

ネタバレ確信度を考慮した 試合実況データセット構築と分析手法の検討

白鳥 裕士^{†1, a} 牧 良樹^{†1, b} 阿部 和樹^{†1, c} 中村 聡史^{†2, d}

†1 明治大学大学院先端数理科学研究科 †2 明治大学総合数理学部

a) swany181@gmail.com b) tekkanomaki01@gmail.com

c) cs182001@meiji.ac.jp d) satoshi@snakamura.org

概要 スポーツの録画視聴を楽しみにしている人にとって、Twitter 上などで意図せず遭遇してしまうネタバレ情報は問題であり、自動的にネタバレを検出・遮断できるようなシステムを構築することが重要である。本研究では、そうしたシステムの構築のための足がかりとして、ネタバレツイートを「試合の最終結果が高い確信度で予測できてしまう投稿」と定義し、Twitter 上でのサッカーのネタバレに関するデータセットを構築した。また、ネタバレの文章特性の調査を行い、試合状況とネタバレ内容の連動性を確認した。さらに、構築したデータセットを用いてネタバレ判定実験を行い、3つの判定手法を比較した結果、SVM+試合展開手法が優れていることが明らかになった。

キーワード ネタバレ防止, 機械学習, スポーツ, サッカー, SNS, Twitter

1 はじめに

スポーツは筋書きのないドラマであるため、どちらが勝つか分からないという緊張感や予想もしない試合展開に対する驚きを味わうことができる。そのため、リアルタイムで観戦したいと考えている視聴者は少なくない。しかし、仕事や学業などの時間の関係で、スポーツの試合をリアルタイムで視聴観戦することが困難な場合があり、あらかじめ録画予約をしておき、時間に余裕があるときに改めて視聴するというのも珍しくない。ここで、録画視聴を楽しみにしている人が、視聴前にそのスポーツの試合結果を知ってしまうと、緊張感や驚きが失われてしまう可能性がある。こうした緊張感や驚きを大事にしている視聴者にとって、試合のスコアや勝敗といった「ネタバレ情報」は避けたいものであるため、視聴するまでの間、情報遮断を積極的に行っている。しかし Twitter や Facebook などの SNS はふとしたコミュニケーションをとるために気軽にアクセスすることが多く、その際に「本田ごおおおおおおる！」のようなネタバレ情報を目にしてしまうことも少なくない。

こうしたネタバレ情報を遮断することを目的とした研究は現在盛んに行われている。例えば、中村らは、ネタバレ情報を動的にフィルタリングすることを目的とし、ユーザの意図に基づいてウェブページに含まれるユーザの興味のある情報をフィルタリングする手法[1]や、ユーザとネタバレとの関わり方に注目し、情報の提示を工夫することでネタバレを回避する手法[2]について提案している。しかし、これまで行われてきた研究では、ユーザと

ネタバレとの関わり方といったインタラクションに注目しており、ネタバレはどのような文章特性を持っているのか、また、高精度に判定するにはどうしたら良いのかといった点については取り組まれていなかった。また、我々はこれまで、スポーツにおけるネタバレデータセットを構築し、ネタバレ実験判定を行うことで、試合状況別に SVM モデルを切り替える手法の有用性を明らかにしてきた[3]が、データセットの構築について、ツイート (Twitter 上での投稿) のネタバレの基準を明確にしていなかったために精度に問題があり、また分析や判定が複雑になりすぎている。

そこで本研究では、サッカーの試合に関する Twitter 上でのネタバレツイートについて定義を行い、データセットを再構築する。また、スポーツのネタバレの文章特性について分析し、それらを高精度に判定するための手法についての検討も行う。

本研究の貢献は以下の2点である。

1. これまで構築したデータセットの問題点から、「試合の最終結果が高い確信度で予測できてしまう投稿」をネタバレツイートとして明確化し、サッカーに関する Twitter 上の投稿についてのデータセットを再構築した。
2. ネタバレ判定実験を行い、4つの手法を比較することで、SVM+試合展開手法の有用性を明らかにした。また、確信度帯別の判定結果の分析や前処理方法の検討も行った。

