

Masterarbeit im Fach Informatik
RHEINISCH-WESTFÄLISCHE TECHNISCHE HOCHSCHULE AACHEN
Lehrstuhl für Informatik 6
Prof. Dr.-Ing. H. Ney

Unsupervised Learning of Neural Network Lexicon and Cross-lingual Word Embedding

20. September 2018

vorgelegt von:
Autor Jiahui Geng
Matrikelnummer 365655

Gutachter:
Prof. Dr.-Ing. H. Ney
Prof. B. Leibe, Ph. D.

Betreuer:
M.Sc. Yunsu Kim

Erklärung

Geng, Jiahui

Name, Vorname

365655

Matrikelnummer (freiwillige Angabe)

Ich versichere hiermit an Eides Statt, dass ich die vorliegende ~~Arbeit/Bachelorarbeit/~~
Masterarbeit* mit dem Titel

Unsupervised Learning of Neural Network Lexicon and Cross-lingual Word Embedding

selbständig und ohne unzulässige fremde Hilfe erbracht habe. Ich habe keine anderen als die angegebenen Quellen und Hilfsmittel benutzt. Für den Fall, dass die Arbeit zusätzlich auf einem Datenträger eingereicht wird, erkläre ich, dass die schriftliche und die elektronische Form vollständig übereinstimmen. Die Arbeit hat in gleicher oder ähnlicher Form noch keiner Prüfungsbehörde vorgelegen.

Aachen, 20. September 2018

Ort, Datum

Unterschrift

*Nichtzutreffendes bitte streichen

Belehrung:

§ 156 StGB: Falsche Versicherung an Eides Statt

Wer vor einer zur Abnahme einer Versicherung an Eides Statt zuständigen Behörde eine solche Versicherung falsch abgibt oder unter Berufung auf eine solche Versicherung falsch aussagt, wird mit Freiheitsstrafe bis zu drei Jahren oder mit Geldstrafe bestraft.

§ 161 StGB: Fahrlässiger Falscheid; fahrlässige falsche Versicherung an Eides Statt

(1) Wenn eine der in den §§ 154 bis 156 bezeichneten Handlungen aus Fahrlässigkeit begangen worden ist, so tritt Freiheitsstrafe bis zu einem Jahr oder Geldstrafe ein.

(2) Straflosigkeit tritt ein, wenn der Täter die falsche Angabe rechtzeitig berichtigt. Die Vorschriften des § 158 Abs. 2 und 3 gelten entsprechend.

Die vorstehende Belehrung habe ich zur Kenntnis genommen:

Aachen, 20. September 2018

Ort, Datum

Unterschrift

Abstract

In recent years, cross-lingual representation of words enable us to explore the structures of different languages in a shared space. Based on word similarity in a cross-lingual space. we can realize bilingual dictionary induction, information retrieval and knowledge transfer between resource-rich and resource-lean languages.

In this thesis, I propose a novel training method of cross-lingual word embedding which is data-driven and supported by the language model. We better utilize the context information to find the We compare the results with that using BPE as the translation unit in our unsupervised model, find BPE does not improve the performance.

Contents

| | |
|--|-----------|
| Abstract | v |
| 1 Introduction | 1 |
| 1.1 Related Work | 2 |
| 1.1.1 Word Embedding | 2 |
| 1.1.2 Unsupervised Machine Translation | 3 |
| 1.2 Outline | 4 |
| 1.3 Notation | 4 |
| 2 Machine Translation | 5 |
| 2.1 Statistical Machine Translation | 5 |
| 2.1.1 Word-Based Model | 5 |
| 2.1.2 Phrase-Based Model | 8 |
| 2.1.3 Decipherment | 8 |
| 2.2 Neural Machine Translation | 9 |
| 2.2.1 Encoder-Decoder Framework | 10 |
| 2.2.2 Attention Mechanism | 11 |
| 2.2.3 Transformer | 13 |
| 2.2.4 Neural Unsupervised Machine Translation | 15 |
| 3 Word Embedding | 19 |
| 3.1 Monolingual Embedding | 19 |
| 3.1.1 CBOW and Skip-gram Model | 19 |
| 3.1.2 Enriching Embedding with Subword Information | 21 |
| 3.2 Cross-lingual Word Embedding | 21 |
| 3.2.1 Training Algorithm | 21 |
| 4 Unsupervised Learning of Cross-lingual Word Embedding | 27 |
| 4.1 Initialization | 27 |
| 4.2 Vocabulary-based Training | 27 |
| 4.2.1 Iterative Procrustes Analysis | 27 |
| 4.2.2 Adversarial Training | 27 |
| 4.3 Corpus-based Training | 27 |

| | | |
|----------|--|-----------|
| 5 | Sentence Translation | 29 |
| 5.1 | Context-aware Beam Search | 29 |
| 5.1.1 | Language Model | 29 |
| 5.1.2 | Beam Search | 30 |
| 5.2 | Denoising Neural Network | 30 |
| 5.2.1 | Denoising Auto-encoder | 30 |
| 5.2.2 | Noise Model | 31 |
| 6 | Experiments | 35 |
| 6.1 | Corpus Statistics | 35 |
| 6.2 | Word Translation | 36 |
| 6.2.1 | Baseline | 36 |
| 6.2.2 | Data-driven Approach | 36 |
| 6.3 | Sentence Translation | 38 |
| 6.3.1 | Comprehensive Table | 38 |
| 6.3.2 | BPE vs Word | 38 |
| 6.3.3 | Artificial Noise | 39 |
| 6.3.4 | Phrase Embedding | 39 |
| 6.3.5 | Vocabulary Cutoff in Translation | 40 |
| 7 | Conclusion | 41 |
| A | Appendix | 43 |
| | List of Figures | 45 |
| | List of Tables | 47 |
| | Bibliography | 49 |

Chapter 1

Introduction

Building a good machine translation system requires a enormous collection of parallel data. NMT systems often fail when the training data is not enough. However parallel corpora are costly because they require huge human labor, time and expertise of corresponding languages. Several approaches has been proposed to alleviate this issue, for instance, triangulation or semi-supervised learning techniques; these systems however still need a strong cross-lingual signal. Unsupervised machine translation uses only monolingual data of the source and target languages to train the model and such monolingual corpora are readily available.

Recently cross-lingual word embedding draws more and more attention since it helps to exploit the multilingual semantic information. For instance, it enables computation of cross-lingual word similarities, which can be directly applied to bilingual dictionary induction or cross-lingual information retrieval. Several methods are proposed to learn the cross-lingual without any parallel data. This correspondingly inspires work on unsupervised machine translation.

In this thesis, I first make a survey of the current training algorithm of cross-lingual word embedding, find the significant training skills among them and their respective advantages and disadvantages. Then I propose a novel data-based unsupervised training method, the experiments the method can achieve considerable results in comparison with the start-of-the-art method. I further explore the application of the cross-lingual word embeddings in the machine translation domain. I develop a total unsupervised machine translation system starting from the simple word-by-word translation principle, to over the shortcomings to such more, I design the context-aware beam search to find the best translation candidate. To handle to re-ordering, I implement a denoising network with artificial noises, which are to mimic the true noise in the word-by-word translation, The results demonstrate that such simple but efficient translation system can surpass the most unsupervised neural machine translation system with costly iterative training.

1.1 Related Work

1.1.1 Word Embedding

Traditional language processing systems treat words as discrete atomic symbols, which assign for each word a specific ID number. Such encodings are arbitrary; it does not provide any information about the relations that may exist between individual symbols. In comparison, distributed representations of words in a high dimensional continuous space help learning algorithm to achieve better performance in natural language processing. According to the distributional hypothesis, similar words tend to occurs with semantic neighbors, having similar word representations. At the beginning, in order to explore word similarity, word embedding training involves dense matrix multiplication and complex matrix analysis.

Mikolov et al. [2013a] propose two models for learning word embeddings with neural networks, namely skip-gram and continuous bag of words (CBOW) models. Neural architectures makes the training extremely efficient. In further work, Mikolov et al. [2013c] discuss about learning the embedding for phrases. Pennington et al. [2014] propose a regression model which combines the global matrix factorization and local context window methods, making the training process more interpretable. Bojanowski et al. [2016] represent each word as a bag of character n -grams, and assign for each character n -gram a distinct embedding. In this way, they can exploit subword information in word embedding training.

Mikolov et al. [2013b] again notice that continuous embedding spaces exhibit similar structures across languages, even when for distant language pairs, for example, English-Vietnamese. With separately learned word embeddings they learn a linear mapping from source embedding space to target embedding space. To do this, they employ a parallel vocabulary of five thousand words as anchor points and evaluate their approach on a word translation task. Xing et al. [2015] show that results can be improved by enforcing an orthogonal constraint on the linear mapping.

Recently, Artetxe et al. [2017a] raise an iterative method that aligns the word embedding spaces gradually. In their work, the vocabulary size for learning can be reduced to 25 word pairs and even with Arabic numerals. The performance is comparable to the learning with large vocabulary dictionary Cao et al. [2016] propose a distribution-based model, which is a modified CBOW model minimizing the distribution dissimilarity between source and target embeddings; here, distribution information refers the mean and variance. Zhang et al. [2017] put forward a method using adversarial training without any parallel data. The discriminator tries to distinguish if the embedding is from the source side, the generator aims to learn the mapping from source embedding space to the target one. These methods are totally unsupervised but the performances are far worse than the supervised training. Conneau et al. [2017] simplifies the adversarial training structure with dif-

ferent loss functions for the generator and discriminator, improving the performance significantly. They also propose to use cross-domain similarity local scaling (CSLS) to handle hubness and a validation criterion for unsupervised model selection.

1.1.2 Unsupervised Machine Translation

As mentioned previously, the lack of parallel corpora motivates people to use monolingual data to improve the machine translation system. Some researchers use triangulation techniques (Cohn and Lapata [2007]) and semi-supervised approaches (Cheng et al. [2016]) to alleviate this issue. But these methods still require parallel corpora.

