

Pregatirea lotului national de informatica

Alba Iulia, 2004



SUFFIX ARRAYS

Prof. Emanuela Cerchez

Liceul de Informatica “Grigore Moisil” Iasi

Problema

Fie $A = a_0 a_1 \dots a_{N-1}$ un sir de lungime N peste alfabetul S .

Query (W) :

Fie $W = w_0 w_1 \dots w_{P-1}$ un sir de lungime P peste alfabetul S . Sa se determine toate aparitiile sirului W in sirul A .

Observatie

Sirul A este fixat si trebuie sa suporte numeroase interogari.

Solutii

- ⌘ **Cautare naiva** - $O(N \cdot P)$
- ⌘ **KMP** - $O(N + P)$
- ⌘ **Suffix trees** - $O(N)$ complexitate timp pentru constructia arborelui si $O(N \cdot |S|)$ complexitate spatiu, apoi $O(P)$ complexitate timp pentru query => complexitate spatiu prohibitiva. Prin compresie - $O(N \log |S|)$ complexitate timp pentru constructie arbore si $O(N)$ complexitate spatiu, apoi $O(P \log |S|)$ complexitate timp pentru query.

Suffix arrays

- ⌘ Este o structura de date proiectata pentru a raspunde eficient la interogari pentru un sir A fixat.
- ⌘ Complexitate spatiu: $O(N)$
- ⌘ Complexitate timp pentru interogari in cazul cel mai defavorabil $P + \lceil \log_2(N-1) \rceil$
- ⌘ Complexitate timp pentru constructie suffix-array: $O(N \log N)$ sau $O(N)$, folosind structuri de date suplimentare.

Cautare

- ⌘ Notam A_i = sufixul lui A care incepe la pozitia $i = a_i a_{i+1} \dots a_{N-1}$.
- ⌘ Asociem sirului A un vector Pos , de lungime N . $Pos[i] =$ pozitia de inceput in A al celui de-al k -lea cel mai mic sufix dintre $\{A_0, A_1, \dots, A_{N-1}\}$ (k -smallest), relativ la relatia de ordine lexicografica.
- ⌘ $A_{Pos[0]} < A_{Pos[1]} < \dots < A_{Pos[N-1]}$ (unde $<$ indica relatia de ordine lexicografica).

Exemplu

- ⌘ Fie $A = \text{aabaabab}$ un sir de lungime 8.
- ⌘ Cele 8 sufixe ale lui A sunt: $A_0 = \text{aabaabab}$,
 $A_1 = \text{abaabab}$, $A_2 = \text{baabab}$, $A_3 = \text{aabab}$,
 $A_4 = \text{abab}$, $A_5 = \text{bab}$, $A_6 = \text{ab}$, $A_7 = \text{b}$.
- ⌘ Sufixele in ordine lexicografica sunt: $\{A_0, A_3, A_6, A_1, A_4, A_7, A_2, A_5\}$.
- ⌘ Vectorul Pos este $(0, 3, 6, 1, 4, 7, 2, 5)$.

Notatii

- ⌘ Fie x un sir. Notam x^p prefixul de lungime p al lui x (adica $x_0x_1 \dots x_{p-1}$).
- ⌘ Fie x si y doua siruri $x <_p y$ daca si numai daca $x^p < y^p$ (se compara lexicografic prefixele de lungime p).
- ⌘ Analog se pot defini $\geq_p, \leq_p, >_p$.

Observatie. Toate sufixele care au p -prefixe egale apar in Pos pe pozitii consecutive.

Alte notatii

$$\text{⌘ } L_w = \min \{ k \mid W \leq_P A_{\text{Pos}[k]} \} \text{ sau } N$$

$$\text{⌘ } R_w = \{ \max \{ k \mid A_{\text{Pos}[k]} \leq_P W \} \text{ sau } -1$$

⌘ Deoarece Pos este ordonat si dupa relatia \leq_P deducem ca W apare incepand cu pozitia i in A ("se potriveste" cu $a_i a_{i+1} \dots a_{i+p-1}$)
daca si numai daca $i = \text{Pos}[k]$ unde $k \in [L_w, R_w]$. Deci numarul de aparitii ale lui W in A este $R_w - L_w + 1$.

Algoritmul de determinare L_w si R_w

⌘ Vectorul POS fiind ordonat in raport cu relatia \leq_P , putem determina L_w si R_w modificand algoritmul de cautare binara. Prin urmare sunt necesare $O(\log N)$ comparatii de siruri de lungime P , deci complexitatea este $O(P \log N)$.

