译性能在前3~5个轮次显著提升,并且通常在10个轮次后变得稳定。将互训练SMT最后一轮获得的翻译结果的BLEU值作为基准线,在图中由一条直线表示。从图2互训练SMT+NMT以及互训练SMT+NMT+QKD的曲线可以看出,两者均有效提升了翻译性能。



(a) 英-中



(b) 中-英

图2 SA-UDA的验证集BLEU趋势曲线

Fig. 2 BLEU convergence curve of SA-UDA on validation set

6 拓展应用

SA-UDA能有效改善电商产品翻译质量,其在无监督领域适应中也非常有意义。因此,我们不仅在电子商务领域的机器翻译上进行相关探索和研究,还将SA-UDA应用于公开的低资源领域数据集Ted, Law和Medical OPUS,这些数据集被视为领域内数据,WMT-14被视为领域外数据。

表5分别列出了Ted, Law和Medical OPUS上各种对比方法以及SA-UDA的实验结果,其中SMT_{out}与第3节中SMT_{out}设置相同的参数,在外领域WMT-14数据上训练得到翻译模型。NMT相关实验室基于Transformer模型,其实验过程的设置与DAFE实验相同。互训练NMT+QKD中,

式(5)中的 α 对于英-德、德-英均设置为0.1。

表5 SA-UDA应用于不同领域的BLEU值

Table 5 Experiment results of SA-UDA applied to different domain evaluated by BLEU

| | 德-英 | | | 英-德 | | |
|------------------------|-------|-------|-------|-------|-------|-------|
| | Ted | Law | Med | Ted | Law | Med |
| SMT _{out} | 25.70 | 24.74 | 24.32 | 18.95 | 16.62 | 18.46 |
| NMT _{out} | 28.15 | 24.61 | 26.75 | — | — | — |
| Back-DAFE+DAFE | 34.89 | 31.46 | 38.79 | — | — | — |
| NMT _{out} +BT | 32.87 | 38.05 | 40.49 | 25.5 | 28.59 | 35.82 |
| 互训练NMT | 33.52 | 40.38 | 44.70 | 28.86 | 30.25 | 38.20 |
| 互训练NMT+QKD | 33.73 | 40.64 | 45.01 | 29.13 | 30.91 | 38.95 |

从表5中SMT_{out}和NMT_{out}的结果可知,NMT模型在这3个领域的效果都优于SMT,因此本文在这3个领域的数据上实现了互训练NMT。实验结果表明,互训练NMT在迭代训练中可以利用领域内和领域外数据,拟知识蒸馏的方法能够进一步提高翻译性能。由于三者在职域内的风格较为统一,因此在拟知识整理过程中,给予老师模型较小的权重。

结束语 在电子商务中翻译产品信息有两大挑战,一个挑战是无监督的领域自适应,它将在资源丰富的领域外训练的翻译模型适配到领域内的数据上,即电子商务领域,在该领域中没有平行训练数据并且难以构建平行训练数据。另一个挑战是解决领域内非平行数据的两种语言之间的风格差异。为了解决这两个问题,本文提出了风格感知的无监督领域适应方法,即SA-UDA。为了有效利用领域内非平行数据和领域外平行数据,在SA-UDA中应用了互训练过程。同时为了使翻译过程中维持各自的风格,提出了拟知识蒸馏。在电子商务产品翻译上的实验表明,所提方法的性能明显优于各种无监督领域适应方法,在单词级和句子级翻译上都得到了显著的改进。此外,本文提出的方法不仅适用于电子商务领域,在其他领域的无监督领域适应翻译中也具有积极的作用。本文在中英语言的电子商务领域机器翻译上进行了相关的研究探索,在未来工作中我们将在拓展电子商务领域的其他语言的同时,进一步研究情感分析^[21]及语义角色标注^[22]等多任务学习方法,以提升电商机器翻译性能。

参考文献

- [1] CURREY A, BARONE A V M, HEAFIELD K. Copied Monolingual Data Improves Low-Resource Neural Machine Translation[C]// Proceedings of the Second Conference on Machine Translation. Denmark: Association for Computational Linguistics, 2017: 148-156.
- [2] SENNRICH R, HADDOW B, BIRCH A. Improving Neural Machine Translation Models with Monolingual Data[C]// Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics. Germany: Association for Computational Linguistics, 2016: 86-96.
- [3] DOU Z Y, HU J J, ANASTASOPOULOS A, et al. Unsupervised domain adaptation for neural machine translation with domain-aware feature embeddings[C]// Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Lan-

