

Twitter での Retweet 情報を利用した情報拡散予測

Predicting Information Diffusion in Twitter through Retweeting Behaviors

村上 明子^{*1*2} 鈴木 秀幸^{*3}
Akiko Murakami Hideyuki Suzuki

^{*1}日本アイ・ビー・エム株式会社 東京基礎研究所
IBM Research - Tokyo

^{*2}東京大学大学院 学際情報学府
Graduate School of Interdisciplinary Information Studies, The University of Tokyo

^{*3}東京大学 生産技術研究所
Institute of Industrial Science, The University of Tokyo

In social media, resharing contents is one important way for users to propagate information to the other users. In this paper, we investigate retweeting behaviors of Twitter users; our investigation shows that retweet activities are highly affected by their preferences, such as topics of interest. Based on this investigation, we propose a probabilistic model for retweeting behavior that depends on tweet contents and other information. This probabilistic model can be used with information diffusion models for predicting information diffusion.

1. はじめに

近年、インターネット上でのコミュニケーションは Twitter や Facebook などリアルタイムに情報や体験を共有する形に変化してきている。一般的にソーシャルメディアでは他の人から来た情報の一部または全部を再共有するという形で情報が伝播していくが、リアルタイム型は興味が新鮮なうちに情報を共有相手に伝えることができるため、より高速に情報が伝播したり、大きな範囲で情報が広がって行くと言える。

この情報の拡散は、バイラル・マーケティングに代表されるような“良い”意味での拡散だけではなく、間違っている情報であったり、風評や企業の悪い評判といった“悪い”情報も拡散してしまうことが多々ある。2011年の3月に起こった東北地方太平洋沖地震では、多くのデマや未確認情報が拡散されていることが観測された。デマや風評を打ち消すためにはいち早く拡散すると思われる内容を把握し、それを否定するための情報を複数のメディアにうたななければいけない。そのため情報の拡散を早期に把握することが重要であるといえる。

拡散の予測は情報を媒介するユーザーをノードとし、情報の共有元から供給対象へ向きのあるリンクを張った情報の伝播経路のネットワークでモデル化をすることができる。Twitterの場合は、ユーザーからフォロワーにリンクを張ったフォロワーネットワークがこれにあたる。このネットワークを用いて、各々のノードが情報を伝達するか(再共有するか)どうかを確率的なモデルで捉えることで、このネットワーク上の情報の拡散を推測することができる。

boydらの研究[1]では、Twitterの返答や再共有といった行動は通常の会話の要素が含まれているとし、発言に同意したり、内容を肯定的あるいは否定的な自分の意見と共に広めたいという要求によって情報の再共有を行っているとして述べている。本稿ではTwitterにおける実際のRetweetによって、再共有がコンテンツならびにユーザーの行動に依存したものであるか検証し、情報拡散のモデルにおける情報の再共有の確率モデルを提案する。

2. 関連研究

ソーシャルメディア上での情報伝播の研究ではその人気度とAPIなどの充実度からTwitterが対象とされることが多い。Twitterは発言をTweet、情報の再共有をRetweet(RTと略される)と呼び、Retweetにはシステムがサポートをしている公式RTと、ユーザーが自主的に行っている非公式RTがある。非公式RTは共有者の発言とみなされ、コメントを入れることも、発言の内容を変更することもできる。一方、2009年より提供されている公式なRetweetは改変されない元の発言のまま誰がRetweetしたかの情報と共に発信される。2010年までの研究は非公式RTの内容を解析し構造化したものが主であったが、近年はRTの情報は公式のAPIでもサポートされているため、近年では公式RTを分析した研究も見受けられる。

Retweetの行動からユーザーの嗜好を推定する太田ら[2]の研究では、Retweetしている対象をその人の好む話題とみなし、多くRTを共有しているユーザーほど興味が似ているとしてRetweetを視覚化し、フォロー対象者を発見するための支援を行っている。またZamanら[3]はユーザーごとにどのTweetをRetweetしているかという情報を集め、Retweetするかどうかを確率的協調フィルタリングモデルで推定している。さらに過去にRetweetをした時間帯や、誰の内容をリツイートしていたか、過去の自分の発言とどれくらい内容が近いかなどの指標を用いてある発言をRetweetするか否かを判定する研究もある[4]。

