

NewsML エディタのための SubjectCode 階層に基づく記事分類システムの試作

An Article Classifier based on a SubjectCode Hierarchy for a NewsML Editor

児玉政幸[†] 大園忠親[†] 新谷 虎松[†]

Masayuki Kodama, Tadachika Ozono, Toramatsu Shintani

1 はじめに

近年、ニュース管理／配信フォーマットの標準化が進められており、国際新聞電気通信評議会 (IPTC) が策定した XML ベースの NewsML が普及しつつある [1]。NewsML フォーマットによる記事（以下、NewsML 記事）は様々なメタデータを包含可能であり、メタデータを利用することで一般的なキーワードベースの検索よりも詳細な検索を行うことが期待される。NewsML に含まれるメタデータの一つに、NewsML 記事をカテゴリ毎に分類するために使われる SubjectCode がある。SubjectCode は Subject（大分類）、SubjectMatter（中分類）、SubjectDetail（小分類）といった3つの階層構造によって表現される。IPTC で定義されている SubjectCode の総数は1,371種にも及ぶ。これだけの数があれば、記事を従来より詳細に分類することができる。しかし、これだけの数の中から編集した記事に適切な SubjectCode を付加することは記者にとって非常にコストがかかる作業である。また、記事の主題は多岐に渡ることが多く、複数の SubjectCode を適切に付加するとなると、さらに作業負担が強いられる。上記問題を解決するためには、記事の主題に沿った SubjectCode を自動分類してくれるサポート機構が記者にとって必要である。そうでなくとも、SubjectCode 候補を数件（ここでは3～5件を想定している）に絞り込むだけでも、記者を支援できることを考えている。

本研究の目的は、NewsML エディタ上で記者が編集した記事を SubjectCode に基づき分類し、記者の SubjectCode 付加作業を支援することである。具体的には、NewsML 記事集合の学習データを用いて、SubjectCode 特徴ベクトルを各々作成し、編集された記事から SubjectCode 候補を記者に提示することで支援を図る。本稿では、SubjectCode 階層構造や概念類義語を用いることで、SubjectCode 特徴ベクトルの各次元の重みを調整する手法について述べる。

2 NewsML エディタのための記事分類

既存の研究では、テキストを決められたクラスタに分類する研究は盛んに行われている。福本ら [3] は文書中に現れる語に対し、重み付けの学習を行った結果を用いて文書の自動分類を行う手法を提案している。Buckley ら [2] は情報間の類似度をキーワードの重みを要素とするベクトルで表現し、ベクトル間の距離を計算することにより、情報間の類似度を計算している。本研究でも同様に、重み付けを利用し特徴ベクトルを作成していくが、記事の構造や NewsML 特有の SubjectCode の階層構造を利用することで重み付けを行っている。ニュースの特徴を十分に用いることで、ニュースに特化した特徴ベクトルを作成できると考えている。

また、猪野ら [4] は WordNet から抽出した共通概念を特徴語として利用し、テキスト分類に有効であることを示している。本研究においても、共通概念による特徴語の拡張は必要であると考えており、EDR 概念辞書を利用して特徴ベクトルを作成している点において関連した研究である。

本研究では、分類システムの利用として NewsML エディタへの組み込みを想定している。記者にとって SubjectCode 付加は、記事を編集する上で余分かつ煩雑な作業である。したがって、NewsML エディタ設計の前提として記事編集過程

をできるだけ妨げない工夫が必要である。記事編集過程を妨げないために、SubjectCode 分類は記事編集中に行う。具体的には、記事編集者が1文を入力し終えた時点で、その都度、編集記事が分類システムにかけられる。記事編集者はシステムによって数種に分類された SubjectCode を選別するだけで良いので、記事編集中においても SubjectCode 作業負担が軽くなる。

3 SubjectCode の決定

3.1 NewsML 記事からの語抽出

SubjectCode を分類するために、各 SubjectCode 毎に SubjectCode 特徴ベクトルを作成する。SubjectCode 特徴ベクトルは、NewsML 記事から抽出した語を用いて作成する。まず、各 SubjectCode 毎にそれを含む NewsML 記事を全て収集し、リスト化する。日本語係り受け解析器 Cabocha¹を用いて収集した記事から語を抽出する。抽出する際のルールとして、人物名は性と名を一つに、組織名とそれに続く未知語は一つに、連続する名詞は一つにする。人物名や組織名は Cabocha によって解析された語をそのまま用いる。

