

Algorytmy Tekstowe

Tomasz Waleń

MIM UW

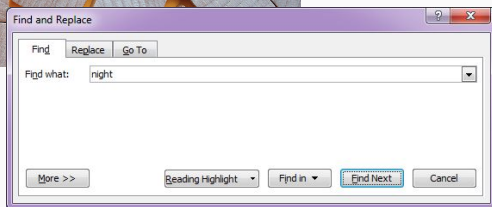
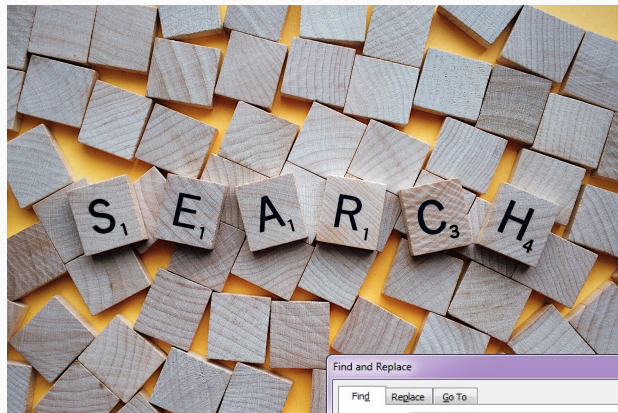
Algorytmy i Struktury Danych, 2018-12-03

slajdy: <http://bit.ly/aisd-20181203>

Licencja: Creative Commons Attribution-NonCommercial



Problem wyszukiwania wzorca



Problem wyszukiwania wzorca

Definicja

Dla tekstu T i wzorca P ($|T| = n$, $|P| = m$)
wyznacz wszystkie wystąpienia wzorca w tekście:

$Occ(P, T) = \{i : \text{wzorzec } P \text{ występuje w tekście } T \text{ na pozycji } i\}$

Problem wyszukiwania wzorca

Definicja

Dla tekstu T i wzorca P ($|T| = n$, $|P| = m$)
wyznacz wszystkie wystąpienia wzorca w tekście:

$$\text{Occ}(P, T) = \{i : T[i, \dots, i + m - 1] = P\}$$

Przykład

P : abaaba

0 1 2 3 4 5 6 7 8 9 10 12 14 16 18 20 22
 T : ababaabaababaabaabaaba
abaaba abaaba abaaba
abaaba

Occ: 2, 5, 10, 17

Algorytm 1: NaiveSearch

Dane: T – tekst, P – wzorzec

Wynik: pozycje w tekście w T na których występuje wzorzec P

- 1 $Occ := \emptyset$
 - 2 **foreach** $i \in \{0, \dots, |T| - |P|\}$ **do**
 - 3 **if** $T[i, \dots, i + |P| - 1] = P$ **then**
 - ▷ zgłoś wystąpienie P na pozycji $T[i]$
 - 4 $Occ := Occ + \{i\}$
 - 5 **return** Occ
-

Lemat

W pesymistycznym przypadku algorytm naiwny wykonuje $O(nm)$ operacji.

Dowód.

Dla $T = a^{n+1+m}$ i $P = a^m b$, algorytm naiwny stara się dopasować wzorzec na n pozycjach, niestety każde dopasowanie wymaga $m + 1$ porównań.

Stąd algorytm wykonuje $O(nm)$ operacji. □

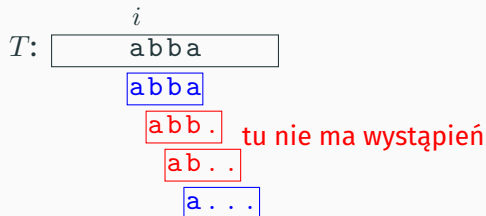
Gdzie tkwi nieefektywność algorytmu naiwnego?

Za każdym razem przesuwamy obliczenia tylko o jedną pozycję w prawo i każde kolejne dopasowanie ignoruje informację otrzymaną we wcześniejszych krokach.

Gdzie tkwi nieefektywność algorytmu naiwnego?

Za każdym razem przesuwamy obliczenia tylko o jedną pozycję w prawo i każde kolejne dopasowanie ignoruje informację otrzymaną we wcześniejszych krokach.

Przykład

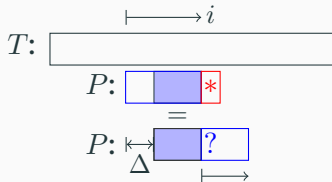


Jak mógłby wyglądać lepszy algorytm?