As to fully unsupervised machine translation, which only takes advantage of monolingual data, Ravi and Knight [2011] first consider it as a deciphering problem, where the source language is considered as ciphertext. They also propose iterative expectation-maximization (EM) method and Bayesian decipherment to train this unsupervised model. To solve the bottleneck of such model – mainly the huge memory required to store the candidates search space, – Nuhn et al. [2012] limit search candidates according to the word similarity. Nuhn and Ney [2014] limit the search space by using and preselection beam search. For the same reason, Kim et al. [2017] enforce the sparsity with a simple thresholding for the lexicon, also initializing the training with word classes to efficiently boost the performance. Although initially not based on distributional semantics, recent studies show that the use of word embeddings can bring significant improvement in statistical decipherment Duong et al. [2016].

For neural machine translation, He et al. [2016] show a bidirectional, iterative way of effectively using monolingual data. More concretely, they train two agents to translate in opposite directions (e.g. French \rightarrow English and English \rightarrow French) and make them teach each other through a reinforcement learning process. While promising, this approach still requires a small parallel corpora for a warm start. Artetxe et al. [2017b] and Lample et al. [2017] propose two bi-directional unsupervised machine translation model which totally rely only on monolingual corpora in each language. Both models need to use cross-lingual word embedding to initialize the MT system and train the sequence-to-sequence system with denoising autoencoder. They turn the unsupervised training into a supervised one by introducing back-translation techniques. The most important principle for these two works is the shared representations. Sentences from different languages are encoded into a shared space and then translated into specific languages with different decoders.

1.2 Outline

The remainder of this thesis is structured as follows. In Section 1.3 we introduce the notations for this thesis. Chapter 2 introduces the development of machine translation systems from statistical models to neural models. Some basic techniques and principles for unsupervised machine translation are also explained. Chapter 3 gives a survey of training and applications details of cross-lingual word embedding. Chapter 4 discusses the unsupervised learning of cross-lingual word embedding and our efforts on data-driven training. We describe our unsupervised machine translation system in Chapter 5, which contains mainly context-aware search with the help of language model and denoising autoencoder method aimed at reordering. We demonstrate experimental results of cross-lingual word embedding model and unsupervised machine system in Chapter 6, comparing the performance with the state-of-the-art model. In Chapter 7, we summarize our work and draw conclusions.

1.3 Notation

In this thesis, we use the following notations:

- Source sentence: $f_1^J := f_1 \cdots f_j \cdots f_J$
- Target sentence: $e_1^I := e_1 \cdots e_i \cdots e_I$
- Single character for each language: f, e
- Source and target word embedding: \mathbf{f}, \mathbf{e}
- Source and target embedding matrices: \mathbf{F}, \mathbf{E}
- Source and target Corpora: \mathcal{F}, \mathcal{E}
- Mapping default from source to target embedding space: W
- Mapping with specific direction: $W_{f \rightarrow e}, W_{e \rightarrow f}$
- Internal states of decoder: \mathbf{s} and \mathbf{s}_i
- Internal states of encoder: \mathbf{h} and \mathbf{h}_i
- Context vectors in attention mechanism: \mathbf{c} and \mathbf{c}_i

Chapter 2

Machine Translation

This chapter describes the classical or start-of-the-art machine translation (MT) models from statistical MT to neural MT. The related research work on unsupervised MT is also included.

2.1 Statistical Machine Translation

Statistical MT had achieved success in the early era of MT, until neural MT emerged. The initial statistical models for MT are based on words as atomic units that may be translated, inserted, dropped and reordered. In statistical MT, we use both a translation model and a language model which encourages fluent output. Later, the statistical MT prefers to use translation of phrases as atomic units. These phrases are any contiguous sequences of words and not necessarily linguistic entities. In this approach, the input sentence is broken up into a sequence of phrases and are mapped one-to-one to output phrases, which may be reordered.

2.1.1 Word-Based Model

Noisy-channel model

Noisy-channel model is based on the notion of a noisy channel from Shannon's information theory, which had been applied to many language processing problems. Assuming that the source sentence is a distorted message omitted from the target sentence, we have a model on how the message is distorted (translation model $Pr(f_1^J|e_1^I)$) and also a model on which original messages are probable (language

model $Pr(e_1^I)$. Their task is to find the best translation \hat{e}_1^J for an input foreign sentence.

$$\begin{aligned}\hat{e}_1^J &= \operatorname{argmax}_{e_1^I} \{Pr(e_1^I | f_1^J)\} \\ &= \operatorname{argmax}_{e_1^I} \frac{Pr(f_1^J | e_1^I) Pr(e_1^I)}{Pr(f_1^J)} \\ &= \operatorname{argmax}_{e_1^I} \{Pr(f_1^J | e_1^I) \cdot Pr(e_1^I)\}\end{aligned}$$

The task of word alignment is an artifact of word-based MT. Alignment models target at the reordering problem for word-based translation.

IBM Models

Based on the noisy-channel model, IBM word-based translation makes the model more complicated by adding submodels. Starting from lexical translation, absolute alignment model, fertility etc. are added step by step.

The advances of the five IBM models are:

- IBM Model 1: lexical translation;
- IBM Model 2: adds absolute alignment model;
- IBM Model 3: adds fertility model;
- IBM Model 4: adds relative alignment model;
- IBM Model 5: fixes deficiency.

Alignment model

This is an explicit model for reordering words in a sentence. More often than not, words that follow each other in one language have translations that follow each other in the output language. In detail, the position j in the source sentence is aligned with the position i in the target sentence when translating, denoted as $i = a_j$. Alignment model is a global reordering model. For whole sentence, we denote the alignment as:

$$a_1^J := a_1 \cdots a_J$$

There is not probabilistic model for alignment in IBM Model 1, it treats probability of all alignments equally. IBM Model 2 addresses the issue of alignment with an explicit model for alignment based on the positions of the input output words and sentence lengths $a(i|j, I, J)$.

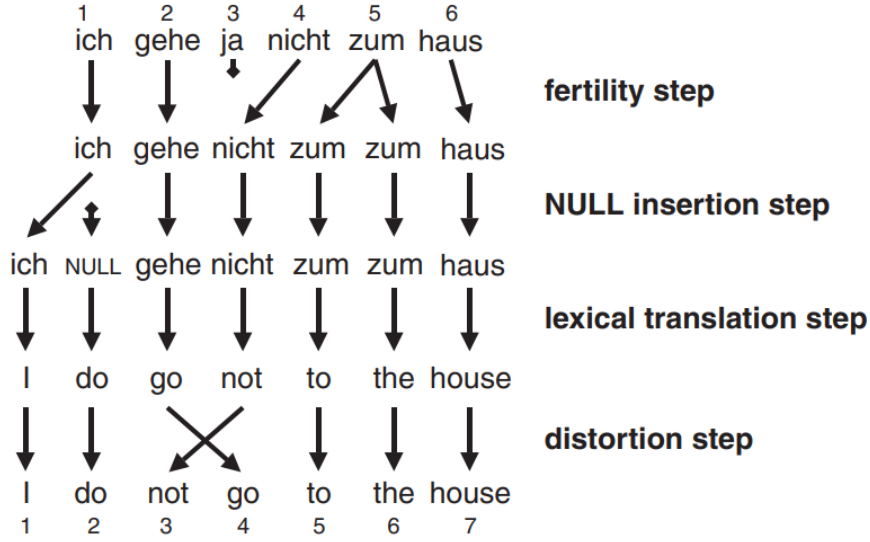


Figure 2.1: Illustration of a translation process using IBM-3 model (Koehn [2009])

Fertility model the specific number of output words in the output language. Generally one word in one language just translate into one single word in the other language. But some words produce multiple words or get dropped. we denote the fertility model as $\phi(e)$, so the length of translation sentence

$$J = \sum_{i=0}^I \phi(e_i)$$

The word-based translation process can be described as the following figure, we use IBM-3 model for example: IBM-3 model contains four steps:

Weighted Model

To overcome shortcomings in modeling, introduce scaling exponents λ_m as in speech recognition. For example

$$Q(f_1^J, e_1^I; a_1^J) = P(J|I)^{\lambda_1} \cdot \prod_i P(e_i|e_h)^{\lambda_2} \cdot \prod_j [p(a_j|a_{j-1}, I, J)^{\lambda_3} \cdot p(f_j|e_{a_j})^{\lambda_4}]$$

where the four probabilities corresponds to length model, language model, alignment model and lexical model separately. e_h denotes the history words in N -gram model

2.1.2 Phrase-Based Model

Actually when translating, words may not be the best candidates for the smallest units for translation, sometimes one word in a foreign languages should be translated into two English words, or vice versa. Word-based models often break down in these cases.

Phrase-based models typically do not strictly follow the noisy-channel approach proposed for word-based models, but use a log-linear framework This model allow s straightforward integration if additional features.

Log-linear Model Combination: Consider arbitrary models ("feature functions"):

$$q_m(f_1^J, e_1^I; a_1^J) > 0 \quad m = 1, \dots, M$$

$$\begin{aligned} Q(e_1^I, f_1^J; a_1^J) &= \prod_{m=1}^M q_m(f_1^J, e_1^I; a_1^J)^{\lambda_m} \\ &= \exp\left(\sum_{m=1}^M \lambda_m \log q_m(f_1^J, e_1^I; a_1^J)\right) \end{aligned}$$

In this frame work, we view each data point as a vector of features and the model as a set of corresponding feature functions, these functions are trained separately and combined assuming that they are independent of each other.

Components for log-linear model can be such as language model, phrase translation model, reordering model are used as feature functions with appropriate weights.

- Phrase translation model can be learned from a word-aligned parallel corpus, alternatively we also can use expectation maximization algorithm to directly find phrase alignment for sentence pairs.
- Reordering model for phrase case is typically modeled by a distance based reordering cost that discourage reordering in general. Lexicalized reordering model can be introduced for specific phrase pair.

The model can be further extended by components like: bidirectional translation probabilities, lexical weighting, word penalty and phrase penalty.

2.1.3 Decipherment

Ravi and Knight [2011] frame the MT problem as a decipherment task, treating the foreign text as a cipher. They also propose iterative expectation-maximization

method (EM) and Bayesian decipherment to train this unsupervised model. Many concepts in decipherment are the same as the SMT.

