

1 これまでの修学内容

私は現在、龍谷大学先端理工学部電子情報通信課程に在学しており、データマイニング関連の研究室に所属している。3年次には、Cookpadのレシピデータから調理手順中の利用食材と調理動作の順序性に注目することでレシピの人気分析を行った。卒業研究では、双曲空間を用いたデータマイニング研究に着手する予定である。

2 NAISTで取り組みたい研究

2.1 はじめに

奈良先端科学技術大学院大学で取り組みたい研究テーマは「日本語の同音異義語誤り箇所をラティス構造を用いたニューラル機械翻訳の頑健性向上」である。本稿では、研究テーマの背景・関連研究、提案手法、提案手法の評価について述べる。

2.2 背景・関連研究

近年、NMTモデルは翻訳精度の向上において目覚ましい発展を遂げているが、ほとんどのNMTモデルではノイズを含んだ入力は翻訳文に悪影響を与えてしまう[1]。一般的な入力ノイズの一種として同音異義語ノイズが挙げられる。同音異義語ノイズは私たちが文章作成時に誤った漢字変換をしてしまった際や、音声入力時などに発生する。このようなノイズは入力文が膨大になるにつれ人手での検出は困難になり、翻訳時に理想的な出力を生成するための障壁となる。

このような問題を解消するために昨今、ノイズを含んだ入力文から正確な翻訳文を獲得する研究が行われている。Qinら[2]は、中国語から英語の翻訳にて、同音異義語の疑似誤りデータを用いて同音異義語誤り検出器を学習させ、これを用いて入力文中の同音異義語誤り率の高い単語を音節に置き換えた系列にて学習することで、翻訳性能とノイズに対する頑健性が向上することを示した。

本稿で取り扱う日本語に関する研究では、同音異義語の誤りを検出・訂正する手法において三品ら[3]は、従来のn-gramのみの手法に比べ、n-gramと確率的潜在意味解析法を組み合わせることで精度が改善することを示した。また、藤井ら[4]はBERTとZero-shot学習を組み合わせた検出を試みた。この他にも様々な同音異義語誤りを検出・訂正する手法が提案されているが、日本語の同音異義語誤りに頑健なNMTモデルの提案は未だ行われていない。

他方、NMTモデルの性能向上を図る取り組みの一例として、入力系列に複数の分割候補を同時に表現できるラティス構造を取り入れる研究が行われている。しかし、近年NMTモデルとして非常に高い精度を達成したTransformerでは、位置符号化を用いて各トークンの系列内位置情報を捉えており[5]、ラティス構造をそのまま入力することは困難である[6]。そこで、Xiaoら[6]はラティス構造を構成するノードの1文字目を絶対位置とし、ラティス構造を単一の系列で表現することでTransformerへの入力を可能とした。同時にラティ

ス同士の距離を適切に扱えるself-attentionを用いたラティス空間のエンコーダを提案し、中国語から英語と英語からドイツ語の機械翻訳タスクにて従来のTransformerエンコーダーに比べ優位性を示した。Liら[7]は中国語NERタスクにて、ラティス構造を構成する各トークンに始点と終点の位置を割り当て、スパンからなる平坦な構造に変換することでTransformerへの入力を可能とし、従来法に比べ性能・効率ともに優れていることを示した。

このようにラティス構造を活用した研究は幅広く行われており、本稿で取り扱う問題の同音異義語誤りについて、誤り箇所をラティス構造を用いることで、他の同音異義語候補を複数持つ系列を表現することが可能となる。Liら[7]のラティス構造をNMTへ入力する手法に加えて、Douら[8]が提案した入力文書とそれを補助する情報をそれぞれエンコードする2つのエンコーダを持つガイド付き要約フレームワーク(GSum)を踏襲することで、NMTモデルが翻訳文に基づいて複数の候補から適切な語を注視するように学習することが期待できる。

