推薦論文

希少性および影響力に着目した マイクロブログからの犯罪関連投稿抽出

高橋 柊^{1,a)} 菊地 悠¹ 落合 桂一¹ 深澤 佑介¹

受付日 2016年12月9日, 採録日 2017年5月16日

概要:近年, Twitter などのマイクロブログに代表される SNS サービスのユーザ数増加により, 実世界における情報がリアルタイムに web 上にアップロードされるようになった。そのため, マイクロブログ上の情報をセンシングすることで, 実世界における事象を検知する研究が活発となっている。本研究では, マイクロブログ上に投稿される犯罪関連投稿に着目する。マイクロブログより, リアルタイムな犯罪関連投稿を抽出することができれば, 犯罪事象に対し短時間で適切な防犯対策が可能となる。また, 警察官によるパトロールなど既存のセンシング手法では抽出困難であった犯罪事象がマイクロブログ固有の情報より抽出可能となる。マイクロブログにおける投稿抽出には投稿テキスト情報を利用する手法があるが, 犯罪という希少な事象を投稿テキスト情報から抽出するのは困難である。提案手法ではユーザの投稿内容, 投稿関連位置情報および関連ステータス情報を用いることで, 投稿内容の希少性および影響力について特徴量を生成し, 投稿者が経験あるいは目撃した犯罪関連投稿を抽出する。Twitter の日本語投稿データに対し提案手法を適用したところ, 投稿テキスト情報のみを利用した既存手法 (AUC=0.6146) に対し, 提案手法を適用したところ, 投稿テキスト情報のみを利用した既存手法 (AUC=0.6146) に対し, 提案手法を適用したところ, 投稿テキスト情報のみを利用した既存手法 (AUC=0.6146) に対し, 提案手法

キーワード:Twitter, イベント抽出, データマイニング

Exraction of Criminal Related Posts from Microblogs based on Rarity and Influence

SHU TAKAHASHI^{1,a)} HARUKA KIKUCHI¹ KEIICHI OCHIAI¹ YUSUKE FUKAZAWA¹

Received: December 9, 2016, Accepted: May 16, 2017

Abstract: Recently, with increase of the SNS users such as Twitter, real world information is uploaded on the web in real time. By using the proposed method, it is possible to extract criminal related posts experienced or witnessed based on rarity and influence from users' posted contents, associated location information and status information. Crime events that were difficult to detected by existing methods such as patrol by police officers can be extracted from unique information of microblogs. By using the proposed method, it is possible to extract criminal related posts experienced or witnessed from users' posted contents, associated location information and status information. The proposed method applied on Japanese Twitter data shows higher AUC score (0.7183) compared to the existing methods which only use text information (0.6146).

Keywords: Twitter, event detection, data mining

1. はじめに

近年, Twitter などのマイクロブログに代表される SNS

サービスのユーザ数増加により、実世界における情報がリアルタイムに web 上にアップロードされるようになった。リアルタイムに市民が情報を投稿するソーシャルメディア

本論文の内容は 2016 年 5 月の第 79 回モバイルコンピューティングとパーベイシブシステム研究会にて報告され,同研究会主査により情報処理学会論文誌ジャーナルへの掲載が推薦された論文である.

¹ NTT ドコモ NTT DOCOMO, INC., Yokosuka, Kanagawa 239–8536, Japan

a) syuu.takahashi.us@nttdocomo.com

は、既存のセンシングデバイスでは抽出困難なリアルタイム情報を多く含むため、社会的な事象が大きく反映される。これらの研究は Human Probe(HPB)と呼ばれ、社会およびビジネス分野への応用が進んでいる。

HPB における代表的な研究分野としてイベント検出がある. Chen ら [1] は画像投稿コミュニティサイトである Flickr*1の投稿情報を用いることで、イベントを検出する手法について検討している. しかし、投稿数の出現回数に依存するため、ごく少数の投稿からイベントを抽出することは困難である. 山田ら [2] は、Twitter*2からイベント情報を自動的に抽出する手法について検討している. Twitterを情報源として活用することで、多様かつ大量なイベント情報の抽出が可能となることを示している.

本研究では、Twitter に投稿される犯罪関連投稿に着目する。平成24年度の法務省調査では、アンケート調査より得られた個人犯罪被害のうち、暴行・脅迫では56.8%、性的事件では74.1%が警察機関などに届出がされていない*3. Twitter の投稿情報より投稿者が経験あるいは目撃した犯罪関連投稿を抽出することができれば、警察によるパトロールなど既存のセンシングでは顕在化困難であった犯罪事象についての知見を得ることができる.

日本において市民が犯罪にあう経験はごく希少な確率で発生すると考えられる。提案手法ではごく希少な確率で発生するという犯罪の特性を利用することで,犯罪関連投稿を抽出する。Twitter はリアルタイム性の高いソーシャルメディアであり,投稿者の投稿頻度が他のソーシャルメディアと比べ高い。そのため,犯罪事象が起こることで投稿者の投稿内容あるいは位置における語の出現分布に違いが起こることが期待できる。本研究では希少性および投稿の影響力を特徴量化することで,マイクロブログから犯罪関連投稿を抽出する手法を提案する。

本研究の貢献は以下のとおりである.

i) 希少性および影響力の特徴量化

犯罪という希少性の高い投稿を抽出するための投稿者 希少性および位置希少性を提案する.投稿者のステー タス情報を用い,対象投稿の影響力の特徴量化を行う. また,提案特徴量が犯罪関連投稿の分類に有効である ことを示す.

ii) 提案特徴量による犯罪関連投稿抽出

本研究で提案した特徴量を用いて、投稿者が経験あるいは目撃した犯罪関連投稿をマイクロブログより抽出する.投稿テキスト情報のみを用いた既存手法と提案手法との分類精度を比較し、提案手法が有効であることを示す.

本稿では、次章にて関連研究について述べ本研究との差

分について説明する.次に、3章にて提案手法の詳細を説明する.4章では実データを用いた実験を行い、既存手法と提案手法の精度について検証する.最後に、5章にて本研究のまとめおよび今後の課題について述べる.

2. 関連研究

本研究と関連する研究には SNS から現実世界の事象を検出するソーシャルセンシングと、地域の特性を統計データから推定する Geographic profiling がある。本章では、これらの関連研究について説明するとともに、本研究との差異について述べる。

Twitter などの SNS データから現実世界における事象の 説明を試みる研究は、1章であげたイベント検出に関する研 究をはじめ多く行われている [3], [4], [5]. Watanabe ら [6] は位置情報の少ない小規模なローカルイベントの検出につ いて検討している. 位置情報のついていない投稿に対して 投稿テキストより位置情報を推定することで, ローカルイ ベントを検出する手法を提案している. Sakaki ら [7] は災 害時に被災地とその他の地域で Twitter の利用に差がある ことを示し, 災害の発生場所を早期に推定する手法を提案 している. これらの研究は任意の災害やイベントに関連し た多数の投稿を集約し、投稿内容の変動を利用することで イベント検出を行う. しかし, 犯罪関連事象は投稿者が単 体で経験あるいは目撃することが考えられ、必ずしも複数 人が同一の犯罪関連事象を経験あるいは目撃するとは限ら ない. そのため、規模の大きいテロなどの犯罪関連事象に ついては適用可能であると考えられるが、多くの犯罪事象 を説明することは困難である. 本研究では, 多数の投稿に よる検出ではなく、犯罪という特有の事象が持つ希少性お よび影響力を考慮することで犯罪関連投稿の抽出を行う.

