

3G-10

## 動的シソーラスを用いた連想検索 ～リンク重みの導入～

巖寺俊哲 木本晴夫  
NTT情報通信網研究所

### 1. はじめに

近年、コンピュータは著しい技術進歩、普及を遂げている。また、種々の情報のデータベース化も進展し、それらのデータベースは、ネットワークを介して、容易に利用可能になっている。さらに、専門家だけでなく、いままでコンピュータと無縁であった一般ユーザにとってもコンピュータやデータベースが身近に、また、必要不可欠なものになってきており、自らそれを操作する必要が生じている。

従来、データベース検索を行う際には、多くの場合、キーワードをAND/ORで組み合わせた検索式を使用していた。しかし、目的の情報を検索するためには、適切なキーワードを発見し、それらを適切に組み合わせた検索式を作成する必要があり、容易には、目的の情報を検索することができなかった。

現在、我々は、目的の情報を高精度で、容易に、検索することを可能にするために、動的シソーラスを用いた連想検索法の検討を進めている[1, 2]。

本稿では、連想検索法の検索精度をさらに向上させるために導入した動的シソーラスのリンク重みとこれを用いたキーワードの関連度の算出法について報告する。また、実験により本方法が適合率の向上に有効であることを示す。

### 2. 動的シソーラスを用いた連想検索

連想検索法は、以下の4点を特徴とする。

- (1) サンプル文書の学習
- (2) 動的シソーラス
- (3) 連想キーワードの生成
- (4) 検索結果の順位付け

サンプル文書の学習では、ユーザから提供された検索の目的とする文書の見本(サンプル文書)から検索の目的に関して特徴的なキーワードやキーワードの関係を学習する。さらに、それを動的シソーラス上に反映する。

動的シソーラスは、ノードとリンクとから構成される。ノードは、検索に使用され得る語句を示す。リンクは、その語句間の関係を示す。ノードは、その語句の重要度を示すノード重みを持つ。また、リンクは、語句間の関係の強弱を示す本稿で述べるリンク重みを持つ。

連想キーワードの生成では、ユーザの入力したキーワードから動的シソーラスを使用してユーザの検索目的により適したキーワード(連想キーワード)を生成する。

検索結果の順位付けでは、連想キーワードを使用して検索された検索結果文書を動的シソーラスと照合する。さらに、その動的シソーラスとの適合の程度に応じて順位付けする。

従来の連想検索法では、入力キーワードに対応するノードと関係の強弱にかかわらずリンクで結合しているノードを連想キーワード候補とし、ノード重みを使用して連想キーワードとして生成するか否かを決定していた。この方法では、キーワード間の関係の強弱がまったく考慮されていない。このため、ノード重みを併用せず、リンクのみを連想キーワード生成に使用した場合、再現率は向上するが、適合率は低下する。これは、リン

クが、適切なキーワードの生成には有効に機能しているが、不適切なキーワードも生成してしまっているためである。さらに検索精度、特に、適合率を向上させるためには、より適切な関係にあるキーワードのみを生成する必要がある。

### 3. リンク重みの導入

前述したように、適切な連想キーワードを生成するためには、より適切なリンク、すなわち、関係がより強いリンクを使用する必要があり、その関係の強弱を表現する尺度が必要である。このリンクに付与される関係の強弱を示す値をリンク重みとよぶ。

このリンク重みは、次の条件を満たす必要がある。

(1) リンクは、検索目的に関して、必要な連想キーワードを生成している(再現率向上に寄与)[1]。この機能を損なうことがない。

(2) リンク重みを導入することによって必要なキーワードのみを生成可能にする(適合率の向上)。

上記条件を満たすリンク重みを連想キーワードの生成に使用することで検索精度を向上させることが可能である。

#### 3.1 リンク重みの付与方法

連想検索法では、リンクの生成に共起関係を使用している[1]。この共起関係の頻度をその強弱ととらえ、これを使用してリンク重みを付与する。

リンクの生成と同様にサンプル文書中でのノードに対応するキーワードどうしの共起頻度を使用して算出する。共起頻度とは、ある2つのキーワードが共起したサンプル文書数である。ここで全サンプル文書数がDN、ノードiとノードjとに対応するキーワードの共起頻度がDN<sub>ij</sub>のとき、ノードiとjを結合するリンクのリンク重みLW<sub>ij</sub>は、次式により算出される。

$$LW_{ij} = \frac{DN_{ij}}{DN}$$

上式により算出されたリンク重みは、必ず共起するノード間のリンクでは1に、まったく共起しないノード間のリンクでは、0になる。

#### 3.2 リンク重みを使用した連想キーワードの生成方法

##### (1) 連想キーワード生成手続き

以下に述べる手順で連想キーワードを生成する。

- ① ユーザの入力キーワードを後述するサンプル文書中での出現確率で点数付けする。
- ② 入力キーワードに対応するノード(生成開始ノード)からあらかじめ指定した距離までリンクで結合しているノードを連想キーワードの候補とする。ここで距離は、生成開始ノードからあるノードまでの経路上にあるリンクの本数である。
- ③ 各連想キーワード候補の入力キーワードとの関連度をリンク重みを用いて計算する。関連度は、後述する方法で求める。こ

