

リプライを用いたバズツイートの分類

羽田 優太 松本 和幸 吉田 稔 北 研二

徳島大学

{matumoto;mino;kita}@is.tokushima-u.ac.jp

1 はじめに

近年 SNS の発達により、多数のユーザ間で口コミの拡散や共有が盛んにおこなわれるようになった。これによりインターネット上でコンテンツの流行が度々起こるようになった。ユーザ同士のコミュニケーションによっておこる流行のことを「バズる」と呼ぶ。この「バズる」ことを利用したマーケティングが重要視されており、いち早くバズに気づき取り入れることが今後のマーケティングで重要になってくると考えられる。

既存研究において、時系列での拡散のされ方などから流行を検知するもの [1][2] があるが、それらの手法では拡散された要因までは知ることができないという問題点がある。本研究では、Twitter においてバズったツイートに寄せられたリプライをベクトル化し、ツイートに対するリプライ全体を特徴ベクトルとすることで、「バズっている」特徴の定義を目指すとともにバズツイート内におけるタイプ別の分類分けの可能性について検討する。

2 提案手法

2.1 データ収集

Twitter API を用いて、バズったツイートおよびそれに対するリプライ、バズっていないツイートとそれに対するリプライの収集を行う。

学習用データ

リプライを特徴ベクトルに変換する際に、ツイート単位での特徴抽出手法が必要となる。本研究では、ツイート全般において、既存の単語の頻度ベクトルを用いる Bag of Words 特徴量では未知の単語への対応が困難であるため、より汎用性の高い手法として、単語分散表現を文単位に拡張した、文の分散表現を採用する。具体的には、

Doc2Vec の実装の一つである Sentence2vec と、リカレントニューラルネットワーク (RNN) に基づく AutoEncoder、畳み込みニューラルネットワーク (CNN) に基づく AutoEncoder を用いる。分散表現モデルの学習用データとしては、リプライのみを収集した。表 1 に実際に収集したリプライデータの一例を記載する。Twitter Streaming API を用いることで、リプライが付きやすいユーザ (芸能人、有名人、政治家、etc.) を対象として、それらのユーザのフォロワーが投稿したリプライに限定して収集した。

表 1: リプライデータの一例

	リプライ
1	すごい!! ステキ!!
2	ナイスプレー!!
3	まっ無理だろうね。
...	...

バズツイートのデータ

バズツイートに対して投稿されたリプライを収集する。バズツイートの特徴ベクトルを求めるために、リプライのテキストを先述したそれぞれの手法でベクトル化する。表 2 に実際に収集したバズツイートについたリプライ例を記載する。

バズっていないツイートのデータ

バズっていないツイートに対して投稿されたリプライを収集する。バズツイートとの特徴量の差を見つけるために利用する。

2.2 ベクトル化

リプライ文をベクトル化する。ベクトル化には複数の手法を用い、どの手法が最適かを実験結果

表 2: バズツイートに対して投稿されたリプライの実例

元ツイート	娘を小児科に連れて行った時の事。診察室から『ギャー!助けてー!やめろー!』と幼い男子の悲鳴が響き渡った。そしたらお母さんが『コラ!何て口の利き方なの!』と一蹴。一瞬の沈黙の後、『ギャー!やめて下さいませー!お医者さまー!助けて下さいませー!ギャー!』院内が笑いに包まれた
リプライ 1	リアルに飯吹いた w
2	コピペ乙
3	仮面ライダー誕生のシーンもこんな感じだったんだろうか
...	...

の比較によって明らかにする。図 1 に、リプライのベクトル化までの実験の流れを示す。本稿では、RNN(Recurrent Neural Networks)に基づくベクトル化では、LSTM(Long Short-Term Memory)とGRU(Gated Recurrent Unit)[3]を隠れ層にそれぞれ持つ AutoEncoder, 畳み込みニューラルネットワーク CNN(Convolutional Neural Networks)[4]に基づくベクトル化では、複数の畳み込み層とプーリング層を持つ AutoEncoder を構築した。入力には、文字の one-hot 表現を用い、抽出する隠れ層の次元は 64 次元とした。

また、文の分散表現ベクトルの学習手法である Sentence2vec[5]を用い、日本語形態素解析器 MeCab によって単語分かち書き表現にした後、単語列を入力として文の分散表現を学習した。分散表現ベクトルの次元数は 100 とした。

本研究で単語頻度ベクトルを特徴として用いない理由として、リプライは短文であることが多く、さらに人によって語彙の差があることから、十分な特徴が見いだせない可能性が高いと考えたからである。以下、本稿で利用するそれぞれの手法の概要について述べる。

