

2010 年度冬の LA シンポジウム [S10]

文法圧縮に基づいた圧縮データの自己索引構造化の提案

馬場 雅大* 丸山 史郎† 坂本 比呂志‡ 定兼 邦彦§ 山下 雅史¶

1 はじめに

データを圧縮保存している場合、文字列検索や部分復元といった操作を行う際に全体を復元するのは不利益が大きい。それを回避する手段として自己索引 (self-index) 構造が提案されている。

本稿で扱うのは、文脈自由文法 (CFG) に基づいた圧縮と構文木を圧縮した全二分木の巡回によって文字列検索を行うものである [4]。このデータ構造は長さ u の文字列 S から変換された CFG G が持つ異なる規則 (変数) の数を n 、 G の構文木の高さを h としたとき、 $1.5n \log n + n \log h + h \log n + o(n \log n)$ ビット領域で表し、長さ m のパターン P の出現回数を求める計算量は $O(m \log^2 n + occ_c(m \log n + h))$ である。CFG の最適解を n_* とすると $n = O(n_* \log u)$ 、 $h = O(\log u)$ である。 occ_c とは構文木中に現れるコアの出現回数を示す。コアとは P に含まれる部分文字列を十分に長く符号化した変数である。構文木で P が存在するところにはコアが存在するので検索時の必要条件とする。 P が長い場合 occ_c の数は少なくなり検索時間は小さくなる。

一方でこの手法は、文字列検索以外の操作をサポートしていない。自己索引構造で効率よくサポートすべき操作を以下に定義する。

- $count(S, P)$: S 中の P の出現回数
- $locate(S, P)$: S 中の P の位置
- $access(S, i)$: S 中の i 番目の文字

2 準備

以下ではデータ構造で用いる簡潔データ構造などについて説明する。

簡潔データ構造は、対象となるデータ構造をそのデータの情報理論的下限に漸近的に一致する一致する簡潔表現と、データに対して何らかの操作を効率よく行い、漸近的に小さい簡潔索引とで成り立つ。情報理論的下限とは、位数 L に対して $\lg L$ ビットである。なお \lg は底が 2 の対数を表している。長さ n のビット列であれば、 $L = 2^n$ 通りのパターンが考えられるので情報理論的下限は n ビットである。

2.1 Rank/select 演算

$|A| = \sigma$ のアルファベット A 上の長さ n の文字列 S を考える。ここで S に対する操作を以下のように定義する。

- $rank_c(S, i)$: $S[..i]$ 中の c の出現回数
- $select_c(S, i)$: S 中で先頭から i 番目の c の位置
- $access(S, i)$: $S[i]$

このような操作を効率的にサポートするデータ構造を rank/select 辞書という。

2.1.1 ビット列の簡潔データ構造

長さ n のビット列のサイズは 1 の数が少なければ圧縮できる。1 の数が m 個のビット配列での操作を定数時間で行う $B(n, m) = \lceil \lg \binom{n}{m} \rceil + O(n \lg \lg n / \lg n) = m \lg \frac{n}{m} + \Theta(n) + O(n \lg \lg n / \lg n)$ ビットのデータ構造がある [6]。これを完全索引付辞書 (FID) と呼ぶ。

*九州大学大学院システム情報科学府
†九州大学大学院システム情報科学府
‡九州工業大学大学院情報工学研究院
§国立情報学研究所
¶九州大学大学院システム情報科学府

2.1.2 文字列の簡潔データ構造

アルファベットサイズが2よりも大きい場合にはウェーブレット木 [2] を用いる。これはテキスト S を表現する深さ $\lg \sigma$ の二分木であり、各ノードはビット列を格納している。ある深さにある列の長さを合計するとテキスト長 n となる。各列に対して rank/select 辞書を追加したサイズは $n \lg \sigma(1+o(1))$ ビットとなる。rank_c などの操作の計算量は $O(\lg \sigma)$ である。

2.2 全二分木の簡潔データ構造

順序木の簡潔データ構造として BP 表現 [5] と DFUDS 表現 [1] が挙げられる。両者は全体としてバランスする開き括弧と閉じ括弧の括弧列で構成され、ノード数を n とすると長さは $2n$ である。

2.2.1 BP, DFUDS で共通するデータ構造

括弧列を用いて表現する BP, DFUDS では木の巡回を行うために、以下の操作を $o(n)$ ビットで定数時間で行う補助データ構造が提案されている [5]。対象となる括弧列を P とすると以下の通りである。

- $findopen(P, x)$: $P[x]$ にある開き括弧に対応する閉じ括弧の位置を返す。
- $findclose(P, x)$: $P[x]$ にある閉じ括弧に対応する開き括弧の位置を返す。
- $enclose(P, x)$: $P[x]$ にある括弧とそれに対応する括弧を囲う最小の括弧対の開き括弧の位置を返す。

これらと括弧列の rank/select 索引などにより木の巡回を実現する。

2.3 全二分木の簡潔データ構造

全二分木は情報理論的下限が $n - \Theta(\log n)$ ビットなので n ビットの別表現 [8] を用いる。

全二分木を前置順に巡回して内部ノードであれば開き括弧、葉ノードであれば閉じ括弧を順番に並べる。括弧全体をバランスするために先頭に開き括弧を配置する。この括弧列 F から以下のような操作を

定数時間で行うことができる。なお x は前置順に並べたノードを表現しているとする。

- $isleaf(x)$: x は葉であるか。
- $parent/sibling(x)$: x の親・弟
- $left/rightchild(x)$: x の長男・次男
- $childrank(x)$: x がその親の左から何番目の子か
- $leaf_rank/select$: 葉ノードの rank/select
- $inner_rank/select$: 内部ノードの rank/select

