

係り受け情報を用いた全文検索とその評価

新美和彦, 兵藤安昭, 池田尚志

岐阜大学工学部

〒501-11 岐阜県岐阜市柳戸 1-1

Tel: 058-293-2740

Fax: 058-293-2740

E-Mail: {kazuhiko,hyodo,ikeda}@ikd.info.gifu-u.ac.jp

概要

大量の電子化文書に容易にアクセスできる環境が整ってくるにつれて、その中からユーザが必要とする情報のみを正確に検索する技術がますます必要不可欠となってきている。従来の検索手法では、文書中に単語が出現するか否かに関するブール演算による絞り込みが主として用いられてきた。そのほか、単語間の関係による絞り込みとしては近接演算が用いられてきたが、近接演算では単語間の正確な関係を表現することは出来ない。本論文では、単語間の係り受け情報を用い、高精度な絞り込みが可能な全文検索システムについて述べる。特許データを対象とした検索実験で、係り受け関係を用いた検索精度は適合率92.11%, 再現率96.01%を示し、近接関係を用いた検索と比べ良好な結果を得た。また、インデックス容量の増加も27%程度に抑えることが出来た。

キーワード

全文検索, 係り受け情報, 骨格構造解析, 近接演算

Full-text retrieval using dependency structure and its evaluation

Abstract

Today we can easily access a lot of large scale electronic documents, and with these advance the eager wish for high precision text retrieval is increasing. In this paper we propose, to improve the precision, the full-text retrieval using dependency relation between words instead of proximity operation. The proximity relation has been used as a substitution for syntactic relation so far, because the syntactic analysis is still a difficult task for a computer. We apply our method of skeletal syntactic analysis for Japanese to full-text retrieval and evaluate the index size, response time, accuracy of retrieval and others verifying usefulness of this method.

Keyword

Full Text Retrieval, Dependency Structure, Skeletal Syntactic Analysis, Proximity Operation

1. はじめに

大量の電子化文書に容易にアクセスできる環境が整ってくるにつれて、その中からユーザが必要とする情報のみを正確に検索する技術がますます必要不可欠となってきた。従来の検索手法では、文書中に単語が出現するか否かに関するブール演算による絞り込みが主として用いられてきた。そのほか、単語間の関係による絞り込みとしては近接演算が用いられてきたが、近接演算では単語間の正確な関係を表現することは出来ない。

検索精度向上のためには言語情報の活用が効果的である。[兵藤 96] は構文解析情報を活用した翻訳支援のための類似用例検索について述べている。本論文では、『ある単語がある単語に係る』という係り受け情報を利用した高精度な全文検索システムについて述べる。

係り受け情報は形態素解析及び構文解析を施すことにより作成されるが、長文に対する安定した構文解析はまだ困難である。本システムにおいては表層的情報のみを用いて行う骨格構造解析法 [兵藤 95] によって、文書データベースの係り受け解析を行った。

特許データを対象とした検索実験で、係り受け関係を用いた検索精度は適合率92.11%、再現率96.01%を示し、近接関係を用いた検索と比べ良好な結果を得た。また、インデックス容量の増加も27%程度に抑えることが出来た。

2. 全文検索システム

2.1 システム概要

本システムは、図1に示すように文書データベース、係り受け解析部、インデックス部、照合部、インタフェース部から構成される。インデックス部は、1次記憶上の単語エン트리と、2次記憶上の単語出現情報及び係り先情報とから成る。これらのインデックスは文書データベースに対する係り受け解析処理の結果から作成される。照合は2段階に分けて行われる。ユーザが、係り受け関係を含む検索パターンを入力すると、まず始めに、1次記憶上の単語エント리를検索し、単語が出現する文を抽出する。次に検索されたすべての文を対象として、検索パターンと係り受け構造が一致するか否かの照合を行い、検索結果をインタフェース上に表示する。検索システムはサーバ上にあり、ユーザはWeb上のインタフェースを通して検索する。

