

卒業論文

ソーシャルメディアにおける
ユーザーコミュニティの情報を
用いたバースト予測に関する研究

03-130925 石塚 淳

指導教員 坂田一郎 教授

2015 年 2 月

東京大学工学部システム創成学科知能社会システムコース

概要

近年, Twitter や Facebook などのソーシャルメディアと呼ばれるサービスが普及し, 一般のユーザーが容易に情報を発信できるようになった. ソーシャルメディアの中でも Twitter はリアルタイム性が高く, リツイートなどの機能によって情報伝播が起こりやすいという特徴を持っている. Twitter の普及とともに炎上というウェブ上の現象にも注目が集まっている. Twitter が普及する以前の炎上は, 芸能人のブログに対して, 閲覧者の批判的なコメントが集中的に集まるというような事態を指すことが多かったが, Twitter というソーシャルメディアの普及によって, 炎上は芸能人などの一部の人間だけではなく, 我々の身にも起こりうる身近なものとなっている. 本研究は Twitter 上には異なる特徴を持った様々なコミュニティが存在し, コミュニティによって炎上への関わり方が異なっているのではないかという仮説に基づいて, ユーザーのコミュニティに着目する. 本研究の目的は, この仮説に基づき, Twitter のネットワークデータやツイートデータ, プロフィールデータを分析することにより, 情報伝播や炎上に影響を与える要素について知見を得ることである. 情報伝播や炎上についての知見は, 炎上の弊害を回避し, 個人として炎上に巻き込まれないために役立つと考えられる. Twitter 上で炎上が起きている最中では, あるツイートが数多くリツイートされる現象がよく起こる. 本研究ではその現象をバーストと定義し, バーストをある種の情報伝播であると捉え, ツイートがバーストするかどうかを予測する二値分類器を構築する. 日本国内の全 Twitter ユーザーから, ネットワークを構築してクラスタリングを行うことによって, ユーザーのコミュニティを抽出する. その後, 従来の情報伝播予測やリツイート予測でも使われていたユーザー属性とツイート属性に加えて, 抽出したコミュニティの情報をを用いたコミュニティ属性を特徴量として使用することでバースト予測の精度の改善が見られた.

キーワード ソーシャルメディア, SNS, Twitter, 情報伝播, 炎上, バースト

目次

第 1 章	はじめに	1
1.1	背景	1
1.2	目的	2
1.3	各章の構成	3
第 2 章	関連研究	4
2.1	ソーシャルメディアのコミュニティ分析	4
2.2	ソーシャルメディアの情報伝播と炎上分析	5
2.3	本研究の位置付け	5
第 3 章	ソーシャルメディアにおけるバースト予測	7
3.1	コミュニティ情報を用いたバースト予測	7
3.2	コミュニティ抽出	8
3.3	バースト予測モデルの構築	9
第 4 章	コミュニティ抽出と炎上事例分析	14
4.1	実験	14
4.2	抽出されたコミュニティの結果と考察	16
4.3	コミュニティ情報とツイート時間帯分析	17
4.4	コミュニティ情報とリツイート分析	19
4.5	まとめ	22
第 5 章	バースト予測モデルの構築	26
5.1	実験	26
5.2	特徴量の組み合わせによるバースト予測精度の結果と考察	31
5.3	予測に有効な素性の結果と考察	33
5.4	個々の事例による予測精度の結果と考察	34
5.5	まとめ	35
第 6 章	結論	37

iv 目次

謝辭 39

参考文献 40

目次

1.1	Twitter のユーザーホーム画面	2
3.1	提案手法における全体のフレームワーク	8
3.2	SVM におけるマージン最大化	12
3.3	ロジスティック曲線の例	13
4.1	コミュニティごとの 1 日あたりのツイート数 (生活保護)	18
4.2	コミュニティごとの 1 日あたりのツイート数 (献血)	18
4.3	コミュニティごとの 1 時間あたりのツイート数 (美味しんぼ)	19
4.4	一定回数以上リツイートされたツイート件数の割合	20
4.5	コミュニティごとの 1 日あたりのリツイート数 (生活保護)	21
4.6	コミュニティごとの 1 日あたりのリツイート数 (献血)	21
4.7	コミュニティごとの 1 時間あたりのリツイート数 (美味しんぼ)	22

表目次

4.1	3つの事例におけるツイートの収集条件	15
4.2	3つの事例における収集したツイートの基礎データ	15
4.3	コミュニティ抽出におけるパラメータ	16
4.4	コミュニティごとの3つの炎上事例への参加人数とコミュニティ情報	24
4.5	コミュニティごとのTwitterの基本データの中央値	25
4.6	コミュニティごとのクラスター係数と各中心性の平均値	25
4.7	3つの事例におけるリツイート回数の分散値	25
5.1	バースト予測モデルの構築に用いたデータ	27
5.2	実験に用いた特徴量	27
5.3	3つの事例の全データにおけるバースト予測モデルの特徴量の組み合わせ	28
5.4	6つの事例の個々のデータにおけるバースト予測モデルのテスト結果	28
5.5	機械学習による予測結果の分類	29
5.6	6つの事例におけるツイートの収集条件	30
5.7	6つの事例における収集したツイートの基礎データ	30
5.8	3つの事例におけるある回数以上のリツイートされたツイート件数	31
5.9	それぞれの事例におけるSVMのパラメータ	31
5.10	3つの事例の全データにおけるバースト予測モデルの実験結果	32
5.11	P値が0.1以下のロジスティック回帰モデルにおける特徴量	33
5.12	6つの事例の個々のデータにおけるバースト予測モデルのテスト結果	35

第 1 章

はじめに

1.1 背景

ソーシャルメディアと呼ばれるサービスが注目を浴びてから久しい。ソーシャルメディアとは、一般のユーザーが情報を発信することができるサービスであり、双方向のコミュニケーションができるメディアである。その代表例とも言える Facebook や Twitter というサービスは全世界で流行し、日本でも多くのユーザーが利用している。

Twitter はソーシャルメディアの中でもリアルタイム性が高く、情報伝搬が起こりやすいという特徴を持つ。ツイートと呼ばれる 140 文字の短文をウェブサイトに掲載したり、利用者間でコミュニケーションを取ることができるサービスで、Twitter のユーザーホーム画面は図 1.1 のようになっている。Twitter にはユーザーをフォローすると、そのフォローしたユーザーのツイートが自身のフィードであるタイムラインに表示されるようになるフォローという仕組みが存在する。このフォローの関係をを用いて、Twitter ではソーシャルネットワークを構築している。2012 年の 3 月時点で、アクティブユーザーは世界では 1 億 4000 万人、日本では 1400 万人に達し^{*1}、多くの人々にとって欠かせない存在となっている。Twitter によって容易に情報を発信できるようになり、リアルタイムな情報を簡単に手に入れることができるようになった。その一方で、個人の知られたくない情報が拡散され、流出してしまったり、信憑性の低い情報が出回るなどのデメリットも存在する。

Twitter というソーシャルメディアが普及するに伴い、炎上というウェブ上の現象にも注目が集まっている。炎上は、なんらかの不祥事をきっかけに爆発的に注目を集める事態、または状況として定義される^{*2}。Twitter が普及する以前の炎上は、芸能人のブログに対して、閲覧者の批判的なコメントが集中的に集まるというような事態を指すことが多かった。しかし近年では、Twitter というサービスが普及したことによって、炎上は芸能人などの一部の人間だけではなく、我々の身にも起こりうる身近なものとなっている。実際に、ある一つのツイートが拡散され、その話題について短期間で爆発的に議論が行われるといった炎上現象も頻繁に起き

^{*1} 2012 年 3 月時点。平成 24 年度版情報通信白書（総務省）

^{*2} Wikipedia ([http://ja.wikipedia.org/wiki/炎上_\(ネット用語\)](http://ja.wikipedia.org/wiki/炎上_(ネット用語))) 25 Dec. 2014, UTC 10:49)

2 第1章 はじめに



図 1.1. Twitter のユーザーホーム画面

ている。NAVER まとめ^{*3}や Together まとめ^{*4}などのキュレーションサイトでは、実際のツイートと共に炎上事例がまとめられ、多くの記事が作成されている。またユーザーが反社会的行動を Twitter 上に投稿し、リツイートなどによって拡散されるバカッターと呼ばれる現象も広まり、2013 年新語・流行語大賞の候補語 50 語にノミネートされている^{*5}。常に炎上が起こりうる状況となった現在では、企業の危機管理においても炎上は重要な存在になってきている。NHK ニュースの特集ページでは食品メーカーの商品に虫が混入されていたという事例を交えて、ネット炎上への迅速な対応が必要となっていることを伝えている^{*6}。こうした Twitter における炎上はユーザーにとって、情報過多や意見の偏り、信憑性などの問題となる。例えば情報過多という点には、短期間で爆発的に議論が行われるので、ユーザーが知りたい情報や正確な情報、全体像を適切に把握することが困難になるという問題が存在する。また意見の偏りという点には、Twitter のフォローという特性上、自分と似たようなユーザーをフォローする傾向があるので、自分の周りの意見は偏ってしまうことが多いなどの問題が存在する。

1.2 目的

近年では Twitter のソーシャルネットワークとしての特性を活かし、コミュニティに着目した研究が数多く行われている [1]。榊らの研究ではツイートとユーザーの双方をクラスタリングすることによって、ツイートのコンテンツによって情報伝播のされ方が異なることを示した [2]。それらの研究から、Twitter には異なる特徴を持った様々なコミュニティが存在し、炎上の初期段階から関与する炎上を引き起こす原因となっているコミュニティや、炎上が盛り上がった段階から興味本位で炎上に関与するコミュニティなどのように、コミュニティによって炎上への関わり方が異なっているのではないかという仮説に至った。本研究の目的は、この仮説に

^{*3} <http://matome.naver.jp/> 05 Feb. 2015, UTC 05:24

^{*4} <http://togetter.com/> 05 Feb. 2015, UTC 05:24

^{*5} 2013 年ユーキャン新語・流行語大賞候補語 50 語より

^{*6} 企業を襲う インターネットの“炎上”(NHK ニュース) <http://www.nhk.or.jp/ohayou/marugoto/2014/04/0416.html> 30 Jan. 2015, UTC 11:24

基づき、Twitter のネットワークデータやツイートデータ、プロフィールデータを分析することにより、情報伝播や炎上についての知見を得ることである。情報伝播や炎上についての知見は、炎上の被害を回避し、個人として炎上に巻き込まれないために役立つと考えられる。また近年では、企業が自社製品の宣伝や採用活動、認知度の向上を目的として Twitter や Facebook のアカウントを作成することも珍しくない。そのため本研究で得られた情報伝播の知見は、企業のマーケティング活動などにも応用が期待できる。

Twitter 上で炎上が起きている最中では、あるツイートが数多くリツイートされる現象がよく起こる。本研究ではその現象をバーストと定義し、バーストをある種の情報伝播であると捉え、ツイートがバーストするかどうかを予測する二値分類器を構築する。日本国内の全 Twitter ユーザーから、ネットワークを構築してクラスタリングを行うことによって、ユーザーのコミュニティを抽出する。その後、従来の情報伝播予測やリツイート予測でも使われていたユーザー属性とツイート属性に加えて、抽出したコミュニティの情報をを用いたコミュニティ属性を特徴量として使用することでバースト予測の精度の改善が見られた。

1.3 各章の構成

この節では、本研究における各章の構成を説明する。まず 3 章で研究の全体のフレームワークや提案手法について説明する。4 章ではコミュニティ抽出の実験を行い、得られたコミュニティの差異を確認するため、生活保護不正受給、群馬県議員による献血に関するツイート問題、漫画「美味しんぼ」における風評被害問題の 3 つの炎上事例を横断した分析を行った。5 章ではツイートがバーストするかどうかを判定する二値分類器を構築し、そのモデルの性能のテストを行った結果と考察について述べる。

