

接尾辞配列とその応用

nPk

Abstract-

Keywords: suffix array, in-place

1 記法

文字列は S とかき、部分文字列は Rust の区間記法にしたがう。つまり、表 1 の通りにする。接尾辞配列 SA についても同じようにする。

| | |
|---------------------|--------|
| S, T, \dots | 文字列 |
| $ S $ | 文字列の長さ |
| $S[i]$ | 文字 |
| $S[i..j], S[i..=j]$ | 部分文字列 |
| $S[i..]$ | 接尾辞 |

2 接尾辞配列

文字列 S の接尾辞 $S[i..]$ は頭文字の添え字 i で指定できる。接尾辞（に対応する添え字）を辞書順に並べたものを接尾辞配列という。接尾字配列上で二分探索を行うと、任意の文字列 T の全部の出現個所を $O(|T| \log |S|)$ で求めることができる。検索文字列についてオンラインに処理できる点で接尾辞配列は Z アルゴリズムや KMP 法よりも優れている。

2.1 文字

Rust 言語において文字は $T: \text{Ord}$ と抽象化できる。接尾辞配列においてすべての接尾辞の頭文字がソートされていることを考えれば、計算量は最悪 $\Omega(|S| \log |S|)$ であることが分かる。実際に、 $O(1)$ の作業メモリを使って、文字列 $\&[T]$ の接尾辞配列を $O(|S| \log |S|)$ で計算するアルゴリズムが存在する。

文字列を座標圧縮することを考える。このとき、 $T = \text{usize}$ である。後述するように、文字列 $\&[\text{usize}]$ の接尾辞配列は $\Theta(|S|)$ で求めることができる。座標圧縮の計算量は $\Theta(|S| \log |S|)$ であるから、これによって最悪計算量が悪化することはない。以下では、文字は usize であるとする。¹

¹より小さな数値型を用いても良いが、アルゴリズムが複雑になってしまう。本稿で紹介する実装では $\text{Vec}\langle T \rangle$ の最大容量が $\text{isize}::\text{MAX}$ であることを利用した最適化を施しており、オリジナルと比べてアルゴリズムが単純になった。最上位ビットをフラグに使っているだけなので、 2^{31} 文字以下ならば u32 でも十分だが、コンパイル時の保証は得られない。

2.2 理論

2つの接尾辞 $S[i..]$ と $S[i+1..]$ は長さが異なるので、大小関係が決まる。そこで、接尾辞およびその頭文字の型を次のように定義する。

定義 2.2.1: $S[i..] \geq S[i+1..]$ が成り立つとき、 $S[i]$ を L 型・S 型の接尾辞という。ただし、末尾の番兵は S 型であるとする。L 型・S 型の接尾辞の頭文字も L 型・S 型とする。

定理 2.2.1: $S[i] \geq S[i+1]$ のとき、 $S[i]$ は L 型・S 型である。 $S[i] = S[i+1]$ のとき、両者の型は一致する。

証明 前者は自明。後者については、 $S[i+k] \neq S[i+k+1]$ が成り立つ最小の正整数 k を考えればよい。□

接尾辞配列 SA は当然頭文字についてソートされる。同じ頭文字 $S[i]$ をもつ接尾辞たちが記録される部分をバケット $S[i]$ と呼ぶことにする。

系 2.2.1: 各バケットは L 型接尾辞と S 型接尾辞で区切られている。

証明 頭文字を固定して考える。この頭文字をもつ最小の S 型接尾辞を $S[i..]$ とかくと、 $S[i] < S[i+1]$ が成り立つ。同様に最大の L 型接尾辞を $S[j..]$ とかくと、 $S[j] > S[j+1]$ が成り立つ。したがって、 $S[i+1] > S[j+1]$ である。□

L 型接尾辞 $S[i..]$ の性質より、 $S[i..] > S[i+1..]$ が成り立つ。したがって、バケット $S[i]$ に書き込まれる L 型接尾辞の 2 文字目以降からなる接尾辞はバケットの左側にある。同様に S 型接尾辞の 2 文字目以降からなる接尾辞は対応するバケットよりも右側にある。これらの観察より次の定理を得る。