Algorithm 1

```
if ( $W \leq_P A_{\text{Pos}[0]}$ )  $L_w = 0$ ;  
else  
if ( $W >_P A_{\text{Pos}[N-1]}$ )  $L_w = N$ ;  
else  
  {  $L = 0$ ;  $R = N - 1$ ;  
    while ( $R - L > 1$ )  
      {  $M = (L + R) / 2$ ;  
        if ( $W \leq_P A_{\text{Pos}[M]}$ )  $R = M$ ;  
          else  $L = M$ ;  
      }  
     $L_w = R$ ; }
```

Exemplu

⌘ Sa consideram sirul A din exemplul precedent ($A = aabaabab$) si $W = ab$. Vectorul Pos determinat este $(0, 3, 6, 1, 4, 7, 2, 5)$.

⌘ $L_w = \min \{ k \mid ab \leq_2 A_{Pos[k]} \} = 2$.

⌘ $R_w = \max \{ k \mid A_{Pos[k]} \leq_2 ab \} = 4$

⌘ Numarul de aparitii ale lui W este $4 - 2 + 1 = 3$ si aparitiile incep pe pozitiile $Pos[2] = 6$, $Pos[3] = 1$ si $Pos[4] = 4$.

Optimizari



⌘ Timpul de executie a algoritmului precedent poate fi imbunatatit, observand ca nu este necesar ca de fiecare data sa executam toate comparatiile aferente testului \leq_P de la inceput. Putem utiliza informatiile din comparatia curenta pentru a optimiza comparatiile urmatoare.

Optimizare 1

- ⌘ Fie v si w doua siruri. Notam cu $l_{cp}(v, w)$ lungimea celui mai lung prefix comun al sirurilor v si w . Observati ca $l_{cp}(v, w)$ poate fi determinat pe parcursul compararii sirurilor v si w .
- ⌘ Putem modifica algoritmul 1., retinand in variabila $l = l_{cp}(W, A_{Pos[L]})$ si in variabila $r = l_{cp}(W, A_{Pos[R]})$
- ⌘ Initial l este setat in urma compararii dintre W si $A_{Pos[0]}$, iar r este setat in urma compararii dintre W si $A_{Pos[N-1]}$.
- ⌘ In urma comparatiei efectuate intre W si $A_{Pos[M]}$ va fi actualizat l sau r .
- ⌘ Retinand valorile l si r , nu mai este necesar sa efectuam primele $h = \min\{l, r\}$ comparatii de simboluri atunci cand comparam W cu $A_{Pos[M]}$, deoarece $A_{Pos[L]} =_l W =_r A_{Pos[R]}$ implica faptul ca $A_{Pos[k]} =_h W$, pentru orice $L \leq k \leq R$. Aceasta modificare va reduce numarul de comparatii de simboluri, dar in cazul cel mai defavorabil complexitatea ramane $O(P \log N)$.

Optimizare 2

- ⌘ Pentru a reduce numărul de comparații de simboluri la $P + \lceil \log_2(N-1) \rceil$ vom precalcula informațiile referitoare la prefixe.
- ⌘ Să considerăm toate tripletele de forma (L, M, R) care pot apărea în algoritmul 1 ($0 \leq L < M < R \leq N-1$). Există $N-2$ triplete, fiecare având un unic punct de mijloc $0 < M < N-1$. Fie (L_M, M, R_M) tripletul care are mijlocul M .
- ⌘ Vom utiliza doi vectori suplimentari L_{lcp} și R_{lcp} având $N-2$ componente (cate una pentru fiecare triplet).
- ⌘ $L_{lcp}[M] = lcp(A_{Pos[M]}, A_{Pos[LM]})$
- ⌘ $R_{lcp}[M] = lcp(A_{Pos[M]}, A_{Pos[RM]})$

Exemplu

- ⌘ Sa consideram sirul din exemplu precedent $A = \text{aabaabab}$ cu vectorul $\text{Pos} (0, 3, 6, 1, 4, 7, 2, 5)$.
- ⌘ Cele 6 triplete (L, M, R) care pot aparea in timpul executiei algoritmului 1 sunt: $(0, 3, 7)$, $(0, 1, 3)$, $(3, 5, 7)$, $(1, 2, 3)$, $(3, 4, 5)$, $(5, 6, 7)$.
- ⌘ $\text{Llcp}[M] = \text{lcp}(A_{\text{Pos}[M]}, A_{\text{Pos}[LM]})$ (pentru $M=1, 6$).
- ⌘ $\text{Llcp}[1] = \text{lcp}(A_{\text{Pos}[1]}, A_{\text{Pos}[0]}) = 4$
- ⌘ $\text{Llcp}[2] = \text{lcp}(A_{\text{Pos}[2]}, A_{\text{Pos}[1]}) = 1$
- ⌘ $\text{Llcp}[3] = \text{lcp}(A_{\text{Pos}[3]}, A_{\text{Pos}[0]}) = 1$
- ⌘ $\text{Llcp}[4] = \text{lcp}(A_{\text{Pos}[4]}, A_{\text{Pos}[3]}) = 3$
- ⌘ $\text{Llcp}[5] = \text{lcp}(A_{\text{Pos}[5]}, A_{\text{Pos}[3]}) = 0$
- ⌘ $\text{Llcp}[6] = \text{lcp}(A_{\text{Pos}[6]}, A_{\text{Pos}[5]}) = 1$

