

# バンディットアルゴリズムを用いた 特定地域から発信されたツイートの収集

上田 紗希<sup>†</sup> 山口 祐人<sup>††</sup> 北川 博之<sup>†††</sup>

<sup>†</sup> 筑波大学システム情報工学研究科 〒 305-8573 茨城県つくば市天王台 1-1-1

<sup>††</sup> 国立研究開発法人産業技術総合研究所 〒 135-0064 東京都江東区青海 2-3-26

<sup>†††</sup> 筑波大学システム情報系情報工学域 〒 305-8573 茨城県つくば市天王台 1-1-1

E-mail: <sup>†</sup>braose@kde.cs.tukuba.ac.jp, <sup>††</sup>yamaguchi.yuto@aist.go.jp, <sup>†††</sup>kitagawa@cs.tsukuba.ac.jp

**あらまし** 本研究では Twitter から特定地域より発信されたツイートを収集する手法を提案する。ツイートの中にはユーザの現在の状況や周囲での出来事などユーザの現在位置と結びついている情報が多く含まれる。しかし Twitter はツイート収集のための API に利用制限があり、また発信地を示すジオタグが付与されたツイートの割合も非常に少ない。そこで本研究では、収集対象地域からツイートを発信する可能性の高いユーザを発見・フォローすることで対象地域のツイートを収集する手法を提案する。限られた時間の中で目的とするツイートを収集するために、ユーザの探索と発見したユーザからのツイートの収集のトレードオフを考慮する必要がある。そのため提案手法では、バンディットアルゴリズムによりフォローするユーザを決定する。また、ユーザの位置情報に関する報酬は直接得ることが出来ないため、ツイート内容から推定する。また実データを用いた評価実験により、提案手法の有効性を示す。

**キーワード** Twitter, Focused Crawling, バンディットアルゴリズム

## 1. 序 論

### 1.1 研究の背景と目的

今日では、ソーシャルメディアを通じて多くの人々が多様な情報発信を行うようになった。Twitter<sup>(注1)</sup> は主要なサービスの一つであり、ユーザはツイートと呼ばれる 140 文字以内の短文を投稿することができる。Twitter には携帯端末からも手軽にアクセスすることができるため、投稿されたツイートの中にはユーザの現在の状況や周囲での出来事など、ユーザの現在位置と深く結びついた情報が含まれていることも多い。そのため、特定の地域から発信されたツイートを収集することができれば、その地域に対するマーケティング、イベント検出 [12] や災害分析 [14] などに応用することが可能であると考えられる。

ツイートを収集する際に、大きく分けて 2 つの問題点がある。1 つ目は収集者の計算機的、および時間的資源に限られていることである。Twitter には常に大量のツイートが世界中から投稿されており、それらのツイートを取捨選択することなく収集・保存し続けることは収集者にとって計算機資源上の大きな負担である。また一般的に収集にかかる時間は有限であるため、限られた時間の中で目的のツイートをできるだけ多く収集することが求められる。2 つ目は Twitter はツイートを収集するための API を公開しているが、その枠組みの中で特定地域から発信されたツイートを収集することは難しいことである。Twitter の API を用いる場合、時間あたりの利用回数や取得ツイート数等には制限が設けられている。例えば投稿されている

ツイートをリアルタイムに収集できる Streaming APIs<sup>(注2)</sup> の中の statuses/sample は、投稿された全ツイートのうちランダムに選ばれた 1% を取得できる。しかし statuses/sample によってランダムに収集されたツイートから収集者が目的とするツイートを探すことは非効率である。statuses/filter は、パラメタにより設定した条件に合致したツイートを収集することができる。follow パラメタでは、一度に最大で 5000 ユーザまでのユーザ ID を指定することで、それらのユーザのツイートを取得することができる。また locations パラメタでは、経度と緯度の組み合わせによる任意の範囲を指定することで、その範囲内のジオタグが付与されたツイートを収集することができる。しかし先行研究 [5] によると、Twitter に投稿されたツイートのうち、ジオタグが付与されたツイートは全体のわずか 0.42% であり、多くのユーザはプライバシー等の問題から自身の位置情報を明らかにしていない。従って先に述べたような理由から Twitter の既存の枠組み内では、限られた時間の中で特定地域のツイートを多く収集することは困難であると考えられる。

本研究では収集対象地域から発信されたツイートを継続的に収集するために、その地域からツイートを投稿する可能性の高いユーザを発見・フォローする。ユーザをフォローすることによって特定地域からのツイート収集を行う場合、まず膨大な Twitter ユーザの中から収集対象地域からツイートを投稿する可能性の高いユーザを探索することが必要である。またユーザの探索と同時に、発見したユーザをフォローしそのツイートを収集することも求められる。一度にフォローできるユーザ数や収集期間は限られているため、収集者はこの二つの行動のトレ

(注1) : <https://twitter.com/>

(注2) : <https://dev.twitter.com/streaming/public>

ドオフを考慮しながらツイートの収集を進めなければならない。そこで本研究では、バンディットアルゴリズム [13] を用いて特定地域からのツイートを収集する手法を提案する。バンディットアルゴリズムはデータの分布が未知である環境下で、得られる累積報酬が最大になるよう行動の選択を繰り返し行うアルゴリズムである。バンディットアルゴリズムでは多くの報酬を得ることのできる選択肢を探す「探索」と、多くの報酬が得られると判明した選択肢を選んで実際に報酬を獲得する「活用」の二つのトレードオフを考慮しながら累積報酬の最大化を目指す。本研究では収集対象地域とフォロー対象の候補となるユーザ集合を入力し、その地域から発信されたツイートを出力として受け取るという状況を想定する。提案手法ではバンディットアルゴリズムにおける選択をユーザのフォロー、報酬をユーザから得られる対象地域から発信されたツイート数と位置付けることで、ツイート収集問題にバンディットアルゴリズムを適用する。本研究における累積報酬の最大化は、対象地域から発信されたツイートの最終的な収集数を最大化することに該当する。

一般にバンディットアルゴリズムにおける報酬は、行動を選択することによって明確に与えられる。しかし本研究の場合、ジオタグが付与されたツイートの割合は非常に小さく、ツイートから直接位置情報に関する情報を得て報酬を計算することはできない。そこで収集したツイートのテキストを用いて発信地推定を行うことでユーザが対象地域からツイートを投稿した確率を計算し、それを報酬とする。

