

周辺語に基づく有効期限を表す時間表現の判定

長島 弘昂[†] 田島 敬史^{††}

[†] 京都大学大学院情報学研究科 〒606-8501 京都府京都市左京区吉田本町

^{††} 京都大学大学院情報学研究科 〒606-8501 京都府京都市左京区吉田本町

E-mail: [†]nagashima@dl.soc.i.kyoto-u.ac.jp, ^{††}tajima@i.kyoto-u.ac.jp

あらまし 明確な有効期限がある文書は、その有効期限を表す時間表現を含むことが多い。一方で、明確な有効期限のない文書にも、同様の時間表現は含まれることがある。その為、自動的に有効期限を推定して文書の取捨選択をしたい際には、これらを区別する技術が必要となる。そこで、本論文では、文書中の時間表現がその文書の有効期限を表すものかどうかを判定する手法を提案する。本研究では、まず、有効期限のある文書が多く存在する Twitter において、ツイート内に現れる時間表現と、その時間表現が示す時点におけるリツイート数の変化に基づき、有効期限を表していると考えられる時間表現と有効期限を表していないと考えられる時間表現を多数抽出し、それらの周辺語に対し有効期限を表す時間表現と共起する確率を求める。この確率を用いて、ある文書中の時間表現が有効期限を表しているかを推定する。有効期限があると推定された文書に複数の時間表現が含まれる場合、より具体的な時間を示す、より先の時間表現を有効期限として選定する。実験の結果、提案手法は F 値 0.580 の性能が得られ、有効期限の選定では 0.897 の正解率が得られた。

キーワード Twitter, 情報検索, 情報要約, 自然言語処理, 機械学習

時間表現以外の情報から、これらを区別する手法が必要となる。

1 はじめに

近年、特にソーシャルメディアの分野で、リアルタイムに流れる情報の量が增大している。このようなソーシャルメディアには、有効期限のある情報が多数含まれていることから、近年は有効期限のある情報の量も増大傾向にある。そうした状況により、有用な情報を効率的に収集する為の、有効期限に基づく情報フィルタリング手法、及び、情報検索手法に対するニーズは、益々高まっている。

文書は、有効期限に関して、二つの基準から分類できる。一つ目の基準は、有効期限の有無である。例えば、「今日が締切」という文書は、有効期限のある文書と考えられる一方で、「今日も猫は可愛い」という文書は、有効期限のない文書と考えられる。二つ目の基準は、有効期限のある文書における、明示的な時間表現の有無である。例えば、先の「今日が締切」という文書は、「今日」という、明示的な時間表現を含む有効期限のある文書と考えられる一方で、「ご飯食べに行かない？」という文書は、明示的な時間表現を含まない有効期限のある文書と考えられる。

本研究では、情報フィルタリング及び情報検索への応用を想定し、対象を明示的な時間表現を含む文書に限定して、その時間表現がその文書の有効期限を表すものかどうかを判定する手法を提案する。つまり、判定においては、有効期限のない文書にも、同様の時間表現が含まれる可能性を考慮しなければならない。先の例であれば、有効期限のある文書「今日が締切」にも、有効期限のない文書「今日も猫は可愛い」にも、同一の明示的な時間表現「今日」が含まれている。よって、ある時間表現がその文書の有効期限を表すものかどうかを判定する為には、

提案手法では、ある時間表現がその文書の有効期限を表すものかどうかを判定する為に、時間表現の周辺語を用いる。機械学習によって、周辺語とその文書の有効期限の有無の関係を学習させ、このモデルを用いて、周辺語から先の判定を行う。従って、その為のトレーニングデータが必要となるのであるが、本論文では、併せて、文書に対する外部からのリアクションの時間変化に基づき、トレーニングデータを自動的に収集する手法を提案する。

具体的には、本論文では、まず、有効期限のある文書が多く存在する Twitter において、ツイート内に現れる時間表現と、その時間表現が示す時点におけるリツイート数の変化に基づき、有効期限を表していると考えられる時間表現とその他の時間表現を多数抽出する。つまり、ツイートに含まれる時間表現が示す時点において、リツイート数に一定の伸びの鈍化が認められるかを判断材料とし、一定の伸びの鈍化が認められる場合には、その時間表現を有効期限として推定する。その後、この推定によって得られた、有効期限を表していると考えられる時間表現とその他の時間表現の周辺語に対し、有効期限を表す時間表現と共起する確率を求める。最終的に、この確率を用いて、ある文書中の時間表現が有効期限を表しているかを推定する。

実験においては、時間表現を 1 以上含む 25,046 のツイートのツイート内容及びリツイート数の変化、そのツイートに対してツイート内容のみによって人手で付与した有効期限のデータを使用した。結果として、人手で付与した有効期限のデータをトレーニングデータとして使用した場合に一定の性能が得られ、又、ツイート内容及びリツイート数の変化から推定した有効期限のデータをトレーニングデータとして使用した場合にも、先の性能に近い性能が得られた。この二つの結果は、周辺語に基

づく有効期限を表す時間表現の判定が一定の性能を得られること、そして、そのトレーニングデータがツイート内容及びツイート数の変化より推定したものであっても同様であることを示している。

