

マイクロブログのコンテキストを用いた 行動予測への確率過程モデルの適用と評価

高山 翼¹ 山上 慶² 齊藤 裕樹¹ 戸辺 義人³ 鉄谷 信二¹

概要: GPS 機能を備えた携帯端末等による位置情報取得技術の普及により、人々の訪れた場所同士を結ぶ行動履歴を解析する研究が活発に行われている。経路情報の蓄積によって得られる行動履歴は、人々の興味と移動との関係を解析する手段として注目されている。一方、Twitter をはじめとするマイクロブログはテキストと共に GPS センサによる位置情報を付与し、公開することが可能なサービスとして世界的に利用されている。本論文では、マイクロブログ上に蓄積された人々の行動履歴を基に利用者の行動とコンテキストを解析し確率過程モデルに適用させることで、将来の利用者の行動を予測する手法の提案を行う。また、マイクロブログの実データを用いた行動予測精度の評価実験より、単純統計を用いた従来手法より高い予測精度が得られることを確認した。

Application and Evaluation of Stochastic Model for People Behavioral Prediction Using Contexts of Microblog Services

TSUBASA TAKAYAMA¹ KEI YAMAGAMI² HIROKI SAITO¹ YOSHITO TOBE³ NOBUJI TETSUTANI¹

Abstract: The advance of GPS-enabled portable devices such as PDAs and smart phones facilitates people to record their location histories. Location trajectories imply human behaviors and preferences related for their interests. On the other hand, microblog services such as Twitter enable us to publish text messages (e.g. Tweets) and location-tags (e.g. Geo-tags) to subscribers. This paper proposes a schema for predicting user behavior by analyzing location trajectories and contexts by applying a stochastic model. And, we confirm the effectiveness of our schema through experiment using the actual data obtained from microblog service.

1. はじめに

GPS 機能を備えた携帯端末等による位置取得技術の普及により、人々が訪れる場所同士を結ぶ行動履歴を解析する研究が活発に行われている [1], [2], [3], [4]。経路情報を蓄積することによって得られる行動履歴の解析は、人々の興味と移動との関係を抽出する手段として注目が集まっている。一方、Twitter をはじめとするマイクロブログサー

ビスは、テキストとともに GPS センサによる位置情報を付与し公開することが可能なサービスとして世界的に利用されている。

本論文では、マイクロブログサービスで発信されるテキストと位置情報から、テキストと移動の関係を確率過程モデルに適用し行動解析を行う。また、解析により得られた情報を蓄積し行動モデル化を行い、利用者の未来の行動を予測する手法を提案をする。提案手法の有効性を示すために Twitter 上の実データを用いた行動予測精度の評価実験の結果、単純な行動履歴から求めたものよりも有効性が高いことを示す。

本論文の構成は次のとおりである。まず第 1 章に引き続き、第 2 章でマイクロブログサービスを用いた行動予測の概要について示す。次に第 3 章で提案の中核となる確率過

¹ 東京電機大学未来科学部情報メディア学科
Department of Information Systems and Multimedia Design,
Tokyo Denki University

² 東京電機大学大学院未来科学研究科情報メディア専攻
Graduate School of Science and Technology for Future Life,
Tokyo Denki University

³ 青山学院大学理工学部情報テクノロジー学科
Department of Integrated Information Technology,
Aoyama Gakuin University

程モデルへの適用, 移動確率計算手法, 行動予測手法について述べる. 第4章では, 提案手法の有効性を確認するために行ったマイクロブログサービスの実データを用いた評価実験とその結果について述べる. さらに第5章では, 関連研究と本研究の位置づけを示し, 第6章では, 本論文の内容をまとめ結論づける.

2. マイクロブログサービスを用いた行動予測

GPS機能を備えた携帯端末の普及により, 多くの位置情報サービス (Location Based Service: LBS) が iPhone や Android といった様々なプラットフォーム上で実現されている. 初期の位置情報サービスでは, 地図上でのナビゲーションを行うことや, 位置を指定した最近傍の施設やサービスの検索などが提供されたが, 位置情報は利用者の携帯端末内で用いられるのみで, 他の利用者と共有は行われていなかった. しかし, 最近の位置情報サービスでは, 人やものの位置情報を蓄積し利用者間で情報共有することで, 位置情報に新たな価値を生み出す動きが活発に行われている. また, GPSによって取得された人々が訪れた場所同士を結んだ行動履歴を蓄積することで, 街中の人々の動線を解析することやナビゲーションシステムなどへの応用が期待されている.