2 関連研究

2.1 コンテンツ体験へのネタバレの影響

ユーザのコンテンツ体験へのネタバレの影響調査については、Leavittらが小説に着目し、ネタバレ情報の提示の有無によって、ユーザのコンテンツの楽しみ方などのような差があるのかを実験により調査している[4]。実験の結果、ネタバレ情報はコンテンツの面白さを落とさないと主張しているが、最初にネタバレされても名作であれば最後まで読書すると面白いという結果を示しているだけであった。また、どちらかといえばあらすじ提示によって、内容や人物関係の理解を支援し、結果的に小説が読みやすくなるということを示唆するものであった。これについては、Rosenbaumらが、小説を読み慣れていない人はネタバレをされた方がストーリーを面白いと感じ、読み慣れている人はネタバレをされない方がストーリーを面白いと感じるということを示している[5]。こうした研究から、ネタバレを防止していく必要性は十分にあるといえる。本研究では、スポーツにおいて、ネタバレを自動的にフィルタリングできるようなシステムを構築するための足がかりとして、スポーツのネタバレに関するデータセットを構築し、高い精度でネタバレを判定できるような手法の検討を行うものである。

2.2 ネットバレルの防止

TwitterのようなSNSでのネタバレを問題視した研究として、田島らはテレビアニメのようなストーリーコンテンツにおいて、放送時間差によって、SNS上でネタバレをされてしまうことを問題としており、致命的なネタバレとなる「生死」「勝敗」などのトピックに対し、人物名一般化などの事前処理を施してから機械学習を行うことでネタバレを判定する可能性について明らかにしている[6]。こうした研究に対し、スポーツのネタバレは試合結果に関するものが多く、これらの研究で扱われているネタバレとは内容が異なる。

同じくSNS上でのネタバレを問題視した研究として、Golbeckは放送時差によりTwitter上でドラマやスポーツに関するネタバレがされてしまうことを問題としており、視聴対象のコンテンツに関するワードリストを生成することにより、関連するツイートをミュート可能としている[7]。これに対して、我々は対象のコンテンツに関するツイートをすべて遮断するのではなく、ネタバレツイートのみを高精度に判定する手法を検討している。また、JeonらはTwitter上でのコメントに対して、「固有表現」や「頻繁に使用される動詞」「時制」などに注目した機械学習を用いてネタバレ検出をする手法を提案している[8]。同研究では、テレビ番組に関するコメントを用いて実験を行うことで、これまでのキーワードマッチングやLDA(Latent Dirichlet Allocation)を用いた手法に比べて高い適合率でネタバレを検知することを可能とし、有用性を示して

いる。また、彼らはスポーツのネタバレについても判定実験を行っている。しかし、スポーツについては1試合しか実験を行っていないうえ、ラベル付けを著者自身で行っているため、ネタバレの定義に疑問が残る。さらに、未来時制に注目した手法を用いているが、未来時制が存在しない日本語には対応できない。

本研究では、スポーツの試合に対するツイートについて実際にラベル付けを行ってもらうことで、ネタバレデータセットを構築し、日本語のネタバレを高い精度で判定できるような手法についての検討を行うものである。

3 ネットバレルデータセットの構築

本章では、ネタバレを防止していくための足がかりとして、どういった情報がネタバレとなるのか、ネタバレの特性を分析する。ここで、視聴者がネタバレにうっかり遭遇する媒体として多いTwitterのようなSNSに注目する。Twitterでは、アクセスするだけで友人の現在の状況を知り、気軽にコミュニケーションを取ることができるため、何気なくアクセスするユーザが多く、その際にネタバレ情報も目にしてしまう可能性が高い。そこで、本章ではスポーツの試合に対するTwitter上の投稿(ツイート)を収集し、ネタバレデータセットを構築することで、ネタバレの特性を分析していく。