2.2 係り受け解析部

対象とする文書には、形態素解析および係り受け解析を施す。係り受け解析には骨格構造解析 [兵藤 95] を用いた。骨格構造解析とは必ずしも完全な係り受けの構造を求めるものではなく、並列構造の解析など意味に立ち入らなければ解析できない部分は曖昧なブロックとしてそのまま残し、文全体の構造を把握しようとするものである。解析例を図2に示す。

2.3 インデックス部

インデックスは、1次記憶上の単語エン트리と、2次記憶上の単語出現情報及び係り先情報から成る (図3参照)。

単語エント리는パトリシア構造を用いて構築している。現在のところ、数字・記号を除くすべての自立語を登録している。

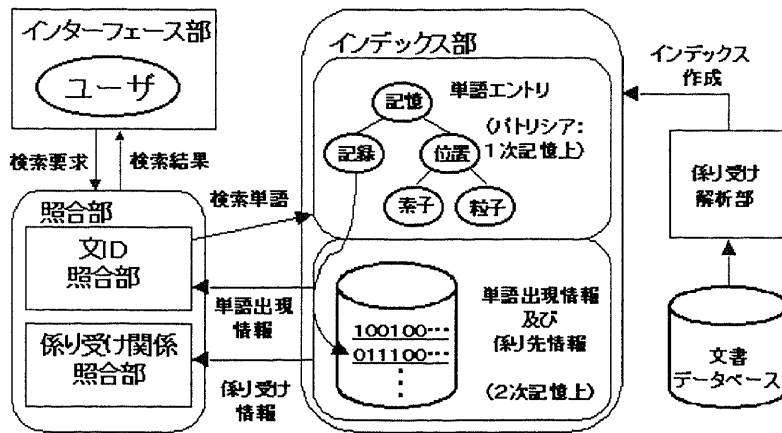


図1：システムの概要

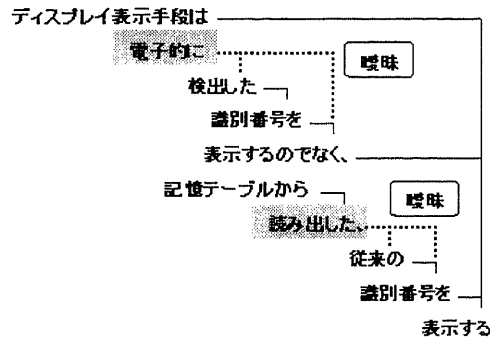


図2：解析例

単語出現情報は文IDと、文内での単語位置情報から構成されている。文IDは文書番号、文書内の項目番号、項目内での文番号から成り、[泓田97]の手法による階層化ビットベクトル用いて実装した。また、単語位置情報と係り先情報は、それぞれ文IDごとに単語番号列として登録する。単語番号は文内でのインデックス対象語を一意に表した番号である。

2.4 照合部

検索は2段階に分けて行う。ユーザが、係り受け関係を含む検索パターンを入力すると、まず始めに、1次記憶上の単語エンタリを検索し、2次記憶上の文IDおよび単語位置情報と係り先情報を読み込む。そして、読み込まれた文IDベクトル(階層化ベクトル)の論理積を実行し、指定した単語が出現する文を検索する(1次検索)。次に検索されたすべての文を対象として、検索パターンと係り受け構造が一致するか否かの照合を行う(2次検索)。係り受け構造の照合は、単語位置情報と係り先情報をビットベクトルに展開し論理積を実行することで行う。これにより、1つの単語が複数の位置に出現している場合や、係り先が特定出来ず複数の解析結果が得られている場合でも高速に照合が可能である。(図4参照)