2.2.1 LSTM AutoEncoder

リカレントニューラルネットワークとは隠れ層からの出力を再び隠れ層に入力するニューラルネットワークである。LSTM は長期的な時系列のデータを学習することができる。

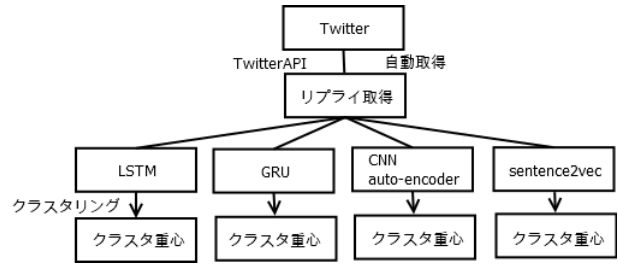


図 1: 実験の流れ 1

2.2.2 GRU AutoEncoder

GRU は LSTM の忘却ゲートと入力ゲートを組み合わせ更新ゲートという単一のゲートにしたニューラルネットワークである。LSTM よりシンプルなモデルを得ることができ、LSTM よりも高速な学習が可能である。

2.2.3 CNN AutoEncoder

CNN は畳み込みニューラルネットワークであり、入力を畳み込み演算により特徴以外の情報を圧縮することで入力の特徴を検出するニューラルネットワークである。図 2 に簡単な CNN の隠れ層の処理の流れを記載する。

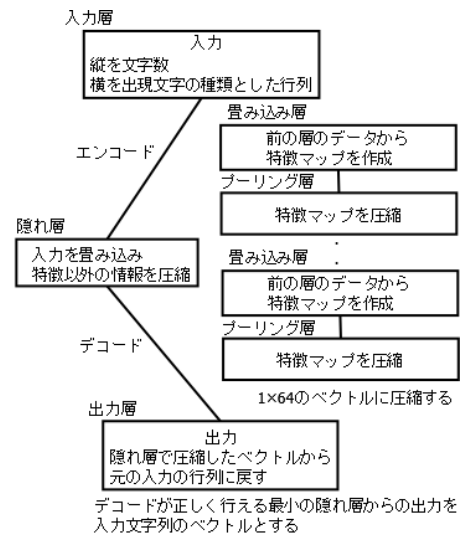


図 2: CNN の隠れ層の処理

2.2.4 Sentence2vec

Word2vec の登場により、単語を分散表現ベクトル化することで単語間の関連度をベクトル間の演算によ

り求めることが可能になり、単語の意味をとらえることが以前よりも容易になった。Sentence2vec は文を複数の単語ベクトルの集合とし、文のベクトルを求める手法である。これにより文同士の類似度を求めることも可能である。

2.3 クラスタリング

ベクトル化したリプライデータをクラスタリングする。リプライデータにクラスタリングを行うことで、リプライ先ツイートにどのようなタイプのリプライがどれだけ投稿されているかがわかるため、それをもとに特徴ベクトルを作成することで、バズツイートの特徴をとらえることができる。本研究では、リプライ毎にクラスタを決定し、どのクラスタのリプライがどの程度投稿されたかに基づくクラスタ頻度ベクトルをバズツイート毎に作成する。本稿では、K 平均法を用いてクラスタリングを行う。図 3 にクラスタ頻度ベクトルを求めるまでの実験の流れを示す。

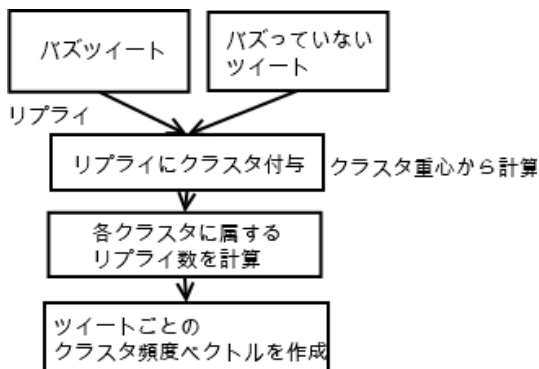


図 3: 実験の流れ 2

2.3.1 K 平均法

K 平均 (k-means) 法とは非階層型クラスタリングのアルゴリズムである。階層型クラスタリングとは異なり、あらかじめクラスタ数を設定するクラスタリングである。本研究ではクラスタ頻度ベクトルを求めるという点でクラスタ数が事前に決まっているほうが望ましいため、この方式を採用した。

