低資源言語におけるニューラル機械翻訳の 精度向上のための文ベクトルの変換

Transformation of sentence vector to improve translation quality of neural machine translation in low-resource language

北海学園大学 大学院工学研究科 教授

或前谷 博



1996 年北海学園大学大学院工学研究科修士課程修了。博士(工学)。2013 年~現在北海学園大学大学院工学研究科教授。機械翻訳の研究に従事。アジア太平洋機械翻訳協会(AAMT)/ Japio 特許翻訳研究会委員。

echi@hgu.jp

011-841-1161(内線:7863)

(1) はじめに

ニューラル機械翻訳や大規模言語モデルの進展により、 機械翻訳の精度が向上している。特に日本語英語間等の 資源の豊富な言語間の翻訳精度の向上は著しい。日本語 英語間のような言語間では対訳コーパス等の言語資源を 比較的容易に収集できるため、大規模な学習データを用 いた精度の高い翻訳モデルの構築が可能である。

それに対して、言語資源の乏しい低資源言語において は十分な学習データの収集が困難なため、構築するモデ ルの精度も不十分となる。このような低資源言語に対す るニューラル機械翻訳の精度向上のためのアプローチと しては、データ拡張のアプローチと複数の言語モデルを 利用するアプローチに大きく分類される。データ拡張 のアプローチでは、例えば Zhang ら [1] は機械翻訳シ ステムにより得られた訳文を原言語文に対する目的言 語文とすることで疑似対訳コーパスを生成した。また、 Sennrich ら ^[2] は元の対訳コーパスに逆翻訳により生 成した疑似対訳コーパスを加えることで翻訳精度の向上 を試みた。さらに、疑似対訳コーパスを大規模言語モ デルにより生成する手法 [3] も提案されている。しかし、 このような疑似対訳コーパスを用いるアプローチではそ の精度が問題となり、翻訳精度に必ずしも良い影響を与 えるとは限らない。特に低資源言語においてはニューラ ル機械翻訳や大規模言語モデルより得られる訳文の精度 は低くなる可能性が高く、その結果、翻訳精度の向上の 妨げとなる。

また、複数の言語モデルを利用するアプローチとして、

Zoph ら [4] は低資源言語を用いて構築した子モデルの性能向上のために豊富な言語資源より構築された親モデルのパラメータに基づき子モデルの初期化を行った。しかし、この場合、子モデルと親モデルの構築に用いるデータの言語は類似していることが前提となり、類似性が低い言語間においては親モデルのパラメータは子モデルにおいて有効に機能しない可能性がある。

そこで本報告では、低資源言語においては逆翻訳により得られる疑似対訳コーパスの精度が著しく低下する可能性を考慮し、複数の言語モデルを用いるアプローチによる、新たなニューラル機械翻訳を提案する。提案手法では、Transformer^[5] に基づくニューラル機械翻訳として提案された BERT-fused NMT^[6] において、事前学習したモデルが生成する文ベクトルを Encoder の単語埋め込みベクトルとして用いる。このモデルは原言語の文ベクトルに対応する目的言語の文ベクトルの情報を有した文ベクトルを生成するため、翻訳精度の向上が期待される。本報告では、提案手法とその有効性をアイヌ語日本語間、ベトナム語日本語間、そして、アッサム語英語間による翻訳実験に基づき述べる。

(2) 提案手法の概要

提案手法に基づくニューラル機械翻訳の概要図を図1 に示す。提案手法によるニューラル機械翻訳では原言 語文の文ベクトルを目的言語文の文ベクトルとの間で 対応関係にある文ベクトルに変換し、それを単語埋め 込みベクトルとして用いる。まず、原言語文を多言語の