第 2 章

関連研究

この章では、まず一般的なソーシャルメディアのコミュニティ分析を行った研究を紹介し、次にリツイートなどに代表される、本研究と関連の深い、ソーシャルメディアの情報伝播と炎上分析についての研究を紹介する。そして最後に本研究の位置付けとして、それらの研究との差異を述べる。

2.1 ソーシャルメディアのコミュニティ分析

Twitter というソーシャルメディアの研究はサービスが開始して以降、盛んに行われている。Twitter は 2006 年 10 月にサービスが開始されたマイクロブローギングサービスであるが、その一年後には Twitter を対象としたネットワーク分析の研究が発表されている [3]。Java らによるネットワーク分析の研究では、フォローネットワークを分析し、フォロー数とフォロワー数によって分割されたコミュニティ間ではユーザーの使用目的や使用方法が異なることを明らかにした。例えば Twitter を情報収集のツールとして利用していて、自らはほとんどツイートしないユーザーが存在することを示した。

Kwak らは、フォローとフォロワーの関係を分析し、Twitter では約 8 割が一方向のフォローであり、残りの約 2 割しか相互フォローは存在しないことを明らかにした [4]。さらに 6 割強のアカウントはフォロワーが存在せず、誰からもフォローされていないことを示した。この研究も Twitter を情報収集のツールとして利用するユーザーが多いことを示している。

また Huberman らによる研究では、フォローネットワークではなく、相互メンションネットワークが提案された [5]。Twitter には、「@ (ユーザー名)」の形式をツイートに含ませることによって、あるユーザー宛にツイートをするメンションという機能がある。彼らの研究によると、相互にメンションが行われたユーザー間にリンクを張ることによって構築されるネットワークにおけるリンク数の方が、フォロー数やフォロワー数よりもツイート数に影響を及ぼすことを明らかにした。

ソーシャルネットワークとしての Twitter という側面に着目した丸井らの研究では、ユーザーのプロフィール情報とメンションによるネットワーク情報、言葉遣いを利用してユーザーのコミュニティを抽出する手法を提案した [1]。また抽出したコミュニティを用いて、ネット

ワーク情報や言葉遣いの類似しているコミュニティを分析し、ユーザーのプロフィールにおける興味や属性に基づいた結果が得られるということを明らかにした。

2.2 ソーシャルメディアの情報伝播と炎上分析

次に本研究と関連の深い、ソーシャルメディアでの情報伝播と炎上分析に関する研究を紹介する。情報伝播に関する研究は、ユーザーの行動により情報伝播のモデルを構築することが目的となっているものが多い。またリツイート（他のユーザーのツイートを引用形式で自分のアカウントからツイートすること）を情報伝播と捉え、モデル化に応用したり、リツイート数を予測する研究も存在する。

情報伝播力の高いインフルエンサーを発見する研究では、フォロー関係にあるユーザーが URL を投稿したとき、同じ URL を自分もツイートした場合に情報伝播が起きたと仮定し、情報伝播をモデル化している [6]。その後、各々のユーザーに対して情報伝播から影響力と定義したスコアを計算し、フォロー数やツイート数などのユーザー属性から影響力の予測を行っている。

榊らの研究では人工知能の表紙における問題を対象に、ユーザーとツイートに含まれる URL の双方をクラスタリングすることによって、情報伝播の可視化を提案している [2]。ユーザーコミュニティとツイートクラスターを用いることで、コミュニティやツイートのコンテンツによって情報伝播のされ方が異なることを示した。ユーザーのクラスタリングにはモジュラリティを最大化する手法である Louvain 法を用いている [7]。

また、Yang らによる研究では、ユーザーがトピックを含むツイートをした後、そのユーザーへトピックを含んだツイートをリプライしたときに情報伝播が起きたとモデル化している [8]。この研究では、情報伝播が最初に起きるまでの時間、そのユーザーから直接情報伝播の起きる人数、情報伝播が人を介して続く人数を Cox の比例ハザードモデル [9] によって、ユーザー属性とツイート属性から予測している。

リツイート予測に関する研究では、オンライン学習器である PassiveAggressive[10] という手法を用いることによって、Twitter から取得したストリームデータをリアルタイムに学習し、予測するモデルが提案されている [11]。予測のための特徴量は Yang らの研究と類似しており、ユーザー属性とツイート属性を使用している。

2.3 本研究の位置付け

本研究は抽出したコミュニティを用いて複数の炎上事例を横断して分析し、そのコミュニティの情報を特徴量として加えることによって、バースト予測モデルを構築する手法を提案している。

前述したようにソーシャルメディアにおけるコミュニティ分析は数多く行われており、榊らの研究では実際の炎上事例を対象にユーザーをクラスタリングしてコミュニティを抽出し、分析を行っている [2]。しかし榊らの研究では、発生した炎上事例に関わったユーザーからコミュ

6 第2章 関連研究

ニティを抽出しているので、事前にコミュニティを抽出している本研究とは異なり、コミュニティ情報を予測に活かすことができない。また本研究では事前にコミュニティを抽出することによって、同一のコミュニティを用いて複数の事例を対象に分析することができるという点も異なる。

またツイートのバースト予測モデルの構築では、Yang らや Petrovic らによる既存研究での予測で使われているユーザー属性とツイート属性を用いて、機械学習による分類を行う [8][11]。しかしながら、特徴量にコミュニティ属性を用いるといった研究はいまだ行われておらず、本研究は前述した RT 予測や情報伝播予測の発展と位置付けられる。

第3章

ソーシャルメディアにおけるバースト予測

本章では研究の着想や提案手法における全体のフレームワークについて述べた後、具体的な提案手法について述べる。提案手法における全体のフレームワークは、Twitter の相互メンションデータからユーザーのコミュニティを抽出し、ツイート、ユーザー、コミュニティの情報から特徴量を生成し、それらのデータに機械学習の分類モデルを適用し、ツイートのバーストを予測する二値分類器を作成するという流れになっている。

3.1 コミュニティ情報を用いたバースト予測

3.1.1 研究の着想と全体のフレームワーク

榊らの研究ではツイートとユーザーの双方をクラスタリングすることによって、コミュニティやツイートのコンテンツによって情報拡散のされ方が異なるを示した [2]。本研究ではその研究を始めとする情報伝播についての研究から、Twitter には異なる特徴を持った様々なコミュニティが存在し、炎上の初期段階から関与する炎上を引き起こす原因となっているコミュニティや、炎上が盛り上がった段階から興味本位で炎上に関与するコミュニティなどのように、コミュニティによって炎上への関わり方が異なっているのではないかという着想を得ている。研究のアイデアは、従来の情報伝播予測 [8] やリツイート予測 [11] で特徴量として用いられていたツイートの属性やユーザーの属性に加え、ユーザーのコミュニティの属性を加えることで、バーストの予測モデルの精度を向上させることができるのではないかという仮説に基づいている。

提案手法における全体のフレームワークの図を図 3.1 に示す。まず始めに Twitter の相互メンションデータからユーザーの相互メンションネットワークを構築する。次に構築した相互メンションネットワークをネットワーククラスタリングしてコミュニティを抽出する。その際、クラスター係数や中心性などのネットワーク指標も計算しておく。そしてツイート、ユーザー、コミュニティの情報から特徴量を生成し、それらのデータに機械学習の分類モデルを適用して、

ツイートのバーストを予測する二値分類器を作成するという流れになっている。

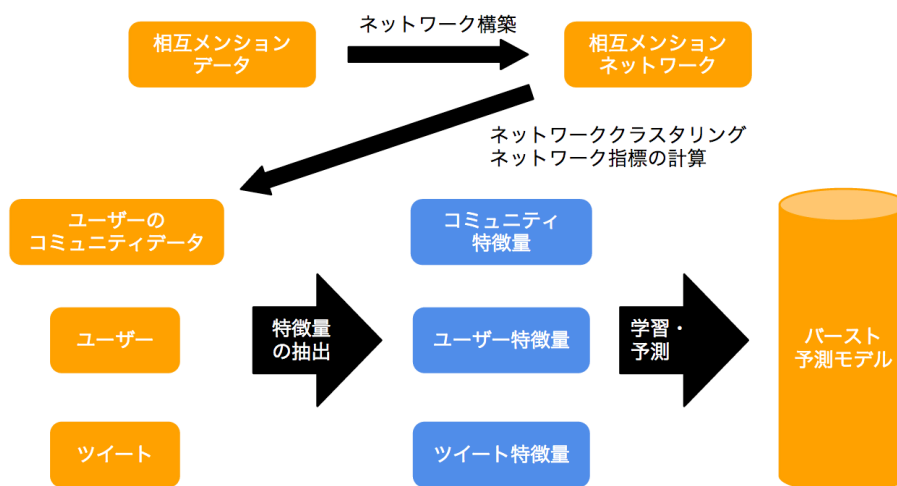


図 3.1. 提案手法における全体のフレームワーク

3.1.2 炎上とバースト

Twitter 上で炎上が起きている場合、その話題についてのツイート数が急激に多くなる。本来 Twitter におけるバーストとは、ある語を含むツイート数が急激に増加することと定義されることが多い [12]。例として「天空の城ラピュタ」が地上波で放送された日に「バルス」という語を含むツイート数が急増した事例が挙げられる^{*1}。しかしながら、この定義では個々のツイートに対してバーストが定義できない。また「バルス」の例のようにテレビなどの Twitter の外部から情報伝播が起こっていることも多いので、Twitter 上の情報伝播かどうかを判断することが難しい。したがって本研究ではツイートのバーストを、そのツイートが数多くリツイートされる現象と定義した。Twitter 上で炎上が起きている最中では、ツイート数が急増するのと同様に、リツイートも急増する。

3.2 コミュニティ抽出

それぞれの事例で統一されたコミュニティを利用するため、日本国内の全ユーザーからネットワークを構築し、ネットワーククラスタリングをすることによりコミュニティの抽出を行う。その後、抽出されたコミュニティを解釈するために、コミュニティを特徴付ける語を取得する。コミュニティごとのユーザーのプロフィール文から tf-idf 値を計算し、特徴語の抽出を行った。

^{*1} 「バルス」が 2013 年最高の TPS (Tweets Per Second) を記録 <https://twitter.com/twitterjp/status/363494742518013952> 04 Feb. 2015, UTC 18:08

3.2.1 ネットワークの構築とネットワーククラスタリング

まず日本国内の全ユーザー間の相互メンションデータを使ってリンクを張り、ネットワークを構築する。次に構築したネットワークを、コミュニティにクラスタリングする。手法はネットワーククラスタリングにおいて、最も代表的な手法の一つである Louvain 法を使用する [7]。Louvain 法はモジュラリティ Q を最大化するクラスタリング手法である Newman 法を改良した手法であり、Newman 法よりも高速かつ精度が良いことが知られている。モジュラリティ Q は以下の式で表され、各コミュニティの結合度を示す。

$$Q = \frac{1}{2m} \sum_{ij} \left[A_{ij} - \frac{k_i k_j}{2m} \right] \delta(c_i, c_j) \quad (3.1)$$

またモジュラリティが最大になるようにクラスタリングを行うと、コミュニティのメンバー数が数十万となり、コミュニティと呼ぶにはメンバー数が多すぎるコミュニティも生成される。そのため、本研究ではコミュニティのメンバー数が 1 万ユーザー以下となるように、再度そのコミュニティの中で Louvain 法を適用することによりサブクラスタリングを行い、コミュニティを抽出している。

3.2.2 コミュニティの特徴語の抽出

抽出したコミュニティを特徴づけるために、tf-idf 値による特徴語の抽出を行った。tf-idf 値は文書中の単語の重みを表す値であり、以下の 2 つの式の積で表される。