Liら [8] は結婚、就職、出産などのプライベートなライフイベントを Twitter 上の投稿情報から抽出する手法について検討している。プライベートなライフイベントは、パブリックなイベントと比較し関連投稿数の増加が少ないため抽出が困難である。そこで、祝辞やなぐさめなどライフイベントに対する反応が含まれる返信投稿に着目することで、高精度にライフイベントを抽出することが可能であることを示している。一方で、多種のライフイベントについては返信内容を利用したイベント抽出が困難であるとし、投稿テキストより抽出した特徴に対し、ライフイベントごとに人手でラベル付けした情報を用いる手法を提案している。そのため、一部の犯罪事象に関連した投稿を抽出することは可能であるが、多種の犯罪事象を抽出するには犯罪事象ごとにラベリング作業が発生してしまうため抽出が困難である。

古川ら [9] は Twitter 上の犯罪関連投稿を経験・公的・ 参照投稿に分類し、経験投稿の自動抽出手法について検討 している。一般市民の投稿者が現場で経験したものを経験

^{*1} Flickr, https://www.flickr.com/

^{*2} Twitter, https://twitter.com

^{*3} 平成 24 年版 法務省:犯罪白書 第 5 編/第 3 章/第 2 節/2

投稿,公共組織などのアカウントによる投稿を公的投稿, ニュースやブログの引用をともなう投稿を参照投稿とし, 経験投稿が他の犯罪情報と異なる性質を持っていること を明らかにし,既存センシングでは抽出困難な犯罪事象が Twitter上に存在していることを示している。また,投稿 テキストの Bag-of-Words を用い経験投稿と他の投稿を機 械学習により分類しているが,単純な単語素性のみでは判 別が困難であることを示している。そのため研究目的は本 研究と同じであるが,希少性や影響力など投稿テキスト以 外の情報を考慮した本研究の提案手法とは試みが異なる。

位置および犯罪事象との関連性を分析する Geographic profiling [10] は、犯罪の防止および捜査に利用されている. Geographic profiling により、犯罪の空間的な動作を把握し、重点的に治安維持活動を行う地域を特定することができる。また、ソフトウェアを用いることで、多量のデータから効率的に Geographic profiling を行うアルゴリズムおよびシステムが提案されている [11]、[12]、[13].

Geographic profiling では、一般に捜査活動より得られ たデータや過去の犯罪統計情報を利用することで位置およ び犯罪事象との関連性を分析する. ソーシャルメディア上 のデータより犯罪関連情報を抽出することができれば、よ り多様なデータから高度な Geographic profiling が可能と なる. White ら [14] は Twitter 投稿情報およびソーシャル グラフを用いた Geographic profiling の有効性について検 討している. あらかじめ設定された犯罪行為者のソーシャ ルグラフを用いた分析や、キーワードの出現数変化をとら えることで、Twitter が Geographic profiling に有効であ ることを示している. また, Wang ら [15] は LDA および 線形モデルを用いることで、Twitter の投稿情報からひき 逃げ事件の予測が可能であることを示している.しかし, これらの研究は犯罪関連投稿が経験投稿であるか公的・参 照投稿であるかについては判別していないため, 投稿者が 犯罪を経験あるいは目撃した犯罪関連投稿を抽出するとい う本研究の目的とは異なる.

提案手法

本研究の目的は、Twitter の投稿情報より投稿者が犯罪を経験あるいは目撃した投稿を抽出することで、警察官のパトロールなど既存のセンシングでは顕在化が困難であった犯罪関連投稿を抽出することである。図1に提案手法における分析ステップを示す。

本提案手法は以下6ステップにより構成される.

i) 犯罪語マッチ処理

Twitterより取得した投稿および、犯罪関連語を含む 犯罪語辞書を用い、投稿テキスト中に犯罪語辞書の単 語を含む投稿を抽出する.

ii) 投稿者希少性を考慮した特徴量計算

犯罪語マッチより得られた投稿の投稿者が、過去に同

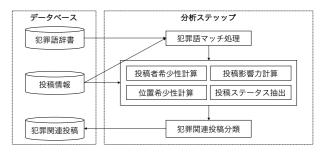


図1 分析ステップ

Fig. 1 Analysis step.

様の犯罪事象について投稿しているか否かを計算する.

iii) 位置希少性を考慮した特徴量計算

犯罪語マッチより得られた投稿における関連位置において,過去に類似した投稿がされているか否かを計算する.

iv) 投稿影響力を考慮した特徴量計算

犯罪語マッチより得られた投稿の影響力について計算 する.

v) 投稿ステータス抽出

投稿テキスト中に含まれる URL リンク内容を抽出する.

vi) 犯罪関連投稿の分類

計算された特徴量を利用し、犯罪関連投稿を抽出する 分類器を構築する.

3.1 犯罪語マッチ処理

投稿内容が犯罪関連語を含んでいる投稿を抽出したい. 投稿テキスト中に犯罪関連語を含むか否かを判定するため に,犯罪語辞書を作成した.**表1**に作成した犯罪語辞書を 示す.日本において強盗や殺人などの重犯罪は,注目度が 高く,警察庁などが統計情報を公表している.本研究では, ソーシャルメディア特有の犯罪関連投稿を抽出したい.

個人犯罪被害のうち、警察機関に届けられていない割合が高いものに暴行・脅迫および性的事件がある。そのため、本研究では比較的刑事罰が弱く、また通報されることが少ないと考えられる性的事件や暴行・脅迫事件および迷惑行為を対象とした。犯罪語辞書に登録する犯罪語には、「痴漢」「ストーカー」などを検索語として抽出された投稿に含まれる形態素のうち、頻出度が高く、犯罪関連投稿の検索に有効であると思われる30語を登録した。

投稿情報および、犯罪語辞書を用い、犯罪関連語を含む 投稿情報のみを抽出する. Twitter の投稿テキストでは、 一般に漢字で表記されるような語についても、ひらがなや カナカナを用いて投稿されるなど、表記ゆれが存在する.

図 2 にキーワードマッチングの実施例を示す. たとえば,「痴漢, ちかん, チカン」はすべて同様の犯罪事象について投稿しているが,「痴漢」のみをキーワードマッチングの条件とすると,「ちかん, チカン」については, 除外

表 1 犯罪語辞書

Table 1 Dictionary of criminal related words.

