

HEŠOVÁNÍ: RŮZNÉ DRUHY TABULEK

Příklady rodnin hešovacích funkcí

- skalární součin: $\mathbb{Z}_p^k \rightarrow \mathbb{Z}_p$, $h_{\vec{a}}(\vec{x}) = \vec{a} \cdot \vec{x}$ pro $\vec{a} \in \mathbb{Z}_p^k \dots$ 1-univerzální
 $h_{\vec{a},b}(\vec{x}) = \vec{a} \cdot \vec{x} + b$, $\vec{a} \in \mathbb{Z}_p^k, b \in \mathbb{Z}_p \dots$ (2,1)-nezávislý
- lineární: $\mathbb{Z}_p \rightarrow [m]$, $h_{a,b}(x) = ((ax+b) \bmod p) \bmod m$, $a, b \in \mathbb{Z}_p, a \neq 0 \dots$ 1-univerzální
 přípustně-li i $a=0 \dots$ 2-univerzální, (2,1)-nezávislý
- multiply-shift: $[2^w] \rightarrow [2^k]$, $h_{a,b}(x) = (ax+b) \langle w-l, w \rangle \dots$ 2-univerzální
 $a, b \in [2^w]$, a liché
 $h_{a,b}(x) = (ax+b) \langle t-l, t \rangle$, $a, b \in [2^t]$, a liché, $t \geq w-l \dots$ 2-nezávislý
 ← bity na pozicích $w-l$ až $w-1$
- polynomy: $\mathbb{Z}_p \rightarrow [m]$, $h_{\vec{a}}(x) = \left(\sum_{i=0}^{d-1} a_i x^i \bmod p \right) \bmod m$, $\vec{a} \in \mathbb{Z}_p^d, p \geq dm \dots$ d -nezávislý
- tabelace: $[2^{kl}] \rightarrow [2^t]$, nagenerují tabulky $T_1 \dots T_k : [2^l] \rightarrow [2^t]$
 $h(x) = \bigoplus_{i=1}^k T_i[x \langle (i-1)l, il \rangle]$ } čas $O(c)$
 prostor $O(c \cdot U^{1/c})$
 slov délky t
 \dots je 3-nezávislé, ale není 4-nezávislé

Separované řetězce pro $m = \Theta(n) \dots$ zvětšují/zmenšují tabulku dle potřeby (jako náduky, pole)

- $E[C_t] \in O(1)$
 \uparrow délka řetězce v t -té přílohdce
 - $\text{var } C_t \in O(1)$
- } stačí c -universalita

Dle: $\text{var } C_t = E[C_t^2] - E[C_t]^2$

$$L = \frac{1}{m} \sum_{s=1}^m E[C_s^2] = \frac{1}{m} \sum_{i \neq j} \Pr[h(x_i) = h(x_j)] \leq \frac{cn^2}{m^2} \in O(1).$$

$\leq c/m$

- pro \vec{a} zcela náhodnou funkci h : $\max_t C_t \in \Theta\left(\frac{\log n}{\log \log n}\right)$ s velkou pravděpodobností
 $\exists \epsilon > 0 \forall n$ (až na konečné množiny výjimek)
 $\Pr[\text{platí pro } n] \geq 1 - n^{-c}$
- též pro $\Omega\left(\frac{\log n}{\log \log n}\right)$ -nezávislý systém
- a pro tabelaci

Lineární přidávání

- vše v jednom poli velikosti m
- Insert hledá první volnou pozici z $h(x), h(x)+1, h(x)+2, \dots \pmod{m}$
- Find se zastaví na první prázdné pozici
- Delete prvky "škrtná", když je škrtnutých moc, přebudujeme tabulku.

Příklad hodují podle posl. číslice, vkládám 75, 36, 14, 42, 24, 95, 17:

(2)

		42		14	75	36	24	95	17
0	1	2	3	4	5	6	7	8	9

Analýza pro $m \geq (1+\epsilon)n$: Insert trval v průměru

- $O(1/\epsilon^2)$ pro úplně náhodnou fci h
- $O(1/\epsilon)$ pro $(\log n)$ -nezávislý systém
- $O(1/\epsilon^{1.5})$ pro 5-nezávislý systém
- $\Omega(\log n)$ pro některé 4-nezávislé
- $\Omega(\sqrt{n})$ pro některé 2-nezávislé, $\Theta(\log n)$ pro Multiply-Shift
- $O(1/\epsilon^2)$ pro tabulaci

Perfektní hashování FKS (Fredman, Komlós, Szemerédi) ¹⁹⁸⁴

Chceme najít funkci $f: U \rightarrow [m]$, která na dané množině $X \subseteq \binom{U}{n}$ nemá žádné kolize, = je "perfektní"

