

# DS v sublineárním prostoru

- vstup se nevejde do paměti → streaming model (1 nebo více sekvencí/průchodů)
- typicky jenko aproximace, často randomizované

## Bloomův filtr - reprezentace množiny s 1strannou chybou

- pole  $k$  bitů  $B[k] = B[k-1] +$  hashovací funkce  $h(x) \in U \rightarrow [k]$
- Insert( $x$ ):  $B[h(x)] \leftarrow 1$
- Find( $x$ ) testuje  $B[h(x)]$
- Delete neumíme

$x \in M \Rightarrow$  odpovíme AND  
 $x \notin M \Rightarrow$  možná také AND

Rozbor:  $Pr[\text{chyba} | x \notin M] = 1 - \prod_{a \in M} (1 - Pr[h(a) = h(x)]) = 1 - (1 - 1/k)^n \approx 1 - e^{-n/k}$  ... sblácíme pod  $1/k$  volbou  $k = cn$

→ pro  $\delta = 1/2$  [Pr chyby] potřebujeme prostor  $\Theta(n)$  bitů a čas  $O(1)$  na op.   
 ↑ pozor, silné požadavky na uzavřenost  $h$ !

Zesílení: Poradíme si  $t$  filtrů paralelně s různými  $k$ ,  $Find(x) = \bigwedge_{i=1}^t Find_i(x)$ .  
 $\Rightarrow \delta \leq \alpha^t$ , čas  $O(t)$ , prostor  $\Theta(tn)$ .

[třeba nechat všechny  $h_i$  hashovat do téhož pole  $aport$ ]

## Majoneta - cílem najít prvek, který se ve streamu vyskytl více jak $n/2$ -krát

- udržujeme kandidata  $k$  a jeho počet  $c$

Init:  $k \leftarrow \emptyset, c \leftarrow 0$   
 Process( $x$ ):  
 Pokud  $k = \emptyset$  &  $k \neq x$   
 Pokud  $x = k$  &  $c++$   
 Jinak  $c--$   
 Pokud  $c = 0$  &  $k \neq \emptyset$

Věta: Pokud prvek  $m$  je majoritní, pak na konci  $k=c$ .  
 (dalším průchodem možno ověřit)

De: Rozdělíme vstup na bloky, blok končí resetem  $k$



V každém bloku má kandidát přesně polovinu zastoupení → zde musí mít nadpoloviční ⇒ je to kandidát

⇒  $O(n)$  času,  $O(\log n + \log U)$  prostoru

## Odhad frekvencí (zejm. pro hledání prvků vyskytujících se více jak $(n/k)$ -krát) - Misra, Gries 1982

- udržujeme slovník počítadel  $A[x]$ , třeba v BVS

Init:  $A \leftarrow \emptyset$   
 Process( $x$ ):  
 Pokud  $x \in \text{keys}(A) : A[x]++$   
 Jinak je-li  $|\text{keys}(A)| < k-1 : A[x] \leftarrow 1$   
 Jinak:  $\forall a \in \text{keys}(A) : A[a]--$   
 Pokud klaso na 0, odstraníme  $A[a]$ .

pro analýzu: představíme si, že  $A[x]$  vždy zvyšuje a pak můžeme zkusit hledat shráně

Output: Pro  $x \in \text{keys}(A)$  vrátíme celkový  $A[x]$ , jinak vrátíme 0.

obecně platí:  $A[a] \leq 0$ , když  $a \notin \text{keys}(A)$

Paměť:  $O(k \cdot (\log n + \log U))$   
 Čas:  $O(n \cdot \log k)$

známení:  $f_a =$  skutečná frekvence  $a$   
 $\hat{f}_a =$  náš odhad, tedy  $A[a]$  malouci výřtu

Věta:  $\forall a : f_a - n/k \leq \hat{f}_a \leq f_a$   
 ↑ ↑  
 řízení řízení

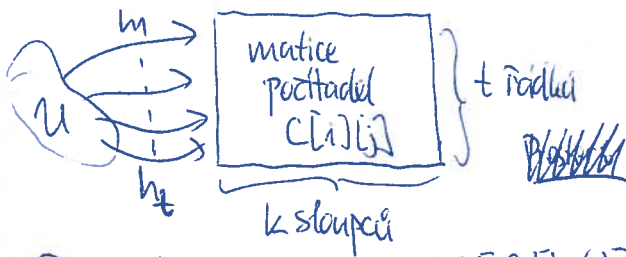
De: když  $k$  snižují  $A[a]$ , snižím i  $k-1$  dalších počítadel  
 ↓  
 stane se to max.  $n/k$ -krát