Twitter上での情報の拡散の分析についても多くの研究がある。風間ら[5]は返信、Retweetによって関連するTweetをグループ化し、その各グループ内での情報の流れを有向グラフとして可視化している。内容の類似している返答によって情報伝播があったとみなし、情報伝播の速度や広がりの大きさなどを推測している研究[6]や、同じ話題を取り上げる際に使われるハッシュタグに着目して、どのように情報が広がるのかをみた研究[7]もある。

本研究で扱うようなネットワーク上の情報拡散は、感染症の

感染伝播と同様のモデル化が可能である。情報拡散の時間変化を考えずに最終的に情報が拡散する範囲だけを考える場合は、ボンドパーコレーションによるモデル化が可能である。これらの理論的性質は当初格子のような単純なネットワークで研究されてきたが、近年発展が著しい複雑ネットワークの分野においても調べられている。また、Twitter のフォロワーネットワークのような有向グラフ上のボンドパーコレーションについても研究されている [8]。また情報拡散の時間変化を考える場合には、感染症モデルである（ネットワーク上の）SIR モデルによるモデル化も提案されている [9]。

3. Twitter 上での情報拡散モデル

感染症のダイナミクスを記述する数理モデルである SIR モデルでは、ネットワーク上の各頂点が感受性 (S)、感染性 (I)、免疫獲得 (R) の 3 状態のいずれかを取る。情報拡散のモデルとしては、情報を知らない状態が未感染で免疫を持たない状態 S、情報を受け取って他人に話す状態が感染状態 I、そして噂は知っているが他人に話すことがなくなった（興味を失った）状態が免疫獲得状態 R であると考えることができる。状態 S が状態 I の頂点と隣接しているとき、状態 S の頂点は状態 I の頂点から一定の確率で情報を得て状態 I へと遷移（感染）する。また、状態 I の頂点は、一定の確率で興味を失い、状態 R へと遷移する。

この情報の拡散モデルを Twitter に当てはめて考えることもできる。情報を受け取り他の人に再共有するノードが I であり、再共有が S と I の接触と考えてよい。通常の SIR モデルでは隣接ノードに感染するかどうかは確率的に考えるため、ユーザーが情報を再共有するかどうか、確率的なモデルで表すことが必要となる。

しかし、Twitter での Retweet 行動に対して単純に SIR モデルを適用するのは適切ではないと考えられる。Twitter において情報を受け取ったあと、他の人に情報を伝える状態（感染状態 I）になるか、情報を知っているが発信しない状態（免疫獲得状態 R）か、いずれかに変化すると考え、これを Twitter における情報拡散として考える。これはつまり、自分がフォローしているユーザーから情報を受け取った後、それを Retweet するかしないか、という 2 つの状態と対応する。つまり、Retweet するという行為は感染状態であり、しないというのは免疫獲得状態であるといえる。そこで、受け取った情報を Retweet するか否かが、情報の拡散モデルにおける感染確率であるといえる。

本稿では、Retweet という行動がユーザーの興味をモデル化する一つの手がかりになると考える。次の章では、その仮説が正しいか実際にユーザーが Retweet した Tweet を検証する。

4. ユーザーごとの Retweet の特徴

ユーザーの Retweet が人の興味に依存しているのならば、ユーザーごとに Retweet に明らかな特徴があるといえる。それを検証するため、本章では Twitter 上での実際のデータを元に検証する。検証のため、2 人のユーザーのタイムライン（自分がフォローしている人の Tweet がリアルタイムで表示される場所）を 2012 年 1 月から 3 月まで取得した。その中には、ユーザーが公式に Retweet したものも含まれる。

