

# BERTopicを利用したTwitter（現X） における くまモン・ファンダムの言論分析

畠山 真一

## 1 はじめに

本研究では、畠山（2023）で提示したデータを再分析し、くまモンに関する「Twitter（現X）」への投稿をもとに、コアなファン層の特質について明らかにすることを目的とする。本研究で対象とするデータは、2021年9月29日～2022年1月5日の期間に収集されたものであるため、現在Xと呼ばれているSNSプラットフォームを旧来の「Twitter」と呼ぶ。本研究では、投稿のトピックを深層学習を応用したBERTopicと呼ばれる手法を用いて分析し、その結果としてくまモン・ファンの成長過程とくまモンを実体化の間に関連性があることを明らかにする。

## 2 使用したデータについて

本分析のベースとなるデータは、「Twitter」社から提

供されているAPIを利用して、「くまモン」を本文に含む投稿を収集したものであり、その基本的な情報は次のとおりである。

- ・収集した期間…2021年9月29日～2022年1月5日

- ・本文に「くまモン」を含む投稿のすべてから、次の条件に当てはまるものを削除した。

- 「送料無料」という広告目的のキーワードを含む投稿を削除した<sup>(1)</sup>。

- 同一内容の複数投稿については、1つのみを残した。

- 投稿から絵文字、顔文字、URLを削除した。

- neologd<sup>(2)</sup> モジュールを利用して正規化した<sup>(3)</sup>。

投稿そのものに加えて、上述のデータには、次のような情報が含まれている。

・ ユーザ名、ユーザのプロフィール欄、投稿日時、返信 (reply) ならば誰宛か、RT 数など  
以後、上述のデータをくまモンデータと呼ぶ。

### 3 くまモンデータの基本情報

くまモンデータの基本的な情報は次のとおりである。

- ・ 投稿の総数は、66,214 である。
- ・ くまモンデータを構成するユーザの総数は、27,935 である。

当該期間 (2021 年 9 月 21 日 ~ 2022 年 1 月

5 日) の投稿数を横軸

に、その数の投稿を行ったユーザの人数を縦軸にとったヒストグラムを作ると図 1 のようになる (縦軸は対数)。

図 1 が示すように、当該期間に 1 ~ 2 個の投稿をするユーザが圧倒

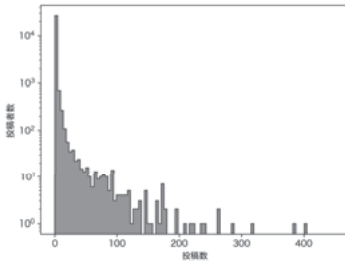


図 1 投稿数とユーザ数の連関

的に多いものの、この期間に 403 ポストしたユーザもあり、ある程度ロングテールの様相を見せている。

続いて、基本的な形態素解析結果について述べる。

本研究では、次のようなセッティングで投稿内容の形態素解析をおこなった。

- ・ 形態素解析器として、mecab を使用した。
- ・ システム辞書は、mecab-ipadic-NEologd を使用するとともに (佐藤・橋本・奥村 2017)、独自のくまモンデータに含まれる独自の表現 (たとえば、「おはくま」など) を含んだ個人辞書を使用した。

- ・ 形式名詞など、一般的な stop-word については除外した。

このようなセッティングのもとでくまモンデータに含まれる投稿内容に対して形態素解析を行った結果、投稿に出現する名詞、動詞、形容詞、間投詞の上位 10 語として表 1 が得られた。

表 1 単語の出現数

単語	出現数
くまモン	66619
さん	8363
今日	6928
くま	5053
熊本	5052
おほくま	4399
モンちゃん	3919
笑	2893
ありがとう	2875
ちゃん	2758

## 4 BERTopicを用いた投稿からのトピック抽出

### 4.1 コア・ユーザーとライト・ユーザー

前節で述べたように、くまモンデータは、極めて活発にポストを行うユーザーと1、2度の言及に留まるユーザーが併存しているのが特徴である。ここで、全投稿数をおおまかに3分割し、投稿数とユーザー数の関係を見てみると、次の表2が得られた（カッコ内が当該ユーザーの投稿数である）。

このように、約2%にすぎない550ユーザーが、全投稿のおよそ40%弱をポストしている。これらの活発なユーザーを、これ以降、コア・ユーザーと呼ぶことにしよう。

反対に、圧倒的多数を占める（76%程

表 2 投稿数と投稿者数

	投稿数 1 回	投稿数 2 回以上 12 回未満	投稿数 12 回以上
ユーザー数	21298	6087	550
投稿数	21298 投稿	18700 投稿	26216 投稿

