

BERT を用いたフィルタリングによる Twitter からの 教師データ作成手法

金澤 滉典[†] 櫻井 義尚[‡]

明治大学先端数理科学研究科[†] 明治大学総合数理学部[‡]

1. はじめに

企業が SNS をマーケティング分析の場として利用する「ソーシャルリスニング」は、消費者の声が直接得られる一方で、意見情報のみを効果的に抽出することは非常に困難である。これに対し、教師あり機械学習によって作成された意見抽出モデルを用いた判定方法では、学習する際の教師データが不均衡であるために抽出精度が落ちることが確認されている。不均衡緩和を目的としてアンダーサンプリングを適用する場合、データ数が減少することによる過学習が問題となっていた。

野崎らは、ツイート集合に対して評価表現辞書に基づいた段階的なフィルタリングをかけることで、教師データ数を確保しつつ意見として判定できるようなデータを抽出することが可能な手法 PSSA[1]が提案された。一方で、辞書を用いたフィルタであることによって、文脈的に意見だと判断できる文章が収集できない可能性が指摘されている。

本研究では、辞書を用いたフィルタリングを、前後の文脈を扱う自然言語処理モデルを用いたフィルタリングに置き換えることで、辞書に基づいた意見ツイートと文脈として判断できる意見ツイートを抽出できる教師データ作成手法を提案する。

2. 関連研究

2.1. 不均衡緩和に関する研究

機械学習における不均衡な教師データを効率的に用いるための議論は多い。少数データを何らかの形でかさ増しするオーバーサンプリングに対して焦点が当たる研究が多く、画像分類などで効果を発揮している。

紺野ら[2]が提案した深層学習中の特徴量を用いた擬似特徴量生成手法も、画像分類に対する効果を示している。画像分類ほど効果的な手法は少ないが、自然言語処理においても単語・文節の入れ替えによる少数データのオーバーサンプリングを澤崎ら[3]が提案している。

2.2. 意見抽出に関する研究

SNS のテキスト集合は意見抽出を目的としているわけではないため、意見抽出器の教師データとして活用できる文章は非常に少ない。評価表現辞書を用いる意見抽出手法を提案した立石ら[4]は、その辞書を元に評価表現の判定と抽出を行った。

本研究では、辞書に含まれていない表現ではあるものの文脈上では評価表現と同義である文章も抽出可能なフィルタを作成し、意見抽出を目的とした教師データの効果的なサンプリング手法の確立を目指す。

3. 提案手法

本論文では、意見抽出器の教師データ作成を目的とした PSSA の評価表現辞書によるフィルタリングを BERT による実装に置換することで、辞書に含まれないが文脈で評価表現を示す文章への反応性向上を図る。PSSA のパラメータは先行研究と同様に 3 ブロック、ラベル比は 5:5 を目標とした。

3.1. PSSA; Prefilter based Stepwise Sampling for Annotation

ランダムサンプリング時に特定の分類クラスのデータが少なくなる場合、少数データが増加するような機械的なフィルタを適用することで、データの不均衡を緩和しつつデータ数も確保するアルゴリズムである。このとき、機械的フィルタを適用する都合上、抽出後のデータが単純化する可能性があるため、絞り込み効果の小さな弱フィルタから順に適用することでサンプリング後のデータセットの複雑性を担保している。

3.2. 評価表現辞書を用いた PSSA でのサンプリング

辞書には小林が策定した品詞なし評価表現辞書を用いる。本辞書は評価表現の可能性のある約 5200 表現が存在し、ドメイン横断を許容している。以下の 3 段階に分けてサンプリングすることで、不均衡を緩和したデータセットを作成する

- 第 1 段階

辞書を利用してフィルタリング。抽出されたものとされなかったものが 7:3 となるようにサンプリング。

- 第 2 段階

辞書に対し MeCab を用いて形態素解析を行い、形容詞・副詞助

Creating Training Data from Twitter using BERT-based Filtering Method

[†] Hironori Kanazawa [‡] Yoshitaka Sakurai
Meiji University

詞類・助動詞・名詞副詞接続・名詞形容動詞語幹に該当する表現のみを用いてフィルタリング。抽出されたものとされなかったものが8:2となるようにサンプリング。

- 第3段階

辞書に対し MeCab を用いて形態素解析を行い、形容詞・副詞助詞類・助動詞・名詞副詞接続に該当する表現のみを用いてフィルタリングし、その中から 71 文字以上のツイート抽出。抽出されたものとされなかったものが5:1となるようサンプリング。