本論文の以降の構成としては、まず、2 節において関連研究について述べる。3 節では提案手法の詳細、続く 4 節では実験の詳細を述べ、5 節で総括を行う。

2 関連研究

本研究は、大まかに下記の二段階の推定によって構成されている。本研究の、周辺語に基づく有効期限を表す時間表現の判定は、下記 1 の結果をトレーニングデータとした、下記 2 によるものである。本節では、主に、下記 2 に対する関連研究を採り上げる。

(1) 時間表現とリアクションの時間変化に基づく有効期限の推定

(2) 周辺語に基づく有効期限を表す時間表現の判定

2.1 有効期限に基づく情報フィルタリングに関する研究

情報フィルタリングに関する研究は数多く存在する。その一方で、情報の有効期限に基づく情報フィルタリングに関する研究は少ない。情報の有効期限に基づく情報フィルタリングに関する研究として、本研究に最も関連すると考えられる研究は、竹村ら [1] による、ツイート分類に関する研究である。竹村らは、Twitter のタイムラインにおいて、重要なツイートを見逃さない為の仕組みとして、ツイートの有効期限によってツイートを分類する手法を提案している。最終的にツイートは、今読むべき、読むのは後でも良い、既に有用でない、の 3 カテゴリに分類されるが、その分類の為に、ツイートの有効期限は 6 クラスに分類される。有効期限の推定には、決定木、及び、SVM を採用している。

本研究とは異なり、竹村らは、対象を、明示的な時間表現を含む文書に限定しておらず、有効期限を 6 クラスとして扱っている。又、分類器の学習において、竹村らは、人手で付与した有効期限のデータを使用しており、大規模なトレーニングデータの確保が困難となっている。

加えて、竹村らは、ツイート中の時間表現、各語の時間依存性を特徴量として含んでおり、この点については本研究と類似している。一方で、竹村らはこれらの他に、文字数、ツイートの種類、短期間で急速に拡散された語 (bursty keywords) の含有、返信の間隔、URL 及び画像の含有を特徴量として採用している為、ツイート以外の文書に対しては、有効期限の同推定手法が適用不可能である。

2.2 有効期限に関する研究

2.1 節で述べた、有効期限に基づく情報フィルタリングに関する研究の他、文書の有効期限を推定する研究として、Almquist ら [2] は、機械学習による推定手法を提案している。しかし、竹村らの研究と同様に、Almquist らも、対象を、明示的な時間表現を含む文書に限定しておらず、有効期限を 5 クラスとして

扱っている。又、モデルの学習において、Almquist らも、人手で付与した有効期限のデータを使用しており、大規模なトレーニングデータの確保が困難となっている。

又、Twitter のようなマイクロブログにおける短い文書の分類手法として、Sriram ら [3] は、文書著者のプロフィールテキストに基づいてマルチクラス分類する手法を提案している。論文において Sriram らは、ツイートを、ニュース、イベント、意見、取引、プライベートメッセージに分類している。確かに、特にイベント等に分類されたツイートが、有効期限のあるツイートである確率は高いと推測されるが、その他クラスであっても、有効期限のあるツイートである可能性はあり、又、提案手法では、その有効期限がいつかまでは分からない。

イベントの発生や隆盛に関する研究としては、樺ら [4] は、Twitter におけるリアルタイムイベントに対するユーザの相互作用を調査することによって、リアルタイムイベントを検知する手法を提案しており、又、Mathioudakis ら [5] は、トレンド検出の手法を提案しているが、有効期限については焦点が当てられていない。

3 提案手法

2 節で述べた通り、本研究は、大まかに二段階の推定によって構成されている。本節においては、一段階目の推定と二段階目の推定とを分け、各詳細を 3.1 節と 3.2 節にて述べる。

3.1 時間表現とリアクションの時間変化に基づく有効期限の推定

3.2 節における、周辺語に基づく有効期限を表す時間表現の判定に当たっては、学習に十分なデータセットが必要となる。従って、本節においては、人手で文書に対して有効期限を付与する代わりに、機械的に大規模なデータセットを確保する手法を提案する。

なお、本論文においては、Twitter を推定環境に設定した。時間表現とリアクションの時間変化に基づく有効期限の推定において、リアクションの時間変化は、時間表現のいずれかが有効期限を表しているかを判定する目的で用いられる。従って、これを満たす情報が得られる環境であれば、推定環境は問わない。一方で、Twitter は、各文書が短く、且つ、大量の文書が日々生成され、又、多種多様な文書が混在しているという、データセットの構築においては極めて有効な複数の要因を満たしており、大規模なデータセットの確保においては高効率であると考えられる為、ここでは Twitter を推定環境に設定した。

従って、本推定は、ツイート内容から時間表現を抽出する段階、リアクションの時間変化からその時間表現のいずれかが有効期限を表しているかを判定する段階で分けられる。一段階目のツイート内容からの時間表現の抽出においては、時刻情報正規化 API¹を用いる為、特筆すべき事項はない。従って、本節においては、リアクションの時間変化からその時間表現のいずれかが有効期限を表しているかを判定する段階について、詳細