本研究の貢献は以下に示す 2 点である。

(1) テキストから推定した報酬を用いるバンディットアルゴリズムにより、特定地域から発信されたツイートを継続的に収集する手法を提案する。

(2) 実際の Twitter データを用いた評価実験により、提案手法によるツイート収集の有効性を示す。

提案手法の有効性を示すため本研究で行う評価実験は以下の 2 つである。

- **報酬の推定方法の評価**：提案手法における報酬の推定方法が適切なものであるかを検証する。複数の報酬推定方法を比較し、提案手法で用いる推定方法が対象地域から発信されたツイートを最も多く収集できることを示す。

- **提案手法と他手法の比較**：提案手法とベースライン手法の比較を行い、提案手法が対象地域から発信されたツイートを最も多く収集できることを示す。

## 2. 関連研究

### 2.1 特定トピックに関するツイート収集手法

ある特定のトピックに関するツイートを Twitter から収集する研究として、Gisselbrecht らは WhichStreams? [7] を提案した。この手法はバンディットアルゴリズムを用い、あるトピックと関係が深いユーザの投稿を収集する手法である。Gisselbrecht らの手法では本研究と同様にバンディットアルゴリズムを用いており、収集したツイートから得られる報酬に従ってフォローユーザの再選択を行う。しかし本研究では多くの場合ツイートの発信地が未知であるため、Gisselbrecht らの手法と違い、収

集したツイートそのものから直接報酬を得ることはできない。そのため本研究では収集したツイートのテキストを用いて発信地推定を行うことで報酬を推定し、推定された報酬に従ってバンディットアルゴリズムを用いてフォローユーザの選択を行う。

Li ら [10] はトピックに関するツイートを収集するために効果的なキーワードを自動的に検出する手法を提案した。Li らの手法ではツイートを収集するための検索キーワード集合を決定し、収集したツイートに含まれる語句をもとに一定時間ごとに検索キーワードを更新する。本研究はキーワードによる検索ではなくユーザのフォローによってツイートを収集する点において Li らの手法とは異なっている。

### 2.2 ユーザおよびツイートの位置推定手法

Twitter における位置推定手法は、1) ユーザの居住地推定手法、2) ツイートの発信地推定手法の 2 種類に分けることができる。

#### 2.2.1 居住地推定手法

ユーザの居住地推定に関わる手法は多く提案されている。居住地推定に関する既存手法として、投稿されたツイートの内容を用いるコンテンツベース手法と、ユーザのソーシャルグラフを用いるグラフベース手法がある。

コンテンツベース手法の一つとして、Cheng ら [5] はローカルワードを用いてユーザの居住地推定を行う手法を提案した。ローカルワードとはその単語を投稿したユーザの居住地に偏りがあるような単語のことである。また、Yamaguchi ら [16] はソーシャルストリームから検出された地理的局所性を持つイベントを用いてユーザの居住地推定を行う手法を提案した。

グラフベース手法として、Backstrom ら [3] や Clodoveu [6] は、ソーシャルメディアユーザの友人関係に着目して居住地を推定する手法を提案した。これらの手法ではソーシャルメディア上で友人関係にあるユーザ同士は現実の居住地も近い可能性が高いとし、ユーザの居住地を推定している。

しかし、これらの手法の目的はユーザの居住地を推定することであり、特定地域から発信されたツイートを収集するという本研究の目的とは異なっている。

#### 2.2.2 発信地推定手法

ユーザの発信地推定に関わる手法はこれまでもいくつか提案されている。Kinsella ら [9] はジオタグ付きツイートから各地域の言語モデルを構築し、この言語モデルからツイートの発信地を推定する手法を提案した。Priedhorsky ら [11] は混合ガウス分布を用いる手法を提案した。Ikawa ら [8] は、ユーザの過去の投稿内容からツイートの位置情報とキーワードの関連付けを行う手法を提案した。

しかし、これらの手法の目的は個々のツイートの発信地を明らかにすることである。本研究では一つの地域に焦点を当てその地域から発信されたツイートを収集するため、対象の地域以外から発信されたツイートは必要としていない。そのため、これらの既存手法では対象地域外に推定された大半のツイートの収集データや収集に要した時間が無駄になるという問題点がある。本研究ではバンディットアルゴリズムを用いて対象地域からツイートを投稿する可能性の高いユーザを探索・フォローし、

効率的にデータを収集する。

### 3. バンディットアルゴリズム

バンディット問題 [2] はデータの分布が未知である環境下での選択の最適化を考える問題である。バンディット問題では毎回の試行で有限の集合内から選択を行い、その選択に対して報酬を受け取る状況を考える。このとき試行回数には制限があり、各選択に対する報酬の分布は事前に分かっていないものとする。バンディット問題の最終的な目的は、このような制約下で最終的に得られる累積報酬を最大化することである。

累積報酬を最大化するために、選択者は「探索」と「活用」という二種類の行動を使い分ける。得られる報酬の分布は未知であるため、選択者は複数回の試行を行って報酬の分布を調べる必要がある。この行動が探索である。一方で、報酬の少ない選択を繰り返し行うことは損であり、より良い選択枝を探索すると同時に現時点で多くの報酬を得られる選択枝を何度も選ぶことも求められる。この行動を活用と呼ぶ。選択者はこの「探索」と「活用」の二つの行動のトレードオフを考慮しながら、累積報酬を最大化することを目指す。

バンディット問題に関するアルゴリズムは、UCB アルゴリズム [2] や  $\epsilon$ -greedy アルゴリズム [15], softmax アルゴリズム [4], Thompson sampling アルゴリズム [1] など数多く提案されている。UCB アルゴリズムは各選択枝の現時点までの試行回数を考慮しながら探索と活用のバランスをとるアルゴリズムである。 $\epsilon$ -greedy アルゴリズムはある一定の確率  $\epsilon$  でランダムな選択を行い、確率  $1 - \epsilon$  で現時点で最も報酬の期待値が高い選択を行うアルゴリズムである。softmax アルゴリズムは  $\epsilon$ -greedy アルゴリズムの探索において、選択によって得られる報酬が明らかになっている場合でも、報酬の大小に関わらず全ての選択を等しい確率で行ってしまうという欠点を解決する手法である。softmax アルゴリズムでは、探索の際に報酬の比によって確率的に選択を行うことで、より大きな報酬が得られる選択を高い確率で試行できる。Thompson sampling アルゴリズムは全ての選択枝の報酬に対して予め事前分布を仮定しておき、その後は得られた報酬によって更新された事後確率に従って選択を行うアルゴリズムである。