3.3. BERT を用いた PSSA でのサンプリング

評価表現辞書を用いた機械的フィルタリング部分を、BERT による文章分類に置換する。サンプリング比などは対照実験のため変更しない。

4. 実験

4.1. ツイートデータの収集と前処理

2019年5月1日から2020年4月30日までの、テーマパークや都内5つ星ホテルなどの19施設がキーワードとして含まれるツイートを、各施設15000件を上限としてRTとリプライを除き収集した。その後、収集した全てツイートに対し以下の前処理に掛けた。

- URL は該当部分を<URL>に置換
- 全ての日本語は全角に置換
- 全英数字は半角小文字に置換

以上が完了後、MeCab を用いてツイート本文の分かち書きを行った。このとき、MeCab の辞書は標準の IPADIC に加えて、インターネット上の新語に対応している IPADIC-NEologd を利用した。

4.2. BERT によるフィルタ機能の実装

4.2.1. フィルタ代替モデルの構築

本実験では、前処理が完了した約190,000件のツイートに対して評価表現辞書を用いた PSSA を適用し、各段階のフィルタに抽出されたか否かを positive/negative のラベルとして設定し教師データとした。その後、BERT モデルに全層結合層加えた自然言語モデルを用いて、文章分類のタスクとして実装した。使用ライブラリは Pytorch の Transformers、日本語の事前学習モデルは東北大学乾研究室が作成した bert-base-japanese-whole-word-masking を使用した。テストデータは、検索条件はそのままに、収集期間を2018年5月1日から2019年4月30日までとするツイート約1万件とした。

4.2.2. フィルタ機能の学習

BERT による辞書フィルタの予測結果を示す (Table 1)。分類

モデルのため、評価指標は Accuracy, Precision, Recall, F1-score とした。

Table 1 BERT による辞書フィルタの予測結果

PSSA 段階	Accuracy	Precision	Recall	F1-score
第1段階	0.76	0.88	0.76	0.78
第2段階	0.85	0.90	0.85	0.86
第3段階	0.99	0.99	0.99	0.99

段階が増えていくにつれ、モデルでの予測と辞書を用いた抽出との差異がなくなっていった。フィルタが厳しくなれば対象となる文章が少なくなるため、BERT の判断が適切に行われたと思われる。また、浅い段階での Recall の値が Precision と比べて少し低いことがわかるが、これは辞書の特性だと思われる。評価表現辞書には漢字1文字で1表現として存在する場合があります、当辞書をそのまま用いて抽出する第1段階などは文章に過剰にマッチングしやすい。BERT モデルでは、文字ではなく文脈を優先するため、このような結果になったと考えられる。

5. まとめ

本研究では、教師データの不均衡を緩和する手法である PSSA の辞書に依存するという問題点を解決するため、辞書によるフィルタを BERT に置き換える手法を提案し、比較検証を行った。実験結果より、第1段階、第2段階という浅い段階における辞書フィルタは、Recall の低さが目立った。これは辞書内に存在する評価表現が、漢字1文字で1表現を示す場合があります、辞書フィルタでの過剰なマッチングが起きる可能性の高さが示唆されている。

今後の課題として、実際に BERT フィルタ用いて意見抽出タスクに向けた教師データを作成し、意見抽出モデルの精度を検証する必要がある。また、従来の教師データ作成手法との比較検証をすることで、本研究の有用性を検証する。

[参考文献]

- [1] 野崎雄太, 櫻井義尚. 「Twitter からの意見抽出モデル構築のための教師データ作成手法」研究報告数理モデル化と問題解決(MPS)2020-MPS-127.9(2020):1-6.
- [2] 紺野友彦, 藤井秀明, 岩爪道昭. "深層学習抽出特徴量から生成した擬似特徴量を用いた不均衡データ多クラス画像分類." 人工知能学会全国大会論文集 第32回全国大会(2018). 般社団法人 人工知能学会, 2018.
- [3] 澤崎夏希, et al. "量的不均衡データに対する学習精度改善のための文書かさ増し手法." ARG WI2 No.11 (2017)
- [4] 立石健二, 石黒義英, and 福島俊一. "インターネットからの評判情報検索." 情報処理学会研究報告自然言語処理(NL) 2001.69 (2001-NL-144) (2001): 75-82.