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}} \quad (3.2)$$

$$idf_i = \log \frac{|D|}{|d : d \ni t_i|} \quad (3.3)$$

$n_{i,j}$ は単語 t の文書 d における出現回数、 $|D|$ は総文書数、 $|d : d \ni t_i|$ は単語 t を含む文書数である。ある文書における単語の tf はその文書中に出現回数が多くなるほど高くなり、ある文書における単語の df は単語の出現する文書数が少ないほど高くなる。ここでは、各コミュニティごとにユーザーのプロフィール文を結合し、それを 1 文書として、単語の tf-idf 値を算出した。多くの文書に出現する一般的な単語を除くため、7 割以上の文書に出てくる単語はストップワードとして除去し、一つの文書でしか出現していない単語も汎用的な語ではないために除去した。そして各コミュニティの文書ごとの tf-idf 値が上位の単語をその文書の特徴語とする。

3.3 バースト予測モデルの構築

この節では、前節で得られたコミュニティの情報を利用した、バースト予測モデルの構築手法について説明する。まずデータから作成した特徴量について述べた後、SVM によるバースト

予測モデルの構築手法について述べる。最後にロジスティック回帰モデルによる、モデル構築後の特徴量を解釈する手法について述べる。

3.3.1 ツイートベースの特徴量

ツイートベースの特徴量は、大きく分けて2つの特徴量に分けられる。1つ目はツイートの中でも明確な数値で表現される属性による特徴量で、ハッシュタグ数、メンション数、URL 数、文字数、リプライかどうかの5つである。ハッシュタグとは、#記号とキーワードから成る文字列で、同じイベントに参加していた人や同じ興味を持つ人がツイートを検索しやすくすることに用いられる機能である。またメンションとリプライは混同されることが多いが、リプライは@とユーザー名がツイートの先頭にあるツイートのことを指し、メンションは文中にある@とユーザー名のことを指す。リプライはツイートユーザーとリプライ先のユーザーの双方をフォローしているユーザーのタイムラインにしか表示されないが、メンションはツイートユーザーをフォローしているユーザーのタイムラインに表示されるという違いがある。

2つ目はツイートの内容による特徴量で、ポジティブ率、主観率の2つである。ポジティブ率、主観率は0から1の少数値であり、東北大学の乾・岡崎研究室の公開している日本語極性辞書^{*2}から単語のマッチングを行い、評価している。ポジティブ率はツイートに含まれるマッピングした感情語がポジティブである確率、主観率はツイートに含まれるマッピングした感情語が主観的である確率と定義する。例えば、ツイートにおける辞書にマッピングした語がすべてポジティブであればポジティブ率は1、全てネガティブであれば0、ポジティブとネガティブが同数かマッピングした語がない場合は0.5となる。

3.3.2 ユーザーベースの特徴量

ユーザーベースの特徴量は、フォロー数、フォロワー数、お気に入り数、ツイート数の4つである。フォロー数はユーザーがフォローしているユーザー数で、フォロワー数はユーザーがフォローされているユーザーの数である。

3.3.3 コミュニティベースの特徴量

コミュニティベースの特徴量は、大きく分けて2つの特徴量に分けられる。1つ目は自身が属するコミュニティ内での影響度を示す特徴量で、コミュニティ内でのクラスター係数、次数中心性、近接中心性、媒介中心性、固有ベクトル中心性、PageRank, HubScore, AuthorityScoreの8つである。クラスター係数と基本的な中心性である次数中心性、近接中心性、媒介中心性の定義を以下に示す。

クラスター係数

^{*2} 乾・岡崎研究室 日本語極性辞書 <http://www.cl.ecei.tohoku.ac.jp/index.php?Open%20Resources%20Japanese%20Sentiment%20Polarity%20Dictionary>

クラスター係数はあるノードの隣接ノード同士が隣接ノードである確率を表し、ノード数を N 、ノード i の隣接している k 個のノードの中のエッジ数を E_i としたとき、以下の式で定義される。この値が高いほど、ネットワークは密であると言える。

$$C = \frac{1}{N} \sum_i \frac{E_i}{k C_2} \quad (3.4)$$

次数中心性

単純にノードに張られているエッジの本数が多いものが中心性が高いとする最もシンプルな考えに基づく次数中心性 C_d はノード i に張られているエッジの本数を d_i とすると以下の式で定義される。

$$C_d = d_i \quad (3.5)$$

近接中心性

他のノードとの距離が近いほど中心性が高いとする、ノード間の距離に着目した近接中心性 C_c は、ノード数を N 、ノード i とノード j の距離を $d(i, j)$ とすると以下の式で定義される。

$$C_c = \frac{N - 1}{\sum_{i \neq j} d(i, j)} \quad (3.6)$$

媒介中心性

そのノードを通る経路が多いと中心性が高いとする、ノード間の経路に着目した媒介中心性 C_b は、ノード数を N 、ノード i を通る経路数の総和 E_i とすると以下の式で定義される。

$$C_b = \frac{E_i}{N - 1 C_2} \quad (3.7)$$

2つ目は自身の属するコミュニティ自体の特徴量で、コミュニティのメンバー数、特徴語の2つである。コミュニティの特徴語はコミュニティの抽出で説明した tf-idf 値の上位の単語を使用し、ある単語がコミュニティの特徴語に出現したら 1、出現しなければ 0、というようなバイナリーデータとする。この手法だと単語の次元は何万、何十万といった次元になるため、疎行列になる。その高次元の疎行列を低次元に圧縮するために潜在意味解析 (Latent Semantic Analysis, LSA) を用いる [13]。圧縮次元数はパラメータとして指定することができる。LSA では特異値分解 (SVD) を用いて、 m 行 n 列の単語文書行列 A を、 m 行 m 列のユニタリ行列 U 、 m 行 n 列の実対角行列 Σ 、 n 行 n 列のユニタリ行列の随伴行列 V^* の3つの行列に分解する。以下にその式を示す。

$$A = U \Sigma V^* \quad (3.8)$$

低次元に変換されるので計算の効率が良くなるのはもちろん、単語の関連性を保って低次元に圧縮できるので、単語間の類似度を適切に評価できるようになる。例えば従来の手法では「猫」

と「ネコ」という同じ意味をもつ2つの語は全く違う単語と判断されるが、LSAを使うことで2つの語の類似性を保つことができる。

3.3.4 SVMによるバースト予測モデルの構築

本研究では上記の特徴量の中からモデルに採用する特徴量と圧縮次元数を、機械学習を用いた実験により決定することでバースト予測モデルを構築する。機械学習は、入力となる多数の訓練データから学習を行い、新しい入力データに対して出力を予測する教師あり学習と、クラスタリングなどに代表される訓練データを用いない教師なし学習に分けられる。サポートベクターマシン（SVM）は1995年に発表されて以降、最も認識性能が優れた教師あり学習を用いる分類モデルの一つとして知られている[14]。サポートベクターマシンでは、入力事例をもとに正例と負例を分ける分離平面を構築する。分離平面は図3.2のようにマージン（分離平面に最も近い正例と負例への距離）を最大化するように構築される。訓練データ i の特徴ベクトルを \mathbf{x}_i 、ラベルを y_i とすると、マージン $\|\mathbf{w}\|$ はLagrangeの未定乗数法より、 $\alpha_i \geq 0$ を用いて最小化できる。ラグランジュ関数は以下の式になる。

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^l \alpha_i [y_i \{(\mathbf{w} \cdot \mathbf{x}_i) + b\} - 1] \quad (3.9)$$

また、カーネル関数を用いて入力事例を高次元空間に写像し、その高次元空間で分離平面を構築することにより非線形な分離を可能にするカーネル法という仕組みも使われている。

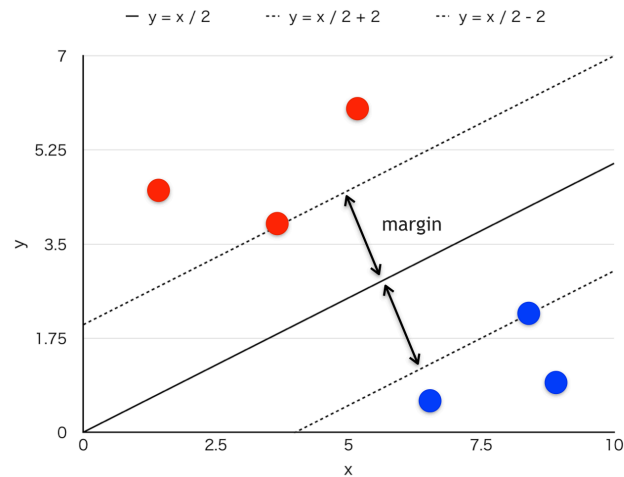


図 3.2. SVM におけるマージン最大化

3.3.5 ロジスティック回帰モデルによる特徴量の解釈

SVMにより考慮する特徴量と圧縮次元数を決定した後、ロジスティック回帰モデルを使用して特徴量の解釈を行う。ロジスティック回帰は図3.3のように、二値分類のデータを直線では

なくシグモイド曲線で回帰する手法である。一般に汎化性能は SVM に劣ると言われるが、モデル化された曲線の式の変数の偏回帰係数からオッズ比を求めることができるので、モデルを解釈しやすい。ロジスティック回帰のモデル化の式は以下のように定義される。

$$p = \frac{1}{1 + \exp(-b_0 - \sum_{i=1}^n b_i x_i)} \quad (3.10)$$

b_0 は定数, b_i は特徴量 x_i における偏回帰係数である。偏回帰係数の指数をとったもの $\exp b_i$ をその特徴量のオッズ比といい, その特徴量の値に 1 を加えたときの確率 p の上昇比を示す値である。本研究ではこのオッズ比や有意確率を使用して特徴量の解釈を行う。

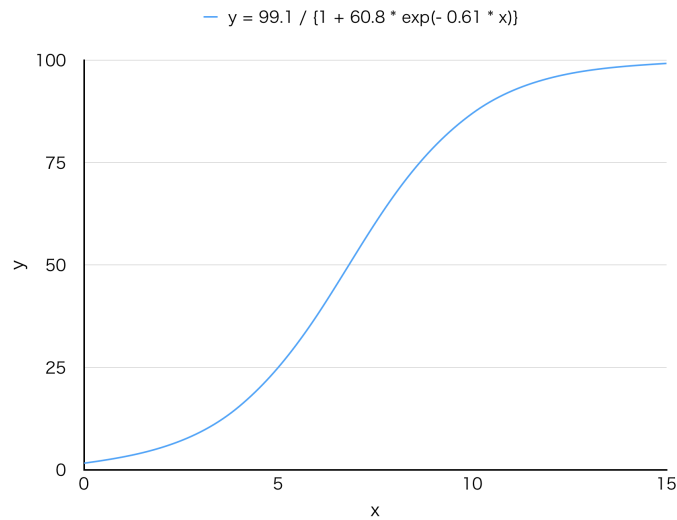


図 3.3. ロジスティック曲線の例

第 4 章

コミュニティ抽出と炎上事例分析

本章ではコミュニティ抽出と炎上事例分析についての実験、結果、考察について述べる。コミュニティ抽出を行った後、抽出したコミュニティごとに Twitter の基本データやネットワーク特徴量の比較を行い、コミュニティごとの特徴を分析する。この章では Twitter の基本データをツイート数、フォロワー数、フォローワー数、お気に入り数として実験を行う。後半の炎上事例分析では、抽出したコミュニティを用いて、ツイート時間帯やリツイートについて分析する。

4.1 実験

4.1.1 目的

本実験の目的はソーシャルメディアのネットワークデータから、5 章のバースト予測に用いることができるコミュニティを抽出することである。さらに抽出されたコミュニティごとの Twitter の基本データ、ネットワーク情報、ツイートやリツイートの行動についての特徴を得ることも目的の一つである。

4.1.2 手順

実験の手順は、3.2 節で説明したコミュニティ抽出と同様である。Twitter の相互メンションデータからネットワークを構築した後、ネットワーククラスタリングを行い、コミュニティを抽出する。次にコミュニティごとのプロフィール情報を用いて、単語の tf-idf 値を計算し、tf-idf 値が上位の単語をそのコミュニティの特徴語とする。その後、分析対象事例への参加人数が多い 10 コミュニティについて、Twitter の基本データやネットワーク特徴量の比較を行う。実装は大規模ネットワーク分析用のライブラリである SNAP ツール [15] を用いて、Louvain 法やネットワーク特徴量の計算を行っている。