### 4. 問題定義

本章では、本研究で用いる用語および本研究で扱う問題を定義する。本研究では連続した時間を一定時間ごとに区切り、一つ一つの時間帯をタイムウィンドウと呼ぶ。各タイムウィンドウにおいてユーザをフォローし、ツイートを収集する。それぞれの投稿におけるユーザ  $u$ , タイムウィンドウ  $t$ , 発信地  $l$  を合わせてツイート  $p = (u, t, l)$  と定義する。ユーザ, ツイート, 発信地の集合をそれぞれ  $U, P, L$  で表す。ツイートは発信地  $l$  が Twitter 上で公開されているツイート  $p \in P^L$  と、公開されていないツイート  $p \in P^U$  とに分けられ、 $P = P^L \cup P^U$  となる。またタイムウィンドウ  $t$  においてフォローするユーザをフォローユーザ集合  $U^t$  で表す。収集対象地域を  $\hat{l}$  で表す。これらの記号の定義を表 1 にまとめる。

表 1 記号と定義

記号	定義
$u$	ユーザ
$t$	タイムウィンドウ
$l$	ツイートの発信地
$\hat{l}$	収集対象地域
$p$	ツイート
$U$	ユーザ集合
$U^t$	タイムウィンドウ $t$ におけるフォローユーザ集合
$P$	ツイート集合
$P^L$	発信地が公開されているツイート集合
$P^U$	発信地が公開されていないツイート集合
$L$	発信地集合

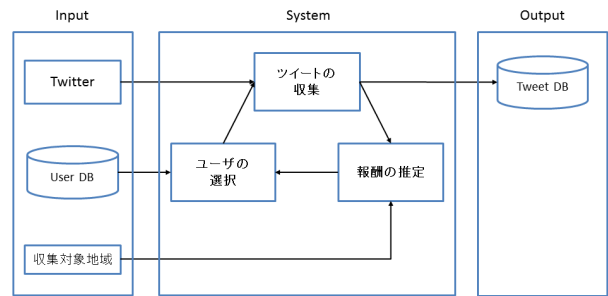


図 1 提案手法全体図。提案手法の処理の流れを示す。提案手法ではユーザ集合と収集対象地域が与えられたとき、バンディットアルゴリズムを用いて対象地域から発信されたツイートの収集を行う。

このとき、対象地域から発信されたツイート収集問題を次のように定義する。ユーザ集合  $U$ , 収集対象地域  $\hat{l}$ , 収集期間  $T$  が与えられたとき、以下の二つの制約下で最終的なツイート収集数  $|\{p_i \in P \mid l_i = \hat{l}\}|$  が最大となることが期待される各タイムウィンドウ  $t$  においてフォローするユーザ集合  $U^t$  の系列  $\{U^1, U^2, \dots, U^T\}$  を出力する。

- 毎タイムウィンドウでユーザ集合  $U$  から  $K$  ユーザを選択してフォローする。
- 収集期間  $T$  の間、ユーザからツイートを収集する。

## 5. 提案手法

### 5.1 手法概要

本研究は、対象地域から発信されたツイートを収集することを目的とする。提案手法では対象地域からツイートを発信する可能性の高いユーザをフォローすることによってそのユーザのツイートを収集する。そこでツイート収集問題にバンディットアルゴリズムを適用し、ユーザのフォローを選択、ユーザから得られる対象地域から発信されたツイート数を報酬ととらえることで、一定時間ごとにフォローユーザの選択を繰り返し行う。これは、過去に頻繁に対象地域からツイートを投稿しているユーザは今後も対象地域を訪れてツイートを投稿する可能性が高いという仮定に基づいている。また、発信地を示すジオタグが付与されているツイートはごく少数であり、ツイートの位置情報を得て報酬を計算することは多くの場合不可能であると

いう問題点がある。そのため、ツイートのテキストを用いて発信地推定を行うことでユーザが対象地域からツイートを投稿した確率を計算し、そのユーザから得られる報酬を推定する。

提案手法の流れを図 1 に示す。提案手法ではフォロー候補ユーザ集合と収集対象地域が入力として与えられる。提案手法の処理は大きく以下の三つのステップに分けられる。この三つのステップをタイムウィンドウごとに繰り返す。

(1) **ユーザの選択**: タイムウィンドウ  $t$  におけるフォローユーザ集合をバンディットアルゴリズムを用いて決定する。

(2) **ツイートの収集**: フォローユーザ集合に含まれるユーザのツイートをタイムウィンドウ  $t$  の期間内、収集する。

(3) **報酬の推定**: 収集したツイートを元に、各ユーザから得られた報酬を推定する。

## 5.2 ユーザの選択

ユーザの選択では、タイムウィンドウ  $t$  においてフォローするユーザ集合を決定する。バンディットアルゴリズムは複数存在するが、本研究ではその中の一つである  $\epsilon$ -Greedy アルゴリズム [15] を用いる。  $\epsilon$ -Greedy アルゴリズムは通常は最も報酬の高い選択枝を選択する活用を行いながら、同時に確率  $\epsilon$  で全選択枝の中からランダムに選択枝を選ぶ探索を行うアルゴリズムである。探索を行うことで、エージェントは現時点で最も報酬の高い選択よりもより良い選択の可能性を探ることができる。

提案手法の場合では、確率  $\epsilon$  で報酬とは無関係にランダムなユーザ選択を行い、確率  $1 - \epsilon$  で最も報酬の期待値が高いユーザを選択する。タイムウィンドウ  $t$  時点でのユーザ  $u$  の報酬の期待値  $Q_{u,t}$  は以下の式で与えられる。

$$Q_{u,t} = \begin{cases} 0 & (F_{u,t-1} = 0) \\ \frac{1}{F_{u,t-1}} \sum_{i=1}^{t-1} g_{u,i} & (\text{otherwise}) \end{cases} \quad (1)$$