IBM model tries to maximize the probability with hidden alignment model

$$\operatorname{argmax}_{\theta} \prod_{e,f} P_{\theta}(f|e) = \operatorname{argmax}_{\theta} \prod_{e,f} \sum_a P_{\theta}(f, a|e)$$

While for unsupervised case, we train

$$\operatorname{argmax}_{\theta} \prod_f P_{\theta}(f) = \operatorname{argmax}_{\theta} \prod_f \sum_e P(e) \cdot P_{\theta}(f|e)$$

for hidden alignments:

$$\operatorname{argmax}_{\theta} \prod_f \sum_e P(e) \cdot \sum_a P_{\theta}(e, a|e)$$

Since the model is very complicated, more assumptions are added to the model. The model accounts for word substitutions, insertions, deletions and local reordering during the translation process but does not incorporate fertilities or global re-ordering.

The generative process:

1. Generate an English sentence e with probability $P(e)$.
2. Insert a NULL word at any position in the English sentence with uniform probability.
3. For each English word token e_i (including NULLs), choose a foreign word translation f_i , with probability $P_{\theta}(f_i|e_i)$, the foreign word may be NULL.
4. Swap any pair of adjacent foreign words f_{i-1}, f_i , with probability $P_{\theta}(swap)$.
5. Output the foreign string f_1^M , skipping over NULLs.

However, this method is limited to rather short sentences and simplistic settings.

2.2 Neural Machine Translation

Neural MT (NMT) has recently become the dominant method for MT task. In comparison to the traditional statistical MT (SMT), NMT systems are trained end-to-end, taking advantages of the continuous representation of the hidden states

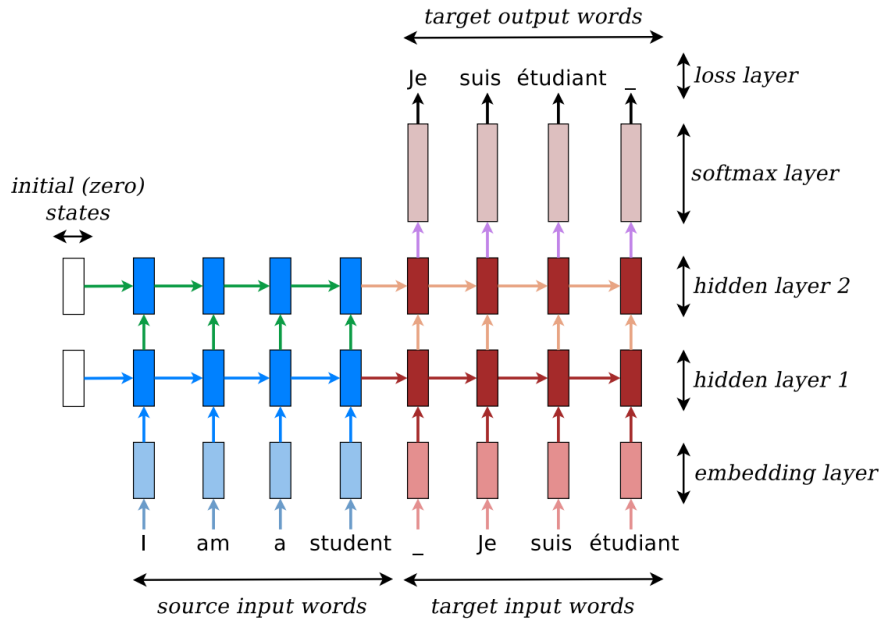


Figure 2.2: Neural MT – example of a deep recurrent architecture (Luong et al. [2015a])

that greatly alleviate the sparsity problem and make better use of more contextual information. Besides this, traditional phrase-based MT, which consists of several models that are tuned separately, neural MT tries to build a more general neural network model which can directly output translations given input, it contains only one single model, and only one single training criterion.

Sutskever et al. [2014] first use a multi-layer Long Short-Term Memory (LSTM) to train the MT system, attention mechanism has lately been used to improve neural MT by selectively focusing on parts of the source sentence during translation. The inherently sequential nature precludes parallelization within training examples. Convolutional sequence to sequence (ConvSeq2seq) and Transformer architectures are brought forward for significantly more parallelization during training to better exploit the GPU hardware and these models can reach a new state of the art in translation quality.

2.2.1 Encoder-Decoder Framework

The most common network structure is the encoder-decoder framework, in which two recurrent neural networks work together to transform one sequence to another.

An encoder condenses an input sequence into a vector \mathbf{c} , and a decoder unfolds that vector into a new sequence. The model can be expressed as:

$$p(e_1^I | f_1^J) = \prod_i p(e_i | e_0^{i-1}, f_1^J) = \prod_i p(e_i | e_0^{i-1}, \mathbf{c})$$

The RNN decoder predicts the next word on the basis of the internal states and the previous word. In more detail, one can parameterize the probability of decoding each word e_i as:

$$p(e_i | e_0^{i-1}, \mathbf{c}) = \text{softmax}(g(\mathbf{s}_i))$$

with g being the transformation function that outputs a vocabulary-size vector. Here \mathbf{s}_i is the RNN hidden state, updated as:

$$\mathbf{s}_i = f(\mathbf{s}_{i-1}, \mathbf{c})$$

It is common in neural MT systems to use a beam-search to sample the probabilities for the words in the sequence output by the model. The wider the beam width, the more exhaustive the search, and, it is believed, the better the results.

Drawbacks:

1. The model compresses all information from the input sentence into a hidden vector \mathbf{c} , while ignores the length of input sentence, when the length of input sentence get very long, even longer than the training sentences, it becomes harder to extract specific information for predicting the target word, the performance will get worse.
2. It's not suitable to assign the same weight to all input words, one target word corresponds usually to one or several words in the input sentence. Treating all words equally does not distinguish the source information and influence the performance badly.

2.2.2 Attention Mechanism

Alignment is the problem in MT that identifies which parts of the input sequence are relevant to each word in the output, whereas translation is the process of using the relevant information to select the appropriate output. Bahdanau et al. [2014] introduce an extension to the encoder-decoder model which learns to alignment and translate jointly. Each time the model predicts the next target word, it softly searches for a set of positions in a source sentence where the most relevant information is concentrated. The model then predicts a target word based on the context vectors associated with these source positions and all the previous generated target words.

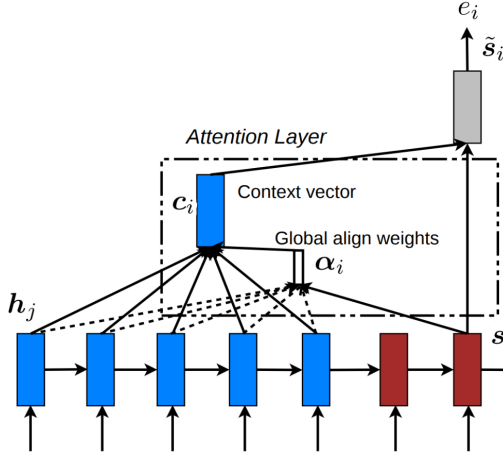


Figure 2.3: Global attention

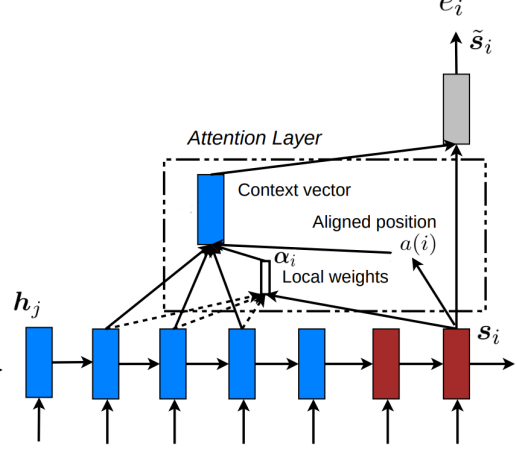


Figure 2.4: Local attention

Described in formula, attention mechanism derives a context vector \mathbf{c}_i that capture the input information to help to predict the target word at the position i . Given the target hidden state \mathbf{s}_i and the source-side context vector \mathbf{c}_i , we can compute the hidden state $\tilde{\mathbf{s}}_i$ by combining the current hidden state \mathbf{s}_i and the context vector \mathbf{c}_i :

$$\tilde{\mathbf{s}}_i = \tanh(W_c[\mathbf{c}_i; \mathbf{s}_i])$$

Then the target word is correspondingly predicted by softmax function:

$$p(e_i | e_0^{i-1}, f_1^J) = \text{softmax}(W_s \tilde{\mathbf{s}}_i)$$

Global attention & Local attention

Global attention attends all the input words, weighted average of source hidden states (word representations):

$$\mathbf{c}_i = \sum_j \alpha_i(j) \cdot \mathbf{h}_j$$

where $\alpha_i(j)$ is the normalized attention weight that output at position i is aligned with input position j , and it can be calculated as the following softmax-like function:

$$\alpha_i(j) = \text{align}(\mathbf{s}_i, \mathbf{h}_j) \quad (2.1)$$

$$= \frac{\exp(\text{score}(\mathbf{s}_i, \mathbf{h}_j))}{\sum_{j'} \exp(\text{score}(\mathbf{s}_i, \mathbf{h}_{j'}))} \quad (2.2)$$

Since score function is referred as a content based function, different forms can be considered:

$$\text{score}(\mathbf{s}_i, \mathbf{h}_j) = \begin{cases} \mathbf{s}_i^T \mathbf{h}_j & \text{dot} \\ \mathbf{s}_i^T W_a \mathbf{h}_j & \text{general} \\ \mathbf{v}_i^T \tanh(W_a[\mathbf{s}_i; \mathbf{h}_j]) & \text{concat} \end{cases} \quad (2.3)$$

For global attention, for each target word we need to attend the whole input sentence, it is very expensive and impractical to translate longer sentences. Luong et al. [2015a]) proposed the local attention, the local attention is actually the trade-off between soft and hard attention.

Local attention first predicts an aligned position $a(i)$ for each target word at position i . Then the context vector \mathbf{c}_i is a weighted sum within the window $[a(i) - D, a(i) + D]$, D is selected empirically. The model predict the aligned position p_i as followed:

$$a(i) = S \cdot \text{sigmoid}(\mathbf{v}_p^T \tanh(W_p \mathbf{s}_i))$$

W_p and \mathbf{v}_p are the model parameters which will be learned to predict positions. S is the source sentence length. As a result of sigmoid function, $a(i) \in [0, S]$. To favor the words near position $a(i)$, they place a Gaussian distribution which centered at $a(i)$.

$$\alpha_i(j) = \text{align}(\mathbf{s}_i, \mathbf{h}_j) \cdot \exp\left(-\frac{(j - a(i))^2}{2\sigma^2}\right)$$

2.2.3 Transformer

The RNN encoder-decoder models have achieved state of the art in sequence modeling and MT problem. However such RNN models also have some disadvantages, because of their inherently sequential computation which prevents parallelization across elements of the input sequence. That means it is more difficult to fully take advantage of the modern computing devices like GPU and TPU. Convolutional Gehring et al. [2017] and fully-attentional feed-forward architectures like Transformer Vaswani et al. [2017] models are proposed as alternatives for RNNs.