Exemplu - continuare

⌘ $Rlcp[M] = lcp(A_{Pos[M]}, A_{Pos[LM]})$ (pentru $M=1, 6$).

⌘ $Rlcp[1] = lcp(A_{Pos[1]}, A_{Pos[3]}) = 1$

⌘ $Rlcp[2] = lcp(A_{Pos[2]}, A_{Pos[3]}) = 2$

⌘ $Rlcp[3] = lcp(A_{Pos[3]}, A_{Pos[7]}) = 0$

⌘ $Rlcp[4] = lcp(A_{Pos[4]}, A_{Pos[5]}) = 0$

⌘ $Rlcp[5] = lcp(A_{Pos[5]}, A_{Pos[7]}) = 1$

⌘ $Rlcp[6] = lcp(A_{Pos[6]}, A_{Pos[7]}) = 2$

Utilizarea vectorilor L_{lcp} si R_{lcp}

- ⌘ Constructia vectorilor L_{lcp} si R_{lcp} se poate face in etapa de sortare a sufixelor.
- ⌘ Sa consideram tripletul (L, M, R) corespunzator iteratiei curente, $h = \max(l, r)$.
- ⌘ In analiza noastra sa presupunem fara a restrange generalitatea ca $r \leq l = h$. Comparand h cu $L_{lcp}(M)$ obtinem urmatoarele cazuri:

Cazuri

⌘ **Cazul 1:** $Ll_{cp}(M) > l$

In acest caz, $A_{Pos[M]} =_{l+1} A_{Pos[L]} \neq_{l+1} W$ deci W trebuie sa fie in jumatatea din dreapta iar l ramane neschimbat.

⌘ **Cazul 2:** $Ll_{cp}(M) = l$

In acest caz stim ca primele l simboluri din $A_{Pos[M]}$ si W sunt egale; deci trebuie sa comparam simbolurile $l+1, l+2, \dots$ pana cand gasim simbolul $l+j$ a.i. $W \neq_{l+j} A_{Pos[M]}$. Simbolul $l+j$ va stabili daca LW este in partea stanga sau in partea dreapta. In oricare dintre cele doua situatii stim noua valoare pentru l sau pentru $r(l+j)$.

⌘ **Cazul 3:** $Ll_{cp}(M) < l$

In acest caz W se potriveste cu primele l simboluri din L si $< l$ simboluri din M , deci L_W va fi in partea stanga, iar noua valoare a lui r va fi $l_{cp}[M]$.

⌘ **Observatie.** Pentru cazul $l < r$ se analizeaza in acelasi mod, utilizand Rl_{cp} .

Cautare in $O(P + \log N)$ pentru LW

```
l=lcp(W, A_Pos[0]); r=lcp(W, A_Pos[N-1]);
if (l=P || W[l]<=A[Pos[0]+1]) LW=0;
else
if (r<P || W[r]<=A[Pos[N-1]+r]) LW=N;
else
  {L=0; R=N-1;
   while (R-L>1)
     {M=(L+R)/2;
      if (l>=r)
        if (Llcp[M]>=l)
          m=l+lcp(A_Pos[M]+1, W_l)
        else m=Llcp[M]
```

Cautare in $O(P + \log N)$ pentru LW

- continuare

```
    else //l<r
        if (Rlcp[M] >= r) m=r+lcp(APos[M]+r, Wr)
            else m=Rlcp[M];
    if (m==P || W[m] <= A[Pos[M]+m]) {R=M; r=m;}
        else {L=M; l=m;}
}
LW=R;
}
```

Sortarea sufixelor

- ⌘ Sortarea va fi realizata in $\lceil \log_2(N+1) \rceil$ etape, prin distribuire (*buckets sort*). Numerotam etapele 1, 2, 4, 8, ...
- ⌘ Vom utiliza $|S|$ galeti (*buckets*), cate una pentru fiecare simbol din alfabet.
- ⌘ In prima etapa, vom distribui sufixele in galeti in functie de prima litera (sufixul A_i va fi plasat in galeata corespunzatoare lui $A[i]$). Deci dupa prima etapa, fiecare galeata va contine sufixele care incep cu acelasi simbol (deci sunt sortate in raport cu \leq_1).
- ⌘ Rezultatul sortarii este memorat in vectorul POS . In plus vom utiliza un vector boolean BH care va indica modul de distribuire a sufixelor in cele m_1 galeti, astfel $BH[i]=1$ daca si numai daca $POS[i]$ indica primul sufix dintr-o galeata.

