

1R-11

入力文字情報を最大限に活用した 文字列検索アルゴリズムの提案

大曾根 匡、佐藤 創

(専修大学経営学部情報管理学科)

1. はじめに

近年、文献情報等の文書情報のDB化が急速に進められているのに伴い、文書情報処理の高速化のニーズが極めて高まっている。このような状況の中で、文書情報処理のうち最も基本的かつ高負荷な処理の一つであるストリング・サーチ処理の高速化は重要な課題である。その高速化を目的としたアルゴリズムとして、KMP法やAC法、BM法などが著名である。また、AC法とBM法をハイブリッドしたアルゴリズムもいくつか提案されている[1][2]。しかし、これらのアルゴリズムは、入力した文字の情報を十分に活用しきっていないとはいえないものであった。そこで、本稿では、入力文字情報を最大限に活用した文字列検索アルゴリズムを提案する。また、性能実験により、その有効性についても検証する。

2. 提案アルゴリズム

高速なアルゴリズムとして著名なBM法は、パターン末尾の文字から照合を始めることにより、高速化を図っている。例えば、パターンが「ABCD」の場合、「D」から照合を始める。そのために「D」に対応する文字をテキストから入力する。もしこの入力文字が「A」だったらパターンを右へ3文字分シフトし、再び、シフトしたパターン末尾の文字に対応する文字を入力する。しかし、BM法では、パターンを右へシフトした際、それまでに入力した文字の情報を全て忘れてしまうようになっている。

そこで、提案するアルゴリズムでは、パターンを右へシフトした際にもそれまでに入力した文字情報を覚えておくようにした。前述の例では、現在「A_ _ *」という照合状況であるとして覚えておくようにする。ここで、_ は未照合文字、*は次の照合文字を表す。すなわち、この照合状況は、パターンの先頭文字の照合は済んでおり、つぎに4文字目の文字の照合を行うことを示している。

提案アルゴリズムでは、パターンから考えられる全ての照合状況をまず生成する。パターンが「ABCD」の場合の全ての照合状況を表1に示す。次に、[1][2]のアルゴリズムと同様に、状態 i のときに文字 c が入力されたら、次の状態が何になるかという情報が書かれた状態遷移テーブル $T[i, c]$ と、次に何文字先の

文字をテキストから入力すればよいかという情報の書かれたスキップテーブル $S[i, c]$ を作成する。表2と3に、パターンが「ABCD」の場合の状態遷移テーブルとスキップテーブルを示す。そして、パターンの検索は、この2つのテーブルの参照を繰り返すことにより行われる。その動作例を図1に示す。初期状態は0とし、テキストの4文字目から文字を入力し始める。この例の場合、提案法では、25文字のテキストに対し10文字の入力で検索を終了している。すなわち、平均スキップ幅は2.5文字である。一方、BM法では14文字の入力を要する。一般に、提案法は、原理的にBM法やKMPの性能を下回らないことに注意されたい。

3. 性能実験

アルゴリズムの性能は、平均スキップ幅によって表現できる。そこで、提案法とKMP法、BM法の平均スキップ幅を定量的に比較するために、各種の実験を行った。そのうちの顕著な差の認められた結果を以下に示す。図2は、パターンを「ABCD」とした場合、テキスト上の「C」の出現確率に対し平均スキップ幅がどのように変化するかを示した図である。実験では、アルファベットの文字種は32文字とし、長さ10000文字のテキストをランダムに生成し、そのテキストに対する平均スキップ幅を求めた。ここで、「C」以外の文字の出現確率は一律とした。この試行を100回行い、その平均を求めた。これより、「C」の出現確率の増加に伴い、提案法の性能とBM法の性能との差がだんだん大きくなっていくことがわかった。しかし、「C」以外の文字の出現確率が増加に対しては、BM法と提案法の性能にほとんど差が出ないという結果が得られた。図3は、パターンを「AAAA」とした場合の、テキスト上の「A」の出現確率の変化に対する平均スキップ幅である。これより、「A」の出現確率の増加に伴い、BM法の性能が著しく劣化するのに対し、提案法の性能はKMP法の性能を下回らないことが確認された。図4は、パターンを「10101」とした場合の、テキスト上の「1」の出現確率の変化に対する平均スキップ幅である。ここで、アルファベットの文字種は「1」と「0」の2文字とする。このように、特にアルファベットの文字種の少ない場合に、提案法が有効であることがわかる。

A String Searching Algorithm with full use of Information obtained in Text Scanning

Tadashi OHSONE and Hajime SATO

Department of Information Management, Senshu University

表1. 状態の定義

| 状態 | 照 合 状 況 |
|----|---------------|
| 0 | __ * _ _ _ |
| 1 | _ _ * D _ _ _ |
| 2 | _ * C D _ _ _ |
| 3 | * B C D _ _ _ |
| 4 | A B C D _ _ * |
| 5 | _ _ C * _ _ _ |
| 6 | _ B * _ _ _ _ |
| 7 | A _ * _ _ _ _ |
| 8 | _ B * D _ _ _ |
| 9 | A _ * D _ _ _ |
| 10 | A * C D _ _ _ |