Założenia szybszego algorytmu

1. maksymalne wykorzystanie zdobytej wiedzy – nigdy nie cofamy się w tekście,
2. sprytniejsze kolejne dopasowania (Δ) – nie mogą być za małe (żeby algorytm nie był zbyt wolny), ani zbyt duże (żeby nie zgubić żadnego wystąpienia).

Przykład



Proste definicje

Dla słowa $w = a_0, \dots, a_{n-1}$:

- ▶ (właściwym) prefiksem słowa w nazywamy słowa postaci $\text{pref}(w, \ell) = a_0, \dots, a_{\ell-1}$, gdzie $0 \leq \ell < n$,
- ▶ (właściwym) sufiksem słowa w nazywamy słowa postaci $\text{suf}(w, \ell) = a_{n-\ell}, \dots, a_{n-1}$, gdzie $0 \leq \ell < n$,
- ▶ puste słowo ϵ jest zarówno prefiksem jak i sufiksem dowolnego niepustego słowa.

Przykład

$w = \text{alamakota}$

- ▶ prefiksy: $\epsilon, a, al, ala, \dots$
- ▶ sufiksy: $\epsilon, a, ta, ota, kota, \dots$

A może przeddefiniujemy problem?

Definicja

Dla tekstu T i wzorca P ($|T| = n$, $|P| = m$) należy wyznaczyć tablicę L długości n :

$$L[i] = \max\{j : \text{pref}(P, j) \text{ jest sufiksem } T[0, \dots, i]\}$$

A może przeddefiniujemy problem?

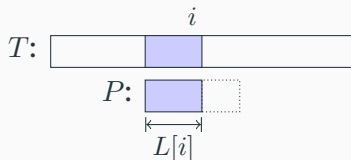
Definicja

Dla tekstu T i wzorca P ($|T| = n$, $|P| = m$) należy wyznaczyć tablicę L długości n :

$$L[i] = \max\{j : \text{pref}(P, j) \text{ jest sufiksem } T[0, \dots, i]\}$$

A prościej?

Dla każdej pozycji $T[i]$ w tekście chcemy się wiedzieć jaki jest **najdłuższy** (początkowy) fragment wzorca który kończy się na $T[i]$.



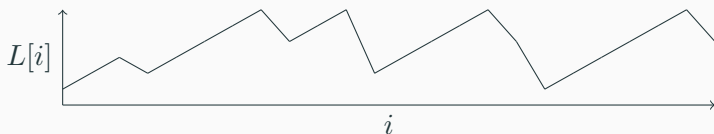
Rozszerzony problem – przykład

P : a b a a b a

T : a b a b a a b a a b a b a a b a a b a a b a a

$L[i]$: 1 2 3 2 3 4 5 **6** 4 5 **6** 2 3 4 5 **6** 4 1 2 3 4 5 **6** 4

a b a a b a



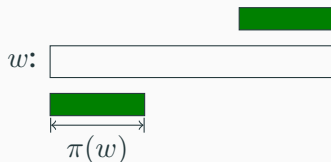
Funkcja prefikso-sufiksowa

Definicja

$$\pi(w) = \max\{\ell : \text{pref}(w, \ell) = \text{suf}(w, \ell)\}$$

Długość najdłuższego prefiksu w , który jest jednocześnie sufiksem.

Dodatkowo definiujemy $\pi(\epsilon) = \epsilon$.



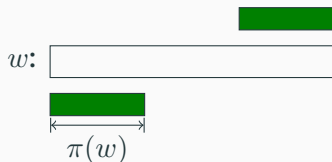
Funkcja prefikso-sufiksowa

Definicja

$$\pi(w) = \max\{\ell : \text{pref}(w, \ell) = \text{suf}(w, \ell)\}$$

Długość najdłuższego prefiksu w , który jest jednocześnie sufiksem.

Dodatkowo definiujemy $\pi(\epsilon) = \epsilon$.



Ale właściwie dlaczego?

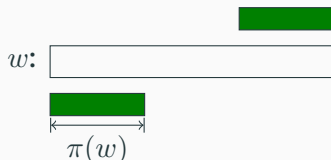
Funkcja prefikso-sufiksowa

Definicja

$$\pi(w) = \max\{\ell : \text{pref}(w, \ell) = \text{suf}(w, \ell)\}$$

Długość najdłuższego prefiksu w , który jest jednocześnie sufiksem.

Dodatkowo definiujemy $\pi(\epsilon) = \epsilon$.



Ale właściwie dlaczego?

- ta funkcja mówi o ile można przesunąć obliczenia!

