ニューラル機械翻訳におけるコーパスフィルタリング に関する固有表現に注目した分析

Corpus Filtering Focusing on Named Entities for Neural Machine Translation

本間広樹 Hiroki Homma

山岸駿秀 Hayahide Yamagishi 松村雪桜 Yukio Matsumura 小町守 Mamoru Komachi

首都大学東京

Tokyo Metropolitan University

Some parallel corpora include sentences that disturb learning of machine translation systems. By removing such noisy sentences like containing many out-of-vocabulary from the training corpus, it is expected to makes better translations. In this paper, we focus on the sentences containing named entities because most of the named entities fall into out-of-vocabulary due to low-frequencies. We propose two kinds of filtering methods, using byte pair encoding and using named entity recognition. By removing noisy sentences from a training corpus on Japanese-English language pair, BLEU scores improve statistically significantly by 0.5 points in both proposed methods. Analysis revealed that both our methods overcome the mistakes such as suffix of the noun, determiner, and sentence lengths.

1. はじめに

2013年に提案されたニューラル機械翻訳 [Kalchbrenner 13] は現在,コンピュータを用いて自動的に翻訳を行う機械翻訳の 主流の手法となっている。自然言語処理の分野では学習用の データに含まれない単語を未知語と呼び,機械翻訳システムな どの言語生成システムでは一般的に,未知語を生成できないと いう問題がある.さらに,ニューラル機械翻訳は計算量の関係か ら実際に使用する語彙サイズを絞る必要があり,既知の単語で あっても低頻度のものは未知語と同様に扱われてしまう.これ らの,システムが扱うことのできない単語を out-of-vocabulary (OOV)と呼ぶ.

ー般的な単語は、構成性のある単語と構成性のない単語の 2 種類に分類される.ここで構成性のある単語とは、語義を 持つサブワードに分割することが可能な単語を指す.例えば unkingly という単語は構成性があり、un(~でない),king (王),ly(~らしい)の3つのサブワードに分割された場合、 「王らしくない」という意味を推測できる可能性がある.逆に 構成性のない単語は語義のあるサブワードに分割できない単語 を指す.例えば Maldives という単語は構成性を持たず、この ような単語が OOV であった場合は元の意味を復元するのは困 難である.

このように、単語をさらに小さなサブワードという単位に分 割することで OOV の問題を解決する手法として Byte Pair Encoding (BPE) モデル [Sennrich 16] が提案された. この手法 では OOV を高頻度の細かいサブワード,最も細かいもので文 字単位に分割することで,語彙サイズの制限による OOV の問 題を解消する.

また,機械翻訳の手法の多くはモデルの学習のために原言 語文と目的言語文のペアを集めたパラレルコーパスを必要と する.特にニューラル機械翻訳においては他の手法と比較して 大規模なデータを要する.現在公開されているパラレルコーパ スの多くは,複数の言語で書かれた論文や新聞記事などの中か ら,対応の取れる文対を機械的に集めることで対訳文としてい る.しかし,正しい対訳になっている文対であってもその文対 の必要性について検討されていない場合も多く,これらのコー

連絡先:本間広樹, homma-hiroki@ed.tmu.ac.jp

パスは翻訳モデルの学習に悪影響を及ぼす文を含有している ことが考えられる.このような文を本稿では以降ノイズ文と呼 ぶ.パラレルコーパスにさまざまなノイズを与えてその影響を 調べた研究 [Khayrallah 18] ではニューラル機械翻訳は以前の 主流な機械翻訳手法である統計的機械翻訳に比べてノイズに 弱いことがわかった.ノイズ文を学習データから取り除くこと で,より精度の高い機械翻訳システムが構築されることが見込 まれる.

固有名詞などに代表される固有表現は一般的な語彙に比べて 出現頻度が低い.このため,固有表現を多く含む文は計算量の 制約で語彙サイズが制限されるニューラル機械翻訳において, 多数の OOV を持つことになる.OOV の占める割合が高い文 は,文としての意味が大きく欠けている場合が多く,入力文と 参照文のどちらに用いても翻訳モデルの学習に対して悪影響を 与えると思われる.また,前処理として BPE などのサブワー ド化を行う場合を考えると,出現頻度の低い固有表現はある程 度の長さのある高頻度のサブワードで構成されていないことが 多く,他の構成性のある単語よりも細かいサブワードに分割さ れる.このように細かく分割された固有表現を多く含む文は翻 訳結果の改善にほとんど寄与せず,逆に文長の増加や語義を持 たないサブワードの増加により正しい翻訳を妨げる学習につな がってしまうと考えられる.

