

Permuted Longest-Common-Prefix Array

Demian Banakh

1 Notacja

- $a * b$ znaczy katenację słów a i b .
- Tekst wejściowy oznaczam $t[1..n]$
- Sufiks i znaczy $t[i..n]$
- Tablica sufiksowa dla t to $SA[1..n]$

$$t[SA[1]..n] < t[SA[2]..n] < \dots < t[SA[n]..n]$$

- Funkcja $lcp(x, y)$ zwraca najkrótszy wspólny prefiks słów x i y . Dla wygody będę oznaczał

$$lcp(i, j) := lcp(t[i..n], t[j..n])$$

- Tablica $LCP[1..n]$ zdefiniowana przez t i SA to

$$LCP[j] = lcp(SA[j - 1], SA[j])$$

- Tablica $PLCP[1..n]$ jest permutacją tablicy LCP zdefiniowaną wzorem

$$PLCP[SA[i]] = LCP[i]$$

2 Algorytmy i dowody

Kluczowa własność $PLCP$:

Lemma 2.1. $PLCP[i] \geq PLCP[i - 1] - 1$ dla każdego $i > 1$.

Dowód. Niech j, j' takie że $\text{SA}[j] = i$ oraz $\text{SA}[j'] = i - 1$. Mamy z definicji

$$\begin{aligned}\text{PLCP}[i] &= \text{LCP}[j] = \text{lcp}(\text{SA}[j - 1], i) \\ \text{PLCP}[i - 1] &= \text{LCP}[j'] = \text{lcp}(\text{SA}[j' - 1], i - 1)\end{aligned}$$

Rozważmy następujące sufiksy

$$\begin{aligned}A &= t[\text{SA}[j - 1] - 1..n] = t[\text{SA}[j - 1] - 1] * t[\text{SA}[j - 1]..n] \\ B &= t[\text{SA}[j' - 1]..n] = t[\text{SA}[j' - 1]] * t[\text{SA}[j' - 1] + 1..n] \\ C &= t[i - 1..n] = t[i - 1] * t[i..n] \\ B &\leq C \text{ (poprzedni leksykograficznie)}\end{aligned}$$

Oczywiście $\text{lcp}(\text{SA}[j' - 1] + 1, i) \geq \text{lcp}(\text{SA}[j' - 1], i - 1) - 1$, przy czym równość jest wtw gdy $t[\text{SA}[j' - 1]] = t[i - 1]$. Skoro nie ma sufiksów pomiędzy $\text{SA}[j - 1]$ a $\text{SA}[j] = i$ w porządku leksykograficznym, to

$$t[\text{SA}[j' - 1] + 1..n] \leq t[\text{SA}[j - 1]..n] < t[i..n]$$

Zatem

$$\text{lcp}(\text{SA}[j - 1], i) \geq \text{lcp}(\text{SA}[j' - 1] + 1, i) \geq \text{lcp}(\text{SA}[j' - 1], i - 1) - 1$$

□

Ten lemat prawie wprost prowadzi do wydajnego Algorytmu A obliczenia PLCP: mając obliczoną wartość $\text{PLCP}[i - 1]$, odejmujemy 1 i porównujemy sufiksy znak po znaku dopóki są równe. Oczywiście musimy wiedzieć który to jest sufiks $\text{SA}[j - 1]$; do tego celu liczymy dodatkową tablicę $\Phi[\text{SA}[j]] = \text{SA}[j - 1]$.

Algorithm 1 Algorytm A

Require: SA and text

Ensure: PLCP

```

for  $i = 1$  to  $n$  do
     $\Phi[\text{SA}[i]] = \text{SA}[i - 1]$ 
end for
 $l \leftarrow 0$ 
for  $i = 1$  to  $n$  do
     $s \leftarrow \Phi[i]$ 
    while  $\text{text}[i + l] = \text{text}[s + l]$  do
         $l \leftarrow l + 1$ 
    end while
     $\text{PLCP}[i] \leftarrow l$ 
     $l \leftarrow \max(l - 1, 0)$ 
end for
```

Łatwo widać, że złożoność czasowa Algorytmu A jest $O(n)$, bo zmniejszamy l o 1 co najwyżej n razy, a maksymalna wartość l jest co najwyżej n . Z punktu widzenia pamięci, potrzebuje dodatkowej tablicy Φ rozmiaru n .

W dość naturalny sposób można ulepszyć kontrolę nad space-time trade-off modyfikując algorytm tak, żeby produkował tylko co q -ty element PLCP (wystarczy obliczyć tylko co q -ty element Φ ; naiwne liczenie PLCP[i] zaczynamy od (PLCP[$i - 1$] - q) - wynika z Lematu 1), a obliczenie dowolnego elementu PLCP na bazie co q -tego było w czasie amortyzowanym $O(q)$ (dla dowolnego i szukamy najbliższych takich k , $k + q$ na których wartość PLCP jest znana; wtedy obliczamy naiwnie PLCP[i] w interwale zadanym przez PLCP[k] i PLCP[$k + q$], którego długość jest średnio $O(q)$).

Definition 2.1. Nazywamy wartość PLCP[i] = $lcp(i, \Phi[i])$ redukowalną, gdy $t[i - 1] = t[\Phi[i] - 1]$.

Lemma 2.2. Jeśli PLCP[i] jest redukowalna, to PLCP[i] = PLCP[$i - 1$] - 1.