ただし、 $g_{u,t}$  はタイムウィンドウ  $t$  におけるユーザ  $u$  の報酬、 $F_{u,t-1}$  はタイムウィンドウ  $t-1$  までにユーザ  $u$  がフォローされた回数とする。あるタイムウィンドウにおいて一度にフォローできるユーザ数を  $K$  とするとき、この選択を  $K$  回繰り返し行うことで、フォローするユーザ集合を決定する。なお、1 度目の選択については全ユーザの累積報酬が 0 であるため、フォローユーザ  $K$  人をランダムに選択する。

## 5.3 ツイートの収集

選択したフォローユーザ集合のツイートをタイムウィンドウ  $t$  の期間内、収集する。ツイートの収集には Streaming APIs 中の statuses/filter より follow パラメータを用いる。

## 5.4 報酬の推定

タイムウィンドウ  $t$  の終了後、フォローユーザから得られた報酬の推定を行う。まず、タイムウィンドウ  $t$  間に収集したツイートから、各ツイートの発信地が対象地域である確率を計算する。確率を計算するために、提案手法ではナイーブベイズ分類器 (Naive Bayesian Classifier) を用いる。入力として特徴ベクトル  $X$  (後述) が与えられたとき、クラスラベルが  $L$  である確率は以下の式で計算される。

$$P(L|X) = \frac{P(L)P(X|L)}{P(X)} \propto P(L)P(X|L) \quad (2)$$

学習データには事前に収集したジオタグ付きツイートをを用い、クラスは発信地が対象地域  $\hat{l}$  であるか否かの二通りとした。ナイーブベイズ分類器にはラプラススムージングを用いた。本研究では特徴ベクトル  $X$  をタイムウィンドウ  $t$  間に収集したツイート毎に作成する。まず、位置情報が既知であるツイートをを用い、特徴語として名詞を抽出する。そして抽出された特徴語から bag-of-words に従ってツイートを特徴ベクトルで表す。bag-of-words は、ベクトルの各次元を一つの単語に対応付け、各文書内のその単語の頻度を値とするベクトル表現方法である。特徴語集合を  $W = \{w_1, w_2, \dots, w_n\}$  とし、ツイート  $p$  内の単語  $w_i$  の出現回数を  $n(w_i, p)$  で表すとき、ツイート  $p$  の特徴ベクトル  $X$  は以下の式で表される。

$$X(p) = (n(w_1, p), n(w_2, p), \dots, n(w_n, p)) \quad (3)$$

しかし実際のデータにおいては、対象地域外から投稿されたツイートの数は対象地域から投稿されたツイートの数を大きく上回るという問題がある。このようなクラスの分布の偏りによる影響を解消するため、class mass normalization (CMN) [17] を用いる。CMN は学習データのクラスの分布と現実のデータのクラスの分布とが大きく異なるような場合に分類結果の補正を行う。label 1 と label 0 の二つのクラスが存在し、サンプル  $i$  が label 1 に分類される確率を  $f(i)$  とするとき、通常のカテゴリではサンプル  $i$  は一般に以下の式に従ってクラス  $l_i$  に分類される。

$$l_i = \begin{cases} \text{label 1} & (f(i) > \frac{1}{2}) \\ \text{label 0} & (\text{otherwise}) \end{cases} \quad (4)$$

一方、CMN を用いる場合、サンプル  $i$  は以下のように分類される。

$$l_i = \begin{cases} \text{label 1} & (q \frac{f(i)}{\sum_i f(i)} > (1-q) \frac{1-f(i)}{\sum_i (1-f(i))}) \\ \text{label 0} & (\text{otherwise}) \end{cases} \quad (5)$$

このとき、 $q$  は理想的な label 1 の割合である (実験ではデータから算出する)。

タイムウィンドウ  $t$  において、ナイーブベイズ分類器から推定されるユーザ  $u$  のツイート  $p_i$  の発信地が対象地域  $\hat{l}$  である確率を  $P(\hat{l}|p_i)$ 、タイムウィンドウ  $t$  内でユーザ  $u$  から発信されたツイート集合を  $P_u^t$  とする。このときタイムウィンドウ  $t$  でのユーザ  $u$  の報酬  $g_{u,t}$  は、 $P_u^t$  に含まれる全ツイートの  $P(\hat{l}|p_i)$  の CMN による補正值の和として以下の式で与えられる。

$$g_{u,t} = \sum_{p_i \in P_u^t} \frac{c_{p_i}^p}{c_{p_i}^p + c_{p_i}^n} \quad (6)$$

ただしこのとき、 $c_{p_i}^p$ 、 $c_{p_i}^n$  は式 5 における  $q \frac{f(i)}{\sum_i f(i)}$ 、 $(1-q) \frac{1-f(i)}{\sum_i (1-f(i))}$  に該当するため、

$$\begin{cases} c_{p_i}^p = q \frac{P(\hat{l}|p_i)}{\sum_u \sum_{p_j \in P_u^t} P(\hat{l}|p_j)} \\ c_{p_i}^n = (1-q) \frac{1-P(\hat{l}|p_i)}{\sum_u \sum_{p_j \in P_u^t} (1-P(\hat{l}|p_j))} \end{cases} \quad (7)$$

となる。

推定された報酬を用いて、タイムウィンドウ  $t+1$  で再びフォローするユーザ集合を選択する。

## 6. 評価実験

### 6.1 実験設定

提案手法の有効性を検証するため、実際のツイートを用いて特定地域から発信されたツイートの収集を行った。実験は全て64bit Windows8.1, Intel Core i7-4820K@3.70GHz CPU, 32GB RAMで構成されたPC上で行った。提案手法のプログラムはPythonで記述し、ライブラリはscikit-learn, Numpyを用いた。

まず、Streaming APIsのstatuses/filterを用い、2016年5月30日から2016年6月7日の間に日本国内のジオタグ付きツイートを収集した。次に収集したジオタグ付きツイートをユーザごとに集計した後、収集数の多い上位2万ユーザを選出し、実験に用いるため、これらのユーザのツイートを2016年7月16日16時から2016年8月10日16時までの期間収集した。収集したツイートのうち、ユーザがリツイートしたツイートを除いたものをテストデータセットに含めた。テストデータセット内に含まれる最終的なツイート総数は7,992,311件であった。

