

ソフトウェア製品ニュースからの開発傾向の抽出

徳永秀和† 天雲 勇作† 青江順一‡

近年、インターネット上に多量の新製品ニュース記事が存在する。これらの記事の有効利用のひとつとして、ソフトウェア製品ニュース記事から最近の開発傾向を探るシステムの開発を行うこととした。そのためには、製品の機能を示すキーワードを抽出することが必要となる。この抽出を、サポートベクターマシンで行った。さらに、抽出もれを補うため、全文検索と重み付け統計を取ることによってソフトウェア製品の開発傾向が得られることを報告する。

Extraction of The Development Tendency from Software Product News Article

Hidekazu TOKUNAGA† Yusaku TENKUMO† Jun-ichi AOE‡

At present, a lot of new product news article exist on the Internet. So, it is effective to develop the system which explores the latest development tendency from a software product news article. It is necessary to extract the keyword which expressing the function of a product. We performed keyword extraction using the support vector machine. Furthermore, in order to compensate the omission in extraction, the full-text search and statistics of weighted keywords was performed. Consequently, we report that the development tendency of a software product was acquired.

1. はじめに

現在、コンピュータの急速な進化とネットワーク環境の整備により、一般人でも www にネットワーク接続しているパソコンさえあれば多くの情報を閲覧可能である。しかし、多くの情報が得られる反面その中から目的の情報を探す労力がより多

く必要になる。例えば、多くのソフトウェア製品ニュースから最近の開発傾向を知ろうとすると、一つ一つのニュース文を見ていく必要があり、ユーザ側に大変な労力を強いる。そこで、ニュース記事より製品の機能的特徴を抽出し、開発傾向をユーザに分かりやすく提示するシステムを開発することが有用である。

本論文ではセキュリティ分野に限って、SVM (サポートベクターマシン) を用いて、製品の機能的特徴を示すキーワードをラベル付けしたニュース記事を学習させることによって、未知のニュース記事からキーワードを抽出する。得られたキ

† 高松工業高等専門学校

Takamatsu National College of Technology

‡ 徳島大学知能情報工学科

Dept. of Information Science & Intelligent Systems, University of Tokushima

ワードを全文検索と重み付け統計を用いて、開発傾向を提示するシステムに関する研究結果を示す。

2. SVM (サポートベクターマシン)

サポートベクターマシン^[1]は、特にパターン認識の能力について、現在知られている中でもっとも優秀な学習モデルのひとつである。SVMというのは「線形しきい素子」というモデルに「マージン最大化」という基準で学習させ、更にカーネルトリックという工夫を加えて性能を上げたものである。

2.1 線形しきい素子

SVMの基礎となる線形しきい素子は、ニューロンを非常に単純化したモデルで、入力ベクトル \mathbf{x} に対して、2値(-1, 1)の出力 y を識別関数 (1) に従って出力する。

$$y = \text{sgn}\left[\sum_{i=1}^n \langle \mathbf{w}_i, \mathbf{x} \rangle - h\right] \cdot \dots \cdot (1)$$

$\langle \mathbf{w}, \mathbf{x} \rangle$ はベクトルの内積をあらわしており、

$\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_d) \in \mathbf{R}^d, h \in \mathbf{R}$ である。

また、 $\text{sgn}[u]$ は $u > 0$ で 1, $u \leq 0$ で -1 をとる符号関数である。このモデルの働きを幾何学的に説明すると、入力空間である \mathbf{R}^d を $\langle \mathbf{w}, \mathbf{x} \rangle - h = 0$ で定義される超平面で二つに分け、一方に 1 をもう一方に -1 を対応させることにあたる。線形しきい素子の学習は、与えられた例題に対して、 \mathbf{w} や h というパラメータを調整することにより行われる。

