

# Datové struktury I

## 10. přednáška: Sufixová pole

Jirka Fink

<https://ktiml.mff.cuni.cz/~fink/>

Katedra teoretické informatiky a matematické logiky  
Matematicko-fyzikální fakulta  
Univerzita Karlova v Praze

Zimní semestr 2024/25

Licence: Creative Commons BY-NC-SA 4.0

## Hledání jehly v kupce sena

- Máme daný velmi dlouhý text (seno) a krátký text (jehla)
- Úkolem je najít všechny výskyty jehly v kupce sena
- Příklad: v seně bananas se jehla ana vyskytuje hned dvakrát
- V seně anna se tatáž jehla nevyskytuje vůbec

## Algoritmy na hledání jedné jehly ve velmi velkém senu

- Knuth, Morris, Pratt
- Aho, Corasick: Více předem daných jehel
- Robin, Karp: Randomizovaný algoritmus

## Opačný přístup

Máme pevně dané seno a chceme hledat různé jehly.

## Abeceda

- $\Sigma$  je konečná množina znaků (písmen)
- $\Sigma^*$  je množina konečných posloupností (slov, řetězců) nad  $\Sigma$
- $\varepsilon$  je speciální prázdné slovo
- $\$$  je speciální znak značící konec slova

## Pro slovo $\alpha \in \Sigma$ značíme ①

- $|\alpha|$  délku  $\alpha$
- $\alpha[k]$  je  $k$ -tý znak slova  $\alpha$ , počítáno od 0 do  $|\alpha| - 1$
- $\alpha[k : l]$  je podslovo začínající  $k$ -tým znakem a končící těsně před  $l$ -tým
- $\alpha[: l]$  je prefix prvních  $l$  znaků
- $\alpha[k : ]$  je sufix začínající  $k$ -tým znakem

## Motivace práce se všemi sufixy

Výskyt jehly znamená, že nějaký sufix sena začíná jehlou

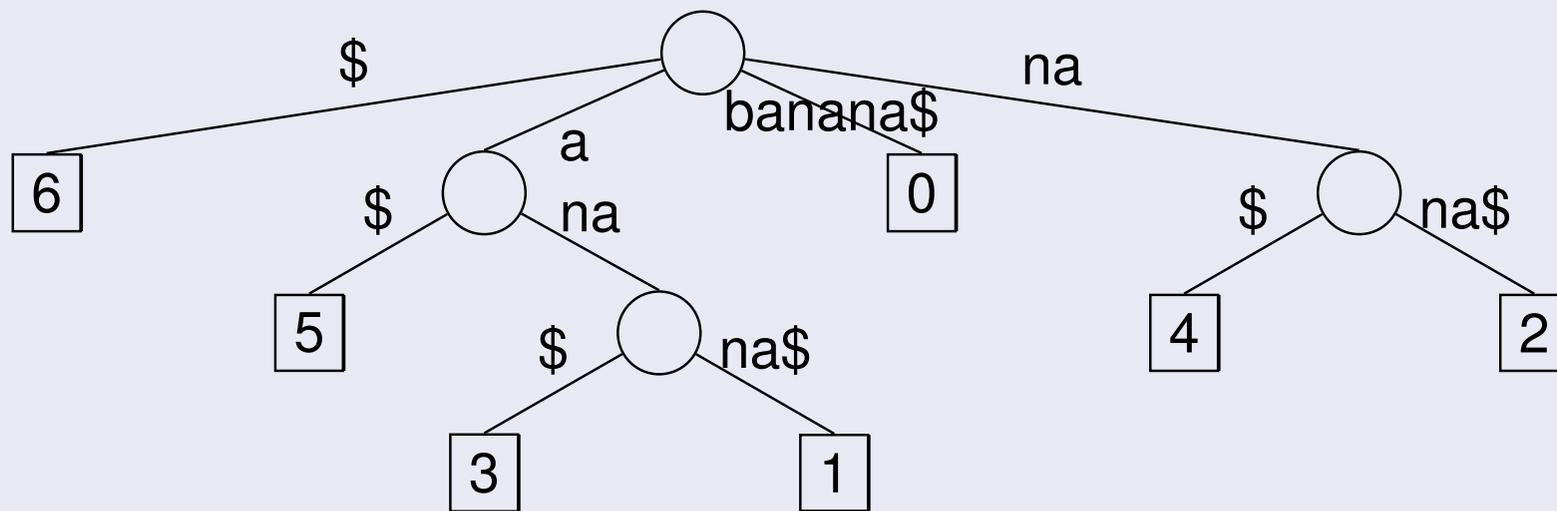
1 Značení je podobné jako v Python.