一方, Twitterをはじめとするマイクロブログサービスは, 短いテキストを発信するサービスであり, その場の状況に関するメッセージを即座に発信できる特徴がある. 特に, スマートフォン等の携帯端末上で用いられることにより, 即時性と臨場感の高い情報発信が可能である. また, テキストを発信する際に GPSによる位置情報を付加し, 実世界と直接リンクした情報を扱うことも可能である. このように, 人がセンサとなり現在の状況や行動に関するテキストを位置情報と共に情報発信すること可能である. これにより, 人々の日常活動や目標指向の行動を知ることが可能になることから, 新たな価値のある情報を作り出す研究が活発である [5], [6].

本研究では, まず人々がマイクロブログサービスに発信したメッセージと位置情報から, 発言コンテキストと移動経路を抽出し行動履歴として蓄積を行う. 蓄積した行動履歴を基に利用者のコンテキストから未来の行動を予測する手法を検討する. 具体的には, 人の行動において, 移動先は移動元の情報のみに影響されるという仮定に基づき, 人の行動にマルコフ連鎖を適用し行動モデルを作成する. また, 人々がマイクロブログ上で発信するメッセージと行動には一定の関係があるとの仮定に基づき, メッセージと行動の組み合わせを条件付確率による定式化を行う. これらにより, 利用者の未来の行動予測を行うものである.

次に, 予測の結果から利用者に次の行動に対する適切な情報提示を行う推薦システムの提案を行う. 図1に提案システムのシステムモデルを示す. 本システムは, 二つの機

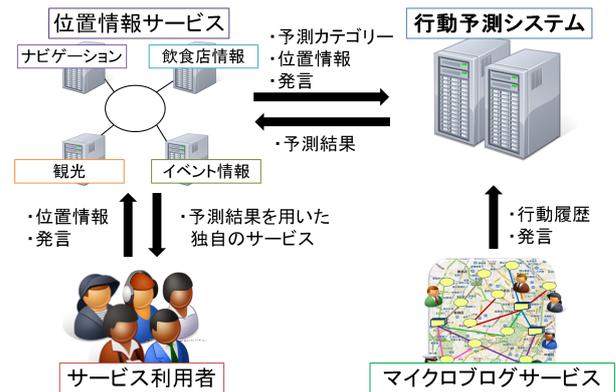


図1 マイクロブログのメッセージと位置情報を用いた行動予測システム

Fig. 1 User Behavior Predicting System among Text Messages and Location Histories in Microblog Services.

能を核として設計する. 一つ目は, 行動履歴をマイクロブログサービスから取得し, 行動解析を行い蓄積する機能である. 二つ目は, 利用者の発言と位置情報から行動予測を行う機能である.

3. マイクロブログの行動履歴を用いた行動予測手法

本章では, まず行動推定の基本的方針を述べ, 人の行動に対して確率過程モデルの適用を行い, 移動確率の計算手法と行動予測手法について述べる.

3.1 行動予測の基本方針

マイクロブログは, その場その瞬間に見たことや感じたことを気軽に短いテキストとして発信できることから, 利用者の興味関心を反映するメディアであると言える. マイクロブログの発言に付与された位置情報により, 利用者が興味関心を持った場所を知ることができる. また, 利用者の位置情報を時間軸上に並べることにより移動経路を抽出する. さらに, 移動中に行った発言を集約することで, 利用者の行動と移動の意味や目的を知ることができる. 例えば, 鎌倉を観光で訪れた人々の場合, 鶴岡八幡宮や建長寺, 高徳院などの位置情報とともに, 鶴岡八幡宮でお参りしたことや, 高徳院で鎌倉大仏を見たこと, 建長寺の半僧坊からの景色に関する発言をした可能性が高いと考えられる. このような行動履歴を多くの人から集約することにより, 移動経路とその移動の意味を知ることが可能であると考えられる. そこで, 本論文では, 過去に蓄積された行動履歴に対して, 確率過程モデルを適用することにより行動解析を行い, 移動経路と発言を考慮した行動モデルを作成する. また, 作成した行動モデルを用いて未来の利用者の行動を予測する手法を提案する.

