

スプログの収集における HTML 構造の類似性およびアフィリエイトの分析

森尻惇宜史[†] 片山 太一^{††} 石井 聡一^{†††} 宇津呂武仁^{††} 河田 容英^{††††}
福原 知宏^{†††††}

[†] 筑波大学理工学群工学システム学類 〒 305-8573 茨城県つくば市天王台 1-1-1

^{††} 筑波大学大学院システム情報工学研究科 知能機能システム専攻 〒 305-8573 茨城県つくば市天王台 1-1-1

^{†††} 東京電機大学大学院未来科学研究科 〒 101-8457 東京都千代田区神田錦町 2-2

^{††††} (株) ナビックス 〒 141-0031 東京都品川区西五反田 8-3-6

^{†††††} 独立行政法人 産業技術総合研究所 サービス工学研究センター 〒 135-0064 東京都江東区青梅 2-3-26

あらまし アフィリエイト収入を得ることを目的とするスプログの検出タスクにおいて、これまで、HTML 構造の類似性やアフィリエイト ID を用いることにより、一定の範囲のスプログが検出できることが知られている。これらの手法は単独で用いた場合の適用範囲が十分ではなく、両者の手がかりを併用する必要がある。これに対して、本論文では、HTML 構造の類似性、アフィリエイト ID のいずれによっても検出できないスプログに対して、SVM を適用することにより、高適合率の検出が可能であることを示す。

キーワード スパムブログ検出, HTML 構造, アフィリエイト, 機械学習

Analyzing Similarity of HTML Structures and Affiliate in Automatic Collection of Splogs

Akihito MORIJIRI[†], Taichi KATAYAMA^{††}, Soichi ISHII^{†††}, Takehito UTSURO^{††}, Yasuhide
KAWADA^{††††}, and Tomohiro FUKUHARA^{†††††}

[†] College of Engineering Systems, School of Science and Engineering, University of Tsukuba

^{††} Graduate School of Systems and Information Engineering, University of Tsukuba

^{†††} Graduate School of Science and Technology for Future Life, Tokyo Denki University

^{††††} Navix Co., Ltd.

^{†††††} Center for Service Research, National Institute of Advanced Industrial Science and Technology

Abstract Spam blogs or splogs are blogs hosting spam posts, created using machine generated or hijacked content for the sole purpose of hosting advertisements or raising the number of in-links of target sites. It has been shown that splogs can be detected based on similarity of HTML structures and affiliate IDs. The similarity of HTML structures of splogs is effective in splog detection, and the identity of affiliate IDs extracted from splogs can identify spammers much more directly than similarity of HTML structures, although it is not easy to achieve high coverage in extracting affiliate IDs. The coverage of the intersection of the two clues, similarity of HTML structures and affiliate IDs, is relatively low, and it is necessary to apply them in a complementary strategy. This paper studies how to detect splogs which cannot be detected based on either similarity of HTML structures nor affiliate IDs. We apply SVMs to this task and show that splogs of above type can be detected with high precision.

Key words spam blog detection, HTML structure, affiliate, machine learning

1. はじめに

ブログには個人の意見情報が記されており、市場の動向を推

測するための手掛かりや製品についての意見調査をする上で有益であるとして、近年注目を集めている。そのため、従来からあるインデクシングのみを行う検索エンジンとは異なる、プロ

グ特有の情報検索サービスが出現している。具体的には、ブログ解析サービスとして、Technorati, BlogPulse [2], kizasi.jp などが存在する。多言語ブログサービスとしては、Globe of Blogs が言語横断ブログ記事検索機能を提供している。また Best Blogs in Asia Directory がアジア言語ブログの検索機能を提供している。一方で、ブログのウェブコンテンツの作成と配信は非常に容易になっており、そのことが引き金となって、アフィリエイト収入を得ることを目的とするスパムブログ (以下、スプログ) が急増している [3], [9], [10], [12]。スプログにおいては、通常、広告主への誘導または対象サイトの被リンク数を増加する目的のもとで、機械的な文書作成や他サイトの引用という手段を用いて自動的に記事を生成し、大量のリンクを有するブログを機械的に自動生成する。文献 [10] は英語ブログにおいて、約 88% のブログサイトがスプログであり、それは全ブログポストの 75% を占めると報告している。このことから、文献 [11] に述べられているように、スプログは情報検索品質の低下やネットワークと格納資源の多大な浪費などといった問題を起す要因となる。そのため、近年、スプログの分析や検出を目的とした研究が進められている。いくつかの既存研究 [9], [10], [12] はスプログの重要な特性を報告している。文献 [12] では、TREC Blog06 データコレクションを用いて、スプログのピング時系列特性、入力度数/出力度数の分布特性、典型的な単語群を分析している。また、文献 [9], [10] は、BlogPulse データセットを用いたスプログ分析の結果を報告している。一方、文献 [4], [8], [11], [13] 等においては、言語情報、リンク情報、HTML タグ情報、時間情報といった多様な特性を手がかりとしてスプログを検出する技術を提案している。

