

# タイトルと画像が一致しないニュース記事によるクリックベイトの分析

関 喜史<sup>†</sup>

<sup>†</sup> 株式会社 Gunosy 〒106-6125 東京都港区六本木 6-10-1 六本木ヒルズ森タワー 25 階  
E-mail: †yoshifumi.seki@gunosy.com

あらまし クリックベイトは釣りとも呼ばれ、ニュース記事のタイトルを過剰に表現することでユーザの閲覧行動を誘発する行為のことであり、近年問題になっている。ユーザがクリックベイトなニュース記事を閲覧するとユーザが不快になり、配信したニュース配信サービスの満足度を招く恐れがある。このようにクリックベイトなニュース記事を分析し、特定することは社会的にもビジネス的にも重要な課題ある。本研究では画像を想起するニュース記事におけるクリックベイトを対象に分析を行う。ニュース記事の中には、画像を想起するものがあり、特にエンターテインメント分野では多い。一方で想起した画像が記事にないことがあり、その場合にユーザが不満を持ってしまう恐れがある。クリックベイトなニュース記事に関する研究では、タイトルと本文の一致性が中心に議論されているが、タイトルと画像の不一致もクリックベイトのひとつであり、重要である。クリックベイトなニュース記事を特定するために、本研究ではニュース記事に対する質問を複数の回答者に回答させることによって、ニュース記事のタイトルと画像が一致しているかのデータセットを作成することを目指す。ニュース記事に対して「タイトルから画像があるべきだと思うか」「このタイトルに対して画像は適切だと思うか」という設問を用意し、各ニュース記事に複数の回答者からの回答を集め、回答結果の分析を行った。結果として対象としたニュース記事の中で 35%程度が画像とタイトルが一致しないニュース記事であり、少なくない数のニュース記事がクリックベイトであることが示唆された。

キーワード データ分析、ニュース分析、ジャーナリズム

## 1. はじめに

近年ニュースの閲覧行動は大きく変化している。2016 年 3 月の LINE 株式会社による「世代間のニュースサービス利用に関する意識調査」によれば、84%がニュース閲覧時にスマートフォン利用しており、テレビの 61%を大きく上回った [1]。50 代、60 代は 40%以上が新聞を利用すると回答したものの、10 代、20 代では 20%前半に留まる。このように若者を中心に多くの人々がニュースを読む際にスマートフォンを利用している。

スマートフォンでのニュース閲覧の増加に伴い、スマートフォンでの閲覧に特化した新興のニュース媒体も数多く誕生している。スマートフォンのニュース閲覧では、Yahoo ニュース<sup>(注1)</sup>、グノシー<sup>(注2)</sup>、SmartNews<sup>(注3)</sup>のようなニュースアプリや、Twitter<sup>(注4)</sup>、Facebook<sup>(注5)</sup>、LINE<sup>(注6)</sup>のような SNS(Social Networking Service) がよく用いられている。新興のニュース媒体社はニュースアプリや SNS からの流入を効果的に獲得することで PV(ページビュー) を急激に増加させている。以前は PV が収益に直接つながらないこともあったが、近年では広告技術の進歩により媒体社は PV を高めることで収益を得やすくなった。このような背景から新興のニュース媒体が急成長して

いる。

一方で PV を高めることを目指した媒体運営のために問題も起こっている。アメリカ大統領選やイギリスの EU 離脱問題で注目されたフェイクニュースは代表的な問題である [2]。フェイクニュースを配信する媒体は虚偽の情報を事実であるかのように発信し、SNS で多くシェアされることで PV を獲得することを狙っている。また最近問題になったキュレーションメディアでは安価に大量の PV を獲得するために、専門知識がない編集者が大量にニュース記事を執筆することによって不正確な情報が発信されてしまっていた。本研究で取り上げるクリックベイトもこうした問題の 1 つである。

