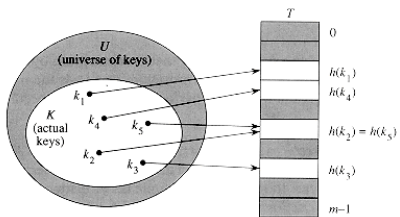


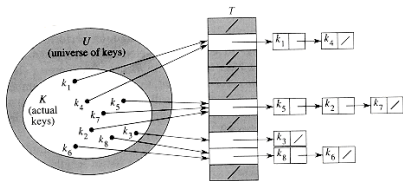
# Hešovanie

kuko

20.10.2020

Vybrané partie z dátových štruktúr





priemerná čas. zlož.



očekávaná čas. zlož.

rodina hešovacích funkcií  $\mathcal{H}$

$$h \in_R \mathcal{H}$$

$\mathcal{H}$  je univerzálna ak

$$\forall x_1 \neq x_2 : \Pr_{h \in_R \mathcal{H}} [h(x_1) = h(x_2)] \leq 1/n$$

$\mathcal{H}$  je 1-nezávislá, ak

$$\forall x, y : \Pr_{h \in_R \mathcal{H}} [h(x_1) = y] = 1/n$$

$\mathcal{H}$  je 1-nezávislá, ak

$$\forall x, y : \Pr_{h \in_R \mathcal{H}} [h(x_1) = y] = O(1/n)$$



$\mathcal{H}$  je  $k$ -nezávislá, ak

$$\forall \underbrace{x_1, \dots, x_k, y_1, \dots, y_k}_{\text{rôzne}} : \Pr_{h \in_R \mathcal{H}} [h(x_1) = y_1 \wedge \dots \wedge h(x_k) = y_k] = O(1/n^k)$$

$\forall x, y : \Pr_h[h(x) = y] = O(1/n)$  a pre rôzne  $x_1, \dots, x_k$  sú náhodné premenné  $h(x_1), \dots, h(x_k)$  skoro nezávislé

$\mathcal{H}$  je  $k$ -nezávislá, ak

$$\forall \underbrace{x_1, \dots, x_k}_{\text{rôzne}}, y_1, \dots, y_k : \Pr_{h \in_R \mathcal{H}} [h(x_1) = y_1 \wedge \dots \wedge h(x_k) = y_k] = O(1/n^k)$$

$\forall x, y : \Pr_h[h(x) = y] = O(1/n)$  a pre rôzne  $x_1, \dots, x_k$  sú náhodné premenné  $h(x_1), \dots, h(x_k)$  skoro nezávislé

- $x \mapsto (ax \bmod p) \bmod n$  je univerzálna
- $x \mapsto (ax) \gg (w - k)$  (ak  $n = 2^k$  a register má  $w$  bitov) je univerzálna
- $x \mapsto ((ax + b) \bmod p) \bmod n$  sú 2-nezávislé ( $a, b \in \mathbb{Z}_p$ ,  $a \neq 0$ )
- všeobecne  $x \mapsto ((a_k x^k + \dots + a_1 x + a_0) \bmod p) \bmod n$  sú  $k$ -nezávislé
- jednoduché tabulačné hešovanie je 3-nezávislé
  - vygenerujeme si tabuľky  $T_1, \dots, T_c$  veľkosti  $u^{1/c}$  s úplne náhodnou hešovacou fn.;
  - na  $x \in U$  sa pozeráme ako na vektor  $x = x_1 \dots x_c$  a  $h(x) = T_1(x_1) \oplus \dots \oplus T_c(x_c)$

- $x \mapsto (ax \bmod p) \bmod n$  je univerzálna
- $x \mapsto (ax) \gg (w - k)$  (ak  $n = 2^k$  a register má  $w$  bitov) je univerzálna
- $x \mapsto ((ax + b) \bmod p) \bmod n$  sú 2-nezávislé ( $a, b \in \mathbb{Z}_p$ ,  $a \neq 0$ )
- všeobecne  $x \mapsto ((a_k x^k + \dots + a_1 x + a_0) \bmod p) \bmod n$  sú  $k$ -nezávislé
- jednoduché tabulačné hešovanie je 3-nezávislé
  - vygenerujeme si tabuľky  $T_1, \dots, T_c$  veľkosti  $u^{1/c}$  s úplne náhodnou hešovacou fn.;
  - na  $x \in U$  sa pozeráme ako na vektor  $x = x_1 \dots x_c$  a  $h(x) = T_1(x_1) \oplus \dots \oplus T_c(x_c)$