これらのことから,固有表現を多く含む文を学習データか ら取り除くことで,より正しい翻訳が行われる機械翻訳システ ムが構築されることが見込まれる.この仮定のもと,本研究で は2種類のコーパスフィルタリング手法を提案する.一つは 固有表現抽出を用いたフィルタリング(NERF)で,文中に含 まれる固有表現数が多い文を取り除く手法である.この手法は 直接固有表現抽出ツールを用いる必要があるため,ツールが存 在しない言語にも容易に適用できる手法として,BPEを用い たフィルタリング(BPEF)も考案した.これはサブワード化 したときに文長が大きく増加する文,つまり細かいサブワード に分割される単語を多く含む文を取り除く手法である.

本研究では、これら2種類の手法でパラレルコーパスから ノイズ文を取り除く実験とその分析を行った。日英翻訳の実験 ではノイズ文を取り除くことで、どちらの手法でも機械翻訳 の自動評価尺度である BLEU スコアが 0.5 ポイント改善した. 同実験において人手評価も実施し,NERF に関してはスコア の向上が見られた.分析から,BLEU スコアを上昇させる要 因として複数形の有無や冠詞の間違いなどが改善されたほか, 文長がより正しく学習できていることが分かった.

2. フィルタリング手法

パラレルコーパス内に含まれるノイズ文を取り除く方法と して,固有表現抽出を用いる手法とBPEを用いる手法の2種 類を考える.

2.1 NERF: 固有表現抽出によるフィルタリング

固有表現抽出はテキストから固有表現を抽出するタスクで ある.現在では計算機によって大量のテキストから固有表現を 自動的に抽出する研究が一般的であり、ツールとして様々なも のが公開されている.固有表現抽出手法は数多く提案されてい るが、本研究では機械学習を用いた系列ラベリングによる固有 表現抽出の結果を用いる.

固有表現抽出による学習データのフィルタリングは以下の流 れで行う.

- 1. 固有表現抽出を用いて原言語側または目的言語側の学習 データの各文で固有表現数をカウント
- 2. 固有表現数の多い上位 n 文の文番号を取得*1
- 3. 2. で取得した文番号の文を原言語側と目的言語側の両方 から削除

2.2 BPEF: BPE によるフィルタリング

Byte Pair Encoding (BPE)[Sennrich 16] は、低頻度の単 語を高頻度のサブワードに分割することで機械翻訳の OOV の 処理を行う手法であり、現在、ニューラル機械翻訳では一般的 に使用されている. BPE の計算方法を以下に示す.

- 1. 単語分割された文の集合 *2 を入力 "Quickly, accurately, and, effectively"
- 2. 文の集合に含まれる各単語に対して文字を要素としたリスト化
 [Q, u, i, c, k, l, y, @], [a, c, c, u, r, a, t, e, l, y, @],
- [a, n, d, @], [e, f, f, e, c, t, i, v, e, l, y, @]*3
 3. リストの要素 2-gram の頻度を計算
- $\{ly: 3, y@: 3, Qu: 1, ui: 1, ic: 1, ... \}$
- 最も高頻度の2要素を結合してリストを更新 [Q, u, i, c, k, ly, @], [a, c, c, u, r, a, t, e, ly, @], [a, n, d, @], [e, f, f, e, c, t, i, v, e, ly, @]
- 語彙サイズを決定するパラメータη回だけ 3. と 4. を繰り返し実行 [Quick, ly@], [accurat, ely@], [and@], [effectiv, ely@]
- 6. 使用した全ての文字とサブワードを語彙として使用 {Q, u, i, c, k, ..., ly, ly@, ely@, ..., accurat, effectiv}
- *1 nは削除文数を決定するパラメータ

ここで,リストの要素 2-gram とはリストに含まれる要素を隣 同士で結合したものである。例えば,リスト [*a*, *b*, *c*] の要素 2-gram は {*ab*, *bc*} となる。BPE によって使用する語彙サイ ズが削減され,低頻度の単語は複数のサブワード*⁴ に分割さ れるようになる。BPE によってサブワード化された文の例を 以下に示す。

サブワード化前(5 トークン)

HDPE における メルトフラクチャー 特性 。

サブワード化後(8 トークン) HDPE における メルト フラ クチャ ー 特性 。

BPE の結合回数を語彙サイズの制限以下にすることで全て の単語をサブワードで表現できるため,OOV がなくなる.そ の一方で,構成性のない単語は高頻度の大きなサブワードに分 割できず,そのような単語を多く含む文は,とても細かく分割 されたサブワードの増加と文長の増加により学習に悪影響を与 える文になる可能性がある.つまり,BPEの適用で文長が大 きく増加した文を学習データから取り除くことで,NERFと 同様により正しい翻訳が行われる機械翻訳システムが構築され ることが見込まれる.この仮定のもと BPEF では BPE の適 用で文長が大きく増加した文を学習データから取り除く.