人による Retweet の傾向の違いを見るため、ユーザー 2 人が Retweet した Tweet から名詞を抽出した。抽出された名詞のうち頻度上位 10 件を表 1 に示す。調査の対象とした Retweet の数はユーザー A が 189 件、ユーザー B が 185 件である。

表 1: Retweet における一般名詞の頻度上位 10 件

順位	ユーザー A		ユーザー B	
1	データ	13	人	22
2	人	10	日本	10
3	論文	8	自分	10
4	話	8	IBM	10
5	IBM	8	女性	8
6	現在	7	Google	8
7	更新	7	今日	7
8	情報	6	研究	7
9	日本	6	仕事	7
10	原発	6	日本語	5

2 人の興味の違いはユーザー A に「原発」、ユーザー B に「女性」という語から少し読み取れるが、そのほかは一般的な単語であったり、2 人に共通に出現する「IBM」や「研究」に関連する語であり、頻度だけでは分別するのは難しいといえる。

より強くユーザーの興味をモデル化するため、ユーザーの目にするタイムラインの Tweet と、Retweet を行った Tweet との比較を行う。すべての Tweet の集合と Retweet した Tweet の集合にわけ、名詞の中でより Retweet の集合により相関が高いものを抽出する。相関が高い名詞の上位 10 位を相関値^{*1}と共に表 2 に示す。

表 2: Retweet における一般名詞の相関上位 10 件

順位	ユーザー A		ユーザー B	
1	放射線モニター	17.3	有名人	9.9
2	(URL ^{*2})	17.3	及川	7.1
3	停電	10.1	ロマンス	5.0
4	解釈	9.7	(URL ^{*3})	5.0
5	三部作	9.7	オープンハウス	4.16
6	等高線	9.7	女性	3.1
7	減少傾向	9.7	グローバル	3.0
8	全般	9.7	Google	2.5
9	各地	8.9	デモ	2.4
10	放射線量	7.0	梓	2.4

ユーザー A は原発・放射線関係に、ユーザー B は就職活動や採用について多く Retweet していることがこのリストより分かる。ユーザーによる Retweet の嗜好が、Retweet した Tweet とタイムラインに現れる語を比較することによってより正しくモデル化できるといえるだろう。

*1 ここで示している相関値は IBM Content Analytics によって算出された値である。一般的な相関値とは違い、全体での頻度が低いものに高い相関値が割り振られるのを防ぐため、頻度の少ない単語に対しては値を低く見積もっている。http://www-01.ibm.com/software/ecm/content-analytics/bundle.html

*2 放射線関係の URL

*3 就職活動関係の URL

5. Retweet の確率モデル

前章で実際の Retweet のデータじゃら、ユーザーごとに Retweet する内容に特徴があることがわかった。この特徴を用いて、ユーザーがある Tweet を Retweet するか否かを示すモデルを作成する。

ある情報の再共有を行うためには、下記の 2 つの条件がある。

1. 共有された情報を受信し、再共有可能な状態となること
2. 共有された情報を、自分のフォロワーに対して再共有しようと思うこと

最初の条件は、情報が自分に向けて発信されたときにリアルタイムにその情報を受信する必要があるということを行っている。Twitter においては、自分のタイムラインに対象の Tweet が来たときに、Twitter を利用しているかどうかということと等しい。

2 つ目の条件で示されている再共有の理由は、面白い内容だから、あるいはフォロワーが興味を持ちそうだから、など人それぞれ異なる。その結果として再共有した内容にユーザーごとの特徴が出ることとなる。

これらのことから、ある Tweet の Retweet されやすさを示す、情報の再共有の確率モデルは、ユーザーとその Tweet の組み合わせで決定される。そのため、まずはユーザーの興味をモデル化し、それをを用いてある Tweet がどれくらい再共有されやすいかを、計算する必要がある。