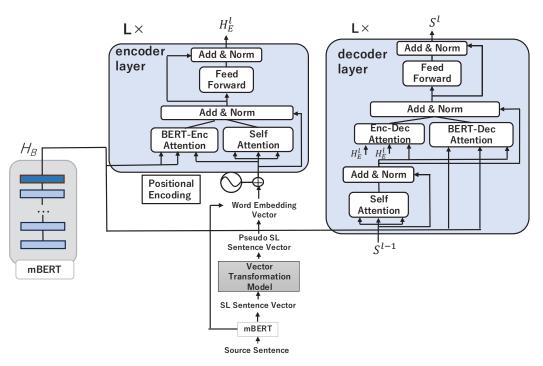


図1 提案手法に基づくニューラル機械翻訳 Vec-trans NMT の概要図

Encoder モデルである mBERT[7] を用いて文ベクトル(以 降、原言語文ベクトルと呼ぶ)を得る。原言語文ベクト ルは mBERT により得られる単語ベクトルの平均を求め ることで得られる。次いで、事前学習されたニューラル ネットワークモデルを用いて目的言語の文ベクトル(以降、 目的言語文ベクトルと呼ぶ)に対応する疑似的な原言語文 ベクトル(以降、疑似原言語文ベクトルと呼ぶ)を生成す る。そして、その疑似原言語文ベクトルを単語埋め込み ベクトルとして BERT-fused NMT の Encoder に用い る。文ベクトルから単語ベクトルの変換処理は単語ベク トルの要素ごとの平均が文ベクトルの要素と一致するよう に単語ベクトルの各要素に任意の固定値を付与することで 行う。本報告では疑似原言語文ベクトルを得るための事 前学習済みのニューラルネットワークモデルをベクトル 変換モデルと呼ぶ。また、このベクトル変換モデルを用 いたニューラル機械翻訳を Vec-trans NMT と記す。

3 ベクトル変換モデル

3.1 ベクトル変換モデルの構築の概要

ベクトル変換モデルは原言語文ベクトルとその目的言語文ベクトル、そして、目的言語文に対する評価スコアを用いて学習を行う。原言語文ベクトルと目的言語文ベクトルは共に mBERT より得られる単語べ

クトルの平均ベクトルを求めることで得る。評価スコアは mSBERT^[8] より得られる原言語文ベクトルと目的言語文ベクトルとの間のコサイン類似度を用いる。mSBERT は文ベクトル間の類似度を求めるためにシャムネットワークにより学習された SBERT^[9] に対して、知識蒸留により多言語に拡張した Encoder モデルである。ベクトル変換モデルの目的は入力された原言語文ベクトルとその目的言語文ベクトルの間のコサイン類似度が評価スコアと等価になるように原言語文ベクトルを変換することにある。

ここで対応関係にある原言語文ベクトルとその目的言語文ベクトル間のコサイン類似度は基本的には最大値の1.0に近しい値になることが理想ではあるが、mBERTは文ベクトル間の類似度を得ることを目的とはしていないため、文ベクトル間の意味的な類似度を得ることは困難と考えられる。そこで、ベクトル変換モデルにおいて目的言語文ベクトルとの間で意味的な類似度が得られるように原言語文ベクトルを新たな文ベクトル、すなわち、疑似原言語文ベクトルに変換する。そして、その文ベクトルをBERT-fused NMTの単語埋め込みベクトルとして利用することで翻訳精度の向上を図る。

また、学習時に必要な正解データには目的言語文ベクトルとの間のコサイン類似度が評価スコアと等価になるように変換された原言語文ベクトル(以降、正解文ベク

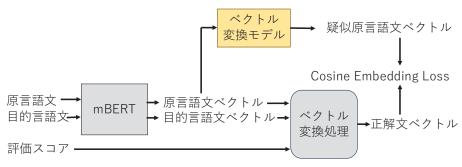


図2 ベクトル変換モデルの構築の概要

トルと呼ぶ)を用いる。この正解文ベクトルの生成方法については次節の3.2で詳細を述べる。正解文ベクトルとベクトル変換モデルより得られる疑似原言語文ベクトルとの間のLossの計算はCosine Embedding Lossを用いる。図2にベクトル変換モデルの構築の概要を示す。