全二分木の表現は、パトリシアトライや DFUDS 圧縮 [3] などが提案されている。主要項では同サイズであるが、パトリシアトライは右の子への対応に結局索引を追加する他、[3] は複雑であり、簡潔索引そのものは $2n$ ビットの括弧列に対するものなので括弧列そのものを n ビットとしている本稿の表現の方が有利である。

3 文法圧縮に基づく圧縮索引の自己索引構造化

本節では、edit sensitive parsing に基づいてテキスト S から構築された構文木 T に対する表現と操作について簡単に説明する。構文木とは $X \rightarrow AB$ といった生成規則を持つ CFG と等価な順序木のことを指す。このとき X から AB を求めるためのデータ構造を辞書 D とし、逆に AB から X を求める逆引き辞書を D^R と表記する。

パターン P を検索する際に、 P の S における全ての出現に対して、 P の部分文字列を十分に長く符号化した極大な変数 (規則) をコアとする。 T 内に構文木が存在し、かつ P がある程度長ければ、 S から作られた D^R を用いて P をパーシングすることでコアを決定できる。

文字列を列挙して検索する際には、構文木中でパターンのコアを列挙していき、各コアの周辺で照合を行う。すなわち $P = X_1 X_2 \dots X_k$ といった変数列において、そのコアを X_t すると構文木中での X_t にラベル付けされたノード v を列挙した上で、各々について v に隣接する部分木の中に他の $X_1 X_2 \dots X_k$ が

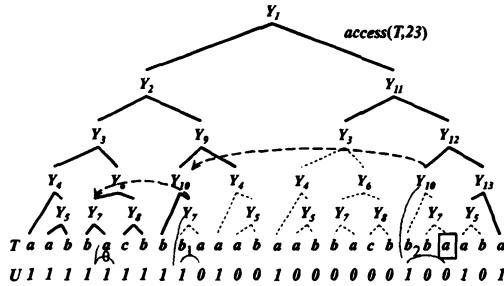


図 2: 配列 U と $access$ 操作の概要

1. 以下を構文木での葉，すなわち変数ではなくアルファベットを示す葉に至るまで繰り返す。
2. $S[i]$ の祖先にあたる葉の左からの順番を $k = rank_1(U, i)$ とする。
3. $j = select_1(U, k)$, $s := i - j + 1$
4. k について $lmost_occ$ を参照. 参照先を根とする部分木の s 文字目を求めるため再帰操作へ.

繰り返しは高々 h 回であり，左分木が明示的ではない葉ノードは $lmost_occ$ の参照を行うので計算量は $O(h \lg n)$ である。 $S[i]$ を求めた上で，そこから右側に展開していけば部分文字列復元となる。このとき PT の各葉は配列 L で隣接関係にある。葉は 1 文字以上カバーしているので復元する長さを l とすると，辿る枝の本数は $O(l)$ 本である。計算量は $O((h+l) \lg n)$ となる。

$locate(S, P)$ 操作: $locate$ は $count$ に付随して行う。

1. パタン圧縮時にコアがパタン先頭から j 文字離れていることを確認
 2. ヒット時にコアが示す部分文字列先頭位置 k
 3. $k - j$ を返す。
1. についてはパタン圧縮に並行して行うことができる。 2. についても同様に $count$ での照合作業に並行して行うことができるので， $locate$ の計算量は本来の $count$ 操作の計算量を超えることはない。

定理 1 サイズ $1.5n \lg n + B(u, n) + o(n \lg n + u)$ ビットのデータ構造を用いて $access$ 操作は $O(h \lg n)$ 時間， $locate$ 操作にかかる時間は $count$ 操作にかかる時間を増やすことはない。

4 おわりに

本稿では，文法圧縮に基づいた索引構造に対して，新たに索引を加えることで検索だけではなく検索位置の特定や部分復元ができることを示した。部分復元の計算量についてはウェーブレット木に変わるデータ構造を利用することで改善できる可能性がある。

参考文献

- [1] D. Benoit, E. D. Demaine, J. I. Munro, R. Raman, V. Raman, S. S. Rao : “Representing trees of higher degree.” *Algorithmica* 43(4), 275–292, 2005.
- [2] R. Grossi, A. Gupta, and J.S. Vitter. “High-order entropy-compressed text indexes.” In *SODA04*, pages 636.645, 2004.
- [3] J. Jansson, K. Sadakane, and Sung, W.-K.: “Ultra-succinct representation of ordered trees.” In *SODA (2007)*, N. Bansal, K. Pruhs, and C. Stein, Eds., SIAM, pages. 575–584.
- [4] S. Maruyama, H. Sakamoto, M. Baba, H. Ono, K. Sadakane, M. Yamashita. : “Searching Long Patterns from Grammar-Based Compression.” Submitting.
- [5] J. I. Munro, V. Raman : “Succinct Representation of Balanced Parentheses and Static Trees.” *SIAM Journal on Computing*, 31(3):762–776, 2001.
- [6] R. Raman, V. Raman, S. S. Rao : “Succinct Indexable Dictionaries with Applications to Encoding k -ary Trees and Multisets.” In *Proc. ACM-SIAM SODA*, pages 233–242, 2002.
- [7] H. Sakamoto, S. Maruyama, T. Kida and S. Shimozone : “A space-saving approximation algorithms for grammar-based compression.” *IEICE Trans. on Information and Systems*, E92-D(2):158–165, 2009.
- [8] 馬場雅大, 小野廣隆, 定兼邦彦, 山下雅史. “全二分木の簡潔な表現.” 情報処理学会研究報告, Vol.2010-AL-129 No.1, pages 1-8, 2010.