

Datové struktury 2 - LS 2015/2016 - Martin Mareš

1

- Poprvé podle nové akreditace - dost jiná přednáška než dříve (snad je to zrušena k lepšímu?)
- Kontakty: mj@ucw.cz, <http://mj.ucw.cz/lyuka/ds2/> → repríza 2017/2018 (jen mírné úpravy)
- Cvičení nejsou, ale můžete si domluvit konzultaci.
- Přibližný plán:
 - statické slovníky
 - celočíselné DS
 - dolní odhady
 - cache-oblivious DS
 - DS pro stromy a obecné grafy
 - geometrické DS
 - ústorné DS
 - streaming algoritmy

Literatura dosti kusá, budeme přidávat odkazy na web + scanovalé stránky + videozáznamy z r. 2016

• Požadavky ke zkoušce: znát odpřednesené, umět to aplikovat a upravovat

VÝPOČETNÍ MODEL

- kdybychom studovali poly. vs. exp, na modelu nezáleží
- u DS ale potřebujeme rozlišovat $\log n / \log \log n / O(1) / \dots \Rightarrow$ model musíme specifikovat
- Budeme používat Word-RAM (neřekneme-li jinak)

- w -bitová celá čísla - slova
- na slovech umíme počítat v konstantním čase (jako v Cěčku ...)
- aritmetika: $+, -, *, /, \%$
- logické operace: $\&, |, ^, \ll, \gg, \sim$
- porovnávání: $=, <, >$
- paměť je pole slov indexované slovy \rightarrow potřebujeme $w \geq \log_2 n$
- vstup a výstup předáváme v paměti
- čas = # provedených instrukcí
- prostor = rozptyl mezi min. a max. adresou políček paměťové banky

BůHO dokážeme počítat i s $O(w)$ -bit. slovy

všechny logaritmy budou nadále implicitně dvojkové

STATICKE MAZEŘINY

universum (treba slova RAMu)

• Chceme pro n -prvkovou $S \subseteq U$ vybudovat DS, která bude umět rychle odpovídat na dotazy "x ∈ S?"

	Build	Member	
• Co už umíme:	$O(n \log n)$	$O(\log n)$	vyhledávací strom [v porovnávacím modelu nete lépe]
	$O(n)$ průměrně	$O(1)$ w.c.	kukaččí hesování (potřebuje $\log n$ -nezávislou rodinu fci)
• Ukážeme:	$O(n)$ průměrně	$O(1)$ w.c.	perfektní hesování FKS (stati & nezávislost)
	$O(n)$ průměrně	$O(1)$ w.c.	... jiný přístup...
	$O(n \log n)$ w.c.	$O(1)$ w.c.	derandomizace

PERFEKTNÍ HESOVÁNÍ

FKS = Fredman, Komlos, Szemerédi 1984

(2)

Opakování: Def: Systém \mathcal{H} hesovacích funkcí $U \rightarrow [m]$ je c-universální ($c > 0$)
 $\equiv \forall x, y \in U, x \neq y: \Pr_{h \in \mathcal{H}} [h(x) = h(y)] \leq c/m$.

Navíc obvykle chceme "hezkou parametrizaci" - tedy aby náhodný výběr $h \in \mathcal{H}$ šlo provést rovnoměrně náhodným výběrem $O(1)$ parametrů a pomocí nich pak $h(x)$ vyhodnotit v čase $O(1)$.

Pro $S \subseteq U$ a funkci $h \in \mathcal{H}$ počítáme kolize: $\{x, y\} \in \binom{S}{2}$ t.j. $h(x) = h(y)$.

nastane s $\text{řstí} \leq \frac{c}{m}$

Lemma $\mathbb{E}[\# \text{kolizí}] = \sum_{\{x, y\}} \mathbb{E}[C_{xy}] \leq \binom{n}{2} \cdot \frac{c}{m} \leq \frac{n^2 \cdot c}{2m}$
 ↑
 indikátor kolize

Pro $m = \lceil n^2 \cdot c \rceil$ je $\mathbb{E}_{h \in \mathcal{H}} [\# \text{kolizí}] < \frac{1}{2}$, takže podle Markovovy nerovnosti je

$$\Pr_h [h \text{ koliduje na } S] = \Pr [\# \text{kolizí} > 2 \cdot \mathbb{E}[\# \text{kolizí}]] < \frac{1}{2}$$

