

NHKニュース記事からのキーワードの抽出と記事の分野判別実験

浦谷 則好 畠田 のぶ子

NHK 放送技術研究所

{uratani,hatada}@strl.nhk.or.jp

1. はじめに

文書を検索する場合、分野やキーワードで指示することが通常行なわれている。キーワード抽出や文書分類に関して研究が種々なされてきた。しかし、文書を代表するキーワードをどのように選んだら良いのか、そのキーワードを用いて文書がどの分野に判別できるのかについて、効果的な基準や手法がいまだに存在しない。

NHKではニュース記事を、「政治」、「経済」、「社会」、「国際」、「スポーツ」、「各局」、「首都圏」の7つに分けて管理している。そこで、記事中からキーワードを選び出し、それらのキーワードが分野分類にどの程度効力があるのかを調べたので、それについて報告する。

2. キーワードの抽出

通常、キーワードは対象とする文章を形態素解析して、自立語の中から選ばれる。しかし、その場合、形態素解析の精度が問題となり、単語の長さを短単位にとるか長単位にするかは形態素解析プログラムに依存することになる。しかも、自立語、付属語の区別は必ずしも明確ではない。そこで、形態素解析を行なわず、字面処理でキーワードを選定することにした。

われわれは、以前字面から表現パターンを抽出するために3つの基準を比較している¹⁾。その結果、翻訳ユニット（機械翻訳のための表現パターン）の抽出のためにはエントロピー基準が一番適していることが判明した。（このエントロピー基準は下畠らのエントロピー基準²⁾とは違って文字列全体の分布のみから計算されるので計算コストは著しく小さい。）

しかし、エントロピー基準は長いパターンが優

●仕事量基準

$$n \times f_a$$

●平均情報量基準

$$\log (N / f_a) / I_n$$

●エントロピー基準

$$\log (N / f_a) / H_n$$

ここで、

α ：長さ n の文字列（例えば defgh）

f_a ：文字列 α の頻度

A_n ：長さ n の文字列全体の集合

N ： A_n に含まれる文字列の頻度の総和

f_n ： A_n に含まれる文字列の頻度の平均

I_n ： A_n 中で平均頻度を持つ文字列の情報量

$$= -\log (f_n / N)$$

H_n ： A_n のエントロピー ($= -\sum p_i \log p_i$)

先される傾向があるので頻度が小さくなり、できるだけ少数のキーワードで分野分類したいという観点から見て、キーワード抽出にはあまり適していないように思われる。仕事量基準は短いパターンが優先される傾向があり、分野に特有の複合的な語や言い回しを捉まえられない恐れが強い。そこで、今回は平均情報量基準を用いてキーワードの抽出を行なった。（手法の詳細は文献1を参照）対象とした文書は96年9月のNHKニュース記事1カ月分で、約90万文字（記事数4,038、文数21,078）である。その結果、頻度5以上で40,175個の文字列を抽出した。

3. 分類のためのキーワードの選定

2. で抽出した文字列から分野分類に有効と思われるキーワードを以下のようにして選定した。

どの分野にも同様に出現するものであれば、その文字列には分野を特定する能力がないと断定できる。逆にもしある文字列が特定の分野にだけ出

現するならば、その文字列を含んでいるだけで分野を特定できる。そこで、分野毎にどのような出現の偏りがあるかを調べるために、各々の文字列の分野毎の出現頻度を調べた。つまり、

$$p = M/N \quad M: \text{ある分野での出現頻度}$$

$$N: \text{全分野での出現頻度の総和}$$

としたとき、 p が大きいほど分野分類のキーワードとして適切だということができる。しかし、 p が大きくてもその出現頻度が小さければ実際上はあまり分野分類に寄与するとは思われない。また、 p が 0.5 のもの 2 つが出現するより p が 1.0 のもの 1 つ出現する方が明らかに分野を確信を持って限定できると考えられるので、分別力が p に比例すると考えるのは不適切だと思われる。そこで、

$(1.0 + p \log_2 p + (1.0 - p) \log_2 (1.0 - p))$ をキーワードの信頼度 (Q) とし、 $M \times Q$ が大きいものを分類のためのキーワードとして採用することにした。(ただし、 $p < 0.5$ のときは 0 とする)

$M \times Q$ の大きいものを表 1 に上げる。

文字列は特定の 1 分野にだけ偏っていると言えなくとも、特定の 2 分野を取れば明らかに偏った

表 1 $M \times Q$ の大きい文字列 (上位 15)