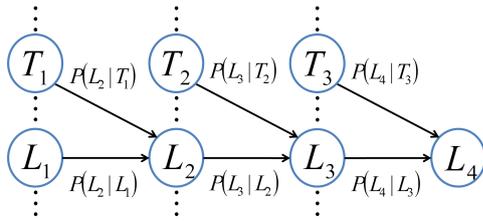


図 2 移動と発言コンテキストの状態遷移モデル

Fig. 2 State Transition Model based on Movements and Contexts.

提案手法は、具体的には以下のような方針に基づく。ある移動経路を通った人々のうち特定の発言をした人の割合を、条件付確率で計算すると、移動経路に対して関係の強いキーワードが高い割合を占めることとなる。このことを元に、発言から未来の移動経路の確率を求める。例えば、過去に建長寺から若宮大路を経て鶴岡八幡宮に行くルートを通った人は、高い割合で歴史的な建物や風情や景色に関する発言をしたと考えられる。このような割合を元に、現在建長寺で風情や景色を楽しんでいる発言をする人に対して、推定される目的地である鶴岡八幡宮の観光情報を提示することや目的地へのナビゲーションを行う。

以下、確率過程を用いた行動のモデル化および行動予測の手法について説明する。

3.2 人の行動の確率過程モデルへの適用

利用者は GPS 機能を備えた携帯端末を持ち、マイクロブログサービスへのメッセージ送信とともに「現在地の座標-タイムスタンプ」の組合せ $(L_k, time_k)$ を記録する。また、発信したメッセージは、単語ごとに分解し重み付けし、ベクトル空間上に表す。この情報とタイムスタンプの組合せ $(T_k, time_k)$ を記録する。システムでは、この情報を蓄積しユーザごとに移動に伴う位置 L の時系列 L_1, L_2, \dots, L_n 、および発言コンテキストの時系列 T_1, T_2, \dots, T_n を保持する。

人の行動において、位置 L_k から位置 L_{k+1} への移動は、蓄積された過去の移動情報だけに依存していると仮定すると、2 地点間の移動確率はマルコフ過程によって表すことが可能である。3 地点間以上移動も 2 地点間の移動の連なりとして表し、多地点の行動もマルコフ連鎖により表すことが可能である。したがって、 N 地点間の移動確率 $P(L_1, \dots, L_N)$ は式 (1) で得ることができる。

$$P(L_1, \dots, L_N) = P(L_1) \prod_{n=2}^N P(L_n|L_{n-1}) \quad (1)$$

次に、位置 L_k で発生した発言コンテキスト T_k について考える。図 2 は、滞在地点を 1 つの状態と見なし移動を状態遷移として表した図である。位置 L_k で発生した発言コンテキスト T_k は、位置 L_{k+1} への移動に関係性があり、移動確率に影響すると仮定する。これにより、 N 地点間の位置 L と発言コンテキスト T の同時分布確率は式 (2) で得

ることができる。

$$P(L_1, \dots, L_N, T_1, \dots, T_N) = P(L_1) \left[\prod_{n=2}^N P(L_n|L_{n-1}) \right] \prod_{n=2}^N P(L_n|T_{n-1}) \quad (2)$$

3.3 条件付確率を用いた 2 地点間移動確率の計算

次に、これまでの行動履歴が既知であり、現在地から次の移動先への移動確率を求める方法を検討する。前節で得られた式 (2) を漸化式で表すと式 (3) のようになる。

$$P(L_1, \dots, L_N, T_1, \dots, T_N) = P(L_1, \dots, L_{N-1}, T_1, \dots, T_{N-1}) P(L_N|L_{N-1}) P(T_N|T_{N-1}) \quad (3)$$

このことから、これまでの行動履歴が既知であり、位置 L_k 上に存在する利用者の位置 L_{k+1} への移動確率は、式 (4) によって得られる。

$$P(L_{k+1}|L_k, T_k) = P(L_{k+1}|L_k) P(L_{k+1}|T_k) \quad (4)$$

式 (4) のうち、位置 L_k を条件とした位置 L_{k+1} への移動確率 $P(L_{k+1}|L_k)$ と、発言コンテキスト T_k を条件とした位置 L_{k+1} への移動確率 $P(L_{k+1}|T_k)$ をそれぞれ求める方法を以下に示す。