上記の既存研究とは異なり、HTML 構造の類似性やアフィリエイト ID を手がかりとして、スプログの検出を行なう研究もある。スプログにおいては、一人の作成者が複数のスプログを機械的に生成していると考えられる。そこで、文献 [7] では、同一の作成者によって作成されたスプログの組において、HTML 構造が類似している場合があり、その特性を利用したスプログの検出を行なうことができることを報告している。これは、図 1 における領域「1」, 「3」を適用範囲とする手法になる。また、文献 [5] においては、アフィリエイト ID をスプログから自動抽出し、複数のブログサイトに含まれるアフィリエイト ID に着目して、スプログを収集、分析する手法を提案している。これは、図 1 における領域「2」, 「3」を適用範囲とする手法になる。これらの 2 つの手がかりについて、文献 [6] では、それぞれの手がかりの適用範囲の違いを示し、これらの 2 つの手がかりは相補的に用いる必要があることを報告している。以上をふまえて、本論文では、HTML 構造の類似性、アフィリエイト ID のいずれによっても検出できないスプログ (図 1 における領域「4」に対応する) に対して、機械学習の一つである Support Vector Machines (SVMs) [16] を用いることで、高適合率でスプログが検出できることを示す。

2. ブログの HTML 構造の類似性の測定

2.1 HTML ファイルからの DOM 系列の抽出

本論文では、文献 [18] で提案されたブロック抽出の方式をふまえて、HTML 文書から DOM 系列を抽出する [7]。まず、図 2 に示すように、HTML 文書 s 中の全ての HTML タグを木構造で表現する。次に、この HTML タグの木構造に対して、ブロックレベル要素として用いられるタグのうち、P タグおよび DIV タグによって木構造を分割し、これらのタグの下位にあるタグを取り込むことによって、個々のブロックを構成する。ここで、一般に、ブロックレベル要素としては、P タグおよび DIV タグ以外のタグも用いられるが、本論文では、簡単化のために、P タグおよび DIV タグに限定する。また、文献 [18] と同様に、BODY タグも、P タグおよび DIV タグと同様に扱い、BODY タグの位置において、HTML タグの木構造の分割を行う。さらに、文献 [18] では、ブラウザにレンダリングされない SCRIPT と STYLE の二タグ及びその下位ノードはブロック内に含まないとしているが、本論文では、ブロックの中身の詳細を区別するために、これらのタグ以下もブロック内に含める。次に、ブロックにまとめあげられた HTML タグの木構造を横型探索することにより、ブロックのリスト構造を形成し、HTML 文書 s の DOM 系列 $dm(s)$ とする。

2.2 DOM 系列の差分の割合

HTML 文書 s および t に対して、それぞれから抽出された DOM 系列 $dm(s)$, および $dm(t)$ の差分を DP マッチングによって求める。DP マッチングの際、挿入および削除のコストを 1、置換のコストを 2 とし、DP マッチングにより求まる編集距離 (レーベンシュタイン距離) を edit distance ($dm(s), dm(t)$) とする。次に、抽出された DOM 系列 $dm(s)$ の要素数を $|dm(s)|$ とし、以下の式で s, t の DOM 系列の差分の割合 $Rdiff(s, t)$ を計算する。

$$Rdiff(s, t) = \frac{\text{edit distance } (dm(s), dm(t))}{|dm(s)| + |dm(t)|}$$

また、HTML 文書 s に対して、HTML 文書集合 T の要素 $t \in T$ との間で、DOM 系列の差分の割合 $Rdiff(s, t)$ が最も小さいものを求め、その差分の割合の最小値を $\text{MinDF}(s, T)$ と定義する。

$$\text{MinDF}(s, T) = \min_{t \in T} Rdiff(s, t)$$

3. アフィリエイト ID を用いたスプログの分析

図 3 にアフィリエイト ID の抽出例を示す。アフィリエイトリンクには、そのアフィリエイトリンクを生成したアフィリエイトの アフィリエイト ID や、広告主の ID、商品 ID などが含まれており、我々はその中からアフィリエイト ID の抽出を行った。本論文では、特に、文献 [5] にしたがって、ASP (アフィリエイト・サービス・プロバイダ) のうち実際にアフィリエイト ID の抽出が可能な 10 社^(注1)を対象としてアフィリエイト ID の抽