5.1 ユーザーの興味のモデル化

ユーザーがどのような内容の Tweet を Retweet しやすいかは、過去にそのユーザーが Retweet した情報でモデル化することができる。これはユーザーのある情報 (Tweet) に対する感度といえる。この感度は、ユーザーの行動 (Retweet) に対して重みを考慮する。ここでいう行動に基づく重みは、Retweet された発言に含まれる単語が重視されるようにすると同時に、Retweet されていない発言に含まれる単語に対しても重みを考慮したものではなくてはならない。そのため、読まれたが Retweet されなかったものと、そもそも読まれずに Retweet をする判断を下されなかったものとを区別する。

また、実際のデータで示したように、単に Retweet した内容だけではなく、Retweet しなかったデータとの比較で、よりユーザーの特徴を得ることができるといえる。そのため、ユーザーの興味をモデル化するステップは以下のように考える。

1. ユーザー u のタイムライン (フォローしている人の Tweet) を集める。各 Tweet を m とする
2. 各々の Tweet から単語を抽出する。抽出された単語のベクトルを $f_j^{(u,m)}$ とする。
3. Tweet のユーザーの行動に基づき、各発言 m に対する重み $r^{(u,m)}$ を決める
 - Retweet した Tweet は $r^{(u,m)} = 1$
 - Retweet しなかった Tweet で、読んで RT しなかったものは $r^{(u,m)} = 0$ 、読んでいないものは $0 < r^{(u,m)} < 1$
4. 発言全体の頻度と、行動による重み付きの頻度の比で、単語ごとの重みを決定する

Retweet をされないときの $r^{(u,m)}$ は、その Tweet が読まれたか読まれていないかに依存する。Tweet が読まれたかどうかは直接測定することはできないので、ユーザーのその他の行動 (Tweet) や、それまでの行動パターン (Twitter 上での活動のサイクルなど) によって推定して、値を決定する。

したがって、ユーザー u のモデル化された興味は単語の重みつきベクトル $f_j^{(u)}$ として以下のように表される。 j は単語の ID, c は定数とする。

$$f_j^{(u)} = \frac{\sum_m r^{(u,m)} f_j^{(u,m)}}{\sum_m f_j^{(u,m)} + c} \quad (1)$$

5.2 ある Tweet に対する Retweet 確率

まず、前節でモデル化されたユーザーの興味と比較するために Tweet t に含まれる単語を抽出する ($f^{(u,t)}$)。ここで、出現した単語と、過去の Retweet による興味のモデル化である $f_j^{(u)}$ との比較を行う。これにより、特徴キーワード j が出現したとき再共有しない確率は

$$1 - f_j^{(u)} \cdot f_j^{(u,t)} \quad (2)$$

で表される。今、各単語の出現を独立事象と仮定するとユーザー u の、ある Tweet t を Retweet する確率 $P(t, u)$ は、以下のように記述することができる。

$$P_{(t,u)} = 1 - \prod_j (1 - f_j^{(u)} \cdot f_j^{(u,t)}) \quad (3)$$

この Retweet する確率は、3 章で述べたような Twitter 上での情報拡散モデルにおいて、フォロワーに対する情報を伝播するための確率であるといえる。この確率を用いることで、フォロワーネットワーク上での情報の拡散を予測することができる。

6. まとめと今後の課題

本稿では、実際のデータを用いて情報の再共有の行動がユーザーの興味に基づいていることを示し、それに基づいた情報再共有の確率モデルを提案した。今後は、実際のデータを用いて再共有の確率モデルの妥当性を検証したい。

ユーザーの情報再共有の確率モデルには改善の余地がまだ数多くある。例えば、表 2 において、ユーザー B に「有名人」や人や会社名などの固有名詞が出てきているが、これは 2011 年 1 月に放映されたテレビ番組に関する Tweet を頻繁に Retweet していたことによる。このような時期的に偏ったトピックをある時期にまとめて Retweet するというのはよく見られる現象である。Retweet のモデルではこのような時期に依存するトピックがノイズになることも考えられるため、時間による忘却などを考慮する必要があるだろう。また、本稿で提案したある Tweet に対する Retweet 確率では、その Tweet を実際に目にするかどうかは考慮にいていないため、モデル化するにあたり考慮する必要があるといえる。