まず、位置 L_k を条件とした位置 L_{k+1} への移動確率は、位置 L_{k+1} が位置 L_k の状態に依存するため、式 (5) のとおりベイズの定理によって求めることができる。

$$P(L_{k+1}|L_k) = \frac{P(L_{k+1}) P(L_k|L_{k+1})}{P(L_k)} \quad (5)$$

次に、発言コンテキスト T_k は、実際には複数の単語の集合 $message_k = \{word_1, \dots, word_n\}$ で構成される。本研究では、複数の単語の生起確率は独立したものとし bag-of-word の手法をとる。よって、位置 L_{k+1} と発言コンテキスト T_k を構成する $\{word_1, \dots, word_n\}$ の条件付確率は、単純ベイズ確率モデルによって計算することが可能である。以上から、条件付確率 $P(L_{k+1}|T_k)$ は式 (6) で得ることができる。

$$P(L_{k+1}|T_k) = P(L_{k+1}) \prod_{i=1}^n P(word_i|L_{k+1}) \quad (6)$$

3.4 人の移動先の予測手法

前節の移動確率計算を行い、行動履歴 $\{L_k, T_k\}$ から最も移動する可能性の高い地点を選択する。そのために、すべての位置 $\{l_1, \dots, l_n\}$ についての移動確率を計算し、式 (7) に示す移動確率行列 $M_{L_k \rightarrow L_{k+1}}$ を生成する。

$$M_{L_k \rightarrow L_{k+1}} = \begin{pmatrix} P(l_1|L_k, T_k) \\ P(l_2|L_k, T_k) \\ \vdots \\ P(l_n|L_k, T_k) \end{pmatrix} \quad (7)$$



図 3 地図上にプロットしたジオタグ付き発言*1
Fig. 3 Tweets with Geo-tags plotted on map.*1

式 (7) より位置 $\{l_1, \dots, l_n\}$ への移動確率を得た上で、最も移動確率の高い位置を選択し、予測結果とする。

また、複数地点の一連の移動を1つの状態と見なすことにより、経路と経路上の発言コンテキスト $\{L_1, \dots, L_k, T_1, \dots, T_k\}$ を用いた行動予測が考えられる。これを式 (7) に適用すると、移動確率行列 $M_{\{L_1, \dots, L_k\} \rightarrow L_{k+1}}$ は以下のとおりとなる。

$$M_{\{L_1, \dots, L_k\} \rightarrow L_{k+1}} = \begin{pmatrix} P(l_1 | L_1, \dots, L_k, T_1, \dots, T_k) \\ P(l_2 | L_1, \dots, L_k, T_1, \dots, T_k) \\ \vdots \\ P(l_n | L_1, \dots, L_k, T_1, \dots, T_k) \end{pmatrix} \quad (8)$$

同様にして、位置 $\{l_1, \dots, l_n\}$ への移動確率のうち最も移動確率の高い位置が予測結果となる。

4. 行動予測手法の評価実験

提案手法の有効性を示すために評価実験を行った。本実験では、マイクロログサービスより得た行動履歴を基に、従来手法と提案手法で予測した移動先について、正解である実際に利用者のとった行動と比較することで正答率を求める。

4.1 実験データセット

本実験では、マイクロログサービスである Twitter より、神奈川県鎌倉市内の 9km 四方の空間でアップロードされた位置情報付き発言の取得を行った。取得には Twitter の StreamingAPI を用いた。この API は、実世界の緯度、経度を指定することにより、指定領域内でアップロードさ

*1 地図データ ©2013 Google, ZENRIN.