2.2 マージン最大化

すべてのサンプルに対して正しい出力を出す超平面は一意的ではない、そこでそのような超平面のうち汎化能力の点において最適なものはあるはずである。直感的に考えると訓練パターンすれすれを通る超平面よりは多少余裕をもって分ける超平面のほうが良いと思われる。SVMではその余裕をマージンという量で測り、マージンをできるだけ大きくするような超平面を求める。マージン

は超平面と一番近くにある訓練パターンとの距離の最大値である。

3. システム概要

3.1 システム構成

まず学習用に Web 上から取得したセキュリティ関連のニュース記事を形態素解析し、形態素単位で抽出 ($y = 1$) と非抽出 ($y = -1$) のラベル付けを人手で行い学習用ファイルを作る。そして次節で示す特徴ベクトルを作成し、SVMで学習させる。その学習済みのSVMを用いて未知のニュース記事からキーワードを抽出する。そして抽出もれを補うため、得られたキーワードを未知の記事全体に渡って全文検索を行い、重み付け統計によって開発傾向を示す。

3.2 特徴ベクトルの各要素

本研究での特徴ベクトルは、山田らの論文^[2]を参考にし、必要最小限の要素をいくつか試した結果一番精度が良い傾向にあった次のようなものとした。

1. 対象単語の品詞
2. 対象単語の隣にある左右の単語の品詞
3. 対象単語の文字種 (英字, カタカナ, 平仮名, 漢字, 数字, 記号)
4. 対象単語の隣にある左右の単語の文字種
5. 対象単語が記事全体の何文目か
6. 対象単語の文の手前部分に特定単語 (ソフト, ソフトウェア, ツール, システム, パッケージ) がいくつ存在するか
7. 対象単語の文の後ろ部分に特定単語 (上記) がいくつ存在するか

3.3 特徴ベクトルの例

表1は元記事の例として取り上げたニュース記事である。このニュースはDELETE MASTERという個人情報保護ツールについての記事である。

まずこの元記事を、文章の各単語の形態素を得るために、形態素解析ソフト chasen^[3]にかけて各単語の形態素を得る。得られた結果は表2の通りで、この出力結果は左から単語そのもの、単語の

表 1 元記事の例

メディアビジョン, 個人情報保護ツール (DELETE MASTER) を発売
2003 年 7 月 29 日
(株)メディアビジョンは 29 日、個人情報保護ツール『DELETE MASTER (デリートマスター)』、『DELETE MASTER Personal (デリートマスター パーソナル)』を
⋮

表 2 Chasen の処理結果

メディアビジョン. メディアビジョン. メディアビジョン. 17 0 0
⋮ 79 0 0
個人. コジ. 個人. 2 0 0
⋮

表 3 ラベル付け済みのデータ

×. メディアビジョン. メディアビジョン. メディアビジョン. 17 0 0
×. ⋮ 79 0 0
○. 個人. コジ. 個人. 2 0 0
⋮

読み, 単語の原型, 最後に品詞とその活用を数値で表したものである。

そして, この段階で学習用データには人間の手でラベル付けを行う。この例では「個人情報保護ツール」というキーワードに正解の○をつけ, それ以外は×をつける。そのようにして正解のラベル付けを施したものが表 3 である。

次に, 表 3 のデータを読み込み特徴ベクトルを作るプログラムを実行し, 特徴ベクトルを作成したものが表 5 である。これの各要素を説明すると, 左からラベル (1.0 が抽出対象, 0.0 がそれ以外である。), 残りは要素のインデックス (表 4) (そのデータがどの要素のものか示すもの) と値が交互に入っている。表 5 で一番上の行で例を挙げると左から最初はラベルで次が 0 なので表 4 のイ

ンデックス関連表より左側の単語の品詞だが単語自体存在しないので 0 とする, 次は 3 だがこれも左側に単語が存在しないことから文字種も 0, 次は 1 で対象単語の品詞を表すここでは chasen の出力より未知語などを表す 17 という数値が入っている。次は 4 で単語の文字種の値 2 (カタカナ) が次に入っている。このように表 4 を参照しながら表 5 を見ることで特徴ベクトルがどのようなものか理解できる。