ベクトル変換モデルは原言語文ベクトルを目的言語文ベクトルと対応関係にある疑似原言語文ベクトルに変換する。変換された疑似原言語文ベクトルに対する正解文ベクトルはベクトル変換処理より原言語文ベクトルと目的言語文ベクトルとの間のコサイン類似度が評価スコアと等価になるように原言語文ベクトルを変換した文ベクトルである。この正解文ベクトルと疑似原言語文ベクトルとの間で Cosine Embedding Loss に基づきベクトル変換モデルの学習を行う。

3.2 正解文ベクトルの生成

ベクトル変換処理における正解文ベクトルの生成について述べる。ベクトル変換モデルの学習に用いる正解文ベクトルは原言語文ベクトルの任意の 1 つの要素に変数 x を加えた際、目的言語文ベクトルと変数 x を付与した原言語文ベクトルの間のコサイン類似度が評価スコアとなるような x を求め、その値を加えることにより得る。変数 x を付与した n 次元の原言語文ベクトルを S=(s1, …,si+x,…,sn)、目的言語文ベクトルを T=(t1,…,tn) とした場合のコサイン類似度の計算式を以下の式 (1) に示す。

$$\frac{\sum_{j=1}^{i-1} (s_j \times t_j) + (s_i + x) \times t_i + \sum_{j=i+1}^n (s_j \times t_j)}{\sqrt{\sum_{j=1}^{i-1} s_j^2 + (s_i + x)^2 + \sum_{j=i+1}^n s_j^2}} = 評価スコア$$
式 (1)

式 (1) では si に変数 x が付与されている。この計算式は x の 2 次方程式になるため、x の解を求めること

で値が得られ、xの値を Si に加算した原言語文ベクトルが正解文ベクトルとなる。

また、式(1)の評価スコアは先に述べたように Encoder モデルである mSBERT より原言語文ベクトルと目的言語文ベクトルの間のコサイン類似度を求める ことで得られるが、mSBERT に含まれない低資源言語 が対象となった場合には代替言語を利用する。代替言語には低資源言語と類似した言語を用い、資源が豊富な言語文を Google 翻訳より代替言語文に翻訳した文を mSBERT に入力することで代替言語文に対応する文ベクトルを生成する。その概要を図3に示す。

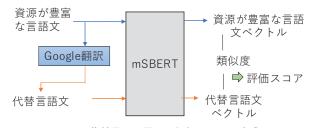


図3 代替言語を用いた評価スコアの生成

低資源言語が mSBERT の学習データに含まれていない場合、その言語の文ベクトルを得ることは困難となるため、低資源言語と文法的に近く、mSBERT の学習データに含まれている言語を用いる。具体的には、Google翻訳を用いて資源が豊富な言語文を翻訳し、代替言語文を得る。そして、それを mSBERT より代替言語文ベクトルに変換し、資源が豊富な言語文ベクトルとの間のコサイン類似度を得る。ただし、得られたコサイン類似度が閾値以上の場合にのみベクトル変換モデルの学習データである評価スコアに用いる。その結果、より精度の高いベクトル変換モデルが得られると考えられる。

4

性能評価実験

本報告では提案手法の有効性を確認するためにアイヌ 語日本語間、ベトナム語日本語間、そして、アッサム語 英語間の翻訳実験を行った。

4.1 実験データ

実験データはベクトル変換モデル及び Vec-trans NMT を構築するためのデータに大きく分かれる。さらに Vec-trans NMT においてはモデル構築のための学習データと翻訳に用いる評価データの2種類に分類される。