犯罪語	ルビ
警察	ケイサツ
喧嘩	ケンカ
ナンパ	ナンパ
変態	ヘンタイ
キャッチ	キャッチ
警官	ケイカン
すり	スリ
絡ま	カラマ
酔っ払い	ヨッパライ
騒ぎ	サワギ
ストーカー	ストーカー
痴漢	チカン
盗撮	トウサツ
パトカー	パトカー
通報	ツウホウ
スカウト	スカウト
サイレン	サイレン
覗き	ノゾキ
泥棒	ドロボウ
客引き	キャクヒキ
けが人	ケガニン
掴まれ	ツカマレ
万引き	マンビキ
野次馬	ヤジウマ
ひったくり	ヒッタクリ
変質者	ヘンシツシャ
付き纏	ツキマト
セールス	セールス
置引	オキビキ
ストーキング	ストーキング

板橋駅 の ホーム で 痴漢された		
形態素解析		
板橋駅 / の / ホーム /	。 で/痴漢/さ/れ/た	
	ルビ抽出	
イタバシエキ/ ノ / ホーム	/ デ / チカン / サ / レ / タ	

図 2 キーワードマッチングの実施例

Fig. 2 Example of keyword matching.

されてしまうという問題がある. そこで, 投稿テキストを JTAG [16] を用い形態素解析し, 各形態素をルビに変換したのち, キーワードマッチングを行う.

Twitter に投稿されるツイートには、多くの場合位置情報が付与されていない。そこで、投稿内容に関連した位置情報を新たに付与する必要がある。投稿者が犯罪を経験あるいは目撃したことについて投稿した場合、投稿テキスト中に投稿関連位置に関する語が出現することが期待できる。地名、施設名などの Point-of-Interest (POI) に関する名称および緯度経度との対応関係を POI-DB として持ち、対象テ

板橋駅 / の / ホーム / で / 痴漢 / さ / れ / た POI-ID 000001 NAME 板線駅 lat.lng 1,747,788 1,951,77578

図 3 位置情報抽出の実施例

Fig. 3 Example of location information extraction.

キストを形態素解析し、形態素と POI-DB のマッチングを行う. テキストからの POI 抽出では、Geo/Non-geo 曖昧性 および Geo/Geo 曖昧性が問題となる [17]. Geo/Non-geo 曖昧性とは地名と同じ表記で地名以外の意味を持つものであり、Geo/Geo 曖昧性とは表記が同じ地名が複数の地理的な場所に存在するものである.

本研究では、Geo/Non-geo 曖昧性および Geo/Geo 曖昧性を解決するために、落合ら [18] による手法を利用する. 落合らは Geo/Non-geo 曖昧性では地名と共起する関連語が対象テキスト中に出現するかを、Geo/Geo 曖昧性では地名と共起する地名が投稿テキスト中に出現するかを比較することで、投稿テキストから適切な POI を抽出可能であることを示している。図 3 にツイートに対する位置情報の付与例を示す。POI-DB に登録する POI 情報には、市区町村名 1,918 件*4、駅名 9,172 件*5、観光スポット名約 3 万件*6を利用した。

3.2 投稿者希少性を考慮した特徴量計算

日本において、同一人物が数カ月間に複数回犯罪を経験あるいは目撃することは少ないと考えられる。そのため、あるユーザが複数回犯罪に関する投稿をしている場合、それはユーザの経験あるいは目撃した投稿ではなく、ニュースなどの伝聞情報やノイズである可能性が高い。たとえば、「泥棒」という単語を高頻度で投稿している投稿者は、泥棒を経験あるいは目撃しているのではなく、ゲームや映画など非現実世界における事象について投稿している可能性が高い。また、「警察」という単語を高頻度で投稿している投稿者は、警察を目撃したのではなく、警察関連組織における、公式アカウントである可能性が高い。一方で、普段「泥棒」という単語を投稿していない投稿者が「泥棒」という単語を投稿していない投稿者が「泥棒」という単語を投稿していない投稿者が「泥棒」という単語を投稿した場合、投稿者が泥棒を経験あるいは目撃した可能性が高いといえる。

そこで、ある犯罪事象に関連した投稿をしていた投稿者が一定期間中に同様の事象について投稿しているか否かを 投稿者希少性スコアとして計算する。投稿者が過去同様の 犯罪事象に関して投稿していないほど、投稿内容は投稿者 の経験あるいは目撃した内容である可能性が高いといえる。

投稿者希少性は、 $tf \cdot idf$ 法を用いて計算する。犯罪語マッチより得られた投稿を対象投稿とする。対象投稿における投稿者の過去投稿を抽出する。対象投稿を d_t , 投稿者

^{*4} 統計局ホームページ, http://www.stat.go.jp

^{*5} 駅データ.jp, http://www.ekidata.jp

^{*6} ドコモ地図ナビ、https://www.nttdocomo.co.jp/service/map_navi/

の過去投稿群を D,対象投稿 d_t における犯罪語 T の出現 頻度を $tf(d_t,T)$,犯罪語 T を含む投稿数を df(D,T) とし たとき,投稿者希少性 $user_score(T,d_t)$ は式 (1) より計算 される.

$$user_score(T, d_t) = tf(d_t, T) \cdot \log \frac{|D|}{df(D, T)}$$
 (1)

3.3 位置希少性を考慮した特徴量計算

任意の場所においてある犯罪事象に関する投稿が過去多く検出される場合、その情報は社会的に知られている可能性が高く、また犯罪関連投稿と考えにくい。たとえば、「痴漢」という単語がある任意の地域において過去多く投稿されている場合、「痴漢」が恒常的に発生しているか、「痴漢」に関連する犯罪事象以外の事象が恒常的に発生している可能性が高いといえる。このように、ある地域における恒常的な事象は、地域における投稿中に含まれる語の分布によって示される。

位置希少性は、犯罪語マッチより得られた対象投稿における投稿位置および投稿日より、投稿位置と紐付いた、投稿日中の投稿と過去の投稿を比較することで計算する.投稿日の投稿群を犯罪日投稿群、過去の投稿を過去投稿群とする.犯罪日投稿群における単語の出現率が、過去投稿群における単語の出現率と異なるほど、犯罪日における位置希少性が高いといえる.

犯罪日投稿群および,過去投稿群の単語の出現率を確率分布とし,Jensen-Shannon divergence を用い確率分布の類似度を計算する.犯罪日投稿群における投稿位置の単語出現分布をP,過去投稿群における投稿位置の単語出現分布をQ,P(i) およびQ(i) をそれぞれ分布P およびQ に従って選ばれた値がi である確率としたとき,

犯罪日希少性 $poi_score(P,Q)$ は式 (2) より計算される.

