

E-012

ニュース検索のための格構造を用いたユーザの興味表現と分類手法

A Representation and Clustering of User's Interests by Using Case Structures for Retrieval of News Documents

東原 智幸†
Tomoyuki Higashihara

渥美 雅保†
Masayasu Atsumi

1. まえがき

インターネット上の文書数の増大によりユーザ自身の求める文書を検索することが難しくなっている。近年、その問題を解決するため、ユーザが興味を示す文書をシステムが自動的に検索し、提示する研究が多く行われている[1][2]。また、ブログなどの文書から話題(トピック)を抽出するサービスも提供されている。それに伴って、ブログからのトピック抽出について研究されており、[3]では文書ベクトルに基づいて、クラスタリングを行い、得られたクラスタをトピックとしている。その際、興味や文書の表現としてtf・idf法が多く用いられているが、文の係り受け構造や格情報などの情報が利用できない点に問題がある。本研究では、格構造の集合からユーザの興味概念構造を表現し、それらをニュース検索に適用する方法について提案する。ユーザがニュースに興味を示した場合、そのニュースのすべてに興味として保持することは興味とは関係のないデータもふくまれるため、検索率の減少の原因となる。本システムでは、ユーザが興味を示したニュースに対して、展望台システム[4]を適用しニュースの主題を選択する。その結果からトピック構造集合を抽出する。次に、その集合をクラスタリングし、ユーザの興味概念構造を作成する。

最後に、興味概念構造を使用した検索実験を行い、ユーザの興味に適應した検索が行えるか評価を行う。

2. 興味構造

興味構造はトピック構造集合をクラスタリングすることにより構築する。以下でこれらについて説明を行う。

2.1. トピック構造

トピック構造は、ニュース中の主題と推測される文章内に含まれる格構造集合である。

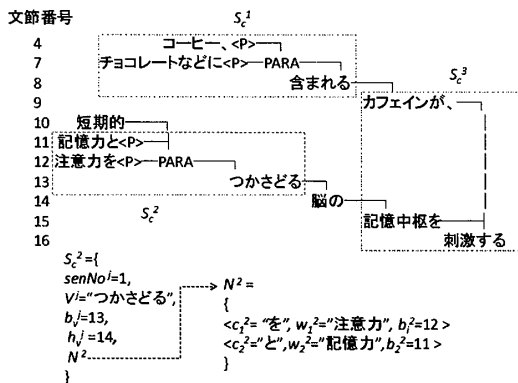


図1 構文解析結果と格構造(一部略)

ニュース中のj番目のトピック構造 S_c^j (図1)は $S_c^j = \{senNo^j, V^j, b_j^j, h_j^j, N^j\}$ (1) と表現される。 $senNo^j$ は、ニュース中の文番号、 V^j は動詞、 b_j^j は動詞の文節番号、 h_j^j は動詞の係り先文節番号である。 N^j は動詞 V^j に係るM個の名詞集合で、表層格 c_m^j 、その格の要素(単語または複合語) w_m^j 、その要素の文節番号 b_m^j の組の集合か

らなり、 $N^j = \{<c_1^j, w_1^j, b_1^j>, \dots, <c_m^j, w_m^j, b_m^j>, \dots, <c_M^j, w_M^j, b_M^j>\}$ (2) と表される。

例えば、ニュース中の1番目の文章に“カフェインが記憶中枢を刺激する”が存在する場合は、格構造は $\{1, \text{“刺激する”}, 16, -1, \{<“が”, \text{“カフェイン”}, 9>, <“を”, \text{“記憶中枢”}, 15>\}$ と表現される。動詞に係る文節がない場合には、係り先文節番号は-1となる。

2.2 トピック構造のクラスタリング

クラスタリングでは、トピック構造類似度と単語間の概念類似度を用いる。

2.2.1 トピック構造間類似度

トピック構造間類似度 $s(x, y)$ はトピック構造 $x=S_c^i$ と $y=S_c^j$ 間の概念的、格構造的な類似度を求める式で、以下の式により定義される。

$$s(x, y) = sim(V^i, V^j) + avr(N^i, N^j)$$

$sim(V^i, V^j)$ は動詞の概念類似度と $avr(N^i, N^j)$ は名詞集合の概念類似度平均である。名詞集合間の概念類似度平均は、

$$avr(N^i, N^j) = \frac{1}{K} \sum_{k \text{ s.t. } c_k^i=c_k^j} sim(w_k^i, w_k^j)$$

