

犯罪予告の早期自動発見手法に関する研究

中村健二[†] 田中成典[‡] 寺口敏生[†] 大谷和史[‡] 山本雄平[‡]

関西大学大学院総合情報学研究科[†] 関西大学総合情報学部[‡]

1. はじめに

近年、Web の掲示板では、犯罪予告の投稿が増加している。それに伴い、予告の実行事例も増加[1]しており、早急な対策が求められている。犯罪予告の早期発見を支援する既存サービスに予告.in がある。予告.in では、Web でユーザが発見した犯罪予告を投稿するシステムと、巨大匿名掲示板群である 2ちゃんねるの投稿から、犯罪予告の可能性のある投稿を自動的に抽出するシステムを提供している。しかし、予告.in はユーザ参加型メディアであるため、投稿者の確保が難しいという問題点がある。また、犯罪予告の自動抽出では、監視対象が限定される問題点がある。そこで、本研究では、予告.in に掲載された犯罪予告から抽出した単語を基に犯罪予告語辞書の構築を行う。そして、犯罪予告語辞書を用いて、犯罪予告の判定を行う。また、掲示板がもつ特有の URL 構造と HTML 構造[2]のパターンを認識し、登録していない掲示板を抽出対象とする。そして、抽出した掲示板を自動監視[3]し、犯罪予告である可能性の高い新規投稿[4][5]を自動的に判別することで、犯罪予告の自動発見を支援する手法を提案する。

2. システムの概要

本研究では、Web の掲示板から犯罪予告である可能性が高い投稿の抽出手法を提案する。システムの概要を図 1 に示す。本システムは、犯罪予告文書と人手で作成した犯罪予告語辞書を入力とし、1) 犯罪予告学習機能、2) 掲示板検出機能、3) 犯罪予告文書抽出機能を通じて犯罪予告文書を出力する。

2.1 犯罪予告学習機能

本機能では、予告.in の報告文書に掲載された犯罪予告文書とその文書を基に、人手で作成した犯罪予告語辞書を入力し、犯罪予告文書の特

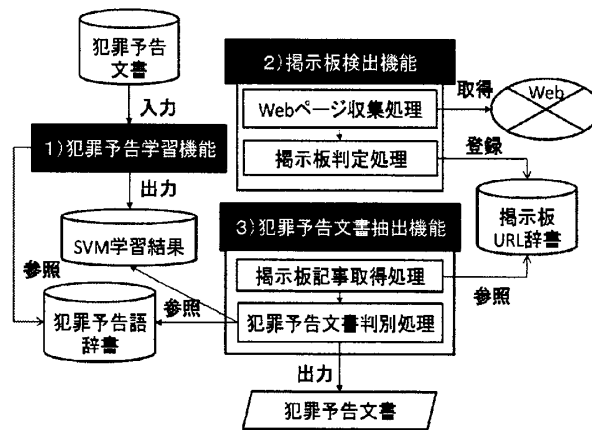


図 1 システムの概要

徴から 15 次元の特徴量を抽出したものを SVM (Support Vector Machine) の学習データとする。このとき、犯罪予告語辞書の作成においては、犯罪予告をより詳細に判別するため、犯罪予告語を殺人、放火、爆破、自殺、性犯罪と犯罪全般の 6 カテゴリに分類し、SVM の学習データとする。

2.2 掲示板検出機能

本機能では、Web ページ収集処理と掲示板判定処理の 2 つの処理を行う。Web ページ収集処理では、Web クローリングにより無作為に Web ページを収集する。掲示板判定処理では、Web ページ収集処理で取得した Web ページの URL を解析する。URL の解析は、URL をドメインとディレクトリに分割し、掲示板に特有の構造をもつ URL を特定することでを行い、検出した掲示板の URL を URL 辞書に登録する。

2.3 犯罪予告文書抽出機能

本機能では、掲示板記事取得処理と犯罪予告文書判別処理の 2 つの処理を行う。掲示板記事取得処理では、掲示板検出機能で作成した掲示板 URL 辞書から複数の掲示板を取得し、それらを解析して掲示板のページの全記事を抽出する。犯罪予告文書判別処理では、犯罪予告学習機能で作成した学習済みの SVM と犯罪予告語辞書を用いて掲示板の記事を判別し、犯罪予告文書を抽出する。

Research for Automatic Detection of Crime Warning on Web
[†] Kenji Nakamura, Toshio Teraguchi,
 Graduate School of Informatics, Kansai University, 2-1-1
 Ryouzenji-cho, Takatsuki-shi, Osaka 569-1095, Japan
[‡] Shigenori Tanaka, Kazufumi Otani, Yuhei Yamamoto
 Faculty of Informatics, Kansai University, 2-1-1 Ryouzenji-
 cho, Takatsuki-shi, Osaka 569-1095, Japan

