

Bloom filter

- chceme udržovat $S \subseteq U$, ale v prostoru $o(|S|)$
- budeme opět hashovat, ale hash tabulka bude mít velikost $o(|S|)$
- a připustíme false positive

$$\text{pro } x \in S \quad \Pr[\text{alg řekne } "x \in S"] = 1$$

$$\text{pro } x \notin S \quad \Pr[\text{alg řekne } "x \in S"] > 0$$

-
- máme m -bitové pole $A[0, \dots, m-1]$ zpočátku vynulované
 - budeme používat k nezávislých hash. fcí h_1, \dots, h_k
 - dnes opět zcela náhodně ☺

Insert(x):

1. pro $i \in 1, \dots, k$
2. $A[h_i(x)] \leftarrow 1$

Find(x)

1. return $A[h_1(x)] \wedge \dots \wedge A[h_k(x)]$

Zahashovali jsme n prvků (tedy $n = 181$).

$\Pr[A[i]=0]$ pro nějaké i ?

každá hash se netrefí $A[i]$
to máme pro všech n prvků

$$\left(1 - \frac{1}{m}\right)^k$$

zahashuj se kamkoliv jinam

využijeme $m \gg k$ a $\lim_{a \rightarrow \infty} \left(1 - \frac{1}{a}\right)^a = e$

$$\rightarrow \Pr[A[i]=0] \approx e^{-\frac{kn}{m}}$$

Pst false positive?

pro prvek $x \in Y$ musí platit, že $A[h_1(x)], \dots, A[h_k(x)]$

musí všechny být 1

pst, že nejsou 1 je z předchozího výpočtu $\approx e^{-\frac{kn}{m}}$

$\Pr[\text{false positive}] \approx \left(1 - e^{-\frac{kn}{m}}\right)^k$ ← pro h_1, h_2, \dots, h_k

doplňte $A[h_i(x)]=0$, tedy $A[h_i(x)]=1$

Chceme mít pst. false positive pro prvek y menší než ϵ pro

$\epsilon < 1$. Jak vybrat m a k ?

$$\text{ozn } p = e^{-\frac{kn}{m}}$$

$$\Pr[\text{false positive}] \leq (1-p)^k$$

$$\text{ozn } f = (1-p)^k$$

• představme si, že m je taky pevné

• chceme minimalizovat f v prom. k

• ožn $g = \ln f = k \ln(1-p)$

• $\frac{dg}{dk} = \ln(1-p) + \frac{kn}{m} \cdot \frac{e^{-\frac{kn}{m}}}{1 - e^{-\frac{kn}{m}}}$

• ověřte si, že glob. min je v bodě $k = \ln 2 \cdot \frac{n}{m}$

• např. pro $m = \frac{n}{2}$ chceme $k = 2 \ln 2 = 1.39 \leq 2$

SUFFIX ARRAY

• pro řetězec S délky n máme pole A délky n

• $A[i]$... index i tého nejmenšího sufixu S

• $L[i]$... délka nejdelšího společného prefixu

$S[A[i]:]$ a $S[A[i+1]:]$

Jak spočítat ^{nejdelší} ^{společný} $S[A[i]:]$ a $S[A[j]:]$, pokud
prelix

máme pole L ?

RMQ

• = range min query

• máme pole A délky n

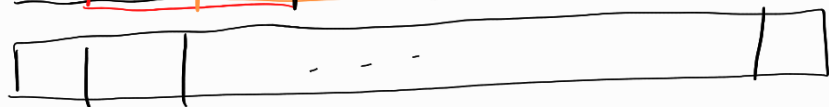
• chceme DS pro rychlé dotazy "min z $A[i], A[i+1], \dots, A[j]$ "



min z celého



min z intervalů délky 2



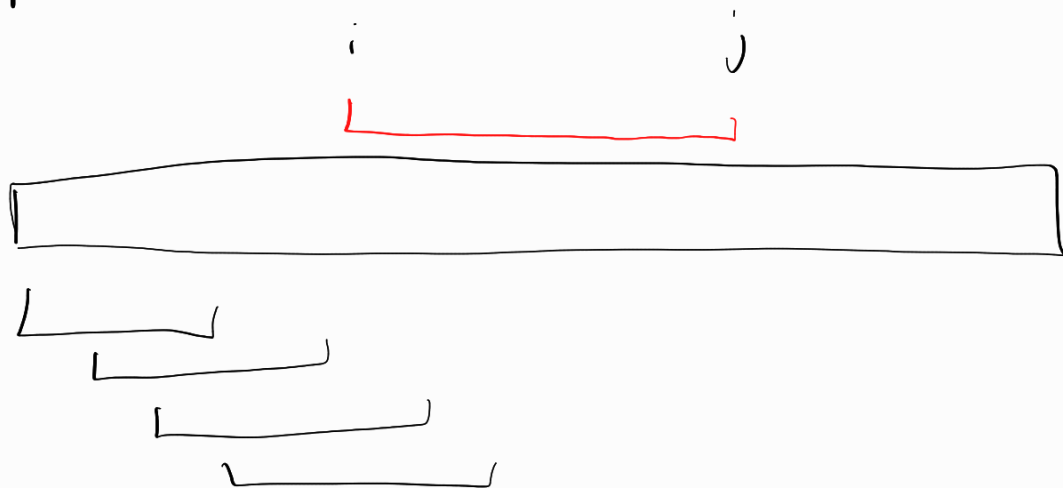
min z intervalů délky 1

- $k = j - i + 1$

- spočítáme p t. z. $2^p \leq k$ a $2^{p+1} > k$

- počítáme se na hladině pro intervaly délky

2^p



- bereme min z červených intervalů

→ 2 dotazy ... konst. složitost

- sice se nám překrývají kusy, ale minimálně

nevadí

- funguje pro max, gcd, ...

- obecně pro operace, kterým nevadí opakování a přibíhání na prvky → aby nevadil ten přibíhání

Jak za pomoci L hledat výsledky patternu P v čase $O(n + \lg m)$?

- opět binární přetváření

- l ... začátek vlevo

- r ... začátek vpravo

- m ... střed

- ozn. $s(i) = T[S[i]:]$

- představme si, že z rekurze máme $\text{lcp}(s(l), P)$ a $\text{lcp}(P, s(r))$

- pomocí sparse table v konst. čase spočítáme $\text{lcp}(s(l), s(m))$ a $\text{lcp}(s(m), s(r))$

- 3 případy, ozn. $k \leftarrow \text{lcp}(s(l), P)$

- i) $k < \text{lcp}(s(l), s(m))$

- $s(l)$ a $s(m)$ mají delší spol. prefix než $s(m)$ a P

$\rightarrow s(l)[k+1] = s(m)[k+1]$ a $s(l) \leq_{\text{lex}} P$

- tím pádem více $s(P)[k+1] > s(m)[k+1]$

- můžeme tedy nastavit $l \leftarrow m$

- ii) $k > \text{lcp}(s(l), s(m))$

- tím pádem $P <_{\text{lex}} s(m)$

- $r \leftarrow m$

iii) $k = \text{kp}(s(l), s(m))$

• pak podle prvního rozdílného znaku $P[k:]$ a $s(m)[k:]$

zjistíme, jestli $P \leq_{\text{lex}} s(m)$ nebo naopak

• jeden znak P nikdy nesrovnáváme s $s(m)$ dvakrát