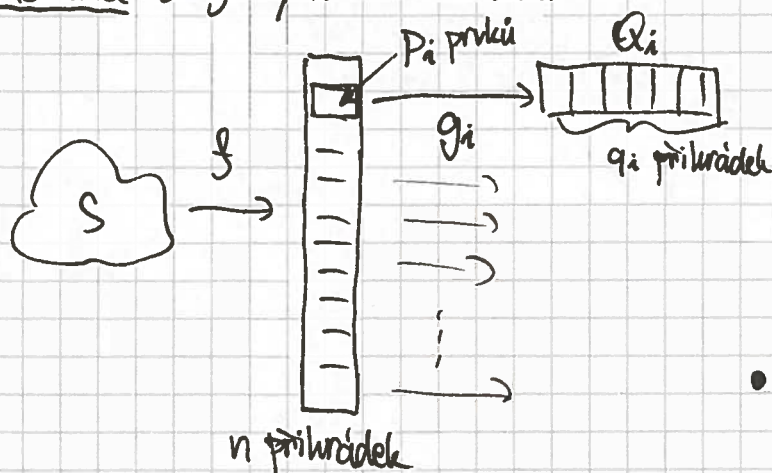
⇒ proto budeme-li volit h náhodně z \mathcal{H} , tak dlouho, než objevíme nějakou perfektní, bude to trvat průměrně ≤ 2 pokusy:

Lemma (O očekávání): Čekáme-li na událost, která nastane s řstí p , pak $\mathbb{E}[\# \text{pokusy, než se dočkáme}] = 1/p$.

... jeden pokus trvá $O(n)$ [pokud kolize detekujeme přívrázkovým tříděním], takže v čase průměrně $O(n)$ perfektní f ci najdeme.

... jenže potřebujeme kvadraticky velkou tabulku ⇒ inicializace stojí $O(n^2)$.

Konstrukce dvojstupňového hesování:



$f \in \mathcal{H}$ hešuje do n přívrádek

$$\mathbb{E}[\# \text{kolizí}] \leq \frac{cn}{2}$$

... umíme efektivně najít f , pro niž $\# \text{kolizí} < cn$.

g_i hešuje do $[q_i]$ pro $q_i = \lceil \frac{P_i^2}{c} \rceil$, takže umíme najít perfektní f ci

Dokážeme, že celková velikost všech Q_i je $O(n)$:

$$\sum_{i=1}^n q_i \leq n + c \cdot \sum_i P_i^2 \leq n + c \left(\sum_i P_i \right) + c \left(\sum_i \binom{P_i}{2} \right) \in O(n)$$

↑
 $P_i \cdot 2 \cdot \binom{P_i}{2}$
 # kolizí v i -té přívrádce

všech kolizí pro $f \leq \frac{cn}{2}$

Spotřeba paměti:

- parametry $f \dots O(n)$
- parametry $g \dots O(n)$
- tabulka indexovaná f (pointery na Q_i) $\dots O(n)$
- tabulky $Q_i \dots O(n)$

} celkem $O(n)$

Čas na konstrukci:

- průměrně $O(n)$ na volbu f
- průměrně $O(q_i)$ na volbu g_i

} celkem průměrně $O(n)$

Čas na dotaz:

- výpočet f
- vhlédnutí do libvni tabulky pro pointer a param. g_i
- výpočet g_i
- vhlédnutí do Q_i

} $O(1)$ w.c.