1 : <https://labs.goo.ne.jp/api/jp/time-normalization>

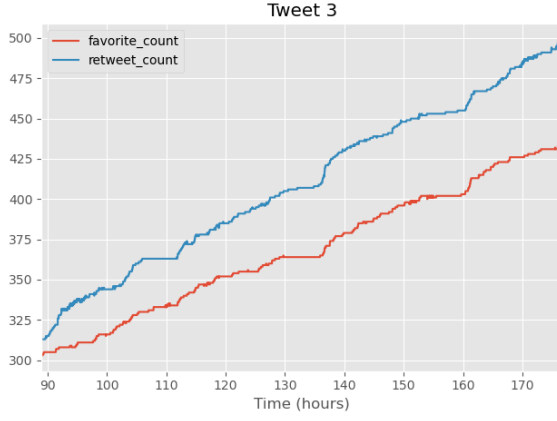


図 1 経過時間に対するリアクションの時間変化の例（有効期限がない場合）

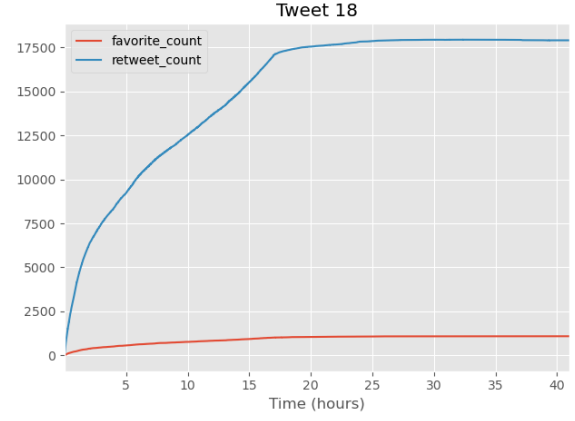


図 2 経過時間に対するリアクションの時間変化の例（有効期限がある場合）

を述べることとする。

3.1.1 Twitter におけるリアクションの時間変化

本論文においては、リアクションの時間変化として、リツイート数の時間変化を採用する。Twitter における有効期限の推定であれば、他にいいね数の時間変化を用いる方法も考えられるが、リツイートが周囲のユーザに対して情報を拡散することを目的とするのに対し、いいねはそのような周知の役割は薄く、ユーザの自己満足を主たる目的とする傾向があると推測される為、本論文においては、リアクションの時間変化としてリツイート数の時間変化を採用する。換言すれば、対象の情報の拡散において、ユーザは、周囲のユーザにとって対象の情報が有用であるかを判断材料としていてと考え、この有用性の時間変化がリツイート数の時間変化となって表出され、ひいては同リアクションの時間変化が有効期限の推定において適切な指標となり得ると判断しているのである。

図 1 から図 3 は、リアクションの時間変化の実例である。横軸がツイートからの経過時間を、縦軸がリツイート及びいいねの回数を表しており、青線がリツイート数を、赤線がいいね数を表している。図 1 においては、リツイート及びいいねの回数は一定の割合で増加し続けているが、図 2、図 3 においては、リツイート及びいいねの回数の伸びは時間の経過に伴って鈍化している。図 1 におけるリツイート及びいいねの回数の時間変化は、対象のツイートに有効期限がないことによって、図 2 におけるリツイート及びいいねの回数の時間変化は、対象のツイートの有効期限によって説明できる。その一方で、図 3 においては、有効期限がないのにも関わらず、リツイート及びいいねの回数の伸びが時間の経過に伴って鈍化しており、これはリツイート数が少ない、即ち、十分に情報の拡散が行われていない、あるいは、ツイートに関係するユーザの絶対数が少ないが為に、有効期限とは無関係にリツイート数が上限に達していることに起因していると考えられる。

3.1.2 リアクションの時間変化に基づくツイートの有効期限の推定

本論文においては、ノイズ低減フィルタ $smooth$ を適用することで平滑化した、リアクションの時間変化 $y(t)$ の二次導関数

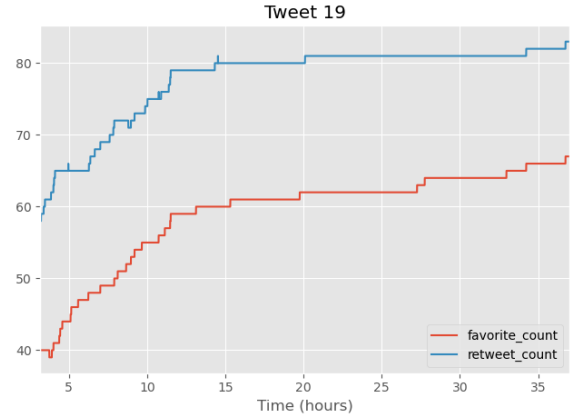


図 3 経過時間に対するリアクションの時間変化の例（有効期限はないが、リツイート数が少ない場合）

に対して、シグモイド関数 $sigmoid$ を乗じることによって有効期限の推定範囲を表現した上で、その負のピークによって有効期限 $t_{expiration}$ を推定する（式 6）。