* : 入力位置

表2. 状態遷移テーブル

| | A | B | C | D | # |
|----|---|---|----|---|---|
| 0 | 7 | 6 | 5 | 1 | |
| 1 | | | 2 | | |
| 2 | | 3 | | | |
| 3 | 4 | | | | |
| 4 | 7 | 6 | 5 | 1 | |
| 5 | 7 | | | 2 | |
| 6 | 7 | 6 | | 8 | |
| 7 | 7 | 6 | 5 | 9 | |
| 8 | | | 3 | | |
| 9 | | | 10 | | |
| 10 | | 4 | | | |

空白は状態0

表3. スキップテーブル

| | A | B | C | D | # |
|----|---|----|----|----|---|
| 0 | 3 | 2 | 1 | -1 | 4 |
| 1 | 5 | 5 | -1 | 5 | 5 |
| 2 | 6 | -1 | 6 | 6 | 6 |
| 3 | 7 | 7 | 7 | 7 | 7 |
| 4 | 3 | 2 | 1 | -1 | 4 |
| 5 | 3 | 4 | 4 | -2 | 4 |
| 6 | 3 | 2 | 4 | -1 | 4 |
| 7 | 3 | 2 | 1 | -1 | 4 |
| 8 | 5 | 5 | -2 | 5 | 5 |
| 9 | 5 | 5 | -1 | 5 | 5 |
| 10 | 6 | 6 | 6 | 6 | 6 |

: その他の文字

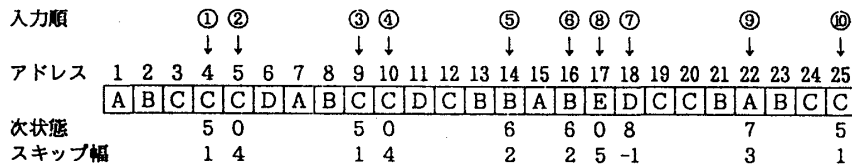


図1. 提案法の動作例

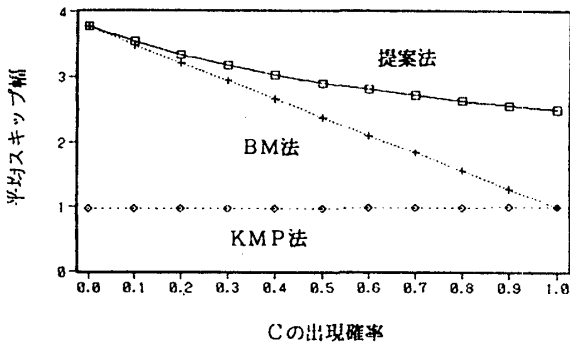


図2. パターン「ABCD」に対する平均スキップ幅

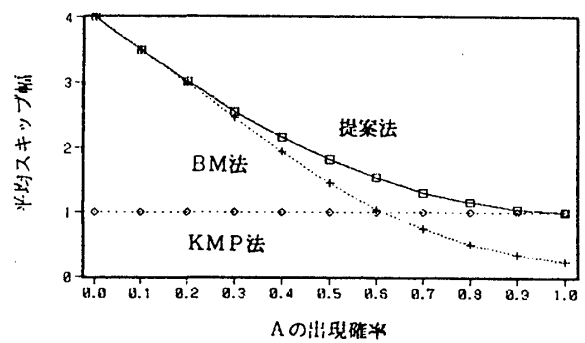


図3. パターン「AAAA」に対する平均スキップ幅

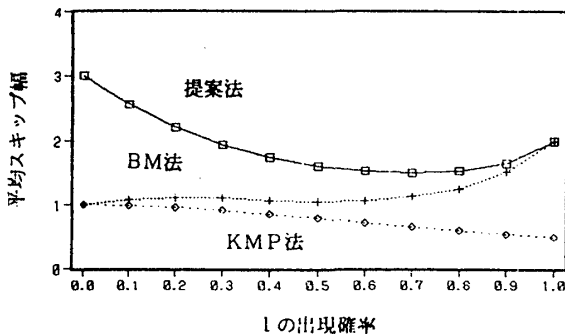


図4. パターン「10101」に対する平均スキップ幅

4. まとめ

入力した文字情報を最大限に活用した文字列検索アルゴリズムを提案した。また、性能実験によりKMP法やBM法より有効であることを定量的に検証した。特に、2種類のアルファベットからなるテキストにおける文字列検索に有効であると考えられる。

参考文献

- [1] 大曾根 他, "高速ストリング・サーチ・アルゴリズムの提案," 情報処理学会第34回全国大会論文集, pp. 463-464 (1987).
- [2] 大曾根, 佐藤, "文字の出現頻度を考慮した文字列検索アルゴリズムの提案," 情報処理学会第45回全国大会論文集(1), pp. 67-68 (1992).