3. システムの実証実験と考察

本システムの掲示板検出機能の実行結果を図2に示す。本提案手法の有用性を実証するために掲示板の検出精度と犯罪予告の検出精度について実証実験を行った。

3.1 実証実験

本実験では、掲示板の検出精度を検証するために掲示板の URL と掲示板でない URL をそれぞれ 30 件ずつ用意し、適合率、再現率と F 値 (F Measure) により掲示板の検出精度を判定する。また、犯罪予告の検出精度を検証するために犯罪予告を含む掲示板の記事 15 件と犯罪予告を含まない掲示板の記事 50 件を用意し、適合率、再現率と F 値により犯罪予告の検出精度を判定する。適合率とは、犯罪予告と判定された件数のうち何件が実際に犯罪予告かを示す基準である。再現率とは、実際の犯罪予告のうち何件が判定できたかを示す尺度である。また、F 値とは、適合率と再現率から算出する値であり、文における特定表現の抽出制度に用いられる。

3.2 結果と考察

掲示板検出精度は、表1に示すように、適合率 0.52, 再現率 0.83 と、F 値 0.64 という結果が得られた。この結果は、再現率が高い値を示しており、本研究の目的として犯罪予告が含まれる掲示板を漏れなく取得するという点から、有用な結果が得られていると考えられる。また、適合率が 0.52 と低い値にとどまった理由は、URL の構造が掲示板と類似している Blog を掲示板と誤判定したためである。次に、犯罪予告検出精度は、表2に示すように、適合率 0.78, 再現率 0.93, F 値 0.85 という結果が得られた。この結果は、犯罪予告を漏れなく取得するという点から、有用な結果が得られていると考えられる。また、適合率が 0.78 と低い値にとどまった理由は、「犯罪予告を起こした者を逮捕したニュース記事」や「ゲームの内容についての掲示板記事」において、犯行予告と同様の情報が含まれていたためであると考えられる。表1と表2の結果から、本提案手法は、Web に散在する犯罪予告を広く検出するという目的に対して、有用であることが実証された。

4. おわりに

本研究では、Web の掲示板から犯罪予告である可能性の高い文書の自動発見手法を提案した。実証実験の結果、犯罪予告である可能性の高い文書を自動的に抽出し、本システムの有用性を証明できた。今後の課題として、本提案手法では掲示板検出機能に URL の構造のみを用いたが、

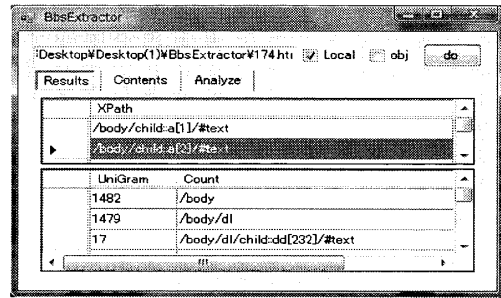


図2 掲示板検出機能の実行結果

表1 掲示板検出精度の検証結果

対象 URL 数	適合率	再現率	F 値
60	0.52	0.83	0.64

表2 犯罪予告検出精度の検証結果

対象記事数	適合率	再現率	F 値
65	0.78	0.93	0.85

対象 URL にアクセスして Web ページを取得し、Web ページの内容を確認することで適合率を向上できると考えている。また、犯罪予告文書抽出機能では学習データの不足により、誤判定が生じたが、犯罪予告として抽出されたものを随時学習させることにより更なる精度の向上が期待できる。今後の発展として、犯罪予告文書から実際に犯罪が行われる時間と場所の特定を行い、よりの確な犯罪防止に繋げる予定である。また、掲示板以外に Blog のコメント欄や Wikipedia など、不特定多数が文章を投稿できる Web ページにおいても犯罪予告が投稿されることがあるため、対象を掲示板から不特定多数が文章を投稿できる Web ページ全般に拡大する必要がある。

参考文献

- [1] 警察庁：平成 20 年上半期のインターネット・ホットラインセンターの運用状況等について，財団法人インターネット協会，2008.10.
- [2] 南野朋之，斎藤豪，奥村学：繰返し構造に基づいた Web ページの構造化，情報処理学会論文誌，情報処理学会，Vol.45, No.9, pp.2157-2167, 2004.9.
- [3] Kumar, R., Novak, J., Raghavan, P. and Tomkins, A.: World Wide Web, Springer Netherlands, 2003.5.
- [4] 毛利隆軌，北川博之：Hidden Web サイトからの新規トピック文書の抽出，情報処理学会論文誌，情報処理学会，Vol.46, No.SIG 5, pp.56-69, 2005.3.
- [5] Yang, Y., Zhang, J., Carbonell, J. and Jin, C.: Topic-conditioned Novelty Detection, International Conference on Knowledge Discovery and Data Mining, ACM, pp.688-693, 2002.7.