まず、リアクションの時間変化 $y(t)$ は、後述するノイズ低減フィルタ $smooth$ によって平滑化される（式 1）。得られたリアクションの時間変化 $y_{smooth}(t)$ に対し、一次導関数を求め、同様に平滑化を行い、 $\dot{y}_{smooth}(t)$ を得る（式 2）。同様にして、 $\dot{y}_{smooth}(t)$ に対しても一次導関数を求め、平滑化を行うことで $\ddot{y}_{smooth}(t)$ を得る（式 3）。

$$y_{smooth}(t) := smooth(y(t)) \quad (1)$$

$$\dot{y}_{smooth}(t) := smooth\left(\frac{d}{dt}y_{smooth}(t)\right) \quad (2)$$

$$\ddot{y}_{smooth}(t) := smooth\left(\frac{d}{dt}\dot{y}_{smooth}(t)\right) \quad (3)$$

その後、有効期限の推定範囲を $[0, \arg \max_t y_{smooth}(t)]$ に限定した上で（式 4）、得られた $\ddot{y}_{smooth}(t)$ に $sigmoid$ を乗じることによって更に有効期限の推定範囲を狭め、この関数を最小化する t を有効期限 $t_{expiration}$ と推定する（式 6）。具体的には、 t が 0 から $\arg \max_t y_{smooth}(t)$ に近づくほど、有効期限は

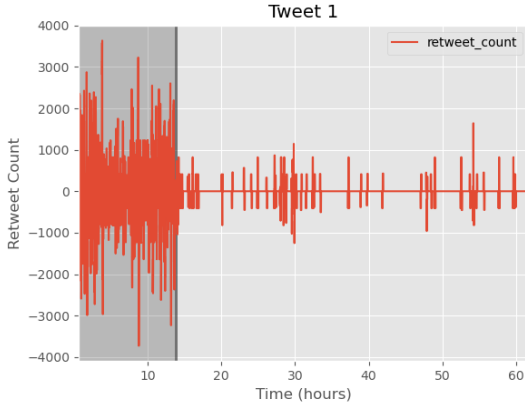


図 4 移動平均幅が 2 の場合のリアクションの時間変化の二次導関数の例

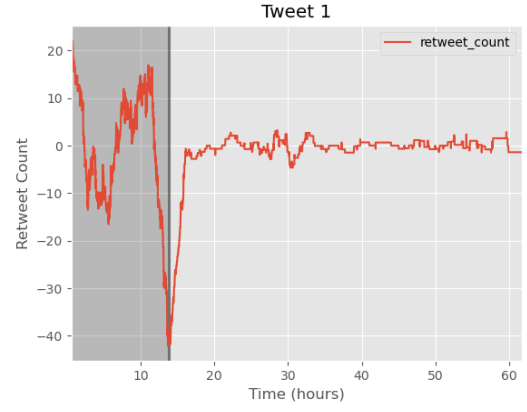


図 5 移動平均幅が 50 の場合のリアクションの時間変化の二次導関数の例

現れやすくなると考え、 t が 0 から $\arg \max_t y_{smooth}(t)$ に近づくほど、大きな重みを与えるよう、*sigmoid* を乗じている。

$$t_{up} = \arg \max_t y_{smooth}(t) \quad (4)$$

$$b'_{sigmoid} = \arg \min_{t \leq t_{up}} |y_{smooth}(t) - b_{sigmoid} y_{smooth}(t_{up})| \quad (5)$$

$$t_{expiration} = \arg \min_{t \leq t_{up}} (\ddot{y}_{smooth}(t) \times \text{sigmoid}(a_{sigmoid}(t - b'_{sigmoid})/t_{up})) \quad (6)$$

ここで、

$$\text{sigmoid}(z) = \frac{1}{1 + e^{-z}} \quad (7)$$

又、 $a_{sigmoid}$ 、 $b_{sigmoid}$ は定数である。

ノイズ低減フィルタ *smooth* の必要性は、図 4 及び図 5 において認められる。図 4 及び図 5 は、同一のリアクションの時間変化の二次導関数について、適用する単純移動平均の幅 w_{ma} を、順に 2、50 と変化させた例である。横軸がツイートからの経過時間を、縦軸がリツイート数を表しており、赤線が移動平均を適用したリツイート数の二次導関数を、灰色の縦線が同ツイート内容から抽出された時間表現が指し示す時間を表している。導関数の次数が増加するにつれて、微分の特性によってノイズも増加することは自明であるが、そのノイズが図 4 において認められる。しかしながら、移動平均の幅 w_{ma} を増加させることで、図 5 には灰色の縦線で示した時間において負のピークが現れている。なお、実際、同ツイートの有効期限は灰色の縦線で示した時間である。

本論文においては、ノイズ低減フィルタ *smooth* として、Savitzky-Golay フィルタを用いる。Savitzky-Golay フィルタは、一定幅のデータについて、最小二乗法によって多項式近似を行うことを繰り返し、ノイズの低減を行う。この際、高次多項式を採用することによって、元となるデータのフィーチャを減衰させることなく、平滑化が可能である。この特徴は、Savitzky-Golay フィルタが、平滑化においてフィーチャ減衰が起こる単純移動平均よりも優れていることを示しており [6]、主にリアクションの時間変化の二次導関数の負のピーク、即ち、