- $x \mapsto (ax \bmod p) \bmod n$  je univerzálna
- $x \mapsto (ax) \gg (w - k)$  (ak  $n = 2^k$  a register má  $w$  bitov) je univerzálna
- $x \mapsto ((ax + b) \bmod p) \bmod n$  sú 2-nezávislé ( $a, b \in \mathbb{Z}_p, a \neq 0$ )
- všeobecne  $x \mapsto ((a_k x^k + \dots + a_1 x + a_0) \bmod p) \bmod n$  sú  $k$ -nezávislé
- jednoduché tabulačné hešovanie je 3-nezávislé
  - vygenerujeme si tabuľky  $T_1, \dots, T_c$  veľkosti  $u^{1/c}$  s úplne náhodnou hešovacou fn.;
  - na  $x \in U$  sa pozeráme ako na vektor  $x = x_1 \dots x_c$  a  $h(x) = T_1(x_1) \oplus \dots \oplus T_c(x_c)$

- $x \mapsto (ax \bmod p) \bmod n$  je univerzálna
- $x \mapsto (ax) \gg (w - k)$  (ak  $n = 2^k$  a register má  $w$  bitov) je univerzálna
- $x \mapsto ((ax + b) \bmod p) \bmod n$  sú 2-nezávislé ( $a, b \in \mathbb{Z}_p$ ,  $a \neq 0$ )
- všeobecne  $x \mapsto ((a_k x^k + \dots + a_1 x + a_0) \bmod p) \bmod n$  sú  $k$ -nezávislé
- jednoduché tabulačné hešovanie je 3-nezávislé
  - vygenerujeme si tabuľky  $T_1, \dots, T_c$  veľkosti  $u^{1/c}$  s úplne náhodnou hešovacou fn.;
  - na  $x \in U$  sa pozeráme ako na vektor  $x = x_1 \dots x_c$  a  $h(x) = T_1(x_1) \oplus \dots \oplus T_c(x_c)$

- $x \mapsto (ax \bmod p) \bmod n$  je univerzálna
- $x \mapsto (ax) \gg (w - k)$  (ak  $n = 2^k$  a register má  $w$  bitov) je univerzálna
- $x \mapsto ((ax + b) \bmod p) \bmod n$  sú 2-nezávislé ( $a, b \in \mathbb{Z}_p$ ,  $a \neq 0$ )
- všeobecne  $x \mapsto ((a_k x^k + \dots + a_1 x + a_0) \bmod p) \bmod n$  sú  $k$ -nezávislé
- jednoduché tabulačné hešovanie je 3-nezávislé
  - vygenerujeme si tabuľky  $T_1, \dots, T_c$  veľkosti  $u^{1/c}$  s úplne náhodnou hešovacou fn.;
  - na  $x \in U$  sa pozeráme ako na vektor  $x = x_1 \dots x_c$  a  $h(x) = T_1(x_1) \oplus \dots \oplus T_c(x_c)$

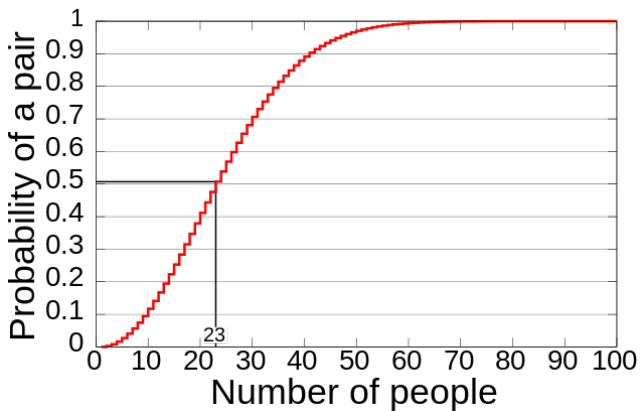
hešovanie so zretazením	$O(1)$	univerzálna rodina $\mathcal{H}$
lineárne sondovanie	$O(1/\varepsilon^2)$	$n \geq (1 + \varepsilon)m$
	$O(1/\varepsilon^2)$	5-nezávislá rodina $\mathcal{H}$
	$O(1/\varepsilon^2)$	tabulačné hešovanie



dá sa vyhľadávanie v  $O(1)$  v najhoršom prípade?  
čo keby sme vedeli všetky prvky dopredu?

