

パターンストリームにおけるデータ検索

山口信[†] 島田 諭[†] 三浦 孝夫[‡]

法政大学[†] 法政大学[†] 法政大学[‡]

1. 前書

情報検索の分野では従来様々な手法によりデータから特定のパターンを発見する手法が提案されている。SuffixArray や SuffixTree は予めテキストデータに索引付けをしておくことにより、高速な検索を実現するアルゴリズムである。しかしながら大量のストリームデータにおいて、このようなデータ構造を構築、維持することは、空間的速度的に非常に高価であり、現実的でない。一方、KMP 法や BM 法といったそのアルゴリズムによる手法は、構造のないデータ検索に適しているといえる。またストリームパターンに対応する方法はあまり知られていない。しかし検索パターンが変動する中で検索を行わねばならない状況は十分に考えられる。例えば、株価は常に変動しているが、その変化の特徴を逐次過去のデータから検索することは、株主にとって有用なことである。本研究手法では KMP 法において、検索パターンが連続的に増加すると仮定し、増加分による検索を、差分情報を用いて高速に文字列検索するアルゴリズムを提案する。

2. KMP 法と BM 法

パターン中の部分文字列の接頭辞と接尾辞の最長一致長から瞬時に SKIP 量を計算し、冗長な検索を省略することができるアルゴリズムが KMP 法または BM 法である。skip 量をまとめた skip 表の構築時間はパターンの線形時間であり、一般にパターンはデータに比べ小さいので、この構築時間は大きくなりにくい。KMP 法では文字列検索はパターンの先頭から行い、BM 法は一般的に文字列検索をパターンの末尾から行う。末尾から比較し不一致点でスキップすることによりテキスト上で一度も比較しない文字が生じ、計算量を抑えることが出来る。ただし BM 法は日本語のように全角文字を用いる文化圏では空間的に高価な可能性がある。

3. パターンストリーム

KMP 法は検索時の一致長内の接頭辞と接尾辞の最長一致がパターンの幅に関与する。そこでそ

の接尾辞と接頭辞の最長一致を MPL として予めテーブルに保持しておけばそれが skip 表の役割を果たすことになる。

MPL の定義

$$\text{pattern}[0, 1, 2, \dots, m] \text{ のとき}$$

$$\text{MPL}(m) = \text{MAX} \{k \mid [\text{pattern}[0\dots k] = \text{pattern}[(m-k), \dots, m]], 0 < k < m\}$$

接頭辞と接尾辞の一部が重なっている場合、そのパターンは、接頭辞の重複を除く接頭辞部分の繰り返し文字列になっている。ただし最後の繰り返し部分は完全な接頭辞になっているとは限らない。また接頭辞と接尾辞の一部が重なっていない場合は以下のようなになる、 $\text{pattern}[m], K = \text{MPL}(m)$ とするとき、

$$k < \frac{m}{2}$$

特に $k=0$ の時は $\text{pattern}(m) \neq \text{pattern}(1)$

かつ $\text{pattern}(m-i, \dots, m) \neq \text{pattern}(0, \dots, i)$ である。

次に、差分的に求める MPL 長を $\text{MPL}(m+1)$ とすると、 $\text{MPL}(m+1)$ は $\text{MPL}(m)$ より 1 文字長くなるか、 $\text{MPL}(m)$ より短くなることを示す。α を $\text{MPL}(m)$ で一致する文字列、z を接頭辞 α の次の文字、x を $\text{MPL}(m+1)$ で追加された文字とする。

α	z	⋯	α	x
---	---	---	---	---

まず $x=z$ のときは、 $\alpha z = \alpha x$ は明らかなので $\text{MPL}(m+1) = \text{MPL}(m) + 1$ はひとまず成り立つが、実はこれ以上長い $\text{MPL}(m+1)$ は存在しない。仮に存在すると仮定すると、より長い $\text{MPL}(m+1) = \beta$ が存在して、下のようになる。

α	x	⋯	x	⋯	α	x
β				⋯	β	

接尾辞 β の末尾は x ゆえ、接頭辞 β は

$\beta = \alpha x \dots x$ の並びになっている。ここで $\alpha x \dots = \beta'$ とおいて接頭辞 $\beta = \beta' x$ としておく。一方接尾辞 β は $\beta = \dots \alpha x = \beta' x$ の並びになっている。つまり “ $\dots \alpha$ ” は β' というより長い $\text{MPL}(m)$ を保持していたことになる。この矛盾から $\text{MPL}(m+1) = \text{MPL}(m) + 1$ が成り立つ。次に $x \neq z$ のときは $\text{MPL}(m+1) \leq \text{MPL}(m)$ となる。仮に $\text{MPL}(m+1) >$