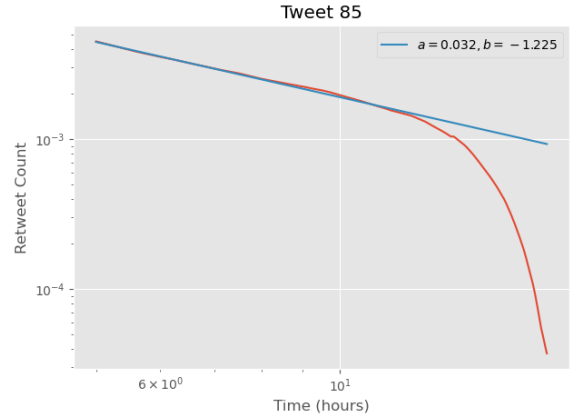


図 6 power law with exponential cutoff に従う単位時間当たりのリツイート数の例

ウェッジを基に有効期限の推定を行う本手法において、有効であると判断した為、同 Savitzky-Golay フィルタをノイズ低減フィルタとして採用した。Savitzky-Golay フィルタの主なハイパーパラメータは、多項式近似を行う幅 w_{sg} 、及び、多項式の次数 o_{sg} である。又、実装においては、SciPy [7] を使用する。

一般に、単位時間当たりのリツイート数は power law with exponential cutoff に従って減衰する [8]。式 9 が power law with exponential cutoff を表す式であり、式中の a 、 b 、 c は定数である。

$$p(t) = at^{-b} \quad (8)$$

$$p_{cut}(t) = at^{-b}e^{-ct} \quad (9)$$

図 6 は、power law with exponential cutoff に従っていると考えられる、単位時間当たりのリツイート数の例である。横軸がツイートからの経過時間を、縦軸がリツイート数を表しており、赤線が単位時間当たりのリツイート数を、青線が power law (式 8) に従う場合のその近似を表している。図 6 からは、 $t = 13$ 付近までは power law に従っているものの、以降においては急速に減衰していることが分かる。これは、power law with exponential cutoff の exponential cutoff、 e^{-ct} の作用に

よるものであると考えられる。

よって、正確にこの法則に従う限り、リアクションの時間変化の二次導関数 $\ddot{y}_{smooth}(t)$ も、cutoff の周辺で唯一のピークを示すことはない。一方、有効期限が存在するケースにおいては、リアクションの時間変化の収束付近において、より急激な伸びの鈍化が起こると予測される為、この現象によって有効期限を推定するのである。急激な伸びの鈍化の例は、図 2 において示した通りである。

最終的に、この手法によりリアクションの時間変化のみから推定される有効期限と、実際にツイート内容から抽出される時間表現を比較し、その誤差が 1 時間以内である場合に、その時間表現が有効期限であると判定する。

3.2 周辺語に基づく有効期限を表す時間表現の判定

本推定では、まず、有効期限を表す時間表現の含有判定を行う。その後、有効期限を表す時間表現が文書に含まれると推定された場合には、有効期限を表す時間表現の選定を行う。本節においては、有効期限を表す時間表現の含有判定と選定とを分け、各詳細を 3.2.1 節と 3.2.2 節にて述べる。

3.2.1 周辺語に基づく有効期限を表す時間表現の含有判定

本論文では、周辺語に基づく有効期限を表す時間表現の含有判定においては、ナイーブベイズ手法の多項モデル及びベルヌーイモデル [9] を用いて、即ち、周辺語の出現情報のみに基づいて、文書の 2 クラス分類を行う。ここで、周辺語とは、同一文書中の、時間表現を除いた全ての語を指す。

クラス分類の際、適切に語のフィルタリングを行う。本手法においては、品詞によるフィルタリング、データセットにおける語の総出現頻度数によるフィルタリングを行う。従って、本手法における主要なハイパーパラメータは、使用する品詞集合 S_{POS} 、データセットにおける語の総出現頻度数に基づく上位採用件数 k である。その他基準による語のフィルタリングについては、Schneider による論文 [10] が詳しい。

又、実装においては、scikit-learn [11] を使用する。

3.2.2 有効期限を表す時間表現の選定

有効期限を表す時間表現が文書に含まれると推定され、且つ、その文書に複数の時間表現が含まれる場合には、有効期限を表す時間表現の選定が必要となる。

本論文で使用する時刻情報正規化 API からは、多様な単位の時間表現が得られる。例えば、「今日の 20 時が締切」という文書からは、今日の日付と対応付けられた「今日」という時間表現と、直前に現れた今日の、20 時の日時と対応付けられた「20 時」という時間表現が抽出される。なお、この文書の最適な有効期限が、後者の「20 時」であることを考慮すれば、前者の「今日」は、有効期限としては不十分である。

従って、本論文では、時刻情報正規化 API・有効期限・日本語の性質を踏まえ、二つの手法を提案する。一つ目の手法は、最後尾に現れる時間表現を有効期限とする手法、二つ目の手法は、日付より日時を優先し、その内、最も先の時間表現を有効期限とする手法である。