ユーザの関心を引くために過度に扇情的なタイトルをニュース記事に付けた結果、ユーザがそのニュース記事を読んだ際に騙されたと感じてしまうことを本研究ではクリックベイトと呼ぶ。クリックベイトはフェイクニュースの一部とされている場合もあるが、本研究ではフェイクニュースは「嘘のニュース記事」であり、クリックベイトは「嘘ではないがユーザに誤解を招くようなニュース記事」であると定義する。クリックベイトなニュース記事は釣り記事とも呼ばれている。新聞ではニュース記事のタイトルは見出しと呼ばれており、見出しは本文を十分に要約したものにすることが推奨されていた。しかし新聞は購読した時点で収益が発生するのに対し、ニュースアプリや SNS においてはニュース記事の閲覧を行わないと PV が発生しないことから、タイトルを扇情的にすることでユーザを惹きつけ、PV 数を増やそうという試みを様々な媒体社が行っているが、その試みが行き過ぎてしまうことがクリックベイトを引き

(注1): <https://news.yahoo.co.jp/>

(注2): <https://gunosy.com>

(注3): <https://www.smartnews.com/>

(注4): <https://twitter.com/>

(注5): <https://www.facebook.com/>

(注6): <https://line.me/>

起こしている。クリックベイトなニュース記事閲覧したユーザは媒体社だけでなく、利用していたニュースアプリや SNS に対して不満を持ってしまう。そのためクリックベイトなニュース記事を特定し、ユーザが閲覧しないようにすることは重要である。

しかしクリックベイトは明確な線引をすることが難しい。フェイクニュースの場合は嘘であるため、どのニュース記事がフェイクニュースかは一意に定まる。しかしクリックベイトの場合は、「騙された」という主観的な基準のため判断は難しい。ニュース媒体側も明確に線引をすることができず、悪意がないままにクリックベイトなニュース記事を配信してしまうこともある。このような背景から我々は、ユーザが「騙された」と感じたことをサービス内での行動から特定することを目指している。

クリックベイトなニュース記事を機械学習を用いて特定しようという試みはすでに報告されているが、どのようなニュース記事をクリックベイトなニュース記事とするかという点は十分に議論されていない [3] [4]。我々は以前クリックベイトなニュース記事を特定することを目的に、滞在時間が短いニュース記事を分析しその類型化を行い [5]、滞在時間が短いニュース記事には画像を中心としたニュース記事が多く、その中にはタイトルと画像が一致していないニュース記事が一部存在することを明らかにした。

本研究ではタイトルと画像が一致しないニュース記事によるクリックベイトを画像想起型クリックベイトと定義する。そしてこの問題に対してより深い分析を行うことを目的とし、タイトルと画像の不一致に関するデータセットを人手で作成し、その分析と評価を行った。その結果、タイトルと画像が一致していないことを人手による評価で判断することは簡単ではないこと、そしてタイトルと画像が一致していないクリックベイトなニュース記事は多く存在することが示唆された。

## 2. クリックベイトなニュース記事の特定に関する関連研究

本章ではクリックベイトなニュース記事を特定する試みについて紹介する。

Facebook はクリックベイトの特定に力を入れており 2014 年、2016 年、2017 年に対策を行っていることを明らかにする告知を行っている [6] [7] [8]。2014 年の告知では滞在時間やそのページから戻ってきたときの行動に着目してニュース記事を排除するようにしたこと [6]、2016 年の告知では (1) タイトルがニュース記事を理解するための重要な情報を保留していないか、(2) タイトルが読者の誤解を招くような誇張をしていないかという 2 点の観点から数万のニュース記事を人手で分類し、テキストの特徴からクリックベイトなニュース記事を特定するシステムを構築したことを明らかにしている [7]。そして 2017 年の告知ではクリックベイトの特定をニュース記事に対してではなく個別の投稿に対して行っていること、(1)、(2) を区別して特定するようにしたこと、英語以外の言語にも適用し始めたことを明らかにした [8]。このように Facebook はクリックベイトなニュース記事の特定に力を入れているが、これらの手法は学術