Prefikso-sufiks – przykłady

w : abaaaab

ab

ab

$$\pi(w) = 2$$

w : abababa

ababa

ababa

$$\pi(w) = 5$$

w : ababcde

|

|

$$\pi(w) = 0$$

Definicja tablicy prefikso-sufiksów

A co by było, gdyby policzyć wartości funkcji π dla wszystkich prefiksów słowa?

Definicja

Dla słowa w definiujemy tablicę prefikso-sufiksów:

$$\pi_w[\ell] = \pi(\text{pref}(w, \ell))$$

Dla $0 \leq \ell \leq |w|$.

Tablica prefikso-sufiksów – przykład

Dla $w = \text{abaaba}$:

| ℓ | $\pi[\ell]$ | |
|--------|-------------|---------------------------------------|
| 0 | 0 | - |
| 1 | 0 | $\pi(\text{a}) = 0 (\epsilon)$ |
| 2 | 0 | $\pi(\text{ab}) = 0 (\epsilon)$ |
| 3 | 1 | $\pi(\text{aba}) = 1 (\text{a})$ |
| 4 | 1 | $\pi(\text{abaa}) = 1 (\text{a})$ |
| 5 | 2 | $\pi(\text{abaab}) = 2 (\text{ab})$ |
| 6 | 3 | $\pi(\text{abaaba}) = 3 (\text{aba})$ |

Tablica prefikso-sufiksów – przykład

Dla $w = \text{abaaba}$:

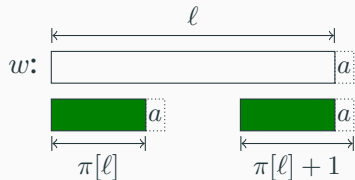
| ℓ | $\pi[\ell]$ | |
|--------|-------------|---------------------------------------|
| 0 | 0 | - |
| 1 | 0 | $\pi(\text{a}) = 0 (\epsilon)$ |
| 2 | 0 | $\pi(\text{ab}) = 0 (\epsilon)$ |
| 3 | 1 | $\pi(\text{aba}) = 1 (\text{a})$ |
| 4 | 1 | $\pi(\text{abaa}) = 1 (\text{a})$ |
| 5 | 2 | $\pi(\text{abaab}) = 2 (\text{ab})$ |
| 6 | 3 | $\pi(\text{abaaba}) = 3 (\text{aba})$ |

Ale jak efektywnie liczyć tablicę π ?

Tablica prefikso sufiksów - przypadki

Jeśli $\ell > 1$ i $w[\pi_w[\ell]] = w[\ell]$, to

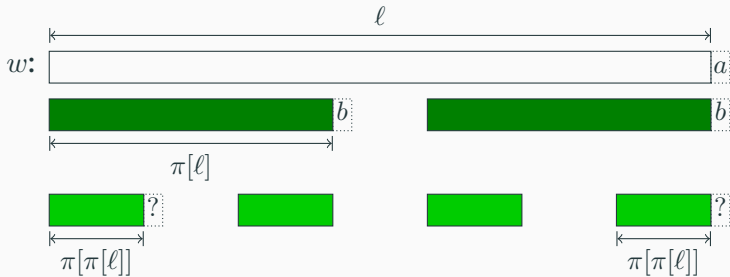
$$\pi_w[\ell + 1] = \pi_w[\ell] + 1$$



Tablica prefikso sufiksów - przypadki

Jeśli $\ell > 1$ i $w[\pi_w[\ell]] \neq w[\ell]$, to

$$\pi_w[\ell + 1] \leq \pi_w[\pi_w[\ell]]$$



Algorytm wyznaczania tablicy prefikso-sufiksów

Algorytm 2: Compute- π

Dane: P , $|P| = m$

Wynik: tablica $\pi_P[\ell]$

```
1  $j = 0$ 
2  $\pi[0] := 0$ 
3 foreach  $\ell \in \{1, \dots, m\}$  do
4   while  $j > 0$  and  $P[\ell - 1] \neq P[j]$  do
5      $j := \pi[j]$  ▷ zmniejsz wartość  $j$ 
6   if  $j < \ell - 1$  and  $P[\ell - 1] = P[j]$  then
7      $j := j + 1$ 
8    $\pi[\ell] := j$ 
9 return  $\pi$ 
```

Lemat

Złożoność czasowa algorytmu $\text{Compute-}\pi$ wynosi $O(m)$.

Dowód.

