

クラメールの連関係数を援用した類似文書検索システムの提案 A Framework for a Similar Documents Retrieval System Using Cramer's Coefficient of Association

樽松理樹†
Masaki Kurematsu

1. はじめに

現在、コンピュータを利用して数多くの文書に容易にアクセスできる環境が整ってきている。それとともに、それらの文書を効率良く処理する技術の開発が活発化[1]している。この分野の課題の一つとして、類似文書検索 [2]がある。類似文書検索の基本的な方法は、①元となる文書を何らかのモデルに変換する。②検索対象となる文書集合の各文書と同じモデルに変換する。③モデル間の類似度を計算し、類似度の高いものを抽出する。というものである。モデルとしては、文書中に出現する語の TF*IDF からなる文書ベクトルや、文書中に出現する語の出現確率に基づく確率モデルなどがある。これらの考えに基づく商用システムなども開発されているが、まだ精度には課題が残っている。そのため、さらなる手法についての研究が進められているのが現状である。

本研究では、このような類似文書検索に対し、カテゴリデータ間の関係を示すクラメールの連関係数[3]を援用するアプローチを検討した。本稿では、本手法を示すとともに、プロトタイプを用いた評価実験結果について報告する。

2. クラメールの連関係数を援用した類似文書検索システム

2.1 本研究で対象とする類似文書検索

本研究における類似文書検索は、特定の文書と文書集合中の全文書とを比較し、類似している文書を検索するというものである。端的に言えば、文書をクエリとした文書検索となる。これを実現するためには、任意の二つの文書の類似度を求める必要がある。この点に対し、クラメールの連関係数を援用する。

2.2 クラメールの連関係数

クラメールの連関係数は、カテゴリデータ間の関連の程度を表す指標の一つであり、二つのカテゴリの連関（独立性）を測る指標である。 k 個の要素からなるカテゴリデータ A と l 個の要素からなるカテゴリデータ B 間のクラメールの連関係数 $C_{A,B}$ は、式(1)によって求めることができる。

$$C_{A,B} = \sqrt{\frac{\chi^2}{n \times \min\{k-1, l-1\}}} \quad \chi^2 = n \left(\sum_{i=1}^k \sum_{j=1}^l \frac{f_{i,j}^2}{f_i f_j} - 1 \right) \quad \dots \text{式(1)}$$

ここで、 n はデータの総数、 $f_{i,j}$ は A の i 番目の要素 A_i と B の j 番目の要素 B_j が一緒に出現したデータ数、 f_i は A_i が出現したデータ数、 f_j は B_j が出現したデータ数を示す。

またクラメールの連関係数は $0 \leq C_{A,B} \leq 1$ の値をとり、1 の時に完全に連関となり、二つのカテゴリデータ間には強い相関があると言える。

2.3 クラメールの連関係数の援用方法

本研究では、クラメールの連関係数を文書の類似度と見立て、援用する。以下、その算出方法を説明する。

- ① 類似度を求めるために、文書を文書ベクトルに変換する。文書ベクトルの各要素は、形態素解析を用いて抽出した名詞および名詞列とその出現回数である。名詞および名詞列の順番は、出現回数（降順）、辞書順（昇順）でソートする。
- ② ①で作成した文書ベクトル A および B をカテゴリデータとみなし、式(1)によってクラメールの連関係数を求める。ここで式(1)中の各値は以下のように求める。
 k = 文書 A 中の名詞および名詞列の個数
 l = 文書 B 中の名詞および名詞列の個数、
 $f_{i,j}$ = 文書 A の i 番目の語句 $W_{A,i}$ の個数
 \times 文書 B の j 番目の語句 $W_{B,j}$ の個数
 \times 語の類似度 ($W_{A,i}, W_{B,j}$)
 語の類似度 (A,B) = $2 \times$ (語句 A と語句 B の共通意味数
 \div (語句 A の意味数 + 語句 B の意味数))
 なお意味数は、計算機可読型辞書(MRD)から求める。
 $f_i = f_{i,1} + f_{i,2} + \dots + f_{i,k}$
 $f_j = f_{1,j} + f_{2,j} + \dots + f_{l,j}$
 $n = f_{i,j}$ の総和
- ③ ②で求めた類似度の高い順に結果をユーザに提示する。

2.4 特許検索システムの構築

以上の提案内容に基づき、処理する文書を特許公報に限定した検索システムを構築した。そのスクリーンショットを図 1 に示す。

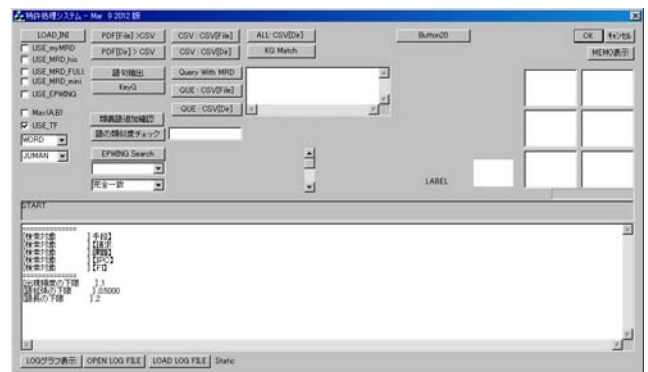


図 1 : 特許検索システム

†岩手県立大学ソフトウェア情報学部

