

ab-experiment: časté chyby

- někteří jste si nevšimli, že $(a, 2a)$ -strany mají konst. amort. počet strukt. změn při insertu a deletu.

- to je rozepsáno ve skriptech, sekce 3.2 v kapitole o (a, b) -stromech

HASHOVÁNÍ

Značení

- m ... velikost hashovací tabulky,

- $\mathcal{U} = \{0, \dots, U\}$... univerzum. Tedy prvky, jež je podmnožina S si chceme pamatovat.

- h, h_1, h_2, \dots

... hashovací fce. Dnes to bude taková fce, že $f(x)$ je uniformně náhodná hodnota z $\{0, \dots, m-1\}$.

Jednoduše se analyzuje ☺

- pokud pro $x, y \in \mathcal{U}$ platí $h(x) = h(y)$, pak nastala kolize.

- pokud útočník zná hashovací fci, pak lze vyvodit, že operace find / insert / delete budou mít složitost $\Omega(n)$.

- proto volíme hashovací fce náhodně

- „pokud ani já nevím, co dělám, třeba to bude vidět útočník“ ☺

- takže např. klád $\Pr[h(x)=h(y)]$... co je náhodného? To hashovací fc.
- úplný zápis by byl $\Pr_{h \in \mathcal{H}} [h(x)=h(y)]$, kde \mathcal{H} je množ. fun. z \mathcal{U} do $\{0, \dots, m-1\}$
 ↖ udáme uniformně náhodně.

Př: Hashujeme n prvků do pole velikosti $m := n^2$. Jaká je pst., že nastane kolize? (Chceme určit nějakou rozumnou horní mez, tj. ne ≤ 1)

• $\Pr[\exists i, j \in \{1, \dots, n\} \text{ t.j. } h(x_i) = h(x_j) = k] \leq \binom{n}{2} \cdot \left(\frac{1}{n^2}\right)^2 \leq \frac{n^2}{2} \cdot \frac{1}{n^4} \leq \frac{1}{2} \cdot \frac{1}{n^2}$

• $\Pr[\exists k : X_k] \leq n^2 \cdot \frac{1}{2} \cdot \frac{1}{n^2} \leq \frac{1}{2}$
 union bound přes přehledky

v tu inkluzi exkluzi každý
 - je větší než následující +

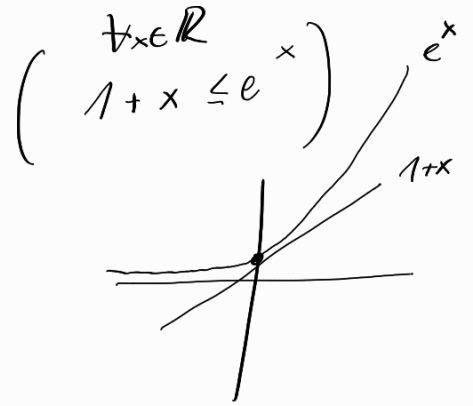
• tedy pst., že nastane kolize je nejvýše $\frac{1}{2}$

Př: (Narozeninový paradox) Hashujeme do tabulky velikosti m . Kolik je potřeba zhashovat prvků, aby nastala aspoň jedna kolize s pstí aspoň $\frac{1}{2}$?

- p_k ... ne nastane kolize, hashujeme-li k prvků
- $p_1 = 1$
- $p_2 = 1 \cdot \left(1 - \frac{1}{m}\right)$
 druhý prvek se smí zhashovat kamkoliv kromě místa, kam se zhashoval první
- $p_3 = 1 \cdot \left(1 - \frac{1}{m}\right) \left(1 - \frac{2}{m}\right)$
 ;

$$\bullet p_k = 1 \left(1 - \frac{1}{m}\right) \cdots \left(1 - \frac{k-1}{m}\right)$$

$$\bullet p_k \leq 1 \cdot e^{-\frac{1}{m}} \cdot e^{-\frac{2}{m}} \cdots e^{-\frac{k-1}{m}} = e^{-\frac{1}{2m}k^2}$$



$$\bullet \text{při jakém } k \text{ platí } e^{-\frac{k^2}{2m}} \leq \frac{1}{2} \quad ?$$

$$\bullet \text{pro } k = \sqrt{2m} : e^{-1} = \frac{1}{e} \leq \frac{1}{2}$$


→ hashujeme-li $\sqrt{2m}$ prvky, pak pst kolize je aspoň $\frac{1}{2}$

Srovnáme s předchozí úlohou

• tam jsme takřka hashovali m prvky, ale pst kolize byla

nejvýše $\frac{1}{2}$???

• dost krátkě jsem ignoroval konstanty ... ty to dost ovlivní

• a celkově jsem počítal jako 

Cv: Mějme hash tabulku velikosti m . Ukážete, že existuje konst. c_1

t.č. pst, že nenastane kolize při hashování $c_1 \sqrt{m}$ prvky, je nejvýše $\frac{1}{2}$.

Podobně ukážete, že existuje konst. c_2 t.č. pst, že nenastane kolize při

hashování $c_2 \sqrt{m}$ prvky, je aspoň $\frac{1}{2}$.

- Pokuste se minimalizovat $|C_1 - C_2|$.
 - V nějakém bodu výpočtu by se mohlo hodit uvažovat o destičkách.
 - Mohlo by se hodit uvidět, že $e^{-x-x^2} \leq 1-x$ pro $x \in [0, \frac{1}{2}]$
 - Tuhle úlohu můžete najít v
M. Mitzenmacher, E. Upfal: Probability and Computing
jako cvičení 5.3.
-

KUKAČČÍ HASHOVÁNÍ

- chceme worst-case konstantní řád, ne jen v průměrném případě
- princip: máme dvě hashovací fce h_1, h_2
- pro $x \in \mathcal{U}$ platí, že v tabulce (či řádě) na indexu $h_1(x)$ nebo $h_2(x)$ a nikde jinde!
- pseudokód insertu je ve skriptech nebo poznámkách.