Zauważmy, że licznik j może zostać zwiększony co najwyżej m razy (tylko raz dla każdej wartości ℓ). Każda iteracja pętli **while** zmniejsza licznik (co najmniej o 1). Stąd sumaryczna liczba iteracji pętli **while** nie może przekroczyć m . \square

Algorytm KMP

Algorytm 3: KMP

Dane: tekst T , wzorzec P , $|T| = n$, $|P| = m$

Wynik: tablica wystąpień Occ

- 1 $\pi = \text{Compute-}\pi(P)$
 - 2 $\ell = 0$
 - 3 $Occ := \emptyset$
 - 4 **foreach** $i \in \{0, \dots, n - 1\}$ **do**
 - 5 **while** $\ell > 0$ **and** $T[i] \neq P[\ell]$ **do**
 - 6 $\ell := \pi[\ell]$
 - 7 **if** $T[i] = P[\ell]$ **then**
 - 8 $\ell = \ell + 1$;
 - 9 **if** $\ell = m$ **then**
 - ▷ znaleziono wystąpienie kończące się na $T[i]$
 - 10 $Occ = Occ + \{i - m + 1\}$
 - 11 $\ell = \pi[\ell]$
 - 12 **return** Occ
-

Lemat

Złożoność czasowa algorytmu KMP wynosi $O(n + m)$.

Dowód.

Analogicznie jak w poprzednim dowodzie zauważamy, że licznik ℓ może zostać zwiększony co najwyżej n razy. Każda iteracja pętli **while** zmniejsza licznik (co najmniej o 1). Stąd sumaryczna liczba iteracji pętli **while** nie może przekroczyć n . □

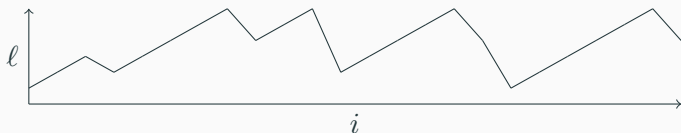
Algorytm KMP – przykład

P : a b a a b a

T : a b a b a a b a a b a b a a b a a a b a a b a a

l : 1 2 3 2 3 4 5 **6** 4 5 **6** 2 3 4 5 **6** 4 1 2 3 4 5 **6** 4

a b a a b a



Problem wyszukiwania wzorca – ciekawostki

- ▶ algorytm KMP jest bardzo tolerancyjny jeśli chodzi o alfabet (może być dowolny, nawet nie musi być sortowalny) i definicję równości słów – można go zaadoptować nawet do bardzo egzotycznych problemów,
- ▶ w praktycznych zastosowaniach zamiast KMP używa się szybszych algorytmów, np. `grep` używa wariantu algorytmu Boyera-Moora,
- ▶ problem wyszukiwania wzorca można rozwiązać w czasie $O(n + m)$ i dodatkowej pamięci $O(1)$,
- ▶ dzięki rozszerzeniom SIMD w nowoczesnych procesorach, wiele standardowych algorytmów można zaimplementować znacznie efektywniej (zwłaszcza dla krótkich wzorców i małych alfabetów)



Definicja

Dla tekstu T należy przygotować strukturę danych, która umożliwia zapytania:

$OccQuery(P)$: lista wystąpień wzorca P w tekście T .

Definicja

Dla tekstu T , *tablica sufiksowa* (Suffix Array / SA) zawiera uporządkowanie leksykograficzne wszystkich sufiksów.

$$SA[i] = \{j : \text{t. że } T[j, \dots] \text{ jest } i\text{-tym co do rangi sufiksem}\}$$

Zauważmy, że dzięki uporządkowaniu, pozycje w tekście zaczynające się od podobnych słów będą w tablicy sufiksowej blisko siebie.

Tablice sufiksowe pełnią podobną rolę, jak drzewa sufiksowe. Często stosuje się je zamiennie.

Przykład

0 1 2 3 4 5 6 7 8 9 10
T: mississippi

| <i>p</i> | <i>T</i> [<i>p</i> ,...] | <i>i</i> | <i>SA</i> [<i>i</i>] | <i>T</i> [<i>SA</i> [<i>i</i>],...] |
|----------|---------------------------|----------|------------------------|--|
| 0 | mississippi | 0 | 10 | i |
| 1 | ississippi | 1 | 7 | ippi |
| 2 | ssissippi | 2 | 4 | issippi |
| 3 | sissippi | 3 | 1 | ississippi |
| 4 | issippi | 4 | 0 | mississippi |
| 5 | ssippi | 5 | 9 | pi |
| 6 | sippi | 6 | 8 | ppi |
| 7 | ippi | 7 | 6 | sippi |
| 8 | ppi | 8 | 3 | sissippi |
| 9 | pi | 9 | 5 | ssippi |
| 10 | i | 10 | 2 | ssissippi |

Zastosowania tablicy sufiksowej

- ▶ wyszukanie wzorca $Occ(P)$ – w czasie $O(|P| + |Occ| + \log |T|)$,
- ▶ zliczanie liczby różnych podstów – w czasie $O(|T|)$,
- ▶ najdłuższe wspólne podstowo – w czasie $O(|T_1| + |T_2|)$,
- ▶ LZ-faktoryzacja (używana do kompresji ZIP),
- ▶ rozwiązanie jednego z zadań laboratoryjnych 😊
- ▶ i wiele wiele innych.