SIR モデルで情報拡散のモデルを考えると、情報を知らない状態 S、情報を受け取って拡散する状態 I、そして興味を失い情報を拡散しない状態 R の 3 状態を考えた。通常の噂の伝播モデルでは、状態 S が状態 I の頂点と隣接しているとき、状態 S の頂点は状態 I の頂点から一定の確率で情報を得て状態 I へと遷移 (感染) し、状態 I の頂点は、一定の確率で興味を失い、状態 R へと遷移すると述べた。しかし、Twitter に

おける拡散モデルを考えるとときには通常の噂の伝播とは違うモデルを考える必要がある。

[10] では Twitter における返答および RT の構造を網羅的に調査しており, Retweet の半数以上が情報が共有されてから 1 時間以内に再共有されているという結果を示した。これは Twitter がリアルタイム型のメディアである特性から, 情報が共有されてから時間が経過すると情報の再共有という行為があまり行われないうことを示している。つまり, Twitter のようなリアルタイム型のメディアにとっては状態 I である時間は短く, 情報を再共有するとすぐに状態 R に遷移してしまう。また, 情報に接触しても興味がなかったり, その情報を他者に伝えようというモチベーションがない場合は, 再共有を行わず状態 S から速やかに状態 R に遷移する。今後, 情報拡散を考えるとときには本稿で述べた情報の再共有の確率モデルを既存の情報拡散に当てはめるだけではなく, リアルタイム型のメディアに即した新しいモデルを提案する必要があるだろう。

謝辞

本研究の一部は, 総合科学技術会議により制度設計された最先端研究開発支援プログラム (FIRST 合原最先端数理モデルプロジェクト) により, 日本学術振興会を通して助成されたものです。

参考文献

- [1] danah boyd, Scott Golder, and Gilad Lotan. Tweet, tweet, retweet: Conversational aspects of retweeting on twitter. In *Proceedings of the 43rd Hawaii International Conference on Social Systems (HICSS)*, 2010.
- [2] 太田侑介, 寺田実, 丸山一貴. Twitter におけるリツイート経路の重ね合わせによるユーザ発見支援. 第 10 回情報科学技術フォーラム (FIT2011), 2011.
- [3] T. R. Zaman, R. Herbrich, J. Van Gael, and D. Stern. Predicting information spreading in twitter. In *Proceedings of Computational Social Science and the Wisdom of Crowds Workshop*, 2010.
- [4] Zi Yang, Jingyi Guo, Keke Cai, Jie Tang, Juanzi Li, Li Zhang, and Zhong Su. Understanding retweeting behaviors in social networks. In *Proceedings of the 19th ACM international conference on Information and knowledge management, CIKM '10*, 2010.
- [5] 風間一洋, 今田美幸, 柏木啓一郎. Twitter の情報伝播ネットワークの分析. 第 24 回人工知能学会全国大会, 2010.
- [6] Jiang Yang and Scott Counts. Predicting the speed, scale, and range of information diffusion in twitter. In *ICWSM*, 2010.
- [7] Geerajit Rattananitnont, Masashi Toyoda, and Masaru Kitsuregawa. A study on characteristics of topic-specific information cascade in twitter. In *Forum on Data Engineering (DE2011)*, 2011.
- [8] 増田直紀, 今野紀雄. 複雑ネットワーク 基礎から応用まで. 近代科学社, 2010.
- [9] S. Boccaletti, V. Latora, Y. Moreno, M. Chavez, and D-U. Hwang. Complex networks : Structure and dynamics. *Phys. Rep.*, Vol. 424, No. 4-5, pp. 175-308, 2006.
- [10] Haewoon Kwak, Changhyun Lee, Hosung Park, and Sue Moon. What is Twitter, a social network or a news media? In *Proceedings of the 19th international conference on World wide web (WWW2010)*, 2010.