また、SubjectCode の意味を SubjectCode 特徴ベクトルに追加することで、SubjectCode の意味を反映した特徴ベクトルを作成できる。しかし、SubjectCode の意味は SubjectCode 名では表現しきれない。そこで、SubjectCode 説明書から SubjectCode を表現する語を抽出する。SubjectCode 説明書は日本新聞協会が配布している SubjectCode に対する簡潔かつ的確な説明がなされている文書である²。SubjectCode を的確に表現している語が多くあるため、特徴ベクトルに付加する。また、抽出した語から EDR 概念辞書³を用いて類義語を抽出する。類義語を用いると言い換え表現や同義語にも対応できると期待できる。また、SubjectCode の日本語名称として使われている語は、SubjectCode を直接表現している語だと考えられる。したがって、日本語名称に含まれる語に重み付けを行う。ここでは、語の出現頻度を2倍とする。

3.2 SubjectCode 特徴ベクトルの語の重み付け

SubjectCode の階層構造を利用した語の重み付けの例として、SubjectCode「自動車事故」の特徴ベクトルの作成例を示す。NewsML 記事では例えば、「自動車事故」は SubjectCode の小分類に属しており、大分類の「災害・事故」、中分類の「事故（交通・運輸）」と共に記述される。これらの要素は概念上階層構造になっている。つまり、上位分類に含まれる概念は下位分類にも継承される概念だといえる。さらに、上位分類の語は下位分類の語よりも一般的な語である可能性が高い。裏を返せば、下位分類の語は大分類の語よりも記事の主題を特徴付けやすい語だと考えられる。このことから、より詳細な SubjectCode を分類することができるように重み付けを行う必要がある。上位分類の語をベクトル要素にそのまま加えると下位分類の語と出現頻度が同じになるため、下位分類の語の出現頻度を高くする。ここでは、下位分類における語の出現頻度を2倍とする。すなわち、小分類の語の重みを1としたとき、中分類の語の重みは $\frac{1}{2}$ 、大分類の語の重みは $\frac{1}{4}$ となる。

¹<http://chasen.org/~taku/software/cabocha/>

²<http://www.newsml.jp/>

³http://www.iijnet.or.jp/edr/J_index.html

[†]名古屋工業大学 大学院情報工学専攻

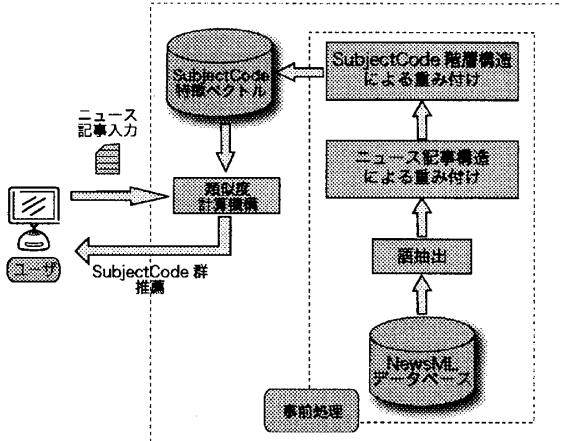


図1: システム概要

また、記事を構造的に見てみると、一般的に記事の大意は先頭のパラグラフに記述される傾向がある。また、タイトルは記事本文全体を要約した文章だと考えられる。つまり、記事本文の先頭になればなるほど、その記事を特定するための語が含まれると考えることができる。Brandowら[5]は、実験からこの考え方方が有効であるとしている。本研究では語の重み付けに上記手法を適用する。まず、記事をタイトルと各パラグラフに分ける。タイトルから順に語の出現頻度の重みを高くする。ここでは、下位パラグラフにおける語の出現頻度を $\frac{1}{2}$ とする。すなわち、タイトルの重みを 1 としたとき、第1パラグラフの重みは $\frac{1}{2}$ 、第2パラグラフの重みは $\frac{1}{4}$ となる。