Zastosowania tablicy sufiksowej – ale jak?

Niestety potrzebujemy dodatkowych struktur danych:

- ▶ tablicy LCP,
- ▶ tablicy SA^{-1} ,
- ▶ często również struktury RMQ.

Definicja

$\text{lcp}(X, Y)$ = długość najdłuższego wspólnego prefiksu X i Y

X :

| | |
|--|-----|
| | x |
|--|-----|

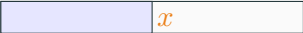
Y :


| | |
|--|-----|
| | y |
|--|-----|

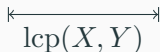
←————→
 $\text{lcp}(X, Y)$

Definicja

$\text{lcp}(X, Y)$ = długość najdłuższego wspólnego prefiksu X i Y

X : 

Y : 


 $\text{lcp}(X, Y)$

Przykład

$\text{lcp}(\text{alamakota}, \text{alamapsa}) = 5$

Definicja

Tablica LCP zawiera informacje o wartościach lcp dla pomiędzy sąsiednimi sufiksami w tablicy sufiksowej.

$$\text{LCP}[i] = \text{lcp}(T[\text{SA}[i - 1], \dots], T[\text{SA}[i], \dots])$$

Dla $i = 0$, $\text{LCP}[i]$ nie jest zdefiniowane.

Tablica LCP – przykład

| i | $SA[i]$ | $T[SA[i], \dots]$ | $LCP[i]$ |
|-----|---------|-------------------|----------|
| 0 | 10 | i | |
| 1 | 7 | ippi | 1 |
| 2 | 4 | issippi | 1 |
| 3 | 1 | ississippi | 4 |
| 4 | 0 | mississippi | 0 |
| 5 | 9 | pi | 0 |
| 6 | 8 | ppi | 1 |
| 7 | 6 | sippi | 0 |
| 8 | 3 | sissippi | 2 |
| 9 | 5 | ssippi | 1 |
| 10 | 2 | ssissippi | 3 |

SA^{-1} – funkcja odwrotna do SA

$$SA^{-1}[p] = \text{RANK}[p] = i \text{ wtw } SA[i] = p$$

| | i | $SA[i]$ | $T[SA[i], \dots]$ |
|----|-----|---------|-------------------|
| 0 | 0 | 10 | i |
| 1 | 1 | 7 | ippi |
| 2 | 2 | 4 | issippi |
| 3 | 3 | 1 | issippi |
| 4 | 4 | 0 | mississippi |
| 5 | 5 | 9 | pi |
| 6 | 6 | 8 | ppi |
| 7 | 7 | 6 | sippi |
| 8 | 8 | 3 | sissippi |
| 9 | 9 | 5 | ssippi |
| 10 | 10 | 2 | ssissippi |

$T: \text{mississippi}$

$SA^{-1}[2] = 10$

Ale jak efektywnie wyznaczyć SA?

Zacznijmy od budżetu i rachunków:

- ▶ naszym celem jest złożoność $O(n)$,
- ▶ idealnie byłoby gdyby algorytm nie był zbyt zależny od alfabetu (ale zakładamy, że alfabet jest sortowalny w $O(n)$),
- ▶ jeśli koncentrujemy się na algorytmach rekurencyjnych to następująca funkcja wygląda interesująco:

$$T(n) = T\left(\frac{2n}{3}\right) + O(n)$$

- ▶ oczywiście $T(n) = O(n)$.

Juha Kärkkäinen, Peter Sanders,

“Simple linear work suffix array construction”

ICALP 2003

~530 cytowań + ~350 (wersja czasopismowa: J. of ACM 2006)

Ciekawostka:

praca zawiera implementację algorytmu (~50 wierszy w C)

Algorytm 4: ComputeSA

Dane: słowo T , $|T| = n$

Wynik: tablica sufiksowa SA słowa T

- 1 podziel sufiksy słowa T na trzy kategorie S_0, S_1, S_2
 - ▷ $|S_i| \approx n/3$
 - 2 wygeneruj słowo X reprezentujące sufiksy z S_1 i S_2
 - ▷ $|X| \approx 2n/3$
 - 3 $SA_X = \text{ComputeSA}(X)$
 - 4 $SA_{1,2} =$ wyznacz z SA_X kolejność leks. sufiksów z S_1 i S_2
 - 5 $SA_0 =$ wyznacz z SA_X i T kolejność leks. sufiksów z S_0
 - 6 $SA = \text{scal } SA_{1,2} \text{ i } SA_0$ ▷ ta część jest dosyć sprytna
 - 7 **return** SA
-

Lemat

Algorytm KS ma złożoność czasową $O(n)$

Dowód.

