単語単位ラベルを考慮したテキストデータ拡張手法の研究

2237007

教育研究センター



研究代表者 (助成金受領者) 共同研究者

同志社大学 大学院 博士後期課程 寺 本 優 香 文化情報学研究科 同志社大学 文化情報学部 教 授 波多野 賢 治 名古屋大学 数理・データ科学 孝 准教授 駒水 裕

[研究の目的]

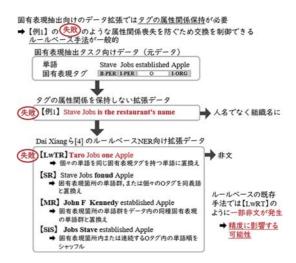
固有表現抽出とは、文中から人名や地名、組織名、日付表現、薬品名、遺伝子名といった特定の属性を持つ単語を抜き出す自然言語処理の基礎技術であり、医療・創薬分野での情報抽出やチャットボットなど様々な分野で利用されている。

このタスクでは教師ありの機械学習モデルが 頻繁に利用されており、近年は深層学習モデル が優れた性能を示している。しかし、その学習 には大規模かつ高品質なデータが必要であり、 特に分野固有の固有表現を扱う場合は、専門知 識を持つアノテータによるデータ作成が不可欠 であり、人的コストが課題となっている。

固有表現抽出タスクで中心的役割を果たしてきたデータセットに CoNLL2003 があるが、その固有表現カテゴリ(人名・地名・組織名)は比較的粗い粒度であり、分野やアプリケーションによってはより細かい粒度が求められる。近年では、Ding らによって提案された Few-NERD のように、政治家や作家など細粒なカテゴリに分けたデータセットも登場している。今後も使用目的に応じて多様な粒度・種類の固有表現カテゴリが求められることは自明である。

このようなデータ準備のコストを削減する方 法の一つにデータ拡張がある。データ拡張は、 小規模な学習データに対して加工や合成などの 処理を施し、新たなデータを生成して学習に用いる技術である。固有表現抽出タスクでは、単純なルールに基づく交換・挿入・削除といった手法が主流であり、これは系列ラベリングにおいて単語の並びの構造的整合性を保持する必要があるためである。本手法の目的は、既存手法のルールベースを発展させ、固有表現抽出のデータ拡張において新たな置換可能条件を明らかにすることである。

[研究の内容. 成果]



既存手法の交換対象は同種のトークン・単語境界範囲内か、同一文中に出現する他のセグメントに限定される。本研究では NLP で伝統的に用いられる構文木と格フレームを手掛かりと

し. より広範囲の境界線を探索する。前者の構 文木は部分木ごとに文構造上のまとまりを取得 でき、後者は固有表現と関連する規則性とみな されている。

着眼点・目標

- グの交換はより制約を追加して行われるべきという仮定
- ・文法上の機能が類似した単語同士での交換拡張により、拡張データの 文法的規則上の破綻を抑制

① POSタグを用いた手法

ひとまとまりの固有表現箇所・またはOタグを持つ単語に対し、以下の交換規則を設定



② SubTreeを用いた手法

- 文法構造を表現する構文木より部分木 (SubTree)を取得
- 同じ機能を持つ部分木同士を交換
- ➡ 同じ深さと属性を持つ部分木の交換による拡張は構造・意味的に合理的[6-8]



構文木の検証では子ノードを二つ以上もつよ うな部分木を交換対象とする。図の例では計3 パタンの部分木が得られる。元データ全体の整 合性を担保するため、 そこから得られた用法の 近い部分木、あるいは名詞同士を置換候補とし た。これにより、固有表現箇所以外の交換につ いても実現が可能となる。上記の両手法につい て、既存の置換手法と置換候補の幅、データ品 質を比較し検証する。

我々の実験において比較に採用する既存手法 および提案手法の略字表記と簡単な説明を以下 に記す。

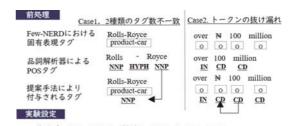
MR:既存手法。固有表現を同種固有表現と交

MR+POS: MR において, 交換条件に POS タグが同じことを追加

OLwTR: ○ タグを持つ個々の単語同士を交換 OLwTR+POS: において交換条件に POS タ グが同じことを追加

MR+OLwTR: MR と OLwRT を同時に適用 MR+OLwTR+POS: MR+OLwTR におい て交換条件に POS タグが同じことを追加 **TreePos**:交換条件に TreePos の深さ属性を

なお、解析器の違いなどによるセグメントの 差に関する取扱い、使用データ、使用ツール、 データ拡張時の実験設定については,以下の図 に記す。



- ・使用データ: Few-NERD 解析ツール: Stanford CoreNLP
- NLP 夕を作成 ^{*電}络データ=1:10 ・元のデータセットと同じタグ比率の3/167スケールデータを作 ・単語・固有表現の交換発生確率0.7 元データ:拡張後ラ

実験の結果、交換に POS タグを考慮した手 法群では、もともとのデータセットにおいて, 複数のトークンから成立する固有表現が多く含 まれるカテゴリの Recall が向上した。使用し たデータセットでは、主に紛争名、災害名など のカテゴリがこれにあたる。一方で、〇タグ にあたるトークンを変更した手法では. 固有表 現・一般名詞両方の用途がある単語に対する正 答率が低下した。

具体的な例としては、「民主党」あるいは「民 主的な(形容詞)]」)の両方の意味を持つ Democratic という固有名詞がこれに該当する。一般 名詞として使用される場合には、特定の文脈の みに出現していたが、ルールを設定したランダ ムな交換によりその規則性が無くなり、文脈情 報が失われたと考えられる。こうした傾向は. とりわけ MR+OlwRT+POS を用いた場合に. Democratic のような複数の意味を持つ単語を 含んだ固有表現箇所で頻繁に発生した。以下の 表は、本研究の結果、得られた成果である。

提案手法は、ごく少数のデータを使用するよ うなカテゴリの分類では、既存手法よりわずか に精度が向上した。一方、オリジナルのデータ が数十件程度確保できる場合には、既存手法に よるデータ拡張のほうがより高精度であった。

proposed macro-AC macro-F1 macro-P macro-R method Original-small 0.915 0.626 0.604 0.650 0.902 0.562 0.601 MR 0.528 MR+POS 0.9020.562 0.527 0.602 OLwTR 0.902 0.570 0.545 0.597 OLwTR+POS 0.903 0.578 0.552 0.606 MR+OLwTR 0.560 0.594 0.901 0.529 MR+OLwTR+POS 0.902 0.569 0.603 0.540 SubTree 0.900 0.536 0.593

また、100件程度の正解データを確保できる場合、提案手法・既存手法ともにデータ拡張はかえって精度の低下につながるということが明らかになった。このことから、生成系モデルなどである程度のデータ量が確保できる一般的なドメインの固有表現抽出タスクでは、現時点では

データ拡張そのものがあまり現実的ではない可能性がある。一方で提案手法は、歴史的な書籍など少数のデータの利活用が必要なドメインにおいて有効である可能性が示唆される。

「成果の発表、論文など]

- 李優香, 駒水孝裕, 波多野賢治: "固有表現タグおよび POS タグによる交換制約付きデータ拡張手法", 第 15 回データ工学と情報マネジメントに関するフォーラム予稿集, No. 1b-6-4, March 2023. (学生プレゼンテーション賞 受賞)

寺本優香, 駒水孝裕, 波多野賢治: "POS タグおよび 構文木の SubTree を用いた NER 向けデータ拡張手 法", 東海関西データベースワークショップ 2023