Důsledek: Uvime najít (na 2 průchody)  $\{a \mid f_a > n/k\}$

potenciálně  $\Sigma$  počítadel

Pravděpodobnostní odhad frekvencí - Count-Min

to je jižite stavila obz. universalita



$h_1 - h_t : U \rightarrow [k]$  náhodně vybereme 2 2-nezávislého systému funkcí

$\forall x, y \in U \forall i, j \in [k] \Pr[h_i(x)=i \& h_i(y)=j] = \frac{1}{k^2}$

Process(x): Pro  $i=1-t : C[i][h_i(x)] +$

Result(x):  $\min_i C[i][h_i(x)]$

Nastavíme  $k = \lceil 2/\epsilon \rceil, t = \lceil \log 1/\delta \rceil$

Chceme  $(\epsilon, \delta)$ -estimátor

$\Pr \left[ \left| \frac{x-OPT}{OPT} \right| > \epsilon \right] < \delta$

(něco jako FPTAS)

Analýza Je opt. frekvence (skutečné), Je výstup algoritmu  $\alpha$  evidentně  $f_a \geq f_a$  ... o kolik víc?

Fixa:  $X_i = C[i][h_i(a)] - f_a$

Pro  $j \in U \setminus \{a\} : Y_{i,j} = [h_i(j) = h_i(a)] \dots X_i = \sum_{j \neq a} f_j \cdot Y_{i,j}$

$E[X_{i,j}] = 1/k \Rightarrow E[X_i] = \frac{1}{k} \sum_{j \neq a} f_j = \frac{f - f_a}{k}$

Podle Markovae:  $\Pr[X_i > \epsilon \cdot \frac{f - f_a}{k}] < \frac{1}{k\epsilon} \leq \frac{1}{2}$

Pak:  $\Pr \left[ \sum_{i=1}^t X_i > \epsilon n \right] = \Pr[\exists i X_i > \epsilon n] \leq 2^{-t} \leq \delta \checkmark$  Je to  $(\epsilon, \delta)$ -estimátor.

prostor:  $O(\log \epsilon \log 1/\delta \log n + \log 1/\delta \log U)$

Odhad # různých prvků - Alon, Matias, Szegedy 1999 [AMS]

Df:  $tz(x) = \max \{i \mid 2^i \mid x\}$  (#trailing zeroes) "pravitková funkce"

Init:  $z \leftarrow 0, h : U \rightarrow U$  2-nezávislého syst. náhodně vybraná

Process(x):  $z \leftarrow \max(z, tz(h(x)))$

Output:  $2^{z+1/2}$

Intuice: každý  $2^i$ -tý prvek má  $tz(x) \geq i$

Analýza  $\forall j \in U \forall r \geq 0 X_{r,j} = [tz(h(j)) \geq r]$  ... toto jsou 2-nezávislé jevy

$Y_r = \sum_{j: f_j > 0} X_{r,j}$

$t$  = hodnota  $z$  na konci výpočtu  
 $d$  = # různých prvků  
 $d$  = náš odhad, tedy  $2^{t+1/2}$

$t \geq r \Leftrightarrow Y_r > 0$

$E[X_{r,j}] = 1/2^r, E[Y_r] = d/2^r$

$\text{var}[X_{r,j}] \leq E[X_{r,j}^2] = E[X_{r,j}] = 1/2^r$

$\text{var}[Y_r] = \sum_{j: f_j > 0} \text{var}[X_{r,j}] \leq d/2^r$   
 díky 2-nezávislosti

Myjí zvolíme a nejmenší celé  $t \geq 2a$ . Pak  $2^{a+1/2} \geq 3d$ . Pak:

$$\Pr[\hat{d} \geq 3d] = \Pr[t \geq a] = \Pr[Y_a > 0] = \Pr[Y_a \geq 1] \stackrel{\text{Markov}}{\leq} \frac{E[Y_a]}{1} = \frac{d}{2^a} \leq \frac{\sqrt{2}}{3} \approx 0.47$$

$\uparrow$  z def.  $Y_a$        $\uparrow$  z celistvosti       $\uparrow$  Markov

Pro opačný směr zvolíme  $b$  největší celé, pro nějž  $2^{b+1/2} \leq d/3$ . Pak:

$$\Pr[\hat{d} < d/3] = \Pr[t \leq b] = \Pr[Y_{b+1} = 0] \leq \Pr[|Y_{b+1} - E[Y_{b+1}]| > \frac{d}{2^{b+1}}] \stackrel{\text{Chebyshev}}{\leq} \frac{2^{b+1}}{d} \leq \frac{\sqrt{2}}{3}$$

$\uparrow$  z def.  $Y_{b+1}$        $\uparrow$  z def.  $Y_{b+1}$        $\uparrow$  Chebyshev

Získali jsme tedy  $(3, \frac{\sqrt{2}}{3})$ -estimator.