Zakładając, że uda nam się wykonać wszystkie kroki związane z obliczaniem słowa X , tablic $SA_{1,2}$ i SA_0 oraz ich scalaniem, w czasie $O(n)$. Otrzymujemy następujące równanie rekurencyjne:

$$T(n) = O(n) + T(2n/3)$$

co rzeczywiście, daje nam rozwiązanie $T(n) = O(n)$. □

Drobna uwaga techniczna

W zależności od wartości $|T| \bmod 3$, uzupełniamy T o 1, 2 lub 3 znaki \$ (zakładamy, że znak \$ jest mniejszy od dowolnego innego).

Podział sufiksów na S_0, S_1 i S_2

Drobna uwaga techniczna

W zależności od wartości $|T| \bmod 3$, uzupełniamy T o 1, 2 lub 3 znaki \$ (zakładamy, że znak \$ jest mniejszy od dowolnego innego).

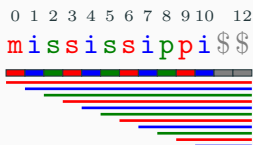
Zazwyczaj najprostsze pomysły są najlepsze, więc po prostu definiujemy:

$$S_j = \{i : i < |T| \text{ oraz } i \bmod 3 = j\}$$

$$S_0 = \{0, 3, 6, 9\}$$


$$S_1 = \{1, 4, 7, 10\}$$

$$S_2 = \{2, 5, 8\}$$



Podział sufiksów na S_0 , S_1 i S_2

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$



$$S_0 = \{0, 3, 6, 9\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$




$$S_1 = \{1, 4, 7, 10\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$



$$S_2 = \{2, 5, 8\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$



Algorytm 5: ComputeSA

Dane: słowo T , $|T| = n$

Wynik: tablica sufiksowa SA słowa T

- 1 podziel sufiksy słowa T na trzy kategorie S_0, S_1, S_2
 - ▷ $|S_i| \approx n/3$
 - 2 wygeneruj słowo X reprezentujące sufiksy z S_1 i S_2
 - ▷ $|X| \approx 2n/3$
 - 3 $SA_X = \text{ComputeSA}(X)$
 - 4 $SA_{1,2} =$ wyznacz z SA_X kolejność leks. sufiksów z S_1 i S_2
 - 5 $SA_0 =$ wyznacz z SA_X i T kolejność leks. sufiksów z S_0
 - 6 $SA =$ scal $SA_{1,2}$ i SA_0
 - 7 **return** SA
-

Słowo X reprezentujące sufiksy S_1 i S_2

$$S_1 = \{1, 4, 7, 10\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$

$$S_2 = \{2, 5, 8\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$

X' : 0 1 2 3 4 5 6 7 8 9 10 12 14 16 18 20

| | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|----|----|---|---|---|---|---|---|---|---|---|
| i | s | s | i | s | i | p | p | i | \$ | \$ | s | s | i | s | s | i | p | p | i |
|---|---|---|---|---|---|---|---|---|----|----|---|---|---|---|---|---|---|---|---|

Słowo X reprezentujące sufiksy S_1 i S_2

$$S_1 = \{1, 4, 7, 10\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$

$$S_2 = \{2, 5, 8\}$$

0 1 2 3 4 5 6 7 8 9 10 12
mississippi\$\$

X' : 0 1 2 3 4 5 6 7 8 9 10 12 14 16 18 20

| | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|----|---|---|---|---|---|---|---|---|---|
| i | s | s | i | s | s | i | p | p | i | \$ | \$ | s | s | i | s | s | i | p | p | i |
|---|---|---|---|---|---|---|---|---|---|----|----|---|---|---|---|---|---|---|---|---|

Niestety słowo X' jest zbyt długie (ma $\approx 2n$ znaków zamiast oczekiwanych $2n/3$).

Słowo X reprezentujące sufiksy S_1 i S_2

A gdyby skompresować X' ?

Słowo X reprezentujące sufiksy S_1 i S_2

A gdyby skompresować X' ?