## Perfektné hešovanie

## Narodeninový paradox



$$1 \times \left(1 - \frac{1}{365}\right) \times \left(1 - \frac{2}{365}\right) \times \cdots \times \left(1 - \frac{n-1}{365}\right)$$

$$e^x \approx 1 + x$$

$$\approx 1 \cdot e^{-1/365} \cdot e^{-2/365} \cdots e^{-(n-1)/365} = e^{-(n(n-1)/2)/365}$$

$$1 \times \left(1 - \frac{1}{365}\right) \times \left(1 - \frac{2}{365}\right) \times \cdots \times \left(1 - \frac{n-1}{365}\right)$$

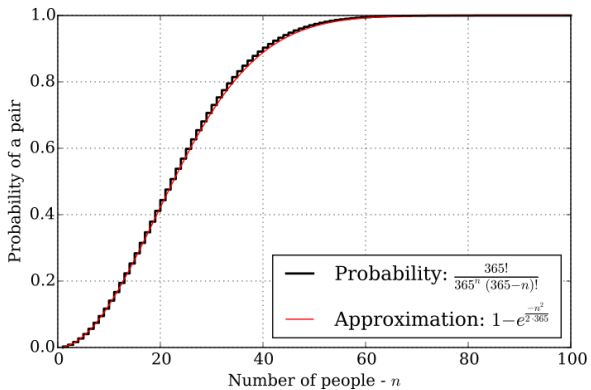
$$e^x \approx 1 + x$$

$$\approx 1 \cdot e^{-1/365} \cdot e^{-2/365} \cdots e^{-(n-1)/365} = e^{-(n(n-1)/2)/365}$$

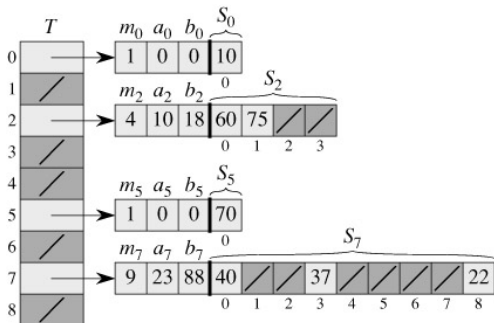
$$1 \times \left(1 - \frac{1}{365}\right) \times \left(1 - \frac{2}{365}\right) \times \cdots \times \left(1 - \frac{n-1}{365}\right)$$

$$e^x \approx 1 + x$$

$$\approx 1 \cdot e^{-1/365} \cdot e^{-2/365} \cdots e^{-(n-1)/365} = e^{-(n(n-1)/2)/365}$$



- (FKS'84) dve úrovne hešovania, namiesto zret'azenia tabuľku kvadratickej veľkosti





- $E[\#\text{kolízií}] = m^2 \cdot O(1/n) \leq 1/2$  pre dost' veľké  $n = \Theta(m^2)$
- žiadna kolízia s pp. aspoň  $1/2$  (Markov:  $\Pr[X \geq k\mu] \leq 1/k$ )
- $E[\sum_t C_t^2] = \sum_t E[C_t^2] = O(E[\#\text{kolízií}]) + O(n)$
- $= m^2 \cdot O(1/n) = O(n)$  pre  $m = \Theta(n)$

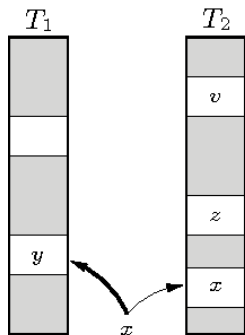
- $E[\#\text{kolízií}] = m^2 \cdot O(1/n) \leq 1/2$  pre dost' veľké  $n = \Theta(m^2)$
- žiadna kolízia s pp. aspoň  $1/2$  (Markov:  $\Pr[X \geq k\mu] \leq 1/k$ )
- $E[\sum_t C_t^2] = \sum_t E[C_t^2] = O(E[\#\text{kolízií}]) + O(n)$
- $= m^2 \cdot O(1/n) = O(n)$  pre  $m = \Theta(n)$