論文として公開されていない。

クリックベイトなニュース記事を特定しようという試みに関する学術論文は、多くはないが存在する。クリックベイトなニュース記事の定義やそのデータセットが存在しないことから、どのようなデータを使うか、どのようにデータを作るかが重要である。

Potthast らの取り組みは機械学習を用いてクリックベイトなニュース記事を特定しようという最初の取り組みであるとされている [11]。Potthast らはよくリツイートされているニュース媒体を 20 個選び、各媒体のリンクを含んだツイートを 150 個ずつサンプリングし、そのツイートのリンクがクリックベイトか否かを 3 人のアノテータによりラベル付を行うことで判定しデータセットを構築し、様々な特徴量を用いて予測を行っている。

Chakraborty らはクリックベイトなニュース記事を多く発信しているとされる 5 つの媒体から 8,069 個のニュース記事を収集し、クリックベイトか否かをラベル付けした [3]。アノテータは 6 人で各ニュース記事は最低 3 人のアノテータに評価されている。これにより 7,623 個のニュース記事をクリックベイトなニュース記事と判定し、そこからランダムにサンプリングした 7,500 記事をクリックベイトなニュース記事、Wikinews のニュース記事からランダムにサンプリングした 7,500 記事をクリックベイトではないニュース記事として、テキストの特徴量を用いて SVM による分類を行った。このデータセットは現在は追加され各ラベル 16,000 記事のデータとなって Github 上で公開されており<sup>(注7)</sup>、このデータを用いた研究も行われている [10] [12]。

また学術論文ではないが Clickbait Classifier というソースコードが Github 上で公開されている<sup>(注8)</sup>。こちらのモデルはインターネット広告へのネガティブ・フィードバックを予測する研究で特徴量として用いられた例が報告されている [13]。これは NewYork Times<sup>(注9)</sup> をクリックベイトではないニュース記事、Buzzfeed<sup>(注10)</sup> と ClickHole<sup>(注11)</sup> をクリックベイトなニュース記事として構築されたニュース記事の分類モデルである。Chakraborty らが構築したデータセットと比較し乱暴なようにも思えるが、Chakraborty らがアノテーションにより排除した指定した媒体の記事が 5%程度であることを考えると、大きな差はないともいえる。

Zheng らは中国の主要なポータルサイトのニュース記事に対して人手でクリックベイトなニュース記事か否かのラベル付を行い、様々な特徴を使ってクリックベイトなニュース記事かどうかを予測するモデルを構築した上で、ユーザがその記事をクリックしたかどうかというユーザ行動を用いてその予測値を修正する方法を提案し予測モデルの精度が向上することを示した。

(注7): <https://github.com/bhargaviparanjape/clickbait>

(注8): <https://github.com/peterldowns/clickbait-classifier>

(注9): <https://www.nytimes.com/>

(注10): <https://www.buzzfeed.com/>

(注11): <http://www.clickhole.com/>

Bourgonje らの試みは Fake News Challenge<sup>(注12)</sup>(FNC1) において 9 位の成績を収めた手法である [4]。FNC1 ではニュース記事の本文とタイトルの関係について *unrelated*, *discuss*, *agree*, *disagree* の 4 つのクラスに分類するというタスクが与えられた。訓練データに *unrelated* とされたニュース記事が 7 割近くあったことから、Bourgonje らはタイトルと本文の類似度を算出し、訓練データから算出したしきい値を下回る場合は *unrelated* とし、下回らなかった場合は構築した 3 クラスの分類モデルを用いてラベルを予測する方法を提案した。これは我々が以前行った分析と類似している [5]。