炎上事例分析では、上記の 10 コミュニティを利用し、時間帯ごとのツイート数やリツイート数を比較する。またリツイートについては、事例ごとのリツイート回数の分布や分散についても比較分析を行う。

4.1.3 データ

本研究を行うにあたり、ソーシャル・ビッグデータ活用を支援するクラウドサービス事業を運営する株式会社ホットリンクから学術的な目的でデータの提供を受けた。本章の実験で使したデータは主に、炎上事例のツイートデータ、ユーザーのプロフィールデータ、ユーザー間の相互メンションデータの三点である。

分析事例として、ここ数年で炎上として認知されていて、ユーザー間の相互メンションデータと発生時期が比較的近い以下の3つの事例について、検索キーワードに基づいて期間内のツイートを収集した。詳しい検索キーワードと収集期間は表4.1に示した通りである。収集したツイートデータは、投稿日時、ツイートID、ツイート内容、投稿ユーザーのアカウント名、リツイートID（リツイート元のツイートID）、リツイート回数から構成されている。基礎データに関しては表4.2の通りで、ツイート数が多い事例では、一日あたり3万から4万ツイートが投稿されている事例も存在する。

1. 生活保護不正受給（2012年）
2. 群馬県議員による献血に関するツイート問題（2012年）
3. 漫画「美味しんぼ」における風評被害問題（2014年）

表 4.1. 3つの事例におけるツイートの収集条件

事例	検索キーワード	検索期間
1	次長課長 河本 不正受給 生活保護 生ポ ナマポ 生ぼ なまぼ	2012/04/01 - 2012/08/31
2	庭山 献血 汚染地域	2012/04/01 - 2012/07/31
3	美味しんぼ スピリッツ おいしんぼ 原因不明の鼻血 風評被害 鼻血描写	2014/04/27 - 2014/04/30

表 4.2. 3つの事例における収集したツイートの基礎データ

事例	ツイート数	RT 数	収集期間
1	35.0 万	10.6 万	121 日
2	13.4 万	2.95 万	153 日
3	14.1 万	5.29 万	4 日

ユーザーのプロフィールデータは、アカウント名（半角英数字とアンダースコアによる15文字以内の一意な名前）、スクリーンネーム（表示名）、プロフィール文、ツイート数、フォロワー数、フォロワー数、お気に入り数、アカウント作成日から構成されている。

ユーザー間の相互メンションデータは、アカウント名のペアで表現されており、無向である。このデータは2012年1月1日から2012年3月31日までの3ヶ月間における日本国内のツ

イートデータから取得した相互メンションデータであり、総リンク数は 36,743,689 本、総ノード数は 5,980,977 個となっている。

4.1.4 パラメータ

用いたパラメータを以下の表 4.3 に示す。Louvain 法ではモジュラリティ Q を最大化するので、メンバー数が数十万人のようなコミュニティと呼ぶには大きすぎるコミュニティを抽出してしまう。そこで本実験ではメンバー数が一定以上のコミュニティに対して、そのコミュニティをさらに Louvain 法で再度クラスタリングを行っている。本実験ではそのメンバー数の上限を 10,000 人として実験を行っている。またコミュニティの特徴語は tf-idf 値の上位の単語を抽出しているが、本実験では 20 単語とした。多くの文書に出現する「です」や「私」などの一般的な単語を除くため、出現する文書数の割合を表す df に対して上限値を設定した。ここでは、df の上限値を全文書数の 0.7 倍と設定しているので、7 割以上の文書に出現する単語を除去している。

表 4.3. コミュニティ抽出におけるパラメータ

パラメータ	値
コミュニティのメンバー数の上限	10,000 人
特徴語の抽出単語数	20 単語
df の上限値	全文書数 * 0.7

4.2 抽出されたコミュニティの結果と考察

コミュニティ抽出の結果として、65535 個のコミュニティが得られたが、その約半数ほどがユーザーを 1 人しか持たないコミュニティであった。メンバー数が 100 名以上のコミュニティに限ると、4366 個のコミュニティが抽出された。

表 4.4 は、3 つの炎上事例への参加人数が多い順に抽出されたコミュニティの上位 10 個のメンバー数や特徴語を示したものである。アンダーバーで接続された数字はサブクラスタリングされる前のコミュニティを表している。例えば、コミュニティ 1_6_1 は始めにクラスタリングされたネットワークの 1 番目（コミュニティのメンバー数の多い順にソートしている）のコミュニティをサブクラスタリングし、その結果得られた 6 番目のコミュニティを再度サブクラスタリングし、その結果得られた 1 番目のコミュニティを表している。特徴語から判断すると、コミュニティ 1 から 7 は全て政治や原発に関するコミュニティとなっている。そしてコミュニティ 8 と 10 は漫画やゲーム、ボーカロイドなどの趣味系、コミュニティ 9 は医療福祉系のコミュニティとなっている。

表 4.5 に、コミュニティごとの Twitter の基本データの中央値を示す。Twitter の基本データとして、ツイート数、フォロワー数、フォロー数、お気に入り数の 4 つのデータの中央値を用

いた。ツイート数やフォロワー数、フォロワー数では外れ値の影響が大きくなるため、平均値ではなく中央値を用いた。また、ネットワークの特徴量としては、クラスター係数、次数中心性、近接中心性、媒介中心性の4つの指標の平均値を用いた。炎上事例への参加人数が多いコミュニティでは、ツイート数やフォロワー数、フォロワー数が概ね多くなっているが、お気に入り数ではその傾向は見られない。しかしながら、趣味系のコミュニティであるコミュニティ10のように炎上事例への参加人数は多くないが、ツイート数は多いものも存在する。

また各コミュニティのネットワーク特徴量として、クラスター係数や中心性を表4.6に示す。Twitterの基本データの中央値におけるツイート数やフォロワー数、フォロワー数と同様に、炎上事例への参加人数が多いコミュニティではクラスター係数や次数中心性は高くなっている。その一方で、媒介中心性や近接中心性ではその傾向は見られない。また趣味系のコミュニティであるコミュニティ10は炎上事例への参加人数は多くないが、クラスター係数や次数中心性が高く、密なコミュニティであると言える。

4.3 コミュニティ情報とツイート時間帯分析

次にコミュニティごとのツイート時間帯を分析する。図4.1、図4.2、図4.3はそれぞれの事例における1日あたり（美味しんぼの事例は1時間あたり）のツイート数の推移である。表4.2の各事例のツイートの収集期間より、生活保護と献血の事例は収集期間がそれぞれ121日、153日であり、長期に渡る炎上の事例である。対して美味しんぼの事例は収集期間が4日であり、短期間に爆発的にツイート数が増加した短期間の炎上である。したがって、ツイート数は生活保護と献血の事例では1日単位で集計し、美味しんぼの事例では1時間単位で集計している。

まず図4.1の生活保護不正受給の事例では、政治系のコミュニティであるコミュニティ1(1.20)のツイート数が全期間にわたって多い。図から5月31日時点でツイート数が急増し、炎上が起きているのがわかるが、コミュニティ1では炎上が起きる以前からツイート数が多く、そのコミュニティでは生活保護不正受給の問題に対する関心が高かったことがわかる。炎上が起きてからの期間では、コミュニティ1と同じく政治系のコミュニティであるコミュニティ4(1.5.1)など他のコミュニティでもツイート数が増加している。

次に図4.2の群馬県議員による献血に関するツイート問題の事例では、原発系のコミュニティであるコミュニティ2(1.6.1)のツイート数が全期間にわたって多くなっている。この事例では炎上が起こったタイミングが3回あるということが、図に存在する3つの山からも見て取れる。最初の炎上のタイミングではコミュニティ1のツイート数が多いが、残り2つの炎上のタイミングではコミュニティ2のツイート数が多い。またコミュニティ2と同じく、特徴語に原発や福島などの単語を含む原発系のコミュニティであるコミュニティ5のツイート数も炎上のタイミングで多くなっている。

最後に漫画「美味しんぼ」における風評被害問題の事例であるが、この事例ではツイート収集期間が4日間となっているため、図4.3におけるツイート数は1時間あたりのツイート数とした。図よりツイート数が多くなっているタイミングが2つ存在することがわかるが、その周辺の時間では他の事例でもツイート数が多かったコミュニティ5、コミュニティ1、コミュニ

18 第4章 コミュニティ抽出と炎上事例分析

ティ4の3つのコミュニティのツイート数が多い。また炎上が起きる直前のツイート数を見ると、コミュニティ5のツイート数が最初に上昇し始めているので、このコミュニティから情報伝播が始まったと考えられる。

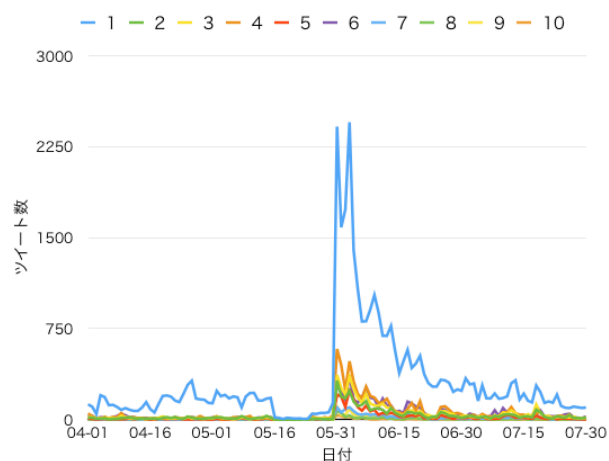


図 4.1. コミュニティごとの 1 日あたりのツイート数（生活保護）

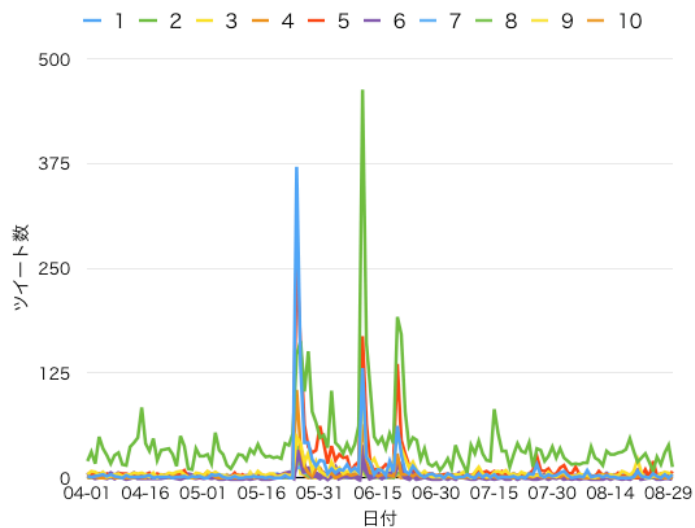


図 4.2. コミュニティごとの 1 日あたりのツイート数（献血）

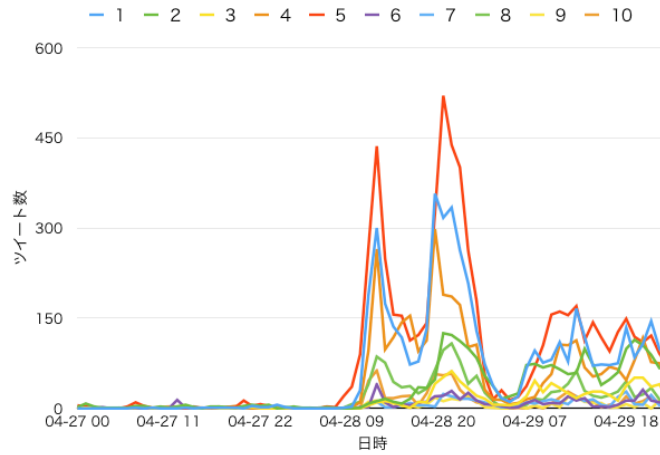


図 4.3. コミュニティごとの 1 時間あたりのツイート数（美味しんぼ）

4.4 コミュニティ情報とリツイート分析

本節では Twitter 上で他のユーザーのツイートをそのままツイートする機能である、リツイートに焦点を当てる。まず始めに各事例ごとのリツイートの分布や分散を比較し、事例ごとのリツイートにおける特徴を分析する。次に、前節のようにリツイートの時間帯分析を行う。

