

# Web 上における単語の意味関係の変化を 把握する手法の提案

関 良太

岸 義樹

茨城大学大学院理工学研究科

茨城大学工学部情報工学科

## 1. はじめに

近年, 単語の意味関係を考慮した Web 検索に関する研究が盛んである. しかし, そこで利用される意味資源は予め作成されたものであり, 更新されることも少ないので, 極めて静的な性質を持つ. 一方で Web は刻々と変化するので, 逆の性質を持つ. この性質の違いから, 上記のような意味資源を Web に利用した場合, その有効性に限界があると考えられる. 意味資源を Web に対して有効なものとするためには, Web 上における単語間の意味関係の変化も知識として把握しておく必要がある.

そこで本研究では, 単語間の意味関係の変化を, 意味構造の 1 つである意味ネットワークの変化から把握する手法を提案して, その有効性を検証した.

## 2. 提案手法

我々が提案する手法は次の通りである.

1. 検索期間  $P = \{p_1, p_2, \dots, p_l\}$  と単語の集合  $T = \{t_1, t_2, \dots, t_n\}$  を与え,  $T$  からクエリの集合  $Q = \{q_1, q_2, \dots, q_n\}$  を作る.
2.  $P$  の要素  $p_i (i=1, 2, \dots, l)$  において,  $Q$  の要素で Web 検索を行い, Web 情報  $w_i$  を取得し, Web 情報の集合  $WI$  に加える.
3.  $WI$  の要素  $w_i (i=1, 2, \dots, l)$  それぞれから, GA で意味ネットワーク  $sn_i (i=1, 2, \dots, l)$  を作り, 意味ネットワーク集合  $SN$  に加える.
4. 3. で作成した  $SN$  中の要素を比較し,  $T$  中の要素の意味関係の変化を把握する.

### 2.1 検索クエリ

$Q$  の要素  $q_{n \times (i-1)+j} (1 \leq i, j \leq n)$  の定義は下式の通りである. 式中の AND は AND 検索を行うことを意味している.

$$q_{n \times (i-1)+j} = \begin{cases} t_i \text{ AND } t_j & (i \neq j) \\ t_i & (i = j) \end{cases}$$

### 2.2 取得する Web 上のデータ

取得する Web 情報  $w_i \in WI$  について説明する. まず,  $Q$  の要素  $q_{n \times (i-1)+j} (1 \leq i, j \leq n)$  それぞれを検索クエリにして検索したときの Hit 件数と上位 50 件分のレスポンスを, Yahoo! API([1])を用いて取得する. さらに, そこから上位 5 件分(1 位から 5 位)と下位 5 件分(46 位から 50 位)の URL を抽出して, それが示す Web コンテンツをダウンロードし, テキスト部分のみを抽出する. 取得した Hit 件数を  $hit_{n \times (i-1)+j}$ , そしてテキスト情報を  $txt_{n \times (i-1)+j}$  として,  $w_i$  を以下で定義する.

$$w_i = \{(hit_1, txt_1), (hit_2, txt_2), \dots, (hit_n, txt_n)\}$$

### 2.3 GA (遺伝的アルゴリズム)

#### 2.3.1 コード化手法

遺伝子  $g_{n \times (i-1)+j} (1 \leq i, j \leq n)$  を以下で定義する.

$$g_{n \times (i-1)+j} = \text{単語 } t_i \text{ と単語 } t_j \text{ の関係の強さ}$$

関係の強さは, 弱い方から強い方へ順に, 0 から 5 までの 6 段階の重みを割り当てることで表現する.

そして個体は, この遺伝子の並びで表現される. すると単語の関係性を全て表現していることになるので, 個体そのものを意味ネットワークと見なすことが出来る.

#### 2.3.2 個体評価手法

構成する意味ネットワークは, Web 上での単語の関係を表現している必要がある. そのためには個体を評価する指標が, Web 上の情報に基づいていなければならない([2], [3]).

個体中の各遺伝子  $g_{n \times (i-1)+j} (1 \leq i, j \leq n)$  における単語  $t_i$  と単語  $t_j$  の関係について 4 つの指標に基づいたスコア付けを  $w_i$  の要素  $(hit_{n \times (i-1)+j}, txt_{n \times (i-1)+j})$  で行い, その結果から個体の適応度を計算する.