このようにクリックベイトなニュース記事を特定するという課題については幾つかの取り組みがあるが、何をクリックベイトなニュース記事とするかという点に関して議論が十分ではない。Facebook の取り組みは (1) タイトルがニュース記事を理解するための重要な情報を保留していないか、(2) タイトルが読者の誤解を招くような誇張をしていないかという 2 点に着目しており、どのようなニュース記事をクリックベイトなニュース記事とするかがある程度明確であるが、本章で紹介した研究の多くは、著者のグループがクリックベイトなニュース記事だとラベル付したか否かによって決まっており、その定義は明確ではない。

本研究ではどのようなニュース記事がクリックベイトなニュース記事なのかを、ニュース配信サービスにおけるユーザ行動分析を元に検討していく。

### 3. 本研究におけるクリックベイトの定義

本章では本研究におけるクリックベイトの定義について述べる。2 章で述べたようにクリックベイトがなにかということは明確ではなく、多くの文献ではクリックベイトなニュース記事をタイトルに情報量が不足しているものと定義している。しかし本研究において我々は異なる立場をとる。

クリックベイトが問題になっているのはユーザを騙しているからである。そしてタイトルから情報量を減らすことがユーザを騙しているかどうかは明らかではない。本研究ではユーザが本文を読んだときに騙されたと感じることをクリックベイトと定義し、どのようなニュース記事がクリックベイトなニュース記事になるのかを明らかにし、クリックベイトなニュース記事を特定することを目指している。

我々はクリックベイトなニュース記事を特定することを目的に、滞在時間が短いニュース記事を対象に分析を行った [5]。滞在時間はウェブにおけるコンテンツの評価指標として良く用いられており、Facebook でも滞在時間が短いニュース記事をクリックベイトなニュース記事としている [6]。しかしグノシーにおいて滞在時間が短いニュース記事を表示されにくくした際に、一部のユーザの満足度が低下する結果となった。そこで滞在時間が短いニュース記事が全てクリックベイトなニュース記事ではない説のもとで、滞在時間が短いニュース記事にはどのようなニュース記事があるのかを調査した。その結果、滞在時間が

短く閲覧数が多いニュース記事の 7 割程度が画像を中心としたニュース記事であり、ユーザが画像を見て満足していることが滞在時間が短い要因の 1 つであることが示唆された。しかし画像を中心としたニュース記事の中には、適切な画像が存在しないことがあることも明らかとなった。画像が中心だと考えられるにも関わらず適切な画像が存在しないニュース記事は、ユーザが騙されたと感じることが予想されることからクリックベイトなニュース記事であるといえる。過去の研究ではタイトルと本文が一致しない点については議論されているものの、画像とタイトルの関係については注目されていない。

本研究ではニュース記事のタイトルと画像が一致しているかを判断するデータセットを作成することを目的に、ニュース記事に対する設問を用意し、その回答を複数人で行ったデータを用意する。この回答結果を分析することで、ニュース記事のタイトルと画像が一致しているかを特定することを目指す。

### 4. データセットの構築

本章ではデータセットの構築方法について述べる。本研究ではスマートフォン向けアプリケーションであるグノシーのニュース記事を利用し、ニュース記事に対してそのニュース記事が画像を想起するか否か、ニュース記事内の画像が想起した画像と一致するかという 2 点についてラベル付を行う。以前の調査では著者が 1 名で少数のニュース記事に対するラベル付けを行い、ニュース記事のタイトルが画像を想起するか、実際の画像とニュース記事のタイトルが一致しているかを調査した [5]。今回は対象とするニュース記事の量を増やすことに加えて、複数のアノテータによってニュース記事の評価を行い、より大規模で、正確なデータセットの構築を目指す。