ベクトル変換モデルの構築においては mSBERT よ り得られる原言語文ベクトルと目的言語文ベクトルとの 間のコサイン類似度の値が閾値を超えている対訳文のみ を学習データに用いた。具体的には、アイヌ語日本語間 の対訳文においては 72,936 対 [10] から閾値 0.95 を 超えていた対訳文 5,787 対を学習データに用いた。そ の際の訓練データ数は 4,787、検証データ数は 1,000 である。また、アイヌ語は mSBERT の学習データに 含まれていないため、アイヌ語に対する代替言語として モンゴル語、トルコ語、そして、韓国語を用いた。評価 スコアはこれら3つの言語それぞれの文ベクトルと日 本語の文ベクトル間のコサイン類似度の平均値を求め、 その平均値が閾値を超えている対訳文をベクトル変換モ デルの学習データに使用した。ベトナム語日本語におい ては対訳文 16.081 対 [11] から文ベクトル間のコサイ ン類似度が閾値 0.9 を超えている対訳文 6,543 対を 学習データに用いた。その際の訓練データ数は5.543、 検証データ数は 1,000 である。そして、アッサム語 英語においては対訳文 54,000 対 [12] から文ベクトル 間のコサイン類似度が閾値 0.9 を超えている対訳文 7,988 対を学習データとしてベクトル変換モデルを構 築した。その際の訓練データ数は6,988、検証データ 数は 1,000 である。また、アッサム語は mSBERT に 存在しないため、アッサム語に対する代替言語としては ヒンディー語、フランス語、そして、ギリシャ語を用い て評価スコアを得た。

ニューラル機械翻訳 Vec-trans NMT の構築においてはアイヌ語日本語では上述した対訳文を含む対訳文91,168 対、ベトナム語日本語では対訳文20,202 対をそれぞれ8:1:1 の割合で訓練データ、検証データ、

そして、評価データに分割して、学習と翻訳を行った。 アッサム語英語においては訓練データ数 50,000、検 証データ数 2,000、そして、評価データ数 2,000 を 用いて学習と翻訳を行った。

4.2 実験方法

実験は BERT-fused NMT をベースラインとして、 提案手法の Vec-trans NMT との比較により行った。 BERT-fused NMT、Vec-trans NMT 共に言語ペア間 の言語方向毎にモデルを構築した。したがって、提案手 法においては3つの言語ペアそれぞれにおいて2つのベクトル変換モデルと2つの Vec-trans NMT のモデルを 構築し、翻訳を行った。評価方法は全て自動評価尺度の SacreBLEU を用いた。COMET^[13]を始めとしたニューラルネットワークベースの自動評価尺度については学習 が行われていない低資源言語の翻訳文に対する評価が不 十分となることが予想されるため使用しなかった。

また、ベクトル変換モデルのアーキテクチャは多層 パーセプトロンとなっており、768 ニューロンの入力 層、1,536 ニューロンの中間層、768 ニューロンの中間層、そして、768 ニューロンの出力層で構成されている。入力層と出力層のニューロンが 768 となっている のは mSBERT による文ベクトルの次元数である 768 と一致させるためである。

4.3 実験結果

表 1 に提案手法による Vec-trans NMT とベースラインである BERT-fused NMT におけるアイヌ語日本語間の BLEU スコアを示す。また、表 2 にベトナム語日本語間の BLEU スコア、そして、表 3 にアッサム語英語間の BLEU スコアを示す。

表 1 アイヌ語日本語の BLEU スコア

	アイヌ語→日本語	日本語→アイヌ語
ベースライン	24. 89	32. 20
Vec-trans NMT	24. 52	33. 36

表 2 ベトナム語日本語の BLEU スコア

	ベトナム語→日本語	日本語→ベトナム語
ベースライン	8. 69	8.60
Vec-trans NMT	9. 42	7. 79



表 3 アッサム語英語の BLEU スコア

	アッサム語→英語	英語→アッサム語
ベースライン	5. 56	14. 34
Vec-trans NMT	5. 09	15. 50

また、ベースラインと提案手法におけるアイヌ語日本 語間の翻訳文の例を表 4 に示す。

表 4 翻訳文の具体例

アイヌ語→日本語			
原言語文	situ turasi e = oman wa		
正翻訳	尾根に沿って川の上流に行って		
ベースライン	尾根の尾根に沿って行って		
Vec-trans NMT	尾根に沿って登って行き		
日本語→アイヌ語			
原言語文	もう山の上に目が傾いた様子で		
正翻訳	tane nupuri ka cup rari kane siran ora		
ベースライン	tane nupuri ka cup rari kane sirki		
Vec-trans NMT	tane nupuri ka cup rari kane sir'an		
	hine		