で表され、表層格の同じ単語同士概念類似度の和の平均を求めたものである。Kは表層格の同じ単語の組数である。

2.2.2 単語間の概念類似度

単語間の概念類似度は、概念体系辞書[5]を用いて概念の深さを検索し、その値を用いて類似度を計算する。

単語 w_0, w_1 の概念体系の深さを d_0, d_1 、共通の深さを d_c とするとき、単語間の類似度は

$$sim(w_0, w_1) = \begin{cases} \frac{2 \times d_c}{d_0 + d_1} & \dots w_0, w_1 \text{ が登録されている} \\ 0 & \dots \text{それ以外} \end{cases} \quad (3)$$

で表される。 w_0, w_1 が辞書には登録されていないが、形態素要素に分割できる場合には、各形態素要素毎に類似度を計算し、それらの最大値を $sim(w_0, w_1)$ とする。

2.2.3 クラスタリング

クラスタリングは抽出されたトピック構造集合TpSetに対して実行される。トピック構造集合TpSetは、主題文1文ごとに形態素解析[6]・構文解析[7]を行い抽出される。

$$TpSet = \{S_c^1, \dots, S_c^j, \dots, S_c^J\} \quad (4)$$

ここで、Jはトピック構造の総数である。 S_c^j の名詞集合 N^j には、 V^j に直接係っている名詞 w_l^j 、その名詞と並列(PARA)関係または同格関係である名詞が含まれる。同格関係とは“長男太郎”などのように複数の分節が対等で同一の対象を示している関係である。

抽出したTpSetを用いて階層併合的クラスタリング[10]によりクラスタリングを行う。基本的なアルゴリズムは以下の通りである。ここで、Gはクラスタ集合、Cはクラスタ数、indexはクラスタ番号である。

1. $S_c^j \in TpSet$ をクラスタ $G_i = \{S_c^j\}$ として割り付ける。C=J, index=C+1と設定する
2. 各データの類似度を以下の式にて計算する $s(G_i, G_j) = \min_{x \in G_i, y \in G_j} s(x, y)$
3. 類似度最大のクラスタ対 G_q, G_r を結合する。Gに $G_{index} = G_q$

†創価大学大学院工学研究科情報システム工学専攻

UG_rを追加し, G_q, G_rをGから除く. C=C-1と更新する. G_{index}を作成する際に, G_q, G_r間共通の概念を求め, G_{index}の概念とする. G_q, G_r間の動詞の組, 表層格が同じ名詞集合の組でそれぞれ共通概念を求め, 共通概念は, 概念体系の共通の深さd_iにある概念である.

- もし, C=1ならば終了. C>1ならばi<indexなるすべてのG_iについてs(G_{index}, G_i)を計算後, index=index+1と変更し, ステップ3に戻る.

図2はクラスタリングされた結果を樹状図で表現した図である. クラスタに概念ID(動詞のみ表示), クラスタ番号, クラスタ間の類似度が付与されていることが分かる.

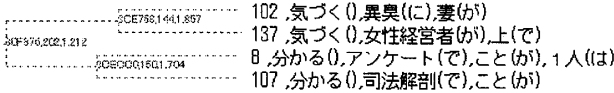


図2 樹状図(一部)

2.3 興味概念構造による検索

上位のクラスタに与えられた概念を検索に用いることにより概念的に広い情報を検索することができる. 例えば, 図2の{「気づく」, «に»;「異臭」, «が»;「妻」}と{「気づく」, «が»;「女性経営者」, «で»;「上」}では上位クラスタの概念は{「気づく」, «が»;「女性」という概念}となり, より広い「女性」という概念に近い主語をもつ「気づく」にも興味があると考えることが可能となる.

3. ニュース検索実験

3.1 テスト環境

本実験では, ニュース提供システム(図3)[8]において, 実験を行った. ニュース提供システムは, ユーザの要求を処理する音声対話部, ユーザの興味概念構造を管理する興味構造管理部, ユーザに提供するニュースを検索するニュース選択部から構成される. 興味構造管理部とニュース選択部の一部を実装している.

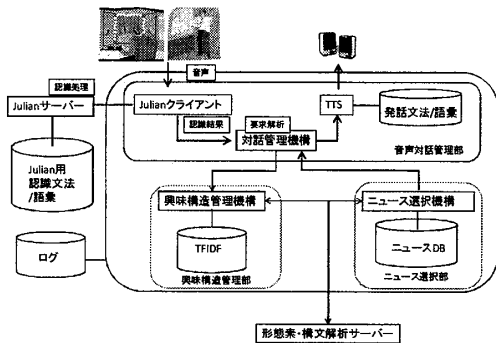


図3 システム構造

3.2 評価データの作成

以下の手順でユーザのニュースに対する評価を収集し, ニュース検索性能に対する評価データを作成する. 使用するニュースは, 毎日新聞2007年版[9]の社会面, 1月1日~17日のニュース500件(2日は0件)である.