今回、人口の異なる複数の地域で検証を行うため、東京23区、京都市、横浜市、つくば市の4地域を対象地域として設定した。ナイーブベイズ分類器の学習データとして用いるため、各対象地域について対象地域内とそれ以外の日本国内地域から、テストデータセットには含まれないジオタグ付きツイートをそれぞれ1万件ずつ収集した。特徴語を抽出する際の形態素解析にはMeCab<sup>(注3)</sup>を用い、辞書はIPA辞書<sup>(注4)</sup>を使用した。提案手法のパラメタについては事前の予備実験より、 $K = 1000$ ,  $\epsilon = 0.5$ , タイムウィンドウの長さは4時間とした。また、CMNのパラメタ $q$ を設定するため、テストデータセットには含まれない日本国内のジオタグ付きツイートを別途10万件収集した。実際に実験に用いたパラメタ $q$ の値について、表2に示す。

表2 パラメタ $q$ の値

対象地域	東京都23区	京都市	横浜市	つくば市
$q$	0.13	0.01927	0.02356	0.00251

### 6.2 評価方法

本実験では対象地域から発信されたツイートの収集数によって評価を行った。実際に収集したツイートの大部分はジオタグが付与されておらず、ツイートそのものからそのツイートの発信地を判断することは難しい。そのため、本実験では以下の方法によってジオタグが付与されていないツイート $p \in P^U$ の発信地を推測し、対象地域から発信されたツイートの収集数を計算した。

タイムウィンドウ $t$ におけるフォローユーザ集合 $U^t$ に含まれるユーザ $x$ は、収集されたツイートに従って以下の三種類のユーザ集合 $U_A^t$ ,  $U_B^t$ ,  $U_C^t$ に分類される。

(A)  $u \in U_A^t$ : タイムウィンドウ $t$ 内でジオタグ付きツイート $p \in P^L$ を1件以上投稿した。

(B)  $u \in U_B^t$ : タイムウィンドウ $t$ 内でツイートを1件以上投稿したが、その全てがジオタグの付与されていないツイート $p \in P^U$ である。

(C)  $u \in U_C^t$ : タイムウィンドウ $t$ 内でジオタグの有無に関わらずツイートを1件も投稿していない。

上記のように分類したユーザ $u$ について、タイムウィンドウ $t$ での対象地域から発信されたツイート収集数 $n_{u,t}$ を以下のように計算し、その合計を評価に用いた。

- (A)のユーザ: ユーザ自身が投稿したジオタグ付きツイートの情報からユーザが対象地域から発信したツイート数を推算する。タイムウィンドウ $t$ においてユーザが投稿したツイート数とジオタグ付きツイート数から対象地域から発信されたツイートの割合を計算する。

$$n_{u,t} = \frac{|\{p_i \in (P^L \cap P_u^t) \mid l = \hat{l}\}|}{|P^L \cap P_u^t|} |P_u^t| \quad (8)$$

- (B)のユーザ: ユーザが投稿したジオタグ付きツイートがないため、(A)のユーザのツイートの情報を用いてユーザが対象地域から発信したツイート数を推算する。タイムウィンドウ $t$ における(A)の全ユーザの $n_{u,t}$ の合計と総ツイート数から対象地域から発信されたツイートの割合を計算する。

$$n_{u,t} = \frac{\sum_{u \in U_A^t} n_{u,t}}{\sum_{u \in U_A^t} |P_u^t|} |P_u^t| \quad (9)$$

- (C)のユーザ: タイムウィンドウ $T$ 内でツイートを1件も投稿していないため、収集数は0である。

$$n_{u,t} = 0 \quad (10)$$

### 6.3 実験結果と考察

提案手法の有効性を、以下の点から検証した。

- (1) 報酬の推定方法の評価
- (2) 提案手法と他手法の比較

#### 6.3.1 報酬の推定方法の評価

本実験では、提案手法で用いた報酬の推定方法が対象地域から発信されたツイートの収集に適しているのかを検証した。今回は、以下の2点を比較した。

- 特徴語の効果: 提案手法では特徴語に名詞を用いた。特徴語にKLダイバージェンスによって抽出された単語を用いた手法と提案手法とを比較する。

- ナイーブベイズ分類器に用いる特徴ベクトルの作成方法の効果: 提案手法ではツイートごとに特徴ベクトルを作成し、ナイーブベイズ分類器による確率計算を行った。同一ユーザのタイムウィンドウ内でのツイートを一つのテキストにまとめることで特徴ベクトルを作成しナイーブベイズ分類器による確率計算を行う手法と提案手法とを比較する。

実験では上記2点の組み合わせによる報酬の推定方法を提案手法で用いた計算式以外に3つ用意し、計4つの推定方法の比較を行った。

(注3): <http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html>

(注4): <https://drive.google.com/uc?export=download&id=0B4y35FiV1wh7MwV1SDBCXZMTXm>

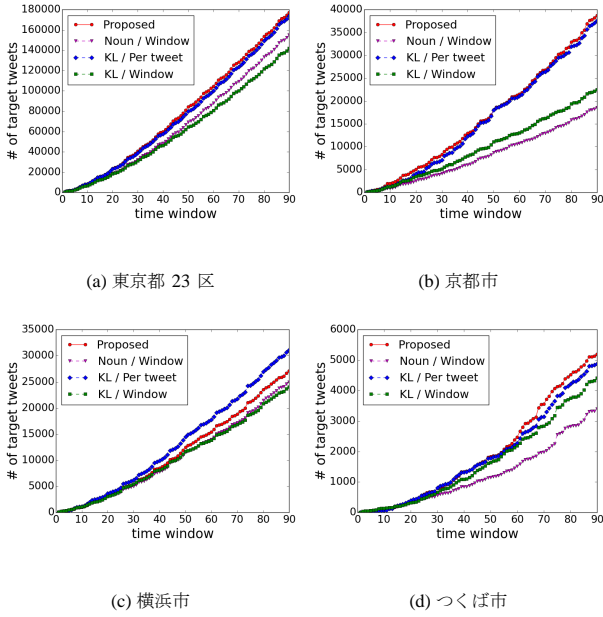


図2 報酬の推定方法の比較。x軸はタイムウィンドウを表し、y軸は各タイムウィンドウまでの対象地域から発信されたツイート数の累計収集数を表す。どの地域の場合でも特徴ベクトルをツイートごとに作成する手法が最も多くのツイートを収集できている。