BPE を用いた学習データのフィルタリングは以下の流れで 行う.

- 原言語側または目的言語側の学習データの各文に対し、 BPEを適用する前のトークン数と適用した後のトークン 数をそれぞれ w と t としてカウント
- 2. *t*/*w* > θ を満たす文の文番号を取得
- 3. 2. で取得した文番号の文を原言語側と目的言語側の両方 から削除

ここで θ は削除する文長比のしきい値を表す. ここでいう文 長比は BPE を適用する前後の文長の増加率である. これによ り,単語単位とサブワード単位で文長が大きく異なるノイズ文 を削除できる.

3. 実験

3.1 実験設定

学習コーパスには科学技術論文の概要から作成された日英パ ラレルコーパスである Asian Scientific Paper Excerpt Corpus (ASPEC)[Nakazawa 16] を用いた. 文の数は学習用, 開 発用, 評価用それぞれ 977,367, 1,790, 1,812 文である. 日本語 側の単語分割には MeCab*⁵ (辞書: IPADic Ver. 2.7.0) を用い, 英語側のトークン化には Moses*⁶ に付属する tokenizer.perl を 用いた. さらに, OOV の処理として日本語側と英語側に BPE を適用した. これは OOV の問題を解決するための BPE の適 用であるため, フィルタリングの有無に関わらず全ての実験で 適用する. このとき, 英語側と日本語側を合わせたコーパス に対して BPE を適用した. BPE の実装には SentencePiece*⁷ を使用した. また, すべての実験において BPE 化の結合回数

*6 http://statmt.org/moses/

^{*2} この例では1文からなる集合を扱う.

^{*3 @}は単語の終わりを表すトークン

^{*4} ここでサブワードとは単語がサブワード化によって複数に分割さ れたときのそれぞれのトークンのことを指し、もとの単語まで結合 されたようなサブワードは含まない

^{*5} http://mecab.sourceforge.net/

^{*7} https://github.com/google/sentencepiece/



図 1: 削除した文中の文字 1-gram と文字 2-gram のサブワードの上位 20 件

表 1: 改善した翻訳結果の一例(置換により BPE をもとに戻した文)

例 1	入力文	感知用と出力用の2基のコイル, 増幅器, 及び位相シフト回路からなる, 位相シフト磁気センサーシステ		
		ムを開発している。		
	参照訳	here was developed a phase shift magnetic sensor system composed of two sets of coils , amplifiers		
		, and phase shifts for sensing and output .		
	ベースライン	the phase shift magnetic sensor system consists of two $\operatorname{\mathbf{coil}}$, $\operatorname{\mathbf{amplifier}}$, and phase shift circuit ,		
		which consists of two coils for sensing and output .		
	NERF	a phase shift magnetic sensor system consists of two $coils$, amplifier , and phase shift circuits for		
		sensing and output .		
	BPEF	a phase shift magnetic sensor system consists of two $coils$, $amplifiers$ and phase shift circuits for		
		sensing and output .		
例 2	入力文	その結果,電気光学特性として,V10=11.8V,V90=18V を得た		
	参照訳	the electro - optical property obtained was $V10 = 11.8V$ and $V90 = 18V$.		
	ベースライン	as a result , we obtained $V10 = 11.8 \text{ V}$, $V90 = 18 \text{V}$ as an electrooptic property .		
	NERF	as a result , the electrical optical properties were obtained with $V10 = 11.8$ V and $V90 = 18$ V.		
	BPEF	as a result , $V10 = 11.8V$, $V90 = 18V$ were obtained as the electrical optical properties .		

は $\eta = 32,000$ とした.各フィルタリング手法において学習 データ内の語彙サイズと文数は表2のようになった.ここで, 削除文数のパラメータは BPEF に合わせてn = 4,030として いる.