4.4.1 各事例におけるリツイートの分布と分散

まず始めに各事例ごとのリツイートの分布を比較する。図 4.4 は、各事例の全ツイートに占める、一定回数以上リツイートされたツイート件数の割合である。リツイート回数が 10 回ほどまでの少ない期間では、生活保護の事例での割合が 1 番高いが、リツイート回数がそれ以上のツイート件数からは美味しんぼの事例での割合が高くなっている。リツイート回数が 20 回以上ほどのツイート件数からは、生活保護よりも献血の事例の割合の方が高くなっていて、生活保護の事例では他の 2 つの事例よりもリツイート回数のばらつきが小さいことがわかる。

また表 4.7 は、各事例におけるリツイート回数の分散値である。美味しんぼの分散は 884.4 と 3 つの事例の中では非常に高い数値となっているが、これは外れ値の影響によるものである。美味しんぼの事例では、リツイート回数が 1 万回を超えるツイートが 1 件存在し、これを除いて分散値を計算し直すと、55.37 となる。したがって、リツイート分布の結果と同様にリツイート回数のばらつきの大きさは美味しんぼ、献血、生活保護の順になっている。

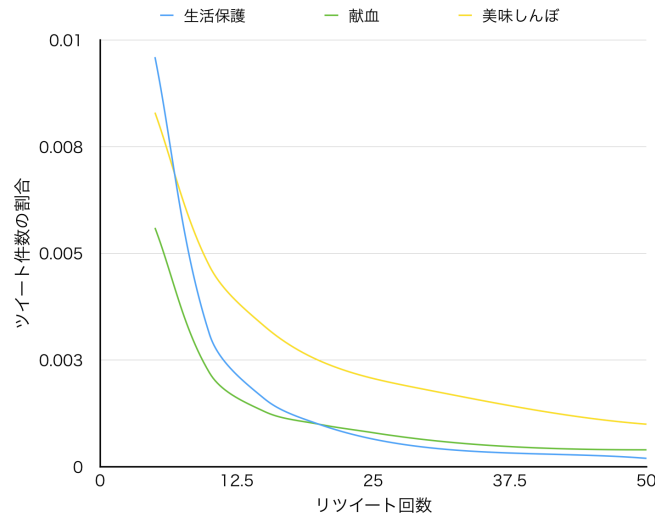


図 4.4. 一定回数以上リツイートされたツイート件数の割合

4.4.2 リツイート時間帯分析

次にコミュニティごとのリツイート時間帯を分析する。図 4.5, 図 4.6, 図 4.7 はそれぞれの事例における 1 日あたり（美味しんぼの事例は 1 時間あたり）のリツイート数の推移である。

図 4.5 の生活保護の事例では、ツイートの時間帯分析のときと同様に、政治系のコミュニティであるコミュニティ 1 (1.20) のリツイート数が全期間にわたって多い。炎上が起きてからの期間で、コミュニティ 1 と同じく政治系のコミュニティであるコミュニティ 4 (1.5.1) など他のコミュニティでもリツイート数が増加していることも同様である。炎上の全期間にわたってリツイートはツイート数の半分程度であることがわかる。

図 4.6 の献血の事例では、政治系のコミュニティ 1 と原発系のコミュニティ 2 のリツイート数が多いのはツイートの時間帯分析と同様であるが、ツイートの多かった原発系のコミュニティ 5 (1.6.2) のリツイートは少なく、原発、政治という特徴語を含むコミュニティ 3 (1.6.3, 黄色) のリツイート数が多くなっている。このことはコミュニティ 5 はリツイートは少なく、自分の意見などをつぶやいているツイートが多かったコミュニティであると考えられる。

図 4.7 の美味しんぼの事例では、コミュニティ 5, コミュニティ 1, コミュニティ 4 の 3 つのコミュニティのリツイート数が多いのはツイート数の分析と同様である。しかし図 4.3 と図 4.6 を比較すると、リツイート数ではツイート数ほど 3 つのコミュニティの間に差が見られなくなっている。したがってコミュニティ 5 はリツイートが少ないコミュニティであるが、反対にコミュニティ 1 とコミュニティ 4 はリツイートの多いコミュニティであることがわかる。またツイート数の分析と同様に、リツイート数が急激に増える直前でコミュニティ 5 のリツイート数が 1 番始めに増加していることがわかるので、コミュニティ 5 から情報伝播が発生したと

考えられる.

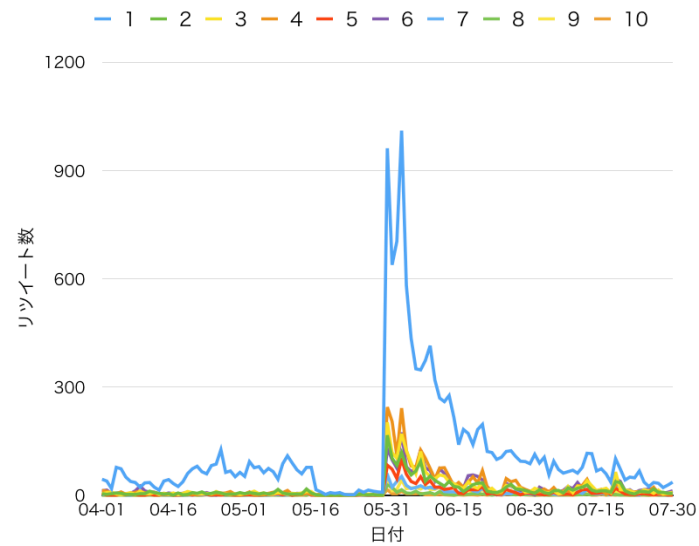


図 4.5. コミュニティごとの 1 日あたりのリツイート数（生活保護）

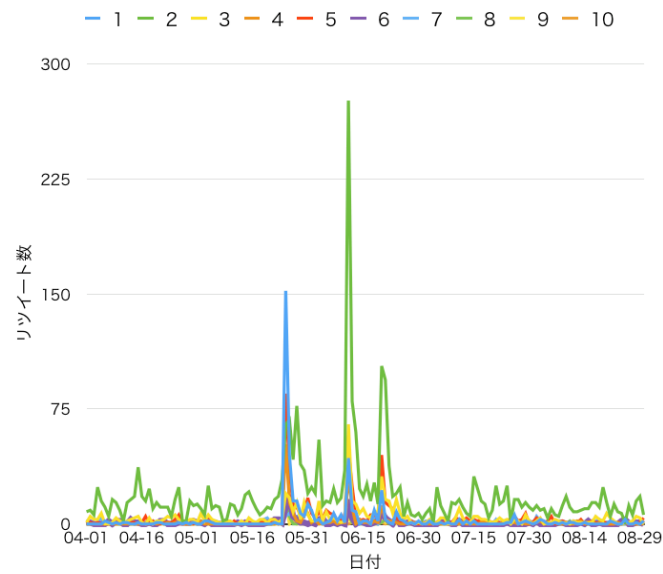


図 4.6. コミュニティごとの 1 日あたりのリツイート数（献血）

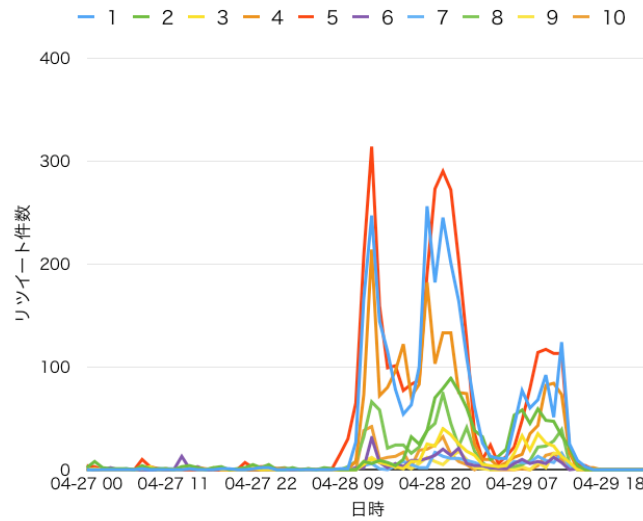


図 4.7. コミュニティごとの 1 時間あたりのリツイート数 (美味しんぼ)

4.5 まとめ

本実験では炎上事例についての知見とコミュニティについての知見の 2 つの知見が得られた。

4.5.1 炎上事例についての知見

今回の実験では生活保護, 献血, 美味しんぼの 3 つの事例を用いたが, それらの炎上事例は複数のタイプに分類されることがわかった。炎上のタイプの分類方法は 2 種類存在すると考えられる。まず一つ目は単純な方法であるが, 炎上の期間による分類である。生活保護と献血の事例はツイートの収集期間がそれぞれ 121 日, 153 日であり, 長期間の炎上事例であるが, 美味しんぼの事例はツイートの収集期間が 4 日であり, 短期間の炎上事例である。二つ目は炎上の原因が Twitter の内部に存在するか外部に存在するかという分類である。美味しんぼの事例では, リツイート回数が 1 万回以上である福島県民のツイートが炎上の引き金となった。献血の事例でも群馬県議員による「献血の車が止まっているけど, 放射能汚染地域に住む人の血って, ほしいですか」というツイートが 400 回以上リツイートされ, それを引用して自分のコメントを加えた著名人のツイートが 700 回以上リツイートされ, 炎上の引き金となった。しかしながら生活保護の事例では最もリツイートされたツイートで 300 回ほどであり, 炎上の原因は Twitter 上ではなく, 芸能人の記者会見であった。このことから生活保護の事例ではリツイート回数の分散が小さく, 散らばりが小さかったと考えられる。

4.5.2 コミュニティについての知見

本実験では炎上事例への参加人数が多いコミュニティで、ツイート数やフォロワー数、フォロワー数などの Twitter におけるデータや、クラスター係数や次数中心性などのネットワーク指標が高くなるという結果が得られた。この結果から Twitter におけるユーザーデータや、クラスター係数や中心性のコミュニティ情報は炎上への参加しやすさに影響すると考えられる。ツイート時間帯やリツイート時間帯分析では、コミュニティごとにツイートのタイミングが違っていたり、それぞれの事例で炎上を引き起こす中心となっているコミュニティが違っていた。また炎上によく参加するコミュニティやリツイートをしやすいコミュニティなど、コミュニティごとに Twitter の使用方法にも違いがあるという結果が得られた。以上の結果により、炎上とコミュニティにはやはり関係があり、5 章のバースト予測にコミュニティ情報は有効であると考えられる。