Dowód. Pokażemy to, kontynuując dowód Lematu 1 z dodatkowym założeniem

$$t[i - 1] = t[SA[j - 1] - 1]$$

W takim razie mamy $A < C$, ponieważ SA[$j - 1$] i SA[j] = i to są kolejne sufiksy. Skoro B i C to też kolejne sufiksy

$$A \leq B < C$$

Skoro pierwsze litery A i C są równe, pierwsza litera B musi być taka sama. Oznaczmy A' , B' , C' odpowiednio A , B , C bez pierwszych liter. Mamy

$$A' \leq B' < C'$$

Ale A' i C' to kolejne sufiksy, więc $A' = B'$, oraz SA[$j - 1$] = SA[$j' - 1$] + 1. W końcu

$$lcp(SA[j - 1], i) = lcp(SA[j' - 1] + 1, i) = lcp(SA[j' - 1], i - 1) - 1$$

ponieważ $t[i - 1] = t[SA[j - 1] - 1] = t[SA[j' - 1]]$. □

Wnioskujemy z tego lematu, że nieredukowalne wartości PLCP[i] jednoznacznie wyznaczają pozostałe, zatem algorytm B obliczenia PLCP najpierw wylicza wszystkie nieredukowalne wartości, potem dopełnia pozostałe.

Algorithm 2 Algorytm B

Require: SA and text**Ensure:** PLCP

```
for  $i = 1$  to  $n - 1$  do
   $j \leftarrow \text{SA}[i]$ 
   $k \leftarrow \text{SA}[i + 1]$ 
  if  $\text{text}[j - 1] = \text{text}[k - 1]$  then
     $\text{PLCP}[k] \leftarrow \text{lcp}(j, k)$ 
  end if
end for
for  $i = 1$  to  $n - 1$  do
  if  $\text{PLCP}[i + 1] < \text{PLCP}[i] - 1$  then
     $\text{PLCP}[i + 1] \leftarrow \text{PLCP}[i] - 1$ 
  end if
end for
```

Lemma 2.3. Suma wszystkich nieredukowalnych wartości lcp jest $\leq 2n \log n$.

Dowód. Niech $l = \text{PLCP}[i] = \text{lcp}(i, j)$, gdzie $j = \Phi[i]$, jest nieredukowalną wartością. Zatem

$$\begin{aligned} t[i - 1] &\neq t[j - 1] \\ t[i..i + l - 1] &= t[j..j + l - 1] \\ t[i + l] &\neq t[j + l] \end{aligned}$$

Dla każdego $0 \leq k \leq l - 1$, będziemy przydzielać koszt 1 do pary $t[i + k] = t[j + k]$ w następujący sposób.

Rozpatrzmy drzewo sufiksowe słowa \bar{t} ; niech v_{i+k} i v_{j+k} będą liśćmi odpowiadającymi prefiksom $t[1..i + k]$ i $t[1..j + k]$. Najniższy wspólny przodek u węzłów v_{i+k} i v_{j+k} odpowiada $t[i..i + k]$ - te prefiksy zgadzają się na dokładnie $(k + 1)$ znaków od końca. Jeśli v_{i+k} jest w mniejszym poddrzewie u , to koszt pary $t[i + k] = t[j + k]$ przydzielamy do węzła v_{i+k} , wpp do v_{j+k} . Gdy v_{i+k} został wybrany, nazwiemy u *drogim przodkiem* węzła v_{i+k} , a v_{j+k} *drogim bratem* węzła v_{i+k} w odniesieniu do u , wpp analogicznie.

Teraz wystarczy pokazać, że każdy liść ma koszt co najwyżej $2 \log n$. W szczególności, pokażmy, że (a) każdy liść ma co najwyżej $\log n$ drogich przodków, i (b) w odniesieniu do każdego z tych drogich przodków dany liść ma co najwyżej 2 drogich braci.

- Rozważmy ścieżkę od v do korzenia. Z konstrukcji, przy każdym drogim przodku węzła v na ścieżce, poddrzewo rośnie przynajmniej w 2 razy. Zatem może być nie więcej niż $\log n$ takich przodków.

- Niech u to drogi przodek węzła v i w to drogi brat węzła v w odniesieniu do u . Mamy, że v, u, w odpowiadają

$$t[1..i+k], t[i..i+k], t[1..j+k] \text{ dla pewnych } i, j \text{ tż } i = \Phi[j] \text{ lub } j = \Phi[i]$$

Bez ograniczenia ogólności, niech $i = \Phi[j]$. Załóżmy, że istnieje inny drogi brat $w' \neq w$ węzła v w odniesieniu do u . Teraz w' musi odpowiadać $t[1..j'+k]$ dla $j' = \Phi[i]$ (ponieważ $i = \Phi[j] \neq \Phi[j']$). Od razu widać, że trzeciego takiego brata nie może istnieć.

Suma kosztów wszystkich liści w drzewie (równa sumie wszystkich nieredukowalnych wartości lcp) jest ograniczona przez $2n \log n$. \square

Z tego lematu wnioskujemy, że złożoność czasowa Algorytmu B jest $O(n \log n)$, natomiast z punktu widzenia pamięci - używa $O(1)$ poza samą tablicą PLCP.

Podobnie jak wcześniej, ten algorytm da się w naturalny sposób zmodyfikować tak, żeby produkował co q -ty element PLCP, a obliczenie dowolnego na bazie co q -tego zajmowało średnio $O(q)$.

Warto dodać, że przedstawione algorytmy obliczenia PLCP są znacznie szybsze w praktyce, niż standardowe algorytmy obliczenia LCP. Zaletą Algorytmu A jest bardzo skuteczne wykorzystanie własności *locality of reference*, drugiego - wymagania pamięciowe (jeśli przedstawić PLCP jako bit-tablicę, potencjalnie można zmieścić się $3n$ bitów pamięci).