まずラベル付の対象となるニュース記事を選択する。本研究では、ユーザからの閲覧数がある程度多く、滞在時間の短く、エンタメカテゴリ (芸能など) であるニュース記事をラベル付の対象とする。クリックベイトなニュース記事は閲覧を誘発しているので、閲覧数が少ないニュース記事はクリックベイトではない。滞在時間が短いニュース記事に絞るのはユーザが画像を見て満足するようなニュース記事を抽出するためである。そしてエンタメカテゴリに絞る理由は、前回の調査で画像を想起する割合が高いカテゴリであったためである。閲覧数と滞在時間は 2017 年 10 月 14 日 0 時から 2017 年 12 月 8 日 24 時までの 8 週間のニュース記事の閲覧データを対象として集計する。滞在時間は期間中に閲覧した全ユーザの滞在時間の中央値である。1 人のユーザが同じニュース記事を数回閲覧した場合は、最大値を代表値とする。これによってラベル付の対象記事を 1,560 件選択した。閲覧数と滞在時間のしきい値は事業上の理由から明らかにすることはできないが、滞在時間のしきい値はニュース記事の長さとの相関が小さく、直帰したと考えられる値に設定しており、閲覧数のしきい値によって期間中の滞在時間が短いエンタメカテゴリのニュース記事の中で閲覧数上位 5% 以内のニュース記事が選択されている。

ラベル付は 12 人のボランティアによって行われた。全員ニュース配信サービスの開発に関わる男性であり、本研究の目

(注12): <http://www.fakenewschallenge.org>

表 1 Q1 で 3 人が一致した回答の件数

Q1 の答え	一致した件数
画像があるべきだと思う	1, 124
どちらともいえない	1
画像がなくてもよい	14
合計	1, 139

表 2 Q2 で 3 人が一致した回答の件数

Q2 の答え	一致した件数
適切だと思う	915
どちらともいえない	11
適切ではないと思う	73
合計	999

表 3 Q1, Q2 共に 3 人が一致した回答の件数

Q1 の答え	Q2 の答え	一致した件数
画像があるべきだと思う	適切だと思う	722
画像があるべきだと思う	どちらともいえない	3
画像があるべきだと思う	適切ではないと思う	63
画像がなくてもいいと思う	適切だと思う	7
合計		795

的について十分に説明を受けている。そのため回答の傾向には偏りがある可能性がある。各ニュース記事に 3 名の回答が集まるようにランダムにニュース記事を振り分け、1 人あたり 380 記事 ~ 400 記事のラベル付を行った。

設問はニュース記事に対して以下の 2 問があり、肯定的な回答、否定的な回答、判別不能の 3 つの選択肢から回答を選択する。

• Q1. このニュース記事のタイトルから、記事には画像があるべきだと思いますか？

- 画像があるべきだと思う
- どちらともいえない
- 画像がなくてもよい

• Q2. 下記の記事タイトルに対して、この画像は適切だと思いますか？画像がない場合は、「画像がありません」ということが適切かどうかを回答してください

- 適切だと思う
- どちらともいえない
- 適切ではないと思う

Q1 ではニュース記事のタイトルが表示され、画像は表示されていない。Q2 ではニュース記事のタイトルと画像が表示されており、画像がない場合は「画像がありません」と表示される。Q1 と Q2 はページが分かれており、Q1 の時点で画像を見ることはできない。

以上のようにしてニュース記事 1, 560 件に対するラベル付きデータセットの作成を行った。

## 5. データセットの分析

本章では作成したデータセットの分析を行う。まず 3 人の回答が一致したものについて述べる。

表 1 に Q1 で 3 人の回答が一致したものについて示す。Q1

表 4 どちらともいえないを 1 つまで含んで一致した Q1, Q2 の回答

Q1 の答え	Q2 の答え	一致した件数
画像があるべきだと思う	適切だと思う	895
画像があるべきだと思う	どちらともいえない	6
画像があるべきだと思う	適切ではないと思う	88
画像がなくてもいいと思う	適切だと思う	24
画像がなくてもいいと思う	適切ではないと思う	2
合計		1, 015

では 1, 139 件が一致しており、全ニュース記事の 73.1%である。一致しているもののうち、1, 124 件が「画像があるべきだと思う」と回答されており、一致している回答のうち 98.6%となり一致したものの殆どが画像があるべきだという回答であった。また全ニュース記事でも 72.1%を占め、今回対象となったエンタメカテゴリの中の滞在時間の短いニュース記事の多くは画像があるべきニュース記事であるといえる。