3.4 投稿影響力を考慮した特徴量計算

犯罪を経験あるいは目撃するという希少な経験に関する 投稿は、その投稿の閲覧ユーザの反応率が、非犯罪関連投稿と比べ高いと考えられる。古川ら [9] は、犯罪関連投稿 のうち、経験投稿における情報の広まりやすさが、非経験 投稿と比較して高いことを指摘している。

Twitter における情報の広まりやすさを示す指標として、 Reply 数と Retweet 数がある. Reply は投稿に対する返信 を意味し、Retweet は投稿の引用を意味している. 単純 な Reply 数、Retweet 数の比較を行うと、つねに Reply や Retweet を受けやすい著名人などの影響が大きい. そのため,対象投稿以前における Reply 数, Retweet 数を考慮し,対象投稿の広まりやすさを評価する必要がある.

対象投稿を d_t , 対象投稿の投稿者における過去投稿集合を D, 投稿 d の Reply 数を Rep(d), 投稿 d の Retweet 数を Ret(d) としたとき,対象投稿の投稿影響 $reply_score$ は式 (3), $retweeted_score$ は式 (4) より計算される.ただし, $\alpha(0 < \alpha < 1)$ は Laplace Smoothing [19] のためのごく小さな定数とする.

$$reply_score(d_t) = \ln \frac{Rep(d_t)|D|}{\alpha + \sum_{d \in D} Rep(d)}$$
 (3)

$$retweeted_score(d_t) = \ln \frac{Ret(d_t)|D|}{\alpha + \sum_{d \in D} Ret(d)}$$
 (4)

Reply 数と Retweet 数は対象投稿の投稿後経過時間によって変化する。本研究では、リアルタイムな犯罪関連投稿を抽出したい。警視庁犯罪情報マップでは、「ひったくり」、「女性に対する声かけ等」、「公然わいせつ」などの犯罪行為および迷惑行為を地理空間上にマッピングしたデータを、2週間に一度の更新頻度で公開している*7。本研究では、警視庁の公開データと比較してよりリアルタイムな情報を抽出することを目指す。そのため、Reply 数と Retweet 数は対象投稿後7日間以内の値を用いる。

3.5 投稿ステータス抽出

Twitter の投稿には、テキストだけでなく、写真や URL リンクなどが存在する。本研究では、投稿者が実際に犯罪を経験あるいは目撃した情報を抽出したい。web サイトの引用リンクなどを含む投稿は、引用情報である可能性が高く、写真情報を含む投稿の場合、投稿者の経験に基づく内容であると考えられる。

投稿テキスト中に URL が含まれるか否かと同時に, URL が写真へのリンクか否かの情報を抽出する. 対象投稿が URL を含むか否かを表す url_flg および, 対象投稿が写真 を含むか否かを表す $photo_flg$ は Algorithm 1 より計算される.

Algorithm 1 url_flg, photo_flgの計算

if 対象投稿が写真 URL を含む then $url_flg = 1$ $photo_flg = 1$ else if 対象投稿が URL を含む then $url_flg = 1$ $photo_flg = 0$ else $url_flg = 0$ $photo_flg = 0$ end if

^{*&}lt;sup>7</sup> 警視庁:犯罪情報マップ, http://www2.wagmap.jp/jouhomap/

3.6 犯罪関連投稿の分類

文章分類問題において、テキストを Bag-of-Words で表現し分類器を構築する手法が広く使われている。Dilrukshiら [20] は Twitter の投稿テキストをラベリングし、Support Vector Machine(SVM)を利用することで、投稿テキストを分類する手法について検討している。SVM を利用することで多量の特徴量を評価した分類器の構築が可能となることを示している。

投稿テキストを Bag-of-Words で表現すると,次元数は数万に及ぶ.次元数が膨大になることで,計算コストが増大してしまうという問題がある. Kireyev ら [21] は Twitter の投稿テキストをトピックモデルの1つである LDA (Latent Dirichlet Allocation) [22] を用いて 300トピックで表現することで,津波や地震などの災害イベント関連投稿を抽出およびクラスタリングをする手法について検討している. LDA を用いることで, Twitter の投稿テキストより潜在トピック情報が可能であることを示している.

本手法では、SVM および RandomForest [23] を用い犯 罪関連投稿を抽出する分類器を構築する. SVM および RandomForest は教師あり学習を用いるパターン認識手法 モデルである. SVM はマージン最大化を特徴とするパター ン認識モデルであり、文章分類問題において広く利用され ている. そのため、提案手法のテキストデータおよび提案 特徴量を用いた文章分類問題においても有効であると考 えられる. RandomForest はアンサンブル学習を特徴とす るパターン認識モデルであり、高い分類精度および汎化性 能が期待できる [24]. また、RandomForest はブートスト ラップ法を用いて学習を行う特徴から, 学習された分類モ デルにおける特徴量の重要度を検証することが可能となる. 投稿テキスト群を LDA で学習することで、投稿テキスト 中に含まれる潜在トピック情報を抽出し, 各投稿テキスト を 300 次元のトピック所属確率で表現する. 分類器を構築 する際,特徴量に位置および投稿者希少性を考慮した特徴 量である poi_score, user_score, 投稿影響力を考慮した特 徴量である reply_score, retweeted_score, 投稿ステータ スより抽出した url_flg, photo_flg を追加することで, 希 少性および影響力を考慮した犯罪関連投稿抽出を行う.

4. 評価

本章では実データに対し提案手法を適用することで、 Twitter から犯罪関連投稿を抽出する. はじめに, 実験条件について述べる. 次に, 対象データに対し提案特徴量を計算し犯罪関連投稿の抽出に有効であるかを明らかにする. 最後に, 既存手法と提案手法とで犯罪関連投稿の分類精度を比較し, 提案手法が犯罪関連投稿の抽出に有効であることを示す.

表 2 ラベリング結果 Table 2 Labeling result.

flg	投稿テキスト (一部)
	池袋でナンパっぽいのされて www しかも全ておじちゃ
	h←
	私事ではありますが…今日初めて変質者 (不審者) を見ま
	した。2回も目の前に現れました。
True	池袋なうだけど電車で前の男の人に痴漢された
	なんかさー、さっきなんてさぁコンビニ袋持ってたら女
	の人に近づかれて無理矢理袋の中開けられたんよ (´・
	-·')
	文の里周辺でおかんが変態に出会ったそうです。
	いつも駅で見る博多女子 (多分) の子が可愛い そして今
	電車で隣に座ってる (俺は変態ではない)
	【不審者出没情報】01:20 函館市桔梗 3 丁目付近で不審
	者が出没した模様です。家の施錠等を確認して下さい。
False	明日、警察に出頭してきます。嘘です。免許証落とした
	から、高島平警察署まで行くだけです…
	こんばんは!寄居警察に確認の電話しました。飼主が現れ
	たかの確認は取れないそうです。
	国会議事堂前に警察車両の「青い壁」―安保法案反対の群
	衆の前に出現 (弁護士ドットコム)

4.1 実験条件

Twitter API*8を利用して Twitter の日本語投稿データを取得した. データ取得期間は 2015 年 1 月 1 日から 2015 年 9 月 31 日までの 9 カ月間であり、日本語投稿および、Twitter 公式のクライアントを利用して投稿されたツイートを対象とした. 犯罪語辞書を用いたキーワードマッチングにて、犯罪語を含む日本語ツイートのうち、20,000 件をランダムサンプリングし、犯罪関連投稿である投稿をTrue、それ以外を False としラベリングした. ラベリングは投稿テキストを 1 件ずつチェックする手法を用い、1 名 10,000 件の計 2 名で実施した. 投稿テキスト中には、テキストが短文であるためラベルの判別が困難である投稿が存在する. 本研究では、投稿テキスト情報のみを用いてラベリングを行う. そのため意味が不明瞭な投稿については除外した.