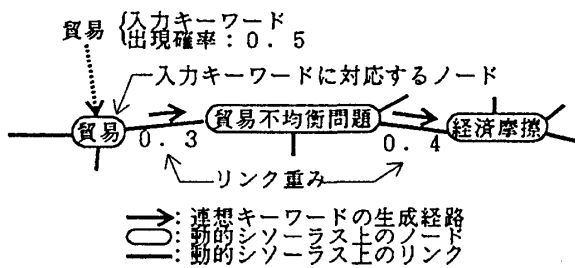


図1 関連度の算出

のとき、複数の入力キーワードから到達し得る候補は、①で求めた点数のより高い入力キーワードとの関連度を採用する。  
④あらかじめ指定した値（関連度しきい値）以上の関連度を持つ連想キーワード候補を連想キーワードとして生成する。  
ユーザの検索目的に関して、関連度の大きい連想キーワードほど入力キーワードと強い関係を持つ。

(2) サンプル文書中での出現確率

あるキーワード  $i$  のサンプル文書中での出現確率  $OP_i$  は、そのキーワード  $i$  が出現したサンプル文書が  $DN_i$ 、全サンプル文書数が  $DN$  のとき、次式により算出される。

$$OP_i = \frac{DN_i}{DN}$$

(3) 関連度

入力キーワード  $i$  と連想キーワード  $j$  の関連度  $RW_{ij}$  は、次式のように、入力キーワードの出現確率  $OP_i$  と連想キーワード  $j$  を生成するために通過した動的シソーラス上のリンク ( $i, \dots, n$ ) のリンク重み  $LW_i, \dots, LW_n$  の積として算出する。

$$RW_{ij} = OP_i \times (LW_i \times \dots \times LW_n)$$

たとえば、動的シソーラスの状態が図1のような場合、入力キーワード「貿易」から生成された連想キーワード「経済摩擦」の関連度  $RW$  は、次のようになる。

$$RW = 0.5 \times (0.3 \times 0.4) = 0.06$$

4. 実験による評価

関連度しきい値の変化とこれに対応する検索結果の再現率、適合率の変化を測定し、リンク重みを用いて算出した関連度の有効性を評価した。

4.1 実験方法

新聞記事161件が格納されている文書データベースを対象とし、二人の被験者によって入力されたキーワードから上記方法で生成された連想キーワードを使用して検索実験を行った。学習させるサンプル文書としては、被験者が検索目的に合致していると認定した文書データベース内の文書の半数を使用した。関連度しきい値は、0~100範囲で変化させ、その値以上

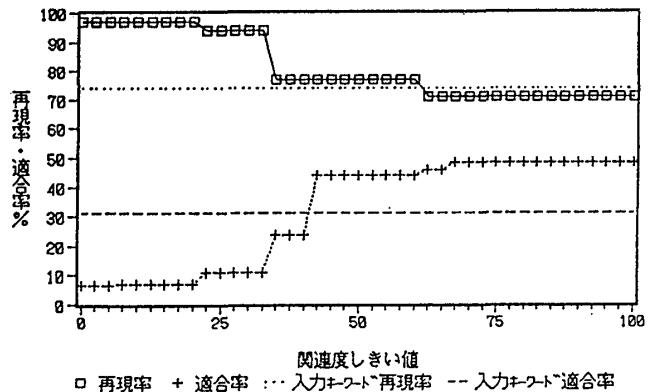


図2 関連度しきい値と再現率・適合率の関係

の関連度を持つ連想キーワードのみを検索に使用した。このとき、前述の算出法で求めた関連度は、あらかじめ0~100の範囲に正規化する。これにより、連想キーワードのもつ関連度の最大値は、100に、最小値は、0になる。  
今回は、リンク重みの有効性を評価するためにノード重みは、使用しない。

4.2 実験結果

実験結果を図2に示す。

実験結果では、リンク重みを使用して算出した関連度のより大きい連想キーワードを使用して検索することによって、再現率は、入力キーワードのみを検索に使用した場合とほぼ同等の値を維持し、かつ、適合率が向上している。

これは、リンクを使用することによって入力キーワードを、ユーザの検索目的に関して、より適切なキーワードへ変換可能であり、そのリンク重みの強いキーワードほど適切な関係にあることを示している。

さらに、上記のことは、共起関係が入力キーワードからユーザの検索目的に対してより適切なキーワードを生成する上で意味があり、その共起関係の強度は、共起頻度から算出し得ることを示している。

5. おわりに

本稿では、より適切なキーワードを生成するために導入した動的シソーラスのノード関係の強弱を示すリンク重みとこれを用いた連想キーワードの関連度の計算方法について報告した。さらに、実験によって関連度の有効性を評価した。その結果、リンク重みが、より適切な連想キーワードの生成に有効であり、検索精度、特に適合率を向上させるのに有効であることが確認できた。さらに、ノード重みと組み合わせて連想キーワードの点数付けに使用することで検索精度をさらに向上させることが期待できる。

[参考文献]

[1] 巖寺, 木本: 「動的シソーラスを用いた連想検索」, 情報学会自然言語処理研究会資料90-NL-75, 1990.  
[2] H. Kimoto, T. Iwadera: "Construction of a Dynamic Thesaurus and Its Use for Associated Information Retrieval," Proc. of 13rd International Conference on Research & Development in Information Retrieval, ACM SIGIR'90, Brussels, 1990.