Sortarea sufixelor - continuare

- ⌘ Sa presupunem ca dupa H etape, sufixele sunt partitionate in m_H galeti, fiecare galeata continand sufixe care au acelasi H -prefix (adica sunt sortate in raport cu relatia \leq_H).
- ⌘ In etapa $2H$ prefixele vor fi sortate in raport cu relatia \leq_{2H} .
- ⌘ Fie A_i si A_k doua sufixe care dupa etapa H sunt in aceeasi galeata, adica $A_i =_H A_k$. Pentru a determina ordinea \leq_{2H} dintre A_i si A_k va trebui sa comparam urmatoarele H simboluri din aceste doua sufixe. Dar urmatoarele H simboluri din A_i sunt primele H simboluri din A_{i+H} , iar urmatoarele H simboluri din A_k sunt primele H simboluri din A_{k+H} (iar pentru aceste sufixe cunoastem deja ordinea \leq_H).

Sortarea sufixelor - continuare

- ⌘ Sa consideram A_i , primul sufix din prima galeata (adica $Pos[0]=i$) si sa consideram A_{i-H} (daca $i-H < 0$, ignoram A_i si trecem la sufixul indicat de $Pos[1]$, s.a.m.d)
- ⌘ Deoarece A_i incepe cu cel mai mic H -prefix, deducem ca A_{i-H} trebuie sa fie primul sufix din galeata sa de la pasul $2H$. Prin urmare vom muta pe A_{i-H} la inceputul galetii sale si vom marca acest lucru.
- ⌘ Pentru fiecare galeata va fi necesar sa cunoastem numarul de sufixe care au fost deja mutate la inceputul galetii (adica sunt ordonate in raport cu relatia \leq_{2H}).
- ⌘ Practic, vom parcurge toate sufixele in ordinea \leq_H , si pentru fiecare sufix A_i vom muta sufixul A_{i-H} (daca exista) pe urmatoarea pozitie disponibila din galeata sa de la pasul H , obtinand astfel sufixele sortate in ordinea \leq_H .

Structuri de date necesare

- ⌘ Vom utiliza 3 vectori de câte N componente, având semnificațiile:
- ⌘ $Pos[i]$ = poziția de început al celui de al i -lea cel mai mic sufix (în raport cu relația \leq_H).
- ⌘ $Prm[i]$ = inversul lui $Pos[i]$ ($Prm[Pos[i]] = i$).
- ⌘ $BH[i] = 1$ dacă și numai dacă $Pos[i]$ indică primul sufix (cel mai din stanga) dintr-o galeată (adică $A_{Pos[i]} \neq_H A_{Pos[i-1]}$).
- ⌘ De asemenea utilizăm doi vectori auxiliari:
- ⌘ $Count[i]$ = numărul de sufixe care au fost mutate la începutul galeții care conține al i -lea sufix.
- ⌘ $B2H[i] = 1$ dacă și numai dacă al i -lea sufix a fost mutat (deci este deja ordonat în raport cu \leq_{2H}).
- ⌘ Putem implementa galețile ca o listă simplu înlantuită alocată static.
 $Bucket[x] = -1$ dacă galeata corespunzătoare simbolului x este vidă sau poziția de început în șirul A a primului sufix din galeată;
- ⌘ $Link[i]$ = poziția de început în A a următorului sufix din galeată.

Descrierea etapei 2H

- ⌘ Initializam vectorul Count cu 0, si vectorul Prm a.i. $\text{Prm}[i]$ indica cea mai din stanga pozitie din galeata continand al i -lea sufix.
- ⌘ Parcurgem vectorul Pos in ordine crescatoare, cate o galeata la un pas. Fie (l, r) limitele stanga-dreapta ale galetii curente.
- ⌘ Sa notam $T[i] = \text{Pos}[i] - H$.
- ⌘ Pentru orice i , $l \leq i \leq r$, incrementam $\text{Count}[\text{Prm}[T[i]]]$, setam $\text{Prm}[T[i]] = \text{Prm}[T[i]] + \text{Count}[\text{Prm}[T[i]]] - 1$ si setam $\text{B2H}[\text{Prm}[T[i]]] = 1$.
- ⌘ Inainte de a trece la urmatoarea galeata, vom pastra valoarea 1 in B2H numai pentru cea mai din stanga componenta din galeata, celelalte componente din galeata le vom marca in B2H cu 0 (astfel in B2H va ramane marcat corect numai inceputul galetii).
- ⌘ Dupa parcurgerea tuturor galetilor, construim Pos (pe baza vectorului Prm) si $\text{BH} = \text{B2H}$.