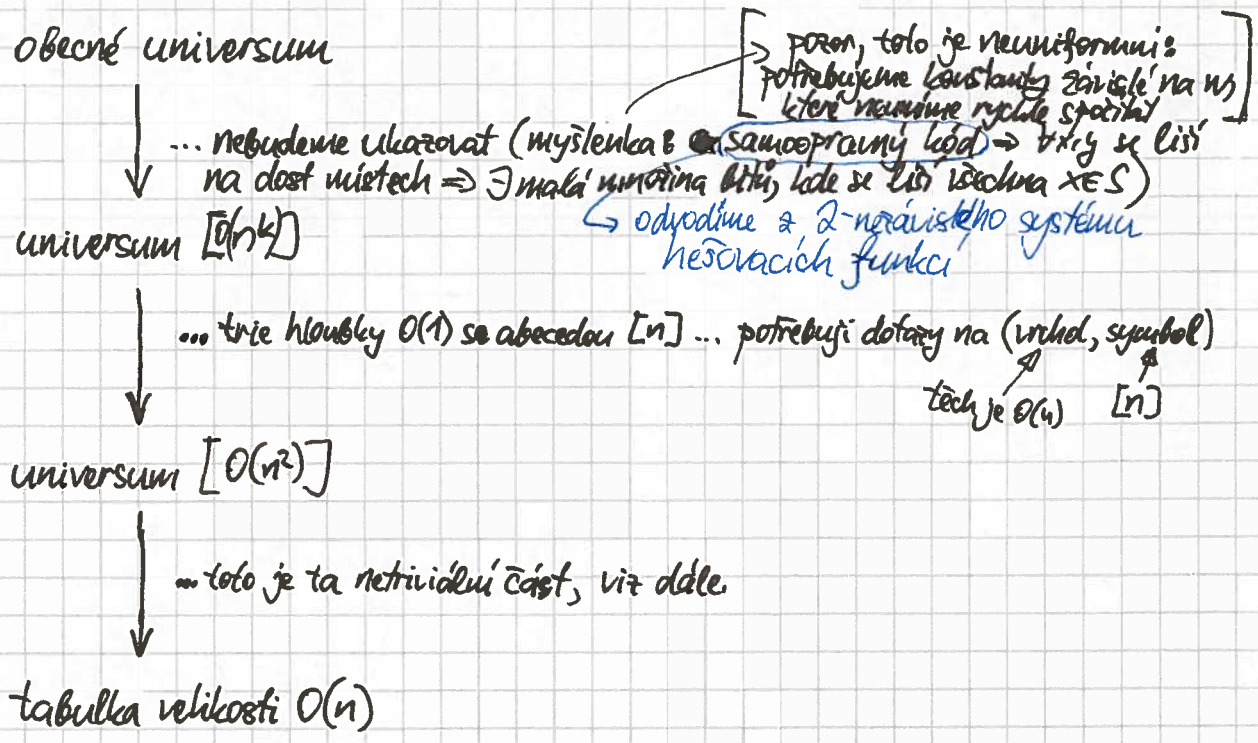
Poznámka: \exists dynamizace s čosem $O(1)$ průměrně amortizovaně na Ins/Del a $O(1)$ w.c. na dotaz.

Odkaz: Trídění reálných čísel vybiraných rovnoměrně náhodně z $[0,1]$.

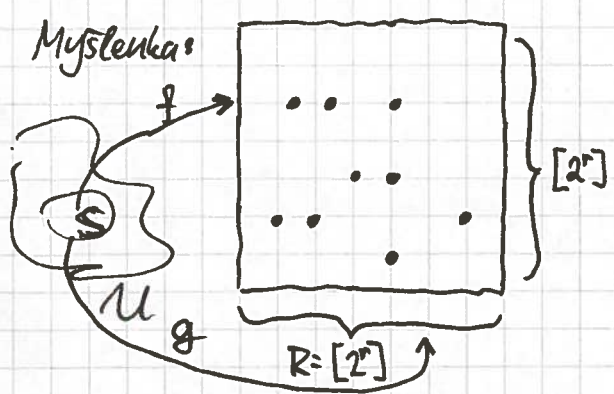
- Rozdělíme $[0,1]$ na n příhradek
 - v $O(n)$ rozmístíme čísla do příhradek ... $\in [\#kolizi] < \frac{n^2}{2}$
 - v každé příhradce dotřídíme bublinkově ... trvá to $O(\text{velikost příhrádky}^2)$, což se seče na
- } průměrně $O(n)$

DETERMINISTICKÉ SLOVNÍKY Hagerup, Miltersen, Pagh 2000

- Nejprve ukážeme ~~deternistickou~~ ^{randomizovanou} verzi, pak ji derandomizujeme.
- Skládáme několik transformací: (všechno to jsou prosté funkce)



• Notace: nerozlišujeme mezi číly z $[O(n^2)]$ a bit. řetěci z $\{0,1\}^{2 \log n + O(1)}$



Prvky z S odpovídají bodům ~~na~~ v mřížce $[2^r] \times [2^r]$.
 Chceme je transformovat tak, aby v řádce byl max. 1 bod.

Krok: $(i, j) \rightarrow (i, j \oplus a_i)$
 ... a pak totéž znovu ve sloupcích

položdě řádce n-krát, sníž #kolizi
omezíme #kolizi ve sloupcích funkcí #kolizi v řádkách

Def: Dvojice funkcí (f, g) z U do R je q -dobrá (pro $q \geq 0$), pokud f má na S nejvýše q kolizí a $x \mapsto (f(x), g(x))$ je na S prostá.

