

1 Suchen in Texten

1.1 Überblick

In Anwendungen, bei denen Texte verarbeitet werden, ist es oft wichtig, in einem langen Text effizient nach dem Vorkommen eines kürzeren Such-Textes (genannt Muster) suchen zu können. Wir betrachten einen Algorithmus, der sehr gut geeignet ist, wenn man in demselben Text nacheinander nach vielen verschiedenen Mustern suchen will. In diesem Fall lohnt es sich, eine etwas aufwendige Vorverarbeitung des Textes zu machen, weil dadurch das anschließende Suchen nach Mustern beschleunigt werden kann.

Ein Wort w ist eine Folge von Zeichen aus einem festen Alphabet Σ (mit konstanter Alphabet-Größe $|\Sigma|$, z.B. ASCII-Zeichen mit $|\Sigma| = 256$). Mit $|w|$ bezeichnen wir die Länge (Anzahl Zeichen) des Wortes w und mit $w[i]$ (für $0 \leq i < |w|$) das i -te Zeichen von w . $w[i..j]$ bezeichnet das Teilwort von w , das mit dem i -ten Zeichen beginnt und mit dem j -ten Zeichen endet. Ein Teilwort der Form $w[i..|w| - 1]$ heißt auch *Suffix* von w . Ein Suffix ist echt (engl. „proper“), wenn $i \neq 0$.

Der gegebene Text ist ein (langes) Wort, das wir mit X bezeichnen. Sei $|X| = n$. Wir nehmen an, daß der Text mit dem speziellen Zeichen '\$' endet, das sonst weder im Text noch in den Suchmustern vorkommt. Man nennt '\$' einen „Sentinel“. Diese Annahme stellt sicher, daß kein Suffix der Anfang (Präfix) eines anderen Suffix von X sein kann. (Ein Suffix, welches das echte Präfix eines anderen Suffixes ist, bezeichnen wir als „eingebettet“ (engl. „nested“).) Mit suffix_i bezeichnen wir das Suffix $X[i..n - 1]$ von X . Ein Suchmuster Y der Länge m kommt in X vor, wenn es ein i gibt mit $X[i..i + m - 1] = Y$. Wir wollen in Zeit $O(n \log n)$ eine Datenstruktur aufbauen, die es uns erlaubt, in X nach einem beliebigen Muster Y der Länge m in Zeit $O(m \log n)$ zu suchen.

Eine dafür geeignete Datenstruktur sind **Suffix Arrays**. Wir wollen hier zeigen, wie man Suffix Arrays direkt und relativ effizient in linear-logarithmischer Zeit ($n \log n$) konstruieren kann. Wir halten uns dabei an den Original-Artikel von Udi Manber und Gene Myers [1]. Die Sortier-Methode könnte man als Recursive-Forward-Bucket-Sort bezeichnen. Wir werden zu einem gegebenen Text der Länge n ein Integer-Array pos der Länge n konstruieren, der die lexikographisch sortierten Indizes der Suffixe enthält.

Text	b	c	c	a	a	b	a	b	a	\$
Index	0	1	2	3	4	5	6	7	8	9
Sort. Index	9	8	3	6	4	7	5	0	2	1
Suffix	\$	a\$	aababa\$	aba\$	ababa\$	ba\$	baba\$	bccaababa\$	caababa\$	ccaababa\$

Abbildung 1: Suffix Array für den Text "bccaababa\$"

1.2 Suchen

Die Suche nach dem Muster Y vollzieht sich in zwei Schritten. Zuerst wird der linke Index i in **pos** gesucht, so dass Y ein Präfix vom **pos**[i]-ten Suffix ist, anschließend wird die entsprechende rechte Stelle gesucht. Dabei verwenden wir binäre Suche.

Zur Suche des linken Index i in **pos** mit Y als Präfix vom **pos**[i]-ten Suffix beginnen wir mit den Indizes $L = 0$ und $R = n - 1$. Als Invariante wollen wir erhalten, dass das **pos**[L]-te Suffix lexikographisch kleiner als Y ist und das **pos**[R]-te Suffix lexikographisch größer oder gleich Y ist.

Falls $R - L = 1$ ist brechen wir ab. In jedem Schritt betrachten wir **pos**[M]. Wenn das **pos**[M]-te Suffix lexikographisch größer oder gleich Y ist, so setzen wir $R := M$, andernfalls $L := M$. Siehe Abb. 2.

<p>Input: Text X, pos, Muster Y. Output: L_Y. Algorithmus: $L := 0$; $R := X - 1$; if $X[\mathbf{pos}[L]]$ ist lex. größer als Y then return 0; else if $X[\mathbf{pos}[R]]$ ist lex. kleiner als Y then return X; else while $R - L > 1$ do $M := \lceil (L + R)/2 \rceil$; if $X[\mathbf{pos}[M]]$ ist lex. kleiner als Y then $R := M$; else $L := M$; return R;</p>
--

Abbildung 2: Binäre Suche nach dem linken Index i in **pos**, so dass Y ein Präfix des **pos**[i]-ten Suffix ist.

