

## Hashování vektorů

- chceme hashovat  $\vec{x} = (x_0, \dots, x_{d-1}) \in \mathcal{U}^d$

Konstrukce hash. fun. z polynomu.

- vybereme prvočíslo  $p \geq |\mathcal{U}|$

- vybereme náhodně  $a \in [p]$

- $h_a(x) = \left( \sum_{i=0}^{d-1} x_i a^i \right) \bmod p$

Věta: Systém  $\mathcal{H} = \{h_a \mid a \in [p]\}$  je  $d$ -univerzální.

Dk:

- uvažme  $\vec{x} \neq \vec{y} \in \mathcal{U}^d$

- $\sum_{i=0}^{d-1} x_i a^i = \sum_{i=0}^{d-1} y_i a^i \iff \sum_{i=0}^{d-1} (x_i - y_i) a^i = 0 \iff a$

je kořen polynomu  $h(a) = \sum_{i=0}^{d-1} (x_i - y_i) a^i$

- zde jsou  $x_0 - y_0, x_1 - y_1, \dots, x_{d-1} - y_{d-1}$  koeficienty polynomu

- základní věta algebry: Nemulový polynom stupně  $d$  má  $\leq d$

kořenů

$\rightarrow$  pro  $\leq d$  hodnot  $a$  se polynom  $h(a)$  vyhodnotí

na 0

$\rightarrow \Pr [h_a(x) = h_a(y)] \leq \frac{d}{p} \implies d$ -univerzální  $\square$

Jak rychle lze vyhodnotit  $h_a(\vec{x})$ ?

$$\begin{aligned} \bullet \quad x_0 a^0 + x_1 a^1 + x_2 a^2 + \dots + x_{d-1} a^{d-1} &= x_0 + a(x_1 + x_2 a + x_3 a^2 + \dots + x_{d-1} a^{d-2}) \\ &= x_0 + a(x_1 + a(x_2 + x_3 a + \dots)) \end{aligned}$$

$\leftrightarrow$   $d$  násobení a sčítání.

• tohle je tzv. Hornerovo schéma

Příkladek  $p \geq 1001$

• ... je celkem drsný

• můžeme zvedkovat  $p$  na nějakí  $m$  jak jsme zvyklí?

• ano, když výslednou hodnotu  $h_a(\vec{x})$  hashujeme novým  
2-nezávislým (např.  $ax+b \pmod{m}$ )

Lem: Necht'  $\mathcal{F}$  je  $c$ -universální hash. systém z  $\mathcal{U}$  do  $[r]$ .

Necht'  $\mathcal{G}$  je  $(2, d)$ -nezávislý hash. systém z  $[r]$  do  $[m]$ .

Pak je jejich složení  $\mathcal{H} = \mathcal{F} \circ \mathcal{G}$  (nejdříve aplikujeme  $\mathcal{G}$ , pak  $\mathcal{F}$ )

je  $(2, c')$ -nezávislý pro  $c' = \left(\frac{cm}{r} + 1\right) d$ .

• pro hashování polynomů: pokud  $p \geq 4dm$ , pak složení polynomů

s  $ax+b \pmod{m}$  je  $(2, \frac{5}{2})$ -nezávislé.

Dle lemmatu:

• co chceme? Pro  $x \neq y$  e $\mathcal{U}$ ,  $i, j \in [m]$  chceme

$$\Pr_{\substack{f \in \mathcal{F} \\ g \in \mathcal{G}}} [g(f(x)) = i \wedge g(f(y)) = j] \leq \frac{c'}{m^2}$$

•  $g \in \mathcal{G}$  je 2-nezávislý, hodovo z definice?

• ne,  $f(x)$  není zameněné různě od  $f(y)$ , abychom mohli aplikovat nezávislost ☹

• ale  $f \in \mathcal{F}$  je  $c$ -universální, takže jsou různě skoro vždy ☺

• ozn:

•  $M$ :  $\text{jev } g(f(x)) = i \wedge g(f(y)) = j$

•  $C$ :  $\text{jev } f(x) = f(y)$

$$\Pr[M] = \Pr[M \cap \bar{C}] + \Pr[M \cap C]$$

$$= \Pr[M | \bar{C}] \Pr[\bar{C}] + \Pr[M | C] \Pr[C]$$

•  $\Pr[M | \bar{C}]$

• zde  $f(x) \neq f(y) \Rightarrow$   $\stackrel{\text{2-nezávislost}}{\leq} \frac{d}{m^2}$

•  $\Pr[\bar{C}] \leq 1$  stačí, takže je těch  $\dots \underbrace{+1}_d$  ve znění

•  $\Pr[M | C]$

•  $f(x) = f(y) =: z$  (došla příměrka, proto  $z$ )

• p-táme se na  $M|C \Leftrightarrow g(x^i) = i \wedge g(x^j) = j$

•  $i \neq j \Rightarrow P_r[M|C] = 0$

•  $i = j \Rightarrow P_r[M|C] \leq \frac{d}{m}$  pokud vidíme cennou  
 konvid přes  $x \in \{0, \dots, m-1\}$ .

