

# Teorie informace a kódování

## L3: dokonalé hašování

Jan Konecny

4. listopadu 2021

# First fit decreasing



A.C. Yao, R.E. Tarjan

Storing a sparse table.

Communications of the ACM. 1979 Nov 1; 22(11):606-11.

- Ukládání statických tabulek o  $n$  položkách, každá z nich je číslo mezi 0 a  $N - 1$ .

Pro jednoduchost budeme uvažovat

- číslo  $N$  je druhá mocnina nějakého celého čísla, t.j.  $N = m^2$  pro nějaké  $m \in \mathbb{Z}$ .
- že  $n \geq \max\{2, m\}$ .

Algoritmus *first-fit-decreasing*

- seřadíme řádky dle počtů nenulových prvků v klesajícím pořadí
- hladovým (greedy) způsobem přiřazujeme posuny.

## Příklad 1

$$A = \begin{pmatrix} 0 & * & 0 & 0 & * \\ * & * & 0 & * & 0 \\ * & 0 & * & 0 & 0 \\ 0 & 0 & 0 & * & 0 \\ * & 0 & 0 & 0 & * \end{pmatrix}$$

	0	*	0	0	*								$r(1) = 1$
*	*	0	*	0									$r(2) = 0$
				*	0	*	0	0					$r(3) = 4$
					0	0	0	*	0				$r(4) = 5$
							*	0	0	0	*		$r(5) = 7$
5	6	1	8	10	4	12	20	18	0	0	24		

Krok 1:

**for**  $i \leftarrow 1$  to  $m$  **do**

    | count( $i$ )  $\leftarrow 0$ ;  
    | list( $i$ )  $\leftarrow \emptyset$

**end**

**for** each nonzero position  $(i, j)$  in  $A$  **do**

    | count( $i$ )  $\leftarrow$  count( $i$ )+1;  
    | list( $i$ )  $\leftarrow$  list( $i$ )  $\cup \{i\}$

**end**

Po kroku 1:

- list( $i$ ) je seznam nenulových sloupců v řádce  $i$
- count( $i$ ) je jejich počet.

Krok 2:

**for**  $l \leftarrow 0$  to  $n$  **do**

| bucket( $l$ )  $\leftarrow \emptyset$

**end**

**for**  $i \leftarrow 0$  to  $m$  **do**

| bucket(count( $i$ ))  $\leftarrow$  bucket(count( $i$ ))  $\cup \{i\}$

**end**

Krok 2 je radix sort řádků podle jejich počtu nenulových prvků.

Krok 3:

```
for  $k \leftarrow 0$  to  $n + m$  do
```

```
  | entry( $k$ )  $\leftarrow$  false
```

```
end
```

```
for  $l \leftarrow n$  downto 0 do
```

```
  | for each  $i$  in bucket( $l$ ) do
```

```
    |  $r(i) \leftarrow 0$  ;
```

```
    | check for overlap: for each  $j$  in list( $i$ ) do
```

```
      | if entry( $r(i) + j$ ) then
```

```
        |  $r(i) \leftarrow r(i) + 1$ ;
```

```
        | goto check overlap
```

```
      | end
```

```
    | end
```

```
    | for each  $j$  in list( $i$ ) do
```

```
      | entry( $r(i) + j$ )  $\leftarrow$  true
```

```
    | end
```

```
  | end
```

```
end
```

Pokud je rozdělení nenulových prvků v řádcích  $A$  „přiměřeně rovnoměrné“, pak *first-fit-decreasing* účinně komprimuje  $A$ .

Toto empirické pozorování můžeme ověřit analýzou metody.

Intuitivně:

- řádky s pouze několika nenulovými prvky neblokují příliš mnoho možných hodnot posunutí pro jiné řádky,
- problémy způsobují pouze řádky s mnoha nenulovými prvky.



Abychom tento jev kvantifikovali, definujeme

- $n(l)$ , pro  $l \geq 0$  jako celkový počet nenulových prvků v řádcích s více než  $l$  nenulovými prvky.

Následující věta ukazuje, že pokud  $\frac{n(l)}{n}$  klesá dostatečně rychle s rostoucím  $l$ , *first-fit-decreasing* funguje dobře.

## Věta 2

Uvažujme že pole  $A$  splňuje *podmínku harmonického rozpadu*:

$$\text{pro libovolné } l, n(l) \leq n/(l + 1).$$

Pak každý posun  $r(i)$  vypočítaný pro  $A$  metodou *first fit* splňuje  $0 \leq r(i) \leq n$ .

## Věta 2

Uvažujme že pole  $A$  splňuje **podmínku harmonického rozpadu**:

$$\text{pro libovolné } l, n(l) \leq n/(l + 1).$$

Pak každý posun  $r(i)$  vypočítaný pro  $A$  metodou first fit splňuje  
 $0 \leq r(i) \leq n$ .

