

# 「フェイクニュース」の理解に向けて: アノテーションスキームの提案と日本語データセット構築

村山太一\* 久田祥平\* 上原誠 若宮翔子 荒牧英治  
奈良先端科学技術大学院大学

{murayama.taichi.mk1, s-hisada, uehara.makoto.ug2, wakamiya, aramaki}@is.naist.jp

## 概要

社会問題となっているフェイクニュースに対して、検出タスクやデータセット構築などの研究が取り組まれている。これらの既存研究はそのニュースが正しいかどうかの事実性ばかり着目する。しかし、複雑な現象であるフェイクニュースを理解するためには、事実性の側面だけでなく、発信者の意図や社会への有害度合い、ニュースの種類など、様々な観点から捉えることが重要である。本研究では、フェイクニュースの様々な側面を捉える、きめ細かなラベリングを行う新しいアノテーションスキームを提案する。そして、このアノテーションスキームに基づいた初めての日本語フェイクニュースデータセットを構築し、<https://hkefka385.github.io/dataset/fakenews-japanese/> で公開している。

## 1 はじめに

ネット上の情報の信頼性を貶める「フェイクニュース」に対して、研究者はフェイクニュース検出タスクに取り組んだり、リソースとして FakeNewsNet [1] や CoAID [2] などのフェイクニュースデータセットを構築している。これらの既存の研究やデータセット構築では、そのニュースが正しいかどうかといったニュースの事実性の側面に着目している。しかし、事実性のラベルだけで、フェイクニュースやそれによって引き起こされる現象を十分に理解することは困難である。フェイクニュースをより理解するためには、発信者の意図や社会への有害性などの様々な視点から、フェイクニュースを捉えることが重要である。

本研究では、「フェイクニュース」の定義や既存のフェイクニュース検出データセットを調査した結果得られた課題を踏まえ、フェイクニュースを様々

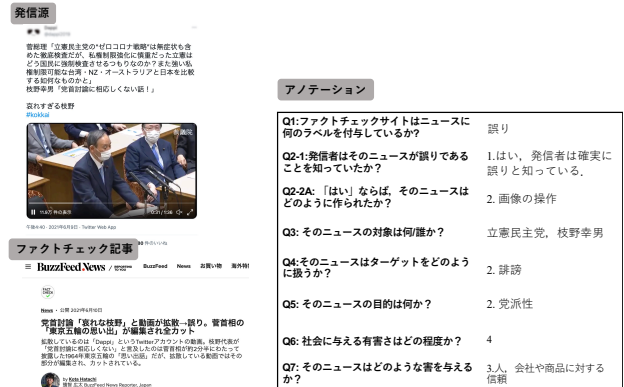


図 1: アノテーション例: 発信源は党首討論の動画をのせたもので、ファクトチェック記事ではこの動画は討論の一部を切り取ったもので野党党首の印象を悪くしていると指摘する。

な観点で捉えるためのアノテーションスキームを提案する。具体的には、ニュースの 1. 事実性, 2. 発信者の意図, 3. 対象, 4. 対象の扱い, 5. 目的, 6. 社会への有害度, 7. 有害性の種類という 7 つの観点でこのスキームは構成される。この詳細なアノテーションは、「発信者がそのニュースを嘘だと知っているかどうかで、返信がどのように変わるのか?」など、フェイクニュースという複雑な現象の深い理解に繋がる。さらに、提案したアノテーションスキームに従って、日本語のフェイクニュースデータセットを構築する (例: 図 1)。これは日本語で構築された初めてのデータセットであり、日本のフェイクニュースの理解や研究に有用である。今後、日本語以外の言語においても、提案したアノテーションスキームに従ってデータセットを構築し、国家間におけるフェイクニュースの性質の比較などを通して、フェイクニュースをより多角的な視点で分析していく予定である。