ベースラインは,エンコーダデコーダモデルにアテンション 機構を付加したモデル [Luong 15] をフィルタリング前のデー タで学習させたものである.学習は 20 epoch 繰り返し,その 中で開発データで測った BLEU スコアが最も高かったものを 用いて評価した.両言語の単語ベクトルには学習用データを 用いて word2vec*⁸ を学習させて得た分散表現を用いた.モデ ルの各種パラメータは,埋め込み層の次元数及び隠れ層の次 元数が 512,レイヤー数が 2 とした.また,最適化手法には AdaGrad を用いて初期学習率は 0.01 とし,ドロップアウト率 0.2 でドロップアウトを行う設定とした.

BPEF では、対象を原言語側とし、サブワード化と同様に 語彙サイズを決定するパラメータを $\eta = 32,000$ とした BPE を用いた。削除文数を決定するパラメータ θ の値は、 $\theta = 1.2$ としたときと $\theta = 1.5$ としたときの結果を比較する予備実験^{*9} を行い、 $\theta = 1.5$ がより有効であったため、本実験ではこのパ ラメータを用いる. NERF では、対象を目的言語側とし、固有 表現抽出の実装に Stanford の CoreNLP (Ver. 3.9.2)^{*10} を 使用した. 固有表現の数は原言語側と目的言語側とで等しいと 考えられるため、どちらを対象としても問題ないとみなした.

各手法の評価には自動評価と人手評価の双方を用いた.自動 評価は,BLEUスコアと METEORスコア及び参照訳の文長 との二乗誤差によって行った.BLEUスコアは n-gram 適合 率に基づいて翻訳の精度を評価する評価尺度で,翻訳の流暢性 などを表し,METEORスコアは類義語,語幹正規化,並べ替 え情報などを考慮した評価尺度で,翻訳の妥当性などを表す. また,BLEUスコアと METEORスコアは機械翻訳の仮説検 定ツール*¹¹を用いて算出し,ベースラインに対する統計的有 意差の有無も合わせて調べた.人手評価は2人の評価者によっ て行い,各手法の出力結果からランダムに抽出した100文に 対し,流暢性と妥当性をそれぞれ1から3の3段階で評価し, その算術平均を用いた.数字が大きいほど良い評価である.

3.2 実験結果

表3に実験結果を載せる.BLEUスコア及びMETEORスコアにおける「*」は有意水準0.05で有意差があることを示す. NERFとBPEFの両方でベースラインよりBLEUスコアが0.5程度改善した.また,METEORスコアがNERFで0.1

^{*8} https://radimrehurek.com/gensim/models/word2vec.html *9 本実験と異なり、ドロップアウトを行わず、レイヤー数を1とし

^{*9} 本実験と異なり、ドロップアウトを行わず、レイヤー数を1とした.また、BPE 化の語彙サイズは $\eta = 32,000$ 及び $\eta = 64,000$ とし、英語側と日本語側で合わせず別々に行った.その他の設定は本実験と同様である.

^{*10} https://stanfordnlp.github.io/CoreNLP/

^{*11} https://github.com/jhclark/multeval

表 2: 各手法で使用する学習データの語彙サイズと文数

コッルカリンガ毛汁	語彙サイズ		√/* ₩/r
ノイルタリンク于伝	原言語	目的言語	人致
ベースライン	26,578	19,755	977,367
NERF	26,526	19,754	$973,\!337$
BPEF	26,453	19,753	$973,\!337$

表 3: 各手法の自動評価及び人手評価

手法	BLEU	METEOR	流暢性	妥当性					
ベースライン	23.3	29.7	2.49	2.31					
NERF	*23.8	29.8	2.54	2.35					
BPEF	*23.8	*30.0	2.48	2.30					

程度, BPEF で 0.3 程度改善した.人手評価では,流暢性と 妥当性の双方において BPEF はベースラインとほぼ変わらず, NERF は 0.05 程度改善した.フィルタリングにより翻訳結果 を向上させることができている.さらに,各手法における出力 の参照文に対する文長の二乗誤差は,ベースライン,NERF, BPEF それぞれ, 33.9, 30.9, 31.5 であった.

4. 分析と考察

4.1 削除した文の分析

削除した文中の文字 1-gram と文字 2-gram のサブワードの 上位 20 件を図 1 に示す. NERF と BPEF の特徴を見るため に学習データから抽出したランダムな 4,030 文と合わせて比較 する. 固有表現抽出と BPE によるフィルタリングでフィルタ リングされたサブワードの数に注目すると, ランダムな 4,030 文に比べて大幅に多いことがわかる. これは, 削除した文の中 には BPE を適用したときに 1 文字か 2 文字の細かいサブワー ドが出るような細かい分割が多く行われていることによる. ど ちらにおいても "." が多いのは, the U.S.A. など, 固有表現 に "." を含む固有名詞が多いからであると考えられる.