(注1): Am 社, At 社, D 社, G1 社, I 社, Lk 社, R 社, St 社, Tr 社, V 社。

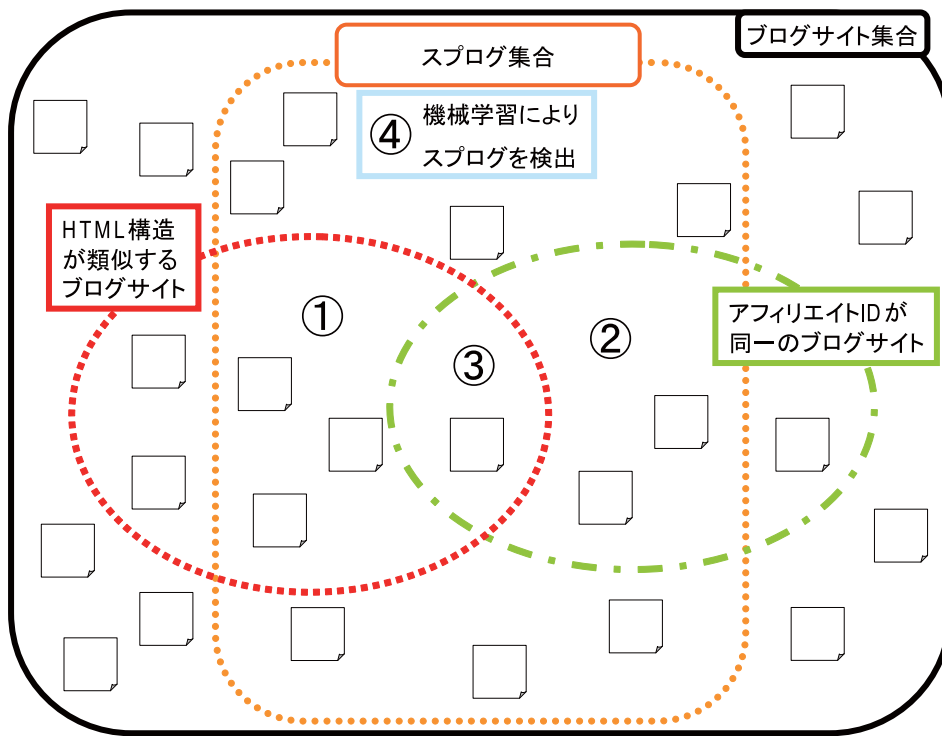


図1 スパログ検出における各手法の適用範囲

出を行った。

日本語ブログ収集にあたり、中国語、日本語、韓国語、英語のブログ記事の収集を行う KANSHIN システム [1] を利用する。このシステムでは、各言語のブログサイトのリストを参照し、ブログサイトの提供する RSS フィードファイルと Atom フィードファイルを取得し、記事をデータベースに蓄積している。このシステムに蓄積されたサイトから、分析対象とした S 社、F 社のブログホスト会社 2 社について、合わせて約 11 万ブログサイトを収集し、S 社の 48,183 サイトのうち 14,352 サイト (約 30%) から、F 社の 60,977 サイトのうち 6,231 サイト (約 10%) から、それぞれアフィリエイト ID が抽出された。

ここで、スパマーは、ASP から得られる報酬を少しでも増やすために、一つのアフィリエイト ID に対して複数のスパログを作成していると考えられる。そのため、複数のブログサイトに出現するアフィリエイト ID は、スパマーがスパログにおいて使用しているアフィリエイト ID である可能性が高くなると考えられる。あるアフィリエイト ID が、スパマーがスパログにおいて使用しているアフィリエイト ID であると推定できれば、そのアフィリエイト ID が出現する全てのブログサイトはスパログと判断することができる。以上の考えに基づき、文献 [5] においては、複数のブログサイトに出現するアフィリエイト ID を分析し、実際にスパマーがスパログにおいて使用しているアフィリエイト ID であると推定できる割合を報告している。本論文においても、文献 [5] の分析結果をふまえて、同一のアフィリエイト ID を含むブログサイト数が多いほど、スパログを自動生成している可能性が高いと考えて、ASP10 社のうちいずれかのアフィリエイト ID が抽出された 20,583 サイト (S 社、F 社の合計) を対象として、10 以上のブログサイト

