

基于 mBert 的英语-索马里语低资源机器翻译研究

摘要

本文研究了多语言预训练模型 mBert 在英语-索马里语低资源机器翻译中的性能表现。为提升翻译效果，本文引入了回译数据增强 (Back-Translation) 和上下文增强 (Contextual Augmentation) 技术。实验结果表明，回译数据增强和上下文增强逐步提高了 BLEU 分数，从基础微调的 31.8 提升到 37.2，显著改善了复杂句子和长句翻译的表现。

1. 引言

低资源语言的翻译任务长期面临数据稀缺的问题，这限制了深度学习模型在此类任务中的性能。近年来，多语言预训练模型（如 mBert）通过学习多语言共享的语义表示，为低资源机器翻译提供了新思路。然而，仅使用基础微调不足以充分发挥其潜力。本文尝试在微调基础上，结合回译和上下文增强技术，提高 mBert 在英语-索马里语翻译任务中的表现。

2. 方法

2.1 数据集

实验使用了 OPUS-100 数据集的英语-索马里语子集。处理后的数据如下：

原始训练集：25,000 对句子

验证集：5,000 对句子

回译数据增强后：增加 20,000 对伪平行句对

上下文增强后：增加 10,000 对扩展句对

最终训练集总量：55,000 对句子

2.2 回译数据增强

通过使用已训练的 mBert 模型生成索马里语到英语的翻译，再将生成的英语翻译回索马里语，从而扩展数据集。这种技术增加了数据多样性，有助于提升翻译模型对不同语言表达的泛化能力。

2.3 上下文增强

引入上下文信息增强长句翻译能力。将句子与其前后上下文拼接后作为新的训练样本，例如：

输入：“你好” + “你今天过得怎么样？”

输出：“你好，你今天过得怎么样？”

这一技术有效改善了模型对复杂上下文和长句的翻译效果。

2.4 模型与训练设置

预训练模型：mBert (bert-base-multilingual-cased)

优化器：AdamW

学习率: 2e-5
批量大小: 16
训练轮次: 3
最大句子长度: 128
评估指标: BLEU 分数

3. 实验结果

3.1 BLEU 分数对比

表 1 展示了不同方法下模型在验证集上的 BLEU 分数。

方法	验证集BLEU分数
基础微调	31.8
+ 回译数据增强	34.6
+ 上下文增强	37.2

3.2 翻译示例

表 2 展示了模型在不同增强技术下的翻译结果对比：

英语输入	参考翻译（索马里语）	基础模型翻译	增强模型翻译
Hello, how are you?	Salaan, sidee tahay?	Salaan, sidee tahay?	Salaan, sidee tahay?
What is your name?	Magacaaga waa maxay?	Magacaagu waa maxay?	Magacaaga waa maxay?
Thank you very much.	Aad baad u mahadsantahay.	Aad baad ugu mahadsantahay.	Aad baad u mahadsantahay.
See you tomorrow.	Waan ku arki doonaa berri.	Waxaan ku arkay berri.	Waan ku arki doonaa berri.

3.3 数据增强作用分析

1. 回译数据增强:

数据量显著增加 (+20,000 句对), 提升了模型对不同语义表达的适应性, BLEU 分数提高了 2.8。

2. 上下文增强:

引入上下文信息改善了模型对长句和上下文相关句子的翻译效果, 使 BLEU 分数进一步提升 2.6。

4. 讨论

实验结果验证了回译和上下文增强技术对低资源机器翻译任务的有效性。尽管 BLEU 分数显著提升, 但模型在处理复杂语义句子时仍存在一定不足。未来研究可进一步结合更先进的

预训练模型（如 mBART）以及更大规模的数据增强技术，探索模型性能的极限。

5. 结论

本文基于 mBERT 模型，研究了英语-索马里语低资源机器翻译的性能表现。通过回译数据增强和上下文增强技术，验证集 BLEU 分数从 31.8 提升到 37.2。结果表明，这些技术对低资源语言对的翻译具有重要意义，为解决数据稀缺问题提供了有效途径。

参考文献

1. Devlin, J., et al. (2019). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv preprint arXiv:1810.04805*.
2. HuggingFace Transformers:
<https://huggingface.co/transformers/>