• chceme lineární paměť, tedy $m = \Theta(n)$

• perf. fci je obecně málo: třeba pro $m=2n$ je

$\Pr[\text{náhodná } f \text{ je perfektní}] = \frac{(2n)!}{(2n)^n} \approx \frac{(2n)!}{n! (2n)^n} \approx \left(\frac{2n}{e}\right)^n / \left(\frac{n}{e}\right)^n (2n)^n = \frac{2^{2n} \cdot n^{2n} \cdot e^n}{e^{2n} n^n \cdot 2^n \cdot n^n} = \left(\frac{2}{e}\right)^n$

klesající mocnina: $a^k = a(a-1)(a-2)\dots(a-k+1)$

$n! \approx \left(\frac{n}{e}\right)^n$

Zvolíme c -univerzální systém \mathcal{H} do $[m]$, $X = \{x_1, \dots, x_n\}$

Pro obecné n, m počítáme $\mathbb{E}[K]$, kde $K = \#\{i, j: i \neq j \text{ a } h(x_i) = h(x_j)\}$ ← # kolizí / seich páří x_i, x_j

$$\mathbb{E}[K] = \sum_{\substack{i, j \\ i \neq j}} \mathbb{E}[\text{indikátor kolize } x_i, x_j] = \sum_{i \neq j} \Pr[h(x_i) = h(x_j)] \leq \frac{cn(n-1)}{m} \leq \frac{cn^2}{m}$$

$\leq c/m \approx \text{univerzality}$

• pro $m \geq 2cn^2$ je $\mathbb{E}[K] < \frac{1}{2}$

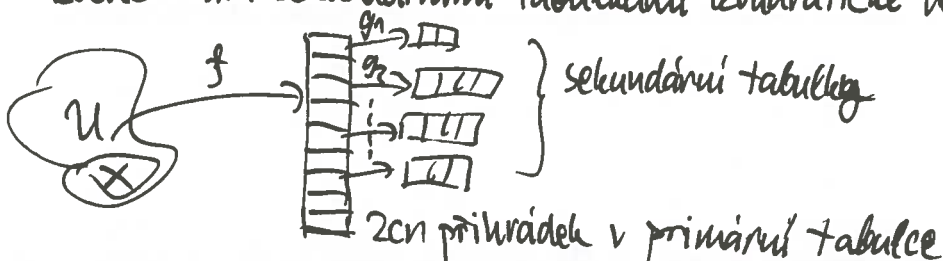
⇒ $\Pr[h \text{ perfektní}] = \Pr[K = 0] \geq \Pr[K < 1] > \frac{1}{2}$ z Markovovy nerovnosti

⇒ průměrný # pokusů, než seřenu perfektní h , je ≤ 2

⇒ průměrný čas $O(n)$... 1 pokus stihnou v $O(n)$ w.c.

• pro $m \geq 2cn$ je $\mathbb{E}[K] < \frac{n}{2}$... umím rychle sehnat h s $K < n$.

- kolize řeším sekundárními tabulkami kvadratické velikosti



- Pamatujeme si:
- parametry funkce f
 - pole $[2cn]$ primárních příhrádek, pro každou
 - počet prvků a_i
 - parametry funkce g_i ↗ perferitív
 - odkaz na sekundární tabulku velikosti $[2ca_i]$

Celková velikost sekundárních tabulek:

$$\sum_{i=1}^n [2ca_i] \leq \sum_{i=1}^n (2ca_i + 1) = n + 2c \left(\sum_{i=1}^n a_i^2 \right) \leq n + 2c \cdot 2n = (4c+1)n \in O(n).$$

$= \# \{i, j \mid f(x_i) = f(x_j)\} = n + \underbrace{\left(\sum_{i,j: i \neq j \text{ a } f(x_i) = f(x_j)} 1 \right)}_{(\# \text{kolikrát fce } f) < n}$

Celkově: prostor $O(n)$, konstrukce v průměrném čase $O(n)$.

Kukačkové hledání (Pagh, Rodler 2004)

Poradíme si 2 hledací funkce $f, g: U \rightarrow [m]$ (hledají buď do společné tabulky T nebo dvou oddělených)

Prvek x hledáme v $T[f(x)]$ a $T[g(x)]$ a uilide finde \rightarrow Find trvá $O(1)$ u.c.

- Insert: pokud $T[f(x)] = \emptyset$, vložíme x do $T[f(x)]$
- jinak prvek $T[f(x)]$ "vykročíme" a pokusíme se ho zalesovat podle opačné funkce ... tím možná vystřelíme nějaký další atd.
 - po $\sim \log n$ krocích vyhlásíme timeout a vše přehodíme s jinými f, g .
 - když se tabulka příliš zaplní, přehodíme.

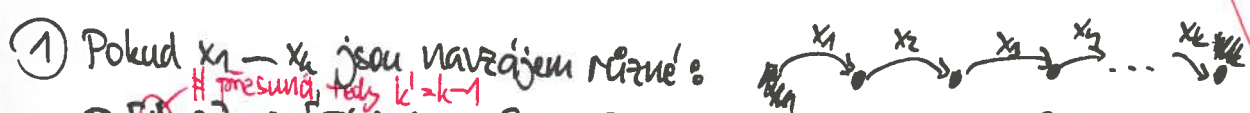
Delete: Skrtáním, občas přehodíme.

Analýza: Timeout nastavíme na $6 \log n$ kroků, předpokládáme $(6 \log n, 1)$ -nezávislý systém.