表 1 10 以上のスパログに含まれるアフィリエイト ID 数およびそれらのアフィリエイト ID を含むスパログの総数

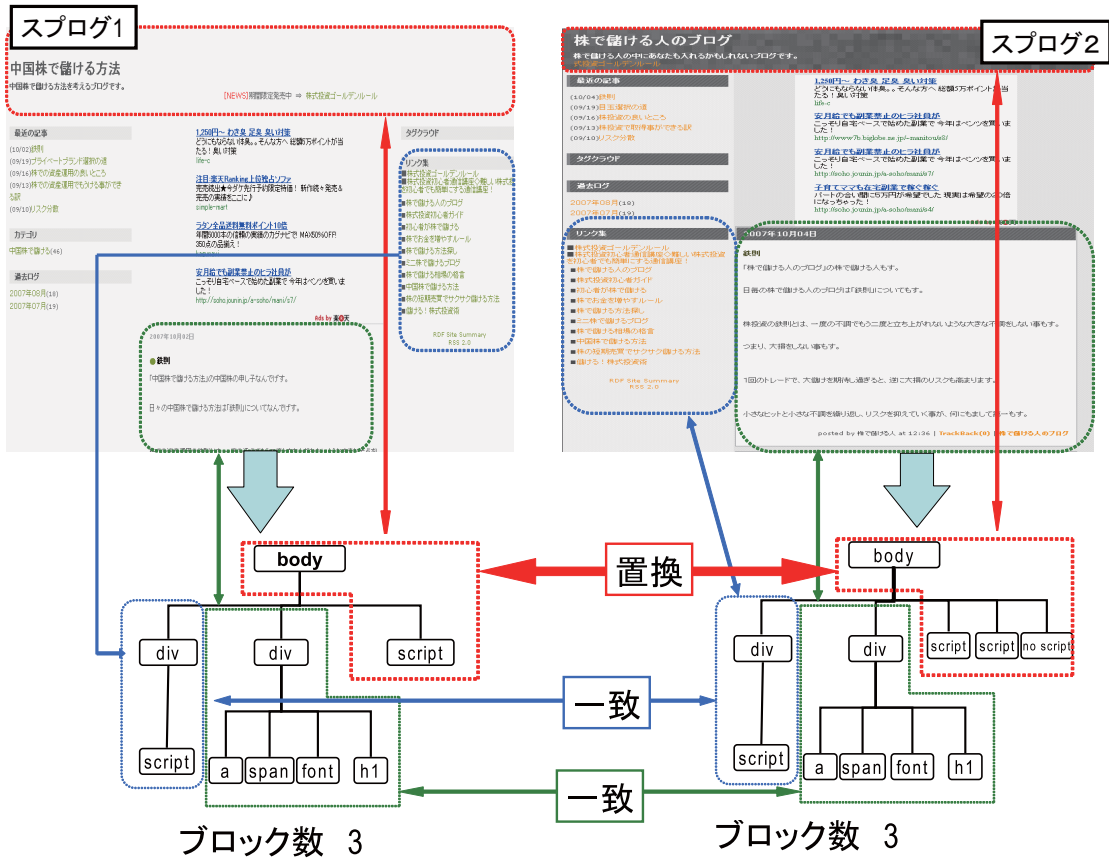
ブログホスト H	アフィリエイト ID 数	$\left \bigcup_x SP_{af}(H, x) \right $
S 社	72	1,101
F 社	56	953

図3 アフィリエイト ID を含むアフィリエイトリンクの例

から抽出されたアフィリエイト ID を分析対象とした。

その結果、129 個のアフィリエイト ID が抽出され、これらのアフィリエイト ID を含むブログサイト数は 2,472 となった。このうち、スパマーがスパログにおいて使用しているアフィリエイト ID であると推定できたアフィリエイト ID は 121 個であり、これらのアフィリエイト ID を含むスパログ数は 2,054 であった (アフィリエイト ID のスパム率は 93.8%, ブログサイトのスパログ率は 83.3%)^(注2)。ここで、この 121 個のアフィ

(注2)：複数のブログサイトに含まれるアフィリエイト ID のうち、スパマーがスパログにおいて使用しているアフィリエイト ID とみなすことができないものを大別すると、自動アフィリエイト作成ツールの作成者のアフィリエイト ID の場合と、ブログホスト会社のアフィリエイト ID の場合とに分けられる。両者とも、スパマーがスパログにおいて使用しているアフィリエイト ID と比較して、



edit distance=2

Rdiff=2/(3+3)=0.33

図2 HTML 文書からの DOM 系列抽出および DOM 系列差分算出の例

リエイト ID の各々を x として、ブログホスト H においてアフィリエイト ID x を含むスプログの集合を $SP_{af}(H, x)$ と定義する。また、各ブログホストに出現したアフィリエイト ID 数、および、いずれかのアフィリエイト ID を含むスプログの総数を表 1 に示す。

4. スプログ検出のための素性

本節では SVM によるスプログ判定において用いる素性について述べる。

4.1 ブラックリスト/ホワイトリスト URL 素性

訓練事例として、スプログ/非スプログが与えられると、その HTML ファイルからアウトリンクとなっている URL を抽出する。その中から以下の条件を満たすものを選定し、ホワイトリスト URL とした。

- i) 訓練事例中のスプログの HTML ファイルのいずれにも含まれない URL である。
- ii) 訓練事例中の非スプログの HTML ファイルの中で、2 回以上出現する URL である。