•  $P_r[C] \leq \frac{c}{r}$  z  $c$ -universality  $\mathcal{F}$ .

$$\rightarrow P_r[M] \leq \frac{d}{m^2} \cdot 1 + \frac{d}{m} \cdot \frac{c}{r} = \frac{1}{m^2} \left( d + \frac{cdm}{r} \right) = \frac{d}{m^2} \left( 1 + \frac{cm}{r} \right)$$

• tabulka taky platí pro  $d' = (c+1)d$

•  $g$  hashuje z  $[r]$  do  $[m]$

$\rightarrow r \geq m \Rightarrow \frac{m}{r} \leq 1$  a dosadíme za  $d'$

## Rabin-Karp alg.

• chceme hledat podřetězec  $\sigma$  v textu  $S$  a máme „malé“  
 prostorn

• nepří. KMP, BM, Suffix Array / Tree využívají paměť

$\Omega(|\sigma|)$

• Zajímá nás  $S[i:i+|\sigma|] = \sigma$  pro  $\approx |\Sigma|$  různých  $i$ :

• kontrola všech pozic je dražší  $\rightarrow O(|\Sigma| \cdot |\sigma|)$

• budeme kontrolovat jen ty pozice, kde  $h(S[i:i+|\sigma|]) = h(\sigma)$

pro vhodnou hash fci  $h$ .

- na první pohled je tohle ještě horší než nánní celý  
nášat musíme navíc hashovat, což stojí  $\Omega(101)$  času.
- ledará bychom uměli počítat hash fci v "konst" čase :-)

• S 

$h(\square\square\square\square\square\square)$

$h(\square\square\square\square\square\square)$

$h(\square\square\square\square\square\square)$

⋮

- chceme každou další  $h$  kromě prvního přičítat v konst. čase

- používáme polynomiální systém z předchozích 4 stran

- máme  $S[j]a^0 + S[j+1]a^1 + S[j+2]a^2 + \dots + S[j+101-1]a^{101-1}$

- chceme  $S[j+1]a^0 + S[j+2]a^1 + \dots$

- tedy stačí odečíst  $S[j]a^0$ ,  
vydělit  $a$ ,  
přičíst  $S[j+101-1]a^{101-1}$ .  
}  $O(1)$  času

a máme to

• alternativně

$$S[j]a^{l-1} + S[j+1]a^{l-2} + \dots + S[j+l-1]a^0$$

a chceme

$$S[j+1]a^{l-1} + \dots + S[j+l-1]a^1 + S[j]a^0$$

• stačí odečíst  $S[j]a^{l-1}$ ,

vynásobit  $a$

přičíst  $S[j]a^0$

• pokud je dělení o hodnotě dražší než násobení, tak alternativa je rychlejší.

• to je většinou pravda, zejména kvůli HW paralelizaci násobení

• pro dělení existuje nějaký trik, ale paralelně

• na starších Pentiumech (Intel) měl těžký bug a

• Intel si ho zapatentoval :-)

Analýza

• první hash stojí  $O(l)$  a každý další  $O(1)$

→ celkem  $O(|S| + l)$  ... ojt :-)

• ozn  $N$  počet pozic i t.č.  $h(S[i: i+|S|-1]) = h(\sigma)$

• celková složitost je

$$O(|S| + |S| + |S| \cdot N)$$

pokud stačí zahltášit první výskyt,  
jinak tam je ještě  
# výskytů  $\cdot |S|$

• kolik je  $N$ ?

• polynomy jsou  $d$ -nezavislý hash systém

$\Rightarrow$  pro  $\sigma \neq \tau$  platí

$$Pr[h(\sigma) = h(\tau)] \leq \frac{d}{p}$$

$$\bullet \mathbb{E}[N] \leq \frac{d}{p} |S|$$

• zvolíme-li  $p \geq d|S|$ , máme vyhráno.

---

Odsud dolů to jde  
možná zlepšit, nestihám

☹