Self-Attention

$$\begin{aligned} \alpha_{ki} &= \frac{\exp(\text{score}(\mathbf{h}_k, \mathbf{h}_i))}{\sum_{j=1}^I \exp((\text{score}(\mathbf{h}_k, \mathbf{h}_j)))} \\ \text{score}(\mathbf{h}_k, \mathbf{h}_i) &= \mathbf{h}_k^\top \mathbf{h}_i / \sqrt{d}, \quad d = \dim(\mathbf{h}_k) \\ \mathbf{h}'_k &= \sum_{i=1}^I \alpha_{ki} \cdot \mathbf{h}_i \end{aligned}$$

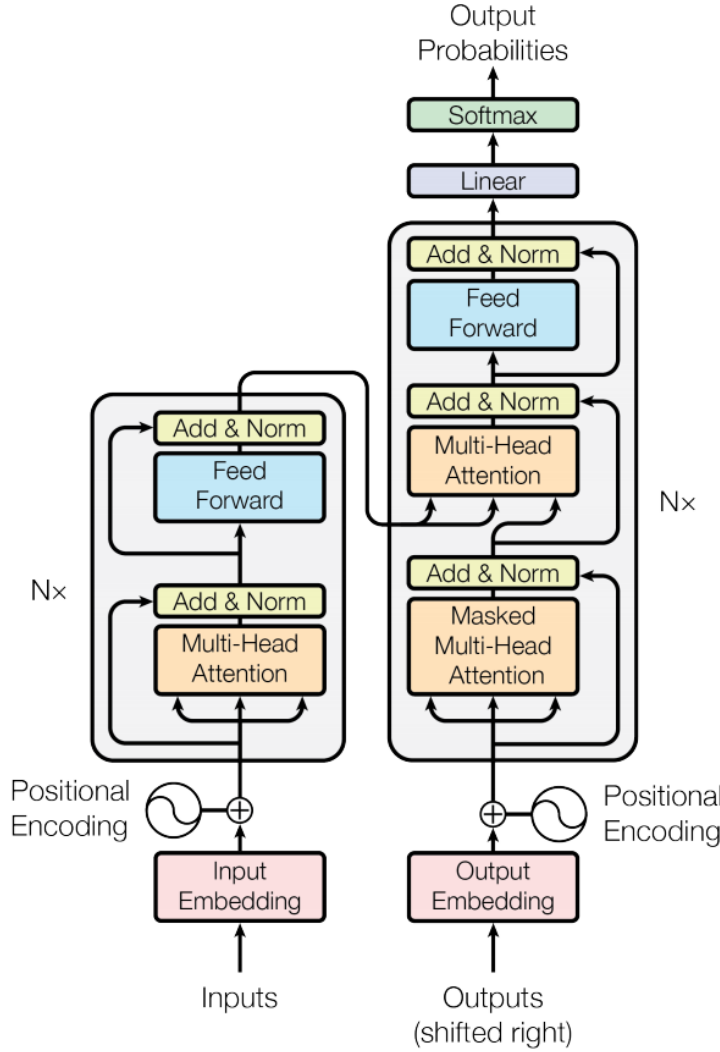


Figure 2.5: The Transformer model architecture (Vaswani et al. [2017])

Multi-head attention

First map Q, K, V into h many lower dimension spaces via linear mapping matrix. Here h is the number of heads. Then apply attention and concatenate the outputs. Suppose the original dimensions of queries, keys, values are d_m . We the number of heads is h , then we set $d_k = d_v = d_m/h$. With $i \in 1 \dots h$ different matrix $W_i^Q \in R^{d_m \times d_k}$ $W_i^K \in R^{d_m \times d_k}$ $W_i^V \in R^{d_m \times d_k}$ and $W^O \in R^{d_m \times d_k}$. The number of attention heads in Transformers impacts their ability to capture long-distance de-

dependencies. Especially many-headed multi-head attention is essential for modeling long-distance phenomena with only self-attention.

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

where $\text{head}_i \in \mathbb{R}^{n_q \times d_k}$, $\text{Multihead}(Q, K, V) \in \mathbb{R}^{n_q \times d_m}$

Encoder-Decoder Attention

$$\alpha_{ki} = \frac{\exp(\text{score}(\mathbf{s}_k, \mathbf{h}_i))}{\sum_{j=1}^I \exp(\text{score}(\mathbf{s}_k, \mathbf{f}_j))}$$

$$\text{score}(\mathbf{s}_k, \mathbf{h}_i) = \mathbf{s}_k^\top \mathbf{h}_i / \sqrt{d}, \quad d = \dim(\mathbf{s}_k)$$

Positional Encoding

Since there is no RNN or CNN structures in transformer model, we need also need to make use of sequence information for seq2seq learning. We need inject position information into the embedding.

$$PE_{(\text{pos}, 2i)} = \sin(\text{pos}/10000^{2i/d_m})$$

$$PE_{(\text{pos}, 2i+1)} = \cos(\text{pos}/10000^{2i/d_m})$$

where pos is the position and i is the dimension. That is, each dimension of the positional encoding corresponding to a sinusoid.

2.2.4 Neural Unsupervised Machine Translation

Although the neural MT systems have achieved near human-level performance on some languages, the lack of large parallel corpora poses a major challenge for many low-source languages. Techniques are raised to alleviate this issue with like triangulation and semi-supervised learning techniques. Monolingual data are leveraged to improve the quality of translation system.

Recently, several neural algorithms with fully unsupervised settings are proposed. The core ideal is very similar to the dual learning in MT. The systems all contain bi-directional models as sub-systems. Pseudo parallel data are created with back-translation. It turns the unsupervised problem into a supervised setting. As better translation, both sub-systems get improved synchronously.

Dual Structure

Dual structure is inspired by the observation: any MT has a dual task: e.g. German-to-English translation (primal) and English-to-German translation (dual). Dual

tasks can form a closed loop and generate informative feedback signals to for both translation models. Based on the feedback signals generated during this process, and leverage language model, we can minimize the reconstruction error of the original sentences. We can iteratively update the two models until convergence. He et al. [2016] trained two agents to translate in opposite directions (e.g. French \rightarrow English and English \rightarrow French), and make them teach each other through a reinforcement learning process. While promising, this approach still requires a parallel corpus for a warm start. Lample et al. [2018] simplified the structure, in the proposed model, the gradients will not be back-propagated through the reverse model, but only make use of the back-translation. The training process as followed:

Algorithm 1 Unsupervised Machine Translation

Language models: Learn language models P_s, P_t over source and target languages;

Initial translation model: Leveraging P_s, P_t , lean two initial translation models in each direction: $P_{s \rightarrow t}^{(0)}, P_{t \rightarrow s}^{(0)}$;

for $k = 1$ to N **do**

Back-translation: Generate source and target sentences using the current translation models $P_{t \rightarrow s}^{(k-1)}$ and $P_{s \rightarrow t}^{(k-1)}$, leverading P_s, P_t ;

Retrain: Train new translation models $P_{t \rightarrow s}^{(k)}$ and $P_{s \rightarrow t}^{(k)}$ using the generated sentences and leverading P_s, P_t ;

end

Initialization

Without enough information to start the dual MT system, it will be hard for the system to catch the meaningful signals and then it will take much more iterations. However if we initialize the system by a naive a word-by-word translation of the sentence, where the bilingual lexicon are derived from the same monolingual data. Though such initial "word-by-word" translation maybe poor if languages or corpora are not closely related, it still preserves some original semantics. Such word-by-word translation can actually achieve several BLEU scores.

Shared Latent Representation

A shared encoder representation actis like an interlingua, which is translated in the decoder corresponding language regardless less of the input source language. Since the only supervision information only comes from the monolingual data. The model learn to translate by learning to reconstruct in both language from this shared space.

There are two different methods to learn the shared feature space:

- Shared encoder The system use only one encoder that is shared by both languages involved. The universal encoder is aimed to produce a language independent representation of the input language but the decoder should separately translate then into corresponding language.
- Adversarial training Train discriminator to classify between the encoding of source and target sentences. The discriminator operates on the output of the encoder, the encoder is trained instead to fool the discriminator.

Optimization

When minimizing the loss function, gradients will not be back propagated through the reverse model which generate the data. Instead the objective function minimized at every iteration is the sum of L^{auto} and L^{back}

- Denoising autoencoder loss

$$L^{auto} = \mathbb{E}_{e \sim E}[-\log P_{t \rightarrow t}(e|\text{noise}(e))] + \mathbb{E}_{f \sim F}[-\log P_{s \rightarrow s}(f|\text{noise}(f))]$$

where s $P_{s \rightarrow s}$ and $P_{t \rightarrow t}$ are the composition of encoder and decoder both operating in the source and target sides, respectively

- Back-translation loss

$$L^{back} = \mathbb{E}_{e \sim T}[-\log P_{s \rightarrow t}(e|u(e))] + \mathbb{E}_{f \sim S}[-\log P_{t \rightarrow s}(f|v(f))]$$

we denote the sentence that translated by intermediate target-to-source translation model as $u(y)$, similarly denote the sentence translated by source-to-target model as $v(x)$, so $u(y)$ should in source language and $v(x)$ should in target language. The pairs $(x, v(x))$, $(u(y), y)$ constitute synthetic parallel sentences.