次に、各ホワイトリスト URL u に以下のように重みづけを行

い、ホワイトリスト URL 素性の値を算出した。

$$\log \sum_u \left(\begin{array}{c} \text{訓練事例全体の中の} \\ \text{非スプログにおける} \\ u \text{ の総出現頻度} \end{array} \right) \times \left(\begin{array}{c} \text{評価事例} \\ \text{における } u \text{ の} \\ \text{出現頻度} \end{array} \right)$$

一方、ブラックリスト URL についても、同様の手順で選定した。

4.2 名詞句素性

文献 [14], [17] の知見より、スプログおよび非スプログ中における単語の分布には異なりがあり、特定の種類の単語は非スプログよりもスプログに現れやすいということがわかっている。

そこで、特定の名詞句とスプログ、非スプログとの間の相関をとらえるために、名詞句素性を導入する。

具体的には、スプログ/非スプログの本文テキストを形態素解析(注3)した結果から名詞句を抽出し、以下の分割表にしたがって、訓練データ中のスプログ/非スプログにおける名詞句 w の出現頻度を用いて、スプログと名詞句 w との間の ϕ^2 統計量を求めた。

アフィリエイト ID が出現するブログサイト数が相対的に多いため、アフィリエイト ID のスパム率よりもブログサイトのスプログ率の方が低くなった。

(注3)：日本語形態素解析器 茶釜 (<http://chasen-legacy.sourceforge.jp/>) および IPA dic 辞書 (<http://sourceforge.jp/projects/ipadic/>) を用いた。

	w	$\neg w$
訓練データ中の スプログ	$\text{freq}(\text{スプログ}, w) = a$	$\text{freq}(\text{スプログ}, \neg w) = b$
訓練データ中の 非スプログ	$\text{freq}(\text{非スプログ}, w) = c$	$\text{freq}(\text{非スプログ}, \neg w) = d$

$$\phi^2(\text{スプログ}, w) = \frac{(ad - bc)^2}{(a+b)(a+c)(b+d)(c+d)}$$

また、評価事例に対しては、この名詞句素性の値として以下の式を用いた。

$$\log \sum_w \phi^2(\text{スプログ}, w) \times \left(\text{評価事例における } w \text{ の出現頻度} \right)$$

4.3 アンカーテキスト名詞句・リンク URL 素性

ブラックリスト/ホワイトリスト URL 素性および名詞句素性よりもより詳細な条件を設定することにより、より有効な性能を示す素性として、アンカーテキストの名詞句およびそのリンク先 URL の (緩い) 組み合わせを用いる。以下、まず、名詞句 w およびブログサイト s に対して、以下の尺度 $\text{Ancf}B(w, s)$ および $\text{Ancf}W(w, s)$ を定義する。

$$\text{Ancf}B(w, s) = \begin{pmatrix} \text{ブログサイト } s \text{ 中で名詞句 } w \text{ が} \\ \text{アンカーテキストに含まれ} \\ \text{そのリンク先がブラックリスト} \\ \text{URL もしくは訓練事例中の} \\ \text{スプログとなっている回数} \end{pmatrix}$$

$$\text{Ancf}W(w, s) = \begin{pmatrix} \text{ブログサイト } s \text{ 中で名詞句 } w \text{ が} \\ \text{アンカーテキストに含まれ} \\ \text{そのリンク先がホワイトリスト} \\ \text{URL もしくは訓練事例中の} \\ \text{非スプログとなっている回数} \end{pmatrix}$$

そして、訓練事例中のスプログ全体の中での総出現頻度 $\sum_s \text{Ancf}B(w, s)$ が 2 以上であるものを選定し、「ブラックリスト URL へのアウトリンクを持つスプログアンカーテキスト名詞句」とする。次に、 w を「ブラックリスト URL へのアウトリンクを持つスプログアンカーテキスト名詞句」として、評価事例 t に対して以下の重みを算出し、評価事例 t に対する「ブラックリスト URL へのアウトリンクを持つアンカーテキスト名詞句素性」の値とする。

$$\log \sum_w \left(\sum_{\text{訓練事例中のスプログ } s} \text{Ancf}B(w, s) \right) \times \text{Ancf}B(w, t)$$

同様に、訓練事例中のスプログ全体の中での総出現頻度 $\sum_s \text{Ancf}W(w, s)$ が 2 以上であるものを選定し、「ホワイトリスト URL へのアウトリンクを持つスプログアンカーテキスト名詞句」とする。次に、 w を「ホワイトリスト URL へのアウトリンクを持つスプログアンカーテキスト名詞句」として、評

価事例 t に対する以下の重みを、評価事例 t に対する「ホワイトリスト URL へのアウトリンクを持つアンカーテキスト名詞句素性」の値とする。

$$\log \sum_w \left(\sum_{\text{訓練事例中のスプログ } s} \text{Ancf}W(w, s) \right) \times \text{Ancf}W(w, t)$$