# Sufixový strom

## Definice

(Komprimovaný) sufixový strom je (komprimovaná) trie obsahující všechny sufixy, kde synové vrcholů jsou lexikografickém pořadí. ①

## Komprimovaný sufixový strom pro slovo banana\$ ②



Čísla v listech značí počáteční pozici sufixu.

## Dotazy

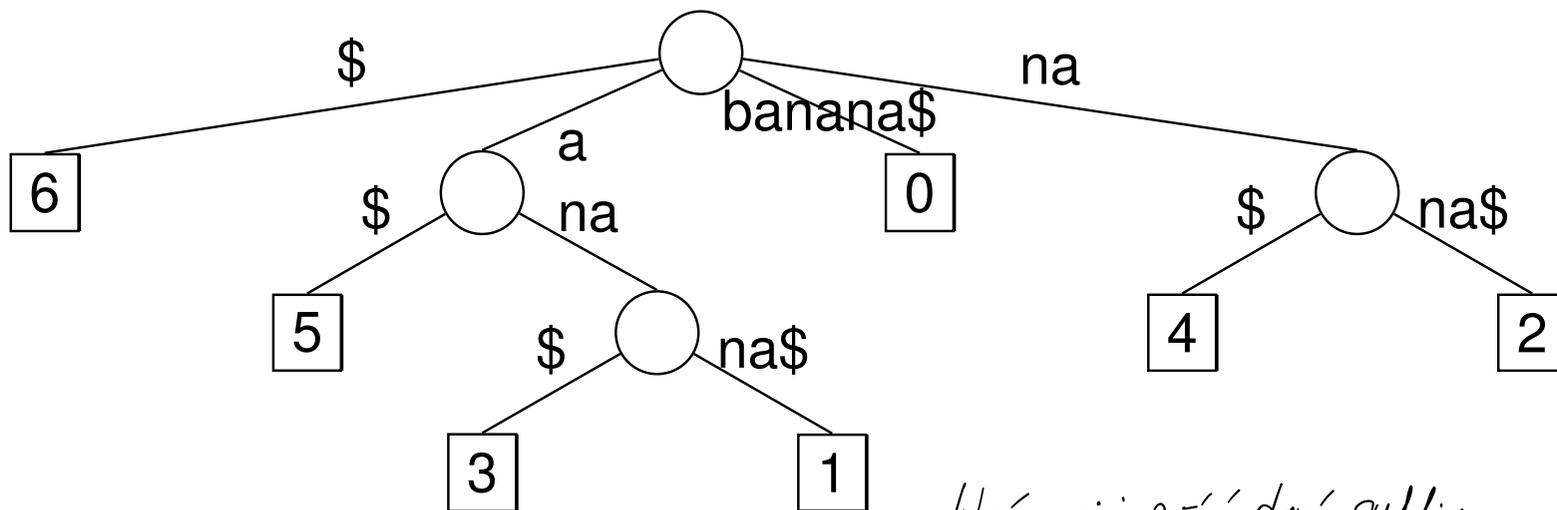
- Nalezení všech výskytů jehly v kupce sena ③
- Nalezení nejdelšího podslova s dvěma (více) výskyty ④
- Nalezení nejdelšího společného podslova dvou slov ⑤

- 1 Hloubkou vrcholu v komprimované trii rozumíme hloubkou odpovídajícího vrcholu v nekomprimované trii, tj. počet písmen na cestě z kořene do vrcholu.
- 2 K uložení komprimovaného sufixového stromu potřebujeme  $\mathcal{O}(|\alpha|)$  paměti, protože každému listu odpovídá sufix a vnitřních vrcholů je méně než listů.
- 3 Vyhledáme jehlu v sufixovém stromu a listy v podstromu udávají všechny výskyty. Chceme-li znát jen četnost, tak si stačí navíc ve všech vrcholech pamatovat počet listů v podstromu.
- 4 Vyhledáme nejhlubší vnitřní vrchol, kde hloubkou rozumíme počet písmen na cestě z kořene do vrcholu.
- 5 Pro slova  $\alpha, \beta \in \Sigma^*$  vytvoříme sufixový strom pro slovo  $\alpha\#\beta$ , kde hledáme nejhlubší vrchol obsahující v podstromu jak sufix obsahující  $\#$ , tak i neobsahující.