表 2 にラベリング結果のうち、True および False を 5 件抽出した結果を示す。ラベリングした結果、20,000 件中 True が 2,375 件、False が 17,625 件となった。ラベル付き データを評価対象データとし、投稿者希少性 $user_score$ 、位置希少性 poi_score 、投稿影響力 $reply_score$, url_flg および、対象投稿中に含まれる url_flg , $photo_flg$ について計算することで、提案手法の有用性について評価する。

4.2 投稿者希少性を考慮した特徴量の有効性評価

本節では、投稿者希少性を考慮した特徴量 user_score が

^{*8} https://www.nttdocomo.co.jp/info/news_release/2011/05/13_01.html

表 3 $user_score$ の高い投稿(降順) Table 3 Posts with a high $user_score$.

flg	投稿テキスト (一部)	user_score
True	佐倉駅の階段で隣登ってたオヤジが上にい	0.9782
True	た JK のスカート凝視してたんだけど	0.9762
True	酔っ払いに絡まれてる間に池袋。	0.9700
False	石切丸様の太ももにすりすりすりすりす	0.9679
	りすりすりしたい	
True	池袋のナンパとキャッチころしたい	0.9675
True	天神川駅に警官とパトカーめっちゃおっ	0.9074
True	た。	0.9074

犯罪関連投稿を分類可能な特徴であるかを,評価対象データを用いて分析する.はじめに評価対象データについて,user_score を計算する.次に正例,負例ごとに計算されたuser_score の頻度分布を生成し,分布間に有位な差があるか確認する.最後に,user_score の高い評価対象データについて,定性的な評価を行う.

評価対象データに対し、過去投稿取得期間を 30 日間に 設定し *user_score* を計算した. 評価対象データにおける 投稿ユーザのうち、過去 30 日間に投稿を持っている投稿 者数は 20.000 投稿中 14.437 投稿となった.

犯罪関連投稿の分類に投稿者希少性が有効であるためには、犯罪関連投稿における投稿者希少性の分布に有意な差異があることが求められる。投稿者希少性の分布が犯罪関連投稿と非犯罪関連投稿とで異なるかを、ノンパラメトリック検定の1つであるBrunner-Munzel検定[25]を用いて検定したところ、有意差があった(p-value < 2.2e-16)。このことから、user_scoreの分布は犯罪関連投稿と非犯罪関連投稿とで異なるため、犯罪関連投稿の抽出に有効な値であるといえる。過去投稿取得期間中に過去投稿が存在しない場合user_scoreを計算することができない。そのため、本研究では分類器の特徴量としてuser_scoreを用いる際に、過去投稿が存在しない場合にはuser_score を用いる際に、過去投稿が存在しない場合にはuser_score = 0 とした。

表 3 に user_score の高い上位 5 投稿を示す. 投稿者が 犯罪を経験, あるいは目撃するという希少な経験に関する 投稿が, user_score を用いることで抽出されていることが 分かる. 一方で非犯罪関連投稿についても, user_score が 高い投稿が存在する. user_score は犯罪後辞書に存在する 犯罪関連語のみについて計算しているため, user_score の みを用いて犯罪関連投稿を抽出するのは困難であるといえる.

4.3 位置希少性を考慮した特徴量の有効性評価

本節では、位置希少性を考慮した特徴量 poi_score が犯罪関連投稿を分類可能な特徴であるかを、評価対象データを用いて分析する。はじめに評価対象データについて、poi_score を計算する。次に正例、負例ごとに計算された

表 4 poi_score の高い投稿(降順) Table 4 Posts with a high poi_score.

flg	投稿テキスト (一部)	poi_score
	ホテルの一室で男が死んだ。部屋を訪れた	
False	のは若い女が一人。彼女が部屋を出たとこ	0.6768
	ろは誰も見ていない。	
False	箱根山:警戒レベル「2」来週にも温泉供給	0.6746
raise	保守再開へ	0.0740
	寝屋川警察署管内において、駅から帰宅途	
False	中の自転車の前カゴからカバンをひったく	0.6616
	られる事件が連続発生しました。	
False	[拡散希望] 兵庫県たつの市富永・たつの市	0.6579
raise	役所近辺で迷子のネコを探しています。	0.6573
	変な男の人が来ちゃって 2 人で乗るの怖	
True	かったからエレベーターのボタンうちが押	0.6450
	したのに乗らずに逃げて来た。	
	•	

poi_score の頻度分布を生成し、分布間に有位な差があるか確認する。最後に、poi_score の高い評価対象データについて、定性的な評価を行う。

評価対象データに対し、過去投稿取得期間を30日間に設定しpoi_scoreを計算した。犯罪関連投稿の分類に位置希少性が有効であるためには、犯罪関連投稿における位置希少性の分布と、非犯罪関連投稿における位置希少性の分布に有意な差異があることが求められる。位置希少性の分布が犯罪関連投稿と非犯罪関連投稿とで異なるかを、Brunner-Munzel 検定を用いて検定したところ、有意差があった(p-value = 1.33e-05)。このことから、poi_score の分布は犯罪関連投稿と非犯罪関連投稿とで異なるため、犯罪関連投稿の抽出に有効な値であるといえる。

表 4 に poi_score の高い上位 5 投稿を示す。位置における希少性が高いと考えられる投稿が抽出されているが、ニュースなどの参照投稿が抽出されていることが分かる。 poi_score はある位置における単語の確率分布を比較した値であるため、語が犯罪関連語であるか否かについては考慮していない。そのため poi_score のみを用いて、投稿者が経験あるいは目撃した犯罪関連投稿を抽出するのは困難であるといえる。

4.4 投稿影響力を考慮した特徴量の有効性評価

本節では、投稿影響力を考慮した特徴量 reply_score および retweeted_score が犯罪関連投稿を分類可能な特徴であるかを、評価対象データを用いて分析する。はじめに評価対象データについて、reply_score および retweeted_score を計算する。次に正例、負例ごとに計算された reply_score および retweeted_score の分布間に有位な差があるか確認する。

評価対象データに対し、 $reply_score$ および $retweeted_score$ を計算した。 パラメータ α はゼロ 頻度問題に対するスムージング処理に用いられる。 そ

表 5 reply_score および retweeted_score
Table 5 reply_score and retweeted_score.