表 2 に Q2 で 3 人の回答が一致したものについて示す。Q2 では 999 件が一致しており、全ニュース記事の 64.1%である。このうち「適切だと思う」が 915 件で一致したものの 91.6%であり、Q2 において回答が一致したものの多くは画像が適切であったといえる。

表 3 には Q1, Q2 両方で 3 人の回答が一致したものを示す。一致したものは 795 件で全ニュース記事の 51.0%である。このうち 722 件が「画像があるべき」であり「画像が適切である」という回答である。また Q1 で「画像がなくてもよい」とされたもので回答が一致したものはすべて「適切である」となった。「画像があるべき」であり「画像が適切ではない」において一致した回答は 63 件である。Q2 の「画像が適切ではない」で 3 人の回答が一致したものは 73 件であり、そのほとんどが Q1 で「画像があるべき」で一致したものであったことがわかる。

まず表 4 に 3 人が一致したものに加えて、「どちらともいえない」を 1 つ含んで残り 2 人の回答が一致したものも含めた際の回答が一致した件数を示す。回答者 3 人のうち、「どちらともいえない」が 1 人で残り 2 人が一致している場合には、3 人が一致するほどではないものの、ある程度合意が取れているものとみなすことができる。全体の件数は増えたものの、各選択肢の比率は大きく変わらない。そのため回答者の合意がある程度取れるものに関しては、全体の 95%程度が画像を想起するニュース記事であり、そのうち 90%程度が画像が適切であり、10%程度が画像が適切でないといえる。

合意が取れるものについて大半がタイトルと画像が一致しているものであった。一方で合意がとれたものは全ニュース記事の 2/3 程度であり、1/3 は合意が取れていない。これは多くのニュース記事が適切な画像を利用していると考えられるが、一方で画像が適切ではないという判断をすることは難しい、もしくは個人によって判断の差があるとも考えられる。そこで次に合意がとれなかったものがどのように回答されていたのかを分析していく。

合意がとれなかったものを分析するにあたり、各選択肢を選んだ人数による回答の分布をみる。

表 5 Q1 で各選択肢を選んだ人数別の記事数

Q1 の答え	1 人が選択	2 人が選択	3 人が選択
画像があるべきだと思う	140	244	1, 124
どちらともいえない	201	46	1
画像がなくてもいいと思う	213	63	14

表 6 Q2 で各選択肢を選んだ人数別の記事数

Q2 の答え	1 人が選択	2 人が選択	3 人が選択
適切だと思う	225	290	995
どちらともいえない	303	71	11
適切でないと思う	218	106	73

表 5 に Q1 の各選択肢を選んだ人数別の記事数を示す。Q1 における「画像があるべきだと思う」という選択肢はほとんどのニュース記事で 3 人に選択され一致しており、選択する人数が少なくなるほど、ニュース記事の数は少なくなっている。一方でこれまでの結果でも明らかであるように「どちらともいえない」、「画像がなくてもいいと思う」という選択肢が 3 人全員に選ばれたニュース記事はほとんどなく、選択する人数が少なくなるほど、ニュース記事の数は多くなっている。

表 6 に Q2 の各選択肢を選んだ人数別の記事数を示す。大まかには Q1 と同じような傾向にある。Q2 における「適切だと思う」という選択肢はほとんどのニュース記事で 3 人に選択され一致しており、選択する人数が少なくなるほど、ニュース記事の数は少なくなっている。一方で「どちらともいえない」、「適切ではないと思う」という選択肢が 3 人全員に選ばれたニュース記事はほとんどなく、選択する人数が少なくなるほどニュース記事の数は多くなっている。

