

Suffixový strom

def:

suffixový strom pro větevec S ,
 $|S| = m$

- kořenový strom s m listy
- vnitřní vrcholy (včetně kořene) mají aspoň 2 potomky
- hrany jsou označeny ^{nepřesahujícími} podřetězcemi slova S
- pro každé vrchol platí: má-li aspoň 2 potomky, listi se označují hranou vedoucí do těchto potomků v prvnímu zvrcholu



Značení:

větevec indexujeme od 1

$S[i]$... symbol na indexu i

$1 \leq i \leq |S|$

$S[i..j]$... podřetězec

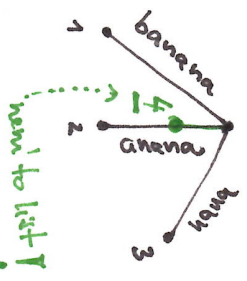
první symbol $S[i]$,

poslední symbol $S[j]$

Př. $S = \text{banana}$
 1 2 3 4 5 6

$S[4..6] = \text{ana}$

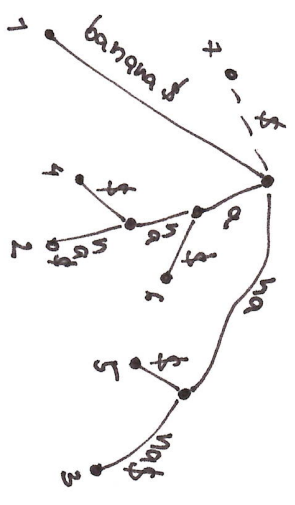
$S[2..4] = \text{ana}$



Proto uvažujeme symboly \$ a, b, c, ... \$
 který není v S a dáváme ho na konec. Dostaneme tedy

$S = \text{banana}\$$

a strom:



Poznámka: pro většinu slova suffixový strom podle původní definice nemusí existovat:

pokud $S = S_1 \dots S_m = S_1 \dots S_k$ pro $S_1, |S_1| = m$

(tj. existují suffixy S_1 , který je prefixem všech suffixů S_i), pak nemáme pro $S_1 \dots S_m$ list!

Algoritmy pro sestavení suffixového stromu

If we have sorted suffixes, we can build suffix tree

Naivní algoritmus pro S^k , $|S|=m$

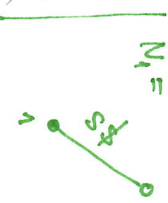
postupností suffix stromů N_1, \dots, N_m , kde

N_i : obsahují suffixy $S[i..m], S[i+1..m], \dots, S[i..i]$

(se započítávají i na konci)

Postup pro $N_i \rightarrow N_{i+1}$:

- Vkládáme $S[i+1..m]$



- čteme $S[i+1..m]$ od kořene s pomocí označení ~~...~~ řádku tak daleko, jak to jde. Zbude x .

Skouška:

→ upřesnění hrany:

→ přidáme na to místo nový ~~...~~ uzel, pokud jeho potomek přidáme nový list s označením $i+1$, hrana do něj označíme x .



pridání 2 uzlů a 1 hrana.

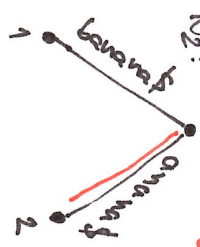
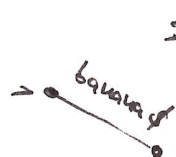
→ ve výsledku:

→ přidáme pomocí uzlů nové listy potomek s označením $i+1$ a hrana do něj označíme x



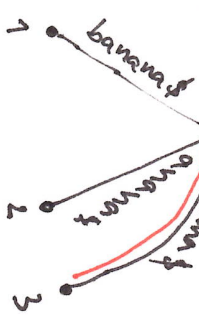
→ Chceme v předchozím křehkém nástroji v řádku: Zalding suffix S^k celý přefixem jeho suffixu S^k .

Př. banana\$

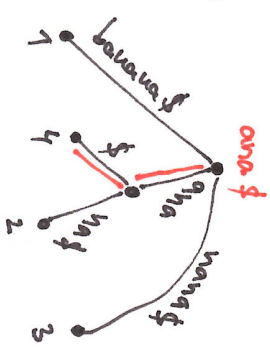


banana\$

N_3 :



N_4 :



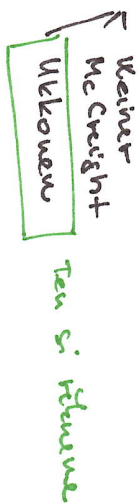
Stabilitost

$$(m-1) + (m-2) + (m-3) + \dots + 1 \in O(m^2)$$

počet porovnání
segmentů = $O(m^2)$

Algoritmy pracující v lineárním čase

- existují i.e.:



Ukkonenův algoritmus

Ukážeme neefektivní formu | kterou optimalizujeme až na DCW). ($|S| = m$).