	平均	分散	対象数
$reply_score(True)$	0.1323	0.5200	2,375
$reply_score({\it False})$	0.2157	0.5948	17,625
$retweeted_score(True)$	0.2302	0.9371	2,375
$retweeted_score({\it False})$	0.2803	0.9148	17,625

表 6 url_flg および photo_flg
Table 6 url_flg and photo_flg.

	True	False	All
$url_flg = 1$	189	5,028	5,217
$url_flg = 0$	2186	12,597	14,783
$photo_flg = 1$	146	2,757	2,903
$photo_flg = 0$	2229	14,868	17,097

のため、 α はごく小さな定数を設定する必要があり、 $reply_score$ 、 $retweeted_score$ ともに 0.01 とした.

対象投稿のうち、Reply および Retweeted が存在しない投稿が多数を占めるため、reply_score および retweeted_score は多くの対象投稿で 0 となった. 表 5 に reply_score および retweeted_score の計算結果を示す.

reply_score および retweeted_score の分布が犯罪関連投稿と非犯罪関連投稿とで異なるかを,Brunner-Munzel 検定を用いて検定したところ,reply_score(p-value = 2.2e-16),retweeted_score(p-value = 9.908e-11)ともに有意差があった。reply_score および retweeted_score は犯罪関連投稿の場合,非犯罪関連投稿と比較し平均値が高い値を示すことが分かる.このことから,reply_score および retweeted_score は犯罪関連投稿を抽出する際に有益な値であるといえる.

4.5 投稿ステータスより抽出した特徴量の有効性評価

本節では、投稿ステータスから抽出される特徴量 url_-flg および $photo_-flg$ が犯罪関連投稿を分類可能な特徴であるかを、評価対象データを用いて分析する。はじめに評価対象データについて、 url_-flg および $photo_-flg$ を計算する。次に正例、負例ごとに計算された url_-flg および $photo_-flg$ に有位な差があるか確認する。

評価対象データに対し、 url_flg および $photo_flg$ 計算した。 url_flg を含む url_flg が含まれるか否かで [pic.twitter.com] あるいは [/photo/1] が含まれるか否かで 判定した。 \mathbf{z} 6 に url_flg , $photo_flg$ の計算結果を示す。有意水準 1%で、 url_flg および $photo_flg$ について、犯罪 関連投稿と非犯罪関連投稿との間に有意な差がみられた。犯罪関連投稿のうち url_flg および url_flg について、犯罪 関連投稿のうち url_flg がみられた。犯罪関連投稿のうち url_flg が分かる。また、犯罪関連投稿のうち写真を含む投稿の割合は url_flg 6.14%であり、非犯罪関連投稿の url_flg 15.64%と比較し少ないこ

表 7 追加した特徴量一覧 Table 7 Proposed feature list.

種別	特徴量
希少性	投稿者希少性を考慮した特徴量 (user_score)
布罗住	位置希少性を考慮した特徴量 (poi_score)
	Reply 数の増加を考慮した特徴量(reply_score)
影響力	Reweeted 数の増加を考慮した特徴量
	$(reweeted_score)$
	投稿中に URL が含まれているかどうかを考慮した
ステータス	特徴量 (url_flg)
A) - 9 A	投稿中に写真が含まれているかどうかを考慮した
	特徴量 (photo_flg)

とが分かる.このことから,犯罪関連投稿の場合,非犯罪 関連投稿と比較しURL,写真URLともに出現率が低いと いえる.

写真を含む投稿は、投稿者の経験に基づく内容であると考えられるため、犯罪関連投稿のうち写真を含む投稿の割合は、非犯罪関連投稿と比較し多いと考えられる。一方で、photo_flg の有効性評価においては、犯罪関連投稿のうち写真を含む投稿の割合は、非犯罪関連投稿と比較し少なくなっている。これは、犯罪を経験あるいは目撃した投稿者は、写真撮影を行うことで犯罪関連事象に巻き込まれる可能性が存在し、撮影が困難であるためだと考えられる。また、犯罪関連事象は突発的に発生することが想定されるため、写真などの付加情報を追加せずに投稿を行っているためだと考えられる。

Boyd ら [26] はランダムサンプリングした 720,000 投稿を分析し、一般に 22%の投稿が URL を含むことを示している。このことから犯罪関連投稿の URL 含有率は一般投稿および非犯罪関連投稿と比較し大幅に低く、 url_-flg および $photo_-flg$ は犯罪関連投稿を抽出する際に有益な値であるといえる。

4.6 犯罪関連投稿の分類精度評価

既存手法における投稿テキストのトピック所属確率のみを利用した分類器を baseline, 投稿テキストのトピック所属確率に提案特徴量を追加した分類器を work とし, 既存手法 (baseline) と提案手法 (work) の分類精度を比較することで, 提案手法の有用性について評価する. 表 7 に提案手法である work に追加した特徴量一覧を示す.

評価対象データのうち、犯罪関連投稿は全体の11.88%と少数であり、ラベルが不均衡であることが分かる。そのため、分類精度はAUC(Area Under the ROC Curve)[27]により評価する。AUC は ROC(Receiver Operatorating Characteristic curve)下の面積を用いて分類精度の良さを評価する指標である。完全な分類が可能な場合には AUC=1、ランダムに分類した場合には AUC=0.5 となり、AUC が1に近いほど優れた分類器であるといえる。また、AUC を

表 8 SVM および RandomForest のハイパーパラメータ **Table 8** Hyperparameter of SVM and RandomForest.

model	param	baseline	work
SVM	C	98	18
S V 1VI	γ	0.001	0.001
	$max_features$	5	45
RandomForest	max_depth	14	11
	$n_estimators$	1,900	1,000

表 9 SVM および RandomForest の AUC Table 9 AUC of SVM and RandomForest.

 分割	baseline		work	
刀削	SVC	RF	SVC	RF
1	0.5894	0.6076	0.6423	0.7248
2	0.5960	0.6098	0.6454	0.7091
3	0.5964	0.6092	0.6413	0.7001
4	0.5835	0.6075	0.6402	0.7201
5	0.6065	0.6219	0.6361	0.7237
avg.	0.5943	0.6111	0.6411	0.7156

用いて精度検証することで,不均衡な分類問題における精度を比較することが可能となる.

検証方法には、交差検定を用いる. 交差検定より生成された各テストサンプルにおける AUC の平均を最終的な評価値とする. また、各サンプルを生成する際には、正例と負例の比率を一定に保つことにした.