Sortarea sufixelor - etapa 1

```
//initializare
for a∈S Bucket[a]=-1;
//distribuim sufixele in galeti, dupa prima litera
for (i=0; i<N; i++) (Bucket[A[i]], Link[i])=(i, Bucket[A[i]])
//construim vectorii Prm, Pos, BH
c=0;
for a∈S
    {i=Bucket[a];
     while (i!=-1) //parcurgem galeata
        {j=Link[i]; Prm[i]=c;
         if (i==Bucket[a]) //primul sufix din galeata
             BH[c]=1;
         else BH[c]=0;
         c++; i=j;}
BH[N]=1;
for (i=0; i<N; i++) Pos[Prm[i]]=i;
```

Sortarea sufixelor - continuare

```
for (H=1; H<N; H*=2) //sufixele sunt ordonate dupa <=H
{
  //parcurgem vectorul Pos galeata cu galeata
  for  $\forall$ galeata[l, r] //l, r sunt capetele galetii
  {
    Count[l]=0;
    for (c=l; c<=r; c++) Prm[Pos[c]]=1;
    d=N-H; e=Prm[d]; Count[e]++; B2H[Prm[d]]=1;
    for  $\forall$ galeata[l, r]
    {
      //distribuim sufixele in 2H ordine
      for (c=l; c<=r; c++)
      {
        d=Pos[c];
        if (d $\in$  [0, N-1])
        {
          e=Prm[d];
          Prm[d]=e+Count[e];
          Count[e]++;
          B2H[Prm[d]]=1;
        }
      }
    }
  }
}
```

Sortarea sufixelor - continuare

```
for (c=l; c<=r; c++)
    {d=Pos[c];
      if (d∈ [0,N-1])
          if (B2H[Prm[d]])
              {e=min{j|j>Prm[j] && (BH[j] || ! B2H[j])}
                for (x=Prm[d]+1; x<e; x++) B2H[x]=0;}
            }
for (i=0; i<N; i++) Pos[Prm[i]]=i;
for (i=0; i<N; i++) BH[i]=B2H[i];
}
```

Calcularea lcp

- ⌘ Cautarea in $O(P + \log N)$ necesita precalcularea vectorilor $Llcp$ si $Rlcp$.
- ⌘ Sa notam A_i si A_j doua sufixe din galeti consecutive.
- ⌘ Initial (cand sufixele sunt distribuite in galeti dupa primul simbol), $lcp(A_i, A_j) = 0$.
- ⌘ Dupa pasul $2H$ sufixele sunt distribuite in galeti dupa primele $2H$ simboluri. Daca A_i si A_j sunt doua sufixe care ajung in galeti adiacente la acest pas (deci erau in aceeași galeata la pasul H), $H \leq lcp(A_i, A_j) < 2H$. Mai mult, $lcp(A_i, A_j) = H + lcp(A_{i+H}, A_{j+H})$.
- ⌘ Se stie ca $lcp(A_{i+H}, A_{j+H}) < H$. Problema este A_{i+H} si A_{j+H} s-ar putea sa nu fie in galeti consecutive.

Calcularea lcp - continuare

- ⌘ Dacă $A_{Pos[i]}$ și $A_{Pos[j]}$ ($i < j$) au $lcp < H$ și Pos este \leq_H ordonat, atunci $lcp(A_{Pos[i]}, A_{Pos[j]})$ este minimul dintre lcp dintre perechile din galetii adiacente dintre $Pos[i]$ și $Pos[j]$:
- ⌘ $lcp(A_{Pos[i]}, A_{Pos[j]}) = \min\{lcp(A_{Pos[k]}, A_{Pos[k+1]}) \mid i \leq k < j\}$
- ⌘ Din păcate, utilizarea formulei precedente necesită prea mult timp. Pentru a determina lcp în $O(\log N)$ puteți utiliza un arbore de intervale.

Bibliografie



- ⌘ Udi Manber & Gene Myers - "Suffix arrays- A new method for on-line string searches."