定理 2.2.2: S 型接尾辞がソート済みならば、L 型接尾辞も $\Theta(|S|)$ でソートできる。逆も成り立つ。

証明 L 型接尾辞をソートすることにする。接尾辞配列を昇順に走査し、初期化済みの要素を探す。これを $SA[i]$ とかく。もし $S[SA[i]-1..]$ が L 型接尾辞ならば対応するバケットに前から詰めて書き込む。同じバケットに属する接尾辞はその 2 文字目以降の部分からなる接尾辞についてもソートされているから、各バケットに書き込まれた L 型接尾辞はソートされて

いる。バケットソートと累積和でバケットの左端を求めることができるので、アルゴリズムは線形時間で動作する。逆も同じように証明できる。□

定理 2.2.2 の証明では S 型接尾辞のうち、1 だけ長い接尾辞が L 型であるもののみを利用している。これより、次の事実を得る。

定義 2.2.2: $S[i-1..]$ が L 型であるような S 型接尾辞 $S[i..]$ をとくに LMS 型 (leftmost S type) という。文字 $S[i]$ についても同様に定義する。

系 2.2.2: LMS 型接尾辞がソートされているとき、接尾辞配列を $\Theta(|S|)$ で構築できる。

証明 LMS 型接尾辞がソートされているとする。定理 2.2.2 より、接尾辞配列 SA を昇順に走査して L 型接尾辞をソートできる。さらに、接尾辞配列を降順に走査して S 型接尾辞をソートすることができる。□

補題 2.2.1: LMS 型接尾辞は高々 $\lfloor \frac{|S|}{2} \rfloor$ 個しかない。

補題 2.2.1 より考えるべき接尾辞の数を半分以下にできましたが、問題サイズは半分になっていません。接尾辞の長さが高々 $|S|$ だからです。LMS 型接尾辞を比較することを考えると、2 つの LMS 型文字に挟まれてできる部分文字列ごとに比較すればよいことが分かります。これを利用して文字列を圧縮し、問題のサイズを半分以下にすることが次の目標です。

定義 2.2.3: 部分文字列 $S[i..=j]$ のうち、 $S[i]$ と $S[j]$ のみが LMS 型であるものを LMS 型部分文字列という。とくに、番兵も LMS 型部分文字列であるとする。

補題 2.2.2: LMS 型接尾辞と LMS 部分文字列の数は一致する。

LMS 部分文字列を順序を保ったまま新しい文字に置き換えるためには、これらをソートする必要があります。逆に、LMS 部分文字列がソートされていれば

隣り合う 2 つの一致判定をとることで、順序を保ったまま改名することができます。

定理 2.2.3: LMS 部分文字列を $\Theta(|S|)$ でソートできる。

証明 バケットソートにより、すべての LMS 文字を線形時間でソートできる。[系 2.2.2](#) の証明と同じようにして、頭文字を除く最初の LMS 文字までソートできることが分かる。これはすべての LMS 部分文字列を含んでいる。□

系 2.2.3: 線形時間で接尾辞配列を構築できる。

証明 時間計算量は $T(n) \leq T\left(\frac{n}{2}\right) + \Theta(n) = \Theta(n)$ である。□

2.3 アルゴリズム

本節では線形時間接尾辞配列構築アルゴリズムの in-place な実装を紹介する。作業メモリを節約するために、バケットソートや区間幅の管理に接尾辞配列 SA 用のメモリを利用する。[2章 1節](#) で述べたように、文字の最上位ビットを活用して接尾辞の添え字とカウンターを区別する。また、文字の型の区別にも利用する。

2.3.1 文字の改名

文字列が座標圧縮されているので、接尾辞配列を使って各文字の登場回数を数えることができる。² 累積和をとると、バケットの右端を得ることができる。[系 2.2.1](#) より、L 型文字をバケットの左端、S 型文字をその右端の添え字に改名してもよい。

