

階層型知識体系を用いたWWW情報の自動カテゴリ推定方法<sup>1</sup>

5 Z - 6

村本 達也 鷲崎 誠司<sup>2</sup>NTT 情報通信研究所<sup>3</sup>

E-mail : {muramoto, suzaki}@isl.ntt.co.jp

## 1 はじめに

近年のインターネットの爆発的普及に伴い、WWW(World Wide Web)上には無数の情報が存在する。この大量の情報の中から利用者に有用な情報の提供を可能にする情報フィルタリング手法が注目されている。情報フィルタリング手法は利用者の興味、趣味、嗜好等を基に情報を提示する方法が多いが、従来方法では利用者に陽に入力を要求するか、閲覧した情報内の特徴語句を抽出することにより推定していた。これに対して我々は、閲覧した情報を既存の階層型知識体系に対応させることにより、利用者興味推定の精度の向上を目指している。本稿では、これを実現するためのWWW情報の自動カテゴリ推定方法とその評価について述べる。

## 2 階層型知識体系

WWW上における階層型知識体系とは、Yahoo Japan[1]やNTT DIRECTORY[2]に代表される登録型のディレクトリ型検索サービスのように、構造が多段階の階層をなしており、各階層が上位概念・下位概念の関係を保持し、階層の各カテゴリにはそれに対応する情報(WWW情報)が保持されているものとする。以下、階層型知識体系の節点をカテゴリ、あるカテゴリから根のカテゴリまでのカテゴリの列をパス、カテゴリに保持されているWWW情報を情報と呼ぶ。また、情報の分類項目は階層型知識体系のパスで表現する(例 スポーツ:サッカー:W杯)。

## 3 自動カテゴリ推定方法

## 3.1 機能構成

本方法は、(1)WWW情報のテキストから形態素解析を用いた単語の抽出、(2)抽出した単語の階層型知識体系への対応付けによるカテゴリの推定、により実現されている。機能構成の外観を図1に示す。

処理(1)では、当該情報からだけでなく、当該情報が参照している情報(リンク先の情報)からも同様に単語を抽出する。これは、当該情報がフレーム組の情報であったり、イメージ情報のみであったりして、テキス

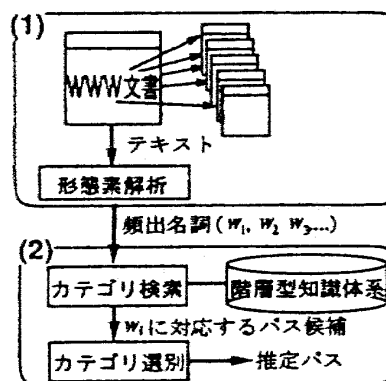


図1: 機能構成の外観

ト情報を全く得られない、または得られても特徴的な語句を抽出できないという事が考えられるからである。より多くの参照情報から単語を抽出することで、より有用な単語の取得が期待できるが、限界効用と計算効率を考慮して、5個の参照情報から抽出することにする。抽出した単語の中から特徴語句になりやすい普通名詞、固有名詞、形容動詞を採用し、出現度数の多いものを使用する。

処理(2)におけるカテゴリ推定では、図1に示すカテゴリ検索装置を用いる。これは、従来のディレクトリ型検索サービスが持つ機能の、検索語句を投入すると、その語句に対応した情報とその情報が格納されている階層型知識体系のパスの組をリストアップする機能である。これを利用して、先に抽出した単語を階層型知識体系に対応させることで、分類先を推定する。具体的な方法は、始めに複数の単語をカテゴリ検索装置に投入して、投入語句に対応するパスと頻度を得る。次に検索で得たパスの頻度の和を各パスについて計算する。そしてパスを頻度の和でソートして上位を採用することにより、推定する。即ち、ここで推定されたパスは、当該情報とその参照情報中の頻出単語を含む情報を多く保持していることになる。このような方法で処理することにより、もし処理(1)で抽出された語句に、特徴的でないもの(例:更新)が抽出されたとしても、この語句に対応するパスは多種になりやすく、頻度の和でソートすると下位になる可能性が高くなる。また、的外れな語句が抽出されたとしても、この語句に対応するパスと他の語句に対応するパスは重複しにくいので、頻度の和でソートすると下位になる可能性が高くなる。

<sup>1</sup>An Automatic Classification of the WWW Information Using a Hierarchical Knowledge-Base

<sup>2</sup>Tatsuya MURAMOTO, Seiji SUSAKI

<sup>3</sup>NTT Information and Communication Systems Laboratories

### 3.2 実行例

幾つかの WWW 情報が含まれるパスを推定した結果を図 2 に示す。これは階層型知識体系に NTT DIRECTORY を用いた場合であるが、各パスは必ずしも互いに素ではないので、図中「JavaFAQ, ICカード」の様に複数のパスに分類される情報も存在する。