4.4 考察

表 1 のアイヌ語日本語においては日本語からアイヌ語の翻訳において BLEU スコアが 32.20 から 33.36 へと比較的大きく向上している。表 2 のベトナム語日本語においてはベトナム語から日本語の BLEU スコアは向上しているが、日本語からベトナム語の BLEU スコアは逆に低下している。しかし、ベトナム語日本語においては、いずれの BLEU スコアも非常に低く、比較するには不十分な結果であったと考えられる。その原因としては学習データが非常に少なく、十分な学習が行われていないためと考えられる。表 3 のアッサム語英語においては、アッサム語から英語の BLEU スコアは低下しているが、英語からアッサム語の BLEU スコアについては 1 ポイント以上向上している。

これらの結果より、アイヌ語日本語とアッサム語英語における低資源言語方向の翻訳においては提案手法により BLEU スコアが大きく向上したことが確認できた。これは提案手法によるベクトル変換モデルは mBERT より得られるベクトルを用いて学習していることから資源の豊富な言語を変換した方が精度の高い文ベクトルを生成できるためと考えられる。そこで、ベクトル変換モ

デルより得られる疑似原言語文ベクトルが目的言語文ベクトルとどの程度類似したものになっているのか調査を 行った。調査方法の概要を図4に示す。

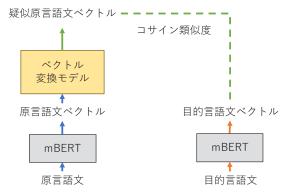


図4 ベクトル変換モデルの精度調査の概要

図4においてmBERTより得られる原言語文ベク トルと目的言語文ベクトル間においては先述したように コサイン類似度を求めても意味的な類似性を得ることは 困難であり、コサイン類似度の精度は低いと考えられ る。それに対して、ベクトル変換モデルより得られる疑 似原言語文ベクトルは目的言語文ベクトルとの間のコサ イン類似度が評価スコアと等価になるように学習されて いるため、目的言語文ベクトルとの間のコサイン類似度 は比較的精度の高い値が得られると考えられる。そこ で、そのことを確認するためにアイヌ語日本語間、ベト ナム語日本語間、そして、アッサム語英語間においてべ クトル変換モデルより得られた疑似原言語文ベクトルと 目的言語文ベクトルとの間のコサイン類似度を求め、比 較を行った。なお、調査に用いた対訳データは3つの 言語ペアのいずれも評価データよりランダムに抽出した 100 データであり、それらより得られた 100 のコサ イン類似度の平均値を求めた。表 5 に疑似原言語文べ クトルと目的言語文ベクトルとの間のコサイン類似度の 平均値を示す。

表 5 文ベクトル間のコサイン類似度の平均値

原言語	目的言語	コサイン類似度の平均値
アイヌ語	日本語	0. 7660
日本語	アイヌ語	0. 7811
ベトナム語	日本語	0. 8563
日本語	ベトナム語	0. 8232
アッサム語	英語	0. 6085
英語	アッサム語	0. 8595

表5において言語ペア毎に双方向でコサイン類似度 の平均値を比較すると学習が不十分なベトナム語日本語 間を除き、アイヌ語日本語間、アッサム語英語間いず れも資源の豊富な言語から低資源言語方向の翻訳のコ サイン類似度の平均値が高いことを確認できる。アイ ヌ語日本語間では日本語から低資源言語のアイヌ語方 向の翻訳時のコサイン類似度の平均値は 0.7811 であ り、アイヌ語から日本語方向の翻訳時の 0.7660 に比 べて高い。また、アッサム語英語間では英語から低資源 言語のアッサム語方向の翻訳時のコサイン類似度の平 均値は 0.8595 であり、アッサム語から英語方向の翻 訳時の 0.6085 に比べて高い。このような調査結果は 表 1 と表 3 の資源の豊富な言語から低資源言語方向の BLEU スコアが低資源言語から資源の豊富な言語方向 の BLEU スコアよりも高くなっている結果と同様であ る。したがって、提案手法におけるベクトル変換モデル の精度向上が翻訳精度、すなわち、BLEUスコアの向 上に直結していることが確認できた。低資源言語から資 源の豊富な言語方向の翻訳精度を向上させるためには、 ベクトル変換モデルによる疑似原言語文ベクトルの生成 の精度向上が重要と考えられる。

