

協調サーチエンジンにおける近接検索

4X-01

上原 稔

東洋大学情報工学科

1. はじめに

インターネットにおける Web ページの増加に伴い AltaVista、Google などの集中型サーチエンジンでは新鮮な情報を検索するのが困難となってきた。そこで、我々は分散型アーキテクチャに基づく協調サーチエンジン (Cooperative Search Engine, CSE)[1]を開発した。CSE では、各サーバに配置された局所的サーチエンジンをメタサーチエンジンで統合し、大域的サーチエンジンを構成する。

膨大な Web ページから適切な文書を検索するには、ブーリアン検索以上に詳細な条件指定が必要である。Lycos[3]などでは豊富な近接演算をサポートすることで高度な検索要求に応えている。近接演算 (proximity operation) とは語順や語間の距離を考慮した演算である。近接演算は絞り込みにも有効であるが、インデックスのサイズが増大するため、小規模なサーチエンジンではサポートしていないことが多い。我々の CSE では局所的サーチエンジンとして namazu[2]を利用可能であるが、namazu も近接演算はサポートしていない。そこで、我々は CSE のために近接演算機能を持つサーチエンジンを開発した。

本章の構成は以下の通りである。2 章では近接演算について説明し、3 章では我々の仕様を述べ、4 章では評価について述べる。最後に結論を述べる。

2. 近接演算

Lycos[3]には豊富な近接演算子が用意されている。

$A \text{ adj } B$ A と B が 1 語以内

$A \text{ near } B$ A と B が 25 語以内

$A \text{ far } B$ A と B が 25 語以上

これら演算子の頭に “o” が付くと語順を考慮する。また、語間距離と無関係なものもある。

$A \text{ before } B$ 語順の A かつ B

STARTS[4]では 2 語 A 、 B の関係として以下の近接演算子を提案している。

$A \text{ prox}[n, X] B$

ここで、 n は AB 間の距離を表す自然数で、 A と B の間の距離が n 以下であることを示す。また、 X は語順を考慮するか否かを示すブール値である。STARTS の演算子を用いると far 以外の Lycos の近接演算子を表現できる。

$A \text{ adj } B = A \text{ prox}[1, F] B$

$A \text{ near } B = A \text{ prox}[25, F] B$

$A \text{ before } B = A \text{ prox}[\infty, T] B$

なお、STARTS に無限大の記法はないが、ここでは理解を容易にするため採用した。

STARTS では、近接演算の被演算子をキーワードに限定している。そのため、 $A \text{ oadj } (B_1 \text{ or } B_2) \text{ oadj } C$ のようなパターンを表現することはできない。

3. CSE の近接演算

CSE では、Lycos の近接演算をすべて表現するため以下のような演算子を採用した。

A prox[m,n,X] B

これは、ある 2 語 A と B が最低 m 語、最大 n 語離れることを表し (m, n は非負数、 $m < n$)、X は語順の有無を表す。これを用いて他の演算子は以下のようにあらわされる。

$$A \text{ prox}[n,X] B = A \text{ prox}[1,n,X] B$$

$$A \text{ far} B = A \text{ prox}[25,\infty,F] B$$

通常の検索におけるキーワード検索の結果は、文書 ID とスコアの対の集合となるが、近接演算におけるキーワード検索では、文書 ID に対応するのは単一のスコアではなく、語の出現位置とスコアの集合である。

したがって、A prox[m,n,X] B は該当する各文書に対して以下の値を返す。

$$\{(\min(P_a, P_b), f(S_a, S_b)) \mid P_b \text{ in } Z\}$$

$$Z = [P_a + m, P_a + n] \quad \text{if } X = T,$$

$$Z = [P_a - n, P_a - m] \cup [P_a + m, P_a + n] \quad \text{if } X = F.$$

ここで、 P_a , P_b は A, B の位置、 S_a , S_b は A, B のスコア、 f はスコア関数である。なお、 $m = \infty$ のときは、A not B を意味する。

$$A \text{ and } B = A \text{ prox}[1,\infty,F] B$$

$$A \text{ not } B = A \text{ prox}[\infty,\infty,F] B$$

また、近接演算子の or と組み合わせて A oadj (B_1 or B_2) oadj C のようなパターンも表現できる。今後、繰り返しが表現できれば正規表現として使うことができる。

4. 評価

検索時間は、文書数 n と語数 m に対して $O(nm)$ であった (表 1 参照)。また、インデックスサイズを図 1 に示す。いずれも 1000 程度の文書数では問題とならない。CSE ではインデックスを分割し、並列・分散検索できるため、それほど性能を落とさず、精度を向上させることができる。

表 1. 検索時間

文書数*1	100	1000	10000
検索時間	0.17s	1.6s	17s
文中語数*2	10	100	1000
検索時間	0.26s	1.8s	16s

*1 文中語数は 1000、*2 文書数は 10000 で固定

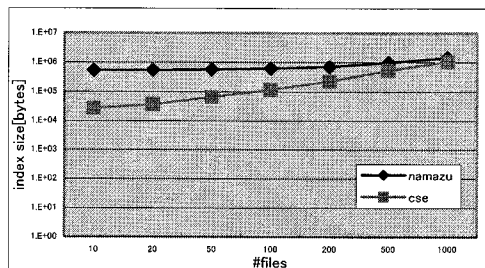


図 1 インデックスサイズ

5. まとめ

既存近接演算のほとんどを表現可能な近接演算子を提案し、論理演算と自由に組み合わせることのできる近接演算を実現した。分散型サーチエンジンでは近接演算の問題点を回避できるため実用の可能性がある。

参考文献

- [1] N. Sato, M. Uehara, Y. Sakai, H. Mori, "Fresh Information Retrieval in Cooperative Search Engine", In Proceedings of SNPD'01, pp.104-111, (2001.8.20)
- [2] Namazu, <http://www.namazu.org/>
- [3] Lycos, <http://www.lycos.co.jp/search/help.html>
- [4] Luis Gravano, et al.: STARTS - Stanford Protocol Proposal for Internet Retrieval and Search -, Stanford University, 1997