**KL ダイバージェンスによる特徴語の抽出** 特徴語の効果の比較では、特徴語に KL ダイバージェンスによって抽出された単語を用いる手法を新たに作成した。KL ダイバージェンスは、2つの確率分布の差異を計る尺度である。特徴語の抽出では、単語  $w$  の離散確率分布を  $Q_w$ 、ラベルの離散確率分布を  $P$  として、各単語について KL ダイバージェンスを計算する。いずれかのラベルの値をとる確率変数  $l$  があるとき、 $Q_w(l)$  は単語  $w$  のラベル  $l$  での発生確率を表す。また、 $P(l)$  はラベル  $l$  の発生確率である。このとき、 $Q_w$  からみた  $P$  の KL ダイバージェンスは以下の式で計算される。

$$D_{KL}(Q_w \parallel P) = \sum_l Q_w(l) \log \frac{Q_w(l)}{P(l)} \quad (11)$$

本実験ではまず位置情報が既知であるツイートから名詞を抽出し、ラベルを対象地域から発信されたツイートと対象地域外から発信されたツイートの2つに分けることで各単語の KL ダイバージェンスを計算した。計算された各単語の KL ダイバージェンスのうち、値が大きいものから上位 500 語を特徴語として用いた。

**特徴ベクトルの作成** ナイーブベイズ分類器に用いる特徴ベクトルの作成方法の効果の比較では、報酬推定の際、同一ユーザのタイムウィンドウ内での全ツイートから一つの特徴ベクトルを作成し、ナイーブベイズ分類器による確率値の計算を行う手法を新たに作成する。特徴ベクトルの作成ではまず、同一ユーザがタイムウィンドウ内で投稿した全てのツイートのテキストを繋げて一つの文章にとらえ、テキスト  $s$  とする。その後テキスト  $s$  から特徴語を抽出し、提案手法と同様に bag-of-words に従って特徴ベクトルを作成する。その後発信地が対象地域

である確率  $P(\hat{l}|s_{u,t})$  をナイーブベイズ分類器によって計算し、CMN による補正を行う。

また、報酬の推定時にはユーザが投稿したツイート数を考慮する。タイムウィンドウ  $t$  におけるユーザ  $x$  の報酬  $g_{u,t}$  は以下の式で与えられる。

$$g_{u,t} = \frac{c_{s_{u,t}}^p}{c_{s_{u,t}}^p + c_{s_{u,t}}^n} \log(|P_u^t| + 1) \quad (12)$$

ただしこのとき、

$$\begin{cases} c_{s_{u,t}}^p = q \frac{P(\hat{l}|s_{u,t})}{\sum_{u \in U^t} P(\hat{l}|s_{u,t})} \\ c_{s_{u,t}}^n = (1-q) \frac{1-P(\hat{l}|s_{u,t})}{\sum_{u \in U^t} (1-P(\hat{l}|s_{u,t}))} \end{cases} \quad (13)$$

とする。

### 6.3.2 比較手法

提案手法と比較するために実際に用いた報酬の推定方法は以下の3つである。

- **計算方法1 (KL / Per Tweet)**: 特徴語に KL ダイバージェンスによって抽出された単語を用いる。特徴ベクトルの作成は提案手法と同じ手法を用い、ツイートごとに作成する。報酬の計算式は提案手法と同様に式 (6) を用いる。

- **計算方法2 (Noun / Window)**: 特徴語に名詞を用いる。特徴ベクトルは同一ユーザのタイムウィンドウ内でのツイートを一つの文章にまとめたテキスト  $s$  から作成する。報酬の計算式は式 (12) を用いる。

- **計算方法3 (KL / Window)**: 特徴語に KL ダイバージェンスによって抽出された単語を用いる。特徴ベクトルは同一ユーザのタイムウィンドウ内でのツイートを一つの文章にまとめたテキスト  $s$  から作成する。報酬の計算式は式 (12) を用いる。

### 6.3.3 実験設定・結果

実験ではテストデータセットより、2016年7月16日16時から2016年7月31日16時までの計360時間分のツイートデータを用いた。実験に用いたツイート数は計4,874,303件である。実験は5回行い、その平均値を結果として示す。実験結果を図2に示す。

実験結果から4地域全ての場合で、提案手法を含む特徴ベクトルをツイートごとに作成する手法が最も多くのツイートを収集できていることが分かる。これは、特徴ベクトルをタイムウィンドウ内でのツイートを一つの文章にまとめたテキストから作成する手法では、途中でユーザの対象地域外への移動やその地域に関係のないツイートが含まれていた場合に、その影響を大きく受けてしまうからだと考えられる。ナイーブベイズ分類器では特徴語の出現確率の積によって、ツイートの確率値を計算している。そのため、出現確率が低い単語が含まれていた場合、その影響によって最終的に算出される確率値は小さくなってしまふと考えられる。タイムウィンドウ内の複数のツイートを一つにまとめて特徴ベクトルを作成する場合、そのツイート全てにこのような影響が及んでしまう。Twitterではユーザの移動やツイートが言及する話題が頻繁に変わることが考えられるため、1つのツイートの内容の影響を他のツイートが大きく受けないよう考慮することが効果的であると考えられる。

また、東京都 23 区、京都市、つくば市の 3 地域では特徴語に名詞を用いる提案手法が最も良い結果となっている。ツイートは 140 字以内という制限があるため、その中で用いられる単語は多くない。KL ダイバージェンスによって抽出された単語を用いる場合、特徴語となる単語数が名詞全てを用いる場合よりも少なくなる。そのため一つのツイート中に含まれる特徴語が少ない、または全く特徴語が含まれないツイートが増え、結果的に推定の手がかりが少なくなってしまうからだと考えられる。

#### 6.4 他手法との比較実験

本実験ではベースラインとなる比較手法と提案手法とを比較し、提案手法が効果的に対象地域からのツイートを収集できているかを評価した。ベースラインとして用いたのは以下の 3 つの手法である。