表 4 インデックスの関連表

0	左隣の単語の品詞
1	対象単語の品詞
2	右隣の単語の品詞
3	左隣の単語の文字種
4	対象単語の文字種
5	右隣の単語の文字種
6	その単語が存在する文において、左側に存在する特定の単語の数
7	その単語が存在する文が何文目に存在するか
8	その単語が存在する文において、右側に存在する特定の単語の数

表 5 学習用特徴ベクトル

0.0	0	0	3	0	1	17	4	4	2	79	5	5	⋯
0.0	0	17	3	4	1	79	4	5	2	2	5	6	⋯
1.0	0	79	3	5	1	2	4	6	2	2	5	6	⋯
													⋮

3.4 学習と抽出

表 5 の特徴ベクトルのパターンを SVM に学習させる。そして学習済み SVM を用いて未知の記事に対して, 同じように各単語の特徴ベクトル (ただし, 表 5 の左端の正解のラベルはない) を作成し各単語をラベル付けする。最後に SVM の出力結果より正解のクラスラベルがついている単語を取り出すことで抽出対象を取り出す。

3.5 全文検索と重み付け統計

SVM でのキーワード抽出には抽出漏れが多いためその抽出漏れを防ぐために全文検索と重み付

け統計という手法を用いる。SVMでの抽出において同じ単語がキーワードである場合でも記事によって抽出ができたり、キーワードの一部分しか抽出できないことが頻繁にある。全文検索はそのような場合を防ぐために得られたキーワードの記事全体に渡って全文検索にかけ該当する単語を拾い出すことによりSVMで抽出できなかった単語を抽出しようという試みである。だが、全文検索を用いることによって問題も発生する。SVMで望むべきものではないキーワードを抽出した場合にそのキーワードをも全文検索にかけてしまうことになる。そこで、抽出される不要キーワードが1単語だけである(複合語でない)ことが多いのと、キーワードが1単語で構成されていることがほとんどないことを利用し、記事をSVMで抽出された単語のみ表示し、その中で複合語となっている場合のみその複合語をキーワードとして抽出することで解決する。

だが、2単語以上で構成される不要なキーワードの抽出や、本来その記事の主題ではない部分に存在したキーワードの抽出などの問題が残る。そこで、なるべく影響を押さえる対策として取り入れたのが重み付け統計という手法である。この手法は抽出したキーワードそれぞれに記事をヒントに重みをつけるのである。ただ、この手法を有効に使うには、必要なキーワードの重みを増やし不要なキーワードの重みを減らすようなルールをうまく考える必要がある。これにより統計の結果も大きく変わってくる。本研究では次のようなルールを用いる事によって重み付け統計を行う。

[以下の場合重みを増やす]

- ・ そのキーワードがある文に発売、発表などの比較的ソフトウェアの特徴を表す語と同時に存在しやすい単語がある場合。
- ・ キーワードが記事の始めのほうに存在する場合。

[以下の場合重みを減らす]

- ・ キーワードが一つの記事に複数含まれる場合。

・ キーワードが記事の終わりのほうに存在する場合。

4. 実験

4.1 実験方法

[実験1]

ASCII24.com よりセキュリティ関係のソフトウェアニュース記事を30記事(23種類のキーワード、キーワード出現回数97回)を集めて、その記事を形態素解析にかけ、単語単位でラベル付けする。ここで正解ラベル付けする単語は、その記事が何に関してのものであるかというもの(例えばウイルス対策ソフト、個人情報保護ツール、総合パッケージなどその記事のソフトウェアがどの種類のソフトウェアであるか判断可能であるもの)である。そのラベル付け済みの30記事分の特徴ベクトルを作成しSVMに学習させる。次に、学習済みのSVMより未知の50記事(キーワード種類39種類、キーワード出現回数179回)からキーワード抽出がどの程度可能であるかを調べる実験を行う。

[実験2]