- guage Processing. Hong Kong: Association for Computational Linguistics, 2019: 1417-1422.
- [4] HU J, XIA M, NEUBIG G, et al. Domain Adaptation of Neural Machine Translation by Lexicon Induction[C]// Proceedings of the 57th Conference of the Association for Computational Linguistics. Italy: Association for Computational Linguistics, 2019: 2989-3001.
- [5] SHEN Y, LEONARD D, PAVEL P, et al. Word-based Domain Adaptation for Neural Machine Translation[C]// Proceedings of the International Workshop on Spoken Language Translation. Belgium, 2019.
- [6] ZHANG Z, LIU S, LI M, et al. Joint Training for Neural Machine Translation Models with Monolingual Data[C]// Proceedings of the AAAI Conference on Artificial Intelligence. USA: AAAI Press, 2018: 555-562.
- [7] HE D, XIA Y, QIN T, et al. Dual Learning for Machine Translation[J]. Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems, 2016(12): 820-828.
- [8] ZHEN Y, WEI C, FENG W, et al. Unsupervised Domain Adaptation for Neural Machine Translation[C]// 24th International Conference on Pattern Recognition. China: IEEE Computer Society, 2018: 338-343.
- [9] ZHENG Z, ZHOU H, HUANG S, et al. Mirror-generative neural machine translation[C]// 8th International Conference on Learning Representations. Ethiopia: ICLR, 2020.
- [10] NIU X, RAO S, CARPUAT M. Multi-task neural models for translating between styles within and across languages[C]// Proceedings of the 27th International Conference on Computational Linguistics. USA: Association for Computational Linguistics, 2018: 1008-1020.
- [11] KOEHN P, OCH F J, MARCU D. Statistical Phrase-Based Translation[C]// Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics. Canada: Association for Computational Linguistics, 2013: 48-54.
- [12] ARTETXE M, LABAKA G, AGIRRE E. An Effective Approach to Unsupervised Machine Translation[C]// Proceedings of the 57th Conference of the Association for Computational Linguistics. Italy: Association for Computational Linguistics, 2019: 194-203.
- [13] ARTETXE M, LABAKA G, AGIRRE E. Unsupervised Statistical Machine Translation[C]// Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing. Belgium: Association for Computational Linguistics, 2018: 3623-3642.
- [14] LUONG T, PHAM H, MANNING C D. Effective Approaches to Attention-based Neural Machine Translation[C]// Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing. Lisbon, Portugal: Association for Computational Linguistics, 2015: 1412-1421.
- [15] ASHISH V, NOAM S, NIKI P, et al. Attention is all you need [C]// Advances in Neural Information Processing Systems 30. USA: NIPS, 2017: 5998-6008.
- [16] YOON K, ALEXANDER M R. Sequence-Level Knowledge Distillation[C]// Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing. USA: Association for Computational Linguistics, 2016: 1317-1327.
- [17] SUTSKEVER I, ORIOL V, QUOC V L. Sequence to sequence learning with neural networks[C]// Advances in Neural Information Processing Systems 27. Canada: NIPS, 2014: 3104-3112.
- [18] SENNRICH R, HADDOW B, BIRCH A. Neural Machine Translation of Rare Words with Subword Units[C]// Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics. Germany: Association for Computational Linguistics, 2016: 1715-1725.
- [19] PAPANENI 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. USA: Association for Computational Linguistics, 2002: 311-318.
- [20] CHRIS D, VICTOR C, NOAH A S. A Simple, Fast, and Effective Reparameterization of IBM Model 2[C]// Proceedings of the North American Chapter of the Association for Computational Linguistics. USA: Association for Computational Linguistics, 2013: 644-648.
- [21] HU D M, ZHU C G, HU C, et al. Multilingual Text Emotional Analysis with Pre-trained Model and Attention Mechanism[J]. Journal of Chinese Mini-Micro Computer Systems, 2020, 41(2): 278-284.
- [22] QIAO B W, LI J H. Neural Machine Translation Combining Source Semantic Roles[J]. Computer Science, 2020, 47(2): 163-168.



NING Qiu-yi, born in 1995, postgraduate. Her main research interests include machine translation and domain adaptation.



DUAN Xiang-yu, born in 1976, Ph. D, professor. His main research interests include machine translation and cross-language information processing.