## 2 予備的調査

\* equal contribution

## 2.1 「フェイクニュース」の定義

情報科学の研究において、フェイクニュースには広義と狭義の定義が用いられる。広義の定義は、“fake news is fabricated information that mimics news media content in form but not in organizational process or intent. [3, 4]”である。この定義は、情報の真偽を重視し発信者の意図を考慮しないもので、風刺やパロディなどの種類のニュースもカバーするものである。しかし、この定義を採用した研究は多くない [5, 6, 7]。一方、狭義の定義は、“a news article that is intentionally and verifiably false. [8, 9]”というもので、発信者の意図を重視する。この定義を採用する研究は多い [10, 11, 12] が、多くのデータセットは、ニュースの真偽に関するファクトチェック機関の判定のみをラベルとして構築されている。

このようなフェイクニュースという言葉の曖昧性に対して批判の声が上がっている。例えば、イギリス政府はフェイクニュースという言葉の曖昧性から公式文書には使用しないことを決定した [13]。また、ネット上における誤った情報の対抗に向けたプロジェクトである First Draft のリーダーの Claire Wardle 氏は、フェイクニュースという言葉は情報の信頼性の問題を十分に説明できないとし、Misinformation, Disinformation, Malinformation の3つの観点で説明すべきと述べている [14]。このように、フェイクニュースの概念は曖昧であるため、そのニュースはフェイクかどうかではなく、様々な観点から説明していくことが重要である。

## 2.2 フェイクニュース検出データセットの課題

ニュースの真偽を評価するフェイクニュース検出タスクのために、数多くのデータセットが構築されている。本節では、51のフェイクニュース検出データセット<sup>1)</sup>を検証し、解決すべき4つの課題を示す。

**意図** 2.1節で説明したように、ほとんどのデータセットは狭義の定義を採用しているにも関わらず、各ニュースの事実性に着目したラベルのみ持って広義の定義に基づいて構築されている。これは、技術開発における定義と言葉の定義の間に乖離があることを意味する。また、意図を持って作成されたニュースは社会的影響力が高く、誤ったニュースが広まる一要因となってい

る [15]。このような背景や、説明可能性の高い検出モデルを構築するためにも、発信者の意図をアノテーションすることが重要である。

**社会への有害性** フェイクニュースは社会に悪影響を与えるが、その規模や種類は様々である。例えば、パロディなどの明らかに嘘であるニュースは社会への害は少ないが、選挙に関する誤ったニュースは人々の意思決定に影響を与え、有害性が大きい。既存のデータセットではこの観点をカバーできていないが、これはファクトチェックの優先度の判断にも有用である。

**言語** フェイクニュースの言語的特徴や拡散パターンは国や言語によって異なる可能性があるにも関わらず、ほとんどのデータセットは米国社会のイベントを対象とし英語で構成されている。一方で、COVID-19の世界的な Infodemic の影響から、英語以外のフェイクニュース検出データセットが増加している。具体的には、51件中、英語以外の言語を含むものは11件で、そのうち8件がCOVID-19に関するものであった。COVID-19以外の様々なトピックを対象とした英語以外のデータセットを構築していくことは、言語横断的な分析や言語に依存しない特性の特定が期待できる。

**ラベル** 機械学習モデルの適用の容易さから、51件中33件のデータセットには、フェイクかどうかの二値ラベルが付与されている。その他のものは、ファクトチェック組織の評価に基づいて構築されており、データセットごとに分類基準が異なり、利用しづらい。「意図」と「社会への有害性」も関連するが、より一般的で頑健な検出モデルを構築するためにも、fine-grainedで一貫したアノテーションスキームが必要である。

## 3 アノテーションスキーム

上記の議論や得られた洞察に基づき、7つの質問からなるフェイクニュースのためのアノテーションスキームを構築した。Q1-Q5は細かなラベルを付与することを目的とした質問で、「意図」と「ラベル」の課題をカバーするものである。Q6とQ7では、COVID-Alam [16]を参考に、社会への有害性の特定を試みている。我々のアノテーションスキームは、ファクトチェック組織の評価に依存しないラベリングが可能となり、解釈性の高い検出モデルの構築が期待できる。ファクトチェック記事や発信源と

