

教育コンテンツ向け分類キーワード生成方法の提案とその実証

2S-6

梅村 武久*1*2、安田 孝美*1*3

通信・放送機構*1 NTT西日本 名古屋支店*2 名古屋大学 情報文化学部*3

1. はじめに

今後、教育シーンでのインターネット利用は必要不可欠なものとなることが予想される。しかし、限られた授業時間内に欲しい情報の発見や入手が困難であるといった問題がある。

この問題を解決するために、筆者らは情報の検索結果を教育向けカテゴリ分類体系に整理して表示することにより、検索結果の絞込みを支援するWWW検索システム(図1)を、通信・放送機構岡崎公共システム開発リサーチセンタに構築し、実用にも耐えうるWWW情報の分類体系化の基盤を確立した。また、学内ネットワークの普及によりイントラ型データベースへの検索利用度も増すため、WWW情報の分類体系化を、SQLやOracleに代表される汎用のデータベースサーバに対して適用させる方法を提案した[1]。

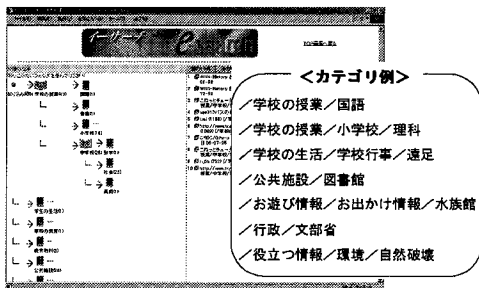


図1 WWW検索システム画面

2. 研究の目的

教育シーンでのインターネット・イントラネットの利用が増すにつれ、システムを運用管理するIT技術者が必要となり、その移動も増大していくなどの問題が発生する。よって、運用管理にかかる人手を最小限に押さえるシステムが要求される。

構築したWWW検索システムは、ロボット型(自動収

Proposal of method of generating classification key word for

Education contents and its Proof

Takehisa Umemura(NTTWEST Nagoya branch),

Takami Yasuda(Nagoya University)

集)であり、カテゴリ分類表示も分類キーワード(自動分類定義ファイル)により自動化されているが、分類キーワードの作成や情報更新に伴うメンテナンスについては、単語頻度などを基に人手に頼る点が多くある。そこで、本稿ではTF-IDF法(図2)を応用した分類キーワードの自動生成方法とその実証結果について報告する。

◆ 語の出現頻度から文書内の語の重要性を測る

$$TF(d, t) = \frac{\text{文書 } d \text{ における語 } t \text{ の出現回数}}{\text{文書 } d \text{ に現れる全語数}}$$

$$IDF(t) = \log\left(\frac{\text{全文書数}}{\text{語 } t \text{ を含む文書数}}\right) + 1$$

$$TF \cdot IDF = TF(d, t) \cdot IDF(t)$$

┌──────────┐
語 t に対するテキスト d の重要性

図2 TF-IDF法

3. 分類キーワード生成方法

WWW情報の分類手法としては、URL名による分類やハイパーリンクの共起による分類[2]など様々な手法が考えられる。また、情報検索分野ではキーワードの抽出方法としてTF-IDF法などを用いた研究[3]も進められている。本研究において、WWW情報分類にキーワードによる手法を選定したのは、これらの様々な手法と組み合わせることが可能となるなど、幅広く応用できることがある。また、分類キーワードを格納した定義ファイルについても、図3のように汎用的なCSV形式のテキストファイルを使用している。

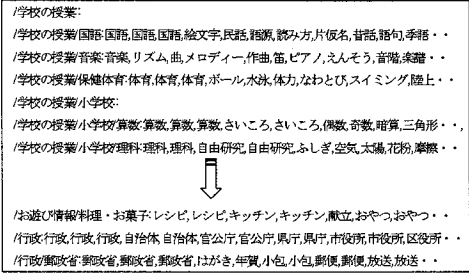


図3 自動分類定義ファイル

次項にて、本定義ファイルの分類キーワードをTF-IDF法を応用して自動抽出する方法について述べる。

3.1 分類キーワードの自動生成処理フロー

IDF (inverse document frequency: 逆文書頻度) 値を各単語の文書に対する特徴度とし、ある一つのカテゴリ内の文書に対する特徴度 a と全カテゴリ内の文書に対する特徴度 b からカテゴリ適合度 (特徴度 b - 特徴度 a) を算出、カテゴリ適合度の大きな単語 (キーワード) を分類キーワードとして抽出する。すなわち、ある一つのカテゴリ文書に対する特徴度が小さく、かつ全カテゴリ文書に対する特徴度が大きい単語を分類キーワードとして抽出する。また、カテゴリ適合度に TF (term frequency: 単語頻度) 値を乗算することでキーワードの重要度を算出し、その大小に応じて分類キーワードの重み付けを行う。初めに、極少数の Web ページを収集し人手を使って分類する必要はあるが、その後は自動 (プログラミング) 化することが可能となる。図4に、自動生成の処理フローを示す。