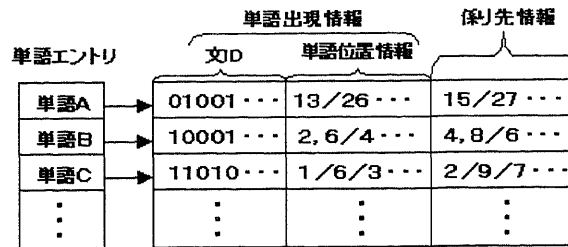


図3: インデックス

2. 5 インタフェース部

インタフェース部はJ A V Aを用いて構築しており、W e b上で使用できる。検索単語の入力、係り受け関係の指定は、G U I上で容易に行うことができる。図5にインタフェース画面を示す。

3. 検索実験

係り受け情報を用いた検索の有効性を実証するため、ブール検索、近接関係を用いた検索との比較を行った。検索対象には表1に示す公開特許公報の一部を用い、サーバには、S P A R C S t a t i o n 2 0 (C P U : S u p e r S P A R C Ⅱ, 7 5 M H z, メモリ: 6 4 M B y t e, O S : S u n O S 4. 1. 4)を使用した。

3. 1 インデックス容量、検索実行時間

係り受け検索と近接関係を用いた検索及びブール検索との間で検索実行速度、インデックス容量の比較を行った。結果を表2に示す。

検索速度は、係り受け検索で1件あたり約21.45(MS)を要し、近接関係を用いた検索より高速で、ブール検索と比べても検索実行時間の増加を約12%に抑えることができた。インデックス容量はブール検索で使用するインデックスの約1.67倍、近接関係を用いた検索で使用するインデックスの約1.27倍の増加となった。

3. 2 係り受け検索の絞り込み精度

本検索システムで用いた骨格構造解析では、意味情報を利用しないと正確に係り受け解析できない部分は、曖昧なまま係り先を特定しないため、検索の際に正しく絞り込みができない可能性がある。図6に誤った絞り込みを行った例を示す。また、誤って解析した場合には、検索洩れを生じる場合もある。係り受け検索の絞り込み精度を評価するため、「[メモリ]が[記憶する]に係る文」のような10件の検索要求に対し適合率と再現率を求めたところ、適合率92.11%、再現率96.01%という結果を得た。

3. 3 近接関係を用いた検索との絞り込み精度比較

3.2で述べた同じ検索要求10件に対し、近接関係を用いた検索との絞り込み精度比較を行った。表3に係り受け検索と近接距離を1~5まで変化させた時の適合率・再現率を示す。

「単語A」→「単語B」を検索

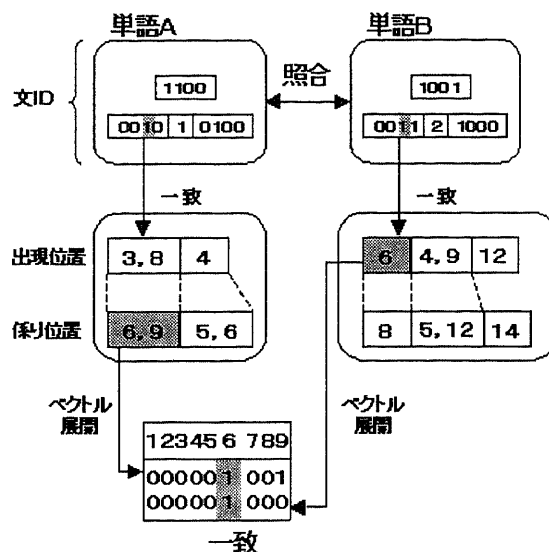


図4：係り受け関係の照合

近接関係を用いた検索では近接距離が短い時は適合率が良いが再現率が低い。図7の例では、近接距離1での検索はAしか検索出来ない。又、近接距離を3にすると、実際には係り受け関係がないCまで検索してしまう。さらに、DやEを検索するため近接距離を長くすると、適合率が低下してしまう。このことから係り受け検索が有用であることが分かる。