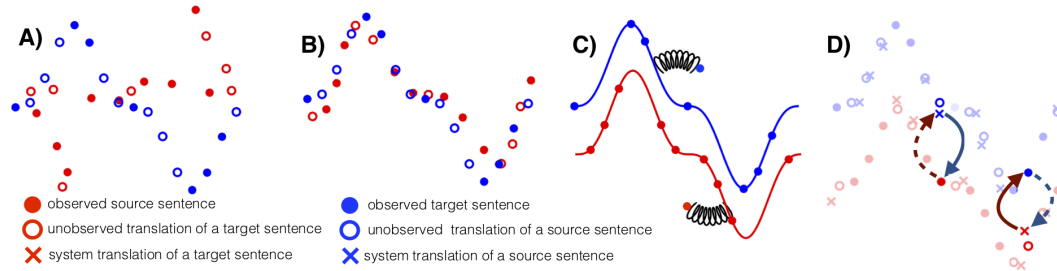


Figure 2.6: Key techniques for neural unsupervised MT: A) the two original monolingual data; B) Initialization, the two distribution are roughly aligned; C) Denoising autoencoder, make the distribution closer to normal one in corresponding language D) Back-translation, improve the both translation model in iterative way (Lample et al. [2018])

Chapter 3

Word Embedding

Word embeddings is distributed representation of words in a vector space. With the learning algorithm it can capture the contextual or co-occurrence information. The word embedding has an interesting and important property: similar words will have similar distribution in the embedding space, with that property, we can find meaningful near-synonyms or Some successful methods for learning word embedding like word2vecMikolov et al. [2013c], Pennington et al. [2014]

3.1 Monolingual Embedding

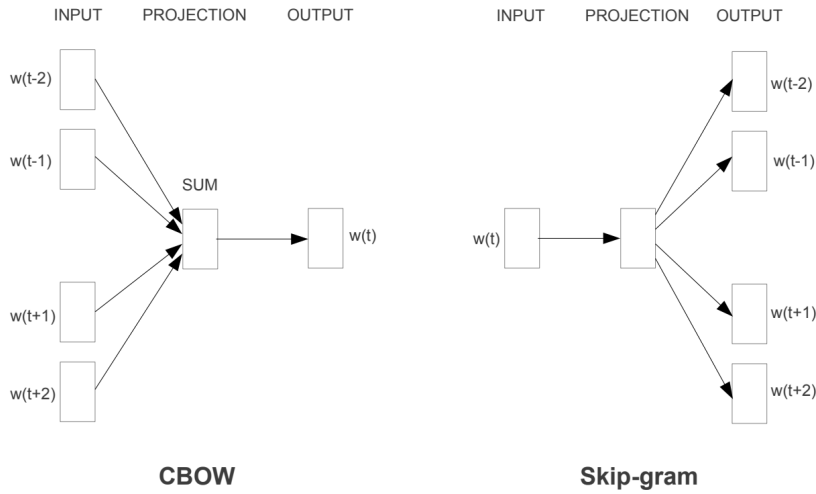


Figure 3.1: Global attention model (Mikolov et al. [2013a])

3.1.1 CBOW and Skip-gram Model

CBOW model and skip-gram model are currently mainstream neural network to learn the word embedding. Algorithmically, CBOW tries to predict the current

word based on the context while skip-gram model is to find better representations that are useful for predicting the surrounding words. Using skip-gram model as example, given a large training corpus represented as a sequence of words e_1, \dots, e_N , the objective of the model is to maximize the following log-likelihood:

$$\frac{1}{I} \sum_{i=1}^I \sum_{-c \leq j \leq c, j \neq 0} \log p(e_{i+j}|e_i)$$

where e_i is the current word and c is the size of training context. We define the loss for skip-gram model as

$$\mathcal{L}_{\text{skip-gram}}^e = -\frac{1}{I} \sum_{i=1}^I \sum_{-c \leq j \leq c, j \neq 0} \log p(e_{i+j}|e_i)$$

Similarly, we define the loss for CBOW model as

$$\mathcal{L}_{\text{CBOW}}^e = -\frac{1}{I} \sum_{i=1}^I \sum_{-c \leq j \leq c, j \neq 0} \log p(e_i|e_{i+j})$$

Furthermore in skip-gram model, suppose given a scoring function s which maps word pairs to score in \mathbb{R} like cosine similarity, the probability can be defined as

$$p(e_{i+j}|e_i) = \frac{\exp\{s(e_{i+j}, e_i)\}}{\sum_{e' \in V^e} \exp\{s(e', e_i)\}}$$

where V^e is the whole vocabulary.

Negative Sampling

However the normalization on the whole vocabulary is very expensive since it is conducted for all words at every training step. Negative sampling is proposed to approximate the softmax to make it computationally more efficient. The problem of predicting words can be considered as an independent binary classification task. Negative sampling trains the model to distinguish a target word from negative samples drawn from a noise distribution P_{noise} .

$$p(e_{i+j}|e_i) = \log Q_{\theta}(D = 1|e_{i+j}, e_i) + \frac{1}{k} \sum_{e' \sim P_{\text{noise}}} \log Q_{\theta}(D = 0|e', e_i)$$

where $Q_{\theta}(D = 1|w_t, w_s)$ is the binary logistic regression probability. According to empirical results, CBOW works better on smaller datasets because CBOW smoothes over a lot of the distributional information while Skip-Gram model performs better when we have larger datasets

3.1.2 Enriching Embedding with Subword Information

The training methods above treat each word as a distinct word embedding, however intuitively we can obtain more information from the morphological information of words. A subword model was proposed to try to fix such problem. The training network is similar, the model design a new presentation of the word: it adds special symbols $<$, $>$ as boundary information at the beginning and the end of a word. Then a normal word is represented as a bag of character n -grams. For example the word "where" and n equals 3, then it can be represented as the following 5 tri-grams:

$$< wh, whe, her, ere, re >$$

Suppose in this way for a specific word $e \in G_e$ the set of character n -grams, we assign for each character n -gram g in G_e a distinct vector \mathbf{v}_g , we calculate the score function using the sum of character-level inner product of vectors:

$$s(e', e) = \sum_{g \in G_e} \mathbf{v}_g^T \mathbf{e}'$$

where \mathbf{e}' the embedding of e'

3.2 Cross-lingual Word Embedding

NLP tasks in multilingual scenarios is receiving increasing interest. The need to represent meaning and transfer knowledge in cross-lingual applications has motivated the works on cross-lingual models which learn the representations in a joint embedding space. Mikolov observe that word embeddings trained separately on monolingual corpora exhibits isomorphic structure across languages, as illustrated in Figure 3.2. That means we can create a connection between source embedding and target embedding even with simple linear mapping. This has far-reaching implication on low-resource scenarios (Adams et al. [2017]), because word embedding requires only plain text to train, which is the most abundant form of linguistic resource.

3.2.1 Training Algorithm

Most training methods can be divided into joint training approaches and mapping-based approaches. The difference is that joint methods train directly the cross-lingual word embedding while the mapping-based approaches first train monolingual word representations on corresponding corpus separately. Then mappings are learned to project different embeddings into the shared space. They learn the transformation from word alignments or bilingual dictionaries.

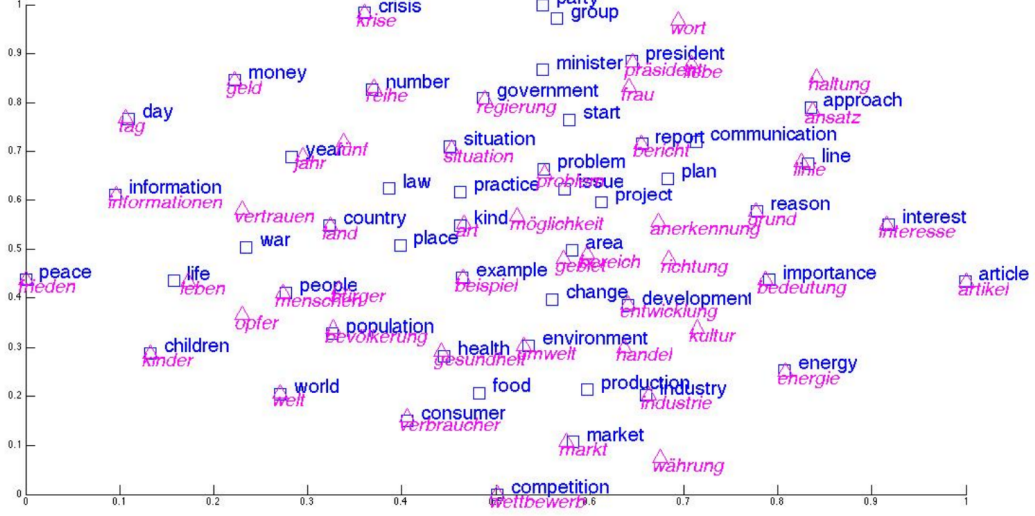


Figure 3.2: A cross-lingual embedding space between German and English (Ruder et al. [2017])

Joint Training Approaches

Joint training means to optimize the source and target embeddings objectives $\mathcal{L}(\cdot)$ jointly with the cross-lingual regularization term $\Omega(\cdot)$.

$$\mathcal{L} = \mathcal{L}_{\text{skip-gram}}^e + \mathcal{L}_{\text{skip-gram}}^f + \lambda \cdot \Omega(\cdot)$$

where the first two losses are exact the same as those in monolingual embedding training. While the $\Omega(\cdot)$ encourages the similar words representations to be similar for words that are related across different languages.

Different cross-lingual regularization terms are proposed to optimize the cross-lingual embeddings (Coulmance et al. [2016], Luong et al. [2015b], Gouws et al. [2015]). For example, Gouws et al. [2015] models it as:

$$\Omega_{\text{BiBOWA}} = \left\| \frac{1}{I} \sum_{e_i \in e_1^I} e_i - \frac{1}{J} \sum_{f_j \in f_1^J} f_j \right\|^2$$

where f_j and e_i are the word embeddings for words in parallel source and target sentences. Instead of relying on expensive word alignments which minimize the distance that aligned to each other, they minimize the distance between the mean of the word representations in the aligned sentences.