れた発言を取得することができる。発言は、テキスト、位置情報、ユーザ ID、タイムスタンプの4つで構成される。本実験では、2012年7月6日から2013年1月25日までに56583件の発言を取得した。取得した発言を MeCab[7] を用いて形態素解析を行い、名詞と固有名詞を抽出し予測に用いる。本実験では、280284個の単語を抽出した。同一ユーザが1日で位置情報を付与した発言を複数行った場合、移動を行ったと見なす。これらにより、43210件分の移動を伴う発言を抽出した。発言のあった位置を地図上に点画したものを図3に示す。

4.2 評価実験方法

まず、取得範囲を格子状に分割し、位置 L を設定する。本実験では、一つの格子を 100m 四方から 1000m 四方まで 100m 刻みで分割した 10 とおりの実験データセットを作成した。分割した一領域あたりの移動を伴う発言の平均回数と、平均単語出現数を表1に示す。以下、一つの格子領域の大きさを変化させたときの正答率の変化を確認する。

実験データセットからランダムに抽出した 90% のデータを学習データとし、残り 10% のデータは正解データとし評価を行う。正答率の比較には、以下の従来手法と 2 とおりの提案手法の結果を用いた。

- 従来手法
発言を考慮しない単純統計による予測結果を求める。この手法では、ある位置 L_k からみた次の移動先 L_{k+1} を、過去に人々がとった行動の人数比のみに基づいて計算を行う。
- 提案手法

表 1 格子領域サイズと取得されたデータ数

Table 1 Grid area size and obtained data.

1 領域あたりの面積 (m^2)	平均移動数	平均単語数
1×10^4	3.36	15.92
4×10^4	13.43	63.69
9×10^4	30.21	143.30
16×10^4	54.41	254.56
25×10^4	83.91	398.07
36×10^4	120.83	573.22
49×10^4	175.14	818.67
64×10^4	217.64	1018.25
81×10^4	271.87	1289.74
100×10^4	335.64	1592.27

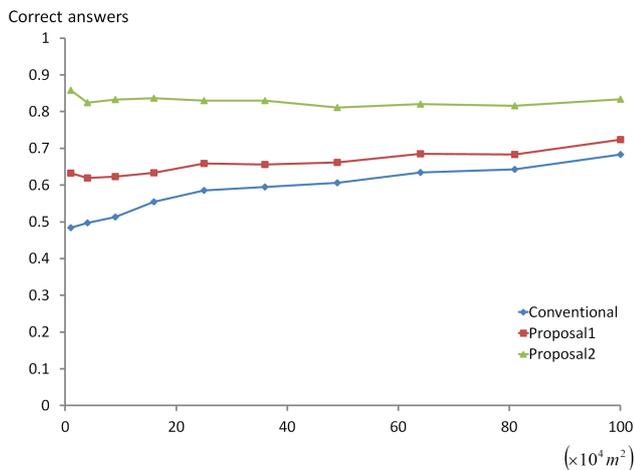


図 4 格子領域面積と正答率の比較

Fig. 4 Comparison of grid area and correct answers.

次の移動先 L_{k+1} について、1 地点前のみの行動履歴 $\{L_k, T_k\}$ を条件とした確率計算を行う手法と、2 地点前の行動履歴 $\{L_{k-1}, L_k, T_{k-1}, T_k\}$ を条件とした計算を行う手法を用いる。

4.3 結果と考察

単純統計による従来手法の正答率 (Conventional) と、提案手法のうち 1 地点前の行動から予測した正答率 (Proposal1)、および 2 地点前の行動から予測した正答率 (Proposal2) を図 4 に示す。格子領域の面積が増えるに従い、すべての手法の正答率が上がることが分かる。これは、一領域あたりの情報量が増え学習の精度が高まったからと考えられる。また、従来手法と比べ、1 地点前の行動からの予測結果の精度が高く、さらに 2 地点前の行動からの予測精度が高いことから、予測に用いる情報量に応じて精度が高まること分かる。

次に、鎌倉駅を含む一つの領域からの 2 つの移動先について、学習した単語の特徴の比較を行う。2 つの移動先の単語ごとのベクトルの比較を図 5 に示す。移動先 1 (Destination1) では「八幡宮」「鶴岡」などが高い値を示すのに対