## Důkaz.

- Pro libovolný řádek  $i$ , uvažujme výběr posunu  $r(i)$ .
- Uvažujme, že  $r(i)$  obsahuje  $l \geq 1$  nenulových prvků.
- Dle harmonického rozkladu je počet nenulových prvků nejvýše  $\frac{n}{l}$ .
- Každý takový nenulový prvek může blokovat nejvýše  $l$  výběrů pro  $r(i)$ .
- Dohromady může být blokováno nejvýše  $n$  výběrů, a tedy  
 $0 \leq r(i) \leq n$ .



## Význam harmonického rozpadu:

Pokud má  $A$  harmonický rozpad,

- musí být alespoň polovina nenulových prvků v  $A$  v řádcích pouze s jedním nenulovým prvkem.
- žádný řádek nesmí mít více než  $\sqrt{n}$  nenulových prvků.

Pokud  $A$  nemá harmonický rozklad, pak *first-fit-decreasing* neudělá dobrou kompresi – aspoň ne v nejhorším případě.

A protože nás zajímá i nejhorší případ, měli bychom se podívat, jak se s ním srovnat.

Potřebujeme najít způsob, jak „uhladit“ rozložení nenulových prvků v řádcích  $A$ , než začneme počítat posuny řádků.

A to tak, že v  $A$  posuneme sloupce:

- každému sloupci přiřadíme posun  $c(j)$  zobrazující každou pozici  $(i, j)$  do nové pozice  $(i + c(j), j)$ .
- to transformuje  $A$  do nového pole  $B$  s více řádky – konkrétně  $\max_j(c_j) + m$  – a se stejným počtem sloupců.

### Příklad 3

$$A = \begin{pmatrix} 0 & * & 0 & 0 & * \\ * & * & 0 & * & 0 \\ * & 0 & * & 0 & 0 \\ 0 & * & 0 & * & 0 \\ * & 0 & 0 & 0 & * \end{pmatrix} \Rightarrow B = \begin{pmatrix} 0 & 0 & 0 & 0 & * \\ * & 0 & 0 & 0 & 0 \\ * & 0 & * & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \\ * & 0 & 0 & 0 & * \\ 0 & * & 0 & 0 & 0 \\ 0 & * & 0 & 0 & 0 \\ 0 & 0 & 0 & * & 0 \\ 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & * & 0 \\ 0 & 0 & 0 & 0 & 0 \end{pmatrix}$$

$$c(1) = 0, c(2) = 5, c(3) = 0, c(4) = 6, c(5) = 0$$

Potřebujeme kritérium pro výběr posunů sloupců.

Tím kritériem bude podmínka **exponenciálního rozpadu** definovaná následovně:

Nechť

- $B_j$  je pole sestávající z prvních  $j$  posunutých sloupců z  $A$ .
- $n_j$  je celkový počet nenulových prvků v  $B_j$ .
- $n_j(i)$  pro  $i \geq 0$  je celkový počet nenulových prvků v  $B_j$ , které se objevují v řádcích  $B_j$  obsahující více než  $i$  nenulových prvků.

$$E_j: \text{Pro libovolné } i > 0, n_j(i) \leq n_j/2^{i(2-n_j/n)}$$

Vybíráme sloupcově posuny tak, aby splňovaly  $E_j$  pro všechna  $j$  s tím, že používáme *first-fit* následovně:

- Vypočítáme postupně – od 1 do  $m$ .
- Vybereme jako posun  $c(j)$  pro sloupec  $j$  nejmenší hodnotu, t.ž.  $B_j$  splňuje  $E_j$ .

# „Univerzální“ hašování

Problém:

- mod je docela pomalá operace



Dietzfelbinger M, Hagerup T, Katajainen J, Penttonen M.  
A reliable randomized algorithm for the closest-pair problem.  
Journal of Algorithms. 1997 Oct 1;25(1):19-51.

- Navržena technika [multiply-shift](#)

# Multiply-shift

Předpokládáme

- $m = 2^M$ .
- $w$  je počet bitů ve slově.
- Hašovací funkce jsou parametrizovány lichými kladnými čísly  $a < 2^w$  (které se vlezou do slova o  $w$  bitech)

$$h_a(x) = (a \cdot x \bmod 2^w) \div 2^{w-M}$$

Neboli, zapsáno v jazyku C:

$$h_a(x) = (\text{size\_t}) (a*x) \gg (w-M)$$



## Příklad 4

$$h_{111}(80) = \dots$$

$a \cdot x$  :

$$111 \cdot 80 = 8880 = 10001010110000_2$$

$(a \cdot x) \bmod 2^w$  :

$$10001010110000_2 = 176$$

pro  $M = 4$ :

$((a \cdot x) \bmod 2^8) \div 2^{8-4}$  :

$$10110000_2 = 11$$

- Až  $4\times$  rychlejší výpočet.
- Toto schema **ne**-splňuje univerzálnost.
- Splňuje volnější podmínku

$$\Pr\{h_a(x) = h_a(y)\} \leq \frac{2}{m}$$

- To je tzv. **2-univerzálnost**.

