

手がかり表現に基づく非論理的な言語推論の学習

張 辰聖子 (指導教員：小林 一郎)

1 はじめに

人間が行っている推論は論理的な含意関係だけではなく、日常的な知識を背景にした常識推論が大きな役割を果たしている。本研究はこの際の知識を、複雑な内容を表現できる自然言語文自体にとらえ、「から」「ため」のような手がかり表現を利用することで、前提と帰結がともに文となる自然言語推論をコーパスから深層学習モデルとして直接学習する。コーパスから抽出した157万文ペアについて学習し、テストデータにおいて、前提から生成した帰結文について人手評価を行ったところ、65.8%が妥当な推論であるとの結果を得た。さらにエラーの原因について考察し、改善と今後の可能性について議論する。

2 自然言語による推論

2.1 研究概要

図1日本語 Wikipedia コーパスから、「理由」を表現する手がかり表現を元に根拠と結論がペアになった文を抽出する。根拠部分を入力、結論部分を出力とし、文を生成する深層学習モデルである T5[2] に学習させる。これにより根拠に相当する文から結論を示す文を生成することを通じて、自然言語推論を実現する。

2.2 データ作成

データの収集には日本語版 Wikipedia コーパスを用いた。以下にデータ作成の手続きを示す。

Step 1. 正規表現で理由節を持つものを抽出 日本語版 Wikipedia 本文から、理由節を持つ1文を正規表現を用いて抽出する。その際、手がかり表現として「から」「ので」「ため」「おかげで」「せいで」の5つを用いた。また、理由節を持つ文の文頭が指示詞（「この」、「その」）であった場合、前文で置換をするという簡単な指示詞の補完を行った。

Step 2. 形態素解析によるフィルタリング 正規表現によって抽出された理由節を持つ文に対し、1文ごとに MeCab を用いて形態素解析を行う。正規表現での抽出では、5つの手がかり表現は、理由を表す以外の意味を持つことがある。そこで、形態素解析により主に理由を表す品詞として使われている場合のみを抽出する。以下に品詞による使われ方の違いの例を示す。

- 信号が青から赤に変わった。(格助詞, 変化の「から」)
- 楽器を持っていなかったからボーカルになった。(接続助詞, 理由の「から」)

例に挙げたように「から」の品詞細分類が接続助詞である時、主に理由を表すことから解析により抽出を行う。上記により、解析の結果、理由節が以下の形態素として含まれていた場合に、抽出を行う。

Step 3. 機械学習によるフィルタリング 形態素解析によって残った文のうち、名詞で用いられる「ため」はいくつかの意味を持つ。そのため、理由を表す「ため」のみを抽出することを目的として、訓練済み日本

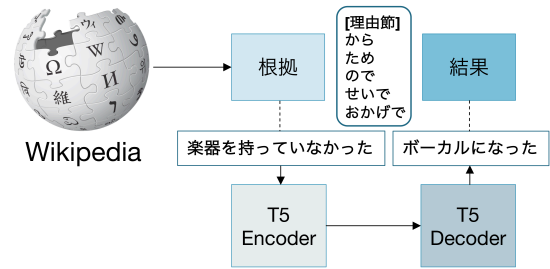


図 1: データ生成と、推論の学習の概要。

語 BERT モデルを用いて二値分類を行った。ここで用いた BERT モデルは Hugging Face¹ の自然言語ライブラリ Transformers² に基づく、東北大学公開の事前学習済み日本語 BERT モデル³ を使用した。「ため」を手がかりとして抽出した文から、ランダムに 300 文を抜き出し、手動で理由を表すか否かのタグ付けをし、二値分類の機械学習を行った。300 文での 5 分割交差検証では 88.7% の精度で理由を表すか否かの分類が可能になった。分類を行った結果、1,288,304 文の「ため」を含む文から、730,133 文 (56.7%) の理由を表す文を抽出することができた。

Step 4. 根拠表現と結果表現の抽出 日本語自然言語処理ライブラリである GiNZA⁴ を用いて、抽出した文に対して係り受け解析を行う。坂地ら [4] の手法を用いて、根拠表現と結果表現の抽出を行った。この手法に、文になっていないものや、文の主語や述語が足りないものを除くために名詞句以外であれば、文節が 2 つ以上のものという制約を加えた。

2.3 推論モデル

収集した文ペアを日本語 T5 モデルに与え、前提から結果を生成する深層学習モデルを学習することで、自然言語による推論を実現する。

3 実験

3.1 実験設定

本研究では、2.2 節で収集した 1,572,956 件の根拠と結果のペア文をデータとして用いる。実験の際は、このデータを訓練:開発:評価=0.95:0.025:0.025 として評価を行う。評価は評価データからランダムで 100 文を抽出し手動で行い、それに伴い BERT-Score[3] によって、生成文と正解文との意味的類似度を検証した。評価指標は 3.2 節に示す。本研究で用いる T5 モデルは Hugging Face の自然言語ライブラリ Transformers に基づく、Isao Sonobe 公開の事前学習済み日本語 T5 モデル⁵ を使用した。