4 実 験

本節では実験の詳細を述べる。

4.1 実験手順

実験は以下の手順で行う。

(1) Twitter において、ツイート内容及びリツイート数の変化を大量に収集する。

(2) 収集したツイートから時間表現を抽出する。

(3) 時間表現を含まないツイートを除外する。

(4) 人手で各ツイートに有効期限を付与する。

(5) 様々なハイパーパラメータ $a_{sigmoid}$, $b_{sigmoid}$, w_{sg} , o_{sg} で、ツイート内容及びリツイート数の変化に基づいて、ツイートの有効期限を推定する。

(6) ツイート及び有効期限を、トレーニングデータとテストデータに、4:1 の割合で分割する。

(7) トレーニングデータを用いて、ナイーブベイズ分類器について 5-分割交差検証を行い、最適なハイパーパラメータ $a_{sigmoid}^*$, $b_{sigmoid}^*$, w_{sg}^* , o_{sg}^* , S_{POS}^* , k^* を決定する。なお、学習には 5 で推定した有効期限を、検証には 4 で付与した有効期限を用いる。

(8) テストデータを用いて、最適なハイパーパラメータにおける性能を測定する。なお、学習には 5 で推定した有効期限を、検証には 4 で付与した有効期限を用いる。

なお、1 の Twitter におけるデータ収集の詳細は以下の通りである。リツイート数の下限は、リツイート数の時間変化から有効期限を抽出しやすくする為に (3.1.1 節)、又、キューは、全ツイートの照会頻度を可能な限り均等にする為に用いている。データ収集においては、Twitter の REST API を使用する。

(1) リツイート数が 50 以上の、日本語のツイートを検索する。

(2) 新規ツイートを収集キュー Q にエンキューする。

(3) Q から 100 件ツイートをデキューし、現在のリツイート数を照会する。

(4) デキューしたツイートを、再び Q にエンキューする。

(5) Q が 10,000 未満ならば 1 へ、そうでなければ 3 へ戻る。

4.2 データセット

実験においては、合計 70,000 のツイートの内、時間表現を 1 以上含む 25,046 のツイートについて、ツイート内容及びリツイート数の変化、及び、ツイート内容のみに基づく人手による有効期限をデータセットとして用いる。25,046 のツイートの内、人手によって有効期限を表す時間表現を含むと判断されたツイートは、6,398 件である。データセットの内、ツイート内容及びリツイート数の変化は学習に、人手による有効期限は検証、及び、周辺語に基づく有効期限を表す時間表現の判定のみの有効性の検証に用いられる。

詳細には、2019 年 12 月 31 日、2020 年 1 月 1 日、2020 年 1 月 11 日、2020 年 1 月 12 日、2020 年 1 月 14 日、2020 年 1

表 1 データセットサイズ

期間	#時間表現を含む	#全体
2019 年 12 月 31 日～2020 年 01 月 29 日	4433	10000
2020 年 01 月 01 日～2020 年 01 月 29 日	4987	10000
2020 年 01 月 11 日～2020 年 01 月 29 日	3065	10000
2020 年 01 月 12 日～2020 年 01 月 29 日	2937	10000
2020 年 01 月 14 日～2020 年 01 月 29 日	3525	10000
2020 年 01 月 15 日～2020 年 01 月 29 日	3302	10000
2020 年 01 月 17 日～2020 年 01 月 29 日	3366	10000

表 2 使用したハイパーパラメータ

	MNB	BNB
$a_{sigmoid}^*$	100	
$b_{sigmoid}^*$	0.9	
w_{sg}^*	101	
o_{sg}^*	0	
S_{POS}^*	$U \setminus \{\text{"動詞"}\}$	
k^*	729	

表 3 周辺語に基づく有効期限を表す時間表現の判定のみの有効性の検証に使用したハイパーパラメータ

	MNB	BNB
S_{POS}^*	$U \setminus \{\text{"動詞"}\}$	$U \setminus \{\text{"フィラー"}\}$
k^*	19683	

月 15 日, 2020 年 1 月 17 日, 各々から 2020 年 1 月 29 日まで, 各期間, データの収集を開始してから, リツイート数が 50 以上のツイートが 1 万件集まるまでのツイートについて, リアクションの時間変化を記録し続けた (4.1 節). なお, 各期間の時間表現を含むツイート数は, 表 1 の通りである. 年末年始においては, 時間表現を含むツイート数が多いことが分かる.

4.3 ハイパーパラメータ

提案手法におけるハイパーパラメータは, $a_{sigmoid}$, $b_{sigmoid}$, w_{sg} , o_{sg} , S_{POS} , k の 6 つである. 6 ハイパーパラメータの内, $a_{sigmoid}$, $b_{sigmoid}$, w_{sg} , o_{sg} は時間表現とリアクションの時間変化に基づく有効期限の推定に, S_{POS} , k は周辺語に基づく有効期限を表す時間表現の判定に用いられる.