## Univerzální hashování vektorů

Hashujeme  $\bar{x} = (x_0, \dots, x_k)$ ,  $x_i$  jsou celá čísla každé má  $w$  bitů.

- Pokud  $\mathcal{H}$  je univerzální, pak také třída hashovacích funkcí

$$h(\bar{x}) = \left( \sum_{i=0}^{k-1} h_i(x_i) \right) \bmod m,$$

kde  $h_i \in \mathcal{H}$  jsou vybrány náhodně a nezávisle, je univerzální.

- Též možno použít multiply-shift:

$$h_{\bar{a}}(\bar{x}) = \left( \left( \sum_{i=0}^{k-1} x_i \cdot a_i \right) \bmod 2^{2w} \right) \div 2^{2w-M},$$

kde  $\bar{a} = (a_0, \dots, a_{k-1})$  jsou náhodná lichá celá čísla na  $2w$  bitech.

## Univerzální hashování řetězců

Předpokládejme, že  $\bar{x} = (x_0, \dots, x_l)$ , kde neznáme horní hranici pro  $l$ .  
Znaky řetězce  $\bar{x}$  budeme uvažovat jako koeficienty polynomu modulo velké  
prvočíslo  $p$ .

$$h_a(\bar{x}) = h_{\text{int}} \left( \left( \sum_{i=0}^l x_i \cdot a^{l-i} \right) \bmod p \right)$$

kde  $a \in \{1, \dots, p-1\}$  je náhodné číslo a  $h_{\text{int}}$  je vybrána z univ. třídy  
hašovacích funkcí na hašování celých čísel  $\{0, \dots, p-1\} \rightarrow \{0, \dots, m-1\}$ .

```
 $h \leftarrow \text{INIT\_VAL};$   
for  $i \leftarrow 0$  to  $x.\text{length}$  do  
  |  $h \leftarrow ((h \cdot a) + x[i]) \bmod p$   
end  
return  $h$ 
```

V praxi můžeme  $\text{mod } a$  vypustit, a povolit přetečení integeru, protože to je v mnoha programovacích jazycích ekvivalentní:  
 $\text{mod } (\text{Max-Int-Value} + 1)$

Hodnoty pro inicializaci  $h$  v některých implementacích:

- djb2:  $\text{INIT\_VAL} = 5381, a=33$
- STLPort 4.6.2:  $\text{INIT\_VAL} = 0, a=5$
- `java.lang.String.hashCode()`:  $\text{INIT\_VAL} = 0, a=31$

## Hash and displace

```
def Construct( $B, y$ ):  
  input :  $S$  – set of keys;  $b$  – number of bins;  $\epsilon > 0$   
  
   $n \leftarrow |S| \cdot (1 + \epsilon)$   
  if  $b < (2 + \epsilon) \cdot n$  then  
    | return -1  
  end  
  repeat  
    | Pick  $f$  at random in UTHF with range  $\{0, \dots, n - 1\}$   
    | Pick  $g$  at random in UTHF with range  $\{0, \dots, b - 1\}$   
    | for  $i \leftarrow 0$  to  $b - 1$  do  
    |   |  $B[i] \leftarrow \{s \in S \mid g(s) = i\}$   
    |   | Find permutation  $P$  s.t.  $|B[P[i]]| \geq |B[P[i + 1]]|$  for  
    |   |    $i = 0, \dots, b - 2$   
    |   |    $m = \max\{i \mid |B[P[i]]| > 1\}$   
    | end  
  until  $|B[P[0]]|^2 + \dots + |B[P[m]]|^2 \leq \frac{n}{1 + \epsilon}$  and  $f$  is 1-1 on each  $B[i]$ 
```

```

for  $i \leftarrow 0$  to  $n - 1$  do
  |  $T[i] \leftarrow \text{false}$  (marks unoccupied slots of the table);
end
for  $i \leftarrow 0$  to  $m$  do
  | repeat
  | | Pick  $d \in 0, \dots, n - 1$  at random
  | until  $T[f(s) + d) \bmod n]$  is false for each  $s \in B[P[i]]$ ;
  |  $D[P[i]] \leftarrow d$ ;
  | for each  $s \in B[P[i]]$  do
  | |  $T[f(s) + d) \bmod n] \leftarrow \text{true}$ 
  | end
end
Set  $E[1], \dots, E[k]$  to distinct values where  $T[E[i]] \leftarrow \text{false}$ ;
for  $i \leftarrow m + 1$  to  $m + k$  do
  | Take (the unique)  $s \in B[P[i]]$ ;
  |  $D[P[i]] \leftarrow (E[i - m] - f(s)) \bmod n$ 
end
return( $f, g, D, n$ )

```