4. おわりに

本論文では、係り受け情報を用いることによる高精度な全文検索システムについて述べた。特許公報に対する、係り受け情報の検索精度については、適合率92.11%、再現率96.01%であり、近接演算を用いた検索より良好な結果を得た。又、インデックス容量の増加は近接関係を用いた検索でのインデックスの約27%に収まった。

参考文献

- [兵藤95] 兵藤安昭, 池田尚志: 表層的情報とN近傍ブロック化手法による日本語長文の骨格構造解析, 情報処理学会論文誌, Vol. 36, No. 9, pp 2091-2101 (1995)
- [兵藤96] 兵藤安昭, 河田実成, 應江黔, 池田尚志: 構文つきコーパスの作成と類似用例検索システムへの応用, 自然言語処理, Vol 3, No. 2, pp 73-88 (1996)
- [泓田97] 泓田 正雄, 溝渕 昭二, 獅々堀 正幹, 青江 順一: 大規模文書データに対する用例文の効率的検索アルゴリズム, 情報処理学会論文誌, Vol. 38, No. 10, pp 2004-2013 (1997)

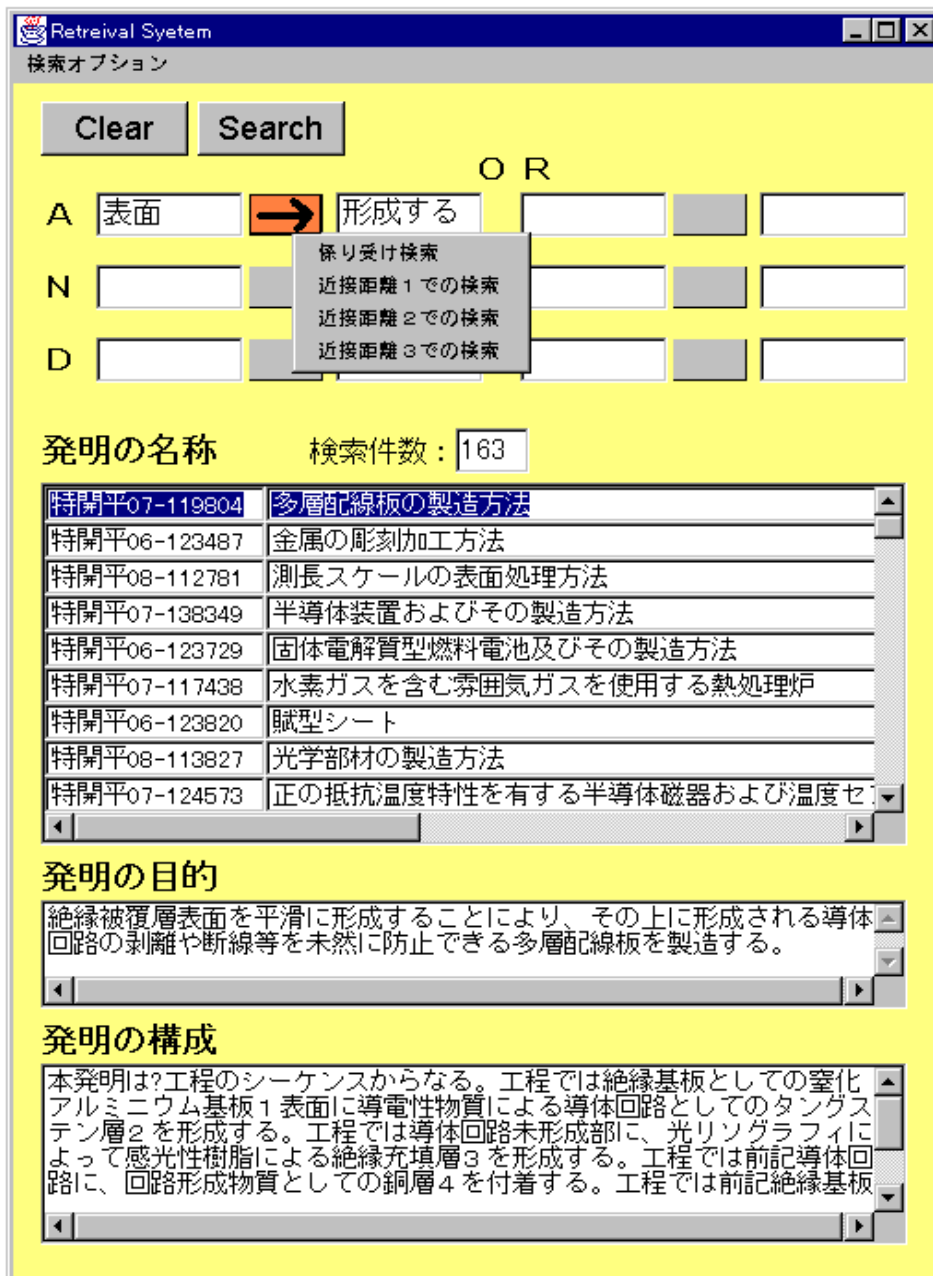


図5：インターフェース画面

対象	公開特許・実用新案中の2項目(「目的」、「構成」)
総特許数	19,870件 (総バイト数：9,021 Kbyte)
総文数	51,070文 (1特許当たりの文数：2.57文)
総単語数	2,034,265個 (1文当たりの単語数：39.83個)
総単語エントリ数	867,496個
単語エントリの出現回数	8.75回/エントリ (MAX：「こと」, 15,789回)

表1：実験に使用した特許文書のデータ