1) データセットを <https://hkefka385.github.io/dataset/fakenews-japanese/#surveylist> に示す。

なったソーシャルメディアの投稿に基づき、これら7つの質問にアノテーターが回答する。

**Q1: ファクトチェックサイトはニュースに何のラベルを付与しているか?**

この質問は、ファクトチェックサイトが付与したラベルである「正確」や「虚偽」などの情報を記載する。「正確」または「ほぼ正確」のラベルが付与された場合、以下全ての質問をスキップする。

**Q2-1: 発信者はそのニュースが誤りであることを知っていたか?**

選択肢: 1. はい、発信者は確実に誤りと知っている、2. はい、発信者はおそらく誤りと知っている、3. いいえ、発信者はおそらく誤りと知らない、4. いいえ、発信者は確実に誤りと知らない

この質問は、課題の一つである「意図」に対応する質問で、4つの選択肢から1つ選択する。1. または2. が選択された場合、既存研究 [3] に基づきそのニュースを「Disinformation」とみなす。3. または4. が選択された場合、発信者は誤ったニュースを流す意図が無いことを意味するため、「Misinformation」とみなす。このような意図のラベリングは、「どのようなニュースを誤りと知らずに拡散しているのか?」など、ニュースの種類によってユーザ行動がどのように異なるかを明らかにする可能性がある。この回答に応じて、誤りが生じた理由に関する質問 Q2-2A もしくは Q2-2B に回答する。

**Q2-2A: 「はい」なら、そのニュースはどのように作られたか?**

選択肢: 1. 創作されたコンテンツ、2. 画像の操作、3. テキストの操作、4. 誤ったコンテキスト

この質問は、意図的に発信されたニュースがどのように作成されたかを理解するためのもので、4つの選択肢を設定した。まず、完全に創作されたニュースか、何かしらのリソースを改竄したニュースかに分類される。前者を「創作されたコンテンツ」とする。後者は、改竄の対象によって3パターンに分け、画像や動画を加工したものを「画像の操作」、ニュースのテキストやメッセージを加工したものを「テキストの操作」、コンテンツ自体は本物だが、誤った文脈情報とともに共有するものを「誤ったコンテキスト」とする。

**Q2-2B: 「いいえ」なら、発信者はどのように誤解したか?**

選択肢: 1. 他の情報源の信頼、2. 不十分な理解、3. ミスリーディング

この質問は、どのように意図しない誤ったニュースの拡散を防いでいくかを考える上で重要である。意図しない誤った情報が発信される背景として、3つの選択肢を設定した。1つ目は、他の情報源を信頼したことによって起こるもの。英語を母国語としないユーザが、英語の記事などを誤訳したり、誤った情報を信用したりすることで起こる。2つ目は、発信者がニュースを十分に読み込まないことによって、不十分な理解や思い込みによって生じるもの。3つ目は、発信者がニュースを十分に理解している可能性はあるが、それを十分に読者に伝えられなかったことで起こるもの。

**Q3: そのニュースの対象は何/誰か? (複数回答可)**

誤ったニュースの対象、つまり主に影響を受ける対象は、ニュースのクラスタリングや検索において有用な情報である。このような情報を特定するタスクはまだ行われていないが、今後フェイクニュースの理解を進める上で重要であると考えられる。

**Q4: そのニュースは対象をどのように扱うか?**

選択肢: 1. 持ち上げ、2. 誹謗、3. どちらでもない

この質問は、そのニュースで扱っている対象について、持ち上げや称賛をしているのか、もしくは誹謗中傷や侮辱をしているのかを選択する。誤ったニュースでも、寄付のような善行か犯罪行為のような悪行かでは読者への印象は大きく異なる。そのため、フェイクニュースが分極化にどのような影響を与えるかを分析する点で重要である。