以上より「ラティス構造を含んだ入力系列に対応したTransformerを用いた同音異義語に頑健な日本語から他言語へのNMTシステム」の提案を行う。

2.3 提案手法

提案手法は(1)同音異義語誤りの検出、(2)同音異義語誤り箇所へのラティス構造の付与、(3)Transformerへのラティス構造の入力、という3つのステップに分けられる。以下でそれぞれの詳細を述べる。

2.3.1 同音異義語誤りの検出

日本語ソース文中の同音異義語誤りを検出する手法について述べる。検出器にはMeCabを用いて形態素解析を行った後に単語をサブワード分割しているBERT base Japanese¹を用いる。また検出器の訓練データとして毎日新聞の記事データを用いる。まず、訓練データに対してMeCabを用いて漢字部の読みごとの同音異義語辞書を作成する。次に、先ほどMeCabによって検出された漢字部を作成した同音異義語辞書を用いて、誤った同音異義語で置換することで、疑似的に誤りを含んだ文を作成する。

また、我々が誤った同音異義語を選択してしまう場面は、同じ同音異義語でも文中での出現頻度が低いものより高いものであるという予想から、同音異義語辞書内に出現頻度でソートし、上位5件からランダムに選択する。加えて、原文の意味を保持するため、修正箇所は3箇所までとし、この数を超えないようにランダムに選択する²。作成された正文と誤り文を用いて各トークンごとに0(誤りなし)と1(誤りあり)の二値分類教師あり学習をBERTに対して行う。

図1は同音異義語誤りの検出手順の概略を示しており、作成した検出器からは①のような出力が得られる。これをMeCab

¹<https://huggingface.co/cl-tohoku/bert-base-japanese-whole-word-masking>

²これらのパラメータについては検討する必要がある

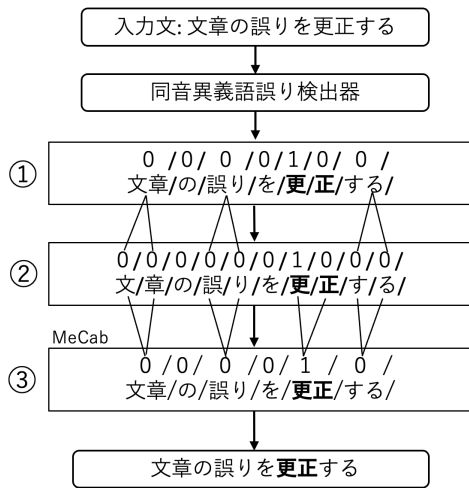


図 1: 同音異義語誤り検出手法の概略図

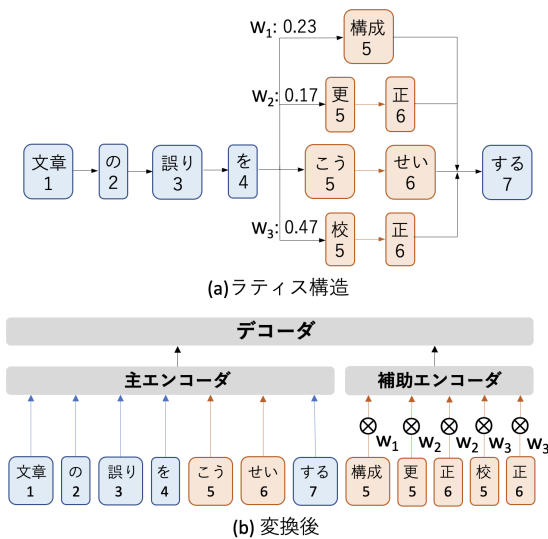


図 2: ラティス構造 (a) と変換後の入力系列 (b)

の分割に当てはめることで同音異義語誤り箇所を特定したいが、出力はサブワードごとになるため MeCab の分割とサブワードの分割が一致するとは限らない。そこで②のように一度文字レベルに 0,1 の割り振りを行う。最後に③で MeCab の分割と検出器の出力を照らし合わせ、各分割ごとに論理和の処理を行うことで同音異義語誤り箇所の検出が可能となる。