アルゴリズム 2.3.1.1:

```
1 fn rename(text: &mut [u8], sa: &mut [u8]) {
2     // 文字の登場回数を数える
3     sa.fill(0);
4     text.iter().for_each(|s| sa[*s] += 1);
5     // 累積和をとり、バケットの右端を求める。
6     (1..sa.len()).for_each(|i| sa[i] += sa[i-1]);
```

rust

² 実際には $S[i] < |S|$ で十分。

```

7   // 改名する。
8   {
9     // 番兵はS型。その1つ前はL型。
10    let mut is_s_type = false;
11    // 番兵は1つしかないので、改名しても`0`のまま。
12    // 直前の要素を型の判別に使うので、更新を遅延する。
13    let mut name = [0, 0];
14    for i in (0..text.len()).rev().skip(1) {
15      if text[i] != text[i+1] {
16        is_s_type = text[i] < text[i+1]
17      }
18      // 遅延を反映する
19      if let Some(s) = text.get_mut(i+2) {
20        *s = name[i%2]
21      }
22      // 反映を遅延する
23      name[i%2] = if is_s_type {
24        sa[text[i]]-1
25      } else {
26        // L型は最上位ビットを立てておく。
27        sa[text[i]-1] | lusize.rotate_right(1)
28      };
29    }
30    if text.len() >= 2 {
31      text[..2].copy_from_slice(&name);
32    }
33  }
34 }
35 /// 文字の型を判定する。
36 fn is_s_type(s: usize) -> bool {
37   s & lusize.rotate_right(1) == 0
38 }
39 /// L型文字ならバケットの左端、S型文字ならバケットの右端の添え字を返す。

```

```

40 fn str_to_ptr(s: usize) -> usize {
41     s & (!0 >> 1)
42 }

```

文字の最上位ビットをフラグに使うことで、文字の型をいつでも定数時間で求めることができるようになった。定理 2.2.1 から文字の型を判別するには S を逆順に走査する必要がある。したがって、系 2.2.2 接尾辞配列を操作する際にはこの手法は利用できない。³ また自由な方向に走査できるので、メモリ IO の回数を削減できるかもしれない。

2.3.2 LMS 文字のソート

LMS 部分文字列をソートするために、LMS 文字をソートする。LMS 文字をバケットの正しい位置に置く際に注意が必要である。LMS 文字はバケットの右端の情報をもつが、愚直に書き込むと $O(|S|^2)$ かかってしまう。

そこで、バケットの右端には次に書き込むべき場所への距離を書き込むことにする。これはすでに書き込んだ文字の個数に 1 を加えたものである。バケット内に LMS 文字が 1 つしかない場合、カウンターは不要なので単に書き込んでしまう。これと区別するために 1 を加えた。

LMS 文字についてソートされていればよいので、バケット内での順序は自由にしてよい。系 2.2.2 では文字ではなく接尾辞の順序を保つので、全体をシフトする必要がある。

アルゴリズム 2.3.2.1:

```

1 fn sort_lms_char(text: &[usize], sa: &mut [usize]) { rust
2     // LMS 文字を数える。
3     sa.fill(0);
4     text.windows(2).for_each(|s| {
5         if !is_s_type(s[0]) && is_s_type(s[1]) {
6             sa[str_to_ptr(s[1])] += 1
7         }
8     });
9
10    // カウンターの最上位ビットを立てて区別する。 実際はグローバルな定数。