¹ <https://huggingface.co>

² <https://github.com/huggingface/transformers>

³ <https://huggingface.co/cl-tohoku/bert-base-japanese>

⁴ <https://megagonlabs.github.io/ginza/>

⁵ <https://huggingface.co/sonoisa/t5-base-japanese>

表 1: 訓練データの例と、それに類似する前提を与えたときの学習済み T5 の文生成出力例。

| | | |
|---------|----|-------------------|
| 訓練データ 1 | 前提 | 小波まで漢字だと固いイメージになる |
| | 結果 | こなみとひらがなにした |
| 訓練データ 2 | 前提 | 漢字ではイメージが固い |
| | 結果 | みずいろと平仮名にした |
| 新しい入力 | 前提 | 漢字ではとっつきにくい |
| | 出力 | ひらがなにした |

表 2: 人手評価されたデータに対する生成文と正解文の BERT-Score の平均。

| 種類 | 件数 | Precision | Recall | F_1 |
|------|-----|-----------|--------|-------|
| full | 100 | 0.678 | 0.700 | 0.688 |
| i | 52 | 0.699 | 0.719 | 0.707 |
| ii | 20 | 0.638 | 0.667 | 0.651 |
| iii | 11 | 0.688 | 0.695 | 0.690 |
| iv | 9 | 0.665 | 0.668 | 0.666 |
| v | 5 | 0.700 | 0.721 | 0.710 |
| vi | 3 | 0.597 | 0.654 | 0.623 |

3.2 評価基準

付録 A の表 4 の生成例 2 にあるように、評価データの正解と比べると全く違う出力であるが、出力内容自体は前提から妥当に導かれるものが多くあった。したがって、生成データセットにおけるモデルの性能は、出力の妥当性を人手で判断することで評価する。評価基準は、自然言語による演繹的な文生成を行うモデルの提案をしている関連研究 [1] を参考にし、以下の 6 つの分類とする。

- i 妥当な推論である 出力は前提から妥当に導かれる結論である。
- ii 妥当な推論であるが文法的な間違いがある 出力は前提から導かれる結論だが、主語などの成分が欠けていたり、動詞の活用に違和感があるなど内容の理解を妨げない程度の文法的な間違いがある。
- iii 前提を繰り返す 出力は、前提文の繰り返し、もしくは言い換えであり前提と同意である。または、前提に含まれる単語を組み合わせているだけである。
- iv 前提から導かれない 出力は文、内容は正しいが、前提から導かれ得る結論ではない。
- v 矛盾している 出力は、前提と矛盾する、もしくは本質的に間違っている。
- vi 結論が理解し得ない 結論が文をなしていない。もしくは前提が文をなしていない。

3.3 実験結果

実験でファインチューニングされたモデルに評価データを与えると、付録 A の表 4 のように文生成が行われた。前提はモデルに与えた文、出力はモデルから生成された文、正解は収集データ中の前提に対する結論文のことである。生成データセットにおけるモデルの性能に対する人手で行った評価結果は、各評価者ごとの評価数を表 3 に、全ての評価者の平均をグラフにしたものを図 2 に示す。人手で評価したデータに対する、生成文と正解文の BERT-Score [3] による類似度を表 2 に示す。BERT-Score とは、BERT の埋め込み (分散表現) を利用することで、文章の類似性を評価するものである。また、訓練データにあるような形の前提であ

表 3: 評価データからランダムに抽出した 100 文に対する人手評価の結果 (件数)。

| 種類 | 評価者 1 | 評価者 2 | 評価者 3 | 評価者 4 | 割合 (%) |
|-----|-------|-------|-------|-------|--------|
| i | 52 | 47 | 45 | 67 | 52.8 |
| ii | 20 | 19 | 13 | 0 | 13.0 |
| iii | 11 | 10 | 9 | 9 | 9.8 |
| iv | 9 | 18 | 21 | 13 | 15.3 |
| v | 5 | 5 | 9 | 7 | 6.5 |
| vi | 3 | 1 | 2 | 4 | 2.8 |

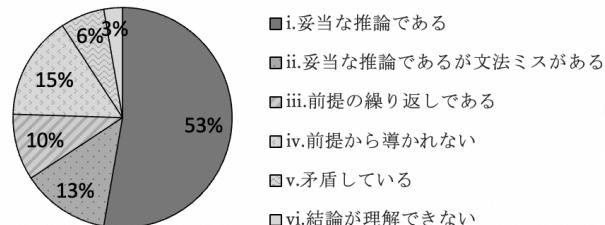


図 2: 評価データからランダムに抽出した 100 文に対する人手評価の結果 (平均)

れば妥当な推論を行えるという仮定の下、訓練データの例と、それに伴うモデルによる出力例を表 1 に示す。出力が妥当な推論である、もしくは妥当な推論であるが文法的な間違いがあると評価された割合は 65.8% となった。