2.3.2 同音異義語誤り箇所へのラティス構造の付与

検出された同音異義語誤り箇所に対して、ラティス構造を適応した系列の作成手順について述べる。まず、藤井ら [4] の手法に基づき、誤り箇所を BERT の [MASK] トークンに置き換えて、その部分を BERT に予測させる処理を行う。予測結果を尤度順に n 語取得し、それらが [MASK] 部分の読みと対応する同音異義語辞書中に存在すれば、上位から最大 3 件抽出する。この制限は attention の計算量が入力系列長の二乗に比例するためである。またその際、抽出した候補に BERT から得られる同音異義語候補すべてで尤度を正規化したものを各同音異義語候補の重みとする。図 2(a) では上記の処理によって得られた同音異義語候補を含んだラティス構

造の例を示している。この際、Qin ら [2] の入力系列中に音節情報を組み込む手法を参考に、提案手法では日本語の読みをラティス構造に組み込む。これによって抽出した同音異義語候補内に正例がなかった場合であっても、NMT モデルが日本語の読みに注視して学習することを期待する。

2.3.3 Transformer へのラティス構造の入力

ラティス構造を NMT モデルへの入力を可能とするために、ラティス構造から単一の系列に変換することで NMT モデルへの入力を可能とする。しかし、このままでは系列中に同一の位置情報が複数存在することになり、基本的な単一のエンコーダからなる NMT モデルへ入力してしまうと、モデルが正確にトークンの位置情報を学習できない可能性がある。そこで、図 2(b) のように Dou ら [8] が提案した GSum を踏襲し、系列を主情報と補助情報の 2 つに分割する。主情報は同音異義語誤り検出部をその読みに置き換えた文とし、補助情報は同音異義語候補のみとする。主情報を主エンコーダに、補助情報を補助エンコーダに入力する。これは同音異義語の読みに対して補助エンコーダの情報をもとに学習することを期待してである。補助エンコーダの位置情報はソース側の同音異義語の位置情報を共有する。加えて、ガイドンス側の入力の同音異義語候補には BERT より得られた正規化された尤度を重みとして掛けることで、モデルが有力な同音異義語に着目することを期待する。

2.4 提案手法の評価

JESC, JPO などの対訳コーパスの日本語入力文に対して、提案した検出器を適応し、NMT モデルを訓練する。

学習した NMT モデルの翻訳性能を BLEU などによって測定し、Transformer のベースモデルとの翻訳性能を比較することで、提案モデルの翻訳性能を評価する。

提案した NMT モデルが、同音異義語誤りに対して効果的か評価する必要があるが、実世界での日本語の同音異義語誤りのみを取り扱った対訳データは存在しない。そこで、Qin ら [2] の手法を参考に、テストセット内に一定の確率で人工的に同音異義語ノイズを混入させたデータを用いて、ベースモデルと提案モデルの翻訳性能を比較することで頑健性の評価を行う。

参考文献

- [1] Y. Belinkov et. al., "Synthetic and natural noise both break neural machine translation", ICLR, 2018.
- [2] W. Qin et. al., "Modeling Homophone Noise for Robust Neural Machine Translation", CASSP, 2021.
- [3] 三品 et. al., "確率的 LSA を用いた日本語の同音異義語誤りの検出・訂正", 情報処理学会論文誌, 2004.
- [4] 藤井 et. al., "BERT を利用した Zero-shot 学習による同音異義語の誤り検出", 言語処理学会 第 27 年次大会, 2021.
- [5] A. Vaswani et. al., "Attention is all you need", NeurIPS, 2017.
- [6] F. Xiao et. al., "Lattice-Based Transformer Encoder for Neural Machine Translation", ACL, 2019.
- [7] X. Li et. al., "FLAT: Chinese NER Using Flat-Lattice Transformer", ACL, 2020.
- [8] Zi-Yi Dou et. al., "GSum: A General Framework for Guided Neural Abstractive Summarization", NAACL, 2021.