```

³別の方法はある。

```

11  const COUNTER_FLAG: usize = 1.rotate_right(1);
12  const COUNT_ZERO: usize = COUNTER_FLAG;
13  const COUNT_ONE: usize = COUNT_ZERO + 1;
14  text.windows(2).enumerate().for_each(|(i, s)| {
15      if !is_s_type(s[0]) && is_s_type(s[1]) {
16          let ptr = str_to_ptr(s[1]);
17
18          if sa[ptr] == COUNT_ONE {
19              // 唯一の要素なので、カウンターは不要
20              sa[ptr] = i+1;
21          } else if sa[ptr-1] == COUNT_ZERO {
22              // バケット内で最初の文字
23              sa[ptr-1] = i+1;
24              // 引き算だとフラグが消えてしまう
25              sa[ptr] = COUNT_ZERO+2;
26          } else {
27              // バケット内で2文字目以降
28              let diff = sa[ptr] - COUNT_ZERO;
29              if sa.get(ptr.wrapping_sub(diff)).is_some_and(|sa|
30                  *sa == COUNT_ZERO) {
31                  // 次に書き込む先はバケット内か、L型用のバケット
32                  sa[ptr-diff] = i+1;
33                  sa[ptr] += 1;
34              } else {
35                  // 隣のS型用のバケットにはみ出してしまう。
36                  // バケット内で最後の1文字なので、カウンターを上書きしてよ
37                  // い。
38                  sa[ptr] = i+1
39              }
40          });
41      // カウンターを除去して、右に1だけシフトする。
42  }

```

```

43     let mut i = sa.len();
44     while i > 0 {
45         i -= 1;
46         if sa[i] > COUNT_ONE {
47             // カウンターを最後のLMS文字で上書きする.
48             let diff = sa[i] - COUNT_ZERO - 1;
49             sa[i] = sa[i - diff];
50             sa[i - diff] = COUNT_ZERO;
51             // LMS文字がある部分はスキップする.
52             i -= diff;
53         }
54     }
55 }
```

2.3.3 誘導ソート

系 2.2.2 の証明で用いたアルゴリズムを誘導ソートという。誘導ソートはこのアルゴリズムの肝であり、ソート済みの LMS 文字から LMS 部分文字列の順序を導き、ソートされた LMS 型接尾辞から接尾辞配列を誘導する。

アルゴリズム 2.3.3.1 (Induced-Sort):

```

1 fn induced_sort(text: &[u8], sa: &mut [u8]) {
2     // LMS型からL型を誘導する
3     sort_l_types(text, sa);
4     // LMS文字を除去する。番兵は誘導できないので、残しておく。
5     sa.iter_mut().skip(1).for_each(|i| {
6         // L型以外の文字は除去してよい。
7         if text.get(*i).is_some_and(|s| is_s_type(*s)) {
8             *i = COUNT_ZERO
9         }
10    });
11    // L型からS型を誘導する
12    sort_s_types(text, sa);
```

```

13 }
14
15 // S型の誘導ソートもほとんど同じだが、SAを降順に走査する。カウンターは
16 // 残らないので、あとから削除する必要はない。
17 fn sort_l_types(text: &[usize], sa: &mut [usize]) {
18     // LMS文字をソートするときとほとんど同じ。違いはSではなくSAを昇順に走
19     // 査すること。
20
21     // L型文字の出現数を数える
22     text.iter().for_each(|s| {
23         if !is_s_type(*s) { sa[str_to_ptr[*s]] += 1 }
24     });
25
26     // バケットに昇順に書き込む。
27     // バケットの中身を1だけシフトしたときには添え字を進めないので、for
28     // ループは使えない。
29     let mut i = 0;
30     while i < sa.len() {
31         // 初期化済みで、1つ長い接尾辞が存在するもの
32         if sa[i] > 0 && sa[i] & COUNTER_FLAG == 0 {
33             let s = text[sa[i]-1];
34             let ptr = str_to_ptr(s);
35             // シフトしても影響がなかった場合、同じ操作はしない。
36             if !is_s_type(s) && sa[ptr] & COUNTER_FLAG != 0 {
37                 if sa[ptr] == COUNT_ONE {
38                     // バケット内で唯一のL型文字
39                     sa[ptr] = sa[i]-1
40                 } else if sa[ptr+1] == COUNT_ZERO {
41                     // バケット内で最初のL型文字
42                     sa[ptr+1] = sa[i]-1;
43                     sa[ptr] = COUNT_ZERO + 2;
44                 } else {
45                     let diff = sa[ptr] - COUNT_ZERO;
46                     if sa.get(ptr + diff).is_some_and(|sa| *sa ==
47                         COUNT_ZERO) {
48                         // 書き込まれていない場合