文字列	$M \times Q$ の値	分 野
選手	668.28	スポーツ
国際	583.49	国際
衆議院	455.04	政治
さきがけ	444.18	政治
気象庁	436.03	社会
橋本總理大臣	303.90	政治
伊豆諸島	277.59	社会
民主党	253.34	政治
リーグ	251.21	スポーツ
社民党	243.69	政治
暴風域に	242.90	社会
イラク	242.63	国際
関東地方	234.75	社会
衆議院選挙	228.59	政治
臨時国会	222.74	政治

出現をしている可能性がある。こうした文字列は単独では分別能力が低いと考えられるが、他のキーワードと一緒に現れるならば分別に寄与することが想像できる。そこで、2 分野を取ったときに $M \times Q$ の大きなものも調べた。この上位を表 2 に上げる。

表 2 2 分野で $M \times Q$ の大きい文字列

文字列	$M \times Q$ の値	分 野
地方	518.62	社会、各局
シ大統領	411.32	国際、政治
大田知事	365.77	政治、各局
千葉県	306.23	社会、首都圏
茨城県	277.42	首都圏、社会
沖縄本島	269.60	社会、各局
大阪	258.68	各局、社会
警察	258.09	各局、首都圏
エイズ	257.31	社会、各局
」と述べ	245.85	政治、国際

$M \times Q$ (2 分野から求めたものは 1/3 として) の大きさの順に文字列をソートした。まず、以下の方法でキーワード候補を絞った。

(a) 文字列を α としたとき、それより下位にその一部となるような文字列 β (例えば $\alpha = \beta \gamma$) があればそれを除く。

(b) 文字列を α としたとき、それより下位に α を含むもの δ (例えば $\delta = \alpha \gamma$) があればそれを除く。

(a) の判断は β より判別に有効なより長いキーワードが存在することを意味しているからである。

(b) は δ が出現するなら α が必ず出現するからである。

次に、 $M \times Q$ の大きい順にキーワードとして次式を満足する個数まで採用することとした。

$$c S_i \leq \sum (M Q_i)^2$$

S_i : 分野 i の記事数の二乗

$M Q_i$: 分野 i に対する $M \times Q$ の値

c は実験では 1.21 とした。ただし、分野毎に最低

30個は採用し、最大は250個に抑えた。結果として、1007個のキーワードが抽出された。分野毎のキーワードの個数は表3に示す。

表3 採用された分野別キーワード数

政治	経済	社会	国際	スポーツ	各局	首都圏
30	213	250	30	30	248	206

4. ニュース記事の分類

3. で選定したキーワードを用いて、キーワード選定の良否と分野判別能力を調べるために、ニュース記事の分類実験を実施した。キーワードを用いた文書分類は種々の方法が提案されている^{3), 4)}が、ここではtf法、tf*idf法⁵⁾、tf*rlb法（提案する手法）、Decision Tree法を採用した。ベクトル法である前3者は、各分野におけるキーワードの出現個数の期待値を用いて、分野の代表ベクターを求めておき、それと各記事のベクターとの内積を求めてその最大となる分野に判別するという方法をとった。（正規化は行なっていない。）

tf*rlb法は

$$tf \cdot df = log(N/df)$$

の替わりに3. で述べた信頼度Qでtfを重み付けする方法である。Decision Tree法ではp（分野を特定したときのキーワードの出現確率で3. のpと同じ）の大きい順に決定木を構成するということにした。

学習サンプル（96年9月のニュース記事）に対するtf*idf法の適合率（P）、再現率（R）、F値の値を表4に示す。ただし、

$$F = 2PR / (P+R)$$

である。記事中に1つもキーワードを含まないもの（未分類となるもの）が54記事（政治：3、経済：2、社会：19、国際：13、スポーツ：6、各局：7、首都圏：4）存在した。これらは全て誤分類と見なした。表3を見ると、一番多くキーワー

表4. tf*idf法の分野別の判定精度
(学習サンプル対象)

分野	記事数	適合率	再現率	F値
政治	652	.8137	.8773	.8443
経済	392	.7955	.8138	.8045
社会	836	.7639	.6423	.6979
国際	584	.7833	.7860	.7846
スポーツ	452	.8950	.9424	.9181
各局	537	.7635	.5773	.6575
首都圏	585	.6925	.8393	.7589
トータル	4038	.7816	.7712	.7764

ドが採用されているにも関らず「社会」に未分類が多いことが分かる。これは、「社会」で取り上げるトピックスが幅広いことを示している。

「スポーツ」、「政治」、「経済」は適合率、再現率とも他と比べて高い。これらの分野はキーワードで明確に特徴付けられることを示している。詳細に見ると「社会」、「各局」、「首都圏」の間で誤判定が多い（社会→首都圏：117、各局→首都圏：80）ことが判明した。