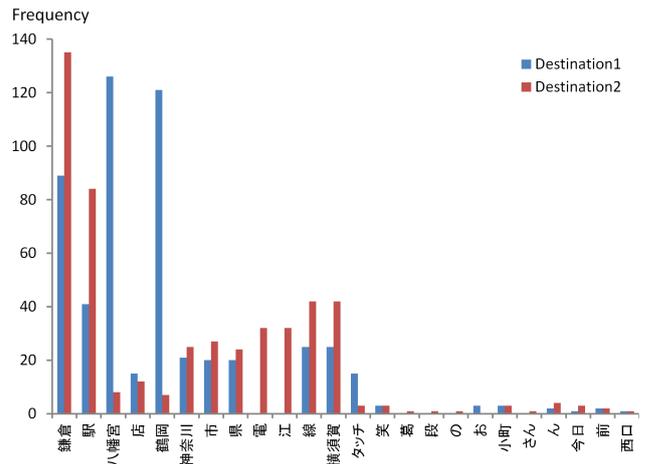


図 5 異なる移動先の単語ごとの頻度分布

Fig. 5 Distinations and frequency distribution of words.

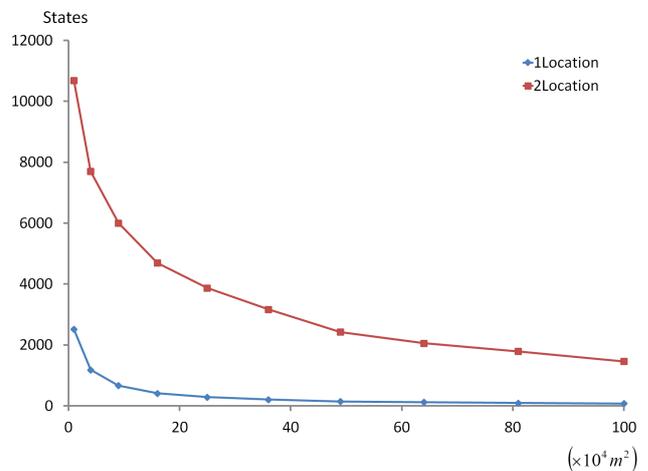


図 6 格子サイズ、地点数と状態数

Fig. 6 Grid size, locations and states.

して、移動先 2 (Destination2) では「鎌倉」「江」「電」などが高い値を示すことが分かる。地図上では、移動先 1 は鶴岡八幡宮を含む領域であり、移動先 2 は江の電沿いの領域であることから、移動に関するコンテキストが学習データに反映されていると考えられる。

提案手法は、予測対象とする領域の格子領域上で識別される行動履歴ごとに状態を保持し、状態ごとにベイズ分類器を生成し処理を行う。領域計算量はベイズ分類器の数に比例するため、一格子あたりの面積と行動履歴に用いる地点数を変化させ、システム上で保持される状態数の評価を行った。図 6 は、格子サイズに対する状態数の変化を 1 地点前の情報のみを用いた場合と、2 地点前の情報を用いた場合について示したものである。格子サイズが小さいと、短い距離で行われる発言でも識別されるため行動履歴の識別数の増加により状態数が多くなり、格子サイズが大きいと、離れた地点の発言も同一の行動履歴と見なされることが増えるため、識別数が減少し状態数が少なくなる。一方、行動予測に用いる地点数が 2 地点の場合、取りうる 2 地点

の組み合わせの数の増加に従い1地点のものより多くなった。2地点の際の状態数は、最大1地点の19倍であるが、累乗分の増加とはならないことが分かる。これは、人々の移動する経路は、隣接する場所や道路上をとることが多く、組み合わせが限定されたためと考えられる。

5. 関連研究

GPSなどの位置計測技術の発展により、位置とタイムスタンプからなる位置履歴を記録することが可能である。この位置履歴を、ユーザの興味や趣向に関する手がかりと見なして分析する研究が行われている。Countsら[1]は、センサによって取得された情報に位置タグをつけ複数ユーザで共有するシステムを提案している。平井ら[2]は、GPSナビゲーションツールを利用して目的地へたどり着く過程を分析し、実世界との知覚的な相互関係を分析している。また、位置履歴からユーザの行動を分析し複数のユーザ間のソーシャルメディアを提案する研究[3]や、データマイニングによってユーザ間のコミュニティを抽出する研究[4]などが存在する。しかし、本研究では、位置の履歴だけを分析するのではなく、マイクロブログ上で人々が送信した発言とジオタグを蓄積し、位置履歴だけでなく同時に発言から得られるコンテキストを用いる。これにより、行動を起こす背景や目的を含めた分析を行うことが可能である。