表 4.4. コミュニティごとの 3 つの炎上事例への参加人数とコミュニティ情報

コミュニティ	メンバー数	参加人数	特徴語
1. 1_20	9504 人	4326 人	OPI,nsc, 勅語, 田母神, 護持, 昌司, 領海 英霊,BKD, 中共, 悪用, 断交, 南北朝鮮, 皇紀 君が代, ポス, 尖閣, 参政, 尊皇,JNSC
2. 1_6.1	7692 人	2596 人	原発, 放射能, こと, 反対, 好き, 子供, 日本 情報, 事故,TPP, 避難,http, 移住, ため 未来, 汚染, 被曝, 福島,11, 自然
3. 1_6.3	4285 人	1795 人	原発, 政治, 日本, 好き, こと,http, 小沢 生活,TPP, 社会, 一郎, 反対, 国民, 趣味 フォロー, 増税, さん, よう, 情報, 音楽
4. 1_5.1	4882 人	1695 人	http, 経済, 政治, 日本, 研究, 軍事, アニメ 歴史, 社会, フォロー,RT, 主義,jp, ゲーム 好き,com, ネタ, 趣味, アカウント, さん
5. 1_6.2	4946 人	1682 人	福島, 好き, こと,http, 原発, フォロー, 出身 在住, 仕事, 最近, 趣味, 大好き, 音楽, 情報 さん, 震災, 日本, もの, 現在,com
6. 1_6.7	3514 人	1365 人	大阪, 日本, 政治, 好き, こと, 原発, 社会 共産党,http, 議員, 維新, 橋下, 趣味, 活動 フォロー, 生まれ, 仕事, 大好き, 在住, 問題
7. 1_6.5	3910 人	864 人	原発,http, 好き, こと, 音楽, 日本,com 活動, 反対, 東京, 自然, 映画, 社会, もの 仕事, 世界, よう, jp,11, たち
8. 1_7.1	6774 人	830 人	漫画, 連載, 発売,http, マンガ, アニメ, 特撮 仕事, ゲーム, 編集, フォロー, まんが, 好き お願い, 現在, コミック,com, さん, 最近,jp
9. 1_6.4	4113 人	746 人	障害, 好き, 自閉症, 発達, こと, 支援, 医療 息子, 福祉,http, 仕事, 社会, 知的, 趣味 大好き, フォロー, 日々, 現在, 看護, 音楽
10. 5_9	7962 人	728 人	SSG, ボカマス,futaba,sP,relations PBP, 崙麗, プロデュース, 宏美,ster, 法被 バンダイナムコゲームス,961, シンデマス,1105 アイマスプロデューサー, 真耶,dic, 東郷, 周子

表 4.5. コミュニティごとの Twitter の基本データの中央値

コミュニティ	ツイート数	フォロワー数	フォロー数	お気に入り数
1. 1_20	4,071	306	302	23
2. 1_6_1	4,269	275	276	122
3. 1_6_3	5,456	418	384	52
4. 1_5_1	6,771	223	194	156
5. 1_6_2	3,295	151	158	39
6. 1_6_7	2,529	252	224	21
7. 1_6_5	2,722	217	204	54
8. 1_7_1	2,942	109	116	23
9. 1_6_4	2,167	143	134	24
10. 5_9	6,723	142	164	28

表 4.6. コミュニティごとのクラスター係数と各中心性の平均値

コミュニティ	クラスター係数	次数中心性	媒介中心性	近接中心性
1. 1_20	0.1852	29.89	0.0031	0.2701
2. 1_6_1	0.1672	39.40	0.0040	0.3373
3. 1_6_3	0.1440	28.34	0.0079	0.2973
4. 1_5_1	0.1608	28.57	0.0102	0.2767
5. 1_6_2	0.1442	23.00	0.0096	0.2678
6. 1_6_7	0.1510	16.96	0.0101	0.2456
7. 1_6_5	0.1518	18.88	0.0070	0.2553
8. 1_7_1	0.1175	16.09	0.0076	0.2391
9. 1_6_4	0.1368	11.20	0.0071	0.2094
10. 5_9	0.2417	31.69	0.0152	0.2676

表 4.7. 3 つの事例におけるリツイート回数の分散値

	生活保護	献血	美味しんぼ
分散値	3.860	16.55	884.4

第 5 章

バースト予測モデルの構築

本章では、個々のツイートのバースト予測の実験に対する結果と考察を行う。4 章で抽出されたコミュニティでは、コミュニティごとにフォロー数などの Twitter の基本データに違いが見られ、ツイートやリツイートといった行動にも違いが見られた。この実験では抽出されたコミュニティ情報を使用することでバースト予測モデルの精度が上昇するかどうかを検証する。まず始めに SVM による特徴量の組み合わせによるバースト予測精度の結果と考察について述べる。次にロジスティック回帰モデルによる予測に有効な素性の結果と考察について述べた後、個々の事例による予測精度の結果と考察を行う。最後にまとめとして、全体の実験結果に対する考察を行い、提案手法の限界や応用についても議論する。

5.1 実験

5.1.1 目的

本実験の目的は、バースト予測モデルを構築し、使用した特徴量を比較することで、コミュニティの情報がバースト予測に有効であるかどうかを検証することである。また構築したモデルを解釈することによって、バーストに有効な特徴量についての知見を得ることも目的の一つである。

5.1.2 手順

表 5.1 にバースト予測モデルの構築に用いたデータを示す。教師あり学習では、正例と負例のデータ数が偏っているデータ（不均衡データ）を訓練データとすると、片方の結果を出力しやすいモデルを構築してしまう。これを回避するために、本研究ではランダムサンプリングを用いて、正例と負例の数を同数にした。ランダムサンプリングの特性上、使用する乱数値によって多少結果にばらつきが出る。そのため 10 回ランダムサンプリングを行い、10 個のデータセットにそれぞれモデルを構築し、テストを行って、その平均の値を取った。またツイートをを行ったユーザーが、相互リプライのリンクを持たない、ネットワークを構築した時期に Twitter を利用していなかったなどの理由で、コミュニティに属していないデータも存在する。本研究

では、それらのデータは欠損値として除去してから正例と負例を同数にしているため、表 5.1 のように使用データ数は元データ数から減少している。

表 5.1. パースト予測モデルの構築に用いたデータ

事例	元データ数	正例データ数	負例データ数	使用データ数
生活保護	45,756	3,357	42,399	4,854
献血	9,978	746	9,232	1,050
美味しんぼ	3,830	1,164	2,666	1,692
ALS	7,432	1,599	5,833	574
人工知能	7771	1,977	5,794	440
STAP 細胞	208,850	4,5976	162,874	10,042

3 章の提案手法で説明した分類モデルの特徴量をまとめ、特徴量のグループを名付けたものが表 5.2 である。TweetA はツイートの明確な数値で表現される属性による特徴量、TweetB はツイートの内容による特徴量、User はユーザーの特徴量、CommunityA はユーザーのコミュニティにおける影響度による特徴量、CommunityB はユーザーが所属するコミュニティによる特徴量である。

表 5.2. 実験に用いた特徴量

特徴量名	使用した特徴量
TweetA	ハッシュタグ数, メンション数, URL 数, 文字数, リプライかどうか
TweetB	ポジティブ率, 主観率
User	フォロー数, フォロワー数, お気に入り数, ツイート数
CommunityA	コミュニティ内でのクラスター係数, 次数中心性, 媒介中心性, 近接中心性 固有ベクトル中心性, PageRank, HubScore, AuthorityScore
CommunityB	コミュニティのメンバー数, 特徴語

まず考慮する特徴量と圧縮次元数を決定し、モデルを決定するため、相互メンションデータと同じ時期である 2012 年の事例を 2 つ含む、比較的古いデータである生活保護、献血、美味しんぼの 3 つの事例を用いて実験を行う。実験における特徴量の組み合わせは表 5.3 に示す。ベースラインは Yang らや Petrovic らの研究でも使われている TweetA と User を考慮したモデルとする [8][11]。次に TweetA と User と TweetB の特徴量の組み合わせ、TweetA と User と CommunityA の特徴量の組み合わせ、TweetA と User と CommunityA と CommunityB の特徴量の組み合わせの 4 つの特徴量の組み合わせを試す。ただし CommunityB の特徴量を追加する際は圧縮次元数を 10, 30, 50, 100, 200, 無圧縮の 6 通りに変更しながら実験を行う。SVM の実装は SVM 用のライブラリである LIBSVM を用いる [16]。

考慮する特徴量と圧縮次元数が決定された後、ロジスティック回帰モデルを構築し、それぞれの特徴量の偏回帰係数、オッズ比、P 値について考察を行う。データはパースト予測モデルの

表 5.3. 3 つの事例の全データにおけるバースト予測モデルの特徴量の組み合わせ

組み合わせ	使用した特徴量
1	TweetA + User
2	TweetA + User + TweetB
3	TweetA + User + CommunityA
4	TweetA + User + CommunityA + CommunityB(圧縮次元数 10)
5	TweetA + User + CommunityA + CommunityB(圧縮次元数 30)
6	TweetA + User + CommunityA + CommunityB(圧縮次元数 50)
7	TweetA + User + CommunityA + CommunityB(圧縮次元数 100)
8	TweetA + User + CommunityA + CommunityB(圧縮次元数 200)
9	TweetA + User + CommunityA + CommunityB(圧縮なし)

構築に用いた生活保護, 献血, 美味しんぼの 3 つの事例のツイートデータを用いた。

次に個々の事例に対するコミュニティ特徴量の影響や次元圧縮の影響を得るために, 6 つの事例それぞれに対してモデルを構築し, 性能をテストした。特徴量は表 5.4 に示される 4 つの組み合わせを用いている。

表 5.4. 6 つの事例の個々のデータにおけるバースト予測モデルのテスト結果

組み合わせ	使用した特徴量
1	TweetA + User
2	TweetA + User + CommunityA
3	TweetA + User + CommunityA + CommunityB(最適圧縮次元数)
4	TweetA + User + CommunityA + CommunityB(圧縮なし)

最後にバーストの予測モデルの評価方法について述べる。教師あり機械学習では訓練データとテストデータを分割するには, 交差検定 (Cross Validation) という手法が用いられる。Kohavi らは, 一般に 10 分割交差検定と呼ばれる手法が最もよいテスト手法であることを示した [17]。この手法は標本群を 10 個に分割し, そのうちひとつをテストデータとし, 残る 9 個を訓練データとする。その後, 10 個に分割された標本群それぞれをテストデータとして 10 回検証を行う。そうして得られた 10 回の結果を平均して 1 つの推定を得るという手法である。

また予測の精度を評価する指標には正解率 (Accuracy), 適合率 (Precision), 再現率 (Recall), F 値 (F-measure) の 4 つがある。予測と実際の結果によって, 予測結果は表 5.5 のように 4 種類に分けられるため, その 4 種類の予測結果を正しく評価するために以下の 4 つの指標が存在する。例えば, 迷惑メールのスパム判定のシステムを作る際, 迷惑メールでないメールが迷惑メールと判定されるのは好ましくないため, 適合率で評価するなど, タスクに応じて使用する評価指標を決定する。今回は単純な正解率である Accuracy を精度の評価指標として用いた。

表 5.5. 機械学習による予測結果の分類

	実際に正例	実際に負例
予測が正例	True Positive (TP)	False Positive (FP)
予測が負例	False Negative (FN)	True Negative (TN)

$$Accuracy = \frac{TP + TN}{TP + FP + TN + FN} \quad (5.1)$$

$$Precision = \frac{TP}{TP + FP} \quad (5.2)$$

$$Recall = \frac{TP}{TP + FN} \quad (5.3)$$

$$F - measure = \frac{2PrecisionRecall}{Precision + Recall} \quad (5.4)$$

5.1.3 データ

本実験では炎上事例のツイートデータとして、4章で用いたデータに加えて、比較的新しいデータである以下の3つの事例のデータを追加した。炎上事例のツイートデータの収集条件である検索キーワードと検索期間は表5.6の通りで、ツイート数やリツイート数、収集期間などの基礎データは表5.7に示す。その他のユーザーのプロフィールデータやユーザー間の相互メンションデータは4章で用いたものと同じである。

1. 人工知能学会紙表紙問題 (2013 年)
2. ALS アイスバケツチャレンジ (2014 年)
3. STAP 細胞に関する問題 (2014 年)

5.1.4 バーストの定義

本研究では、閾値を2つ設定し、リツイート回数が5回以上のものをバーストしたツイート、リツイート回数が1回以上かつ2回以下のものをバーストしていないツイートと定義した。本来リツイートされるはずのないツイートが含まれるのを防ぐため、リツイート回数が0回のツイートは除いた。表5.8は収集した炎上事例3つにおけるツイートのリツイート回数が1, 5, 10回以上のツイート件数を示したものである。事例全体を見ると、バーストの定義であるリツイート回数が5回以上のツイート件数は、リツイートされたツイート数のうちの8.0%で、リツイートされたツイートの中でも上位10%に入るほど多くリツイートされていることがわかる。またバーストの定義をリツイート回数が10回以上とすると、5回よりもリツイートされやすく

