

Datové struktury 2 - LS 2015/2016 - Martin Mareš

1

- Poprvé podle nové akreditace - dost jiná přednáška než dříve (snad je to zrušena k logickým?)
 - Kontakty: mj@ucw.cz, <http://mj.ucw.cz/ujuka/ds2/>
 - Cvičení nejsou, ale můžete si domluvit konzultaci.
 - Přibližný plán:
 - statické slovníky
 - celočíslné DS
 - dolní odhady
 - cache-oblivious DS
 - DS pro strany a obecné grafy
 - geometrické DS
 - listové DS
 - streaming algoritmy
- } literatura dosti kusá,
budeme přidávat odkazy
na web
- 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: &, |, ^, <<, >>, ~
 - 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 použité paměťové buňky
- BůHO dokážeme počítat i s $O(w)$ -bit. slovy
- všechny logaritmy budou nadále implicitně dvojkové

STATICKÉ MAJAZKY

- Chceme pro n -prvkovou $S \subseteq U$ vybudovat DS, která bude umět rychle odpovídat na dotazy "x ∈ S?"
- ↙ universum (treba slova RAMu)

	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 (stačí 2-nezávislost)
	$O(n)$ průměrně	$O(1)$ w.c.	... jiný přístup...
	$O(n \log n)$ w.c.	$O(1)$ w.c.	derandomizace

PERFEKTNÍ HEŠOVÁNÍ

FKS = Fredman, Komlos, Szemerédi 1984

(2)

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

Navíc obvykle chceme "heškovou 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 $h \in U \rightarrow [m]$ počítáme kolize: $\{x, y\} \in \binom{S}{2}$ t.j. $h(x) = h(y)$.

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

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

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

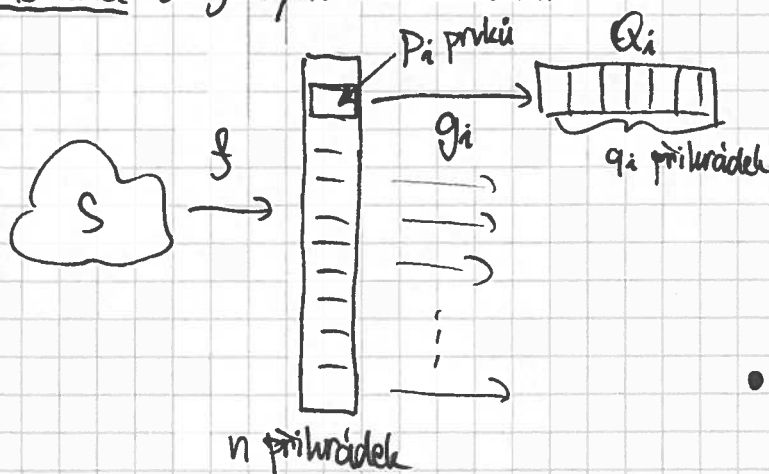
⇒ proto budeme-li volit h náhodně z \mathcal{H} , tak dlouho, než objevíme nějakou perfektní, bude to trvat průměrně ≤ 2 pokusy:

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

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

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

Konstrukce dvojitupňového hešování's



• $f \in \mathcal{H}$ hešuje do n příhradek

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

... umím efektivně najít f , pro něj $\# \text{kolizí} < cn$.

• g_i hešuje do $[q_i]$ pro $q_i = \lceil \frac{p_i^2 \cdot c}{2} \rceil$, takže umíme najít perfektní f_i

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

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

\uparrow $p_i \cdot 2 \cdot \binom{p_i}{2}$ # kolizí v i -té příhradce

$\#$ = # 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 hlavní tabulky pro pointer a param. g_i
- výpočet g_i
- vhlédnutí do Q_i

} $O(1)$ w.c.

