

Обучение машинного перевода без параллельных текстов

Саттаров Тагир

Московский физико-технический институт

Курс: Численные методы обучения по прецедентам
(практика, В.В. Стрижов)/Группа 674, весна 2019

Цель исследования

Цель работы

Улучшить качество перевода на низкоресурсных парах языков

Проблема

При переводе некоторых пар языков необходимо учитывать морфологическое строение слов, что существенно усложняет задачу машинного перевода

Решение

Предлагается построить общее скрытое пространство для обоих языков. Предлагается использовать T2T модель кодировщика-декодировщика вместо используемой seq2seq

- G Lample, A. Conneau, L. Denoyer, M. Ranzato. Unsupervised Machine Translation Using Monolingual Corpora Only. Conference paper at ICLR 2018
- A. Conneau, G. Lample, M. Ranzato, L. Denoyer, and H. Jegou. Word translation without parallel data. ICLR 2018
- A. Zadeh, P. P. Liang, S. Poria, P. Vij, E. Cambria, L-P Morency. Multi-Attention Recurrent Network for Human Communication Comprehension. AAAI-18

Постановка задачи

Дано

Заданы \mathcal{D}_{src} и \mathcal{D}_{tgt} , где \mathcal{D}_{src} – выборка предложений на английском языке, а \mathcal{D}_{tgt} – на французском. Также задана \mathcal{D}_{valid} – предложения на английском, переведенные на французский для проверки качества перевода.

Используются модели

Дискриминатора \mathbf{d} – по вектору из скрытого пространства \mathbf{Z} предсказывает исходный язык

Кодировщика \mathbf{f} – переводит предложение из одного из языков в скрытое пространство

Декодировщика \mathbf{g} – восстанавливает предложение по вектору из скрытого пространства

Общая схема решения

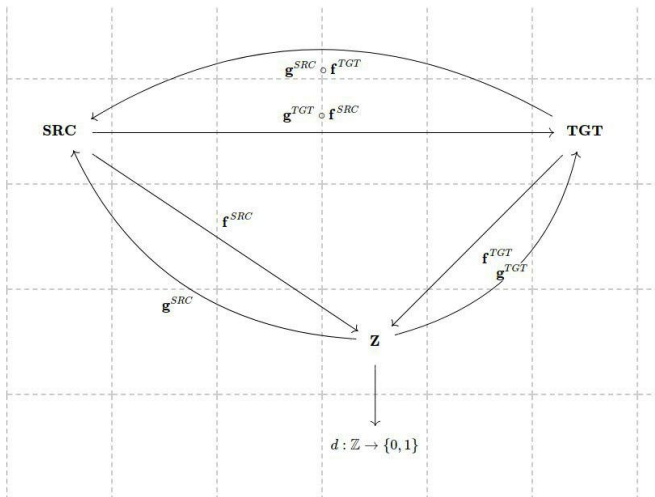


Рис.: Auto-encode

Требования к кодировщику и декодировщику

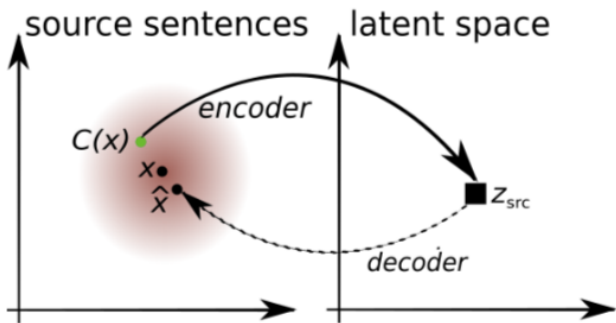


Рис.: Auto-encode

Требования к кодировщику и декодировщику

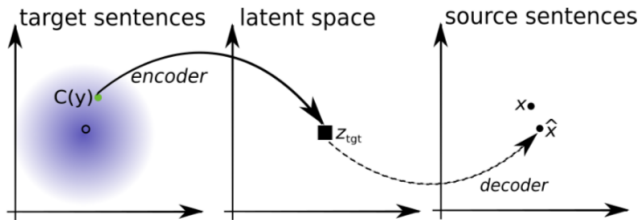


Рис.: Cross-domain

Дано

Один из наиболее распространенных критериев качества – BLEU(bilingual evaluation understudy).

Определение

N -грамма – последовательность из N подряд идущих слов.

Метод

Рассматривается несколько вариантов профессионального перевода и смотрится количество N -грамм лежащих и в машинном переводе, и в данной выборке.

Критерий качества перевода

Пример

Один из наиболее распространенных критериев качества – BLEU(bilingual evaluation understudy).

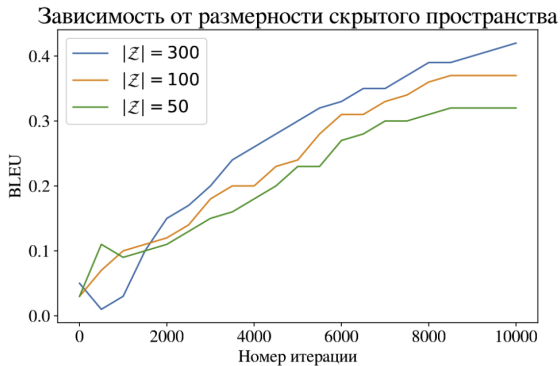
Определение

N -грамма – последовательность из N подряд идущих слов.

Метод

Рассматривается несколько вариантов профессионального перевода и смотрится количество N -грамм, лежащих и в машинном переводе, и в данной выборке.

Вычислительный эксперимент



Пример

Один из наиболее распространенных критериев качества – BLEU(bilingual evaluation understudy).

Определение

N -грамма – последовательность из N подряд идущих слов.

Метод

Рассматривается несколько вариантов профессионального перевода и смотрится количество N -грамм, лежащих и в машинном переводе, и в данной выборке.