4.2 出力結果の考察

表1に出力例を載せる.例1のベースラインの出力を見る と、本来複数形として出力するべき単語を単数形として出力し ている翻訳ミスや、同じフレーズを繰り返し出力している翻訳 ミスが見受けられる.それに対して NERF や BPEF ではこれ らの箇所に改善が見られた.また、例2のベースラインの出 力を見ると、本来 the と出力するべき部分を an と誤って出力 してしまっているが、どちらの提案手法でも正しく the と出力 できていることがわかる.

翻訳改善の理由として、複数形や冠詞以外の意味を持つ、固 有表現がサブワード化されて出現した s や an を含む文を一部 取り除くことで、s や an の複数形や冠詞としての意味をモデ ルがより正しく学習できたことが考えられる.今回の実験結果 では BLEUも METEOR も向上しているが、METEOR の改 善が小さいことがわかる.METEOR スコアは妥当性を測るこ とのできる評価尺度であるため、妥当性より流暢性のほうが改 善していると考えられる.人手評価でも NERF において流暢 性のほうがわずかに上がり幅が大きいことが確認できる.妥当 性が低い出力結果を見ると、strut-screw を but screw や snut screw と誤って出力している例があった.これは strut が str と ut に分かれ、str のほうだけうまく翻訳ができていなかっ たものであり、言語による BPE のかかり方の違いも関わって くると考えられる. また,参照文との文長の二乗誤差を見ると,どちらの提案手 法でもベースラインに比べて小さくなっており,出力する文長 が参照文に近づいていることがわかる.このことから,固有表 現を多く含む文,つまり BPE をかけたときにトークン数が増 加する文を取り除いたことで,翻訳モデルの出力文長がより正 しく学習できたことが考えられる.

5. 関連研究

これまでにも機械翻訳におけるパラレルコーパスのフィル タリングに関する研究が行われている.例えば外れ値検出ア ルゴリズムを用いてパラレルコーパスをフィルタリングする研 究 [Kaveh 11] や,文対の意味の違いを識別することに焦点を 当てたフィルタリングの研究 [Carpuat 17] がある.また,パ ラレルコーパスに様々なノイズを与えてその影響を調べる研究 [Khayrallah 18] も行われている.Khayrallahらはニューラル 機械翻訳が統計的機械翻訳に比べてノイズによる悪影響を受け やすいことを示している.

6. おわりに

本論文では、ニューラル機械翻訳におけるコーパスフィルタ リングに関する分析について報告した.固有表現抽出を用いる 方法と BPE を用いる方法の2種類で学習コーパスのフィルタ リングを行った.フィルタリングを行わない場合とともに日英 翻訳の実験を行ったところ、BLEU スコアの 0.5 ポイントの 上昇を確認した.この理由を分析したところ、複数形の有無や 冠詞の間違い、誤った文長で出力される問題などが一部改善さ れたことが分かった.今後は異なる言語対での実験を行い、言 語ごとの有効性の違いを検証していきたい.

参考文献

- [Carpuat 17] Carpuat, M., Vyas, Y., and Niu, X.: Detecting Cross-Lingual Semantic Divergence for Neural Machine Translation, in WNMT (2017)
- [Kalchbrenner 13] Kalchbrenner, N. and Blunsom, P.: Recurrent Continuous Translation Models, in *EMNLP* (2013)
- [Kaveh 11] Kaveh, T., Shahram, K., and Jia, X.: Parallel corpus refinement as an outlier detection algorithm, in *MT Summit* (2011)
- [Khayrallah 18] Khayrallah, H. and Koehn, P.: On the Impact of Various Types of Noise on Neural Machine Translation, in WNMT (2018)
- [Luong 15] Luong, T., Pham, H., and Manning, C. D.: Effective Approaches to Attention-based Neural Machine Translation, in *EMNLP* (2015)
- [Nakazawa 16] Nakazawa, T., Yaguchi, M., Uchimoto, K., Utiyama, M., Sumita, E., Kurohashi, S., and Isahara, H.: ASPEC: Asian Scientific Paper Excerpt Corpus, in *LREC* (2016)
- [Sennrich 16] Sennrich, R., Haddow, B., and Birch, A.: Neural Machine Translation of Rare Words with Subword Units, in ACL (2016)