5. スプログ検出および信頼度尺度

SVM 機械学習を行うためのツールとして、TinySVM (<http://chasen.org/~taku/software/TinySVM/>) を用いた。カーネル関数としては、線形および 2 次多項式を比較し、2 次多項式の方が性能が良かったため、6. においては、2 次多項式カーネルを用いた場合の結果を示す。また、全ての素性に値がないものは訓練データから除外する。

また、SVM 機械学習での信頼度尺度として、分離平面から各評価事例への距離を用いた [15]。具体的には、スプログとして判定される事例に対する分離平面からの距離の下限 LBD_s および、非スプログとして判定される事例に対する分離平面からの距離の下限 LBD_{ab} をそれぞれ設定する。

6. 評価

6.1 訓練事例

3. で収集したブログサイトが、S 社は 48,183 サイト、F 社は 60,977 サイトあり、これらのうちで、10 以上のブログサイトから抽出されたアフィリエイト ID を含むスプログ数は、表 1 に示すように、S 社が 1,101 サイト、F 社が 953 サイトであった。

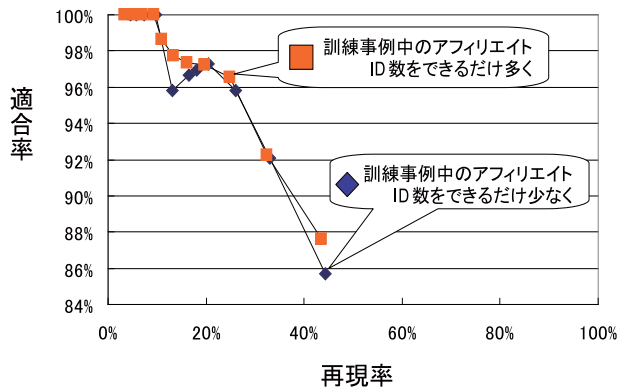
機械学習においては一定数以上の訓練事例が必要である。ここで、アフィリエイト ID を手がかりとして一定数の訓練事例を収集する場合に、以下の二通りの考え方がある。

- 訓練事例に含まれるアフィリエイト ID の数をできるだけ多くなるようにする
- 訓練事例に含まれるアフィリエイト ID の数をできるだけ少なくなるようにする

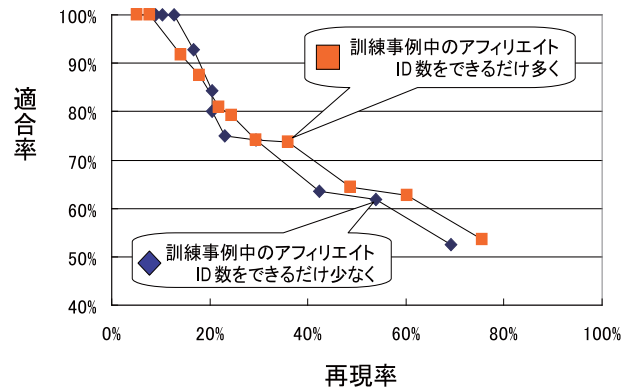
本論文では、この二通りの考え方に基づいた訓練事例の収集をそれぞれ行い、検出性能の比較も行う。

訓練事例に含まれるアフィリエイト ID の数ができるだけ多くなるようにした場合には、訓練事例収集のコストは最大になる。この場合、訓練事例の収集を行う際には、全てのアフィリエイト ID を選択し、S 社については 72 個のアフィリエイト ID を、F 社については 56 個のアフィリエイト ID を、それぞれ選択する。その後、全てのアフィリエイト ID について、それらを含むブログサイトを各アフィリエイト ID について同数程度ずつ選択し、500 サイト (各ホストあたり) を SVM におけるスプログの訓練事例とする。

一方、訓練事例に含まれるアフィリエイト ID の数ができるだけ少なくなるようにした場合には、訓練事例収集のコストは最小になる。この場合、訓練事例の収集を行う際には、より多くのブログサイトに含まれるアフィリエイト ID から順に選択し、S 社については 13 個のアフィリエイト ID を、F 社につい