**Q5: そのニュースの目的は何か?**

選択肢: 1. 風刺/パロディ、2. 党派性、3. プロパガンダ、4. 目的なし/わからない

誤ったニュースの一部は、COVID-19 ワクチンが健康被害を引き起こすなどの自説を広める意図を持って拡散されている。目的を推定できないニュースもあるが、3つの選択肢を設定した。1つ目は読者を楽しませたり、政治を批判したりすることが目的である風刺やパロディである [17]。2つ目は、政治的な文脈で一方的に偏った、党派性をもったニュース。偏っていること自体は誤りではないが、誤りの可能性が高いという報告がある [18, 19]。3つ目はプロパガンダと呼ばれる、イデオロギーや主張を広げたりするために行われる説得の一形態である [20]。

**Q6: 社会に与える有害さはどの程度か?**

この質問は、社会に悪影響を与える可能性のあるニュースを特定するためのもので、0-5の実数値で有害の度合いを回答する。0は社会に害を与えない

パロディなどのニュースを、5は社会に大きな害を与えるニュースであることを示す。

**Q7: そのニュースはどのような害を与えるか?**

選択肢: 1. 害なし (例: 風刺やパロディ), 2. 社会に対する不安, 3. 人, 会社や商品に対する信頼, 4. 政治や社会イベントの正しい理解, 5. 健康, 6. 人種や国に対する偏見, 7. 陰謀論, 8. わからない

誤ったニュースが引き起こす可能性のある被害を7つの選択肢で設定し、さらに「わからない」という選択肢を加えて選択する。この質問は、ニュースによりどのような被害が発生する可能性が高いのかを把握するのに有用である。

## 4 日本語データセット

### 4.1 データセット構築

日本語フェイクニュースデータセットの構築のため、ファクトチェック組織の1つであるファクトチェック・イニシアティブ (FIJ) [21] で2019年7月から2021年10月の間に掲載された検証記事を収集した。アノテーション作業は、ファクトチェック記事と発信源となった投稿・ニュース記事を確認しながら3名のアノテーターによって行われた。アノテーション例を図1に示す。

最終的に307件のニュースにアノテーションが行われた。そのうち、186件のニュースは拡散のきっかけとなった投稿とコンテキスト情報 (Retweet や Reply 情報) を Twitter から収集した。その結果、471,446 ツイート (1 ニュースあたり平均 2,534 ツイート), 277,106 ユーザ (1 ニュースあたり平均 1,489 ユーザ), 17,401 の会話データ (1 ニュースあたり平均 93 会話) が含まれる。ラベルの分布を付録の表1に示す。

### 4.2 日本語データセットの分析

誤りのニュースと正確なニュースの比較を行うのが一般的だが、構築したデータには正確なニュースのサンプルが5例しかない。そのため、本研究では「Q2-1: 発信者はそのニュースが誤りであることを知っていたか?」のはいいいえの回答を元に、つまり Disinformation/Misinformation の観点でツイートデータを分析した。ここでは、特に興味深い結果が得られたユーザアカウントの年齢と Bot 率についてのみ示す。その他の分析は付録Aで示す。

**アカウントの年齢:** 図2にユーザのアカウント作成

Disinformationを投稿したユーザ Misinformationを投稿したユーザ

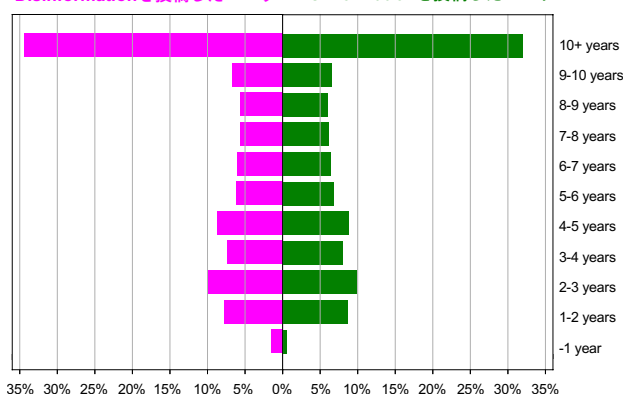


図2: DisinformationとMisinformationを投稿したユーザのアカウント作成日からの経過時間の分布

日からの経過時間の分布を示す。MisinformationとDisinformationの分布は類似しているが、米国においてフェイクニュースを流すユーザ分布の報告 [1] と比較すると、作成から1年未満のアカウントが少なく、10年以上Twitterを利用しているアカウントが多いという傾向がみられた。