データセットの構築にあたって、まずスポーツの試合に対するツイートを用意する必要がある。これについては、以前の研究[3]で我々が収集したツイートをを用いることにした。用意した試合は、日本代表の試合であり、FIFA女子ワールドカップカナダ2015、EAFF東アジアカップ2015、FIFAワールドカップロシア2018、国際親善試合の試合を含む9試合とした。なお、9試合の試合結果は日本代表の5勝1敗3分であった。

我々のこれまでの研究では、ツイートがネタバレであるのか否かを決定する際、「このツイートはネタバレですか?」と、直接問うようなシステムを開発していた。しかし、試合についてほんの些細な情報でもネタバレであると感じてしまう人や、全くネタバレであると感じない人など、人による差が大きく分析が複雑になりすぎてしまった。そこで本研究では、ネタバレツイートを「試合の最終結果が高い確信度で予測できてしまう投稿」と定義し、試合の結果をどの程度の確信度で予測できるかを答えてもらうものとした。ユーザは実験システムを用いて、試合の最終結果が「勝ち」「負け」「引き分け」のどれになるかをページ上に提示されているツイートから予測して対応するボタンを選択したあと、どの程度の確信度で予測できたかについてゲージ(0~100)を移動して値を選択することで、ツイートに対してラベル付けをすることが可能となっている。「勝ち」「負け」「引き分け」のどれにも予測ができない場合(確信度が0の場合)に対応するため、

「勝ち」「負け」「引き分け」に加え、「わからない」というボタンも配置した。

また、ユーザは試合開始からの経過時間との組み合わせで、ある程度のことを予測できてしまう。例えば、「守備固めてもしょうがない、攻めていこう」というツイートが試合開始時点であれば、単なる意気込みと捉えるかもしれないが、試合後半であれば、試合に負けているという状況を考慮した上でのコメントと捉えることも多いであろう。そこで、ツイートの下にそのツイートが投稿された時点での試合開始からの経過時間を表示した(図 1)。



図1 開発したウェブシステム

ここで、ツイートから試合の最終結果が予測できてしまう要素として、独立したツイートや経過時間以外にも、ツイートの組み合わせなども考えられる。しかし、ウェブシステムでの提示方法や予測基準が複雑かつ困難になることを考慮し、本研究では独立したツイートおよび経過時間をネタバレ判断基準とした。また、収集したツイート全てに対して予測してもらうとなると膨大な時間がかかってしまうため、提示するツイートは1試合につき任意の1000件とした。ここで、1試合ずつ提示してしまうと、その試合を視聴したことがある人が予測する場合、1つのツイートから試合の前後関係や内容を思い出し、予測に影響を及ぼしてしまう可能性がある。また、前述したツイートの組み合わせが予測に影響を及ぼしてしまう可能性がある。そこで、9試合分のツイートをランダムに提示することとした。

なお、システムでは9000件のツイートからまだそのユーザによってラベル付けされていない任意の100件のツイートを提示し、ラベル付けを行うものとした。提示可能なツイートが100件に満たない場合には、提示可能な件数分のツイートが提示されるようにした。

このウェブシステムを用いて、実際にツイートに対してラベル付けを行い、データセットを構築した。データセットの構築にあたり、サッカーの試合観戦に興味があり、Twitterを普段から用いている19歳から22歳の大学生17人に協力を依頼した。データセット構築の結果、1ツイートあたり5人以上のラベルが付き、合計45174件のデータを収集することができた。

4 データセットの分析

本章では、前章で構築したデータセットの内容について分析する。まず、構築したデータセットについて、平均確信度別のツイート例を表1に示す。

表1 平均確信度別ツイート例

平均確信度	ツイート	経過時間
0~9	おいしいー	41
	柏木うめえ	31
10~19	長かった。ここから。	29
	あと3分	112
20~29	香川先制点!	9
	森重ごーる	34
30~39	1-0で前半終了。	48
	あと2点ぐらい入ったなあ	51
40~49	アメリカ3点目	14
	ロスタイムに得点されて負けるパターンだろコレ(ーー;))	112
50~59	前半終了 シンガポール 0-2 日本	46
	よかった。ほんとよかった。香川とか岡崎とか原口とか酒井とか。	116
60~69	もうゴールしても喜ばなくなった。	78
	【後半 35分経過】 日本 1 × 1 イラン	98
70~79	勝った勝った	112
	引き分けか…。	114
80~89	5点目きたああああああ	77
	イラン 1-1 日本	113
90~100	試合終了 1-1 引き分け	113
	3-0か～。あと2点くらい欲しかったな。	120