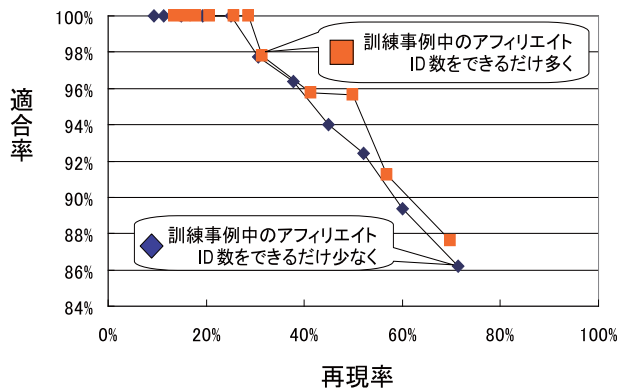


(a) スプログ検出性能

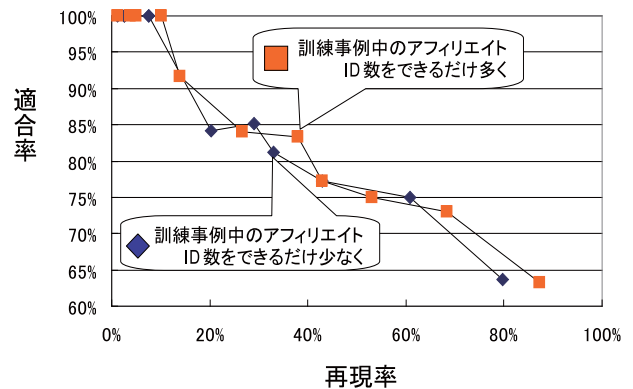


(b) 非スプログ検出性能

図 4 スプログ・非スプログ検出性能 (S 社)



(a) スプログ検出性能



(b) 非スプログ検出性能

図 5 スプログ・非スプログ検出性能 (F 社)

ては 12 個のアフィリエイト ID を、それぞれ選択する。その後、より多くのブログサイトに含まれるアフィリエイト ID を含むブログサイトから順に、これらのアフィリエイト ID を含むブログサイトを 500 サイト (各ホストあたり) を選択して、これを SVM におけるスプログの訓練事例とする。

また、あらかじめ人手で判定した非スプログを 500 サイト用意^(注4)、これを非スプログの訓練事例とした。

6.2 評価事例

3. で収集したブログサイトは、S 社においては 48,183 サイト、F 社においては 60,977 サイトであった。これらのうち、10 以上のブログサイトから抽出されたアフィリエイト ID を含まず、かつ、10 以上のブログサイトから抽出されたアフィリエイト ID を含むスプログと HTML 構造が類似しない ($\text{MinDF} > 0.15$) ブログサイトは、S 社は 47,029 サイト、F 社は 59,982 サイトとなった。このブログサイト集合を評価事例の候補集合とする事で、図 1 における領域「4」のスプログを対象として検出を行う。

この候補集合のうち、分離平面からの距離が分散されるようにスプログ側、非スプログ側の両方から評価事例を選択し、S

社については 283 サイト、F 社については 268 サイトに対して、人手でスプログ/非スプログの判定を付与した。

6.3 評価手順および評価尺度

スプログとして判定される事例に対する分離平面からの距離の下限 LBD_s について、分離平面からの距離が LBD_s 以上となる評価事例に対して、スプログと判定した場合の再現率、適合率を測定する。そして、 LBD_s を変化させた場合の再現率、適合率の推移を、S 社については図 4(a) に、F 社については図 5(a) にプロットした。

また、非スプログについても同様に、非スプログとして判定される事例に対する分離平面からの距離の下限 LBD_{ab} について、分離平面からの距離が LBD_{ab} 以上となる評価事例に対して、非スプログと判定した場合の再現率、適合率を測定する。そして、 LBD_{ab} を変化させた場合の再現率、適合率の推移を S 社については図 4(b) に、F 社については図 5(b) にプロットした。

また、図 4、および、図 5 においては、訓練事例中のアフィリエイト ID 数をできるだけ多くした場合については、「訓練事例中のアフィリエイト ID 数をできるだけ多く」としてプロットし、訓練事例中のアフィリエイト ID 数をできるだけ少なくした場合については、「訓練事例中のアフィリエイト ID 数をできるだけ少なく」としてプロットした。

(注4)：非スプログの収集の際に、本論文では無作為に収集を行なった。しかし、非スプログの作成者は一つしかブログを作成しないという考え方に基づいて、一つのブログサイトにのみ含まれるアフィリエイト ID を手がかりとして、非スプログの訓練事例をより簡単に収集する手法も考えられる。

6.4 評価結果

図 4, および, 図 5 に示すように, 訓練事例中のアフィリエイト ID 数をできるだけ多くした場合と, できるだけ少なくした場合を比較した時, 再現率, 適合率の推移はほぼ同等となった. これにより, アフィリエイト ID を用いて訓練事例を収集する際には, 収集のコストを抑えるために訓練事例に含まれるアフィリエイト ID を少なくしても, 検出の再現率や適合率に与える影響は少ないことが分かった.

また, 図 4(a), および, 図 5(a) に示すように, スプログの検出において, S 社においては再現率約 35%以下の範囲で, F 社においては再現率約 60%以下の範囲で, 90%を超える高い適合率となった. 一方, 非スプログの検出においては, 図 4(b), および, 図 5(b) に示すように, S 社においては再現率約 17%以下の範囲で, F 社においては再現率約 15%以下の範囲で, 90%を超える高い適合率となった.