```

```

44         sa[ptr+diff] = sa[i]-1;
45         sa[ptr] += 1;
46     } else {
47         // 書き込まれているなら、バケット内で最後のL型文字
48         sa[ptr] = sa[i]-1;
49         // 順番を保つ
50         sa[ptr..ptr+diff].rotate_left(1);
51         // 今いる場所がシフトした場合、添え字をインクリメントして
52         // はけない
53         continue;
54     }
55 }
56 i += 1;
57 }
58 }
59 // カウンターを除去して、バケットを左にシフト
60 let mut i = 0;
61 while i < sa.len() {
62     if sa[i] > COUNT_ZERO {
63         // カウンターを初期化して左にシフト
64         let diff = sa[i] - COUNT_ZERO;
65         sa[i] = COUNT_ZERO;
66         sa[i..i + diff].rotate_left(1);
67         // L型文字が書き込まれている部分はスキップ
68         i += diff
69     } else {
70         i += 1
71     }
72 }
73 }
```

2.3.4 LMS 部分文字列のソート

LMS 文字はソート済みなので, LMS 部分文字列をソートできる. 部分問題を構築するメモリを確保するために, LMS 型の文字を SA の末尾に集める. LMS 文字列を改名したら, S での登場順に並び替え, 部分問題を再帰的に解く.

LMS 文字は高々 $\lfloor \frac{|S|}{2} \rfloor$ 個しかないので, 改名後の文字は SA[SA[i]/2] に書き込めばよい. 部分問題は SA の先頭に詰めて書き込まれ, その接尾辞配列は末尾に書き込まれる.

アルゴリズム 2.3.4.1:

```
1  /// LMS部分文字列を改名して, 新しい text と sa を返す. これを rust  
2  // 再帰的に解く.  
3  fn sort_lms_substrings(text: &[u8], sa: &mut [u8]) ->  
4  (&mut [u8], &mut [u8]) {  
5      induced_sort(text, &mut sa);  
6      // LMS文字を末尾に集める.  
7      let n_lms = {  
8          let mut n_lms = 0;  
9          for i in (0..sa.len()).rev() {  
10              // text[sa[i]]はLMS型  
11              if sa[i] > 0 && is_s_type(text[sa[i]]) && !  
12                  is_s_type(text[sa[i]-1]) {  
13                  n_lms += 1;  
14                  sa[sa.len()-n_lms] = sa[i]  
15              }  
16          }  
17          // 不要なノードは初期化  
18          let r = sa.len()-n_lms;  
19          sa[..r].fill(COUNT_ZERO);  
20          n_lms  
21      };  
22      // LMS部分文字列を改名する.  
23      let kind_lms = {  
24          let mut kind_lms = 0;  
25          let i = sa.len()-n_lms;  
26          // これは番兵
```

```

25     sa[sa[i]/2] = 0 // text[sa[i]]
26     for i in i+1..sa.len() {
27         // LMS部分文字列の長さを求める
28         let nl = text[sa[i-1]..].windows(2).take_while(|s| {
29             // 次のLMS型まで進む
30             is_s_type(s[0]) || !is_s_type(s[1])
31         }).count();
32         let nr = text[sa[i]..].windows(2).take_while(|s| {
33             // 次のLMS型まで進む
34             is_s_type(s[0]) || !is_s_type(s[1])
35         }).count();
36         // 一致しないなら次に大きな文字を割り当てる。
37         if text[sa[i-1]..sa[i-1]+nl] != text[sa[i]..sa[i]+nr] {
38             kind_lms += 1;
39         }
40         sa[sa[i]/2] = kind_lms
41     }
42
43     kind_lms
44 };
45 // 改名したLMS部分文字列を前に集める。これはtextに現れた順になっている。
46 {
47     let mut n = 0;
48     for i in 0..sa.len()-n_lms {
49         if sa[i] & COUNTER_FLAG == 0 {
50             sa[n] = std::mem::replace(&mut sa[i], COUNT_ZERO);
51             n += 1;
52         }
53     }
54 }
55
56 let (pre, sa) = sa.split_at_mut(sa.len()-n_lms);
57 let (text, _) = sa.split_at_mut(n_lms);

