

BLOOMOVY FILTRY

[první struktura tohoto typu: Bloom 1970]

Úloha: Chceme reprezentovat množinu $X \subseteq U$, $|X|=n$ a odpovídat na dotazy "a ∈ X?"
... ovšem máme málo paměti ⇒ spokojujeme se s přibližným řešením.

Možné chyby: false positive (tvrdíme, že a ∈ X, pro a ∉ X) } chceme omezit pravděpodobnost
false negative (opačně) [ty nakonec řádně nelidnou]

První pokus: Pořídíme si pole B ^[0...m-1] bitů a kešovací fci $h: U \rightarrow [m]$ z 1-univ. systémem.
začneme se samými 0, postupně nastavíme bity $h(x_1), \dots, h(x_n)$.

- Find(a) vrátí TRUE, pokud $B[h(a)] = 1$
 - pro a ∈ X vždy odpoví správně
 - pro a ∉ X false positive s pstí $\leq \frac{n}{m}$ (z 1-universality)

↑
pro $m = c \cdot n$
je pst. chyby $\leq \frac{1}{c}$

- Insert(a) nastaví $B[h(a)] \leftarrow 1$
- Delete neumíme provést

Vylepšení: Pořídíme si k filtrů s $m = \alpha n$ a různými keš. funkcemi (nezávislými).
Find odpoví TRUE, pokud najdeme ve všech filtrech:
- pro a ∈ X odpovíme správně
- pro a ∉ X je $\Pr[\text{false positive}] \leq 2^{-k}$

Tedy pro chyby ϵ potřebujeme prostor $2n \cdot \log \frac{1}{\epsilon}$.

Podrobnější analýza (předpokládá dokonale náhodné keš. fce)

- fixujeme celkový prostor $M = mk$, hledáme optimální k
- $\Pr[B_i[j] = 0] = (1 - \frac{1}{m})^n = (1 - \frac{k}{M})^n \approx e^{-\frac{kn}{M}}$
_{označíme p}
- $\Pr[\text{false positive}] = (1-p)^k \approx (1-p)^k = (1 - e^{-\frac{kn}{M}})^k = e^{k \cdot \ln(1-p)}$
- $k \cdot \ln(1-p) = (-\frac{M}{n}) \ln p \cdot \ln(1-p)$

↳ má na intervalu [0,1] min. pro $p = \frac{1}{2}$

- chceme tedy $e^{-\frac{kn}{M}} = \frac{1}{2} \dots -\frac{kn}{M} = -\ln 2 \dots k = \frac{\ln 2 \cdot M}{n} \approx 0.69$

- pak $\Pr[\text{f.p.}] = 2^{-k}$

→ opačně: pro $\Pr[\text{f.p.}] \leq \epsilon$ chceme $k := \lceil \log \frac{1}{\epsilon} \rceil$,
tedy $M = \frac{kn}{\ln 2} \approx \log \frac{1}{\epsilon} \cdot n \cdot \frac{1}{\ln 2} \approx 1.44 \cdot n \cdot \log \frac{1}{\epsilon}$
 _{$\log e \approx 1.44$}

Dolní odhad (bez důkazu): Je třeba alespoň $n \cdot \log \frac{1}{\epsilon}$ bitů

⇒ náš B. filtr je ^{max} 1.44-krát horší než optimum.

Dalsí operace: • Máme-li pro $X, Y \subseteq U$ filtry B_x a B_y se stejným m , můžeme snadno spočítat $B_{x \cup y} = B_x \vee B_y$, $B_{x \cap y} = B_x \wedge B_y$

- z počtu 1 ve filtru lze odhadnout velikost množiny:
 - víme, že $\Pr[\text{bit je 0}] \approx e^{-kn/M}$
 - proto $\mathbb{E}[\#0] \approx M \cdot e^{-kn/M}$... navíc #0 je silně koncentrováný kolem své \mathbb{E} (bez důkazu)
- teď lze vše hesovat do 1 společné tabulky (tedy $h_i: U \rightarrow [M]$) ... asymptoticky stejné

Počítací Bloomovy filtry [Fan et al. 2000]

- chceme umět i Delete
- v poli jsou místo 1bit. indikátorů b -bitová počítadla
- Insert zvyšuje o 1, Delete snižuje o 1, Find testuje nenulovost
- POZOR: jakmile počítadlo dopočítá do max. hodnoty $(2^b - 1)$, musí tak navěky zůstat!

$\hookrightarrow \Pr[C_i[j] = t] = \binom{n}{t} \cdot \left(\frac{1}{m}\right)^t \left(1 - \frac{1}{m}\right)^{n-t}$

j-té počítadlo v i-tém filtru

předpokládáme pleš náhodné funkce h_i

$\Pr[C_i[j] \geq t] \leq \binom{n}{t} \left(\frac{1}{m}\right)^t \leq \left(\frac{ne}{t}\right)^t \cdot \left(\frac{1}{m}\right)^t = \left(\frac{ne}{mt}\right)^t \leq \left(\frac{e \cdot \ln 2}{t}\right)^t$

aspoň 1 t-tice se zachová do j-té příhrádky, zbývající prvky mohou kumulovat

std. odhad $\binom{n}{t}$, viz kapitoly z DM

už víme, že chceme $m \geq n \cdot \ln 2$

$\Pr[\exists j: C_i[j] \geq t] \leq m \cdot \Pr[C_i[j] \geq t] \leq m \cdot \left(\frac{e \cdot \ln 2}{t}\right)^t$ pro $t=15$ je to $\approx 3 \cdot 10^{-14}$

\Rightarrow pro takřka jakékoli rozumné m stačí 4-bitová počítadla

"Bloomier filtry" - reprezentace funkcí [původně Fan, Cao, Broder 2000]

[zde upřesněno dle Charles, Chellapilla 2008]

- Máme funkci $f: X \rightarrow [2^r]$ definovanou na $X \subseteq U$
- Chceme sestavit funkci $g: U \rightarrow [2^r]$ takovou, že $g|_X = f$ (mimo X může dávat cokoliv)
- Inspirace kukacím hesováním:
 - pole $R[0 \dots m-1]$ r -bitových hodnot
 - 2 hesovací funkce $h_1, h_2: U \rightarrow [m]$
 - Find vrátí $R[h_1(x)] \oplus R[h_2(x)]$ *XOR*

- Pro danou X najdeme h_1, h_2 t.č. kukacím graf je acyklický (pro $m \geq 2 \cdot (1+\epsilon)n$ to nastane s $pstí$ aspoň konst. \Rightarrow stačí průměrně $O(1)$ pokusů)
- Pak každou komponentu prohlédáme do hloubky. Kořeni nastavíme $R[-] \leftarrow 0$, ostatními vrcholy dopočítáme, aby vyšly XORy.
- Stačí nám $2(1+\epsilon)m$ bitů paměti (konstanty lze zlepšit použitím více hes. fcí).