6.2 で示すように, 今回の評価事例は図 1 における領域「4」を対象としている. これにより, HTML 構造の類似性, アフィリエイト ID のいずれによっても検出できないスプログに対して, 高い適合率での検出を実現することができたことが示せる. しかも, その適合率は, HTML 構造の類似性 [7] やアフィリエイト ID [5] の手がかりを単独で用いてスプログの検出を行なう場合とほぼ同等の高い適合率となった.

7. おわりに

アフィリエイト収入を得ることを目的とするスプログの検出タスクにおいて, これまで, HTML 構造の類似性やアフィリエイト ID を用いることにより, 一定の範囲のスプログが検出できることが知られていた. これらの手法は単独で用いた場合の適用範囲が十分ではなく, 両者の手がかりを併用する必要があった. これに対して, 本論文では, HTML 構造の類似性, アフィリエイト ID のいずれによっても検出できないスプログに対して, SVM を適用することにより, 高適合率の検出が可能であることを示した.

文 献

- [1] 福原知宏, 宇津呂武仁, 中川裕志, 武田英明. 複数の言語で記述されたブログ記事を対象とした言語横断型関心システム. 第 21 回人工知能学会全国大会論文集, 2007.
- [2] N. Glance, M. Hurst, and T. Tomokiyo. Blogpulse: Automated trend discovery for Weblogs. In *Proc. Workshop on the Weblogging Ecosystem: Aggregation, Analysis and Dynamics*, 2004.
- [3] Z. Gyöngyi and H. Garcia-Molina. Web spam taxonomy. In *Proc. 1st AIRWeb*, pp. 39–47, 2005.
- [4] 石田和成. スパムブログの推定と抽出. 日本データベース学会 Letters, Vol. 6, No. 4, pp. 37–40, 2008.
- [5] 石井聡一, 福原知宏, 増田英孝, 中川裕志. アフィリエイト ID を用いたスパムブログの分析. Web とデータベースに関するフォーラム (WebDB2010) 論文集. 情報処理学会, 2010.
- [6] 片山太一, 森尻惇宜史, 石井聡一, 宇津呂武仁, 河田容英, 福原知宏. HTML 構造の類似性およびアフィリエイトを用いたスプログの分析. Web とデータベースに関するフォーラム (WebDB2010) 論文集. 情報処理学会, November 2010.
- [7] 片山太一, 芳中隆幸, 宇津呂武仁, 河田容英, 福原知宏. HTML 構造を利用した類似スパムブログの収集. 第 2 回データ工学と情報マネジメントに関するフォーラム—DEIM フォーラム—論

文集, 2010.

- [8] P. Kolari, T. Finin, and A. Joshi. SVMs for the Blogosphere: Blog identification and Splog detection. In *Proc. 2006 AAAI Spring Symp. Computational Approaches to Analyzing Weblogs*, pp. 92–99, 2006.
- [9] P. Kolari, T. Finin, and A. Joshi. Spam in blogs and social media. In *Tutorial at ICWSM*, 2007.
- [10] P. Kolari, A. Joshi, and T. Finin. Characterizing the splogosphere. In *Proc. 3rd Workshop on the Weblogging Ecosystem: Aggregation, Analysis and Dynamics*, 2006.
- [11] Y.-R. Lin, H. Sundaram, Y. Chi, J. Tatemura, and B. L. Tseng. Splog detection using self-similarity analysis on blog temporal dynamics. In *Proc. 3rd AIRWeb*, pp. 1–8, 2007.
- [12] C. Macdonald and I. Ounis. The TREC Blogs06 collection: Creating and analysing a blog test collection. Technical Report TR-2006-224, University of Glasgow, Department of Computing Science, 2006.
- [13] G. Mishne, D. Carmel, and R. Lempel. Blocking blog spam with language model disagreement. In *Proc. 1st AIRWeb*, 2005.
- [14] Y. Sato, T. Utsuro, T. Fukuhara, Y. Kawada, Y. Murakami, H. Nakagawa, and N. Kando. Analyzing features of Japanese splogs and characteristics of keywords. In *Proc. 4th AIRWeb*, pp. 33–40, 2008.
- [15] S. Tong and D. Koller. Support vector machine active learning with applications to text classification. In *Proc. 17th ICML*, pp. 999–1006, 2000.
- [16] V. N. Vapnik. *Statistical Learning Theory*. Wiley-Interscience, 1998.
- [17] Y.M. Wang, M. Ma, Y. Niu, and H. Chen. Spam double-funnel: Connecting web spammers with advertisers. In *Proc. 16th WWW*, pp. 291–300, 2007.
- [18] 吉田光男, 山本幹雄. 教師情報を必要としないニュースページ群からのコンテンツ自動抽出. 日本データベース学会論文誌, Vol. 8, No. 1, pp. 29–34, 2009.