MPL(m)と仮定すると、MPL(m+1)の接頭辞βは、x ≠ z であるから、この長さは MPL(m)+1 以上になる。

α	z	...	x	...	α	x
β				...	β	

接尾辞βに着目すると $\beta = \dots \alpha x = \beta'x$ となっている。ここで接頭辞βは $\beta = \alpha z \dots x$ であるから $\alpha z \dots x = \beta'x$ これでは、MPL(m)はより長いものが既に存在していたことになり矛盾する。よって $MPL(m+1) \leq MPL(m)$ となる。

4. 提案手法

直前の MPL(m)から MPL(m+1)を差分的に構築する。この時第3章で述べたように追加文字が pattern(MPL(m)+1)と一致した場合は $MPL(m+1) = MPL(m)+1$ となるが、そうでない場合 MPL(m+1)は MPL(m)よりも短い範囲から探せば良いことがわかっている。その際のテーブルは、追加された文字 pattern[i]の直前の出現位置へのポインタが格納されている。こうすることにより比較回数が $\frac{1}{|\Sigma|}$ になることが期待できる。

ここで |Σ| は半角文字の総数である。このリンク表を skip2 と名付け、リンクした位置から先頭まで一致する MPL が MPL(m+1)となる。MPL(M+1)を求めるこの手法の計算量は次の通りである。

$$O(1 + |\alpha| + |\alpha| \times \frac{|\alpha|}{|\Sigma|}) = O(\frac{|\alpha|^2}{|\Sigma|})$$

ここでαは MPL(m)のことである。MPL(m+1)求める際、改めて KMP 法や BM 法を用いることでαの線形時間にするには可能だが、新たに skip 表を作成するコストと、実際にはαは非常に小さいということ考慮し、今回は2乗のオーダーのままにしてある。

5. 実験

実験で用いるコーパスはロイター誌から1996年の記事2586種を抽出する。テキストデータはこの2586種類の記事からランダムに2586回抽出して結合したものとした。検索パターンは2568種類の記事から一つを選択したものである。

実験結果としてパターンの違いにおける性能評価結果を表1に示す。skip 表作成時比較コストは MPL(m+1)を構築する際の文字列比較回数と、表への挿入回数で、skip 表作成時挿入コストは skip 表への値の挿入回数を示し、差分

法の場合は2つの表への挿入回数の合計となっている。skip 表作成コストは上記2つの合計を表し、それと検索時の比較回数を表す検索コストとの合計をとったものが総コストである。

表1 パターンお違いにおける性能差

	差分平均	従来平均
skip 表作成時比較コスト	2150.4	2614346
skip 表作成時挿入コスト	4300.4	2611465
skip 表作成コスト	6450.8	5225811
検索コスト	2490802	2489882
総コスト	2497253	7715693

全ての実験結果で優秀な結果となった。KMP 法の skip 表は動的に作成すればほとんどコストが掛からず平均800倍の速度となった。しかしながら、データ検索のコストを含めると、3倍程度の成果にとどまった。KMP 法では検索速度に課題が残る。また今回はパターンにロイターコーパスを用いたため、英語の文章では大文字と小文字の区別が存在してしまい、MPLは2以上になることは極めて稀な小規模なものとなった。そのため差分構築において2回以上の比較を行うことはまず無いと言える。同じ文字列が繰り返してこないパターンでは差分構築のコストを安く済ませられる反面、検索の効率が上がらず、また今回提案したリンク表である skip2 を十分に活用することができない。

6. 結論

パターンの先頭から MPL をとる KMP 法は、差分構築においてとても優秀な結果を出した。BM 法を用いればより高速な検索が期待できるが、パターン増加に対する差分構築は複雑であり、構築と検索の双方の観点から比較する必要がある。今後は異なる特性を示すコーパスを用いることにより、一層の有用性を示せると考える。

文献

[1] 石畑 清：アルゴリズムとデータ構造

Information retrieval in the pattern stream
 Makoto Yamaguchi · Hosei
 Satoshi Shimada · Hosei
 Miura Takao · Hosei