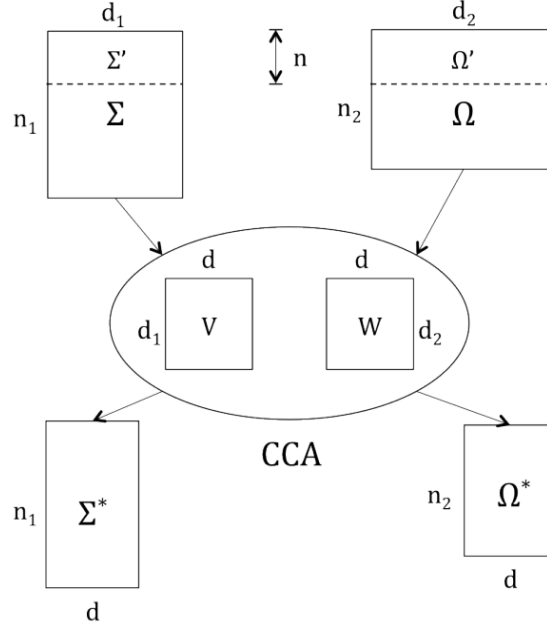


Figure 3.3: Cross-lingual projection with CCA (Faruqui and Dyer [2014])

Mapping-based Approaches:

There are two methods to learn the cross-lingual word embedding, one is to learn the transformation for each language, both monolingual embeddings are mapped into a shared embedding space. The typical model is CCA-based approaches (Faruqui and Dyer [2014], Dhillon et al. [2011] and Lu et al. [2015]). The other method is to learn the mapping from the source embedding space to the target embedding space. The criterion is to minimize the mean squared error between the transformed source embeddings and the target embeddings. The advantage of this method is that it is really fast to learn the embedding alignments. People criticize whether linear or nonlinear can capture the relationship between all words in different languages.

Canonical Correlation Analysis (CCA)

let $F' \in \mathbb{R}^{d_1 \times n}$ and $E' \in \mathbb{R}^{d_2 \times n}$ be the embedding matrices comprised of seed word translation pairs with size n . The W and V are the projection matrices which mapped the original embedding matrices into F^* and E^* , which are in the joint space.

$$F^* = WF' \quad E^* = VE'$$

The objective of CCA is to find the two matrices such that the projections of the parallel matrices are maximally correlated:

$$\operatorname{argmax}_{W \in \mathbb{R}, V \in \mathbb{R}} \frac{\mathbb{E}[(F^*)(E^*)]}{\sqrt{\mathbb{E}[(F^*)^2]} \sqrt{\mathbb{E}[(E^*)^2]}}$$

Once we get the two projection matrices, we can calculate the cross-lingual word embedding for the entire vocabulary.

A linear feature mapping is often not sufficiently powerful to faithfully capture the hidden, non-linear relations within the data. Lu et al. [2015] propose a non-linear extension of CCA using deep neural networks, where two neural networks are used to extract features separately, and trained to maximize the correlations between these features, measured by a linear CCA step with projection mappings the neural network weights and linear projections are optimized together. Different from CCA, the neural network model does not have a closed-form solution, the parameters can be learned via gradient-based optimization.

Cross-lingual Mapping

Mikolov et al. [2013b] notice that the geometric relation that hold between words are similar across languages. This suggest that, we can jump the third embedding space, learn the transformation directly from source embedding space to target space. Given a seed word translation pair set d of size n , the objective is to learn W that minimizes the distance between the mapped source word embedding $W\mathbf{f}_i$ and target embedding \mathbf{e}_i :

$$\operatorname{argmin}_W \sum_{\substack{i=1 \\ (\mathbf{f}_i, \mathbf{e}_i) \in d}}^n \|W\mathbf{f}_i - \mathbf{e}_i\|_2^2$$

This can also be expressed in matrix form as minimizing the squared Frobenius norm of the residual matrix:

$$\operatorname{argmin}_W \|WF - E\|_F^2$$

Xing et al. [2015] showed that the results are improved when we use l_2 normalized embeddings and constrain the W to be an orthogonal matrix. In this case we can use Procrustes analysis which advantageously offers a closed form solution obtained from the singular value decomposition

$$W^* = \operatorname{argmin}_{W \in O_d(\mathbb{R})} \|WF - E\|_F^2 = UV^T, \quad U\Sigma V^T = SVD(EF^T)$$

In order to make learning and inference consistent, Joulin et al. [2018] propose the relaxed CSLS loss (RCSLS), which is inspired from the work of Conneau et al. [2017]

The loss function can be rewritten as:

$$\min_{\mathbf{W} \in \mathcal{O}_d} = \frac{1}{n} \sum_{i=1}^n \left\{ -2\mathbf{f}_i^\top \mathbf{W}^\top \mathbf{e}_i + \frac{1}{k} \sum_{\mathbf{e}_j \in \mathcal{N}_e(\mathbf{W}\mathbf{f}_i)} \mathbf{f}_i^\top \mathbf{W}^\top \mathbf{e}_j + \frac{1}{k} \sum_{\mathbf{W}\mathbf{f}_j \in \mathcal{N}_f(\mathbf{e}_i)} \mathbf{f}_j^\top \mathbf{W}^\top \mathbf{e}_i \right\}$$

Normally, the size n of the dictionary is too small with respect to the entire vocabulary size, in order to incorporate unlabeled data, the unpaired words in the dictionary are used as "negatives samples" in the RCSLS loss, as result, we can optimize the loss function not only the dictionary but the whole vocabulary.

Chapter 4

Unsupervised Learning of Cross-lingual Word Embedding

4.1 Initialization

4.2 Vocabulary-based Training

4.2.1 Iterative Procrustes Analysis

4.2.2 Adversarial Training

4.3 Corpus-based Training

Chapter 5

Sentence Translation

As mentioned previously, training a traditional machine translation system requires large parallel data. With cross-lingual word embedding we can already found ambiguous word translation, in this chapter we propose a simple yet effective method to improve quality of translation which starts from the word-by-word translation. We integrate monolingual model like language model and denoising neural network of the target side to produce meaningful sentence translation. Since all the sub-models are trained on monolingual corpora, our proposed method is total unsupervised. Our system surpasses state-of-the-art unsupervised translation without costly iteratively training.

5.1 Context-aware Beam Search

5.1.1 Language Model

Language models are widely applied in natural language processing, especially machine translation which need frequent queries.

n-gram models

N-gram language models use the Markov assumption to break the probability of a sentence into the product of the probability of each word given a limit history of preceding words.

$$p(w_1^N) = \prod_{i=1}^N p(w_i | w_1, \dots, w_{i-1}) = \prod_{i=1}^N p(w_i | w_{i-(n-1)}, \dots, w_{i-1})$$

The conditional probability can be calculated from *n*-gram model frequent counts:

$$p(w_i | w_{i-(n-1)}, \dots, w_{i-1}) = \frac{\text{count}(w_{i-(n-1)}, \dots, w_i)}{\text{count}(w_{i-(n-1)}, \dots, w_{i-1})}$$

Language model tries to handling sparse data problem because some words or phrases have not been seen yet in the training corpus does not mean they are not impossible. Different smoothing techniques like back-off or interpolation are implemented to assign a probability mass to unseen cases.

5.1.2 Beam Search

The complexity of a search graph is exponential to the length of the given source sentence. Beam search is a heuristic search algorithm that explores a graph by expanding the most promising nodes. At each step of the search process, it will evaluate all the candidates together with the reserved translation results from last step, it will only store a predetermined number (beam width) of translations for next step. The greater the beam width is, the fewer states will be pruned. So it is suggested to prune these word translation candidates as soon as possible to reduce the search space and speed up the translation. According to the similarity of cross-lingual word embedding, we are able to find some meaningful translation candidates for a given word. But there are also words that are noisy in the candidates or obviously incorrect because of grammar checking. With the support of language model, we can select the most probable words from previous word translation candidates.

Given a history h of target word before e , the score of e to be the translation of f is defined as:

$$\hat{e}_1^N = \operatorname{argmax}_{e_1^N} \prod_{n=1}^N p^{\lambda_{LM}}(e_n | e_{n-4}^{n-1}) \cdot q^{\lambda_{emb}}(f_n, e_n)$$

where the lexicon score $q(f, e) \in [0, 1]$ defined as:

$$q(f, e) = \frac{d(f, e) + 1}{2}$$

$d(f, e) \in [-1, 1]$ cosine similarity between f and e

In experiments, we find such lexicon score works better than others, e.g. sigmoid or softmax

5.2 Denoising Neural Network

5.2.1 Denoising Auto-encoder

With the help of language model we have actually improved the quality of word-by-word translation but the results are still far from an acceptable one because of the drawback of the word-by-word mechanism, maybe to some degree we can infer the meaning of sentence, but the sequence of sentences depends on the specific language. We implement the sequential denoising autoencoder to improve the translation output.

An autoencoder is a neural network that is trained to copy its input to its output, autoencoders minimize the loss function like:

$$L(\mathbf{x}, g(f(\mathbf{x})))$$

where L penalizing the difference between the input and output. While a denoising autoencoder (DAE) instead minimizes

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}})))$$

where $\tilde{\mathbf{x}}$ is a noise transformation of \mathbf{x} and denoising autoencoder will try to learn to ignore the noise in \mathbf{x} reconstruct the correct one. Sequential denoising autoencoder will find robust representation of sentences. In practice, denoising autoencoder consists of two parts, namely encoder and decoder. The encoder processes noised data and produces real-valued vectors as an encoding feature of the data. The computational graph of the denosing autoencoder, which attempt to reconstruct the normal input \mathbf{x} from it corrupted version $\tilde{\mathbf{x}}$. The model is trained by minimize the loss.

For our sequential denoising model, the label sequences would be the monolingual data of the target language. However we do not have the noise input. In order to make the model run correctly, we should mimic the noise sentence of word-by-word translation on the target monolingual corpus.

We design different noise types w.r.t. the word-by-word translation. In the experiments, we inject the artificial noise into a clean sentence, the experiment results shows the noise is reasonable and suitable in this case.

5.2.2 Noise Model

We design three types of noise to handle the fertility and reordering problem, namely reordering noise, insertion noise and deletion noise. In experiments, the noise model can improve the sentence translation, but since it actually starts from the word-by-word translation, it can only deal with reordering in limited distance, cannot work for global reordering.