5 まとめ

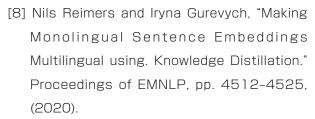
本報告では複数の言語モデルを用いるアプローチの観点より、ニューラル機械翻訳の精度向上を目的とした新たな手法を提案した。提案手法では Encoder モデルより得られる原言語文ベクトルを目的言語文ベクトルの情報を含む疑似原言語文ベクトルに変換するためにベクトル変換モデルを構築した。そして、生成された疑似原言語文ベクトルをニューラル機械翻訳の単語埋め込みベクトルとして利用した。性能評価実験の結果、疑似原言語文ベクトルを利用することで資源の豊富な言語から低資源言語方向の翻訳において、BLEU スコアがベースラインのBLEU スコアよりも向上し、提案手法の有効性を確認した。

今後は、低資源言語から資源の豊富な言語方向の翻訳 精度の向上のためにベクトル変換モデルの改良を行う。 具体的には学習時に用いる Encoder モデルを低資源言語 に対しても有効利用できるようにファインチューニング を行い、ベクトル変換モデルの精度向上を図る予定である。

参考文献

- [1] Jiajun Zhang and Chengqing Zong, "Exploiting Source-side Monolingual Data in Neural Machine Translation," Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pp. 1535–1545 (2016)
- [2] Rico Sennrich, Barry Haddow and Alexandra Birch, "Improving Neural Machine Translation Models with Monolingual Data," Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, pp. 86-96 (2016)
- [3] Seokjin Oh, Su Ah Lee, and Woohwan Jung. "Data Augmentation for Neural Machine Translation using Generative Language Model," arXiv preprint arXiv:2307.16833 (2023)
- [4] Barret Zoph, Deniz Yuret, Jonathan May and Kevin Knight, "Transfer Learning for Low-Resource Neural Machine Translation," Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pp.1568-1575. (2016)
- [5] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention Is All You Need", Proceedings of the 31st Conference on Neural Information Processing Systems, pp. 6000-6010, (2017).
- [6] Jinhua Zhu, Yingce Xia, Lijun Wu, Di He, Tao Qin, Wengang Zhou, Houqiang Li, and Tie-Yan Liu, "Incorporating BERT into Neural Machine Translation," Proceedings of 8th ICLR, pp. 1-18, (2020).
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," Proceedings of NAACL-HLT 2019, pp. 4171-4186, (2019).





- [9] Nils Reimers and Iryna Gurevych, "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks," Proceedings of EMNLP-IJCNLP 2019, pp.3982-3992, (2019).
- [10] 田中蒼大郎, 越前谷博, 荒木健治, "Transformer を用いたアイヌ語 日本語間機械翻訳における精度 と対訳コーパスの関係," 令和4年度電気・情報関係学会北海道支部連合大会講演論文集, pp. 227-228, (2022).
- [11] Ye Kyaw Thu, Win Pa Pa, Masao Utiyama, Andrew Finch, and Eiichiro Sumita, "Introduction of the Asian Language Treebank(ALT)," Proceedings of the Tenth International Conference on Language Resources and Evaluation, pp.1574-1578, (2016).
- [12] Santanu Pal, Partha Pakray, Sahinur Rahman Laskar, Lenin Laitonjam, Vanlalmuansangi Khenglawt, Sunita Warjri, Pankaj Kundan Dadure and Sandeep Kumar Dash, "Findings of the WMT 2023 Shared Task on Low-Resource Indic Language Translation," Proceedings of the Eighth Conference on Machine Translation (WMT), pp. 682-694, (2023).
- [13] Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie, "COMET: A Neural Framework for MT Evaluation," Proceedings of EMNLP, pp. 2685–2702, (2020).