最終的に用いたハイパーパラメータ $a_{sigmoid}^*$, $b_{sigmoid}^*$, w_{sg}^* , o_{sg}^* , S_{POS}^* , k^* を, 表 2 に示す. ここで, $U = \{\text{"その他"}, \text{"フィラー"}, \text{"感動詞"}, \text{"記号"}, \text{"形容詞"}, \text{"助詞"}, \text{"助動詞"}, \text{"接続詞"}, \text{"接頭詞"}, \text{"動詞"}, \text{"副詞"}, \text{"名詞"}, \text{"連体詞"}\}$ である. 又, MNB がナイーブベイズ分類器の多項モデル, BNB がナイーブベイズ分類器のベルヌーイモデルを表している.

又, 周辺語に基づく有効期限を表す時間表現の判定のみの有効性の検証に用いたハイパーパラメータ S_{POS}^* , k^* を, 表 3 に示す.

表 2 と表 3 は, ツイート内容及びリツイート数の変化に基づいて推定した有効期限よりデータセットを構築した場合 (表 2) と, 人手による有効期限よりデータセットを構築した場合 (表 3) とで, 周辺語に基づく有効期限を表す時間表現の判定における最適なハイパーパラメータ S_{POS}^* , k^* が異なることを示し

表 4 人手による有効期限をトレーニングデータとして使用した場合の性能

モデル	P	R	F1	BACC	MCC
MJR	-	0.000	-	0.500	-
RND	0.255	0.508	0.339	0.499	-0.001
MNB	0.731	0.629	0.676	0.773	0.576
BNB	0.750	0.597	0.665	0.763	0.570

表 5 提案手法に基づく推定による有効期限をトレーニングデータとして使用した場合の性能

モデル	P	R	F1	BACC	MCC
MJR	-	0.000	-	0.500	-
RND	0.255	0.508	0.339	0.499	-0.001
MNB	0.781	0.455	0.575	0.705	0.501
BNB	0.823	0.448	0.580	0.707	0.520

ている. 特に, 人手による有効期限よりデータセットを構築した場合の最適なハイパーパラメータ k^* からは, より多くの語を推定に用いる傾向が認められる.

4.4 評価指標

評価指標としては, Precision, Recall, F1-score, Balanced accuracy, Matthews correlation coefficient (MCC) を採用した. 各指標において, 有効期限を表す時間表現を含む文書を正としている. 対象とするデータのクラスサイズがインバランスである為, F1-score の他, MCC も評価指標として採用した. 実際には, 少数データである有効期限を表す時間表現を含む文書を正に設定している為, F1-score のみによる評価も可能であるが, 一般において有効な評価指標として, 又, ランダムにクラス分類を行うモデルに対しても妥当な評価を与える評価指標として, ここでは併記している. MCC, 及び, 各評価指標における MCC の優位性については, Brown による論文 [12], 並びに, Chicco による論文 [13] が詳しい.

4.5 実験結果

実験結果は表 4 及び表 5 の通りである. 表 4 が人手による有効期限をトレーニングデータとして使用した場合の性能を, 表 5 が提案手法に基づく推定による有効期限をトレーニングデータとして使用した場合の性能を表している. 各表において, P が Precision, R が Recall, F1 が F1-score, BACC が Balanced accuracy, MCC が Matthews correlation coefficient, 又, MJR が多数決による分類器 (本実験では, 全文書を有効期限を表す時間表現を含む文書ではないと分類), RND がランダムに分類を行う分類器, MNB がナイーブベイズ分類器の多項モデル, BNB がナイーブベイズ分類器のベルヌーイモデルを表している. MJR 及び RND はベースラインとして記載している.

実験の結果としては, 人手で付与した有効期限のデータをトレーニングデータとして使用した場合に一定の性能が得られ, 又, ツイート内容及びリツイート数の変化より推定した有効期限のデータをトレーニングデータとして使用した場合にも先の性能に近い性能が得られた. この二つの結果は, 周辺語に基づ

表 6 有効期限を表す時間表現の選定における性能

手法	正解率
FIRST	0.435
LAST	0.838
MIN*	0.864
MAX*	0.897

く有効期限を表す時間表現の判定が一定の性能を得られること、そして、そのトレーニングデータがツイート内容及びリツイート数の変化より推定したものであっても同様であることを示している。

3.2.2 節で述べた、有効期限を表す時間表現の選定に関する実験結果は、表 6 に示した通りである。提案手法に基づく推定による有効期限をトレーニングデータとして、ナイーブベイズ分類器のベルヌーイモデルを分類器として使用した場合の、真陽性における正解率を求めた。LAST が最後尾に現れる時間表現を有効期限とする手法を、MAX*が日付より日時を優先し、その内、最も先の時間表現を有効期限とする手法を指している。なお、それぞれの比較として、先頭に現れる時間表現を有効期限とする手法（FIRST）と、日付より日時を優先し、その内、最も前の時間表現を有効期限とする手法（MIN*）の結果を掲載している。

表 6 は、提案手法の有効性、特に、日付より日時を優先する手法の有効性を示している。FIRST と LAST の正解率の差についても、日本語における日付と時間の語順によって説明できると考えられる。又、MIN*と MAX*の正解率の差は、より先の時間表現ほど有効期限となる傾向を示唆している。

