

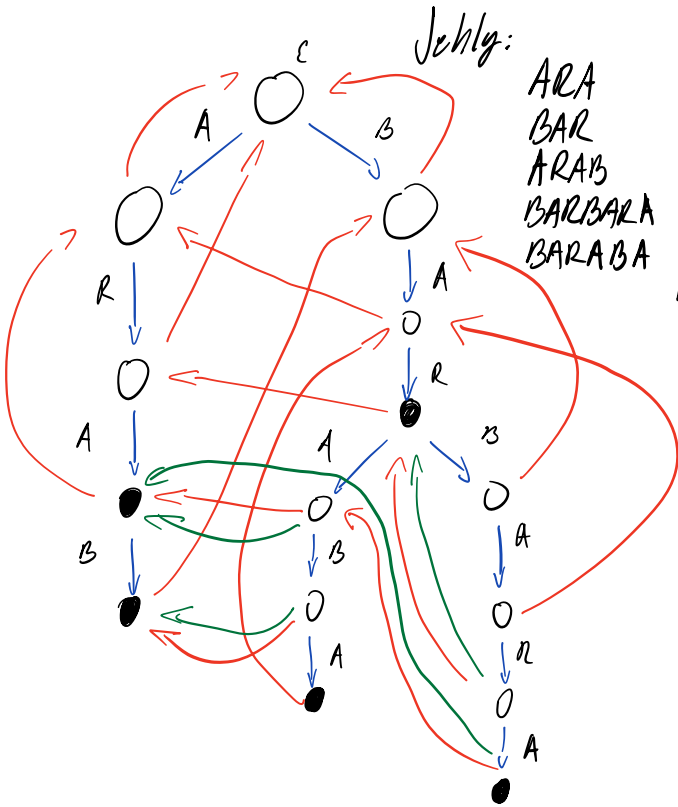
KMP - obecněji (Aho-Corasick)

$i_1 - i_n :=$ jehly

$\sigma :=$ slovo

stavy := prefixy všech jehel

hrany:
 - dopřední → prodlážením prefixu o jednu značku
 - zpětná → zkrácením prefixu na nejdelší vlastní suffix
 - zkratka → cesta po zpětných, jde na největší koncový stav - tedy literou jehly již obsahuje



Reprezentace automatu:

stavy očíslovíme, 0 = kořen (c)

slovo (i) := která jehla končí ve stavu "i"

zpět(i) := kam vede zpětná hrana

zkratka(i) := kam vede zkratková hrana

Dopředn(i, x) := kam vede dopřední hrana pro písmeno x

Vždy maximálně jedna od každé na jednom místě.

Krok (s, x): → To je že pro zadání x xx už není současně věta

Dokud **Dopředn(s, x) = ∅:** → pokud už není kam dál jít, oba skončit zpět

Dokud s = kořen: vrátíme s. → pokud jsem již v kořenu

s = zpět(s) → jímž jde po zpětných

Vrátíme Dopředn(s, x) → už existuje příslušná dopřední hrana

Hledí (σ):

s = kořen

Pro i = 0 — |σ| - 1:

s = Krok(s, σ[i])

t = s - temp. stav

Dokud t ≠ ∅:

Jeli slovo (t) ≠ ∅:

Hlášíme výslyt...

t = zkratka(t) → takže funkce jen dokud je def.

Lemmy: Hledí bývá v čase $O(|σ| + \# \text{výslytů})$

dopředných je nejvýše tolik co |σ|

zpětných je nejvýše tolik co dopředných

tedy $O(|σ|)$

výslytů je tolik, kolikrát se ohlásí.

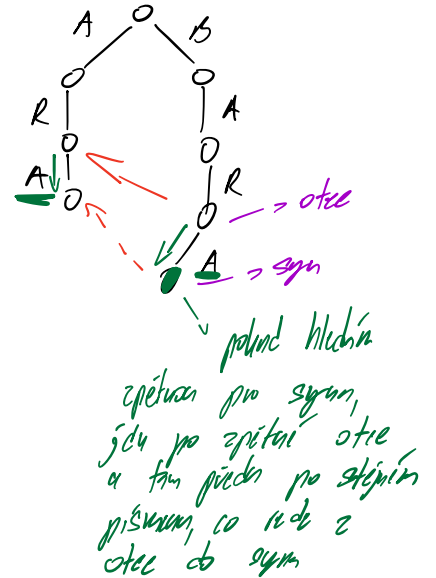
S tím, že poprvé se může projit s výslytem bez ohlášení...

Konstrukce automatu:

Myslel jsem: Stejně jako u UMD, nyní však bude počítat všechny jehly paralelně pro jednoho písmene a strom bude budovat po restřích.

Konstrukce ($i_1 - i_n$):

- 1) Vyberu si tři pro $i_1 - i_n \rightarrow 2$ toho máme dopředu
- 2) Zpět(n) (kořen) = \emptyset , Zhrat(n) (kořen) = \emptyset
- 3) Zpět(n) (\forall syn kořene) = kořen, Zhrat(n) (\forall syn kořene) = \emptyset
- 4) $F =$ frata se syny kořene:
- 5) Dokud $F \neq \emptyset$:
- 6) $V =$ dequeue(F)
- 7) Pro s syny v :
- 8) $Zpět(s) =$ Urok($Zpět(v)$, písmeno na hraně vs)
- 9) $F.Add(s)$
- 10) Pokud Slovo($Zpět(s)$) $\neq \emptyset$:
- 11) $Zhrat(s) = Zpět(s)$
- 12) Jinak $Zhrat(s) = Zhrat(Zpět(s))$.