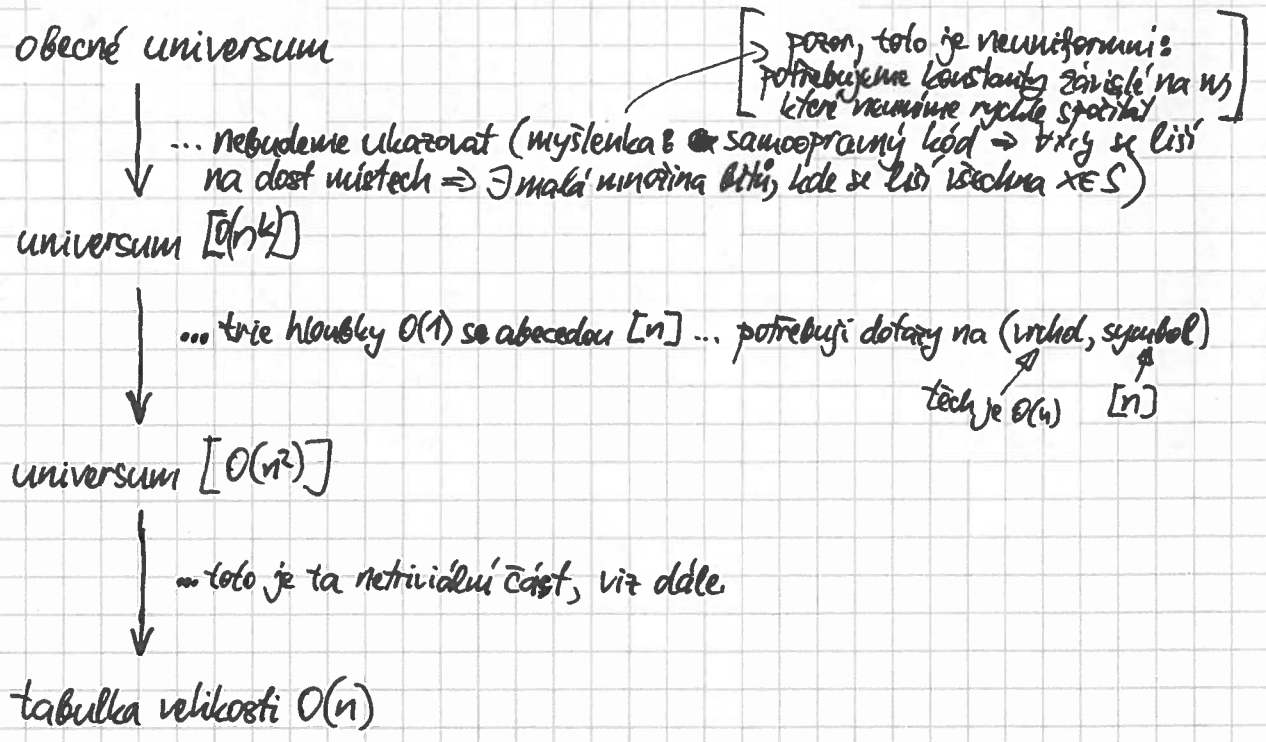
Poznámka: \exists dynamizace s časem $O(1)$ průměrně amortizovaně na Ins/Del a $O(1)$ w.c. na dotaz.

Odkázka: Třídění reálných čísel vybraných rovnoměrně náhodně z $[0,1]$.

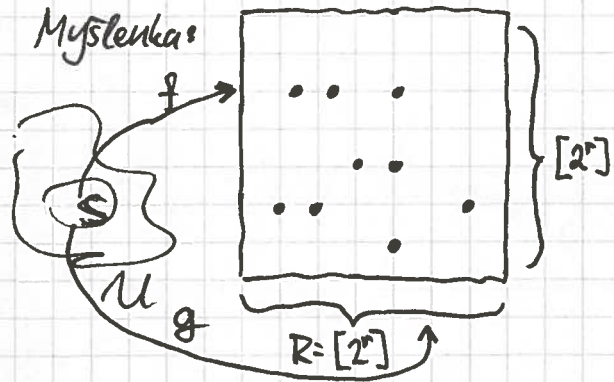
- Rozdělíme $[0,1]$ na n příhradek
- v $O(n)$ rozmístíme čísla do příhradek $\dots \in [\#kolizi] < \frac{n^2}{2}$
- v každé příhradce dotřídíme bublinkově \dots trvá to $O(\text{velikost příhrádky}^2)$, což se seče na $O(n)$

DETERMINISTICKÉ SLOVNÍKY Hagerup, Miltersen, Pagh 2000

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



• Notace: nerozlišujeme mezi čísly z $[O(n^2)]$ a bit. řetězi 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 $(i, j) \rightarrow (i, j \oplus a_i)$
 \dots a pak totéž znovu ve sloupcích

*položení řádek n-krát, sníží #kolizi
 #kolizi v sloupcích funkcí
 #kolizi v řádkách*

Df: Dvojice funkcí (f, g) z U do R je q-dobrá (pro $q \geq 0$),
 pokud f má na S nejvýše q kolizi 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.ž. $(x \mapsto g(x) \oplus a_{f(x)}, f(x))$ je q' -dobrá pro $q' = \lfloor \frac{n}{2^{2^r}} \cdot q \rfloor$
 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 $O(n)$.

