

誰に向けた発言か？：ツイートの指向性推定

清基英則¹ 劉康明¹ 矢田竣太郎¹ 若宮翔子¹ 荒牧英治¹

¹ 奈良先端科学技術大学院大学

{kiyomoto.hidenori.kj5, lieu.kongmeng, s-yada, wakamiya, aramaki}@is.naist.jp

概要

新型コロナウイルスの拡大に伴い、政府や自治体はソーシャルメディアを用いた正確かつ迅速な情報発信が求められている。そのためには、特定の対象（年代や性別など）に向けて発信された情報を、その対象が自分に向けて発信されていると理解できるかどうか、すなわち「指向性」が重要である。情報発信者の属性を特定する研究は多いが、情報が対象とする受信者の属性を特定する研究は見受けられない。本研究では、Twitterにおける雑誌の公式アカウントが発信するツイートは読者層向けに最適化されていると仮定し、各雑誌の対象年齢と性別をラベル付けした指向性ツイートデータセットを用いて、ツイートがどの年齢、どの性別に向けられているものなのか機械学習モデルで分類した。この実験結果を分析し、指向性の定量的測定がもたらす価値を考察した。

1 はじめに

今や情報発信のインフラとしてソーシャルメディアはなくてはならないものとなっている。日本においては東日本大震災を契機として、SNSを活用した情報発信が大きく注目され始めた。その後も、2017年7月九州北部豪雨災害や2018年西日本豪雨災害でも救助を要請するツイートが多数投稿され、政府と市民間での情報発信の有効な手段として利用された[1]。このように、リアルタイム性と拡散性に優れたSNSを用いた情報発信は重要性を増しており、その動きは新型コロナウイルスの流行に伴い、さらに加速的に進んだ。例えば、厚生労働省により組織されたクラスター対策班は、分かりやすい情報を一般市民に届けるために、Twitterアカウントを開設し、情報発信を行った。さらに他の自治体も相次いで、SNSを通じた情報配信を活発化した。従来、SNSに消極的であった公的期間でさえも、SNSを重要なインフラとして捉え始めている。

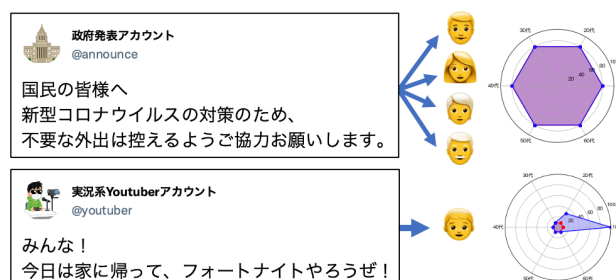


図1 ツイートの指向性推定。政府の発信のツイートには多方面への無指向性が、男子学生が好む話題のツイートには10代男性への超指向性が推定される。

一方で、情報発信における問題点も浮き彫りになっている。2020年6月に報告された新型コロナウイルス感染症に関する情報流通調査[2]において、政府が正しい情報を届けるための工夫を適切に行っていないと評価した人が多かった。その後、新型コロナウイルス感染症対策分科会は、情報の受信者側が関心を持ち、理解を深め、行動を変容させるような“対話ある情報発信”の実現に向けて、情報発信の強化を迅速に進めるよう政府に提言した[3]。しかし、不特定多数に向けた情報発信では、受信者側にうまく伝わらない傾向にあり、様々な年齢・性別・状況の人への情報発信は難しい[4]。

このような背景のもと、リスクコミュニケーションに配慮した情報配信を実現するために、本研究では、ツイートがどの対象（年代や性別など）に向けて発信されているか、この発信者が想定する受信者像を、本研究では「指向性」と呼び、ツイートの「指向性」に着目する。指向性とは、音響学的に、空中に出力された音や電波の伝わる強さが、方向によって異なる性質を意味する[5]。マイクやスピーカの指向性の種類には、全方向に対して発信される無指向性や特定の限られた範囲に発信される超指向性などがある。このアナロジーを情報発信に当てはめると、例えば政府の発表などは広い範囲の無指向性が、敬老会の予定などは高齢者に対する指向性が、カードゲームのイベント情報は若年層の男性に向け