実験1の結果を元に全文検索を行う。実験1でのSVMにより抽出した全てのキーワードを単語に分別する。例えば「ウイルスソフト」を「ウイルス」と「ソフト」に分別し、リストに保存する。そしてリストに存在する単語だけを元記事において表示する。図1に一文だけでこの全文検索を行うとした場合の手順を示す。これは1文だけの場合であるがこれを全記事に渡って行う。この実験ではこの全文検索を使ってキーワードがどの程度取得できるかを検証する。

[実験3]

実験2の全文検索のデータの統計とそれを重み付け統計したものによって、全体の傾向の比較をすることにより重み付け統計の有効性について検証する。ただし、統計は人手によるものであり、類似語(セキュリティソフトとセキュリティツールなど)を一つにまとめたり、実験2で完全に抽出されているキーワード内で部分的なキーワードを取り除いている。例えば「ウイルス対策ソフト」

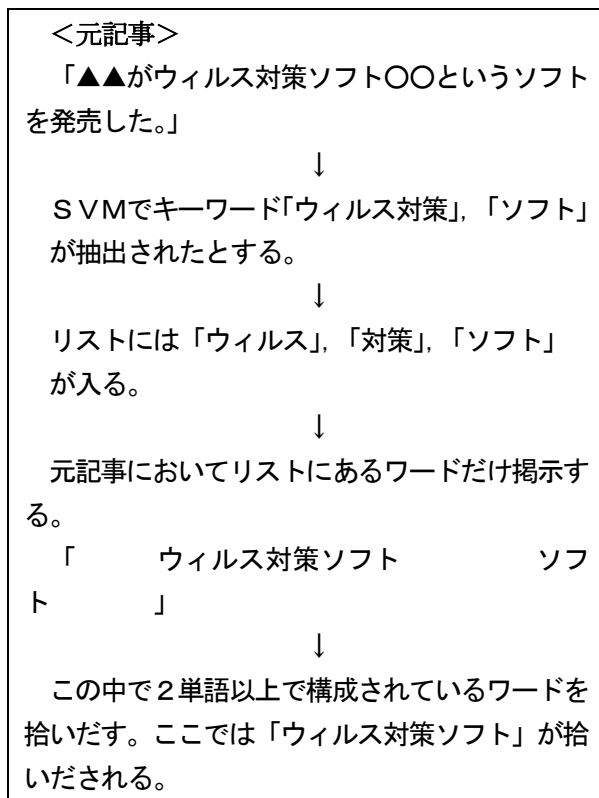


図1. 全文検索の例

が完全なキーワードとすると、「ウイルス対策」、「対策ソフト」は不要なキーワードとして排除している。ただし、記事全体においてその完全キーワード（例では「ウイルス対策ソフト」）の出現数が多い場合に限る。理由としては「管理ソフト」などキーワードでは、詳しく分けると多様なキーワード（個人情報管理ソフト、資産管理ソフト、etc・・・）がある場合別々ではランキングのかなり低い所に位置づけられてしまうが「管理ソフト」とするとランキングの上位にランク付けされ、傾向の把握にも有益な情報であるような場合があるためである。

4.2 実験結果

4.2.1 実験1

学習済みのSVMで未知の50記事（39種類のキーワード、キーワード出現回数179回）からキーワード抽出を行った結果を表6に示す。表6中で部分的、完全の意味は、例えば「ウイルス対策ソフト」の場合「ウイルス」「対策」「ソフト」の3つの単語で構成されておりこれを全部抽出で

きれば完全、どれかが抽出できていれば部分という意味である。

表6 実験1の結果

完全に抽出されたキーワード種類	6種類 (6/39, 15%)
部分的にしか抽出できなかったキーワード種類	25種類 (25/39, 79%)
完全にキーワードが抽出された個数	15個 (15/179, 8%)
部分的にキーワードが抽出された個数	67個 (67/179, 46%)
抽出された不要ワードの種類と個数	8種類 13個
抽出元の記事の特徴を表してなくキーワード種類には存在するもの	9種類 29個