Pr chyby lze snadno zlepšit s použitím paralelně  $t$  estimatorů, vrátíme median ...  
 -- podle Chernova  $\delta < 2^{-\Theta(t)}$ .

$\rightarrow (3, \delta)$ -estimator pro libovolné  $\delta$ , prostor  $O(\log \frac{1}{\delta} \cdot \log U)$ , čas  $O(\log \frac{1}{\delta} \cdot n)$

Lepší odhad # různých prvků, ... tentokrát pro libovolné  $\epsilon, \delta > 0$ . Bar-Yossef et al, 2004 (BJKST)

- Opět  $h: U \rightarrow U$  a počítání  $t_z(h(x))$
- Tentokrát si pamatujeme  $B$  - množinu hodnot s daným  $t_z$ , ale omezíme její velikost.

Init: zvolíme  $h$  náhodně,  $z < 0$ ,  $B < \emptyset$

Process(x): Je-li  $t_z(h(x)) \geq z$ :  
 $B \leftarrow B \cup \{(x, t_z(h(x)))\}$   
 Dokud  $|B| \geq c/\epsilon^2$ :  
 $z \leftarrow z+1$  nastavíme prah  
 odstraníme z  $B$  všechny dvojice s  $t_z < z$

Intuice: cca  $n/2^z$  prvků se objeví v  $B$   
 $\downarrow$   
 $|B| \cdot 2^z$  je dobrý odhad

Output:  $|B| \cdot 2^z$

Prostor:  $O(\frac{1}{\epsilon^2} \cdot \log U)$  [lze zlepšit: místo  $x$  si v  $B$  pamatujeme nějaké hash] do dost velkého prostoru

Analýza:  $X_r$  a  $Y_r$  jako předtím.

$\hat{d} = 2^t \cdot Y_t$  ( $t$  je poslední hodnota proměnné  $z$ )

Pokud  $t=0$ , alg. počítá přesně.

Chyba  $\Leftrightarrow |Y_t \cdot 2^t - d| > \epsilon \cdot d \Leftrightarrow |Y_t - \frac{d}{2^t}| > \frac{\epsilon d}{2^t}$

pro malá a velká  $t$  odhadujeme jinak  
 pro malá  $t$  zvolíme  $s$  tak, aby  
 $\frac{12}{\epsilon^2} \leq \frac{d}{2^s} < \frac{24}{\epsilon^2}$

Myjí:  $\Pr[\text{chyba}] = \sum_{r=1}^{\log n} \Pr\left[|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r} \text{ \& } t=r\right] \leq \sum_{r=1}^{s-1} \Pr\left[|Y_r - \frac{d}{2^r}| > \frac{\epsilon d}{2^r}\right] + \sum_{r=s}^{\log n} \Pr[t=r]$

Líže  $\text{úst. strana}$

$\leq \frac{2^r}{\epsilon^2 d}$

$\Pr[|Y_r - E[Y_r]| > \frac{\epsilon d}{2^r}]$        $\Pr[t \geq s]$   
 $\Pr[Y_{b+1} \geq c/\epsilon^2]$

$$P \leq \sum_{r=1}^{s-1} \Pr\left[\left|Y_r - \frac{d}{2^r}\right| > \frac{\epsilon d}{2^r}\right] + \sum_{r=s}^{\log n} \Pr[t=r]$$

$$\Pr\left[|Y_r - \mathbb{E}[Y_r]| > \frac{\epsilon d}{2^r}\right] = \Pr[t \geq s] = \Pr[Y_{s-1} \geq c/\epsilon^2]$$

↓ Chebyshev

$$\leq \frac{d}{2^r} / \left(\frac{\epsilon d}{2^r}\right)^2 = \frac{d \cdot 2^r \cdot 2^r}{d^2 \epsilon^2 2^r} = \frac{2^r}{\epsilon^2 d}$$

$$\sum_{r=1}^{s-1} \frac{2^r}{\epsilon^2 d} \leq \frac{2^s}{\epsilon^2 d} \leq \frac{\epsilon^2}{\epsilon^2 \cdot n} = \frac{1}{n}$$

geom. řada

voleba s

↓ Markov

~~$$\leq \frac{d}{2^s} / c \epsilon^2 = \frac{2d}{2^s \cdot c}$$~~

$$\leq \frac{d}{2^{s-1}} / c = \frac{2d \epsilon^2}{2^s \cdot c} \leq \frac{2 \cdot 24 \epsilon^2}{\epsilon^2 \cdot c} = \frac{48}{c}$$

→ cca. mírně stáhnout pod 1/2 volba vhodné c

$\Pr[\text{chyba}] \leq 1/6$  -- pod 5 srazíme medicínou trikem.