1.3 Index-Erstellung: Sortieren der Suffixe

Das Sortieren der Suffixe erfolgt in $\log n$ Phasen. In der h -ten Phase (wir beginnen bei 0) werden die Suffixe so sortiert, dass die lexikographische Ordnung bezüglich der ersten 2^h Buchstaben richtig ist. Dazu nutzen wir die Informationen aus der vorherigen Phase und den Fakt, dass wir Suffixe sortieren. Wir teilen die Suffixe in Bereiche („Buckets“) auf, so dass alle Suffixe eines Buckets in den ersten 2^h Buchstaben gleich sind. Wir sortieren dann jeweils die Suffixe eines Buckets. Angenommen wir haben zwei Suffixe **suf** $_j$ und

suf_k aus einem Bucket, die bei den ersten 2^h Buchstaben gleich sind. Wir müssen also ab dem 2^h -ten Buchstaben 2^h weitere Buchstaben vergleichen. Das Ergebnis ist dasselbe, wie wenn wir von den Suffixen suf_{j+2^h} und suf_{k+2^h} die ersten 2^h Buchstaben vergleichen. Wir kennen es schon aus der ersten Phase und wir müssen nur in der Lage sein, die Position der Indizes $j + 2^h$ und $k + 2^h$ in pos zu finden. Dazu bauen wir am Ende jeder Phase ein zu pos reverses Array prm auf (dann ist $\text{pos}[\text{prm}[i]] = i$). Da wir in jeder Phase nur $O(n)$ Aufwand treiben wollen, können wir die Elemente eines Buckets nicht einfach mit einem $n \log n$ -Algorithmus sortieren. Statt dessen gehen wir so vor, dass wir das Suffix Array von vorne bis hinten in seiner groben Sortierung bucketweise durchlaufen und dabei die bestehende Sortierung in die Buckets übertragen. D.h. wenn suf_k ein Element des ersten Buckets ist, so muss suf_{k-2^h} in seinem Bucket vor allen Elementen kommen, die zu Suffixen in höheren Buckets korrespondieren. Schematisch ist das ganze in Abb. 3 dargestellt.

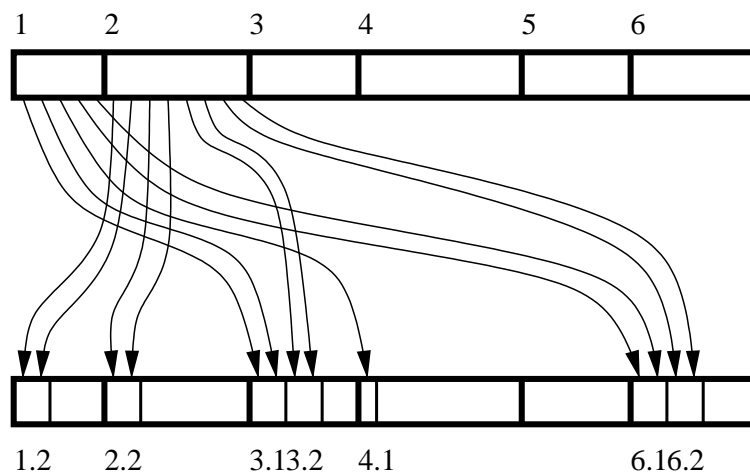


Abbildung 3: Schematische Darstellung der Forward-Sortierung beim Aufbau eines Suffix Arrays

Je nachdem wie geschickt man die Buckets verwaltet, braucht man mehr oder weniger Platz dafür. Wir benutzen noch zusätzlich zwei Arrays \mathbf{BH} und $\mathbf{B2H}$ von Bool'schen Werten und einen temporären Array \mathbf{count} (alle der Länge n). $\mathbf{B2H}$ und \mathbf{count} sind für temporäre Werte. $\mathbf{BH}[i]$ ist 'wahr', wenn an der Stelle i ein neuer Bucket in pos beginnt und sonst 'falsch'.

In einer initialen Phase werden die Suffixe mit einem Radix-Sort nach ihrem ersten Buchstaben sortiert. Nach der ersten Schleife gilt für einen beliebigen Buchstaben d aus Σ , dass in $i := \text{pos}[d]$ das letzte (von hinten) Vorkommen von d in X steht (oder -1). In $\text{prm}[i]$ steht ein Zeiger auf die nächst-letzte Position, wo d in X vorkommt (oder -1). Diese Kette wird dann in der zweiten Schleife abgelaufen und man merkt sich in $\text{prm}[i]$ dann stattdessen die Position an der das i -te Suffix in der lexikographischen Ordnung bezüglich des ersten Buchstabens vorkommt (dies ist $c - c$ wird für jedes Suffix um eins inkrementiert).