- Sufixové pole  $X$  udává lexikografické pořadí sufixů daného slova  $\alpha$ .  $X[i]$  říká, kde v řetězci začíná  $i$ -tý sufix v lexikografickém pořadí.
- Rankové pole  $R$  je inverzní k  $X$ , takže platí  $X[R[i]] = i$ . Hodnota  $R[i]$  říká, kolikátý v lexikografickém pořadí je sufix  $\alpha[i :]$ .
- $LCP(\alpha, \beta)$  udává délku nejdelšího společného prefixu  $\alpha$  a  $\beta$ . ①
- Pole společných prefixů (LCP)  $L$  udává délku společného prefixu mají sufixy sousedící v lexikografickém pořadí.  
Tedy  $L[i] = LCP(\alpha[X[i] :], \alpha[X[i + 1] :])$ .
- Lexikografický následník sufixu začínajícího na pozici  $i$  začíná na pozici  $X[R[i] + 1]$ .

① LCP = longest common prefix

# Souvislost sufixového stromu a pole



na které pozici začíná daný suffix

Stací aboustrnavat jen jedno pole, protože druhé už máme v lin. čase převoditelné.

i	X[i]	R[i]	L[i]	sufix
0	6	4	0	$\epsilon$
1	5	3	1	a
2	3	6	3	ana
3	1	2	0	anana
4	0	5	0	banana
5	4	1	2	na
6	2	0		nana

největší společný prefix i-tého a i+1-tého suffixu

konstrukce  $L[i]$  je hloubka nejbližšího společného předka.

index OutOfRange

## Pozorování

- Sufixové pole udává DFS pořadí listů sufixového stromu
- LCP pole  $L[i]$  udává hloubku nejbližšího společného předka listů  $X[i]$  a  $X[i + 1]$

# Sufixové pole: Příklad pro slovo barokoarokoko

$i$	$X[i]$	$R[i]$	$L[i]$	sufix
0	13	3	0	$\epsilon$
1	1	1	5	arokoarokoko
2	6	12	0	arokoko
3	0	10	0	barokoarokoko
4	11	5	2	ko
5	4	8	2	koarokoko
6	9	2	0	koko
7	12	13	1	o
8	5	11	1	oarokoko
9	10	6	3	oko
10	3	9	3	okoarokoko
11	8	4	0	okoko
12	2	7	4	rokoarokoko
13	7	0	–	rokoko

- Sufixové pole  $X[i]$  říká, kde v řetězci začíná  $i$ -tý sufix v lexikografickém pořadí.
- Rankové pole  $R[i]$  říká, kolikátý v lexikografickém pořadí je sufix  $\alpha[X[i] :]$ .
- Pole společných prefixů  $LCP\ L[i] = LCP(\alpha[X[i] :], \alpha[X[i + 1] :])$ .

## Motivace sufixového pole

- Sufixové pole potřebuje méně paměti

## Cíle

- Vytvořit pole R z X a opačně ①
- Vytvořit sufixový strom z polí X a L a opačně
- Vytvořit pole X a L
- Upravit algoritmy hledání v textu, aby používali pole místo stromu

① Tento krok je triviální, protože jen vytváříme inverzní permutaci.