度）当該期間に「くまモン」という語を含む投稿を1回のみ行った投稿者をライト・ユーザーと呼ぶことにしよう。

これ以降、コア・ユーザーとライト・ユーザーの投稿からBERTopicを用いてトピックを抽出し、比較分析を行う。

### 4.2 BERTopicとは

前節で述べたコア・ユーザーは、どのような投稿を行っているのだろうか。

次節ではこの点を、いわゆる深層学習にもとづく自然言語処理をリードしたと言える言語モデルの一種であるBERT (Bidirectional Encoder Representations from Transformers) をトピック分析に応用したBERTopicを利用した分析を提示するが (Devlin et al 2019; Grootendorse 2022)、本節では、分析ツールであるBERTopicの基本的なフレームワークを提示する。

BERTopicは、おおまかに言えば、次の3段階のステップを踏むことで文章のトピックを抽出する。

1. 文章を埋め込みベクトルに変換し (sentence-BERT が利用される)、類似度を計算できるようにする。

2. 得られたベクトルの次元削減を行った上で (UMAP を使用する)、クラスターリングを行い (類似した内容を持つ文章をグループ化する)、それによって得られたクラスターを一つのトピックとみなす。なお、クラスターリング手法に関しては、HDBSCAN<sup>7</sup>、k-Means などが選択可能である。

3. c-TF-IDF<sup>8</sup> もしくは BM25<sup>9</sup> に基づき、各クラスターの特徴語を抽出する。これが、各トピック (先の各クラスター) を代表する主要語である。

BERTopic は従来の LDA (Latent Dirichlet Allocation) に代表されるようなトピック分析 (Beaulieu et al. 2003) に比べ、語彙や文の意味を分散表現として利用しながら、意味の類似性に基づくクラスターリングを行っており、LDA に比べて高い精度でトピック抽出ができることが報告されている (川原, 2023)。

BERTopic は極めて多彩なパラメータを持つが、本研究では、次のようなセッティングでトピック抽出を行った。

1. sentence-BERT の日本語モデルとして、sentence-bert-base-japanese-v2<sup>10</sup> を利用した<sup>(3)</sup>。

2. 次元削減には、UMAP を、クラスターリングには、HDBSCAN を使用した (ハイパーパラメータについてはデフォルト値を利用した)。

3. 特徴語検出には、文章長の影響を受けにくいとされる BM25 を利用した。

#### 4. 3 コア投稿者言説からのトピック抽出

前節で述べたようなセッティングのもとで、コアユーザの投稿からなる文書集合に対して、BERTopic を用いてトピック抽出したところ、約半数 (13, 075 投稿) の投稿が分類不能となったものの、18 個のトピックが得られた。表 3 は、当該トピックに分類される投稿数の上位 5 グループを表にまとめた

ものである。なお、表3における「ラベル」名は、本稿の筆者がトピック主要語をもとに解釈した結果であり、BERTopicからは得られない要素である。

紙幅の制約のため、トピックへの貢献度についてTopic 0のみグラフ化すると、図2のようになる。

表3・図2から理解されるように、分類可能であった文書集合の半数が、くまモンを擬人化した「くまモンへの手紙」であることは特筆すべき現象と思われる。

同様にTopic 1は、くまモン体操をTopic 2はくまモンスクエアでの出演に関するクラスタであり、コア・ユーザが単なる図像では

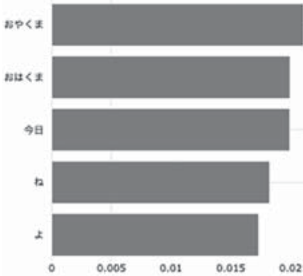


図2 トピックへの貢献度

表3 BERTopic によるトピック抽出 (コア・ユーザ)

Topic	分類された投稿数	トピック主要語	ラベル
0	12688	おやくま・おはくま	くまモンへの個人的な手紙
1	76	体操・ダンス	くまモン体操
2	70	11月・スクエア	くまモン出演情報
3	59	ぼっち・ジョージ	クリスマス
4	29	合唱・絵描き歌	くまモンの絵描き歌