- $E[\#\text{kolízií}] = m^2 \cdot O(1/n) \leq 1/2$  pre dost' veľké  $n = \Theta(m^2)$
- žiadna kolízia s pp. aspoň  $1/2$  (Markov:  $\Pr[X \geq k\mu] \leq 1/k$ )
- $E[\sum_t C_t^2] = \sum_t E[C_t^2] = O(E[\#\text{kolízií}]) + O(n)$
- $= m^2 \cdot O(1/n) = O(n)$  pre  $m = \Theta(n)$

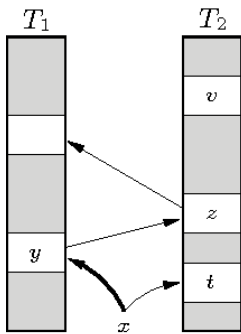
- $E[\#\text{kolízií}] = m^2 \cdot O(1/n) \leq 1/2$  pre dost' veľké  $n = \Theta(m^2)$
- žiadna kolízia s pp. aspoň  $1/2$  (Markov:  $\Pr[X \geq k\mu] \leq 1/k$ )
- $E[\sum_t C_t^2] = \sum_t E[C_t^2] = O(E[\#\text{kolízií}]) + O(n)$
- $= m^2 \cdot O(1/n) = O(n)$  pre  $m = \Theta(n)$

## Kukučie hešovanie

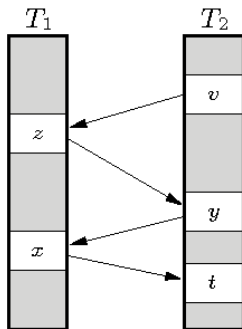


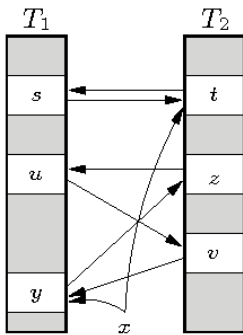


- máme 2 tabuľky  $A$ ,  $B$  dĺžky  $2m$  a 2 hešovacie funkcie  $f, g$
- hľadanie: prvok  $x$  je vždy v  $A[f(x)]$  alebo  $B[g(x)]$









## Riešenie kolízií zret'azením

Koľko prvkov sa zahešuje na tú „najvyťaženejšiu“ pozíciu?

$P_i = \#$ prvkov, ktoré sa zahešujú na pozíciu  $i$

koľko je  $\max P_i$ ?

- nech  $p = 1/n = \Pr[\text{prvok zahešujeme na konkrétnu pozíciu}]$
- nech  $q = 1 - p = \Pr[\text{prvok zahešuje inde}]$

$$\Pr[P_i = k] = \binom{n}{k} p^k q^{n-k}$$

- nech  $p = 1/n = \Pr[\text{prvok zahešujeme na konkrétnu pozíciu}]$
- nech  $q = 1 - p = \Pr[\text{prvok zahešuje inde}]$

$$\Pr[P_i = k] = \binom{n}{k} p^k q^{n-k}$$

- nech  $p = 1/n = \Pr[\text{prvok zahešujeme na konkrétnu pozíciu}]$
- nech  $q = 1 - p = \Pr[\text{prvok zahešuje inde}]$

$$\Pr[P_i = k] = \binom{n}{k} p^k q^{n-k}$$



$$\Pr[P_i = k] = \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k}$$

$$\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$$

$$1/e \leq \left(1 - \frac{1}{n}\right)^{n-k} \leq 1$$

$$\Pr[P_i = k] = \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k}$$

$$\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$$

$$1/e \leq \left(1 - \frac{1}{n}\right)^{n-k} \leq 1$$

$$\Pr[P_i = k] = \binom{n}{k} \left(\frac{1}{n}\right)^k \left(1 - \frac{1}{n}\right)^{n-k}$$