## Pozorování

- Sufixové pole udává DFS pořadí listů sufixového stromu
- LCP pole  $L[i]$  udává hloubku nejbližšího společného předka listů  $X[i]$  a  $X[i + 1]$

## Konstrukce sufixového pole a LCP ze stromu

Obě pole vytváříme při průchodu stromu do hloubky

## Konstrukce sufixového stromu z pole a LCP

Strom vytváříme průchodem do hloubky

- Nejprve vytvoříme list pro slovo  $\varepsilon$
- Po vytvoření listu pro sufix  $X[i]$ 
  - Vynoříme se do vrcholu v hloubce  $L[i]$
  - Přidáme list pro sufix  $X[i]$  v hloubce  $|\text{seno}| - X[i]$

## Časová složitost

Složitost těchto převodů je lineární pro komprimovaný sufixový strom

# Vytvoření LCP pole ze sufixového pole: Triviální postup

## Algoritmus

```
1 for  $i = 0, \dots, |\alpha| - 1$  do ↳ kde je následník
2    $j = X[R[i] + 1]$ 
3    $l = 0$  ↳ dokud jsem ve slově
4   while  $i + l < |\alpha| \ \&\& \ j + l < |\alpha| \ \&\& \ \alpha[i + l] == \alpha[j + l]$  do ↳ shoda znaků
5      $l = l + 1$ 
6    $L[R[i]] = l$ 
```

*Pro vstup „a\*“ strávíme kvadratický čas*

Časová složitost je  $O(|\alpha|^2)$

Algoritmus zrychlíme tak, že si všimneme, že mnoho porovnání provádíme opakovaně.

## Pozorování

Pro každé slovo  $\alpha$  a pro všechna  $i = 0, \dots, |\alpha| - 1$  platí  $L[R[i + 1]] \geq L[R[i]] - 1$ .

*Ten prefix musí být alespoň o jeden delší jako předchozí - 1*

## Důkaz

- Zřejmé pro  $L[R[i]] \leq 1$
- Zřejmé pro  $X[R[i + 1] + 1] = X[R[i] + 1] + 1$  ①
- Jinak platí  $\alpha[i + 1 :] < \alpha[X[R[i + 1] + 1] :] < \alpha[X[R[i] + 1] :]$  ②
- Platí

$$\begin{aligned} L[R[i + 1]] &= LCP(\alpha[i + 1 :], \alpha[X[R[i + 1] + 1] :]) \\ &\geq LCP(\alpha[i + 1 :], \alpha[X[R[i] + 1] :]) \\ &= LCP(\alpha[i :], \alpha[X[R[i]] :]) - 1 \\ &= L[R[i]] - 1 \end{aligned}$$

*suffix, začínající na  $i+1$*   
*jeho následovník*

- 1  $i$  je pozice ve slovu  $\alpha$ ,  $R[i]$  udává pořadí sufixu na pozici  $i$ ,  $X[R[i] + 1]$  udává pozici následníka v lexikografickém uspořádání
- 2 Sufix začínající na pozici  $X[R[[i + 1] + 1]]$  je za sufixem začínající na pozici  $i + 1$ , ale před sufixem na pozici  $X[R[[i] + 1]]$

# Vytvoření LCP pole ze sufixového pole (Kasai, 2001)

## Pozorování

Pro každé slovo  $\alpha$  a pro všechna  $i = 0, \dots, |\alpha| - 1$  platí  $L[R[i + 1]] \geq L[R[i]] - 1$ .

## Algoritmus

```
1 l = 0
2 for i = 0, ..., |α| - 1 do
3     l = max(0, l - 1)
4     j = X[R[i] + 1]
5     while i + l < |α| && j + l < |α| && α[i + l] == α[j + l] do
6         l = l + 1
7     L[R[i]] = l
```

*vždy jde nejvíce jen o luku zpět oproti předchozí iteraci*  
*- díky lemmatu z předchozích slides*

*- vnitřní cyklus se nemohl provést víc jak 2n-krát.*

## Časová složitost je $\mathcal{O}(|\alpha|)$

- Vnější cyklus má složitost  $\mathcal{O}(|\alpha|)$
- Ve vnitřním cyklu se  $l$  vždy zvyšuje o 1
  - Hodnota  $l$  začíná na 0
  - Hodnota  $l$  na konci je nejvýše  $\mathcal{O}(|\alpha|)$
  - Vnější cyklus hodnotu  $l$  sníží dohromady o nejvýše  $\mathcal{O}(|\alpha|)$