X' :

| | | | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 12 | 14 | 16 | 18 | 20 | | | | | |
| i | s | s | i | s | s | i | p | p | i | \$ | \$ | s | s | i | s | s | i | p | p | i |

| | | | | | | | | | |
|-------|----|-----|----|-----|----|-----|----|-----|----|
| i\$\$ | →0 | ipp | →1 | iss | →2 | ppi | →3 | ssi | →4 |
|-------|----|-----|----|-----|----|-----|----|-----|----|

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| 2 | 2 | 1 | 0 | 4 | 4 | 3 |

Sukces!

Algorytm 6: ComputeSA

Dane: słowo T , $|T| = n$

Wynik: tablica sufiksowa SA słowa T

- 1 podziel sufiksy słowa T na trzy kategorie S_0, S_1, S_2
 - ▷ $|S_i| \approx n/3$
 - 2 wygeneruj słowo X reprezentujące sufiksy z S_1 i S_2
 - ▷ $|X| \approx 2n/3$
 - 3 $SA_X = \text{ComputeSA}(X)$
 - 4 $SA_{1,2} =$ wyznacz z SA_X kolejność leks. sufiksów z S_1 i S_2
 - 5 $SA_0 =$ wyznacz z SA_X i T kolejność leks. sufiksów z S_0
 - 6 $SA =$ scal $SA_{1,2}$ i SA_0
 - 7 **return** SA
-

Tablica $SA_{1,2}$

Warto zauważyć, że cała procedura przejścia z T do X jest “odwracalna” i zachowuje porządek sufiksów.

Dla dowolnego sufiksu z S_1 i S_2 możemy wskazać odpowiadający mu sufiks w X .

$\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ X: & 2 & 2 & 1 & 0 & 4 & 4 & 3 \end{matrix}$

 $\begin{matrix} 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 12 \\ T: & m & i & s & s & i & s & s & i & p & p & i & \$ & \$ \end{matrix}$

| i | $SA_X[i]$ | suf. w X | suf. w S | $SA_{1,2}[i]$ |
|-----|-----------|------------|----------------|---------------|
| 0 | 3 | 0443 | i\$\$ | 10 |
| 1 | 2 | 10443 | ippi\$\$ | 7 |
| 2 | 1 | 210443 | issippi\$\$ | 4 |
| 3 | 0 | 2210443 | ississippi\$\$ | 1 |
| 4 | 6 | 3 | ppi\$\$ | 8 |
| 5 | 5 | 43 | ssippi\$\$ | 5 |
| 6 | 4 | 443 | ssissippi\$\$ | 2 |

Algorytm 7: ComputeSA

Dane: słowo T , $|T| = n$

Wynik: tablica sufiksowa SA słowa T

- 1 podziel sufiksy słowa T na trzy kategorie S_0, S_1, S_2
 - ▷ $|S_i| \approx n/3$
 - 2 wygeneruj słowo X reprezentujące sufiksy z S_1 i S_2
 - ▷ $|X| \approx 2n/3$
 - 3 $SA_X = \text{ComputeSA}(X)$
 - 4 $SA_{1,2} =$ wyznacz z SA_X kolejność leks. sufiksów z S_1 i S_2
 - 5 $SA_0 =$ wyznacz z SA_X i T kolejność leks. sufiksów z S_0
 - 6 $SA =$ scal $SA_{1,2}$ i SA_0
 - 7 **return** SA
-

Niestety nie jest możliwe kolejne wywołanie rekurencyjne, żeby obliczyć SA_0 . Stąd musimy wykorzystać informację z $SA_{1,2}$ w celu obliczenia SA_0 .

Każdy sufiks i z SA_0 można opisać jako:

$$(T[i], T[i + 1, \dots])$$

Niestety nie jest możliwe kolejne wywołanie rekurencyjne, żeby obliczyć SA_0 . Stąd musimy wykorzystać informację z $SA_{1,2}$ w celu obliczenia SA_0 .

Każdy sufiks i z SA_0 można opisać jako:

$$(T[i], T[i + 1, \dots])$$

Ponieważ $i + 1 \in S_1$, znamy rangę tego sufiksu:

$$(T[i], SA_{1,2}^{-1}[i + 1])$$

Tak wygenerowane pary zajmują mało miejsca i co najważniejsze możemy je posortować w czasie liniowym. Uporządkowane pary dadzą nam kolejność SA_0 .