表1 雑誌の公式アカウントが発信するツイートの例

対象属性	ツイート
10代女性	みんなどんな夏休みを過ごしてたのかな？
20代女性	韓国で盛んなハズせないファッション
30代女性	季節の変わり目をおしゃれに乗り切るっ！
40代女性	皆さんの興味を引く記事はありますか？
50代女性	みなさん、足元にお気をつけください。
60代女性	ひざ痛は今日みたいに「急な寒さ」で悪化します
10代男性	シンオウ地方で出会えるポケモンを大紹介！
20代男性	このチェキ、カッコいいですよ
30代男性	「0円マイカー」というサービスをご存知でしょうか？
40代男性	この柄に見覚えありますか？
50代男性	100年以上にわたって愛され続ける作品の魅力とは
60代男性	たかがねこ背だなんて軽視していると危険ですよ。

られた超指向性があると考えられる。

我々は、誰に向けたツイートであるかを言語的特徴で捉え、定量的に推定できるようになることで、情報発信者はより迅速かつ効果的に必要な情報を多くの人に届けることができると考えた。本研究では、対象としている層を公開している雑誌の公式アカウントが発信するツイートは指向性を持つと仮定し、雑誌の対象層をラベル付けしたデータセットを用いてツイートの指向性を推定するモデルを構築した。ツイートには30代から40代の男性向けといった複数対象層となっているものが含まれるため、マルチラベル分類タスクとして解き、指向性のある文章の定量的解析を行なった。

2 関連研究

メッセージ送受信に関わる者の属性（年齢、性別、職業など）を推定する場合、これまでよく研究されてきたのは送信者の属性である。Burgerら[6]は、1850万人のユーザが発信した約2億1300万件の多言語ツイートコーパスを構築し、ツイート発信者の性別を分類した。この研究で構築したBalanced Winnow2アルゴリズムがベースの機械学習モデルは、66.5%の精度で性別を分類することができた。またこのモデルが髪や愛に関連する単語、感嘆符や顔文字などの女性がよく用いる用語に強い影響を受けたことを示した。Morgan-Lopezら[7]はTwitterユーザの年齢を推定する際の言語的特徴を捉える研究を行った。ユーザを18から24歳の若年層、13から17歳の若者層、25歳以上の成人層と三つの年齢層に分け、最近発信した200件のツイートを収集し、学習データとして利用した。この研究で構築したロジスティック回帰モデルはF値0.72で分類する

ことができた。若年層は“学校”や“大学”などの用語を特徴として捉えることができたため、高い精度での分類を行うことができたが、成人層の発信するツイートの言語的特徴を捉えることは困難で、低い精度での分類となったことを示した。Abdul-Mageedら[8]は、約160万件のアラビア語のツイートからなるデータセットを構築し、書き手の性別（男性または女性）と年齢グループ（25未満、25-34、35以上）の合計6ラベルの分類を行った。彼らはBERTを用いて、年齢層を51.42%、性別を65.30%の精度で分類できることを示した。さらに、“they ask me”、“you dear”、“good evening”などの語句を分類の根拠としていたことから、性別間で書き方に明確な違いがあることを述べた。

これらの研究に対して本研究の指向性推定は、発信者が想定する受信者像を推定する点で新規である。

表2 指向性ツイートデータセットの統計量

ラベル	対象属性	ツイート数	Mean	SD
F10	10代女性	11,710	1301.1	405.7
F20	20代女性	7,929	881.0	588.1
F30	30代女性	7,230	1032.9	1274.4
F40	40代女性	12,596	1399.6	1209.7
F50	50代女性	17,164	1716.4	1215.1
F60	60代女性	21,087	1405.8	853.8
M10	10代男性	14,653	1971.0	303.8
M20	20代男性	26,926	1941.3	359.2
M30	30代男性	13,568	1356.8	865.5
M40	40代男性	22,118	1701.4	992.3
M50	50代男性	18,393	1839.3	1036.4
M60	60代男性	13,825	1874.7	570.3

