

1 V-9

# 分類された文書集合における 特徴的なキーワードパターンの抽出

小中 裕喜

三菱電機株式会社 先端技術総合研究所

## 1 はじめに

検索対象に対する知識や検索要求が曖昧で適切な検索式が構成できない場合、文書集合を分類/クラスタリングするアプローチが有望である。同一の検索要求に対して閲覧度の高い文書は互いに類似しているという仮説に基づき、類似した文書をまとめることにより、再現率は低いかもしれないが適合率の高い検索が期待されるクラスタが構成される。ここで各クラスタの特徴が何らかの形で提示されていれば、利用者にとって興味深い文書が集まつたクラスタを見い出すことが容易となる。そのような特徴づけを自動的に行う1つの方法は、各クラスタに属する文書のキーワードに類度などに応じて重みを与え、クラスタごとに最も重みを得たキーワードをいくつか提示することである。しかしながら、そのような方法ではキーワードの共起情報が考慮されず、また他のクラスタにも頻出するキーワードが選択されることもあるため、各クラスタの特徴を十分表しているとは言えない。

本稿では、分類/クラスタリングされた文書集合から、各クラスタを他のクラスタと識別するキーワードの組合せを抽出するアルゴリズムを提案する。抽出結果は各クラスタの概要を把握する手がかりとして、検索対象を絞り込むための候補検索式として、また後続の文書をインクリメンタルに分類するための知識として用いられる。

## 2 問題の記述

あらかじめ分類/クラスタリングされた文書集合において、各文書にはキーワード集合が付与されているものとする。 $\mathcal{W} = \{w_1, w_2, \dots, w_m\}$  を全キーワード集合、 $\mathcal{D} = \{d_1, d_2, \dots, d_n\}$  を全文書集合、 $d_i \subseteq \mathcal{W}$  を各文書に対応するキーワード集合とする。また $\mathcal{C} = \{c_1, c_2, \dots, c_l\}$  を全クラスタの集合とする。ただし $c_i \subseteq \mathcal{D}$ 、 $\bigcup_i c_i = \mathcal{D}$  である。パターン $p \subseteq \mathcal{W}$  をキーワードの集合とし、 $|p| = k$  のとき $k$ -パターンと呼ぶ。パターン $p$  を含む文書集合を $\mathcal{S}_p = \{d_i | p \subseteq d_i\}$  とする。

パターン $p$  のクラスタ $c_i$  における支持率 $sup$  と確信度 $cnf$  を $sup = |\mathcal{S}_p \cap c_i| / |c_i|$ 、 $cnf = |\mathcal{S}_p \cap c_i| / |\mathcal{S}_p|$  とする。

これらはパターン $p$  を AND 検索式、クラスタ $c_i$  を適合文書集合とみなした時の再現率、適合率に相当する。

各クラスタにおいて与えられた最小支持率 $minsup$  以上の支持率をもつパターンをラージパターンと呼ぶ。ラージパターンのうち、与えられた最小確信度 $mincnf$  及びその部分パターンの確信度の最大値以上の確信度を持つものを特徴パターンと呼ぶ。また最大確信度 $maxcnf$  が与えられた場合、それ以上の確信度をもつ特徴パターンは十分特徴的であるとみなされ、それを含むパターンは特徴パターンから除外される。

## 3 アルゴリズム

上記の問題を解くアルゴリズムは入力フェーズとマイニングフェーズからなる。

入力フェーズではクラスタごとに各文書のキーワード情報が入力される。このとき、どのクラスタにおいても $minsup$  %以上の文書に含まれていないキーワードは、決して特徴パターンに現れることはなく、結果に影響を与えることもないため、全文書情報から削除される。