Dann beginnt das „rekursive“ Sortieren. $\text{prm}[i]$ zeigt normalerweise auf die Position

des i -ten Suffix in **pos**. Wir iterieren mittels **BH** über die Buckets und setzen **prm**[i] auf das erste Element des Buckets, in dem sich das i -te Suffix befindet. Ausserdem setzen wir **count**[l] = 0 (l ist die linke Bucket-Grenze). Danach durchlaufen wir das Array **pos** Bucket für Bucket und betrachten jeweils das um 2^h verlängerte Suffix $d := \mathbf{pos}[i] - 2^h$, welches sich im Bucket, der bei $k := \mathbf{prm}[d]$ beginnt, befindet. Die ersten **count**[k] Elemente sind schon belegt und das j -te Suffix wird an Stelle $k + \mathbf{count}[k]$ gesetzt und **count**[k] wird inkrementiert. Für jedes Suffix merken wir uns in **B2H**, ob es bewegt wurde. Anschliessend setzen wir mit diesen Informationen in **B2H** die neuen Bucket-Grenzen. In einem letzten Schritt müssen wir **BH** und **pos** wieder aus **B2H** und **prm** richtig setzen. Insgesamt benötigt eine Phase also $O(n)$ Schritte.

Zum Iterieren über die Buckets kann man jeweils zwei Zeiger l und r benutzen. Man setzt l und r auf den Anfang des Buckets und inkrementiert r bis **BH**[r] = **true**. Der nächste Bucket beginnt bei r . Man kann so alle Buckets in linearer Zeit durchlaufen. Abb. 4 zeigt die initiale Phase, das iterative Sortieren die Abb. 5.

```

Input: Text  $X$  der Länge  $n$ .
Output: Suffix Sortierung pos für den ersten Buchstaben,
          Initialisierung von pos, prm, BH, B2H und count.
Algorithmus:
  lege die Arrays pos, prm, BH, B2H und count an;
  //initiale Phase: Radix Sort auf den ersten Buchstaben
  for  $c$  from 1 to  $|\Sigma|$  do
    pos[ $c$ ] := -1;
  for  $i$  from 0 to  $n - 1$  do
     $b := \mathbf{pos}[X[i]]$ ;
    pos[ $X[i]$ ] :=  $i$ ;
    prm[ $i$ ] :=  $b$ ;
   $c := 1$ ;
  for  $d$  from 1 to  $|\Sigma|$  do
     $i := \mathbf{pos}[d]$ ;
    while  $i \neq -1$  do
       $j := \mathbf{prm}[i]$ ;
      prm[ $i$ ] :=  $c$ ;
      if  $i = \mathbf{pos}[d]$  then
        BH[ $c$ ] := true;
      else
        BH[ $c$ ] := false;
       $c := c + 1$ ;
       $i := j$ ;
  //Arrays initialisieren
  BH[ $n$ ] := true;
  for  $i$  from 0 to  $n - 1$  do
    pos[prm[ $i$ ]] :=  $i$ ;

```

Abbildung 4: Initiale Sortier-Phase (Radix Sort)

```

Input: Text  $X$  der Länge  $n$ , Initialisierte Arrays pos, prm, BH,
B2H und count.
Output: Suffix Sortierung in pos.
Algorithmus:
  for  $h$  from 0 to  $\lfloor \log n \rfloor$  do
    for each Bucket  $[l, r)$  do
      count $[l] := 0$ ;
      for  $i$  from  $l$  to  $r - 1$  do
        prm $[\text{pos}[i]] := l$ ;
    for each Bucket  $[l, r)$  do
      for  $i$  from  $l$  to  $r - 1$  do
         $d := \text{pos}[i] - 2^h$ ;
        if  $d < 0$  or  $d \geq n$  then continue;
         $k := \text{prm}[d]$ ;
        prm $[d] := k + \text{count}[k]$ ;
        count $[k] := \text{count}[k] + 1$ ;
        B2H $[\text{prm}[d]] := \text{true}$ ;
      for  $i$  from  $l$  to  $r - 1$  do
         $d := \text{pos}[i] - 2^h$ ;
        if  $d < 0$  or
            $d \geq n$  or
           not B2H $[\text{prm}[d]]$  then
          continue;
         $k := \min\{j : j > \text{prm}[d] \text{ and } (\text{BH}[j] \text{ oder nicht } \text{B2H}[j])\}$ ;
        for  $j$  from prm $[d] + 1$  to  $k - 1$  do B2H $[j] := \text{false}$ ;
    for  $i$  from 0 to  $n - 1$  do
      pos $[\text{prm}[i]] := i$ ;
    for  $i$  from 0 to  $n - 1$  do
      if B2H $[i]$  and not BH $[i]$  then
        BH $[i] := \text{B2H}[i]$ ;

```

Abbildung 5: Recursive-Bucket-Sort-Phase

Literatur

- [1] Udi Manber and Gene Myers. Suffix arrays: A new method for on-line string searches. *SIAM J. COMPUT.*, 22(5):935–948, oct 1993.
- [2] Mohamed Ibrahim Abouelhoda, Enno Ohlebusch, and Stefan Kurtz. The enhanced suffix array and its applications to genome analysis. In *Proceedings of the 2nd Workshop on Algorithms in Bioinformatics*, volume 2452 of LNCS, pages 449–463. Springer, 2002.