# Vytvoření pole $R$ zdvojováním

## Popis algoritmu

- 1  $R_k[i]$  je počet  $j$  takových, že  $\alpha[j : j + k] < \alpha[i : i + k]$  ① ② → když podstav delky  $k$  je lex. menší jak podstav na pozici  $i$
- 2 Cílem je spočítat  $R_k[i] = R[i]$  pro nějaké  $k \geq |\alpha|$
- 3 K výpočtu  $R_1$  jen počítáme počet výskytů jednotlivých písmen
- 4  $\alpha[i : i + 2k] < \alpha[j : j + 2k]$  právě tehdy, když nebo mi to rozhodne oř dvanáct píslka Zelení už známe  
už větší bylo musí  $\alpha[i : i + k] < \alpha[j : j + k] \vee (\alpha[i : i + k] = \alpha[j : j + k] \& \alpha[i + k : i + 2k] < \alpha[j + k : j + 2k])$
- 5  $R_{2k}[i] < R_{2k}[j]$  právě tehdy, když → vychází 2 bodu ④  
 $R_k[i] < R_k[j] \vee (R_k[i] = R_k[j] \& R_k[i + k] < R_k[j + k])$
- 6 K získání  $R_{2k}$  třídíme dvojice  $(R_k[i], R_k[i + k])$  pro  $i = 0, \dots, |\alpha| - 1$

## Časová složitost

- 1 Třídění pro získání  $R_0$  máme v čase  $\mathcal{O}(|\alpha| \log |\alpha|)$  ③
- 2 Následuje  $\mathcal{O}(\log n)$  přihrádkového třídění
- 3 Celková složitost je  $\mathcal{O}(|\alpha| \log |\alpha|)$
- 4 Kärkkäinen, Sanders, 2003: Konstrukce v lineárním čase

- 1 Porovnáváme tedy jen prvních  $k$  znaků sufixů.
- 2 Pokud v zápisu  $\alpha[i : j]$  horní index přesahuje délku slova, pak  $\alpha[i : j]$  značí jen  $\alpha[i : \cdot]$ .
- 3 Pokud je  $|\Sigma| \leq |\alpha|$ , pak můžeme použít přihrádkové třídění a docílit tím času  $\mathcal{O}(|\alpha|)$ .

- rozšíříme podle písmene
- pak podrobdělím podle druhé,
- pak 2-4 písmena
- pak 4-8 písmena
- pak 8-16

proto k-čho budou  
mácinny 2ky

→ což mi povede  
na log. složitost

# Vytvoření sufixového pole zdvočováním

1 Vytvoříme pole  $X[0 \dots n]$  a  $R[0 \dots n]$ .

# Báze: vytvoření  $R_0$

2  $D = \{(\alpha[i], i); i = 0, \dots, n\}$  seříděné lexikograficky

3 **for**  $j = 0, \dots, n$  **do**

4  $X[j] = D[j][1]$

5 **if**  $j = 0$  nebo  $D[j][0] \neq D[j - 1][0]$  **then**

6 |  $R[X[j]] = j$

7 **else**

8 |  $R[X[j]] = R[X[j - 1]]$   $\rightarrow$  protože nás zajímají ostré nové slova

# Indukční krok: vytvoření  $R_{2k}$  z  $R_k$

9 **for** ( $k = 0; k < n; k = 2k$ ) **do**

10  $D = \{(R[i], R[i + k], i); i = 0, \dots, n\}$  seříděné lexikograficky

11 **for**  $j = 0, \dots, n$  **do**

12  $X[j] = D[j][2]$   $\rightarrow$  "který to je suffix?"

13 **if**  $j = 0$  nebo  $(D[j][0], D[j][1]) \neq (D[j - 1][0], D[j - 1][1])$  **then**

14 |  $R[X[j]] = j$

15 **else**

16 |  $R[X[j]] = R[X[j - 1]]$

$\rightarrow$  pozice prouto, abychom na konci věděli, kterému  
dávám suffixu to chceme přiřadit.

$\rightarrow$  nové podstavě délky  $2k$ ,  
jinak znovu vidíme stejný  
podstavě jako v minulém  
kroku  $\rightarrow$  totéž RE...3 ukazuje

- Jak najít nejdelší opakující se podřetězec?
- Jak určit počet různých podřetězců délky  $k$ ?
- Jak najít jehlu v kupce sena?

Hledání bodů v rovině i více-dimenzionálním prostoru.