\rightarrow Vynutím funkci, že se chodí po zpětových a předstírá je máš nalezené a pokračování

Lemna: Konstrukce běží v čase $O(\sum i_i)$

Bca Urok je vyhledávání tří asympt. lineární. Zbytek je BFS.

Jehly procházíme paralelně a přechod jednou jehlou je lineární.

A tedy i celý paralelní přístup je lineární. \rightarrow Dokonce často stav počítán jen jednou a u dalších jehel ho nemusím počítat.

Věta: Nalezení všech výskytů jehel v seně trvá $O(\sum i_i + |S| + \# \text{ výskytů})$.

Důkazy jsou spojeny lemmat.

Robinův - Karpiův alg.

$$h(x_0 - x_{j-1}) \rightarrow \partial \in \mathbb{Z}_H$$



musím umět přepočítat hash v konstantním čase po posunutí okénka jinak bych nic neudělal!

$$\rightarrow h(x_0 - x_{j-1}) :=$$

$$(x_0 r^{j-1} + x_1 r^{j-2} + \dots + x_{j-1} r^0) \bmod H$$

posunutí okénka

$$h(x_1 - x_j) =$$

$$(x_1 r^{j-1} + x_2 r^{j-2} + \dots + x_j r^0) \bmod H$$

Dodíl?

x_0 zmizelo, x_j přibýlo

ostatní se jen vynásobilo r -kem

Řešení:

$$h(x_1 - x_j) = h(x_0 - x_{j-1}) \cdot r - x_0 r^j + x_j$$

A tabule je konstantní $\approx O(1)$

RL alg.

0) zvolíme r z tělesa náhodně

1) $c = h(\text{jehly})$, $a = h(\sigma[0:j])$, spočítám r^j

2) pro $i = 0, \dots, |\sigma| - j$:

3) Pokud $a = c$ and $\sigma[i:i+j] = i$

4) Hlásím výskyt jehly i na pozici i .

5) $a = (a \cdot r - \sigma[i] \cdot r^j + \sigma[i+j]) \bmod H$.

Složitost:

- řešení + počítání hashů $= O(|\sigma|)$

- skutečné výskyty $= O(j \cdot v)$ → výskyty

- falešné výskyty =

\hookrightarrow pro ideální (obdobně náhodnou) h ,

$$P(\text{falešného výskytu}) = \frac{1}{H}$$

$$\rightarrow \text{průměrný čas} = O\left(\frac{|\sigma| \cdot j}{H}\right)$$

$\frac{1}{H}$ pro jednu stránku. Celkem $|\sigma|$ stránek.

Pokud $H \geq j$: $O(|\sigma|)$

Pozor: skutečný výskyt je hodně drahý. Vhodné pro normalizovaný jehel.

Mějme $P(x) := p_0 x^d + p_1 x^{d-1} + \dots + p_{n-1} x^{n-1}$ nad tělesem

Lemna: Pokud $x_1 - x_n$ jsou všechny kořeny polynomu P , pak $P(x) = (x-x_1) \cdot (x-x_2) \cdot \dots \cdot (x-x_n) \cdot Q(x)$

hde $Q(x)$ je polynom bez kořin.

Důsledek: Polynom stupně d má nejvýše d kořin

Uděl by k bylo větší než stupeň P , vyšel by spor.

Lemma: Necht' P a Q jsou polynomy stupně $< n$, x_1, \dots, x_n navzájem různá čísla t.č. $\forall i: P(x_i) = Q(x_i)$.
 Pak $P \equiv Q$.

$R := P - Q$, stupeň opět $< n$ t.č. $\forall i: R(x_i) = 0$ \rightarrow Jediný takový polynom je nulový polynom.
 $R \equiv 0 \Rightarrow P \equiv Q$

Tedy nyní falešný výsledek doplníme:

$$P(r) = h(i)$$

$$Q(r) = h(\sigma[i:i+j])$$

- Ptáme se, že pokud jsou polynomy různé, tak pro konkrétní r jich je pravděpodobnost, že se budou rovnat.

Tedy $P[P(r) \equiv Q(r)] \leq \frac{j}{H} \rightarrow \frac{1}{H}$ opět z velikosti množ. prvků, $\leq j$ -líst, protože to je také pozice, kde se dva polynomy mohou rovnat, aby stále byly identické.

Tedy H by muselo být alespoň j^2 velká.

Celková časová složitost je tedy $O\left(\frac{\kappa \cdot j^2}{H}\right) \rightarrow$ tedy pro $H \geq j^2$ je to $O(\kappa)$

$O(\kappa \cdot j)$:= ověření jehly v seně

$\cdot O\left(\frac{j}{H}\right)$ počet falešných výsledků (kromě toho správného).