- **Statistics NTW**: 対象地域からツイートを発信する可能性が高い  $K$  ユーザを探し、そのユーザをフォローし続けることでツイートを収集する。最初の  $N$  タイムウィンドウまではランダムに選出した  $NK$  ユーザを順番に一度ずつフォローし、 $N + 1$  タイムウィンドウ以降は  $N$  タイムウィンドウまでに得られた報酬上位  $K$  ユーザをフォローし続ける。
- **Number**: ツイート数の多いユーザをフォローすることでツイートを収集する。タイムウィンドウ内で投稿されたツイート数を報酬とする  $\epsilon$ -Greedy アルゴリズムを用いてフォローユーザを選出する。報酬の期待値の計算は提案手法と同じく式 (1) を用いる。
- **Random**: ランダムにフォローユーザを選びツイートを収集する。毎タイムウィンドウ、ユーザ集合の中からランダムに  $K$  ユーザを選出する。

*Number* のパラメタ  $\epsilon$  については  $\epsilon = 0.5$  とした。実験は各手法 20 回ずつ行い、その平均値を結果として示す。実験結果を図 3、図 4 に示す。

図 3 から、4 地域全ての場合で提案手法が全比較手法を最終的なツイート収集数で上回っている。*Statistics* では最初の  $N$  タイムウィンドウの報酬からフォローするユーザを決定するため、 $N + 1$  タイムウィンドウ以降にフォローユーザが対象地域を訪れなくなったり、逆にフォロー外のユーザが対象地域を訪れるようになったとしてもユーザの行動の変化には対応できない。提案手法ではバンディットアルゴリズムを用いることでユーザの行動の変化に合わせてフォローユーザを変更することができ、より効率的に対象地域から発信されたツイートを収集できたと考えられる。また、*Number* ではツイートを多く投稿するユーザをフォローするためにバンディットアルゴリズムを用いた。しかし、ツイートを多く投稿するユーザは必ずしも対象地域からツイートを発信しているとは限らない。提案手法では報酬の推定の際にユーザが対象地域からツイートを投稿した確率を考慮することで、適切なユーザをフォローすることができたと考えられる。

また、図 4 は実験結果のうち、開始から 25 タイムウィンドウ目までを拡大したものである。最終的なツイート収集数では、提案手法の他に *Statistics 18TW* も比較的多くのツイートを収集できているが、収集開始直後に着目すると、提案手法

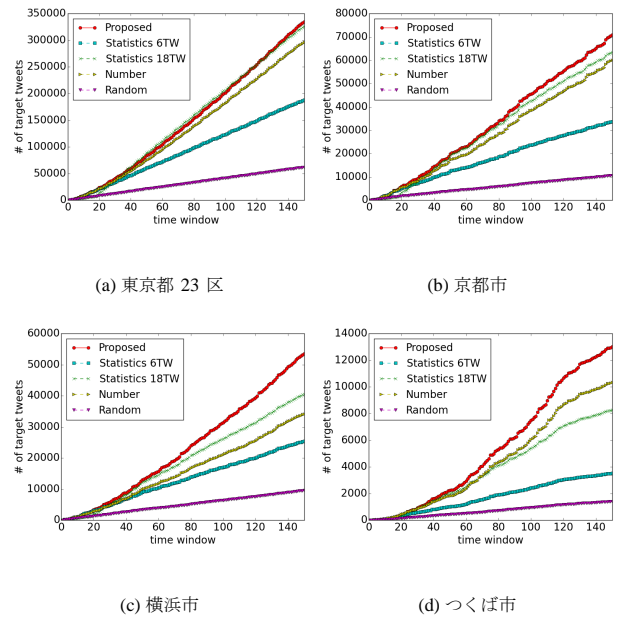


図 3 ベースラインとの比較 (150 タイムウィンドウ). x 軸はタイムウィンドウを表し、y 軸は各タイムウィンドウまでの対象地域から発信されたツイートの累計収集数を表す。全地域で提案手法が比較手法を最終的なツイート収集数で上回っていることが分かる。

と *Statistics 18TW* の差は大きい。*Statistics* の場合、収集開始から  $N$  タイムウィンドウまではユーザを順にフォローしてどのユーザが対象地域からツイートを発信する可能性が高いか調べる必要がある、この期間に仮に対象地域からツイートを発信する確率の高いユーザを発見できたとしても、そのユーザをフォローすることができるのは  $N + 1$  タイムウィンドウ以降になってしまう。提案手法の場合、バンディットアルゴリズムを用いることで対象地域からツイートを発信する可能性が高いユーザを探す「探索」と、実際に可能性が高いとわかったユーザのツイートを収集する「活用」の二つを同時に行うことができるため、収集開始直後からツイートが集まりやすいという利点がある。

さらに、4 地域のうち、人口の少ない京都市、横浜市、つくば市は人口の多い東京都 23 区に比べて提案手法と比較手法間の最終的なツイート収集数の差が大きくなった。人口の少ない地域では膨大なユーザ集合の中からその地域を訪れるユーザを探すことは困難であると考えられるが、提案手法の場合はこのような条件下でもユーザを探索しながら効率的に対象地域からのツイートが収集できている。

## 7. 結論

本研究では Twitter から特定地域より発信されたツイートを収集する手法を提案した。本研究では過去に頻繁に対象地域からツイートを投稿しているユーザは今後も対象地域を訪れてツイートを投稿する可能性が高いという仮定に基づき、バンディットアルゴリズムを用い、ユーザのフォローを選択、ユーザから得られる対象地域から発信されたツイートを報酬ととらえるこ

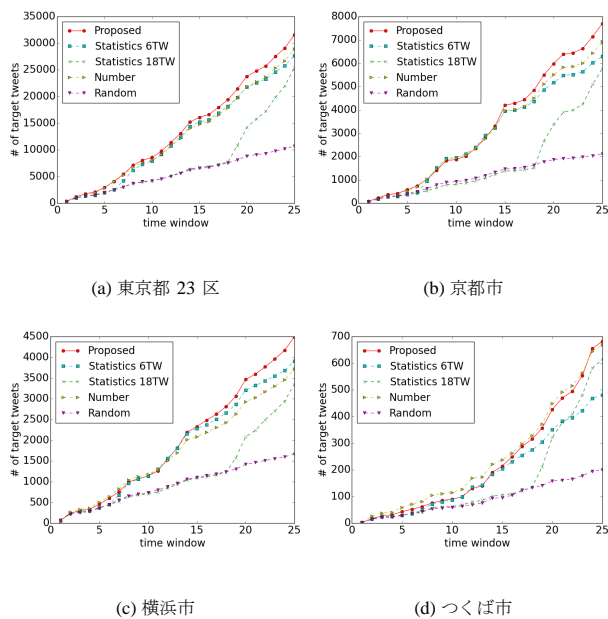


図 4 ベースラインとの比較 (25 タイムウィンドウ) . x 軸はタイムウィンドウを表し, y 軸は各タイムウィンドウまでの対象地域から発信されたツイートの累計収集数を表す. 提案手法と *Statistics 18TW* の差は大きい.