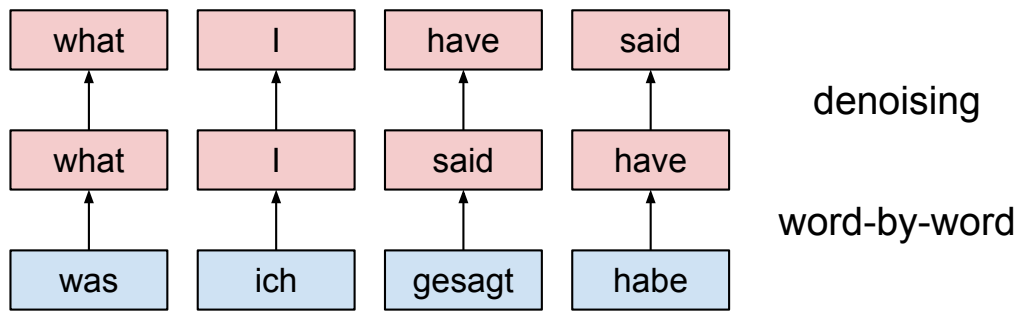


Figure 5.1: Reordering noise

Reordering Noise

The reordering problem is a common phenomenon in the word-by-word translation since the sequence in source language is not exact the sequence in target language. For example, in the grammar of German, the verb is often placed at the end of the clause. "um etwas zu tun". However in English it is not the case, the corresponding translation sequence is "to do something". The verb should always before the noun. In our beam search, language model only assisting in choosing more suitable word from the translation candidates, it cannot reorder the word sequence at all.

For a clean sentence from the target monolingual corpora, we corrupt the word sequence by permutation operation. We limit the maximum distance between the original position and its new position.

The design of reordering noise is as followed:

1. For each position i , sample an integer δ_i from $[0, d_{per}]$
2. Add δ_i to index i and sort $i + \delta_i$
3. Rearrange the words to be in the new positions, to which where indices have been moved

Reordering is actually depends on the specific language pair. However in the experiments we found the performance of the denoising network aimed at such noise is not obvious. The Bleu score before and after the process is close.

Insertion Noise

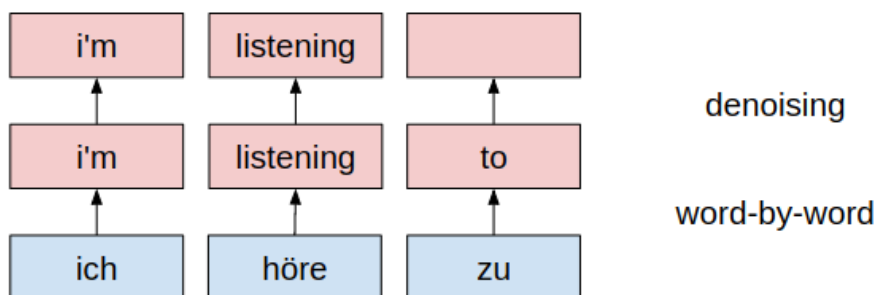


Figure 5.2: Insertion noise

The word-by-word translation system predict the source word at every position of the sentence. However the vocabularies of different systems are not symmetric, for example, in German there are more compound words than that in English. So when

translating cross languages, there are a plenty of cases that a single word will be translated to multiple words and multiple words correspond to a single conversely. We focus on such a case: from a German sentence: "ich höre zu" to "i'm listening". A very frequent word "zu" which corresponds to "to" in English, is dropped from the sentence. The design of reordering noise is as followed:

1. For each position i , sample a probability $p_i \sim \text{Uniform}(0, 1)$
2. If $p_i < p_{ins}$, sample a word e from the most frequent V_{ins} target words and insert it before the position i

We limit the insertion word in a set consisting of the top frequent word in the target language V_{ins}

Deletion Noise

The deletion noise is just a contrary case of insertion noise. Because we are limited to generate only ne word per source word, it is also possible that a target word in the reference is not related to any source word. For example for "eine der besten" the corresponding translation is "one of the best". We need to add an extra preposition in the target sentence. To simulate such situation, we drop some words randomly from a clean target sentence.

1. For each position i , sample a probability $p_i \sim \text{Uniform}(0, 1)$
2. If $p_i < p_{del}$, drop the word in the position i

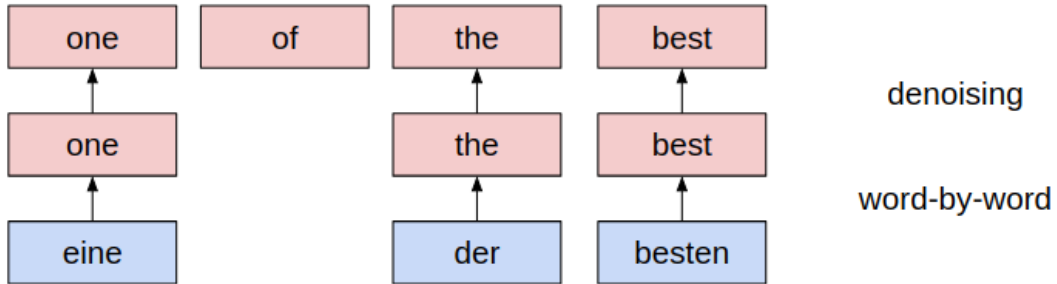


Figure 5.3: Deletion noise

Chapter 6

Experiments

6.1 Corpus Statistics

| Train | | German | English | French |
|-------|---------------|--------|---------|--------|
| | Sentences | 100M | 100M | 100M |
| | Running Words | 1880M | 2360M | 3017M |
| | Vocabulary | 1254k | 523k | 660k |

| Test | | newstest2016 | | newstest2014 | |
|------|-----------------|--------------|-------------|--------------|-------------|
| | | German | English | French | English |
| | Sentences | 2999 | 2999 | 3003 | 3003 |
| | Running Words | 62506 | 64619 | 81165 | 71290 |
| | Vocabulary Size | 11978 | 8645 | 10899 | 9200 |
| | OOV Rates | 4116 (6.6%) | 1643 (2.5%) | 1731 (2.1%) | 1299 (1.8%) |
| | LM perplexity | 211.0 | 109.6 | 51.2 | 84.6 |

Search vocabulary in testing: 50k (src/tgt)

6.2 Word Translation

6.2.1 Baseline

Corpus Size for Monolingual Embedding

| | | De-En | | | | En-De | | | |
|--------------|-----------------|-------|-----|-----|------|-------|-----|-----|------|
| | | 5M | 10M | 50M | 100M | 5M | 10M | 50M | 100M |
| Supervised | Procrustes-NN | | | | | | | | |
| | Procrustes-ISF | | | | | | | | |
| | Procrustes-CSLS | | | | | | | | |
| Unsupervised | Adv-NN | | | | | | | | |
| | Adv-CSLS | | | | | | | | |
| | Adv-Refine-NN | | | | | | | | |
| | Adv-Refine-CSLS | | | | | | | | |

Vocabulary Size for Word Translation

| Accuracy(%) | 50k | 100k | 150k | 200k | No Cutoff |
|-------------|-----|------|------|------|-----------|
| 50k | | | | | |
| 100k | | | | | |
| 150k | | | | | |
| 200k | | | | | |
| No Cutoff | | | | | |

6.2.2 Data-driven Approach

Different Learning Rate Scheduler

Different Learning Rate Scheduler

- Start from Initial Learning Rate
 - German → English

| Accuracy (%) | 100 | 500 | 1k | 2k | 5k |
|--------------|-------|-------|-------|-------|-------|
| 0.01 | 1.00 | 6.01 | 24.21 | 28.37 | 42.71 |
| 0.005 | 0.85 | 5.47 | 7.71 | 17.12 | 38.78 |
| 0.001 | 1.23 | 7.32 | 14.88 | 15.34 | 27.45 |
| 0.0005 | 1.54 | 45.26 | 49.88 | 54.43 | 54.97 |
| 0.0001 | 36.23 | 58.09 | 60.14 | 60.37 | 58.67 |

- English \rightarrow German

| Accuracy (%) | 1k | 2k | 5k |
|--------------|-------|-------|-------|
| 0.01 | 10.02 | 16.50 | 23.36 |
| 0.005 | 47.65 | 22.28 | 24.21 |
| 0.001 | 24.44 | 24.98 | 26.75 |
| 0.0005 | 54.82 | 55.67 | 48.44 |
| 0.00001 | 59.21 | 58.98 | 56.90 |

- Learning Rate Decay

- Inherit Last Learning Rate

- 50k

| Accuracy (%) | 1k | 2k | 5k |
|--------------|-------|-------|-------|
| 0.01 | 7.79 | 4.70 | 6.01 |
| 0.005 | 3.39 | 12.03 | 8.94 |
| 0.001 | 5.40 | 11.49 | 18.89 |
| 0.0005 | 5.78 | 8.33 | 10.25 |
| 0.00001 | 51.50 | 37.63 | 15.27 |

- 50k 2 epochs

| Accuracy (%) | 0.5 | 0.7 | 0.9 |
|--------------|-------|-------|-------|
| 0.01 | 3.78 | 1.39 | 1.31 |
| 0.005 | 4.39 | 3.31 | 14.73 |
| 0.001 | 5.63 | 10.56 | 13.72 |
| 0.0005 | 7.56 | 6.71 | 10.64 |
| 0.0001 | 55.51 | 56.36 | 57.59 |

- 200k

| Accuracy (%) | 0.5 | 0.7 | 0.9 |
|--------------|-------|-------|-------|
| 0.01 | 3.24 | 2.08 | 0.46 |
| 0.005 | 5.32 | 12.72 | 1.23 |
| 0.001 | 4.55 | 8.02 | 17.27 |
| 0.0005 | 6.71 | 6.94 | 6.79 |
| 0.0001 | 57.83 | 57.83 | 57.67 |

6.3 Sentence Translation

6.3.1 Comprehensive Table

| System | De-En Bleu [%] | En-De Bleu [%] | Fr-En Bleu [%] | En-Fr Bleu [%] |
|--|-------------------|-------------------|-------------------|-------------------|
| Word-by-Word | 11.1 | 6.7 | 10.6 | 7.8 |
| + LM (5-gram) + tgt w/ high LM score for OOV | 12.9 | 8.9 | 12.7 | 10.0 |
| + LM (5-gram) + copy from src for OOV | 14.5 | 9.9 | 13.6 | 10.9 |
| + Denoising (RNN) | 16.2 | 10.6 | 15.8 | 13.3 |
| + Denoising (Transformer) | 17.2 | 11.0 | 16.5 | 13.9 |
| Lample et al. [2017] | 13.3 | 9.6 | 14.3 | 15.1 |
| Artetxe et al. [2017b] | - | - | 15.6 | 15.1 |

Table 6.1: Word-by-word translation from German to English

| | Accuracy [%] | Bleu [%] |
|------|--------------|----------|
| 5M | 44.9 | 9.7 |
| 10M | 51.6 | 10.1 |
| 50M | 59.4 | 10.8 |
| 100M | 61.2 | 11.2 |

6.3.2 BPE vs Word

Byte pair encoding (BPE) is a simple data compression technique for word segmentation. It allows for the representation of an open vocabulary through a fixed-size vocabulary of variable-character sequences, making it very suitable word segmentation strategy for neural network models. It helps to reduce the vocabulary size and they eliminate the presence of unknown words in the output translation. We use BPE to represent the mapping between languages.

| Vocabulary | | Bleu [%] |
|------------|------------------------|----------|
| BPE | Merges | |
| | 20k | 10.4 |
| | 50k | 12.5 |
| | 100k | 13.0 |
| Word | Cross-lingual training | |
| | 20k | 14.4 |
| | 50k | 14.4 |
| | 100k | 14.5 |
| | 200k | 14.4 |

As the experiments results demonstrated above, the BPE performs worse than word in this unsupervised learning scenario. It might be difficult to learn the translation relationship between subword units.