このことから Q1 の「画像があるべきだと思う」、Q2 の「適切だと思う」という回答は選択されやすいのに対して、Q1 の「画像がなくてもいいと思う」、Q2 の「適切ではないと思う」という回答は選択されにくいのではないかと見える。今回対象としているニュース記事はグノシーと契約しているニュース媒体が配信しているニュース記事であるため、明らかに不適切と判断できるようなニュース記事が少ないことから適切ではないという選択肢を選びにくいと考えられる。また多くのニュース記事の回答が「画像があるべきだと思う」、「適切だと思う」となっているため回答者が機械的に判断してしまうこともその要因として考えられる。

実際に Q2 で 1 名、2 名のみが「適切ではないと思う」を選んだニュース記事を見てみると、部分的には適切であるともいえるが、ユーザが期待していたような画像ではないとも考えられるニュース記事が多くあった。例えばグラビアアイドルの写真集が発売したことを知らせるニュース記事において、タイトルではその写真集の写真のシーンについて詳しく述べられているが、表示されるのはグラビアアイドルがその写真集を持っている画像というニュース記事や、そのシーンではない画像が用いられているようなニュース記事が多くあった。このようなニュース記事に「適切だと思う」と回答した回答者に理由を聞くと「不適切と断言できるものではないと思った」という

回答が得られた。本研究で特定したいクリックベイトなニュース記事はこのようなニュース記事を含む。そこで Q2 で 1 名でも「適切ではないと思う」が選択されたニュース記事をタイトルと画像が一致しないニュース記事だとして、クリックベイトなニュース記事と考えることとする。その場合 397 件がタイトルと画像が一致しないニュース記事であった。これらすべてが Q1 で「画像があるべきだと思う」で 3 人の回答が一致しているニュース記事であり、画像があるべきニュース記事のなかで、適切な画像がなかったといえるニュース記事は 35.3%であった。

今回対象としたニュース記事はすべてグノシーで配信されているニュース記事であり、ニュース記事の配信元は一定の基準によって審査をうけている。このような状況の中で 35.3%という数値は著しく高いものであると考えられ、クリックベイトがニュース媒体の中で多く発生していることを示唆するものであるといえる。

## 6. おわりに

本研究ではクリックベイトなニュース記事を特定するという目的の元で、タイトルと画像が一致していないニュース記事をクリックベイトなニュース記事と考え、ニュース記事のタイトルと画像が一致しているかどうかを手でラベル付することを試みた。過去クリックベイトなニュース記事を特定しようという研究では、クリックベイトなニュース記事がどのようなニュース記事かという議論が十分に行われていない。そこで我々はクリックベイトなニュース記事をユーザが本文を読んで騙されたと感じるようなニュース記事であると定義し、その中でもタイトルから期待した画像が本文中に存在しない、タイトルと画像が一致していないニュース記事を対象として研究を行った。

まずタイトルと画像が一致していないニュース記事を特定するために、今回は人手でのラベル付を行った。過去の試みからニュース記事中の滞在時間が短いニュース記事は、画像を期待するようなニュース記事であることが多いことが明らかになっている。そこで閲覧数が多く、滞在時間の短い、画像が中心なニュース記事が多いと予想されるエンタメカテゴリのニュース記事を対象にした。これらの記事に対して人手でラベル付を行い、タイトルと画像が一致していないニュース記事の特定を目指した。ラベル付のための質問はそれぞれのニュース記事につき 3 名の回答者がニュース記事のタイトルを見て「このニュース記事に画像があるべきだと思うか」を回答した上で、ニュース記事のタイトルと画像を見て「このニュース記事タイトルにこの画像は適切だと思うか」にそれぞれ肯定的・否定的・判断不能の 3 段階で回答した。

回答を集計して 2 つの知見が得られた。第一に肯定的な選択肢は回答者で合意が取りやすく、否定的な選択肢は合意が取りにくいことである。「画像があるべきだと思う」「適切だと思う」という肯定的な選択肢が 3 人とも一致するニュース記事が多いのに対して、「画像がなくてもいいと思う」「適切ではないと思う」という否定的な選択肢は 3 人が一致するニュース記事は殆どなかった。しかし「このニュース記事タイトルにこの画像は