表 5.6. 6 つの事例におけるツイートの収集条件

事例	検索キーワード	検索期間
1	次長課長 河本 不正受給 生活保護 生ポ ナマポ 生ぽ なまぽ	2012/04/01 - 2012/08/31
2	庭山 献血 汚染地域	2012/04/01 - 2012/07/31
3	美味しんぼ スピリッツ おいしんぼ 原因不明の鼻血 風評被害 鼻血描写	2014/04/27 - 2014/04/30
4	人工知能	2013/12/16 - 2014/04/07
5	ALS アイスバケツ アイスバケツ ”Ice bucket” 氷水	2014/08/01 - 2014/09/02
6	STAP 理研 リケジョ 万能細胞 小保方 オボカタ おぼかた 笹井 若山 バカンティ 野依	2014/01/29 - 2014/06/30

表 5.7. 6 つの事例における収集したツイートの基礎データ

事例	ツイート数	RT 数	収集期間
1	35.0 万	10.6 万	121 日
2	13.4 万	2.95 万	153 日
3	14.1 万	5.29 万	4 日
4	19.0 万	8.88 万	113 日
5	64.1 万	8.10 万	33 日
6	441 万	182 万	153 日

特徴が明確なツイートが抽出されるので、タスクとしての難易度は下がる。しかしリツイート回数が 10 回以上のツイートは、リツイートされたツイート数のうちの 3.1% なので十分な学習データが得られない恐れがあるため、バーストの定義はリツイート回数が 5 回以上とした。また、4 章で説明したネットワークの構築に用いるユーザ間の相互メンションデータにおける総リンク数は 36,743,689 本、総ノード数は 5,980,977 個となっているため、1 ノードあたりの平均リンク数は 6.14 本となっている。したがって、バーストの定義のリツイート 5 回以上は 1 ノードあたりの平均リンク数とほぼ同程度となっているので、単純に考えるとバーストしたツイートは自分から直接リンクが繋がっていないユーザーからもリツイートされることになる。

5.1.5 パラメータ

用いたパラメータを以下の表 5.9 に示す。実験に用いたパラメータは、SVM で使用した正則化パラメータとコストである。カーネルは事前実験により RBF カーネル、線形カーネル、多項式カーネル、シグモイドカーネルの内、最もよい精度を示した RBF カーネルを使用し、正則化パラメータはグリッドサーチにより、正解率を最大にする解を求めた。

表 5.8. 3 つの事例におけるある回数以上のリツイートされたツイート件数

事例	RT1 回以上	RT5 回以上	RT10 回以上
生活保護	50,416	3,357	1,101
献血	10,824	746	292
美味しんぼ	4,428	1,164	657
合計	65,668	5,267	2,050

表 5.9. それぞれの事例における SVM のパラメータ

事例	正則化パラメータ	コスト
3 つの事例	32768.0	2.0
生活保護	8192.0	2.0
献血	512.0	8.0
美味しんぼ	8192.0	2.0
人工知能	8192.0	2.0
ALS	32768.0	0.03
STAP	32.0	8.0

5.2 特徴量の組み合わせによるバースト予測精度の結果と考察

5.2.1 結果

モデルを確定するために、5 章の炎上事例の横断分析で使用した 3 つの事例のデータをまとめて使用してバースト予測の SVM モデルを構築した。これらのデータはネットワークを構築した 2012 年の事例を 2 つ含む比較的古いデータである。実験は 4 章で述べたように、10 分割交差検定で Accuracy を評価して行った。正例負例の数を調整するためにランダムサンプリングを行い、データセットを 10 個作成しているので、Accuracy は 10 回施行した平均値となっている。

特徴量の組み合わせによる実験結果を表 5.10 に示す。ベースラインとして TweetA と User の特徴量を使用した実験では、精度が 59.81% という結果を得た。特徴量を加えた際に、精度が下がったものは TweetB のみであり、他の特徴量では精度は上昇している。中でもユーザーのコミュニティにおける影響度による特徴量を表す CommunityA を加えた際の精度の上昇量が一番高くなっていて、圧縮次元数は 10 が一番高い精度を示した。圧縮次元数が 200 以外のときは、次元圧縮をしないときよりも精度が高くなるという結果が得られた。

表 5.10. 3 つの事例の全データにおけるバースト予測モデルの実験結果

組み合わせ	使用した特徴量	Accuracy
1	TweetA+User	59.81833
2	TweetA+User+TweetB	57.93211
3	TweetA+User+CommunityA	64.53398
4	TweetA+User+CommunityA+CommunityB(圧縮次元数 10)	65.46544
5	TweetA+User+CommunityA+CommunityB(圧縮次元数 30)	65.20869
6	TweetA+User+CommunityA+CommunityB(圧縮次元数 50)	64.97302
7	TweetA+User+CommunityA+CommunityB(圧縮次元数 100)	64.85847
8	TweetA+User+CommunityA+CommunityB(圧縮次元数 200)	64.56618
9	TweetA+User+CommunityA+CommunityB(圧縮なし)	64.64981

5.2.2 考察

組み合わせ 1 の結果からベースラインの Accuracy は 59.81% となっているので、本研究でのバースト予測というタスクの難易度は約 6 割ほどであることがわかった。

また組み合わせ 1 と組み合わせ 2 の結果より、ツイートの内容による特徴量である TweetB は予測に対して有効な特徴量ではないということがわかった。同じツイートによる属性でも TweetA は予測に対して有効で、TweetB は有効でない結果になった理由としては、ユーザーがリツイートする際にはツイートの内容よりも、文字の長さやメンションや URL を含むかなどのツイートの形式や、誰がツイートしているかなどの情報の方を重視しているのではないかと考えられる。今回は複数の事例に対応できる汎用性の高いモデルを構築しようとしたため、ツイート中の単語の情報は使用せず、ポジティブ率と主観率の 2 つの特徴量しか考慮しなかったが、今後の課題として単語の情報を使うことも考えられる。否定的な言い回しや高圧的な態度が炎上を引き起こす原因となることも考えられるので、そういった特徴を表すことができる単語の情報はバースト予測モデルの構築に有効であると考えられる。

また CommunityA と CommunityB の特徴量を追加した実験では精度が上昇しているので、ユーザーのコミュニティ情報はバースト予測に有効な特徴量であったと考えられる。CommunityB の特徴量では次元圧縮を行った方が精度は上昇していて、圧縮次元数は 10 で精度が最大になっている。このことからコミュニティの特徴語は LSA によって次元圧縮をすることで、単語の情報をトピックのような単語の集合に変換することで汎用性を高めることができると考えられる。以上の実験の結果から、バースト予測モデルの特徴量は TweetA, User, CommunityA, CommunityB を使用し、圧縮次元数は 10 と決定した。

5.3 予測に有効な素性の結果と考察

5.3.1 結果

次に決定した特徴量と圧縮次元数を用いてロジスティック回帰モデルによるバースト予測モデルを構築した。データは前節の SVM モデルの構築に用いたデータと同じ 3 つの事例のツイートデータを使用している。表 5.11 に P 値が高い順に P 値が 0.1 以下の特徴量の偏回帰係数, オッズ比, P 値を示す。予測モデルの Accuracy は 53.7% であった。W1~W10 という特徴量はコミュニティの特徴語を LSA により 10 次元に圧縮したものである。表 5.11 を見ると, ツイート属性やユーザー属性, コミュニティ属性からまんべんなく有意な特徴量が得られている。特にオッズ比が高いものは, フォロワー数やツイートの文字数, HubScore, 次数中心性などが挙げられる。逆にオッズ比が低いものは, 近接中心性やツイートに含まれるメンション数, ユーザーのツイート数などが挙げられる。

表 5.11. P 値が 0.1 以下のロジスティック回帰モデルにおける特徴量

特徴量	偏回帰係数	オッズ比	P 値
W9	0.3234744	1.3819208	***
W8	-0.3837588	0.6812957	***
フォロワー数	0.2365654	1.2668905	***
W6	0.1323292	1.1414840	***
近接中心性	-0.6443947	0.5249802	***
文字数	0.2573781	1.2935341	***
HubScore	0.1592978	1.1726872	***
W3	-0.1150489	0.8913225	***
ハッシュタグ数	0.0826386	1.0861492	**
PageRank	0.1249866	1.1331333	**
メンション数	-0.0747095	0.9280130	**
お気に入り数	0.0956376	1.1003602	**
次数中心性	0.1552040	1.1678962	*
W7	-0.0770788	0.9258169	*
W10	-0.1597865	0.8523257	*
ツイート数	-0.0911629	0.9128690	*
クラスター係数	-0.0652864	0.9367991	.
媒介中心性	0.1047521	1.1104353	.

*** P 値 0~0.001 ** P 値 0.001~0.01 * P 値 0.01~0.05 . P 値 0.05~0.1

5.3.2 考察

予測モデルの Accuracy は 53.7% で, SVM の 65.4% と比較すると精度は下がっていて, 予想通り認識性能は SVM の方が優れていた. 有意な特徴量の中でオッズ比が 1 より大きいのは, W9, フォロワー数, W6, 文字数, HubScore, ハッシュタグ数, PageRank, お気に入り数, 次数中心性, 媒介中心性の 10 個である. これは特徴量が増加すると, バーストする確率も上昇することを意味していて, オッズ比が高いほどその上昇量は大きくなる. フォロワー数は自分のツイートがタイムラインに流れるユーザーの数を示すので, オッズ比が大きくなっている. 一方でフォロワー数は有意な特徴量とならなかったのは, フォロー返しを確約する bot や企業用アカウントなどの存在も大きいだろう. またコミュニティの特徴量である HubScore, PageRank, 次数中心性, 媒介中心性もオッズ比が高くなっているため, ツイートするユーザーがコミュニティでどんな立ち位置にいるのかということもバースト予測に有用な特徴量であるといえる.

有意な特徴量の中でオッズ比が 1 より小さいのは, W8, 近接中心性, W3, メンション数, W7, W10, ツイート数, クラスター係数の 8 個である. これは特徴量が増加すると, バーストする確率は減少することを意味していて, オッズ比が低いほどその減少量は大きくなる. メンション数のオッズ比が低くなっているのは, ツイートにメンションを含むと特定のユーザーを名指しすることになるのでツイートの一般性が低くなるからだと考えられる. またクラスター係数のオッズ比が低くなっている理由は, 高校生や友達などの密なコミュニティよりも, Twitter を情報収集のツールとして使っているコミュニティでバーストが起きやすいからではないかと考えられる.

5.4 個々の事例による予測精度の結果と考察

5.4.1 結果

最後に 4 章で述べた 6 つの個々のデータセットを対象に, 決定した特徴量と圧縮次元数を用いて SVM モデルを構築し, 10 分割交差検定によるテストを行った. CommunityA の特徴量による効果や CommunityB の特徴量による効果, 次元圧縮による効果を確認するために, 表 5.4 における組み合わせ 1, 組み合わせ 2, 組み合わせ 3, 組み合わせ 4 の 4 種類の実験を行った. その結果を表 5.12 に示す. 実験 1 と実験 2 の結果を比較すると, どの事例においても精度は上昇しているので, CommunityA のコミュニティ内でのユーザーの影響度を表す特徴量を考慮すると精度が向上している. しかし上昇幅にはばらつきがあり, ネットワークを構築した 2012 年のツイートデータである事例 1, 2 の方が上昇幅は大きくなっている. 次元圧縮の効果は 6 つの事例中 4 つの事例で無圧縮の特徴語よりも高い精度を示しているが, 事例によっては次元圧縮をしない方が精度が高く, 特徴語を入れると精度が落ちる事例も存在する.