~~Nastavíme~~ $m \geq 4n$. ← stačilo by $m \geq 2(1+\epsilon)n$.

Vkládáme prvek x_1 , postupně vyhadzujeme prvky x_2, x_3, \dots, x_k .

Uvažujeme kukaččí graf: vrcholy jsou příhrádky, hrany $\{f(x_i), g(x_i)\}$ ← tedy odpovídají prvkům



① Pokud $x_1 \dots x_k$ jsou navzájem různé:

$$\Pr[k \geq 2] = \Pr[T(x_1) \text{ obsazeno}] = \Pr[\exists y f(x) = g(y) \vee f(x) = f(y)] \leq \frac{2n}{m} = \frac{1}{2}.$$

$$\Pr[k \geq t] \dots \text{indukcí} \leq \frac{1}{2^t}$$

$$E[\# \text{přesunů}] = \sum_t \Pr[\# \text{přesunů} \geq t] \leq \sum_t \frac{1}{2^t} \in O(1).$$

$\Pr[\text{timeout}] \leq \frac{1}{2^{6 \log n}} = n^{-6}$.

↑ každý obecný trik: pro nezápornou celočíselnou veličinu X je $E[X] = \sum_t \Pr[X \geq t]$... rovnou z definice $E[X]$.

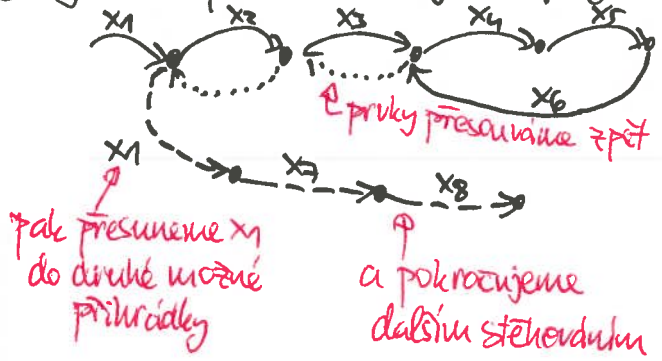
funguje i pro jiné konstanty, ale jsou limi

pro max. $6 \log n$ prvků se chová úplně nahodně

šikovno... ať na milión konstant

pro $(6 \log n, 1)$ -UZ. to platí přemě, jinak ať na konst.

② Vystřídáme prvek x_1 , který už byl jednou vystřídán



a) při 2. příchodu se už nic neopakuje
 • při timeoutu se
 • aspoň 1 typ hran se vyskytne aspoň $\frac{6 \log n}{3} = 2 \log n$ - krát
 podle ① se $\frac{1}{3}$ typ hran vyskytne v průměru $O(1)$ -krát
 $\rightarrow Pr[\text{timeout}] \leq 2^{-2 \log n} \in O(n^{-2})$

b) při 2. příchodu se také nějaký prvek opakuje \Rightarrow určitě nastane timeout

Počítejme konfigurace s t rekurzemi x_i pro fixní x_1 :

#konf. $\leq n^{t-1} \cdot t \cdot t \cdot t \cdot m^{t-1} \in O(n^{t-1} m^{t-1} t^3)$

volby ostatních x_i po kolika krocích skončí 1. cyklus kam se 1. cyklus vrátí kam se vrátí 2. cyklus hodnoty hes.-fci ve vrcholech

Pr dané konfigurace je $\left(\frac{2}{m^2}\right)^t = \frac{2^t}{m^{2t}}$
 ↑ 2 možné orientace hrany

Pr případu b) $\in O\left(\frac{n^{t-1} m^{t-1} t^3 2^t}{m^{2t}}\right) = O\left(\frac{n^{t-1} 2^{2t-2} n^{t-1} t^3 2^t}{2^{4t} n^{2t}}\right) = O\left(\frac{t^3}{2^t n^2}\right)$
 pro konkrétní t $m=4n$

Přes všechna t : $\leq \sum_{t=2}^{\infty} O\left(\frac{t^3}{2^t n^2}\right) = O\left(\frac{1}{n^2} \cdot \sum_{t=2}^{\infty} \frac{t^3}{2^t}\right) = O\left(\frac{1}{n^2}\right)$
 konverguje

Celkem: Pokud nepřehesovávané, Insert stojí průměrně $O(1)$.
 Přehesování nastane s psh $O\left(\frac{1}{n^2}\right)$ a způsobí n Insertů. + příchod řetězcem délky $\Theta(\log n)$

Tedy $E[T_{Ins}] = p \cdot (\Theta(\log n) + n \cdot E[T_{Ins}]) + (1-p) \cdot O(1)$

... řešením rekurence je $O(1)$.

Shrnutí: Pro $(6 \log n)$ -nezávislý systém umíme Find $O(1)$ n.c.
 Ins, Del $O(1)$ očekávaně amortizovaně

- Dále se ví:
- 6-nezávislost obecně nestačí
 - tabelacii hesování stací

umí se sestavit takový, který hesuje v čase $O(1)$ s $O(n^\epsilon)$ prostorem, ale je značně nepraktický [Siegelovy hes. funkce]