適切だと思うか」の設問に 1 名でも否定的な選択肢が選ばれたニュース記事は、合意が取れているニュース記事ほど明らかに不適切ではないが、タイトルに対して画像が十分に適切とは言いがたく、ユーザに不快感を与える可能性があるニュース記事であった。第二にタイトルと画像が一致していないニュース記事がある程度存在することが明らかになったことである。第一の知見を背景に一人でも否定的な回答が得られたニュース記事をタイトルと画像が一致していないニュース記事と判定したところ、対象となったニュース記事の 35.3% がタイトルと画像が一致していないニュース記事と判定された。本研究の対象記事がグノシーが契約を結んでいるニュース媒体が配信したものであると考えると、この比率は決して少ないものではないといえる。

今後の課題としてデータセットの作成に関しては、設問の内容を改善し、問題のあるニュース記事を容易に見つけるようにすること、より多様な回答者によるラベル付を行うことがあげられる。また構築したデータセットを用いて、ユーザ行動と組み合わせた分析や、分類モデルの構築なども行っていきたい。

#### 文 献

- [1] LINE 株式会社, 世代間のニュースサービス利用に関する意識調査, <https://linecorp.com/ja/pr/news/ja/2016/1267>, 2016.
- [2] 藤代 裕之, ネットメディア覇権戦争 偽ニュースはなぜ生まれたか, 光文社, 2017.
- [3] Abhijnan Chakraborty, Bhargavi Paranjape, Sourya Kakarla, Niloy Ganguly, Stop Clickbait: Detecting and preventing clickbaits in online news media, ASONAM, 2016.
- [4] Peter Bourgonje, Julian Moreno Schneider, Georg Rehm, From Clickbait to Fake News Detection: An Approach based on Detecting the Stance of Headlines to Articles, EMNLP Workshop, 2017.
- [5] 関 喜史, 潮 旭, 米田 武, 松尾 豊, クリックベイトなニュース記事の特定に向けたユーザ行動分析, 信学技報, vol. 117, no. 207, NLC2017-27, pp. 65-70, 2017 年 9 月.
- [6] Facebook Newsroom, News Feed FYI: Click-baiting. <https://newsroom.fb.com/news/2014/08/news-feed-fyi-click-baiting/>, 2014.
- [7] Facebook Newsroom, News Feed FYI: Further Reducing Clickbait in Feed, <https://newsroom.fb.com/news/2016/08/news-feed-fyi-further-reducing-clickbait-in-feed/>, 2016.
- [8] Facebook Newsroom, News Feed FYI: New Updates to Reduce Clickbait Headlines, <https://newsroom.fb.com/news/2017/05/news-feed-fyi-new-updates-to-reduce-clickbait-headlines/>, 2017.
- [9] Hai-Tao Zheng, Xin Yao, Yong Jiang, Shu-Tao Xia, Xi Xiao, Boost Clickbait Detection Based on User Behavior Analysis, APWeb-WAIM, 2017.
- [10] Ankesh Anand, Tanmoy Chakraborty, Noseong Park, We Used Neural Networks to Detect Clickbaits: You Won't Believe What Happened Next!, Advances in Information Retrieval. ECIR 2017.
- [11] Martin Potthast, Sebastian Kpsel, Benno Stein, Matthias Hagen, Clickbait Detection, Advances in Information Retrieval. ECIR 2016.
- [12] Md Main Uddin Rony, Naeemul Hassan, Mohammad Yousuf. Diving Deep into Clickbaits: Who Use Them to What Extents in Which Topics with What Effects?. ASONAM, 2017.
- [13] Zhou Ke, Redi Miriam, Haines Andrew, Lalmas Mou-

nia, Predicting Pre-click Quality for Native Advertisements, WWW'16, 2016.