$$\left(\frac{n}{k}\right)^k \leq \binom{n}{k} \leq \left(\frac{ne}{k}\right)^k$$

$$1/e \leq \left(1 - \frac{1}{n}\right)^{n-k} \leq 1$$

$$\Pr[P_i = k] = \underbrace{\binom{n}{k}}_{\leq (ne/k)^k} \left(\frac{1}{n}\right)^k \underbrace{\left(1 - \frac{1}{n}\right)^{n-k}}_{\leq 1} \leq \left(\frac{ne}{k}\right)^k \cdot \frac{1}{n^k} = \left(\frac{e}{k}\right)^k$$

$$\Pr[P_i = k] = \underbrace{\binom{n}{k}}_{\geq (n/k)^k} \left(\frac{1}{n}\right)^k \underbrace{\left(1 - \frac{1}{n}\right)^{n-k}}_{\geq 1/e} \geq \left(\frac{n}{k}\right)^k \cdot \frac{1}{n^k} \cdot (1/e) = \frac{1}{ek^k}$$

$$1/ek^k \leq \Pr[P_i = k] \leq (e/k)^k$$

$$\Pr[P_i = k] = \frac{1}{k^{\Theta(k)}} = \frac{1}{e^{\Theta(k \log k)}}$$

- pre  $k = \ln n / 3 \ln \ln n$  je aspoň  $1/e^{\sqrt[3]{n}}$
- $\implies$  očakávaný počet pozícií, ktoré prekročia  $k$  je  $\Omega(n^{2/3})$

$$1/ek^k \leq \Pr[P_i = k] \leq (e/k)^k$$

$$\Pr[P_i = k] = \frac{1}{k^{\Theta(k)}} = \frac{1}{e^{\Theta(k \log k)}}$$

- pre  $k = \ln n / 3 \ln \ln n$  je aspoň  $1/e^{\sqrt[3]{n}}$
- $\implies$  očakávaný počet pozícií, ktoré prekročia  $k$  je  $\Omega(n^{2/3})$

$$1/ek^k \leq \Pr[P_i = k] \leq (e/k)^k$$

$$\Pr[P_i = k] = \frac{1}{k^{\Theta(k)}} = \frac{1}{e^{\Theta(k \log k)}}$$

- pre  $k = \ln n / 3 \ln \ln n$  je aspoň  $1/e^{\sqrt[3]{n}}$
- $\implies$  očakávaný počet pozícií, ktoré prekročia  $k$  je  $\Omega(n^{2/3})$

$$1/ek^k \leq \Pr[P_i = k] \leq (e/k)^k$$

$$\Pr[P_i = k] = \frac{1}{k^{\Theta(k)}} = \frac{1}{e^{\Theta(k \log k)}}$$

- pre  $k = \ln n / 3 \ln \ln n$  je aspoň  $1/e^{\sqrt[3]{n}}$
- $\implies$  očakávaný počet pozícií, ktoré prekročia  $k$  je  $\Omega(n^{2/3})$



$$\Pr[P_i \geq \ell] = \sum_{k=\ell}^n \underbrace{(e/k)^k}_{\leq (e/\ell)^k} \leq (e/\ell)^\ell (1 + e/\ell + (e/\ell)^2 + \dots) = \Theta((e/\ell)^\ell)$$

$$\Pr[P_i \geq \ell] \leq (e/\ell)^\ell \cdot [1/(1 - e/\ell)]$$

- pre  $\ell = \lceil (3 \ln n) / \ln \ln n \rceil$  dostaneme  $\Pr \leq 1/n^2$

$$\sum_{k=1}^n p_k O(k^2)$$

$$p_k \sim 1/\exp(k)$$

$$\implies \sum p_k O(k^2) = O(1)$$

$$\Pr[P_i \geq \ell] \leq (e/\ell)^\ell \cdot [1/(1 - e/\ell)]$$

- pre  $\ell = \lceil (3 \ln n) / \ln \ln n \rceil$  dostaneme  $\Pr \leq 1/n^2$

$$\sum_{k=1}^n p_k O(k^2)$$

$$p_k \sim 1/\exp(k)$$

$$\implies \sum p_k O(k^2) = O(1)$$