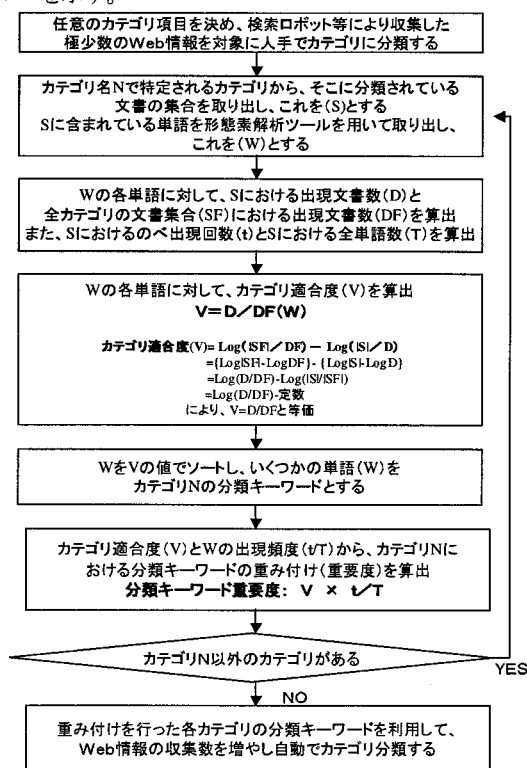


図4 分類キーワード自動生成処理フロー

3.2 実証結果

図4の処理フローが動作するプログラムを、実験運用中のWWW検索システム上に作成し、実証実験を行っ

た。ある一つのカテゴリ(学校の授業/音楽)に分類したWeb情報25ページを対象に、各単語のカテゴリ適合度及び重要度を算出した結果の一部(上位25単語)を表1に示す。

カテゴリ毎ではばらつきはあるものの、表1のようにほぼ妥当な分類キーワードを抽出することができた。抽出した各カテゴリの分類キーワードを用いて約2万ページのWeb情報をカテゴリ分類した結果、全126カテゴリ中ランダムに選択した10カテゴリにおける分類精度は約70%であり、まだサンプル数は少ないがこれまで人手を中心に作成してきた分類キーワードの分類精度とほぼ同様の結果が得られた。

表1 カテゴリ適合度、重要度の算出結果(抜粋)

No	カテゴリ	カテゴリ適合度 >=5	単語頻度 >=5	ページ数/ 該当カテゴリ	ページ数/ 全カテゴリ	重要度 評価値×(単語 頻度/該当 カテゴリ)
1	学校の授業/音楽					
2	楽譜	0.272727273	49	9	33	5.923597679
3	伴奏	0.777777778	9	7	9	3.102836879
4	review	1	7	1	1	3.102836879
5	tempo	1	7	1	1	3.102836879
6	十音	1	6	1	1	2.653574468
7	演奏	0.144977336	36	10	69	2.312673451
8	ピアノ	0.212121212	24	7	33	2.246563664
9	ソロ	1	5	2	2	2.216313357
10	音楽	0.070422535	66	20	284	2.050233743
11	cream	0.5	8	1	2	1.773045643
12	曲	0.191488352	16	3	47	1.348060570
13	バイオリン	0.3	6	3	6	1.329787244
14	タブ	0.25	12	1	4	1.329787244
15	琵琶	0.333333333	6	2	6	1.189083097
16	楽器	0.126666667	13	1	6	1.106150228
17	演奏	0.333333333	6	1	3	0.866524823
18	打ち込み	0.333333333	6	1	3	0.866524823
19	作曲	0.222222222	9	4	18	0.866524823
20	合唱	0.137931034	14	4	29	0.855953001
21	オーケストラ	0.333333333	3	4	11	0.803931057
22	オペラ	0.142857143	12	2	14	0.758779419
23	アンサンブル	0.333333333	3	4	12	0.758779686
24	音	0.093663459	33	8	217	0.532268556
25	オルガン	0.222222222	9	2	9	0.49251379

4. おわりに

教育分野を対象にWWW情報の分類キーワード生成方法とその実証結果について述べた。実証結果より本手法の有効性が把握できたので、今後は更にWeb情報の収集数を増やすと伴に、分類精度の向上を目指した手法の改善を行っていき教育向けカテゴリ分類体系を確立していきたい。

参考文献

- [1] 梅村武久, 安田孝美: 教育分野における WWW 情報の分類体系化, 第60回情処全大, 5M-06, (2000). 教育分野における WWW 情報の分類体系化とその実証, 第61回情処全大, 5S-05, (2000). 教育コンテンツ向け分類キーワード生成方法の提案, 第62回情処全大, 8Y-07, (2001).
- [2] 大久保雅且, 杉崎正之, 田中一男: リンクの共起関係を用いた Web ページ分類法式の検討, 第 59 回情処全大, pp3-81~82, (1999).
- [3] 中川 裕志: 言語メディア論, <http://www.r.dl.itc.u-tokyo.ac.jp/~nakagawa/lmedia/skelton/skelton.html>