全体として、時間経過とともに平均確信度が高くなっている。これは、その時点での試合状況を述べるツイートが多く、試合がそれから動く可能性が時間経過とともに低くなるためだと考えられる。経過時間以外にも「前半終了 シンガポール 0-2 日本」「アメリカ3点目」といったツイートのように、対戦国が明確にされていたり、点差が大きいことがわかったりするものは確信度が高くなっていた。これは、見る人のサッカーへの精通度合いや性格によっては試合結果を悟ってしまうためだと考えられる。平均確信度が50を超えると「よかった。ほんとよかった。香川とか岡崎とか原口とか酒井とか。」といったツイートのように、途中経過ではなく、試合結果が汲み取れてしまうものも見られた。また、平均確信度が70を超えると、直接試合結果を述べているツイートも多く見られ、さらに、平均確信度が80を超えるとスコア情報に関する

ツイートが多く見られた。試合結果について言及しているにも関わらず確信度がそこまで高くないツイートも多かったが、これは、経過時間が試合終了したのかどちらともいえない時間であり、まだ試合が動く可能性があると考えていたためだと推測される。また、人によってはツイート内容を鵜呑みにしない人がいるためだと考えられる。

平均確信度別の分析から、平均確信度が 50 以上のツイートは、その時点での試合状況だけではなく、試合結果も汲み取れてしまうものが多く、致命的なネタバレになり得ると考えた。実際、平均確信度が 50 以上のツイートには「ゴール」「勝つ」「試合終了」など、特定の用語が含まれているものが多かった。また、日本代表が勝っている時(最終的に勝った時含む)は「キター」「嬉しい」、負けている時(最終的に負けた時含む)は「悔しい」「最悪だ」など、試合状況によって内容が異なっていた。

表 2 時間帯別頻出単語

勝ち		負け		同点	
単語	TF-IDF	単語	TF-IDF	単語	TF-IDF
チーム	0.706	チーム	0.805	チーム	0.831
選手	0.596	選手	0.302	選手	0.234
勝つ	0.129	準優勝	0.206	引き分け	0.192
勝利	0.126	する	0.159	試合終了	0.192
する	0.124	おめでとう	0.159	最下位	0.146
前半	0.086	お疲れ様	0.124	勝てる	0.123
後半	0.082	決勝	0.113	draw	0.100
ゴール	0.078	失点	0.096	お疲れ様	0.100
代表	0.071	w 杯	0.080	後半	0.091
試合	0.069	ある	0.075	試合	0.091

そこで、日本代表が勝っている時間帯、負けている時間帯、同点である時間帯別に、平均確信度が 50 未満のツイートと比較して平均確信度が 50 以上のツイートに出現する頻度(TF-IDF[9])が高かった単語上位 10 件を表 2 に示す。単語分割には MeCab を用い、意味が等しい単語の形を揃えるために KAKASHI を用いて半角文字を全角文字に変換した。さらに、選手名やチーム名など意味が等しい固有名詞の試合による分散を抑えるため、選手名を「選手」、チーム名を「チーム」、監督名を「監督」に正規化する処理を行った。なお、意味をもたない単語の出現を減らすため、名詞、動詞、形容詞、副詞、感動詞のみを対象とした。

表 2 より、チーム名と選手名以外は全ての時間帯に共通する単語は存在しなかった。時間帯別で見ると、勝っている時間帯では「勝つ」「勝利」のような試合結果を直接表す単語、「前半」「後半」のような試合の経過を表す単語、「ゴール」のような試合状況の変化を表す単語が頻出していた。負けている時間帯では、「おめでとう」