分類器には SVM および Random Forest を用いる. SVM および RandomForest は scikit-learn [28] にて実装されて いるプログラムを利用し、SVM のカーネルには rbf カーネ ルを用いた. SVM および RandomForest には様々なハイ パーパラメータが存在し、ハイパーパラメータの値により 分類性能は大きく変化する. 本研究ではベイズ最適化を利 用した HyperOpt [29] を用いてハイパーパラメータの最適 化を行う. チューニングの際には、ハイパーパラメータの 探索範囲および探索回数をそろえ, AUC を最大化させるよ う baseline, work それぞれのハイパーパラメータを最適化 した. SVM については、コストパラメータであるCおよび rbf カーネルのパラメータ γ の2つを, RandomForest では 木の生成に利用する特徴量数の制限 max_features, 木の深 さの制限 max_depth および生成する木の数 n_estimators の3つを最適化の対象とした.表8に最適化したハイパー パラメータを示す.

表 9 に 5 分割交差検定をした際の baseline および work の AUC を示す. SVM, RandomForest いずれにおいても, baseline の AUC が高い. baseline と work の予測結果について, AUC の差の検定を DeLong 検定 [30] を用いて行ったところ, baseline と work の差が各分割で有意となった (p-value < 0.01). このことから, 提案手法の特徴量が犯罪関連投稿の抽出に有効であることが分かる.

図 4 に評価対象データのうち 30%をテストデータと

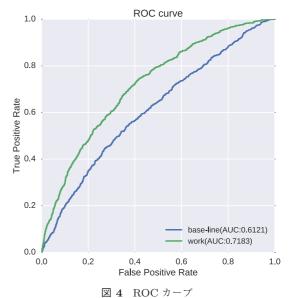


Fig. 4 ROC curve.

表 10 RandomForest (work) の重要特徴量 (降順)
Table 10 Important feature of RandomForest (work).

特徴量	Gini 係数
user_score	0.3235
url_flg	0.1570
poi_score	0.0972
reply_score	0.0631
$topic_292$	0.0227
$photo_flg$	0.0223
$retweeted_score$	0.0221

し、RandomForest で犯罪関連投稿を分類した際の ROC カーブを示す。ROC カーブ下における面積(AUC)が大きいほど、分類器がうまく犯罪関連投稿を抽出している。baseline における AUC が 0.6121 なのに対し、work では 0.7183 となった。RandomForest では不純度の計算過程で算出される Gini 係数の平均値を用いることで、特徴量の重要度を評価することができる。表 10 に work における RandomForest の重要特徴量を降順で示す。重要な特徴量は本研究で提案した特徴量となっていることが分かる。このことから、提案手法における特徴量が有効であるといえる。

表 11 にテストデータの予測結果のうち、予測スコアが高い投稿上位 10 件を示す。予測スコアは 1 に近いほど分類器が正例と判断し、0 に近いほど負例と判断している。予測スコア上位の投稿の多くは投稿者が犯罪を経験あるいは目撃した投稿であることが分かる。投稿者については、参照投稿や警察関連組織や行政の公式アカウントの投稿が排除されている。このことから、本研究の提案手法を用いることで、投稿者が経験あるいは目撃した犯罪関連投稿を抽出することが可能であるといえる。

表 11 予測スコアの高い投稿(降順)

Table 11 Posts with a high prediction score.

flg	投稿テキスト (一部)	予測スコア
False	新宿西口で酔っ払いなう。	0.7577
True	富士急ハイランドいってきた! 財布置き引	0.7535
	きにあった	
True	高駅警察おおい	0.7526
True	伊勢崎駅めっちゃ警察いるじゃん	0.7506
True	京都あんま客引き強ないから三ノ宮の客	0.7487
	引きにびびったおかまバー誘われまくっ	
	たで	
False	警視庁によると、7月6日、世田谷区の東	0.7467
	急電鉄・二子玉川駅のホームで	
True	盗撮騒ぎがある小岩は今日も平和ですね	0.7446
False	新大宮バイパスで変なパトカーを見かけ	0.7443
	た。大宮なのに春日部ナンバーだし	
True	さがみ野前で外人と警官が言い争ってる	0.7432
	笑	
True	酔っ払いがたくさんいる、上野駅	0.7431

5. まとめ

本研究では、投稿の希少性および影響力を考慮し、Twitterより犯罪関連投稿を抽出する手法について検討した。犯罪関連投稿と非犯罪関連投稿において、希少性および影響力を考慮した特徴量について比較した際、分布に差が出ることを示した。

評価対象データを Bag-of-Words で表現したのち, LDA を用い 300 次元に次元縮約を行い, SVM および Random-Forest を用いて犯罪関連投稿および非犯罪関連投稿を分類する分類器を構築した. 分類器の特徴量に Bag-of-Wordsから抽出された特徴量のみを用いた従来手法における分類精度と,提案手法を取り入れた分類精度を比較することで,提案手法の有用性について評価した.

今後の課題として、分類器の精度をより向上させることがあげられる。犯罪関連投稿と非犯罪関連投稿とを比較した際、希少性および影響力の分布がより異なれば、分類器の精度は向上するといえる。そのため、犯罪関連投稿とより関連のある希少性、影響力あるいは他の指標について検討する。また、キーワードマッチングに利用する犯罪関連語について、集合拡張を用いることで犯罪語を半自動的にソーシャルメディア上のデータから抽出する手法についても検討する必要がある。

参考文献

- [1] Chen, L. and Roy, A.: Event Detection from Flickr Data Through Wavelet-based Spatial Analysis, *Proc.* 18th ACM Conference on Information and Knowledge Management, pp.523–532 (2009).
- [2] 山田 渉, 菊地 悠, 落合桂一, 鳥居大祐, 稲村 浩, 太田 賢:マイクロブログを用いたイベント情報抽出技術, 情