マイニングフェーズは図1に示すように各クラスタ $c_i$  ごとに、パターンの構成キーワード数 $k$  を $k_{max}$  まで増やしながら、クラスタ内処理とクラスタ間処理を交互に行っていく。クラスタ内処理は基本的に Apriori アルゴリズム[2]を用いてラージパターン集合 $L_k$  を生成するが、ラージ 2-パターンの候補集合 $C_2$  の生成を効率化するために、ラージ 1-パターン集合 $L_1$  を収集した後に DHP アルゴリズム[3]のハッシング法を適用している。filter() と filter\_and\_count\_other\_support() は、ラージ 1-パターン集合 $L_1$  に現れないキーワードの情報を一時的に削除する。各文書の残りのキーワード情報及び各パターンのキーワード情報はビットベクタで表現し、その後の処理を効率化している。count\_support() はパターン $p$  のクラスタ $c_i$  からの支持 $|\mathcal{S}_p \cap c_i|$  を計算し、count\_other\_support() と filter\_and\_count\_other\_support() は各ラージパターンの他のクラスタからの支持 $|\mathcal{S}_p \cap \bar{c}_i|$  を計算する。ただし前回の反復でラージパターンを支持する文書がなかったクラスタは計算を省略する。apriori\_gen() と gen\_candidate() は[2] 及び[3] と同様であるが、最大確信度に関する変更がなされている。最大確信度が与えられた場合 $maxcnf$  以上の確信度をもつラージパターンが次の反復の候補パターンとして用いられることはない。

```

forall cluster  $c_i \in \mathcal{C}$  do begin
    /* find discriminative 1-patterns */
     $L_1 := \{\text{large 1-patterns}\};$ 
    set all the buckets of  $H_2$  to zero;
     $D_i := \emptyset;$ 
    forall document  $d \in c_i$  do begin
        filter( $d, L_1, \hat{d}$ ); /*  $\hat{d} \subseteq d$  */
         $D_i := D_i \cup \{\hat{d}\};$ 
        forall 2-subsets  $z$  of  $\hat{d}$  do begin
             $H_2[h_2(z)] ++;$ 
        end
    end
     $\overline{D}_i := \emptyset;$ 
    forall cluster  $c_j \in \mathcal{C}, c_j \neq c_i$  do begin
        forall document  $d \in c_j$  do begin
            filter_and_count_other_support( $d, L_1, \hat{d}$ );
             $\overline{D}_i := \overline{D}_i \cup \{\hat{d}\};$  /*  $\hat{d} \subseteq d$  */
        end
    end
    forall pattern  $p \in L_1$  do begin
         $p.cnf := |\mathcal{S}_p \cap c_i| / |\mathcal{S}_p|;$ 
    end
     $G_1 := \{p \in L_1 | p.cnf \geq mincnf\};$ 
     $L_1 := \{p \in L_1 | p.cnf < maxcnf\};$ 

    /* find discriminative 2-patterns */
    gen_candidate( $L_1, H_2, C_2$ );
    count_support( $D_i, C_2$ );
    forall pattern  $p \in C_2$  do begin
         $p.sup := |\mathcal{S}_p \cap c_i| / |c_i|;$ 
    end
     $L_2 := \{p \in C_2 | p.sup \geq minsup\};$ 
    count_other_support( $\overline{D}_i, L_2$ );
    forall pattern  $p \in L_2$  do begin
         $p.cnf := |\mathcal{S}_p \cap c_i| / |\mathcal{S}_p|;$ 
    end
     $G_2 := \{p \in L_2 | p.cnf \geq mincnf,$ 
              $p.cnf > p.maxAncestralCnf\};$ 
     $L_2 := \{p \in L_2 | p.cnf < maxcnf\};$ 

    /* find discriminative k-patterns ( $k \geq 3$ ) */
     $k := 3;$ 
    while ( $|L_{k-1}| \geq 2$  &  $k \leq k_{max}$ ){
         $C_k := \text{apriori\_gen}(L_{k-1});$ 
        count_support( $D_i, C_k$ );
        forall pattern  $p \in C_k$  do begin
             $p.sup := |\mathcal{S}_p \cap c_i| / |c_i|;$ 
        end
         $L_k := \{p \in C_k | p.sup \geq minsup\};$ 
        count_other_support( $\overline{D}_i, L_k$ );
        forall pattern  $p \in L_k$  do begin
             $p.cnf := |\mathcal{S}_p \cap c_i| / |\mathcal{S}_p|;$ 
        end
         $G_k := \{p \in L_k | p.cnf \geq mincnf,$ 
                  $p.cnf > p.maxAncestralCnf\};$ 
         $L_k := \{p \in L_k | p.cnf < maxcnf\};$ 
         $k ++;$ 
    end
    Answer $_i := \bigcup_k G_k;$ 
end

```

図 1: マイニングフェーズ

## 4 予備評価

CD-毎日新聞94年版[1]を対象とした予備評価を行った。各記事にはキーワード集合に加え、RWCテキストデータベース[4]においてUDCコードが付与されている。そこで、UDC主標数の上位2けたが「65」の2822記事(経営関係)を、下位2けたまで分類し、31のクラスタを得た。総キーワード数は71100である。これに対して上記のアルゴリズムをSun Ultra 1の上で適用したところ、 $minsup = 10\%$ ,  $mincnf = 60\%$ ,  $maxcnf = 80\%$ ,  $k_{max} = 3$ として、4629の特徴パターンが約3.1秒で得られた(I/O時間は除く)。具体的な特徴パターンとしては、例えばUDC658.15(私有財務管理、企業の財務管理)に対しては「決算&不良債券」(同19%、78%)、UDC658.16(財務整理)に対しては「救済&信用組合」(支持率17%、確信度100%)、「合併&信用組合」(同17%、85%)などが得られた。このように「不良債券」「信用組合」など単独のキーワードでは十分特徴づけできないクラスタに対しても、キーワードの組合せによって特徴づけられる場合があり、本アルゴリズムの有用性が示されている。一方、1000を超える特徴パターンが生成されるクラスタもあれば、幅広い記事が集まって特徴パターンが見い出せない場合も見られ、適応的なパラメータチューニングやシソーラス援用の必要性が示唆された。

## 5 おわりに

本稿では、あらかじめ分類/クラスタリングされた文書集合から、各クラスタを特徴づけるキーワードの組合せを抽出するアルゴリズムを提案した。現在、本アルゴリズムを応用した検索/分類システムを開発しているが、詳細は別稿にゆずる。今後の課題としては、適応的制御の導入、シソーラスや係受け関係などの考慮、階層的分類への適用などがある。

## 参考文献

- [1] <http://caactus.aist-nara.ac.jp/lab/resource/cdrom/Mainichi/MS.html>.
- [2] R. Agrawal and R. Srikant. Fast Algorithms for Mining Association Rules. In Proc. of VLDB'94, pp. 487-499, 1994.
- [3] J. Park, M.-S. Chen, and P. Yu. An Effective Hash Based Algorithm for Mining Association Rules. In Proc. of SIGMOD'95, pp. 175-186. ACM, 1995.
- [4] Real World Computing Partnership. RWC Text Database: RWCP-DB-TEXT-95-3, 1995. <http://www.rwcp.or.jp/wswg/rwcdb/text/index.html>.