「お疲れ様」のような選手への労いを表す単語や「失点」のような試合状況の変化を表す単語が頻出していた。同点の時間帯では「引き分け」のような試合結果を直接表す単語、「後半」「試合終了」のような試合の経過を表す単語、「お疲れ様」のような選手への労いを表す単語が頻出していた。

このように、試合結果を直接表す単語でも、勝っている時間帯では「勝つ」、同点の時間帯では「引き分け」といったように異なっていたり、勝っている時間帯以外では選手への労いを表す単語が頻出していたりと、時間帯によって特徴的な単語が異なることが示唆された。

5 ネタバレ判定実験

5.1 実験手順

本章では、構築したデータセットを用いて、ネタバレを高精度に判定するための手法についての検討を行った。アルゴリズムについては、我々がこれまで検討してきたパターンマッチ手法、SVM 手法、SVM+試合状況手法[3]に加えて、ランダムフォレスト手法も含めた 4 つの手法を比較した。なお、前章で分析したように、スポーツのネタバレは単語に大きく特徴が表れていたため、本研究では単語を機械学習の特徴量とした。

また、単語分割や半角文字の全角変換、品詞の選定は、前章でデータセット内容を分析した時と同様に行った。ここで、日本語には漢字やひらがな、カタカナといったように、同音同義語を表現する方法が複数あるため、単語をローマ字に変換し、同音同義語を同じ特徴量として扱えば精度が高くなる可能性がある。選手名など試合ごとに異なる固有名詞を正規化する方法についても、4 章で行った単純な正規化の他に、敵味方を分けた正規化や、チーム別に分けた正規化、ポジション別に分けた正規化が考えられる。また、正規化を行わないことで精度が高くなる可能性もある。そこで、ローマ字変換(変換あり、なし)と正規化(単純正規化、敵味方別正規化、チーム別正規化、ポジション別正規化、正規化なし)については、すべての前処理パターン(2 × 5 = 10パターン)で実験を行った。

さらに、どの程度の確信度のネタバレツイートであれば高精度に判定できるのかについても調査を行うため、平均確信度が 50 以上のツイートをネタバレ、49 以下のツイートを非ネタバレとした場合(以降 50-49)、60 以上のツイートをネタバレ、40 以下のツイートを非ネタバレとした場合(以降 60-40)、70 以上のツイートをネタバレ、30 以下のツイートを非ネタバレとした場合(以降 70-30)、80 以上のツイートをネタバレ、20 以下のツイートを非ネタバレとした場合(以降 80-20)、90 以上のツイートをネタバレ、10 以下のツイートを非ネタバレとした場合(以降 90-10)のそれぞれについて判定実験を行った。なお、

いずれの場合においてもネタバレツイートの数が多くなってしまったため、アンダーサンプリングを行ってデータ量を調整した。その結果、50-49 から順に、データ量は 1612 件、1254 件、878 件、616 件、362 件となった。

●パターンマッチ手法:

ネタバレとして出現頻度の高い単語をキーワードとし、キーワードにマッチする単語を含むツイートをネタバレと判定した。ここでは、TF-IDF 値の上位 120 単語をキーワードとした。

●ランダムフォレスト手法:

それぞれのツイートについて、BoW (Bag-of-Words) [10]を特徴量として、学習および判定を行った。また、ハイパーパラメータについてはグリッドサーチ (学習率, 割引率ともに[0.0001, 0.001, 0.01, 0.1, 1.0, 10.0, 100.0, 1000.0]の範囲) を行って設定した。なお、データ量の 8 割を訓練データ, 2 割をテストデータとした。

●SVM 手法:

それぞれのツイートについて、BoW (Bag-of-Words) [10]を特徴量として、学習および判定を行った。また、ハイパーパラメータについてはグリッドサーチ (決定木の数を[50, 100, 200, 300, 400, 500]の範囲) を行って設定した。なお、データ量の 8 割を訓練データ, 2 割をテストデータとした。

●SVM+試合状況手法:

SVM のモデルの作成において試合状況を考慮し、日本代表が勝っている時間帯のツイート, 負けている時間帯のツイート, 同点の時間帯のツイートの学習および判定を分離して行った。例えば、50-49 では 1612 件のツイートがあるが、924 件の勝っている時間帯のツイート, 406 件の負けている時間帯のツイート, 282 件の同点の時間帯のツイートに分離し、それぞれの時間帯のツイートのみで、訓練データとテストデータの準備やハイパーパラメータの選択、テストデータでの判定を行い、最後にそれぞれの時間帯での結果を平均することで精度の算出を行った。SVM のハイパーパラメータ等については SVM 手法と同様に行った。

5.2 実験結果

前処理なしの判定結果について各確信度帯の結果を手法ごとに平均したものを表 3 に示す。また、SVM+試合状況手法における確信度帯別の判定結果を表 4 に示す。

表 3 に示すように、F 値は SVM+試合状況手法が他の手法よりも高い結果となった。また、適合率は SVM+試合状況手法が最も高いが、再現率は SVM 手法が最も高い結果となった。また、表 4 に示すように、確信度帯

別の判定結果は 70-30 の F 値が最も高く、判定しやすいという結果となった。

表 3 ネットバレ判定結果(前処理なし)

手法	適合率	再現率	F 値
パターンマッチ	0.570	0.935	0.708
ランダムフォレスト	0.808	0.857	0.826
SVM	0.809	0.883	0.843
SVM+試合状況	0.831	0.880	0.852

表 4 確信度帯別判定結果(SVM+試合状況手法)

確信度帯	適合率	再現率	F 値
50-49	0.883	0.851	0.866
60-40	0.891	0.856	0.873
70-30	0.873	0.946	0.907
80-20	0.873	0.927	0.893
90-10	0.774	0.890	0.797

さらに、事前処理の有無による結果について、SVM+試合展開手法にローマ字変換を適用した結果を表 5、固有名詞正規化を適用した結果を表 6 にそれぞれ示す。

表 5 に示すように、ローマ字変換を適用すると F 値が向上しており、有用であることがわかった。また、表 6 に示すように、固有名詞は正規化をしないのが最も F 値が高く、どの正規化手法も有用ではないことがわかった。

表 5 ローマ字変換の影響(SVM+試合状況手法)

ローマ字変換	適合率	再現率	F 値
なし	0.831	0.880	0.852
あり	0.844	0.898	0.866

表 6 固有名詞正規化の影響(SVM+試合状況手法)

正規化手法	適合率	再現率	F 値
なし	0.831	0.880	0.852
単純正規化	0.794	0.874	0.828
敵味方別正規化	0.802	0.877	0.836
チーム別正規化	0.824	0.868	0.842
ポジション別正規化	0.815	0.859	0.828

6 考察

ネタバレ判定の結果、SVM+試合展開手法の F 値が最も高く、他の手法よりも優れているという結果となった。特に、適合率が他の手法よりも優れていた。パターンマッチ手法と比較して他の手法の適合率が高かったのは、パターンマッチ手法では特定の単語が存在するだけでネタバレと判断されてしまう場合が多いが、他の手法では単語の数や組み合わせなどの条件が揃うことでネタバレと判断されるためであると考えられる。また、SVM+

試合状況手法の適合率が特に高かったのは、試合状況を考慮しない場合に間違っただけで学習されてしまっていたツイートが、時間帯別にすることによって学習されなくなり、ノイズが軽減されたからだと考えられる。これについては、勝っている時間帯に多いと考えられる「久しぶりにスッキリした」などのツイートで使われるような「久しぶり」と同点時間帯に多いと考えられる「久しぶりにサッカー観る」などのツイートで使われるような「久しぶり」の違いからも理解できる。

また、確信度帯別の判定結果は 70-30 の F 値が最も高くなっていたが、これは、平均確信度が 80 を超えるとスコア情報に関するツイートが多いことから、勝敗に関する単語が少なく特徴量としてあまり学習されなかったためだと考えられる。また、60-40 や 50-49 ではネタバレと非ネタバレのツイート内容の差が小さく、判定難易度が高かったため、F 値がそこまで大きくなかったと考えられる。なお、90-10 ではデータ量が 362 件と少ないため、学習が十分に行えなかった可能性がある。これについては、データ量を増やしていけば結果が変わる可能性があるため、今後の課題とする。