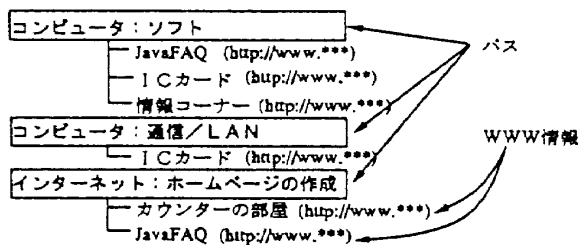


図 2 : 実行例

### 4 評価

実験には 1997 年 9 月に NTT DIRECTORY に登録された情報の中から無作為に抽出した 51 件を用いた。階層型知識体系として、NTT DIRECTORY と Yahoo Japan を用いた。評価尺度として、

$$\text{適合率} = \frac{\text{正解パス数}}{\text{採用したパス数}}$$

$$\text{正当率} = \frac{\text{正解を得た情報数}}{\text{全情報数(51件)}}$$

を用いる。

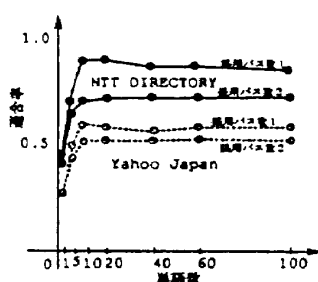


図 3 : 採用単語数の妥当性

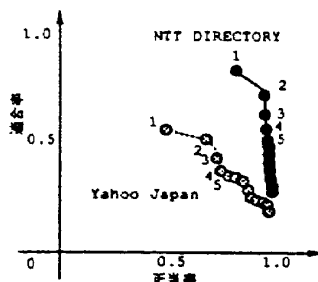


図 4 : 採用パス数の妥当性

図 3 は形態素解析で抽出された頻出単語の上位  $k$  個 (横軸) を用いてパスを検索し、上位 1 個又は 2 個のパスを採用した時の適合率 (縦軸) である。この図から分かるように、NTT DIRECTORY と Yahoo Japan のどちらを階層型知識体系として利用したときも、また採用するパスが 1 個でも 2 個でも、採用する単語の個数が 10 個を越えると適合率の上昇は見られない。よって適合率から判断すると上位 10 個の採用が妥当である。

図 4 は形態素解析で抽出された頻出単語の上位 10 個を検索語句として投入し、検索結果の上位  $k$  個のパスを採用したときの正当率 (横軸) と適合率 (縦軸) の

プロットである。またグラフ中の数字は  $k$  にあたる。この図から分かるように、NTT DIRECTORY と Yahoo Japan のどちらを階層型知識体系として利用したときも、採用パス数を増やしていくと、(1.0, 0.0) に収束していく事が分かる。また、正当率と適合率に同じ重みを持たせて目的関数を 正当率+適合率 とすると、これを最大にするのは  $k=2$  の時である。よって、検索結果からパスを 2 個採用するのが最も有用であるといえる。

上記の実験で得たパラメタである、抽出した頻出名詞の上位 10 個を検索語句として投入し、検索結果の上位 2 個のパスを採用する事により、NTT DIRECTORY で 92%、Yahoo Japan で 73% の正当率を得た。Yahoo Japan はカテゴリの分類が細かく、総カテゴリ数が NTT DIRECTORY の 5 倍以上保持しているため、本方法では良い正当率が得られなかったものと推測される。しかし、Yahoo Japan を階層型知識体系として用いた場合、NTT DIRECTORY を用いた場合よりも詳細に分類されるため、より詳細な情報を得ることが出来る。また、これが両者の構造の善し悪しの基準となるものではない。

### 5 まとめ

本稿では、既存の階層型知識体系を用いた WWW 情報の自動カテゴリ推定方法を提案した。また、階層型知識体系として NTT DIRECTORY, Yahoo Japan を用いて実験を行い、高い正当率でカテゴリが推定できることを確認した。

本方法は、情報フィルタリング手法の利用者興味の自動抽出を念頭に置いているが、現在手作業で行われているディレクトリ型検索サービスの情報の登録の支援や、ロボット型検索サービスの検索結果の再構築等への応用が考えられる。

今後の課題として、単語を抽出する時に、出現度数ではなく  $tf \times idf$  [3] で正規化した頻度を用いた場合との比較等を考えている。また、再現率を計測し、適合率とのバランスをも考察する予定である。

### 謝辞

本システム実現の為に Yahoo Japan のデータの利用を承諾していただいたソフトバンク (株) の方々に感謝いたします。

### 参考文献

- [1] <http://www.yahoo.co.jp/>
- [2] <http://navi.ntt.co.jp/>
- [3] G. Salton: Automatic Text Processing, Addison-Wesley, (1989).