```

```
58     (text, sa)
59 }
```

2.3.5 LMS 接尾辞のソート

LMS 部分文字列の接尾辞配列が得られたので、LMS 接尾辞がソートできた。ここで、部分問題が SA の前半に、その接尾辞配列が SA の後半に詰めて書き込まれているとする。接尾辞配列を正しいバケットに格納するために、SA の後半にソートされた LMS 型接尾辞を詰める。その語、昇順にバケットに書き込む。

アルゴリズム 2.3.5.1:

```
1  fn sort_lms_suffixes(text: &[u8], sa: &mut [u8]) rust
2  {
3      // 部分問題をLMS型接尾辞で登場順に上書きする
4      let mut n = 0;
5      text.windows(2).enumerate().for_each(|(i, s)| {
6          if !is_s_type(s[0]) && is_s_type(s[1]) {
7              sa[n] = i+1;
8              n += 1;
9          }
10     });
11
12     // 部分問題の解から、LMS型接尾辞をソートする。
13     {
14         let l = sa.len() - n;
15         for i in l..sa.len() {
16             sa[i] = sa[sa[i]];
17         }
18         sa[..n].fill(COUNT_ZERO);
19     }
20     n = 0;
21     // 番兵は最小の接尾辞
22     let mut s = 0;
```

```

23     sa[0] = sa[l]
24     for i in l+1..sa.len() {
25         // バケット内でのオフセット
26         if text[sa[i]] == s { n += 1 } else { n = 0 };
27         s = text[sa[i]];
28         // 同じ場所に書き込むことがあるので、初期化してから上書きする。
29         sa[str_to_ptr(s)-n] = std::mem::replace(&mut sa[i],
30                                         COUNT_ZERO);
30     }
31 }
32 }
```

2.3.6 まとめ

アルゴリズムの全体は下記の通り。参照は長さをもつスマートポインターなので、実は $O(\log|S|)$ の作業メモリを使っている。作業メモリを $O(1)$ に抑えるためには生ポインターを使う必要がある。これは unsafe rust の領分なので、本稿では解説しない。

アルゴリズム 2.3.6.1 (Create-Suffix-Array):

```

1  fn suffix_array(text: &mut [usize], sa: &mut [usize]) { rust
2      assert_eq!(text.len(), sa.len());
3      todo!("base case");
4
5      rename(&mut text, &mut sa);
6      sort_lms_char(&text, &mut sa);
7      {
8          let (text, sa) = sort_lms_substrings(&text, &mut sa);
9          suffix_array(&mut text, &mut sa);
10     }
11     sort_lms_suffixes(&text, &mut sa);
12     induced_sort(&text, &mut sa);
13 }
```

2.4 非再帰アルゴリズム

紹介したアルゴリズムの再帰木は鎖型になっている。行きがけと帰りがけでループを2つ用意すると非再帰で書けるはずである。

借用規則を守るために `slice.split_at_mut()` を多用することになる。この可変参照を管理するためにスタックを使うかも？

2.5 制約の緩和

本稿では座標圧縮を仮定したが、文字列の種類が $O(|S|)$ であれば、線形時間で接尾辞配列をもとめるアルゴリズムが存在する。ASCII 文字や Unicode 文字など、コンピューターで利用できる文字の種類は固定なので、接尾辞配列を線形時間で計算できる。

3 LCP 配列

SA 上で差分計算する。RMQ で定数時間クエリ。

4 参考文献

4.1 接尾辞配列の線形時間アルゴリズム

- LI, Zhize; LI, Jian; HUO, Hongwei. Optimal in-place suffix sorting. *Information and Computation*, 2022, 285: 104818. <https://doi.org/10.1016/j.ic.2021.104818>
- NONG, Ge; ZHANG, Sen; CHAN, Wai Hong. Two efficient algorithms for linear time suffix array construction. *IEEE transactions on computers*, 2010, 60.10: 1471-1484. <https://doi.org/10.1109/TC.2010.188>