A proposal of method for understanding change of the semantic relationship between words in Web

Ryota Seki, Yoshiki Kishi  
Ibaraki University

4-12-1 Nakanarusawa, Hitachi, Ibaraki, 316-8511, JAPAN

1. シンプソン係数に基づいたスコア付け;  $sim_k$   
 シンプソン係数([2])は、キーワード X と Y の共起の強さを表す指標である。キーワード X, Y の Hit 件数をそれぞれ  $|X|$  と  $|Y|$ , X と Y の AND 検索での Hit 件数の結果を  $|X \cap Y|$  としたとき、このスコアは次式で与える。

$$sim_k = |X \cap Y| / \min(|X|, |Y|)$$

2. 初出単語距離に基づいたスコア付け;  $FTD_k$   
 $txt_k$  から名詞、動詞、形容詞を抽出し、リスト化する。距離を測りたい単語を X と Y としたとき、初出単語距離([3])はリスト中で最初に現れた X と Y の距離、つまりインデックスの差で表現される。この距離を  $d$  としたとき、このスコアは次式で与える。

$$FTD_k = 1/d$$

3. 最小単語距離に基づいたスコア付け;  $MTD_k$   
 使用リストは初出単語と同じである。最小単語距離([3])は、リスト中の X と Y の距離の中で最小のものである。この距離を  $D$  としたとき、このスコアは次式で与える。

$$MTD_k = 1/D$$

4. 単語の出現密度に基づいたスコア付け;  $AD_k$   
 リスト中の全単語のうち、目的の単語 X, Y が出現する割合のことを表す。このスコアは次式で与える。

$$AD_k = X \text{ と } Y \text{ の出現回数} / \text{リストの総単語数}$$

以上の評価に加え、意味ネットワーク中でリンクを結ぶには相応のコストが必要であると仮定し、結ぼうとするリンクの重みが大きければ大きいほど、コストがかかるように設定した。

これらの要素を考慮して、個体の適合度  $f$  を次式で定義する。

$$f = \sum_{i=1}^n \{ (sim_i + FTD_i + MTD_i + AD_i) \times g_i - \text{コスト} \}$$

## 2.4 変化の把握手法

作成した意味ネットワークから、1 週間ごとの各遺伝子の重みの平均をとる。その平均値が 4 以上であれば強い関係、1 以下であれば弱い関係であるとする。その関係の変化を見ることにより、単語間の意味変化を把握する。

## 3 実験

### 3.1 内容・データについて

平成 13 年度版情報通信白書の文章から「ネットワークインフラ」、「高品質音声サービス」、「次世代携帯電話」、「本格サービス開始」、「ITU 世界

無線通信会議」、「商用サービス」、「試験サービス」、「ジェイフオングループ」、「デジタル方式」、「総務省」という 10 語を抽出して、関係の変化を見た例を示す。Web 情報の収集期間は 2009 年 11 月 15 日から 12 月 12 日までである。GA の交叉は 1 点交叉、選択はルーレット戦略を用い、交叉率は 0.8%、突然変異率は 0.3% としている。コストは、0 から 5 のリンクの重みそれぞれに対して、0.00, 0.01, 0.02, 0.03, 0.04, 0.05 と設定した。

### 3.2 結果

実験で得られた、強い関係、弱い関係を示した遺伝子の数の推移を表 1 に示す。

表 1: 強い関係、弱い関係の遺伝子数の推移

期間 関係	第 1 週	第 2 週	第 3 週	第 4 週
強い	9	9	7	11
弱い	1	0	0	0

また、4 週を通して関係の変化が起きなかった遺伝子数は強い関係を持つもの、弱い関係を持つもの、ともに 0 であった。

## 4 考察

表 1 において強い関係を持つ遺伝子の数の変化や、4 週間ずっと同じ関係を持ち続けた遺伝子は無いことから Web 上における単語の関係が変化していることが分かる。また、弱い関係を持つ遺伝子数が 1 や 0 であり少なかった。これは、弱い関係の基準が厳しく、本来ならば該当すべきものまで除いてしまっていることが考えられる。

## 5 まとめ

本研究では Web 上における単語間の関係の変化を把握する手法を提案して、実際に実験した。その結果、本手法で関係の変化を把握できることが分かった。

## 参考文献

- [1] Yahoo! デベロッパーネットワーク  
<http://developer.yahoo.co.jp/>
- [2] 松尾 豊, 友部 博数, 橋田 浩一, 中島 秀之, 石塚 満: Web 上の情報からの人間関係ネットワークの抽出, 人工知能学会論文誌, Vol.20, No.1E, pp.46-56, 2005
- [3] 田 馳, 手塚 太郎, 小山 聡, 田島 敬史, 田中 克己: 質問キーワードの意味的関連と近接性に着目したウェブ検索の精度改善, 電子情報通信学会第 17 回データ工学ワークショップ (DEWS2006) 論文集, 5A-o1, 2006