Algorytm 8: ComputeSA

Dane: słowo T , $|T| = n$

Wynik: tablica sufiksowa SA słowa T

- 1 podziel sufiksy słowa T na trzy kategorie S_0, S_1, S_2
 - ▷ $|S_i| \approx n/3$
 - 2 wygeneruj słowo X reprezentujące sufiksy z S_1 i S_2
 - ▷ $|X| \approx 2n/3$
 - 3 $SA_X = \text{ComputeSA}(X)$
 - 4 $SA_{1,2} =$ wyznacz z SA_X kolejność leks. sufiksów z S_1 i S_2
 - 5 $SA_0 =$ wyznacz z SA_X i T kolejność leks. sufiksów z S_0
 - 6 $SA = \text{scal } SA_{1,2} \text{ i } SA_0$
 - 7 **return** SA
-

Scalanie SA_0 i $SA_{1,2}$

Został nam już tylko jeden krok — scalanie SA_0 i $SA_{1,2}$.

| i | $SA_0[i]$ | suf. w T |
|-----|-----------|-----------------|
| 0 | 0 | mississippi\$\$ |
| 1 | 9 | pi\$\$ |
| 2 | 6 | sippi\$\$ |
| 3 | 3 | sissippi\$\$ |

| i | $SA_{1,2}[i]$ | suf. w T |
|-----|---------------|----------------|
| 0 | 10 | i\$\$ |
| 1 | 7 | ippi\$\$ |
| 2 | 4 | issippi\$\$ |
| 3 | 1 | ississippi\$\$ |
| 4 | 8 | ppi\$\$ |
| 5 | 5 | ssippi\$\$ |
| 6 | 2 | ssissippi\$\$ |

Żeby scalić te dwa uporządkowane (leksykograficznie) podzbiory sufiksów musimy pokazać, że możemy w czasie $O(1)$ porównywać sufiksy $i \in SA_0$ i $j \in SA_{1,2}$.

Żeby scalić te dwa uporządkowane (leksykograficznie) podzbiory sufiksów musimy pokazać, że możemy w czasie $O(1)$ porównywać sufiksy $i \in SA_0$ i $j \in SA_{1,2}$.

Problem polega na tym, że te sufiksy nie są z tych samych światów:

`cmp(jabłko, ziemniak)`

$$\text{cmp}(i \in S_0, j \in S_1 \cup S_2)$$

Mamy dwa przypadki:

- ▶ $j \in S_1$: porównujemy
 $(T[i], SA_{1,2}^{-1}[i + 1])$ i $(T[j], SA_{1,2}^{-1}[j + 1])$
 $(i + 1 \in S_1, j + 1 \in S_2)$

$$\text{cmp}(i \in S_0, j \in S_1 \cup S_2)$$

Mamy dwa przypadki:

- ▶ $j \in S_1$: porównujemy
 $(T[i], SA_{1,2}^{-1}[i+1])$ i $(T[j], SA_{1,2}^{-1}[j+1])$
 $(i+1 \in S_1, j+1 \in S_2)$

- ▶ $j \in S_2$: porównujemy
 $(T[i], T[i+1], SA_{1,2}^{-1}[i+2])$ i $(T[j], T[j+1], SA_{1,2}^{-1}[j+2])$
 $(i+2 \in S_2, j+2 \in S_1)$

Porównanie – przykład

$$\text{cmp}(\text{sippi}\$\$ \in S_0, \text{ssippi}\$\$ \in S_2)$$

redukujemy do:

$$\text{cmp} \left((s, i, SA_{1,2}^{-1}[\text{ppi}\$\$]), (s, s, SA_{1,2}^{-1}[\text{ippi}\$\$]) \right)$$

| i | $SA_{1,2}[i]$ | suf. w T |
|-----|---------------|----------------|
| 0 | 10 | i\$\$ |
| 1 | 7 | ippi\$\$ |
| 2 | 4 | issippi\$\$ |
| 3 | 1 | ississippi\$\$ |
| 4 | 8 | ppi\$\$ |
| 5 | 5 | ssippi\$\$ |
| 6 | 2 | ssissippi\$\$ |

Dziękuję za uwagę!

Byłbym bardzo wdzięczny za wypełnienie ankiety i uwagi: <http://bit.ly/aisd-1203-ankieta>

Przydatne artykuły:

- ▶ <https://www.youtube.com/watch?v=NinWEPPrkDQ> — wykład E. Demaine dotyczący algorytmów tekstowych
- ▶ <https://www.cs.cmu.edu/~ckingsf/bioinfo-lectures/suffixarrays.pdf> — bardzo dobre slajdy o tablicach sufiksowych
- ▶ <https://www.dmi.unict.it/~faro/papers/conference/faro34.pdf>
— Fast Packed String Matching for Short Patterns
- ▶ https://en.wikipedia.org/wiki/Suffix_array