3.3 SubjectCode 特徴ベクトルを用いた SubjectCode 分類

3.1で抽出した語に3.2で定義した重み付けを行い、特徴ベクトルを作成する。作成された記事ベクトルとSubjectCode特徴ベクトルの類似度を計算することで、与えられた記事のSubjectCodeを分類する。類似度計算はSinghal[6]の方法を利用した。Singhalの方法は、文書の長さによらず適切な類似度を求めることができる方法である。単語ごとに検索語のTF、及び検索対象のTFIDFの積を計算し足し合わせ、文書の長さにより正規化を行う。

図1に本手法を適用したシステムの概要を示す。

4 評価実験

本システムの評価方法として、*N-fold Cross Validation*(交差検定)を用いた。今回はN=3として評価実験を行った。データセットEは4,937件のNewsML記事を用いた。各NewsML記事には1つ、または複数のSubjectCodeが付加されている。EからはSubject(大分類)17種(17種中)、SubjectMatter(中分類)97種(375種中)、SubjectDetail(小分類)33種(979種中)、計147種(1,371種中)のSubjectCode特徴ベクトルを作成することができた。残り1,224種のSubjectCodeは4,937件中のNewsML記事には付加されていなかった。分類精度 $Prec$ は Subject, SubjectMatter, SubjectDetail 毎に式(1), (2)を用いて評価する。ここで、 $Prec_n$ は第n回目の分類精度、Lはテストセットの要素数、 $correct_i$ はNewsML記事l中の正しく分類されたSubjectCode数を示す。ここでの正しく分類されたとは、システムが提示した上位 α 番以内のSubjectCodeに一致するという意味である。Subject、及びSubjectDetailに対しては $\alpha = 5$ とし、SubjectMatterに対しては $\alpha = 10$ とした。記事に含まれるSubjectCode数は

表1: SubjectCode 分類結果

	Subject	SubjectMatter	SubjectDetail
test0	0.84	0.51	0.57
test1	0.81	0.50	0.58
test2	0.83	0.48	0.57
平均	0.83	0.50	0.58

各階層において最大5であったため、上記の設定にした。表1にSubjectCode分類の実験結果を示す。

$$Prec = \frac{1}{N} \sum_{n=1}^N Prec_n \quad (1)$$

$$Prec_n = \frac{1}{L} \sum_{l=0}^L \frac{correct_l}{correct_l + incorrect_l} \quad (2)$$

表1から、SubjectMatterとSubjectDetailがSubjectに比べ精度が低いという結果を得た。これはSubjectCode特徴ベクトルを作成する際の学習データの数が原因である。Subjectは多くの記事に含まれていたが、SubjectMatter、及びSubjectDetailを含む記事は比較的少なかった。今後の課題として、SubjectMatter、及びSubjectDetailの特徴ベクトルを作成する際、より多くの学習データを与える必要がある。

5 おわりに

本稿では、特にSubjectCode階層に着目し語の重み付けを行いSubjectCode特徴ベクトルを作成し、それを用いてニュース記事を分類するシステムを作成した。本システムをNewsMLエディタに組み込むことで、記事編集者は絞られた範囲内からSubjectCodeを選択するだけで良く、編集時ににおけるコストの削減が期待できる。さらにSubjectCodeが適切に付加されると、記事をより詳細なカテゴリで検索することが可能になる。今後、データ量を増やしてSubjectCode分類精度を上げると共に、本システムを実際にNewsMLエディタに組み込み評価する必要がある。

参考文献

- [1] 井上明、猪狩淳一、金田重郎、"ニュース配信のための国際データフォーマットNewsML:その概要と現状について", 情報処理学会論文誌, Vol.2002, No.056, 2002.
- [2] Buckley, C., Allen, J., and Salton, G., "Automatic Routing and Retrieval Using SMART", TREC-2, Infomation Processing Management, Vol.31, No.3, pp.315-326, 1995.
- [3] 福本文代、鈴木良弥、"語の重み付け学習を用いた文書の自動分類", 情報処理学会論文誌, Vol.40, No.4, 1999.
- [4] 猪野陽子、松井藤五郎、大和田勇人、"WordNetからの共通概念抽出によるテキスト分類", 日本ソフトウェア学会, 第22回大会論文集, 2005.
- [5] R. Brandow, K. Mitze, and L.F. Rau, "Information Processing and Management", Automatic condensation of electronic publications by sentence selection., Vol.31, No.5, pp.675-685, 1995.
- [6] Singhal A., Buckley C., and Mitra M., "Pivoted document length normalization", In Proceedings of SIGIR'96, pp.113-126, 1997