手法別の適合率、適合率を表5に示す。

表5. 手法別の精度（学習サンプル）

精度	tf	tf*idf	tf*rlb	DT
適合率	.7357	.7816	.7678	.7947
再現率	.7259	.7712	.7576	.7841

表5を見るとDecision Tree法が一番精度が良いことが分かるがtf*idf法とそれほどの差はない。F値で分野毎の誤分類の傾向を見てみると、誤判定の傾向はどの手法でもほぼ同じであったが、tf*idf法と比べるとtf*rlb法は「経済」で値が低く、Decision Tree法はtf*idf法より「社会」で優れ、「国際」で劣っていた。

テストサンプルを96年10月のNHKニュース記事1カ月分（記事数3,728、文数20,052、文字数829,280）として、同様に適合率、再現率を測定した。結果を表6に示す。当然のことながら全て

表6. 手法別の精度（テストサンプル）

精度	tf	tf*idf	tf*rlb	DT
適合率	.5940	.6449	.6277	.6364
再現率	.5829	.6328	.6159	.6245

の手法で精度が低下している。なお、テストサンプルにおいて1つもキーワードを含まない記事は70であった。（社会：29、スポーツ：11など）

表5、表6の結果を見るとわずか7つしか分野がないにも関わらず分類精度が低すぎるように思われる。この理由は前述した「政治」、「経済」などは内容の分類ではなく、実は「政治部」、「経済部」などの記事を担当した部署の分類に過ぎないからである。特に「社会」、「各局」、「首都圏」の間で誤判別が多いのは、それらが取り扱う話題が極めて似ていることに起因している。そこで、「各局」、「首都圏」を除いて5つの分野で適合率、再現率を測定し直してみた。（学習サンプル：2916、テストサンプル：2549）

結果を表7に示す。これを見ると、tf*idf法が一番優れていることが分かる。Decision Tree法は学習サンプルとテストサンプルで精度の差が大きく、ロバストでないことが分かる。

先に述べたように、「分野」は内容分類になっていない。そこで、テストサンプルにおいてtf*idf法とtf*rlb法で内積の大きかったもの700位まで取って、判別結果を吟味した。700位までの適合率は2つの手法で全く同じで0.9457であった。700位までに含まれる記事のうち538個が一致していた。このうち、正解と異なるものが29個（国際→経済：16、社会→政治：8など）であり、

表7. 手法別の精度（5分野）

	精度	tf	tf*idf	tf*rlb	DT
学習	適合率	.8378	.8656	.8430	.8423
	再現率	.8254	.8529	.8306	.8299
テスト	適合率	.7648	.8013	.7672	.7346
	再現率	.7474	.7831	.7497	.7179

両手法とも全て同じ結果を示していた。内容を見ると、そのうちの1つを除いて、分類としてはむしろ判別結果の方が妥当だと思われた。上位700個のうち、判別結果の方が間違っていると思われるものはtf*idf法で7、tf*rlb法で4しかなかった。つまり、700位までの真の適合率はtf*idf法で0.9900、tf*rlb法で0.9943であると思われる。この結果はキーワードの選定結果と手法の妥当性を示していると考えられる。

5. おわりに

最初に平均情報量基準を用いて、字面処理だけで記事からキーワードの候補を抽出し、文字列の包含関係と信頼度（評価尺度）を用いて、キーワードの選定を行なった。このキーワードを用いて記事の分類実験を行ない、選定されたキーワードと判別手法の妥当性を確認した。実験結果は、記事の内容検索という観点で見た場合、少なくとも評価値（内積）の大きいものについてはもとの分類より判別結果の方が信頼できることを示していた。記事データベースを作成する際には、tf*rlb法（and/or tf*idf法）が、分野の適否を半自動的に検証するのに役立ちそうである。

今後、さらに記事検索や分類に必要な要件を追求し、キーワードの抽出、キーワードの選定、判別手法の改良を図っていく予定である。

参考文献

- 1) 浦谷則好：ニュース原稿データベースからの表現パターンの抽出、情処50全大1R-8, (1995)
- 2) 下畠さよりほか：隣接文字の分散値を用いた定型表現の自動抽出、情処研資NL110-11, (1995)
- 3) 徳永健伸ほか：重み付きIDFを用いた文書の自動分類について、情処研資NL100-5, (1995)
- 4) 岩山真ほか：自動文書分類のための新しい確率モデル、情処研資H33-9, (1995)
- 5) G.Salton and C. S. Yang:On the specification of term values in automatic indexing, *Journal of Documentation*, Vol.29, No.4, pp.351-372(1973)