**Botの割合:** Botometer API [22] を用いてユーザのBot割合を調査したところ、Disinformationで約8%、Misinformationで約6%であった。米国ではフェイクニュースを流すユーザの約22%がBotであるという報告 [1] と比較すると、かなり少ない比率であるといえる。

## 5 おわりに

本研究では様々な観点からフェイクニュースを捉え直す新しいアノテーションスキームを提案した。事実性だけでなく、意図やターゲット、統一的なラベリングを行うことで、フェイクニュースが引き起こす複雑な現象の深い理解が期待できる。そして、このスキームに基づいて、日本語で初めてのフェイクニュースデータセットを構築した。

一方で、この日本語データセットは根拠としたファクトチェック記事の数が少ないため、サンプル数が小さいという問題がある。この問題を解決するため、今後もデータセットを拡張していく。さらに、我々のアノテーションスキームに基づき日本語以外の言語のフェイクニュースデータセット構築にも現在取り組んでいる。これらの取り組みを通して、1. 既存のフェイクニュース検出モデルは付与したラベルをどの程度分類できるか、2. 言語や国によってフェイクニュースの拡散パターンがどのように異なるのか、といった点に注目できる。

## 謝辞

本研究の一部は、JSPS 科研費 JP19K20279 および厚生労働省科学研究費補助金（課題番号：H30-新興行政-指定-004）の支援を受けたものである。

## 参考文献

- [1] Kai Shu, Deepak Mahudeswaran, Suhang Wang, Dongwon Lee, and Huan Liu. Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media. **Big data**, Vol. 8, No. 3, pp. 171–188, 2020.
- [2] Limeng Cui and Dongwon Lee. Coaid: Covid-19 healthcare misinformation dataset. **arXiv preprint arXiv:2006.00885**, 2020.
- [3] Xinyi Zhou and Reza Zafarani. A survey of fake news: Fundamental theories, detection methods, and opportunities. **ACM Computing Surveys (CSUR)**, Vol. 53, No. 5, pp. 1–40, 2020.
- [4] David MJ Lazer, Matthew A Baum, Yochai Benkler, Adam J Berinsky, Kelly M Greenhill, Filippo Menczer, Miriam J Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, et al. The science of fake news. **Science**, Vol. 359, No. 6380, pp. 1094–1096, 2018.
- [5] Soroush Vosoughi, Deb Roy, and Sinan Aral. The spread of true and false news online. **Science**, Vol. 359, No. 6380, pp. 1146–1151, 2018.
- [6] Karishma Sharma, Feng Qian, He Jiang, Natali Ruchansky, Ming Zhang, and Yan Liu. Combating fake news: A survey on identification and mitigation techniques. **ACM Transactions on Intelligent Systems and Technology (TIST)**, Vol. 10, No. 3, pp. 1–42, 2019.
- [7] Zhiwei Jin, Juan Cao, Yongdong Zhang, and Jiebo Luo. News verification by exploiting conflicting social viewpoints in microblogs. In **Proceedings of the AAAI Conference on Artificial Intelligence**, Vol. 30, 2016.
- [8] Kai Shu, Amy Sliva, Suhang Wang, Jiliang Tang, and Huan Liu. Fake news detection on social media: A data mining perspective. **ACM SIGKDD explorations newsletter**, Vol. 19, No. 1, pp. 22–36, 2017.
- [9] Hunt Allcott and Matthew Gentzkow. Social media and fake news in the 2016 election. **Journal of economic perspectives**, Vol. 31, No. 2, pp. 211–36, 2017.
- [10] Eni Mustafaraj and Panagiotis Takis Metaxas. The fake news spreading plague: was it preventable? In **Proceedings of the 2017 ACM on web science conference**, pp. 235–239, 2017.
- [11] Nadia K Conroy, Victoria L Rubin, and Yimin Chen. Automatic deception detection: Methods for finding fake news. **Proceedings of the association for information science and technology**, Vol. 52, No. 1, pp. 1–4, 2015.