前処理については、ローマ字変換が有用であった。これは、Twitter などでは短文形式で文章の程を成していないだけでなく、単語自体にも雑然な表現が多く、そうした単語をローマ字変換することで特徴量として統一されたためだと考えられる。また、固有名詞の正規化がどれも有用でなかったのは、選手名などを正規化してしまうと、ネタバレツイートに多い試合結果に関与する選手と非ネタバレツイートに多い試合結果に関与しない選手が同じ特徴量として学習されてしまったためだと考えられる。これについては、ここで行った正規化処理よりもさらに細かく正規化することで有用になる可能性があるため、今後の課題とする。

7 おわりに

本研究では、ネタバレツイートを「試合の最終結果が高い確信度で予測できてしまう投稿」と定義し、Twitter 上でのサッカーのネタバレに関するデータセットを構築した。データセットを分析した結果、ツイート時点での試合経過時間が長いほど、試合結果への予測確信度が高く、危険なネタバレとなり得ることや、ネタバレの内容が試合状況によって異なることを明らかにした。さらに、ネタバレの判定精度について、パターンマッチ手法、ランダムフォレスト手法、SVM 手法、SVM+試合状況手法で比較した結果、SVM+試合状況手法の F 値が最も高く有用であることや、前処理としてローマ字変換が有用であることを明らかにした。

今後は、学習データを増やしたりデータの前処理をさらに工夫したりすることで、ネタバレの判定精度を向上させていき、実際にネタバレ防止をクライアントなどの形でシステム化することを考えている。また、本研究では独立したツイートをネタバレ判断基準としたが、前後のツイートやツイート数といった要素も視野に入れたネタバレ特性調査やネタバレ判定手法を検討していきたい。さらに、他のスポーツへの適用実験も行っていく予定である。

謝辞

本研究の一部は、JST ACCEL (Grant 番号 JPMJAC1602) の支援を受けたものである。

参考文献

- [1] Nakamura, S. and Tanaka, K.: Temporal filtering system for reducing the risk of spoiling a user's enjoyment, Proc. of IUI'07, pp. 345-348, 2007.
- [2] 中村聡史, 小松孝徳: スポーツの勝敗にまつわるネタバレ防止手法の検討, 情報処理学会論文誌, Vol. 54, No. 4, pp. 1402-1412, 2013.
- [3] 白鳥裕士, 牧良樹, 中村聡史, 小松孝徳: スポーツにおけるネタバレの特性調査と判定手法の検討, 情報処理学会論文誌, Vol. 59, No. 3, pp. 882-893, 2018.
- [4] Leavitt, J. D. and Christenfeld, N. J. S.: Story spoilers don't spoil stories, Psychological Science, Vol. 22, pp. 1152-1154, 2011.
- [5] Rosenbaum, J. E. and Johnson, B. K.: Who's afraid of spoilers? Need for cognition, need for affect, and narrative selection and enjoyment, Psychology of Popular Media Culture, Vol. 5, pp. 273-289, 2016.
- [6] 田島一樹, 中村聡史: ストーリーコンテンツに対するネタバレの基礎調査とその判定手法の検討, 研究報告グループウェアとネットワークサービス(GN), Vol. 2015-GN-96, No. 7, pp. 1-6, 2015.
- [7] Golbeck, J.: The twitter mute button: a web filtering challenge, Proc. of CHI'12, pp. 2755-2758, 2012.
- [8] Jeon, S., Kim, S. and Yu, H.: Spoiler detection in tv program tweets, Information Sciences, Vol. 329, pp. 220-235, 2016.
- [9] R. A. Baeza-Yates. and B. A. Ribeiro-Neto.: Modern information retrieval: the concepts and technology behind Search (2nd Edition), Addison-Wesley Professional, 2011.
- [10] C. D. Manning. and H. Schtze.: Foundations of Statistical Natural Language Processing, MIT Press, 1999.