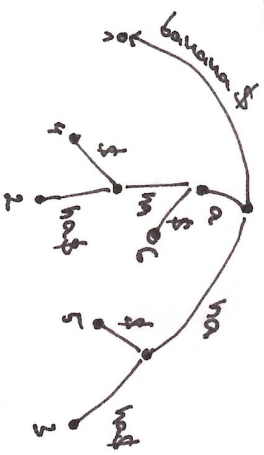
Implicitní suffixový strom pro S

← kde u a v součástí S

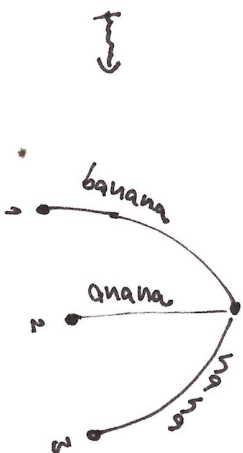
suffixový strom pro S

- odstraníme $\$$ z označení u nebo
- odstraníme hrany bez označení (+ izolované uzly vyloučí)
- odstraníme vnitřní uzly vzhledem k tomu, že dva potomky

Pr:



suffix t. pro banana



impl. s. st. pro banana

Pozn:

- impl. s. st. obsahují všechny suffixy (= obsahují všechny t pomocí označení hran přičtením) a jsou uloženi v listu.

- pokud S se podobá S výsledky S výsledky S

na jiném místě S , ~~suffix~~ je implicitní suffix t . kde! suffix t .
(t no prefix is a prefix of another suffix)

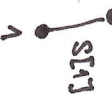
Pozn:

$I_1 \dots$ implicitní suffix strom pro $S[1..j]$

Ukkonenův algoritmus konstruuje posloupnost



strom I_1 je triviální:



Ve fázi $(i+1)$ se konstruuje strom I_{i+1} z I_i pomocí $i+1$ operací expanze. Tj. pro $j = 1 \dots i+1$ zadržíme, abychom byli prefixem suffix $S[j \dots i+1]$ ve strom I_{i+1} .

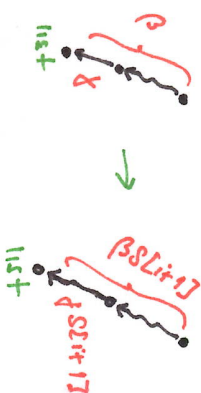
Tj. algoritmus je dvojitý cyklus. První cyklus je přes expanzi. Druhý cyklus je přes kontrakci.

Pravidla pro expanzi

$I_i \rightarrow I_{i+1}$ $3 \leq i \leq i+1$, zapišeme $SI_j \dots SI_{i+1}$ je ve stromu
ozn. $\beta = SI_{i+1} \cdot i$, $(SI_{i+1} | i) = \beta \nabla$

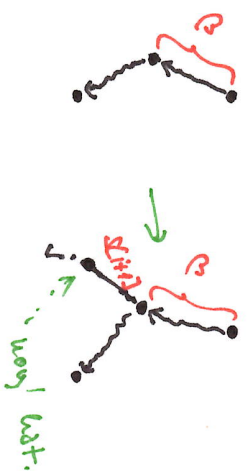
1. Najdeme konec β ve stromu (β patří do $I_i \nabla$)
2. totožto místo předejde vztáčení o symbol SI_{i+1}

Pravidlo 1: β končí v listu
 \rightarrow přidáme SI_{i+1} k označení posloupnosti
řádku na cestě



Pravidlo 2: z konce β upeřeráží cestu symbolem SI_{i+1} ,
cste vjeme v listu).

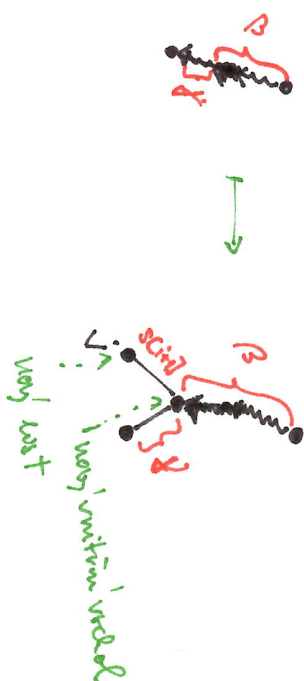
\rightarrow konec je uzel. Přidáme nového potomka a
řádku s označením SI_{i+1} . Potomk je označen i



[TRDY PŘIDÁNÉ LIST,
NEPŘIDÁNÉ VÍTKU]
VRCHOL

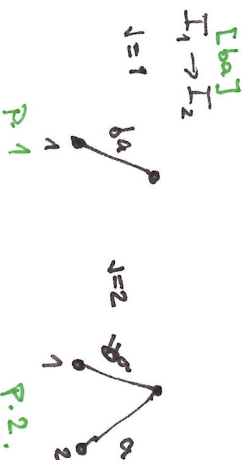
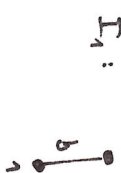
\rightarrow konec není vrchol.