とで, 一定時間ごとにフォローユーザの選択を行った. また, 発信地を示すジオタグが付与されているツイートはごく少数であり, ツイートから位置情報に関する報酬を直接得ることができないという問題点があった. そのため, ツイートのテキストを用いて発信地推定を行うことでツイートから得られる報酬を推定した. 本研究の貢献は, 以下の 2 点である.

(1) バンディットアルゴリズムを用いて, 特定地域から発信されたツイートを継続的に収集する手法を提案した. また, 直接得られない報酬をツイートのテキストから推定する手法を提案した. (第 5 章)

(2) 実際の Twitter データを用いた評価実験により, 提案手法によるツイート収集の有効性を示した. (図 2, 図 3)

今後の課題として以下の 3 つが挙げられる. まず, 提案手法内で用いた報酬の推定方法よりもさらに有効な報酬推定の方法を検討し, より効率的なツイート収集を行うことである. 提案手法では発信地推定にナイーブベイズ分類器を用いたが, 他の分類機の検討により推定の精度を高めることが挙げられる. また提案手法ではツイートごとに発信地推定を行ったが, 投稿間隔の短い前後のツイートの内容を考慮するなどの手法も考えられる. 2 つ目に他のバンディットアルゴリズムの検討が挙げられる. 提案手法では  $\epsilon$ -greedy アルゴリズムを用いたが, バンディットアルゴリズムは他にも様々な手法が存在するため, より効果的な手法の検討が考えられる. 最後に, 今回は最長で 600 時間分のデータを用いて実験を行ったが, より長期的な実験を行った場合にも提案手法はユーザの行動の変化に対応しながらツイートの収集が可能であるか検証することが考えられる.

- [1] Shipra Agrawal and Navin Goyal. Analysis of thompson sampling for the multi-armed bandit problem. In *COLT*, pp. 39–1, 2012.
- [2] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Mach. Learn.*, Vol. 47, No. 2-3, pp. 235–256, May 2002.
- [3] Lars Backstrom, Eric Sun, and Cameron Marlow. Find me if you can: improving geographical prediction with social and spatial proximity. In *Proceedings of the 19th international conference on World wide web*, pp. 61–70. ACM, 2010.
- [4] Nicolò Cesa-Bianchi and Paul Fischer. Finite-time regret bounds for the multiarmed bandit problem. In *ICML*, pp. 100–108. Citeseer, 1998.
- [5] Zhiyuan Cheng, James Caverlee, and Kyumin Lee. You are where you tweet: A content-based approach to geo-locating twitter users. In *Proceedings of the 19th ACM International Conference on Information and Knowledge Management, CIKM '10*, pp. 759–768, New York, NY, USA, 2010. ACM.
- [6] Clodoveu A Davis Jr, Gisele L Pappa, Diogo Rennó Rocha de Oliveira, and Filipe de L Arcanjo. Inferring the location of twitter messages based on user relationships. *Transactions in GIS*, Vol. 15, No. 6, pp. 735–751, 2011.
- [7] Thibault Gisselbrecht, Ludovic Denoyer, Patrick Gallinari, and Sylvain Lamprier. Whichstreams: A dynamic approach for focused data capture from large social media. In *Proceedings of the Ninth International Conference on Web and Social Media, ICWSM 2015, University of Oxford, Oxford, UK, May 26-29, 2015*, pp. 130–139, 2015.
- [8] Yohei Ikawa, Miki Enoki, and Michiaki Tatsubori. Location inference using microblog messages. In *Proceedings of the 21st International Conference on World Wide Web, WWW '12 Companion*, pp. 687–690, New York, NY, USA, 2012. ACM.
- [9] Sheila Kinsella, Vanessa Murdock, and Neil O'Hare. "i'm eating a sandwich in glasgow": Modeling locations with tweets. In *Proceedings of the 3rd International Workshop on Search and Mining User-generated Contents, SMUC '11*, pp. 61–68, New York, NY, USA, 2011. ACM.
- [10] Rui Li, Shengjie Wang, and Kevin Chen-Chuan Chang. Towards social data platform: Automatic topic-focused monitor for twitter stream. *Proc. VLDB Endow.*, Vol. 6, No. 14, pp. 1966–1977, September 2013.
- [11] Reid Priedhorsky, Aron Culotta, and Sara Y Del Valle. Inferring the origin locations of tweets with quantitative confidence. In *Proceedings of the 17th ACM conference on Computer supported cooperative work & social computing*, pp. 1523–1536. ACM, 2014.
- [12] 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*, pp. 851–860. ACM, 2010.
- [13] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, Vol. 1. MIT press Cambridge, 1998.
- [14] Sarah Vieweg, Amanda L Hughes, Kate Starbird, and Leysia Palen. Microblogging during two natural hazards events: what twitter may contribute to situational awareness. In *Proceedings of the SIGCHI conference on human factors in computing systems*, pp. 1079–1088. ACM, 2010.
- [15] Christopher John Cornish Hellaby Watkins. *Learning from Delayed Rewards*. PhD thesis, King's College, Cambridge, UK, May 1989.
- [16] Yuto Yamaguchi, Toshiyuki Amagasa, Hiroyuki Kitagawa, and Yohei Ikawa. Online user location inference exploiting spatiotemporal correlations in social streams. In *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management*, pp. 1139–1148. ACM, 2014.
- [17] Xiaojin Zhu, Zoubin Ghahramani, John Lafferty, et al. Semi-supervised learning using gaussian fields and harmonic functions. In *ICML*, Vol. 3, pp. 912–919, 2003.