検索名	検索実行速度	インデックス容量
ブール検索	19.10(MS)	3,937(Kbyte)
近接関係を用いた検索	22.56(MS)	5,182(Kbyte)
係り受け検索	21.45(MS)	6,594(Kbyte)

表2：検索実行時間とインデックス容量の比較

「指示軸」→「備える」の検索

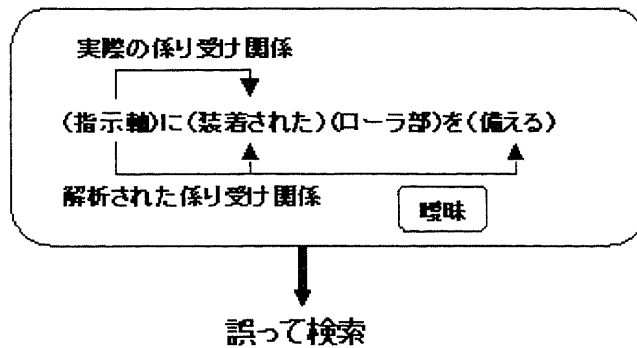


図6：検索誤りの例

検索名	適合率	再現率
係り受け検索	92.11(%)	96.01(%)
近接距離1での検索	96.12(%)	74.25(%)
近接距離2での検索	84.06(%)	81.55(%)
近接距離3での検索	75.92(%)	87.28(%)
近接距離4での検索	69.94(%)	90.52(%)
近接距離5での検索	65.56(%)	93.51(%)

表3：適合率・再現率の比較

検索要求：「表面→形成する」

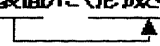
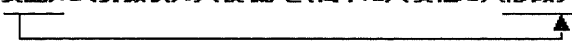
	検索例文	近接距離 ₂	近接距離 ₃	係り受け
A	(基材 表面)に(形成される)(皮膜)である 	○	○	○
B	(ガラス 表面)に(多数)の(溝状 凹部)が(形成された)(ガラス板) 	×	○	○
C	(ZnSe 基盤)の(表面)を(酸化し)て(静電潜像)が(形成される) 	○	×	○
D	(化粧材 表面)に(釘頭状)の(模様)を(簡単に)(安価に)(形成する) 	×	×	○
E	(表面)には、(接着剤層)を(構成する)に(先立)て(印刷 インキ)で (木目 模様)を(印刷)した(印刷層)が(形成され)ている 	×	×	○

図7：検索例