- ニュースをユーザに提示
- 提示されたニュースに対して A (興味がある), B(どちらかといえば興味がある), C(どちらかといえば興味はない), D (興味はない) の4段階で評価する.

3.2 評価結果

興味概念構造の深さの違いが検索性能に及ぼす影響について評価を行った. 興味概念構造の深さは, 葉の部分から根に向かって大きくなっていく. 深さをDとし, D=(0, 1, 2, 5)に変更し, それぞれの深さに対する検索結果(図4)を調査した. 図4は, 上位N位(N=1~20)まで検索した場合の精度を示している. 評価データを5つに分割し, 1つ選択し, 興味概念構造を作成, 残りの4つを評価用データとして使用した. 20位における精度値は, D=0で0.663, D=1で0.630, D=2で0.640,

D=5で0.615となっている. 5位以降では, D=0で18位0.681, D=1で9位0.672, D=2で17位0.653, D=5では7位0.657が最大であった. また, 20位までの平均では, D=0で0.666, D=1で0.652, D=2で0.633, D=5で0.634であった. 深さが大きくなっても精度に大きな違いは見られなかった. 深くなれば検索に用いるクラスタ数も少なくなるが, クラスタは下位のクラスタの概念も含んでいるため検索基準数の影響が少なかったと考えられる.

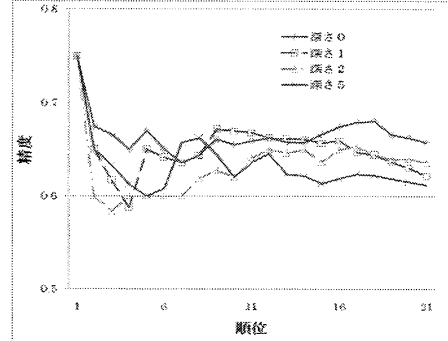


図4 概念の深さと検索精度

表1は, 興味数評価データごとの興味総数の順位において検索された興味ニュースの割合の表である.

表1 興味の割合

データ番号	興味総数(順位)	最小値(D)	最大値(D)
データ1	55	0.605(0)	0.618(2)
データ2	43	0.564(5)	0.605(0)
データ3	24	0.480(5)	0.542(0)
データ4	42	0.464(2,5)	0.482(0)
データ5	57	0.711(1,5)	0.754(2)

4. まとめ

ユーザが興味を示した格構造の集合からユーザの興味概念構造を表現し, ニュース検索において, 興味概念構造の深さを基準とした検索について評価を行った. 今後は, 古い興味の削除や, 興味構造の更新(追加), トピック構造の文章中における重要度を考慮したクラスタリング方法について検討を行う.

参考文献

- 河合, 熊本, 田中: 印象と興味に基づくユーザ嗜好のモデル化手法の提案とニュースサイトへの応用, 日本知能情報ファジィ学会誌, vol.18 No.2, pp.173-183, 2006.
- 野美山, 紺谷, 渡辺他: 個人適応型情報検索システム—個人の興味を学習する階層記憶モデルとその協調的フィルタリングへの適用—, 情報学基礎研究会報告, Vol.96, No70, pp.49-56, 1996.
- 戸田, 黒田, 福田, 石川: ブログにおける多視点からのトピック抽出手法の提案, DEWS2008(第19回データ工学ワークショップ), B4-2, 2008.
- 相良, 砂山, 谷内田: サブトピックを考慮した重要文抽出による報知的要約生成, 電子情報通信学会論文誌 D, Vol.J90-D, No.2, pp.427-440, 2007
- 日本電子化辞書研究所, EDR 電子化辞書(第2版)仕様説明書, TR2-006(改), 2001.
- 黒橋: 日本語形態素解析システム JUMAN ver. 5.1, 東京大学大学院情報理工学系研究科, 2005.
- 黒橋: 日本語構文解析システム KNP ver. 2.0, 東京大学大学院情報理工学系研究科, 2005.
- 東原, 三吉, 渥美: スマートホームにおける音声ニュース提供システムアーキテクチャの構築, FIT2008(第7回情報科学技術フォーラム), E-025, p.193-195, 2008
- CD-毎日新聞(データ集)2007年版
- データマイニングの基礎, 元田ほか, オーム社, 2006