- [12] Martin Potthast, Johannes Kiesel, Kevin Reinartz, Janek Bevendorff, and Benno Stein. A stylometric inquiry into hyperpartisan and fake news. In **Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)**, pp. 231–240, 2018.
- [13] Newsweek. British government bans the phrase ‘fake news’. <https://www.newsweek.com/british-government-bans-phrase-fake-news-1182784>, 2018. [accessed on November 9th, 2021].
- [14] Francesca Giuliani-Hoffman. ‘fake news’ should be replaced by these words, claire wardle says. <https://money.cnn.com/2017/11/03/media/claire-wardle-fake-news-reliable-sources-podcast/index.html>, 2017. [accessed on November 9th, 2021].
- [15] Harvey Leibenstein. Bandwagon, snob, and veblen effects in the theory of consumers’ demand. **The quarterly journal of economics**, Vol. 64, No. 2, pp. 183–207, 1950.
- [16] Firoj Alam, Fahim Dalvi, Shaden Shaar, Nadir Durrani, Hamdy Mubarak, Alex Nikolov, Giovanni Da San Martino, Ahmed Abdelali, Hassan Sajjad, Kareem Darwish, et al. Fighting the covid-19 infodemic in social media: A holistic perspective and a call to arms. In **Proceedings of the International AAAI Conference on Web and Social Media**, Vol. 15, pp. 913–922, 2021.
- [17] John Brummette, Marcia DiStaso, Michail Vafeiadis, and Marcus Messner. Read all about it: The politicization of “fake news” on twitter. **Journalism & Mass Communication Quarterly**, Vol. 95, No. 2, pp. 497–517, 2018.
- [18] Gabriel Hine, Jeremiah Onaolapo, Emiliano De Cristofaro, Nicolas Kourtellis, Ilias Leontiadis, Riginos Samaras, Gianluca Stringhini, and Jeremy Blackburn. Kek, cucks, and god emperor trump: A measurement study of 4chan’s politically incorrect forum and its effects on the web. **Proceedings of the International AAAI Conference on Web and Social Media**, Vol. 11, No. 1, pp. 92–101, 2017.
- [19] Savvas Zannettou, Barry Bradlyn, Emiliano De Cristofaro, Haewoon Kwak, Michael Sirivianos, Gianluca Stringhini, and Jeremy Blackburn. What is gab: A bastion of free speech or an alt-right echo chamber. In **Companion Proceedings of the The Web Conference 2018**, pp. 1007–1014, 2018.
- [20] Garth S Jowett and Victoria O’donnell. **Propaganda & persuasion**. Sage publications, 2018.
- [21] Fact Check Initiative Japan. Fact check initiative japan. <https://fj.info/>. [accessed on November 9th, 2021].
- [22] Clayton Allen Davis, Onur Varol, Emilio Ferrara, Alessandro Flammini, and Filippo Menczer. Botornot: A system to evaluate social bots. In **Proceedings of the 25th international conference companion on world wide web**, pp. 273–274, 2016.
- [23] Vahed Qazvinian, Emily Rosengren, Dragomir R. Radev, and Qiaozhu Mei. Rumor has it: Identifying misinformation in microblogs. In **Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing**, pp. 1589–1599, 2011.
- [24] Amazon Comprehend. Amazon comprehend. <https://docs.aws.amazon.com/comprehend>. [accessed on November 23th, 2021].