5 結 論

本研究では、情報フィルタリング及び情報検索への応用を想定し、対象を明示的な時間表現を含む文書に限定して、その時間表現がその文書の有効期限を表すものかどうかを判定する手法を提案した。判定においては、有効期限のない文書にも、同様の時間表現が含まれる可能性を考慮しなければならず、時間表現以外の情報から、これらを区別する手法が必要となった。

提案手法では、ある時間表現がその文書の有効期限を表すものかどうかを判定する為に、時間表現の周辺語を用いた。機械学習によって、周辺語とその文書の有効期限の有無の関係を学習させ、このモデルを用いて、周辺語から先の判定を行うのである。従って、その為のトレーニングデータが必要となるのであるが、本論文では、併せて、文書に対する外部からのリアクションの時間変化に基づき、トレーニングデータを自動的に収集する手法を提案した。

具体的に、本論文では、まず、有効期限のある文書が多く存在する Twitter において、ツイート内に現れる時間表現と、その時間表現が示す時点におけるリツイート数の変化に基づき、有効期限を表していると考えられる時間表現とその他の時間表現を多数抽出した。つまり、ツイートに含まれる時間表現が示す時点において、リツイート数に一定の伸びの鈍化が認められ

るかを判断材料とし、一定の伸びの鈍化が認められる場合には、その時間表現を有効期限として推定するのである。その後、この推定によって得られた、有効期限を表していると考えられる時間表現とその他の時間表現の周辺語に対し、有効期限を表す時間表現と共に確率を求め、この確率を用いて、ある文書中の時間表現が有効期限を表しているかを推定した。

実験においては、時間表現を 1 以上含む 25,046 のツイートのツイート内容及びリツイート数の変化、そのツイートに対してツイート内容のみによって人手で付与した有効期限のデータを使用した。結果として、人手で付与した有効期限のデータをトレーニングデータとして使用した場合に一定の性能が得られ、又、ツイート内容及びリツイート数の変化から推定した有効期限のデータをトレーニングデータとして使用した場合にも、先の性能に近い性能が得られた。この二つの結果は、周辺語に基づく有効期限を表す時間表現の判定が一定の性能を得られること、そして、そのトレーニングデータがツイート内容及びリツイート数の変化より推定したものであっても同様であることを示している。

今後の課題としては、異なる環境における実験が挙げられる。まず、英語のツイートでも、同様の結果が得られるのか検証し、その後、ニュース等を対象にした実験も行う。

謝 辞

本研究は、JST CREST (JPMJCR16E3)、JSPS 科研費 18H03245 の支援を受けたものである。

文 献

- [1] Hikaru Takemura and Keishi Tajima. Tweet classification based on their lifetime duration. In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management, CIKM '12*, p. 2367–2370, New York, NY, USA, 2012. Association for Computing Machinery.
- [2] Axel Almqvist and Adam Jatowt. Towards content expiry date determination: Predicting validity periods of sentences. In *Proceedings of the 41st European Conference on Information Retrieval*, pp. 86–101, 04 2019.
- [3] Bharath Sriram, Dave Fuhry, Engin Demir, Hakan Ferhatosmanoglu, and Murat Demirbas. Short text classification in twitter to improve information filtering. In *Proceedings of the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '10*, p. 841–842, New York, NY, USA, 2010. Association for Computing Machinery.
- [4] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Earthquake shakes twitter users: Real-time event detection by social sensors. In *Proceedings of the 19th International Conference on World Wide Web, WWW '10*, p. 851–860, New York, NY, USA, 2010. Association for Computing Machinery.
- [5] Michael Mathioudakis and Nick Koudas. Twittermonitor: Trend detection over the twitter stream. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data, SIGMOD '10*, p. 1155–1158, New York, NY, USA, 2010. Association for Computing Machinery.
- [6] J. Guiñón, Emma Ortega, José García-Antón, and Valentín Pérez-Herranz. Moving average and savitzki-golay smoothing filters using mathcad. *Papers ICEE*, 01 2007.

- [7] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, CJ Carey, İlhan Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antônio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. Scipy 1.0: Fundamental algorithms for scientific computing in python. *Nature Methods*, 2020.
- [8] Peter Mathews, Lewis Mitchell, Giang Nguyen, and Nigel Bean. The nature and origin of heavy tails in retweet activity. In *Proceedings of the 26th International Conference on World Wide Web Companion*, WWW '17 Companion, p. 1493–1498, Republic and Canton of Geneva, CHE, 2017. International World Wide Web Conferences Steering Committee.
- [9] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, USA, 2008.
- [10] Karl-Michael Schneider. A comparison of event models for naive bayes anti-spam e-mail filtering. In *Proceedings of the Tenth Conference on European Chapter of the Association for Computational Linguistics - Volume 1*, EACL '03, p. 307–314, USA, 2003. Association for Computational Linguistics.
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, Vol. 12, pp. 2825–2830, 2011.
- [12] J. B. Brown. Classifiers and their metrics quantified. *Molecular Informatics*, Vol. 37, No. 1-2, p. 1700127, 2018.
- [13] Davide Chicco. Ten quick tips for machine learning in computational biology. *BioData Mining*, Vol. 10, No. 1, p. 35, 2017.