なく、「元氣よく動き回り現実世界に身体を持つ」くまモンについて数多く言及していることが理解される。

#### 4. 4 ライト投稿者言説からのトピック抽出

前節と同様の手法でライト・ユーザの投稿からBERTopicを用いてトピック抽出したところ、分類不可能な投稿がコア・ユーザの場合と同じくおよそ半数あった。抽出されたトピックは、144個あり(意味なクラストが144個あったと言い換えられる)、コア・ユーザによる投稿集合と異なり、多様なトピックが出現していることが理解される。次の表4は、当該トピックに分類される投稿数の上位5グループを表にまとめたものである。

これらのトピック集合は、前節のコア・ユーザのものと大きく

表4 BERTopic によるトピック抽出 (ライト・ユーザ)

Topic	分類された投稿数	トピック主要語	ラベル
0	3658	かわいい・好き	くまモンへの愛着
1	552	熊本城・熊本	熊本県の象徴
2	541	グッズ・買った	モノとしてのくまモン
3	454	美味しい・クレープ	パッケージアイコンとしてのくまモン
4	153	歌・fns 歌謡祭	TVの中のくまモン

異なっている。それは、身体性の顕現性に関する差である。ライト・ユーザのくまモンへの愛着を示す「Topic 0」に含まれる投稿を見ると、身体性が感じられないものも多い（感じられるものももちろんあり、コア・ユーザと同様に「くまモンへの手紙」と解釈できるものも数多くある）。

- マスクしてサムズアップして行くくまモン良すぎる

- ひこニャン、くまモンは双璧

- くまモン笑笑たしかに強そう

また、「Topic 1」から「Topic 3」までは、身体を持たないアイコンとしてのくまモンがトピックとなっている投稿と解釈できる。ただし、「Topic 4」については、TVで見るという制約はあるものの、身体を持ち動き回るくまモンをトピックとする投稿集合と解釈できる。

このように、「Topic 4」という例外はあるものの、ライト・ユーザによる投稿に出現するくまモンは身体性があまり感じられず、一種の「アイコン」ととどまっていると解釈するのが妥当であろう。

## 5 コア・ユーザとライト・ユーザの投稿集合の差異について

先述のように、コアユーザは「くまモンへの手紙」という性質を持つ投稿が圧倒的に多く、身体を持ち現実世界に定位するくまモンへの愛着を表現するという傾向を持っている。コア・ユーザも最初の段階ではライト・ユーザであったであろうということを考え合わせると、アイコンとしてのくまモンを愛好するライトなファンから成長した結果、コア・ユーザは、くまモンに身体性を見出すようになったと分析するのが妥当であろう。

ここで問題となるのが、ファンの「成長」を支える環境をどのように構築するかという点である。これに関しては、佐藤（2018）の議論されているファンダムへのアプローチが有効であると考えている。

## 6 おわりに

本レポートでは、「Twitter」における「くまモン」を含む投稿をBERTopicによって分析することを通じて、くまモン・ファンの成長過程を明らかにした。

注

- (1) 畠山 (2023) のツイート削除が充分でないため、本稿では徹底した。
- (2) neologd と <https://pypi.org/project/neologd/> からダウンロード。
- (3) <https://huggingface.co/sonoisai/sentence-bert-base-ja-mean-tokens-v2>

参考文献

- Biel, D.M., Ng, A.Y., and Jordan, M.I.(2003). "Latent Dirichlet allocation." *Journal of Machine Learning Research*, 3, pp. 993-1022.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). "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, Minneapolis, Minnesota. Association for

Computational 6 Linguistics.

- Grootendorst, M. (2022). "BERTopic: Neural topic modeling with a class-based TF-IDF procedure." *arXiv preprint arXiv:2203.05794*.
- 畠山真(2023)「Twitterにおけるくちモンの評判分析ーファンコミュニティの視点からー」尚絅語文第12号, 7-16
- 堀田治(2015)「超高関与消費のマーケティングパター」『AD STUDIES』, 51, 15-20
- 川原一修(2023)「トピックモデルによる市場変動要因の抽出」言語処理学会第29年次大会(NLP 2023), pp.2190-2194, 言語処理学会
- 佐藤一誠(2015)『トピックモデルによる統計的潜在意味解析』コロナ社
- 佐藤尚之(2018)『ファン・ベース』ちくま新書
- 佐藤敏紀・橋本泰一・奥村学(2017)「単語分ち書き辞書mecab-ipadic-NEologdの実装と情報検索における効果的な使用法の検討」言語処理学会第23回年次大会(NLP 2017), pp.875-878, 言語処理学会