Lemma: Necht' (f, g) je q -dobrá a $r \geq \log n + 1$.
 Pak $\exists a_0 \dots a_{2^r-1} \in R$ t.z. $(x \mapsto g(x) \oplus a_{f(x)}, f(x))$ je q' -dobrá pro $q' = \lfloor \frac{2^{3r}}{q} \rfloor \cdot n$

posunu řádky a transponuji mřížku

Navíc všechna a_i lze pro dané S, f, g spočítat randomizovaně v očekávaném čase $O(n)$ a worst-case prostoru $O(n)$.

za předpokladu, že f, g máme vyhodnotit v konst. čase
jelikož U je velké $O(n)$, je $r \leq \log n + O(1)$

Využití: Mějme nějakou $S \subseteq \{0, 1\}^w$. Zvolíme $r > \max(w/2, \log n + 3)$.

0. krok: (f, g) rozkládají $x \in S$ na horních a dolních r bítí (s překryvem, je-li třeba)

- (f, g) je jisté prostá na S
- f má na S nejvýše $\binom{n}{r} < n^2$ kolizí.

par (f, g) je n^2 -dobrá

Lemma (zde je potřeba omezení $q' \leq n$ z tvrzení lemmatu)
 1. krok: (f', g') ... jelikož $2^{3r} < \frac{1}{n}$, musí tento pár být $< n$ -dobrá

Lemma (... a zde naopak druhá část...)
 2. krok: (f'', g'') ... < 1 -dobrá, takže 0 -dobrá $\Rightarrow f''$ je prostá na S .

Výpočet hes. funkce:

1. Rozdělíme x na ~~2^r~~ r -bitové
2. ~~$b \leftarrow b \oplus a_p$~~ $q \leftarrow q \oplus a_p$
3. $p \leftarrow p \oplus b_q$
4. Vyjdáme výsledek

čas $O(1)$

Prostor: Tabulky pro a, b
 + finální hesovací tabulka

$O(n)$ slov

Dle lemmatu: Nejprve očíslovujeme řádky od nejplněnějšího:

$S_i := \{x \in S \mid f(x) = v_i\}$,
 kde $v_1 \dots v_{2^r}$ je permutace na $R = [2^r]$
 taková, že $|S_1| \geq |S_2| \geq \dots \geq |S_{2^r}|$

V tomto pořadí řádkům přidělujeme jejich ~~číslo~~ a_{v_i} .

• Necht' jsme již zpracovali $S_{<i}$:= $S_1 \cup \dots \cup S_{i-1}$ a přidáváme S_i .

Vybereme $a_{v_i} \in R$ náhodně, počítáme, kolik vzniklo nových kolizí:

$$E[\#NK] = |S_{<i}| \cdot |S_i| \cdot 2^r$$

za $O(1)$ pokusů najdu a_{v_i} ,
 pro které $\#NK \leq \lfloor 2 \cdot E[\#NK] \rfloor$ ← je to celé číslo
 $|S_{<i}| \cdot |S_i| \cdot 2^{r-1}$

dvojic $x \in S_{<i}, y \in S_i$
 t.č. $g(x) \oplus a_{g(x)} = g(y) \oplus a_{g(y)}$
 $a_{g(y)}$ pro nějaké $j < i$ = a_{v_i}
 $a_{v_i} = g(x) \oplus g(y) \oplus a_{g(x)}$
 ... to nastane s psti $\frac{1}{2^r}$

• Jak to udělat efektivně? • Udržujeme M_t := kolik bodů jsme už umístili do t -tého sloupce
 ... tedy $|\{x \in S_{<i} \mid g(x) \oplus a_{g(x)} = t\}|$.

• Na počátku seřadíme S podle $(f(x), g(x))$ lexikograficky ... $O(n \log n + 2^r)$
 ↳ poradi $S_1 \dots S_{2^r}$ dalším přírodním tříděním

• Pro každou S_i zvolíme a_{v_i} , pak pro $x \in S_i$ spočítáme pomocí M_x kolize
 ... až najdeme správné a_{v_i} , přepočítáme M_x

$\left. \begin{matrix} O(|S_i|) \\ \times O(1) \text{ přírodně} \\ \text{pokusy} \end{matrix} \right\} O(|S_i|)$
 ↳ v součtu přes všechna i $O(n) + 2^r$ průměrně

2^n -krát

• Jak dobný pár jsme vyrobili?