3 データセット

本稿では、指向性を定量的に推定するために、指向性ツイートデータセット[9]を用いた。指向性ツイートデータセットは、特定の対象に向けて発信された、すなわち、指向性を持つツイートを収集して構築された。特定の対象に向けて発信された、すなわち、指向性を持つツイートを収集し、データセットの構築を行った。指向性を持つツイートとして、雑誌の公式アカウントが発信するツイートは読者層（年齢と性別）向けに最適化されているとみなし、収集対象とした。具体的には、対象年齢と性別ごと

表3 3つの提案モデルによる分類結果 (太字がモデル間の中での最も高い値)

Label	BERT			BERT MASKED PROPER NOUN			BERT MASKED NOUN		
	Precision	Recall	F_1 -score	Precision	Recall	F_1 -score	Precision	Recall	F_1 -score
F10	0.85	0.87	0.86	0.85	0.76	0.80	0.80	0.55	0.65
F20	0.88	0.72	0.79	0.92	0.50	0.64	0.65	0.57	0.61
F30	0.89	0.84	0.87	0.88	0.76	0.82	0.80	0.66	0.72
F40	0.88	0.82	0.85	0.87	0.76	0.81	0.75	0.67	0.71
F50	0.89	0.87	0.88	0.87	0.79	0.83	0.77	0.68	0.72
F60	0.92	0.92	0.92	0.85	0.88	0.87	0.84	0.65	0.73
M10	0.92	0.79	0.85	0.80	0.76	0.78	0.83	0.48	0.61
M20	0.90	0.82	0.86	0.79	0.79	0.79	0.71	0.61	0.66
M30	0.90	0.86	0.88	0.84	0.84	0.84	0.76	0.77	0.76
M40	0.94	0.93	0.93	0.90	0.91	0.90	0.85	0.84	0.84
M50	0.93	0.86	0.89	0.87	0.85	0.86	0.85	0.66	0.74
M60	0.93	0.88	0.91	0.84	0.87	0.86	0.90	0.54	0.68
macro avg	0.90	0.85	0.87	0.86	0.79	0.82	0.79	0.64	0.70

に雑誌をまとめたウェブサイト¹⁾²⁾³⁾を参考に、10代から60代までの男女別の雑誌リストを作成し、Twitterにおける公式アカウント名を取得した。この結果、女性向けと男性向け雑誌のアカウント数はそれぞれ71と35であった。なお、対象とする性別が示されていない雑誌は、男性と女性の両方を対象とする雑誌アカウントとして、幅広い年齢層を対象とする雑誌は、複数の年齢層を対象とする雑誌アカウントとして扱った。例えば、“@safari_online”というアカウントは20-30代の男性向けの雑誌の公式アカウントであるため、20代男性と30代男性の両方に含まれる。ツイート取得方法やノイズ除去処理などの詳細については[9]を参照されたい。データセットにおけるツイートの総数は187199であり、その概要は表2に示す。

4 提案手法

本研究では、ツイートがどの年齢や性別に向けて発信されているのかを、マルチラベル分類タスクとして解く。そのために指向性ツイートデータセットを用いて、文章分類モデルを構築する。分類モデルにはBidirectional Encoder Representations from Transformers (BERT) [10]を日本語コーパスで事前学習したモデル⁴⁾を採用し、指向性ツイートデータ

セットでファインチューニングした。このモデルは12層のエンコーダ、768次元の隠れ層、12個のアテンションヘッドから構成される。学習条件として、最適化手法はAdamW、学習率は 1.0×10^{-5} 、エポック数は5、バッチサイズは32である。各ツイートには複数のラベルが付与されるため、マルチラベル分類モデルを構築した。BERTの最終層の出力を全てのトークンに渡って平均化、その値を線形変換したものを分類スコアとし、スコアが正の値の場合に該当ラベルとした。