- 報処理学会論文誌, Vol.57, pp.123-132 (2016).
- [3] Sankaranarayanan, J., Samet, H., Teitler, B.E., Lieberman, M.D. and Sperling, J.: TwitterStand: News in Tweets, Proc. 17th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems, pp.42–51, ACM (2009).
- [4] Yamaguchi, Y., Amagasa, T., Kitagawa, H. and Ikawa, Y.: Online User Location Inference Exploiting Spatiotemporal Correlations in Social Streams, Proc. 23rd ACM International Conference on Conference on Information and Knowledge Management, pp.1139–1148, ACM (2014).
- [5] Becker, H., Naaman, M. and Gravano, L.: Selecting Quality Twitter Content for Events., *ICWSM*, Vol.11 (2011).
- [6] Watanabe, K., Ochi, M., Okabe, M. and Onai, R.: Jasmine: A Real-time Local-event Detection System Based on Geolocation Information Propagated to Microblogs, Proc. 20th ACM International Conference on Information and Knowledge Management, pp.2541–2544 (2011).
- [7] Sakaki, T., Okazaki, M. and Matsuo, Y.: Earthquake shakes Twitter users: Real-time event detection by social sensors, Proc. 19th International Conference on World Wide Web, pp.851–860 (2010).
- [8] Li, J., Ritter, A., Cardie, C. and Hovy, E.H.: Major Life Event Extraction from Twitter based on Congratulations/Condolences Speech Acts, EMNLP, pp.1997– 2007 (2014).
- [9] 古川忠延,阿部修也,安藤剛寿,岩倉友哉,志賀聡子, 高橋哲朗,井形伸之:Twitterからの犯罪情報抽出の可能 性調査,情報処理学会研究報告,Vol.2011-DD-82,pp.1-6 (2011).
- [10] Eck, J., Chainey, S., Cameron, J. and Wilson, R.: Mapping Crime: Principle and Practice, United States National Institute of Justice (2005).
- [11] Du, B., Liu, C., Zhou, W., Hou, Z. and Xiong, H.: Catch Me If You Can: Detecting Pickpocket Suspects from Large-Scale Transit Records, Proc. 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.87–96, ACM (2016).
- [12] Rich, T. and Shively, M.: A methodology for evaluating geographic profiling software: Final report, Abt Associates Inc., Cambridge (2004).
- [13] Snook, B., Taylor, P.J. and Bennell, C.: Geographic profiling: The fast, frugal, and accurate way, Applied Cognitive Psychology, Vol.18, pp.105–121 (2004).
- [14] White, J.J. and Roth, R.E.: TwitterHitter: Geovisual analytics for harvesting insight from volunteered geographic information, *Proc. GIScience* (2010).
- [15] Wang, X., Gerber, M.S. and Brown, D.E.: Automatic crime prediction using events extracted from twitter posts, International Conference on Social Computing, Behavioral-Cultural Modeling, and Prediction, Springer, pp.231–238 (2012).
- [16] Fuchi, T. and Takagi, S.: Japanese morphological analyzer using word co-occurrence: JTAG, Proc. 36th Annual Meeting of the Association for Computational Linguistics and 17th International Conference on Computational Linguistics, pp.409–413 (1998).
- [17] Amitay, E., Har'El, N., Sivan, R. and Soffer, A.: Weba-where: Geotagging web content, Proc. 27th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval, pp.273–280 (2004).
- [18] 落合桂一,鳥居大祐:時間変化する特徴語によるマイク

ロブログ地名曖昧性解消,情報処理学会論文誌, Vol.7, pp.51-60 (2014).

- [19] Frakes, W.B. and Baeza-Yates, R.: Information retrieval: Data structures and algorithms, Prentice Hall PTR (1992).
- [20] Dilrukshi, I., De Zoysa, K. and Caldera, A.: Twitter news classification using SVM, Computer Science & Education (ICCSE), pp.287–291 (2013).
- [21] Kireyev, K., Palen, L. and Anderson, K.: Applications of topics models to analysis of disaster-related twitter data, NIPS Workshop on Applications for Topic Models: Text and Beyond, Vol.1 (2009).
- [22] Blei, D.M., Ng, A.Y. and Jordan, M.I.: Latent dirichlet allocation, *Journal of machine Learning Research*, Vol.3, pp.993–1022 (2003).
- [23] Breiman, L.: Random forests, Machine learning, Vol.45, pp.5–32 (2001).
- [24] Friedman, J., Hastie, T. and Tibshirani, R.: *The elements of statistical learning*, Vol.1, Springer series in statistics Springer, Berlin (2001).
- [25] Brunner, E. and Munzel, U.: The nonparametric Behrens-Fisher problem: Asymptotic theory and a small-sample approximation, *Biometrical Journal*, Vol.42, pp.17–25 (2000).
- [26] Boyd, D., Golder, S. and Lotan, G.: Tweet, tweet, retweet: Conversational aspects of retweeting on twitter, 43rd Hawaii International Conference on System Sciences, pp.1–10 (2010).
- [27] Fogarty, J., Baker, R.S. and Hudson, S.E.: Case Studies in the Use of ROC Curve Analysis for Sensor-based Estimates in Human Computer Interaction, *Proceedings of Graphics Interface*, pp.129–136, Canadian Human-Computer Communications Society (2005).
- [28] Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M. and Duchesnay, E.: Scikit-learn: Machine Learning in Python, Journal of Machine Learning Research, Vol.12, pp.2825–2830 (2011).
- [29] Bergstra, J., Yamins, D. and Cox, D.D.: Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms, Proc. 12th Python in Science Conference, pp.13–20 (2013).
- [30] DeLong, E.R., DeLong, D.M. and Clarke-Pearson, D.L.: Comparing the areas under two or more correlated receiver operating characteristic curves: A nonparametric approach, *Biometrics*, pp.837–845 (1988).

推薦文

本論文は、マイクロブログ上に投稿される犯罪関連投稿に着目し、警察などに認知されず、顕在化していない犯罪情報を抽出する手法を提案している。著者らの提案する位置と投稿者希少性を考慮した手法により、単純な辞書集合とのマッチングよりも希少性の高い犯罪関連情報の抽出が可能となることを、Twitterの日本語投稿データを用いて示している。精度向上などの課題は残っているが今後の研究でさらなる改善が見込める。また、これまで顕在化されていなかった犯罪情報を防犯対策などに活用できるようになるなど、社会的な貢献が期待できる研究である。以上の

理由により,本論文は推薦に値する.

(モバイルコンピューティングとパーベイシブシステム 研究会主査 稲村 浩)



高橋 柊 (正会員)

2013 年東京都市大学環境情報学部情報メディア学科卒業. 2015 年北陸先端科学技術大学院大学知識科学研究科博士前期課程修了. 同年株式会社NTTドコモ入社. 自然言語処理および機械学習の研究開発に従事.



菊地 悠 (正会員)

2000 年東京大学精密機械工学科卒業. 2002 年同大学院博士前期課程修了. 同 年株式会社 NTT ドコモ入社. SNS お よび位置情報解析の研究開発に従事.



落合 桂一 (正会員)

2006 年千葉大学工学部情報画像工学 科卒業. 2008 年同大学院博士前期課 程修了. 同年株式会社 NTT ドコモ入 社. 2014 年東京大学大学院工学系研 究科技術経営戦略学専攻博士後期課程 入学. SNS および位置情報データ解

析の研究開発に従事. 日本データベース学会, 人工知能学 会各会員.



深澤 佑介 (正会員)

2002 年東京大学工学部卒業. 2004 年 東京大学大学院工学系研究科修士課 程修了. 同年株式会社 NTT ドコモ入 社. 2011 年東京大学大学院工学系研 究科博士後期課程修了. 東京大学人工 物工学研究センターにて客員研究員兼

任. IEEE, 人工知能学会各会員. 博士(工学).