4 まとめ

本研究では、自然言語文そのままの形で推論を行う手法の開発を目的に、自然言語による推論を自然言語文生成として表現する研究に取り組んだ。日本語 Wikipedia コーパスから、理由節を手がかり表現として根拠と結論を抽出した。抽出した因果関係のデータを用いて、根拠部分を入力、結論部分を出力として深層学習を行うことで、推論を行う形の文生成が可能になったことがわかった。モデルによって生成された文に対して人手で評価を行った結果、前提から導かれる推論として妥当だと言える割合は 65.8% となった。今後は、発生したエラーを改善して収集データのノイズを減らし、推論文生成の精度を上げたい。また、作った推論モデルを用いて、多段推論および前提を予測する帰納的推論などに挑戦していきたいと考えている。

参考文献

- [1] Bostrom, K., Zhao, X., Chaudhuri, S. and Durrett, G.: Flexible Generation of Natural Language Deductions, in *EMNLP 2021*, pp. 6266–6278 (2021).
- [2] Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W. and Liu, P. J.: Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer, *Journal of Machine Learning Research*, Vol. 21, No. 140, pp. 1–67 (2020).
- [3] Zhang, T., Kishore, V., Wu, F., Weinberger, K. Q. and Artzi, Y.: BERTScore: Evaluating Text Generation with BERT, in *International Conference on Learning Representations* (2020).
- [4] 坂地泰紀 F 構文パターンを用いた因果関係の抽出, 言語処理学会第 14 回年次大会論文集, 2008, pp. 1144–1147 (2008).

付録 A 推論結果の例

表 4: 学習した T5 深層学習モデルによる推論結果の例.

| 分類種類 | 文章 |
|-----------------------------|---|
| 例 1 妥当な推論である | 前提 一部踏み跡程度の不明瞭な登山道や作業道などがある |
| | 出力 通行には注意を要する |
| | 正解 登山の際には注意が必要である |
| 例 2 妥当な推論である | 前提 最寄の鉄道駅から距離が遠い |
| | 出力 自転車通学が多い |
| | 正解 本校は鉄道を利用して通学する生徒は少ないと思われる |
| 例 3 妥当な推論である | 前提 戦争支持の風潮に背を向けた |
| | 出力 批判された |
| | 正解 非国民のレッテルを張られる |
| 例 4 妥当な推論である | 前提 レーダー類が貧弱である |
| | 出力 対艦攻撃能力は低い |
| | 正解 射撃管制手段は不明だ |
| 例 5 妥当な推論である | 前提 中々仕事がない |
| | 出力 暇をもてあます |
| | 正解 バイトをしている |
| 例 6 妥当な推論である | 前提 帰任先が大阪だった |
| | 出力 関西弁が話せる |
| | 正解 現在大阪在住 |
| 例 7 妥当な推論であるが 文法的な間違いがある | 前提 騎兵隊という部隊名が付いているが馬の輸送が間に合わずまた連発銃の普及により 騎馬のメリットが薄れていた |
| | 出力 格下げされた |
| | 正解 実際は兵士は徒歩だったルーズベルトは黒馬に乗っていた |
| 例 8 妥当な推論であるが 文法的な間違いがある | 前提 光化学系や光化学系の還元力を蓄積させ活性酸素を発生させる |
| | 出力 触媒とし |
| | 正解 植物にとって害となる |
| 例 9 前提を繰り返す | 前提 e は超越数であると証明されている |
| | 出力 e は超越数であると証明されているのだ |
| | 正解 背理法によりは無理数でなければならない |
| 例 10 前提を繰り返す | 前提 マタンサ自体が普通名詞である |
| | 出力 スペイン語圏ではマタンサは普通名詞である |
| | 正解 誤解を防ぐため文章化に当たっては atanzadelcerdo ブタのマタンサなどと表記する |
| 例 11 前提から導かれない | 前提 幼い頃から弟の面倒を見てきた |
| | 出力 家計は裕福だ |
| | 正解 子供の面倒を見る |
| 例 12 前提から導かれない | 前提 咽喉頭の異常感は古代ギリシア時代から記載があるといわれておりヒポクラテスは更年期の 女性が起こす病態と捉え当時子宮は体内を移動する臓器と考えられていた |
| | 出力 ヒポクラテスは子宮の異常感を女性特有の生理現象と捉えた |
| | 正解 これが喉につかえる感じを生じさせていると考えた |
| 例 13 矛盾している | 前提 生息数が減少している可能性が高い |
| | 出力 絶滅のおそれは低いと考えられている |
| | 正解 C は絶滅危惧種としている |
| 例 14 矛盾している | 前提 組織内に寄生したプレロセルコイドに有効な駆虫薬はない |
| | 出力 駆虫薬を併用する |
| | 正解 外科的に摘出することが有効である |
| 例 15 結論が理解し得ない | 前提 1 編成を 1 車両として運用する考え方を基本とした |
| | 出力 これ |
| | 正解 連結器は中間連結部は半永久連結器を使用している |
| 例 16 結論が理解し得ない | 前提 1 の値と独立ではない |
| | 出力 1 |
| | 正解 末尾事象ではない |