さらに、ツイートの話題ではなく、語尾や文体といった言語的特徴を考慮した分類のために、固有名詞と名詞をそれぞれ特殊トークンとして学習させた2つのモデル(BERT MASKED PROPER NOUNとBERT MASKED NOUN)を構築した。固有名詞と名詞の判定には、形態素解析器MeCab [11]を用いた。特殊トークンに置き換えることで、BERTの学習に影響を与える単語を意図的に制限することができる [12]。

5 結果

指向性ツイートデータセットを9(訓練):1(評価)の割合で分割し、BERT, BERT MASKED PROPER NOUN, BERT MASKED NOUNの3つのモデルによるツイートの分類精度を評価した。評価指標には、Macro-F1と各ラベルのF値を用いた。表3に結果を示す。

各モデルのF値の平均はそれぞれ0.87, 0.82, 0.70と高い精度となった。2節で述べたツイート文から

1) <https://www.magazine-data.com/menu/age.html>
 2) <https://www.san-an.co.jp/media/senior.html>
 3) https://en.wikipedia.org/wiki/List_of_manga_magazines
 4) <https://huggingface.co/cl-tohoku/bert-base-japanese-whole-word-masking>

- (1) [CLS]こん##にち##は今日から3月がスタート!!そしてニコラ発売日今##月は中3で思い出づくりしたよぜ##ひチェックしてみてね##-!
- (2) [CLS]こん##にち##は今日から3月がスタート!!そして[Proper noun]発売日今##月は[Proper noun]で[Proper noun]したよぜ##ひチェックしてみてね##-!
- (3) [CLS]こん##にち##は[Noun]から[Noun]が[Noun]!!そして[Noun][Noun][Noun]は[Noun]で[Noun]したよぜ##ひ[Noun]してみてね##-!
- (4) [CLS]まもなくあの3.11から10年です。月刊4月号に寄稿しました。よろ##し##ければお手に取ってみて下さい。
- (5) [CLS]まもなくあの[Proper noun]から10年です。月刊4月号を寄稿しました。よろ##し##ければお手に取ってみて下さい。
- (6) [CLS]まもなくあの[Noun]から[Noun]です。[Noun][Noun][Noun]を[Noun]しました。よろ##し##ければ[Noun][Noun]に取ってみて下さい。

図 2 Attention weight の可視化結果。 (1)(2)(3) は F10 に、 (4)(5)(6) は F60 と M60 に BERT, BERT MASKED PROPER NOUN, BERT MASKED NOUN でそれぞれ正しく分類された例。 マスキングにより、文末や文体の表現に重み付けされている。

書き手の年齢層を推定するモデル [7] の F 値が 0.72 で、一定の有用性を評価しているため、今回の我々の結果もある一定の精度を示したと判断した。

モデルごとに確認すると、BERT の F 値の平均が 0.87 であり、F40 から F60, M30 から M60 の F 値が突出して高い結果となった。ほぼ全てのラベルが F 値が 0.8 を超えていることから、高精度での分類結果となった。BERT MASKED PROPER NOUN については、F 値の平均値は 0.82 であり、BERT に比べてやや低い精度となった。F40 から F60, M30 から M60 の F 値が依然として高いままであったが、他のラベルは下がる結果となった。BERT MASKING NOUN の F 値の平均値は 0.70 となり、他の 2 つのモデルと比べて、適合率、再現率、F 値の全てが下がった。M40 の F 値のみ 0.85 と高いままであるが、他のラベルは BERT に比べて大きく精度が下がった。

6 考察

本節では、分類結果と各モデルの Attention weight の可視化結果をもとに、各モデルの結果と分類の判断根拠について考察する。まず F30 から F60, M30 から M60 のラベルの精度がモデルを問わずに高かった。これは先行研究 [9] では、低い精度であったラベルであり、マルチクラスからマルチラベルにタスクを変更したことで、複数のラベルが付いたツイートの分類精度が大幅に向上したことがわかる。モデルごとに確認すると、通常の BERT の精度が最も高いことから、固有名詞と名詞を含んだ方が分類精度が高くなることがわかった。図 2 は BERT の Attention weight の可視化結果である。図 2 の (1) と (4) より、BERT は、“ニコラ”、“中 3”、“思い出づくり”、“3.11” などの特徴的な言葉を根拠に分類を行っていることがわかる。BERT MASKED PROPER NOUN は、どのラベルも高い精度で推定することができており、固有名詞をマスキングしても、高い精度で分類されていることがわかる。

