

柔軟な文書検索のためのコンパクトなデータ構造 Space-Efficient Data Structures for Flexible Document Retrieval

定兼 邦彦*

Kunihiko Sadakane

1 はじめに

本論文は柔軟な文書検索のためのコンパクトなデータ構造を提案する。現在広く用いられている文書検索システムでは、 tf^*idf スコア [9] に従い文書をランク付けすることで検索の精度を高めている。このスコアは単語 p と文書 d に対するスコア $tf(p, d)$ と $idf(p)$ の積で定義される。前者は文書 d 中の単語 p の出現頻度、後者は総文書数を k 、 p を含む文書数を n_p とすると $\log \frac{k}{n_p}$ と定義される。このスコアを計算するためのデータ構造としては転置ファイル [2] が用いられている。これは各単語ごとにそれを含む文書の番号を列挙したものであり、単語の頻度も格納してあるため tf^*idf スコアを高速に計算できる。ただし任意の文字列に対して計算することができないという欠点がある。

文字列検索の分野では接尾辞木 ([4] 参照) が有名である。これは任意の文字列の出現箇所や頻度を高速に求められるデータ構造であるが、ある文字列 p を含む文書を高速に列挙すること (文書列挙問い合わせ) ができない。Muthukrishnan [6] は接尾辞木と区間最小値問い合わせのデータ構造を用いた最適時間アルゴリズムを考案したが、転置ファイルと比べてデータ構造のサイズが大きいという欠点がある。また、 tf^*idf スコアの計算もできない。

本論文では文書列挙問い合わせおよび任意の文字列に対する tf^*idf スコア計算のためのコンパクトなデータ構造を提案する。検索システム中の文書数を k 、文書の総長を n とすると、主結果は以下の通りである。

定理 1 任意の文字列 p に対し、それを含む q 個の文書を $O(|p| + q \log^\epsilon n)$ 時間 (ϵ は $0 < \epsilon \leq 1$ の定数) で列挙できる。

定理 2 任意の文字列 p と文書 d に対し、 $tf(p, d)$ を $O(|p|)$ 時間で計算できる。また、 p を含む q 個の文書全てに対する $tf(p, d)$ を $O(|p| + q \log^\epsilon n)$ 時間で計算できる。

定理 3 任意の文字列 p に対し、 $idf(p)$ を $O(|p|)$ 時間で計算できる。

これらの問い合わせを実現するデータ構造のサイズは $2|CSA| + 10n + o(n)$ ビットであり、 $O(n)$ 時間で構成できる。ここで、 $|CSA|$ はデータベース中の全ての文書に対する圧縮接尾辞配列 [3] のサイズであり、通常は文書サイズ以下になる [7]。また、文書サイズは通常 $8n$ ビットであるため、このデータ構造のサイズは文書サイズの約 3 倍である。既存手法で用いられている接尾辞木のサイズは文書サイズの 10 倍以上であり、また、Muthukrishnan のデータ構造のサイズはその約 2 倍程度であるため、本論文のデータ構造のサイズはそれらより非常に小さい。また、計算量は高々 $O(\log^\epsilon n)$ 倍になるだけである。このように、本研究のデータ構造は任意の文字列に対する tf^*idf スコアを高速に計算で

*東北大学大学院情報科学研究科 sada@dais.is.tohoku.ac.jp

きるため、転置ファイルなどのデータ構造よりも柔軟な検索を行える。

2 文書列挙問い合わせのためのデータ構造

文書を d_1, d_2, \dots, d_k とし、それを連結した文字列を T で表す。 T の接尾辞木は図 1 のようになる。葉の数字は根からその葉までのパス上の文字列が T 中のどの位置にあるかを表す。これは接尾辞配列 ([4] 参照) と呼ばれ、 SA で表す。配列 D の要素 $D[i] = j$ は接尾辞木で左から i 番目の葉が文書 d_j に含まれることを示す。配列 C の要素 $C[i]$ は $j < i$ かつ $D[j] = D[i]$ となる j のうち最大のものと定義される。そのような j が存在しない場合は $C[i] = -D[i]$ と定義する。

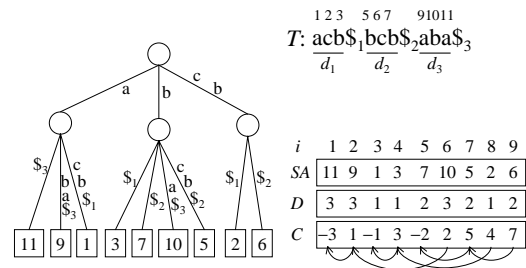


図 1: “acb\$1bcb\$2aba\$3” に対する接尾辞木と文書列挙問い合わせのためのデータ構造。

Muthukrishnan のアルゴリズム [6] はまず p に対応する接尾辞配列の範囲 $[l, r]$ を求め、次に配列 C の $C[l, r]$ の範囲で l 未満の数 $C[i]$ に対応する $D[i]$ の値を列挙する。本論文ではこれらのデータ構造のサイズを圧縮する。まず、配列 $D[i]$ の値は $SA[i]$ から計算できるため D は格納しない。また、 $SA[i]$ は圧縮接尾辞配列を用いて $O(\log^\epsilon n)$ 時間で計算できる。接尾辞配列の範囲 $[l, r]$ は圧縮接尾辞配列を用いて $O(|p|)$ 時間で求まる [8]。

配列 C に関しては、値自体は格納せず、値の大小関係を表す木を格納する (図 2)。配列 $C[1, n]$ の最小値が $C[x]$ のとき、木の根に $C[x]$ を格納し、左の部分木に $C[1, x-1]$ 、右の部分木に $C[x+1, n]$ を再帰的に格納する。このとき、ノードの通りがけ順と接尾辞配列での順序が一致するように、 $C[x]$ を左右の子の間に置く。配列の部分区間 $C[l, r]$ での最小値は、この木のノード間の最近共通祖先 (lowest common ancestor) を用いるデータ構造 [8] を用いて定数時間で求まる。データ構造のサイズは $4n + o(n)$ ビットである。

Muthukrishnan のアルゴリズムでは文書を重複して出力しないために C の値を用いていた。本研究では C の値自身は保存していないため、長さ k のビットベクトルを用い