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 $\binom{w}{r} < n^2$ kolizi.

} pár (f, g) je n^2 -dobrý

Lemma (zde je potřeba omezení $q' \leq n$ z ~~trivial~~ lemmatu)

1. krok: (f', g') ... jelikož $2^{3r} < \frac{1}{n}$, musí tento pár být $< n$ -dobrý

Lemma (... a zde naopak druhá část ...)

2. krok: (f'', g'') ... < 1 -dobrý, takže 0 -dobrý $\Rightarrow f''$ je prostá na S .

Výpočet hes. funkce

1. Rozdělíme x na ~~dvě~~ 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

Dk lemmatu: Nejprve očíslovujeme řádky od nejpnejšího:

$S_i = \{x \in S \mid f(x) = v_i\}$,
 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 ~~tabulky~~ a_{v_i} .

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

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

$$E[\#NK] = |S_{i-1}| \cdot |S_i| \cdot 2^{-r}$$

↓

za $O(1)$ pokusů najdu a_{vi} ,
pro které $\#NK \leq 2 \cdot E[\#NK]$

$$\|S_{i-1}| \cdot |S_i| \cdot 2^{1-r}$$

dvojice $x \in S_{i-1}, y \in S_i$
tož. $g(x) \oplus a_{g(x)} = g(y) \oplus a_{g(y)}$
 $a_{g(y)}$ pro nějaké $j \in S_{i-1}$ $= a_{vi}$

↕

$$a_{vi} = g(x) \oplus g(y) \oplus a_{g(x)}$$

... to nastane s pětí $\frac{1}{2^r}$

• Jak to udělat efektivně? • Udržujeme $M_w :=$ kolik bodů jsme už umístili do w -tého sloupce ... tedy $|\{x \in S_{i-1} \mid g(x) \oplus a_{g(x)} = w\}|$.

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

• Pro každou S_i zvolíme a_{vi} , pak pro $x \in S_i$ spočítáme pomocí M_x kolize

$\left. \begin{matrix} O(|S_i|) \\ \times O(1) \text{ příměrně} \\ \text{pokusy} \end{matrix} \right\} O(|S_i|)$

... až najdeme správné a_{vi} , přepočítáme M_x

↳ v součtu přes všechna i $O(n)$ průměrně

• Jak dobrý 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 $x \mapsto g(x) \oplus a_{g(x)}$:

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

shora omezeno touto

$|S_i| \cdot \sum_{j \in S_{i-1}} |S_j| = \sum_{j \in S_{i-1}} |S_i| \cdot |S_j| \leq \sum_{j \in S_{i-1}} |S_j|^2 \leq 4 \cdot \binom{|S_i|}{2}$
 pokud $|S_j| \geq 2$

$2^{1-r} \leq \frac{1}{n}$, takže $\dots \leq |S_i|$
 $\rightarrow \sum_i \dots \leq n$

$t_i = \max \{i \mid |S_i| > 1\}$ (to je druhá část min (...) z tvrzení lemmatu)

vynecháme členy, pro které $|S_i| = 1$
 ... tehdy $|S_{i-1}| \leq 2^{r-1}$ díky volbě n , takže $\lfloor \dots \rfloor = 0$

Derandomizace

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

Postupně fixujeme části X tak, aby $\mathbb{E}[T | \text{fixovaná část}]$ neklesala.

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 větší z $\mathbb{E}[T | \alpha]$ a $\mathbb{E}[T | \neg \alpha]$.

• Původní randomizovaný krok vypadal takto:

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

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

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

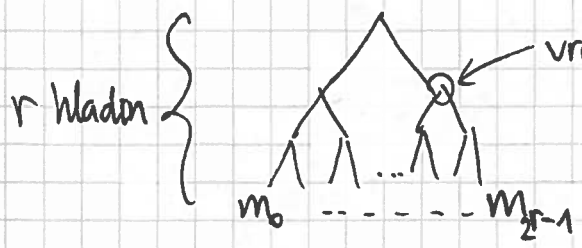
• Postupně fixujeme bity čísla a od nejvyššího a přepočítáváme $\mathbb{E}[\# \text{listů} | \pi_k(a) = A]$

... stále tuto $\mathbb{E}[\dots]$ udržujeme nad 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 Σ 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 $O(nr) = O(n \log n)$.