表 5.12. 6 つの事例の個々のデータにおけるバースト予測モデルのテスト結果

事例	組み合わせ 1	組み合わせ 2	組み合わせ 3	組み合わせ 4	年度
生活保護	57.77297	63.09641	63.33816	63.92872	2012
献血	66.59047	71.43810	70.47666	69.37143	2012
美味しんぼ	68.72338	70.54967	70.53223	68.20330	2014
人工知能	64.86063	65.34843	64.36299	65.43555	2013～2014
ALS	65.22728	67.75000	68.10934	66.20453	2014
STAP	60.32066	62.43277	64.69275	64.85860	2014

5.4.2 考察

実験 1 と実験 2 の結果は、6 つの事例全てにおいてコミュニティでのユーザーの立ち位置を特微量として考慮するとバースト予測の精度が上がることを示している。上昇幅がばらついてしまう理由としては、ソーシャルメディアでのネットワーク情報は絶えず変化してしまうのでリアルタイムなネットワーク情報を得ることが難しいことが挙げられる。また、4 章のデータの前処理で述べたように、コミュニティに属していないユーザーが存在することも理由の一つと考えられる。特に新しいデータである人工知能と ALS の事例では学習データ数が 500 程度となってしまったので、学習が上手くいかなかった可能性も考えられる。

また実験 4 と実験 9 の比較から、次元圧縮については 6 事例中 4 事例で精度が向上した結果となったが、精度が最大となる圧縮次元数は事例によって違っていて、一概に 10 が最適であると言えない結果であった。今回はコミュニティの特徴語を tf-idf 値の上位 20 単語として、出現するかしないかのバイナリーデータとしたが、tf-idf 値をそのまま用いたり、上位 20 単語ではなく、上位 50 単語や全単語として実験を行うと更なる精度の向上が期待できる可能性がある。

5.5 まとめ

実験結果を簡潔にまとめると、ツイートの内容による特微量はバースト予測に有効でなかったが、一方でコミュニティ内でのユーザーの影響度による特微量やコミュニティ自体による特微量は有効であった。この結果によって、ツイートのリツイートやバーストは、どんな内容をツイートするかよりも、ツイートするユーザーやそのユーザーのコミュニティにおける立ち位置やコミュニティの情報が重要であると考えられる。またツイートによる属性の中でも、ツイートの内容は有効な特微量とはならなかったが、文字数やメンション数などの数値で表現される特微量は、ツイートの体裁などの見栄えに関わる部分であるために有効な特徴であったと考えられる。しかしながらツイートの内容や、コミュニティの特徴語による特微量の生成手法にはまだ改善の余地がある。本研究ではツイートの単語はモデルの汎用性を考慮して使用しなかったが、高圧的な態度や否定的な言い回し、ツイート中の言葉遣いなどのバーストに関わりそう

な特徴量を生成する手法として、df が一定以上となる多くの文書に出現する単語のみを用いるなどの方法が考えられる。また次元圧縮は事例によっては有効であることが示されたが、コミュニティの特徴語の特徴量として tf-idf 値をそのまま用いたり、上位 20 単語ではなく、上位 50 単語や全単語として実験を行うなどの方法が考えられる。

提案手法ではコミュニティ情報を用いているが、それによって提案手法の限界も存在する。個々の事例によるバースト予測精度の結果を比較すると、ネットワークを構築した 2012 年の事例のツイートデータを使用した生活保護と献血の事例ではコミュニティによる特徴量を追加したケースの精度の上昇量が多い。その一方で 2014 年のデータである残り 4 つ事例ではコミュニティによる特徴量を追加したケースの精度の上昇量が小さくなっている。手法の限界として、ソーシャルメディアでのネットワーク情報は絶えず変化してしまうのでリアルタイムなネットワーク情報を得ることが難しいことが挙げられる。また、4 章のデータの前処理で述べたように、コミュニティに属していないユーザーが存在し、データが欠損してしまうことも同様に手法の限界として考えられる。

しかしながら本研究で構築されたバースト予測モデルはコミュニティ情報を事前に計算しているので、対象となるツイートとツイートを行ったユーザーのデータがあれば簡単に実行できる。またロジスティック回帰によって得られた特徴量のオッズ比によって、リツイートされやすいツイートに含まれる特徴量やリツイートされづらいツイートに含まれる特徴量が把握できるので、Twitter を利用する際に状況に合わせてツイートの形式を変化させるなどの応用が考えられる。

第 6 章

結論

本研究では Twitter におけるユーザーコミュニティの抽出とそのコミュニティを用いた炎上事例分析、バースト予測モデルの構築を行った。抽出されたコミュニティでは炎上事例への参加人数が多いコミュニティでは、ツイート数やフォロワー数、フォロワー数などのユーザー属性や、クラスター係数や中心性などのコミュニティ属性が高いという特徴が見られた。またコミュニティを用いた炎上事例分析では、事例や時間帯によってコミュニティごとのツイートやリツイートの行動に違いが見られた。炎上事例の分類方法に焦点を当てると、炎上事例は期間が長い事例と短い事例のような分類もできるが、炎上の原因が Twitter 上に存在するか、外部に存在するかという分類方法も考えられ、それによってリツイートの分布やリツイート回数の分散値が異なるということを明らかにした。バースト予測モデルの構築では、コミュニティ情報がバースト予測に有効であったが、ツイートの内容による特徴量はバースト予測に有効でないという結果が得られた。またロジスティック回帰モデルによって得られたオッズ比は、リツイートされやすい、またはされづらいツイートにおける特徴量を示している。これはツイートのバースト以外の応用も考えられ、例えば商品情報等のマーケティングに Twitter を利用する場合、リツイートされやすいように文字数を多くしたり、日頃からフォロワー数を増やす努力をするなどの応用が考えられる。

また本研究では学習データが不足するのを回避するために、バーストの定義はリツイート回数が 5 回以上と定義したが、データ数が十分であれば 10 回以上とする選択肢も考えられる。リツイート回数が 10 回以上のツイートは、5 回以上のツイートよりもリツイートされやすく特徴が明確なのでタスクとしての難易度は下がる。バーストの定義をリツイート 10 回以上として数回実験してみた結果、5 回以上のときよりも精度は 5% から 10% ほど上昇した。応用として実際の利用シーンを想定すると、大きなバーストを確実に検知したい場合はバーストの定義を 10 回以上に設定し、小さなバーストでも取りこぼしのないように検知したい場合はバーストの定義を 5 回以上に設定するなど、使用目的に合わせてバーストの定義を設定すると良い結果が得られるであろう。

炎上は情報過多や情報の偏りなど多くの弊害をもたらすが、炎上が起きてしまった際は組織として炎上に早急に対応できる体制が必要である。現実で沈静できる情報を用意し、ウェブ上に的確なタイミングや形で情報を投入させるなどの早急な対応が必要となる。本研究で得られ

38 第 6 章 結論

た炎上の特徴や、バーストやリツイートされやすい属性についての知見が今後の炎上事例の沈静や企業にとっての組織の整備に繋がれば幸いである。

謝辞

本研究を行うにあたり、非常に多くの方々にご指導、ご鞭撻を賜りましたことを感謝いたします。

指導教官である坂田一郎教授には、お忙しい身でありながら貴重な時間を割いていただき、日頃から多くのご指導、ご鞭撻を賜りましたことを心より感謝します。

森純一郎准教授には研究へのアドバイスから、研究に対する心構えまで色々にご指導くださって感謝しております。また卒論のベースとなる基礎知識から論文執筆に至るまで、ご指導、ご鞭撻を賜りましたことを感謝いたします。

坂田・森研究室秘書の粥川敬子さん、佐藤妙子さん、石原絢さんには、快適に研究を行う環境を整えていただき感謝しています。

株式会社ホットリンク R&D 部門統括兼東京大学工学系研究科客員研究員である榊剛史氏には、研究のテーマ設定から、データの提供、具体的なアドバイスを賜りましたことを心より感謝します。またお忙しい身でありながら、毎週時間を割いてミーティングを設定してくださったこと、論文執筆のアドバイスを頂いたことをこの場を借りて深く感謝します。

博士過程3年の丸井淳己氏には、研究を進めるにあたっての的確なアドバイスを賜りましたこと、サーバーやネットワーク環境の管理まで行っていただき、快適な研究環境を整備してくださったことを心より感謝します。

また研究室の先輩方や同期の仲間には、研究のアドバイスを貰いつつも、楽しい時間を過ごすことができたことを深く感謝します。修士課程2年の澤村氏、早嶋氏には研究の合間に気軽にお声掛けしてくださったり、優しく接していただき、楽しい研究生活を送ることができました。修士課程1年の株田氏、山下氏には毎回の研究会でアドバイスを頂き、卒論の時期には毎回進捗を気に掛けていただいたことを感謝いたします。同期である小林氏、サンタモン氏、河津氏、上子氏とはお互いに切磋琢磨し、研究を進めることができたことを感謝いたします。

改めて皆様に深く感謝すると共に、今後のご多幸をお祈り申し上げます。

参考文献

- [1] Junki Marui, Nozomi Nori, Takeshi Sakaki, and Junichiro Mori. Empirical Study of Conversational Community Using Linguistic Expression and Profile Information. In *Active Media Technology*, Vol. 8610, pp. 286–298. Springer International Publishing, 2014.
- [2] 鳥海不二夫, 榊剛史, 岡崎直観. 「人工知能」の表紙に関するツイートの分析・続報. 第4回 Web インテリジェンスとインタラクション研究会, 2014.
- [3] A. Java, X. Song, T. Finin, and B. Tseng. Why we twitter: understanding microblogging usage and communities. In *WebKDD/SNA-KDD '07 Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 workshop on Web mining and social network analysis*, pp. 56–65, New York, NY, USA, 2007. ACM.
- [4] 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, WWW '10*, pp. 591–600, New York, NY, USA, 2010. ACM.
- [5] B. Huberman, D. M. Romero, and F. Wu. Social networks that matter: Twitter under the microscope. *First Monday*, Vol. 14, , 2009.
- [6] E. Bakshy, J. M. Hofman, W. A. Mason, and D. J. Watts. Everyone’s an Influencer: Quantifying Influence on Twitter. *Proceedings of the Fourth ACM International Conference on Web Search and Data Mining, WSDM ’11*, pp. 65–74, New York, NY, USA, 2011. ACM.
- [7] V. D. Blondel, J.-L. Guillaume, R. Lambiotte, and E. Lefebvre. Fast unfolding of communities in large networks. *Journal of Statistical Mechanics: Theory and Experiment*, Vol. 2008, No. 10, 2008.
- [8] J. Yang and S. Counts. Predicting the Speed, Scale, and Range of Information Diffusion in Twitter. In *Fourth International AAAI Conference on Weblogs and Social Media*, 2010.
- [9] D. Cox, D. R. & Oakes. *Analysis of survival data*. Chapman & Hall, London, 1984.
- [10] K. Crammer, O. Dekel, J. Keshet, S. Shalev-Shwartz, and Y. Singer. Online passive-aggressive algorithms. *The Journal of Machine Learning Research*, pp. 551–585, 2006.
- [11] S. Petrovic, M. Osborne, and V. Lavrenko. RT to Win! Predicting Message Prop-

- agation in Twitter. In *Fifth International AAAI Conference on Weblogs and Social Media*, 2010.
- [12] Qiming Diao, Jing Jiang, Feida Zhu, and Ee-Peng Lim. Finding bursty topics from microblogs. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 536–544, Jeju Island, Korea, 2012.
 - [13] Thomas K Landauer and Susan T. Dutnais. A solution to plato’s problem: The latent semantic analysis theory of acquisition, induction, and representation of knowledge. *Psychological review*, pp. 211–240, 1997.
 - [14] Corinna CORTES. Support-vector networks. *Machine Learning*, Vol. 20, pp. 274–297, 1995.
 - [15] Jure Leskovec and Rok Sosič. SNAP: A general purpose network analysis and graph mining library in C++. <http://snap.stanford.edu/snap>, 2014.
 - [16] Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, Vol. 2, pp. 27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
 - [17] Ron Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, pp. 1137–1143, 1995.