4.2.2 実験2

実験2の結果を表7に示す。これを見ると完全抽出キーワード種類がかなり増えておりSVMで完全に抽出できていなかった部分を補っているのが見て取れる。

表7 実験2の結果

完全抽出キーワード種類	22種類
部分のみ抽出キーワード種類	9種類
抽出された不要ワードの種類	3種類
抽出された不要ワード出現回数	23回

4.2.3 実験3

実験3での結果、統計のランキングを取ったものを表8に示す。A, B, Cはそれぞれ自分での統計、全文検索のみを用いた統計、実験3の重み付けと全文検索を用いた統計の結果である。

4.3 考察

実験1の結果を見ると、キーワードを完全に抽出できているのは僅か10%程度である。ただ部分的に抽出されているキーワードは多く、しかも部分的にとられているのは毎回同じ部分ではなく違う部分である場合が多くあり、部分的なものを

表 8 抽出、検索結果の統計 (上位6ワード)

	A	B	C
少	管理ソフト	ウイルス対策	管理ソフト
	ウイルス対策	管理ソフト	ウイルス対策
	暗号化ソフト	電子メール	セキュリティ
	セキュリティ	セキュリティ	電子メール
	バックアップ	サーバー用ソフ	バックアップ
	メールソフト	バックアップ	サーバー用ソフ ト

A・・・自分で統計した結果

B・・・全文検索のみの統計結果

C・・・全文検索した上で重み付け統計結果

組み合わせると完全になるものが比較的多く見られる。他にもその記事の特徴を表していないが他の記事でのキーワードとなるものを抽出しているものも見受けられる。これは、後の加工の仕方によっては使えるデータが抽出されたとも言えるだろう。また、不要なキーワードが抽出されてしまっているが、ほとんどが1単語で構成されているものであり、全文検索の過程で排除されると考えられる。

次に実験2の結果について考察してみる。実験1の抽出データの傾向からも分かるように表.8を見ると明らかに完全キーワードが増えているのがわかる。実験1の考察で述べた事が原因と考えられる。ただ、不要キーワードの種類が少ない割に出現回数が多いのが問題である。このようなデータは後の統計時に大きく反映されることが予想できる。このように全文検索はSVMでの抽出を補っているが新たな問題点も生じる事がある。

最後に、実験3について考察する。まず表9を見て一番に目につくのが「暗号化ソフト」がない点である。これはSVMによる抽出の時点で抽出できておらず後の処理ではどうしようもないのが現状である。やはり、SVM抽出自体の精度を上げる必要があるだろう。次に目につくのが電子メールがランキングに入っている点である。これはあまり望ましくないキーワードであるのだが、キーキーワードの中に電子メールソフト等があったため全文検索でウイルス対策ソフトなどの記事

に多く含まれる電子メールというキーワードが多く検出されたのが原因である。しかし、重み付け統計されたCを見みると電子メールのランクが落ちており幾分実際の傾向に近いものとなっている。これは記事の後のほうの説明で多くでてきている電子メールというキーワードの重みが下がったためだと思われる。重み付け統計が有効に働いた結果だと言える。

5. 結論

現在の問題点としていくつかある。まず、特定の分野でしか実験を行っていないことから、様々な分野の場合に通じるか分からないというのが現状である。他に統計の際に類似キーワード(例えば「セキュリティソフト」と「セキュリティツール」など)の場合をどうプログラムで判断させるかということが全文検索と統計の自動化を行う上で課題である。

参考文献

- [1] Vapnik, V.N.: Statistical Learning Theory, A Wiley-Interscience Publication, 1998
- [2] 山田 寛康, 工藤 拓, 松本 祐治, "Support Vector Machine を用いた日本語固有表現抽出", 情報処理学会論文誌, Vol.43, No1, PP44-53, 2002
- [3] <http://chasen.aist-nara.ac.jp/>