2.4 類似度計算

クラスタ頻度ベクトルからバズツイートの特徴を定義づけるため、バズったツイートとバズっていないツ

weetのリプライに基づくクラスタ頻度ベクトル間の類似度を求め、差別化可能な特徴を見つける。図 4 にリプライに基づくクラスタ頻度ベクトル間の類似度計算の流れを示す。

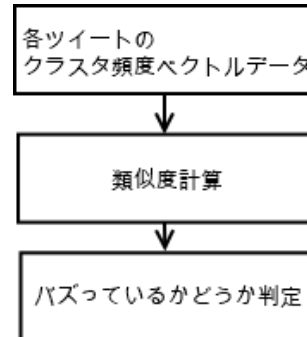


図 4: 実験の流れ 3

3 予備実験結果

本研究における現時点での実験設定と、予備実験結果について記載する。現時点でのデータ取得数は、バズツイート 12 件、そのツイートについてリプライ 1041 件、バズツイート以外も含めたリプライ 32,390 件を取得している。

また、Sentence2vec を用いたベクトル化は、学習パラメータは size=100, window=8 として、学習が完了している。これを 100 クラスタでクラスタリングしたところ、クラスタ頻度ベクトル間の類似度の関係は、図 5 に示すようになった。以降は、これにバズっていないツイートのデータを追加し、再度実験をおこなう。また、CNN や LSTM AutoEncoder を用いてベクトル化を行い、今回の結果との比較をおこなうことで、バズツイートの検出、分類に適した特徴抽出手法について引き続き検討する。

	data2	data3	data4	data5	data6	data7	data8	data9	data10	data11	data12
data1	0.32	0.642	0.443	0.648	0.485	0.517	0.483	0.506	0.5411	0.4546	0.3935
data2		0.399	0.533	0.454	0.484	0.262	0.556	0.432	0.3904	0.4542	0.5367
data3			0.515	0.688	0.619	0.679	0.671	0.719	0.4981	0.4979	0.3647
data4				0.562	0.549	0.381	0.616	0.545	0.5462	0.5864	0.5067
data5					0.583	0.628	0.629	0.672	0.5021	0.4856	0.4413
data6						0.551	0.529	0.542	0.4068	0.5449	0.429
data7							0.549	0.529	0.3402	0.4751	0.3715
data8								0.705	0.619	0.5681	0.4544
data9									0.4234	0.3604	0.2828
data10										0.5545	0.4812
data11											0.4159

図 5: 実験結果のヒートマップ

4 おわりに

本稿では、Twitter 上での流行現象である「バズる」に着目し、バズったツイートに投稿されたリプライに基づく特徴抽出と特徴間の類似度について分析をおこなった。予備実験の結果、ベクトル間の類似性には一定の傾向がみられるため、バズったツイートとそうでないツイートのリプライ特徴をさらに詳しく分析していくことで、バズツイートの検出に本手法が有効かどうかを検証していきたいと考えている。

謝辞

本研究の一部は、科学研究費補助金（18K11549, 15K16077）の補助を受けて行った。

参考文献

- [1] 斎藤翔太, 富岡亮太, 山西健司. ソーシャルネットワークにおける長期間流行する話題の早期検出. 電子情報通信学会技術研究報告. IBISML, 情報論的学習理論と機械学習 = IEICE technical report. IBISML, Information-based induction sciences and machine learning, Vol. 111, No. 480, pp. 77–84, mar 2012.
- [2] 谷季恵, 松村嘉之. Twitter 上の情報拡散がもたらす商品販売効果推定モデルの提案. 精密工学会学術講演会講演論文集, Vol. 2016, pp. 3–4, 2016.
- [3] 五島圭一, 高橋大志, 寺野隆雄. リカレントニューラルネットワークによるボラティリティ変動モデリング. 経営情報学会 全国研究発表大会要旨集, Vol. 2017, pp. 29–31, 2017.
- [4] 大山真司, 山崎俊彦, 相澤清晴. プレゼンテーションスライドのデザインに対する cnn を用いた客観的フィードバックの生成 (画像工学). 電子情報通信学会技術研究報告 = IEICE technical report : 信学技報, Vol. 117, No. 432, pp. 223–228, feb 2018.
- [5] 宮脇克典, 白松俊, 水野創太, 福本加奈恵, 池田雄斗. 科目区分ダイアグラム検索システムにおけるテキスト類似度に基づく科目推薦機構の試作. 人工知能学会全国大会論文集, Vol. 2017, pp. 3N22–3N22, 2017.