6.3.3 Artificial Noise

| d_{per} | p_{del} | p_{ins} | V_{ins} | Bleu [%] |
|------------------|------------------|------------------|------------------|----------|
| 2 | | | | 14.7 |
| 3 | | | | 14.9 |
| 5 | | | | 14.9 |
| 3 | 0.1 | | | 15.7 |
| | 0.3 | | | 15.1 |
| 3 | 0.1 | 0.1 | 10 | 16.8 |
| | | | 50 | 17.2 |
| | | | 500 | 16.8 |
| | | | 5000 | 16.5 |

6.3.4 Phrase Embedding

| Vocabulary | | | No LM Bleu [%] | With LM Bleu [%] | Denoising Bleu [%] |
|------------------------|-----------|------|-------------------|---------------------|-----------------------|
| Word | | | 11.2 | 14.5 | 17.2 |
| Mikolov et al. [2013c] | threshold | 100 | 11.1 | 13.7 | 15.6 |
| | | 500 | 11.0 | 13.7 | 16.2 |
| | | 2000 | 10.7 | 14.0 | 16.5 |
| Top frequent | count | 50k | 12.0 | 15.7 | 16.8 |

6.3.5 Vocabulary Cutoff in Translation

Table 6.2: Word embedding vocabulary cut-off

| Bleu [%] | 20k | 50k | 100k |
|----------|------|-------------|------|
| 50k | 11.1 | 11.3 | 11.2 |
| 100k | 11.2 | 11.2 | 11.1 |
| 150k | 10.9 | 10.9 | - |

Table 6.3: Phrase embedding vocabulary cut-off

| Bleu [%] | 50k | 100k | 150k |
|----------|-------------|------|------|
| 50k | 11.3 | - | - |
| 100k | 11.9 | 11.9 | - |
| 150k | 12.0 | 11.9 | 11.9 |
| 200k | 12.0 | - | - |

Chapter 7

Conclusion

For unsupervised machine translation system, the context-aware beam search language model help at the lexicon choice step.

The denoising networks which aimed at the insertion/deletion/reordering noise in the word-by-word translation sentence works well for fertility and localized reordering problem. Even though BPE is known to be an effective way to overcome the rare word problem in standard NMT, BPE embedding performs worse than word embedding in our case especially when vocabulary is small.

Word-by-word translation based on cross-lingual word embedding depends highly on the frequent word mappings. We found that phrase embedding only helps in word-by-word translation with context-aware beam search, it works not as good as that under the processing of denoising autoencoder.

Appendix A

Appendix

List of Figures

| | | |
|-----|--|----|
| 2.1 | Illustration of a translation process using IBM-3 model (Koehn [2009]) | 7 |
| 2.2 | Neural MT – example of a deep recurrent architecture (Luong et al. [2015a]) | 10 |
| 2.3 | Global attention | 12 |
| 2.4 | Local attention | 12 |
| 2.5 | The Transformer model architecture (Vaswani et al. [2017]) | 14 |
| 2.6 | Key techniques for neural unsupervised MT: A) the two original monolingual data; B) Initialization, the two distribution are roughly aligned; C) Denoising autoencoder, make the distribution closer to normal one in corresponding language D) Back-translation, improve the both translation model in iterative way (Lample et al. [2018]) . . | 17 |
| 3.1 | Global attention model (Mikolov et al. [2013a]) | 19 |
| 3.2 | A cross-lingual embedding space between German and English (Ruder et al. [2017]) | 22 |
| 3.3 | Cross-lingual projection with CCA (Faruqui and Dyer [2014]) | 23 |
| 5.1 | Reordering noise | 31 |
| 5.2 | Insertion noise | 32 |
| 5.3 | Deletion noise | 33 |

List of Tables

| | | |
|-----|---|----|
| 6.1 | Word-by-word translation from German to English | 38 |
| 6.2 | Word embedding vocabulary cut-off | 40 |
| 6.3 | Phrase embedding vocabulary cut-off | 40 |

Bibliography

- Oliver Adams, Adam Makarucha, Graham Neubig, Steven Bird, and Trevor Cohn. Cross-lingual word embeddings for low-resource language modeling. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, volume 1, pages 937–947, 2017.
- Mikel Artetxe, Gorka Labaka, and Eneko Agirre. Learning bilingual word embeddings with (almost) no bilingual data. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 451–462, 2017a.
- Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. Unsupervised neural machine translation. *arXiv preprint arXiv:1710.11041*, 2017b.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. Enriching word vectors with subword information. *arXiv preprint arXiv:1607.04606*, 2016.
- Hailong Cao, Tiejun Zhao, Shu Zhang, and Yao Meng. A distribution-based model to learn bilingual word embeddings. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 1818–1827, 2016.
- Yong Cheng, Wei Xu, Zhongjun He, Wei He, Hua Wu, Maosong Sun, and Yang Liu. Semi-supervised learning for neural machine translation. *arXiv preprint arXiv:1606.04596*, 2016.
- Trevor Cohn and Mirella Lapata. Machine translation by triangulation: Making effective use of multi-parallel corpora. In *Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics*, pages 728–735, 2007.
- Alexis Conneau, Guillaume Lample, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. Word translation without parallel data. *arXiv preprint arXiv:1710.04087*, 2017.

- Jocelyn Coulmance, Jean-Marc Marty, Guillaume Wenzek, and Amine Benhalloum. Trans-gram, fast cross-lingual word-embeddings. *arXiv preprint arXiv:1601.02502*, 2016.
- Paramveer Dhillon, Dean P Foster, and Lyle H Ungar. Multi-view learning of word embeddings via cca. In *Advances in neural information processing systems*, pages 199–207, 2011.
- Long Duong, Hiroshi Kanayama, Tengfei Ma, Steven Bird, and Trevor Cohn. Learning crosslingual word embeddings without bilingual corpora. *arXiv preprint arXiv:1606.09403*, 2016.
- Manaal Faruqui and Chris Dyer. Improving vector space word representations using multilingual correlation. In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 462–471, 2014.
- Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. *arXiv preprint arXiv:1705.03122*, 2017.
- Stephan Gouws, Yoshua Bengio, and Greg Corrado. Bilbowa: Fast bilingual distributed representations without word alignments. In *International Conference on Machine Learning*, pages 748–756, 2015.
- Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tieyan Liu, and Wei-Ying Ma. Dual learning for machine translation. In *Advances in Neural Information Processing Systems*, pages 820–828, 2016.
- Yedid Hoshen and Lior Wolf. An iterative closest point method for unsupervised word translation. *CoRR*, abs/1801.06126, 2018. URL <http://arxiv.org/abs/1801.06126>.
- Armand Joulin, Piotr Bojanowski, Tomas Mikolov, and Edouard Grave. Improving supervised bilingual mapping of word embeddings. *CoRR*, abs/1804.07745, 2018. URL <http://arxiv.org/abs/1804.07745>.
- Yunsu Kim, Julian Schamper, and Hermann Ney. Unsupervised training for large vocabulary translation using sparse lexicon and word classes. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, volume 2, pages 650–656, 2017.
- Philipp Koehn. *Statistical machine translation*. Cambridge University Press, 2009.
- Guillaume Lample, Ludovic Denoyer, and Marc’Aurelio Ranzato. Unsupervised machine translation using monolingual corpora only. *arXiv preprint arXiv:1711.00043*, 2017.

- Guillaume Lample, Myle Ott, Alexis Conneau, Ludovic Denoyer, and Marc'Aurelio Ranzato. Phrase-based & neural unsupervised machine translation. *arXiv preprint arXiv:1804.07755*, 2018.
- Ang Lu, Weiran Wang, Mohit Bansal, Kevin Gimpel, and Karen Livescu. Deep multilingual correlation for improved word embeddings. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 250–256, 2015.
- Minh-Thang Luong, Hieu Pham, and Christopher D Manning. Effective approaches to attention-based neural machine translation. *arXiv preprint arXiv:1508.04025*, 2015a.
- Thang Luong, Hieu Pham, and Christopher D Manning. Bilingual word representations with monolingual quality in mind. In *Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing*, pages 151–159, 2015b.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013a.
- Tomas Mikolov, Quoc V Le, and Ilya Sutskever. Exploiting similarities among languages for machine translation. *arXiv preprint arXiv:1309.4168*, 2013b.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013c.
- Malte Nuhn and Hermann Ney. Em decipherment for large vocabularies. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, volume 2, pages 759–764, 2014.
- Malte Nuhn, Arne Mauser, and Hermann Ney. Deciphering foreign language by combining language models and context vectors. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1*, pages 156–164. Association for Computational Linguistics, 2012.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543, 2014.
- Sujith Ravi and Kevin Knight. Deciphering foreign language. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 12–21. Association for Computational Linguistics, 2011.

- Sebastian Ruder, Ivan Vulić, and Anders Søgaard. A survey of cross-lingual word embedding models. *arXiv preprint arXiv:1706.04902*, 2017.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112, 2014.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.
- Ivan Vulic and Anna-Leena Korhonen. On the role of seed lexicons in learning bilingual word embeddings. 2016.
- Chao Xing, Dong Wang, Chao Liu, and Yiye Lin. Normalized word embedding and orthogonal transform for bilingual word translation. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1006–1011, 2015.
- Meng Zhang, Yang Liu, Huanbo Luan, and Maosong Sun. Adversarial training for unsupervised bilingual lexicon induction. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 1959–1970, 2017.