Proč to funguje?

Def: Kukaččí graf G má vrcholy $V(G) = \{0, \dots, m-1\}$, f_j indexy

tabulky. Pro každý prvek x obsažený v tabulce přidáme hranu $\{h_1(x), h_2(x)\}$, tedy neorientovanou.

? Kdy seže insert?

Pokud pro vkládání prvku y platí, že $h_1(y)$ leží na cyklu kubičtého grafu.

Lem: Necht $c > 1$ je konstanta a $\frac{n}{m} \leq \frac{1}{2c}$ ($\frac{n}{m}$ je tzv. faktor zplnění).

Pro dva vrcholy s a t v kubičtém grafu platí, že pst existence cesty z s do t délky k je nejvýše $\frac{1}{mc^k}$.

Dk:

• in derkeri podle k

• $k = 1$

• union bound přes všech n prvků dáva

$$\Pr[\text{existuje } s \leftrightarrow t \text{ cesta délky } 1] \leq \sum_{t \in V} \frac{1}{m^2} \leq \frac{1}{mc} \quad \checkmark$$

• indukční krok:

• cesta z s do t délky k existuje, pokud existuje u t.ž., existuje

cesta z s do u délky $k-1$, která neprochází přes t , a hrana

z u do t .

$$X_{k-1}(s, u) :=$$

• $\Pr[\text{existuje } s \rightarrow u \text{ cesta délky } k-1] \leq \frac{1}{mc^{k-1}}$ z indukčního předpokladu

$$\Pr[X_1(u, t) \mid X_{k-1}(s, u)] \leq \frac{1}{mc}$$

• ! vrchol t může být na $s \rightarrow u$ cestě, ale pro horní

definice z teorie grafů říká,
že se na cestě nepokoují vrcholy

odhad to může těch $\frac{1}{mc}$ jen snížit.

$$\begin{aligned} \circ P_0 [X_k(s,t)] &\leq \sum_{u \in V} \left(\Pr [X_1(u,t) | X_{k-1}(s,u)] \Pr [X_{k-1}(s,u)] \right) \leq \\ &\leq m \cdot \frac{1}{mc^{k-1}} \cdot \frac{1}{mc} \leq \frac{1}{mc^k} \quad \square \end{aligned}$$

Jak dlouho trvá insert?

Lem: Necht $c > 1$ je konstanta. Střední délka cesty začínající v h_1 vkladáního prvku je $O(1)$.

Dk:

• použijeme předcházející lemma pro všechny cesty délky k a všechny začátky s

• střední délka cesty z $h_1(x)=s$ je nejvýše

$$m \cdot \sum_{k=1}^{\infty} k \frac{1}{mc^k} \leq \sum_{k=1}^{\infty} \frac{k}{c^k} = \frac{c}{(c-1)^2} \in O(1) \quad \square$$

Kolikrát rehashujeme?

Lem: Necht $c > 2$ je konstanta. Střední počet přehastování při vkladání prvku do prázdné tabulky velikosti m je $O(1)$, pokud $\frac{n}{m} \leq \frac{1}{2c}$

Dk:

• aplikujeme předcházející lemma na všechny délky cesty a pro všechna $s=t$.

$$\circ P_0 [\text{existuje cyklus v grafu}] \leq m \cdot \sum_{k=1}^{\infty} \frac{1}{mc^k} \leq \frac{1}{c-1}$$

• psť, že nastane rehash je tedy nejvýše $\frac{1}{c-1}$

• očekávaný počet rehashů je tedy $\leq \sum_{k=1}^{\infty} \left(\frac{1}{c-1}\right)^k \leq 1 + \frac{c-1}{c-2}$ \square

• vkládání má tedy konstantní průměrnou amortizovanou psť \smile

• princip cuckoo hashování se možná nezkouší.

• ale jakmile víte o cuckoím grafu, tak byste ten dík už jistě složili \smile

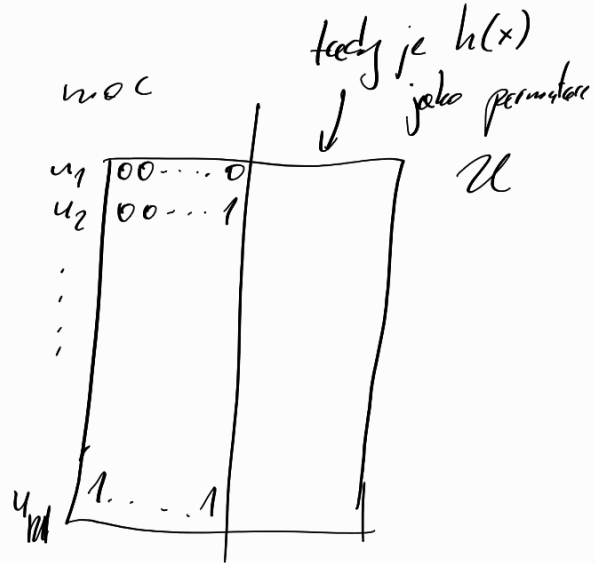
Tabulační hashování

• mit náhodnou ideální hash. fci je jako mít tabulku velikosti

$\lg(|U|!) \cdot \lg|U|$ bitů

To je moc \uparrow tedy je $h(x)$ jako permutace U

př.: 32-bit



• mířeme však mit ideální hash fci pro každý blok 8 bitů

zvlášť a výsledky třeba XORovat

• paměť už je jen $4 \cdot 2^8$ vs. předchozích 2^{32}

• potřebujete v DÚ cuckoo-hash