図 2 の (2) と (5) より、“スタート”、“チェック”、“寄稿” という特徴的な単語を根拠に分類している一方、“みてねー!”, “みて下さい” といった文末の表現により重み付けされていることから、対象となる属性に向けた文体が学習されていることがわかる。BERT MASKED NOUN の精度は他のモデルと比べると下がるが、依然として高い精度で分類ができており、図 2 の (3) と (6) を見ると、“こんにちは”、“したよ”、“ければ”、“みて下さい。” といった文体の表現を根拠としており、より文体的特徴を学習していることが推測される。

7 おわりに

本研究では、特定の年代性別のユーザに向けて発信されたツイートは指向性を持つと仮定し、それを定量的に推定するために、雑誌の公式アカウントによるツイートを収集し、BERT による分類モデルを構築した。実験では、通常の BERT モデルにより一定の精度で推定できることを示し、さらに固有名詞と名詞を特殊トークン化したモデルとの比較により、指向性を持つツイートの言語的特徴を考察した。マルチラベル分類タスクとして解くことで、どのラベルも高い精度で分類することができた。

今後の課題としては、雑誌の公式アカウント以外の指向性ツイートの取得と実用的なツールの開発である。現在、インフルエンサーが発信するツイートを対象として検討している。これにより指向性を雑誌の公式アカウント以外の発信するツイートで確認することができる。また現在、試験的に開発したシステム VOICE2PEOPLE (付録の図 3)⁵⁾ を公開している。VOICE2PEOPLE は今回の実験で用いたモデルを組み込んだ Web アプリケーションで、入力したテキストの指向性を推定することができる。今後、ユーザによる評価実験を行い、指向性の定量的評価の有用性や拡張性について検証する予定である。

5) <https://voice2people.netlify.app/>

謝辞

本研究はAMEDの課題番号JP22mk0101229, JSPS科研費JP20K19932, JP19H01118, JP22K12041, Yahoo株式会社共同研究費の支援を受けたものです。

参考文献

- [1] 佐藤翔輔, 今村文彦. 2018年西日本豪雨災害における「#救助」ツイートの実態:2017年7月九州北部豪雨災害との比較分析. 自然災害科学, Vol. 37, No. 4, pp. 383–396, 2019.
- [2] 総務省. 新型コロナウイルス感染症に関する情報流通調査 報告書, 2020. https://www.soumu.go.jp/menu_news/s-news/01kiban18_01000082.html.
- [3] 新型コロナウイルス感染症対策分科会. “対話ある情報発信”の実現に向けた分科会から政府への提言 令和2年11月12日(木), 2020. <https://www.cas.go.jp/jp/seisaku/ful/bunkakai/seifu.teigen.15.pdf>.
- [4] 公益財団法人東京市町村自治調査会. 誰にも伝わる情報発信に関する調査研究 報告書. 2017. <https://www.tama-100.or.jp/cmsfiles/contents/0000000/672/0.darenimotutawarujouhouhassin.pdf>.
- [5] コトバンク. 第2版, 世界大百科事典内言及日本大百科全書(ニッポニカ), “指向性とは”, 2022. <https://kotobank.jp/word/%E6%8C%87%E5%90%91%E6%80%A7-73018>.
- [6] John D. Burger, John Henderson, George Kim, and Guido Zarrella. Discriminating gender on Twitter. In **Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing**, pp. 1301–1309, Edinburgh, Scotland, UK., July 2011. Association for Computational Linguistics.
- [7] Antonio A Morgan-Lopez, Annice E Kim, Robert F Chew, and Paul Ruddle. Predicting age groups of twitter users based on language and metadata features. **PloS one**, Vol. 12, No. 8, p. e0183537, 2017.
- [8] Muhammad Abdul-Mageed, Chiyu Zhang, Arun Rajendran, AbdelRahim Elmadany, Michael Przystupa, and Lyle Ungar. Sentence-level bert and multi-task learning of age and gender in social media. **arXiv preprint arXiv:1911.00637**, 2019.
- [9] 清基英則, 劉康明, 矢田竣太郎, 若宮翔子, 荒牧英治. 言語的特徴を用いたツイートの指向性推定. 信学技報, Vol. 122, No. 88, pp. 19–24, 2022.
- [10] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In **Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)**, pp. 4171–4186. Association for Computational Linguistics, 2019.
- [11] Taku Kudo. Mecab: Yet another part-of-speech and morphological analyzer. <http://mecab.sourceforge.jp>, 2006.
- [12] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac,

Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers: State-of-the-art natural language processing. **arXiv preprint arXiv:1910.03771**, 2019.

A 参考情報

表4 各ラベルの重要語 (TF-IDF 値)

順位	F10	F20	F30	F40	F50	F60
1	発売 (0.38)	steady(0.44)	さん (0.64)	さん (0.71)	さん (0.71)	さん (0.50)
2	ニコラ (0.33)	さん (0.36)	レシピ (0.24)	家庭画報 (0.22)	おとなの週末 (0.23)	東京人 (0.36)
3	popteen(0.22)	cawaii(0.28)	リンネル (0.24)	九星気学 (0.17)	紹介 (0.21)	和楽 (0.27)
4	さん (0.19)	付録 (0.20)	otona(0.21)	こと (0.17)	こと (0.19)	暮しの手帖 (0.22)
5	本日 (0.19)	spring(0.20)	muse(0.18)	レシピ (0.17)	家庭画報 (0.13)	特集 (0.22)
6	みんな (0.15)	特集 (0.16)	こちら (0.14)	真木 (0.12)	レシピ (0.13)	紹介 (0.19)
7	先生 (0.15)	チェック (0.14)	おすすめ (0.14)	こちら (0.10)	九星気学 (0.10)	こと (0.19)
8	りぼん (0.15)	本日 (0.14)	紹介 (0.14)	クロワッサン (0.08)	皆様 (0.09)	発売 (0.18)
9	コミ (0.12)	mini(0.13)	lee(0.12)	占い (0.08)	人気 (0.09)	記事 (0.12)
10	表紙 (0.12)	こちら (0.13)	mamagirl(0.12)	よう (0.08)	よう (0.09)	歴史 (0.08)

順位	M10	M20	M30	M40	M50	M60
1	発売 (0.42)	発売 (0.44)	さん (0.27)	週刊エコノミスト (0.41)	さん (0.42)	東京人 (0.51)
2	先生 (0.30)	さん (0.34)	こと (0.27)	さん (0.32)	週刊エコノミスト (0.41)	さん (0.42)
3	本日 (0.27)	本日 (0.28)	アイテム (0.18)	こと (0.32)	こと (0.27)	特集 (0.27)
4	コミックス (0.26)	smart(0.19)	out(0.15)	online(0.21)	おとなの週末 (0.26)	発売 (0.21)
5	ジャンプ sq(0.18)	公開 (0.18)	紹介 (0.15)	日本 (0.14)	online(0.19)	記事 (0.15)
6	掲載 (0.16)	掲載 (0.18)	理由 (0.15)	理由 (0.12)	紹介 (0.14)	こと (0.14)
7	公開 (0.16)	連載 (0.14)	人気 (0.15)	よう (0.12)	発売 (0.12)	中央公論 (0.14)
8	少年ガンガン (0.14)	最新 (0.13)	ブランド (0.13)	ため (0.12)	特集 (0.12)	紹介 (0.13)
9	さん (0.14)	こと (0.12)	別注 (0.13)	九星気学 (0.11)	日本 (0.11)	演劇界 (0.12)
10	更新 (0.12)	表紙 (0.11)	go(0.12)	紹介 (0.10)	九星気学 (0.11)	歴史 (0.11)



図3 Voice2People の画面