Přidáme nový vnitřní uzel na
místo konce β a upráníme označení
řádku. Potom postupujeme jako u předchozí
bodě.

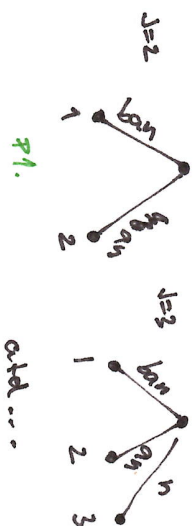
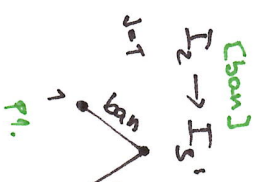


Pravidlo 3: z konce β pokračují cestu
symbolem SI_{i+1}
 \rightarrow udeláme mt, $SI_i \dots SI_{i+1}$ už ve
stromu je

P_i barana



P.2.



atd...

Rychlý pohyb stromem

(Fare 2)

→ více se tímto tam pou-
? úroveň form-
? úroveň form-
? úroveň form-

situace: - jsme ve vrcholku v , (tj. ji to kořen našeho podstromu)
- vlnu $z \in v$ pře generovat (délka k listům) pomocí
úroveň B

(~~to~~ odpovídá situaci, kdy na poslední stránce
nedáme konec B z vrcholku $\bar{S}(v)$, viz obr. 2)

Chceme zjistit, aby složitost vložení konce B závisela
na počtu vrcholů, přes které musíme projít, vlnou na $|B|$.

→ na hranách si pamatujeme délku křídla označení

Jsme-li ve vrcholku w , potom podle definice pravidel symbolů
označení w má vedle w do potomků značení, po které
hraní se y dost. ~~generace~~ ^{generace} ~~uznání~~ ^{uznání} hrany w (viz).

Jeli délka této hrany g a délka B je h (tj. $|B| = h$) pak

- pokud $g > h$: konec B je na B -tém symbolu označení hrany
- pokud $g < h$: konec B naležeme pomocí postupem pro
vrchol z a příslušný suffix B délky $h-g$
- pokud $g = h$: konec B je vrchol z .

Příklad algoritmus tedy zjistit, že složitost vložení
konce B je lin. závislá na počtu vrcholů, které musíme
~~projet~~ navštívit.

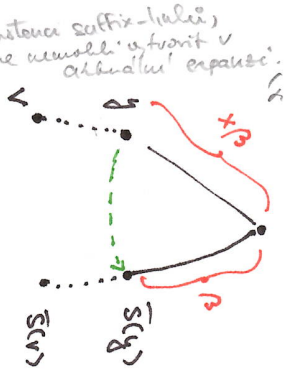
ALGORITHMY

Lemmy: $(v, \bar{S}(v)) \dots$ suffix-link

v čase, kdy půjdeme přes suffix-link
plate depth $(v) \leq \text{depth}(\bar{S}(v)) + 1$.

[depth (x) je hloubka vrcholu x]

Důkaz:



Existence suffix-linků, je možné, existují v abjadu' expanzi: (i kladě)

Potom má g suffix-link do $\bar{S}(g)$
a path-label $\bar{S}(g)$ je B .

Ale protože $x \in B$ je prefix path-label (g) , a tedy
tj. B je prefix path-label $(\bar{S}(v))$ a tedy
 $\bar{S}(g)$ je prefix $\bar{S}(v)$.

Tj. počet prefixů v je vždy
počet prefixů $\bar{S}(v)$, zvěřejně 0 1,
pokud má v prefix s path-labelem x ,
(= jeho symbol) a $B = \epsilon$, suffix-link
tohoto prefixu je kořen.

Použití pozorování:

→ označuji hru dělné pomocí dvojice indexů $[i, j]$, to odpovídá věci $SL[i, j]$

→ zavádíme globální proměnnou $j \in \mathbb{N}$

hru vedeme do kritické body označím

$[i, j] \in J$ (pro vždy index i)

\exists Tím j_i v tom i je možné přejít všechny
expozice pravidla \uparrow pomocí inkrementace \in
 $[0, 1]$

Věsta:

Ukrojení algoritmus s vyhledáním
složitost $O(n)$.

Dukaz:

Uvedeme.

Algoritmus pro Fibi $(i+1)$

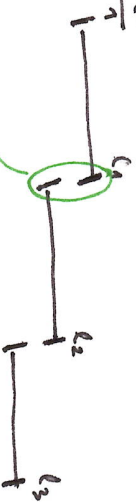
1. inkrementujeme e o 1

2. počítáme expozice L_i až L^* , kde

L^* je první expozice používající pravidlo 3,
nebo poslední expozice v této Fibi $(i+2)$

3. nastavíme L_{i+1} na L^*-1 .

Obrazek:



během řešení expozice \rightarrow Be přejít přes suffix-line.