Pro původní pár (f, g) : $\# \text{kolizí fce } f = \sum_i \binom{|S_i|}{2} \leq q$

Pro nový pár počítáme kolize fce $x \mapsto g(x) \oplus a_{g(x)}$:

$$\# \text{kolizí} \leq \sum_{i=1}^{2^r} \underbrace{[2^{1-r} \cdot |S_i| \cdot |S_{<i}|]}_{\text{už víme}} \leq \sum_{i=1}^t [2^{1-r} \cdot |S_i| \cdot |S_{<i}|] \leq \sum_{i=1}^t [2^{3-r} \cdot \sum_{j < i} \binom{|S_j|}{2}] \leq n \cdot [2^{3-r} \cdot q]$$

shora omezeno tímto

vynecháme členy,
 pro které $|S_i| \leq 1$

... takže $|S_{<i}| \leq 2^{n-1}$
 díky volbě n ,
 takže $\lfloor \dots \rfloor = 0$

Spokud $|S_i| = 0$, je celý součin triviálně 0,
 pro $|S_i| = 1$ víme $|S_{<i}| < 2^{n-1} \Rightarrow \lfloor \dots \rfloor = 0$.

$$|S_i| \cdot \sum_{j < i} |S_j| = \sum_{j < i} |S_i| \cdot |S_j| \leq \sum_{j < i} |S_j|^2 \leq 4 \cdot \binom{|S_j|}{2}$$

jelikož $|S_j| \geq 2$

$2^{1-r} \leq \frac{1}{n}$, takže $\dots \leq |S_i|$
 $\rightarrow \sum_i \dots \leq n$
 (to je druhá část min(...) z horní lemmatu)

Derandomizace

• Obecný trik: Hledáme α : $T(\alpha) \leq \mathbb{E}[T]$

T je obecně nějaká náhodná veličina, tedy funkce, jejíž hodnotu

Postupně fixujeme části α tak, aby $\mathbb{E}[T | \text{fixovaná část}]$ *nerostla*

Využíváme toho, že $\mathbb{E}[T] = P(\alpha) \cdot \mathbb{E}[T | \alpha] + (1 - P(\alpha)) \cdot \mathbb{E}[T | \neg \alpha]$.

V našem případě bude $P(\alpha) = \frac{1}{2}$, takže stačí vybrat *menší* $\mathbb{E}[T | \alpha]$ a $\mathbb{E}[T | \neg \alpha]$.

• Původní randomizovaný krok vypadal takto:

• máme tabulku všech m_x a množinu $X \subseteq R$ *ta hraje roli* $\{m_x g(x) | x \in S_i\}$

• hledáme $a \in R$ t.č. $\sum_{x \in X} m_x \oplus a \leq \lfloor 2^{r-1} \cdot |X| \cdot \sum_{x \in X} m_x \rfloor$ *to je |S_i|*

randomizované jsme to uměli v $O(|X|)$ průměrně, ukážeme, jak deterministicky v $O(|X| \cdot r) = O(|X| \cdot \log n) \Rightarrow$ *sečte se na $O(n \log n)$*

• Postupně fixujeme bity čísla a od nejvyššího a přepočítáváme $\mathbb{E}[\# \text{kolizí} | \pi_k(a) = A]$

... stále tuto $\mathbb{E}[\dots]$ udržujeme *pod* původní $\mathbb{E}[\dots]$, stačí ji umět rychle spočítat. *prefix délky k*

$$\mathbb{E}\left[\sum_{x \in X} m_x \oplus a \mid \pi_k(a) = A\right] = \sum_{x \in X} \mathbb{E}\left[m_x \oplus a \mid \pi_k(a) = A\right] = \sum_{x \in X} \mathbb{E}\left[m_x \mid \pi_k(x) = \pi_k(x) \oplus A\right]$$

průměr všech m_x pro daný prefix $\pi_k(x)$

Budeme udržovat intervalový strom nad všemi m_x :



vrchol si pamatuje součet listů v podstromu

strom uložíme do pole jako haldy

- přepočítání m_i v $O(r)$
- dotaz na \sum listů pro daný prefix v $O(1)$

\Rightarrow • 1 krok derandomizace zvládneme v $O(|X|)$... *prefixy počítáme v $O(1)$ bitovými operacemi*

- celou volbu a zvládneme v $O(|X| \cdot r)$
- pak v $O(|X| \cdot r)$ aktualizujeme strom

} celý algoritmus *běží v $O(nr) = O(n \log n)$.*