また、ブログやSNS、Twitterなどは、その瞬間に見たことや感じたことを気軽にWeb上に公開し共有するなど、世の中の関心を即時的に反映する新たなメディアであることから、これらから有用な実世界情報の抽出を目的とした研究が行われている。Hirutaら[8]は、場所誘因型ジオタグ付きツイートを抽出するため、位置情報メタデータと、発言内容と位置情報の対応付けを解析する手法を提案している。Leeら[9]は、Twitter上の群衆の振る舞いをクラスタ化することにより、実世界のイベントを抽出する試みを提案している。一方、言語解析的なアプローチとしては、ソーシャルメディアでの複数の発言状況から注目する出来事の発生を推定するための、テキストの類似度に基づく情報伝搬経路解析手法[5]や、Twitterの複数発言間の構造をHITSアルゴリズムにより解析し、authorityとhubの2つの尺度を用いて情報取得に有用なユーザを抽出する手法[6]などが存在する。Twitterは発言の文字数が短く、目的のトピックを示す単語が省略されていることや、リツイート元の内容を前提とした発言がされていることが多いため、発言間の情報伝搬を展開してコンテキストを解析することは今後の課題と言える。

6. おわりに

本論文では、マイクロブログサービスから過去の発言と位置情報を取得し、行動解析を行うことにより未来の利用者の行動を予測する手法について提案した。提案手法では、

人の行動に確率過程モデルを適用し解析した統計情報を蓄積することで、未来の利用者の行動予測を行うものである。また、マイクロブログサービスであるTwitterのデータを用いて、提案手法の評価実験を行った。

今後は、移動と発言の関係性について更に研究を行い、予測精度の向上を目指す。また、提案手法では、発言単語のみを考慮しており、発言の意味を考えた予測を行っていない。観光や買い物といった目的別の学習をすることにより、利用者の目的を考慮した予測を行うことが可能となると考えられる。さらに、予測結果を反映させた推薦システムの設計、実装を行う予定である。

謝辞 本研究は、JSPS科研費 若手研究(B) 課題番号24700074の助成を受けたものである。

参考文献

- [1] Counts, S. and Smith, M.: Where were we: communities for sharing space-time trails, *Proceedings of the 15th annual ACM international symposium on Advances in geographic information systems (GIS 2007)*, pp. 1-8 (2007).
- [2] 平井浩将, 森 傑: 経路探索におけるGPSナビゲーションツールの利用とアクション生起との関係-都市空間におけるアクティビティに関する研究-, 日本都市計画学会都市計画論文集, No. 42-3, pp. 541-546 (2007).
- [3] Zheng, Y., Xie, X. and Ma, W.: GeoLife: A Collaborative Social Networking Service among User, location and trajectory, *IEEE Data Engineering Bulletin*, Vol. 33(2), pp. 32-39 (2010).
- [4] Hung, C.-C., Chang, C.-W. and Peng, W.-C.: Mining trajectory profiles for discovering user communities, *Proceedings of the 2009 ACM International Workshop on Location Based Social Networks (LBSN 2009)*, pp. 1-8 (2009).
- [5] Kim, J. W., Candan, K. S. and Tatemura, J.: Efficient overlap and content reuse detection in blogs and online news articles, *Proceedings of the 18th ACM International Conference on World Wide Web (WWW 2009)*, pp. 81-90 (2009).
- [6] 田中淳史, 田島敬史: twitterのツイートに関する分類手法の提案, 日本データベース学会他 第2回データ工学と情報マネジメントに関するフォーラム (DEIM2010) A5-4 (2007).
- [7] *MeCab: Yet Another Part-of-Speech and Morphological Analyzer*, <http://mecab.googlecode.com/svn/trunk/mecab/doc/index.html>.
- [8] Hiruta, S., Yonezawa, T., Jurmu, M. and Tokuda, H.: Detection, classification and visualization of place-triggered geotagged tweets, *Proceedings of the 2012 ACM Conference on Ubiquitous Computing, UbiComp '12*, pp. 956-963 (2012).
- [9] Lee, R. and Sumiya, K.: Measuring geographical regularities of crowd behaviors for Twitter-based geo-social event detection, *Proceedings of the 2nd ACM SIGSPATIAL International Workshop on Location Based Social Networks, LBSN '10*, pp. 1-10 (2010).