

Poradová štatistika

- i - tá poradová štatistika n prvkov (napr. kľúčov tabuľky) je i - tý najmenší prvok v tabuľke.
- i - tá poradová štatistika je prvok, ktorý je väčší ako práve $i - 1$ iných prvkov tabuľky

Predpokladajme jedinečnosť hodnôt kľúčov.

Špeciálne poradové štatistiky:

Minimum: $i = 1$

Maximum: $i = n$

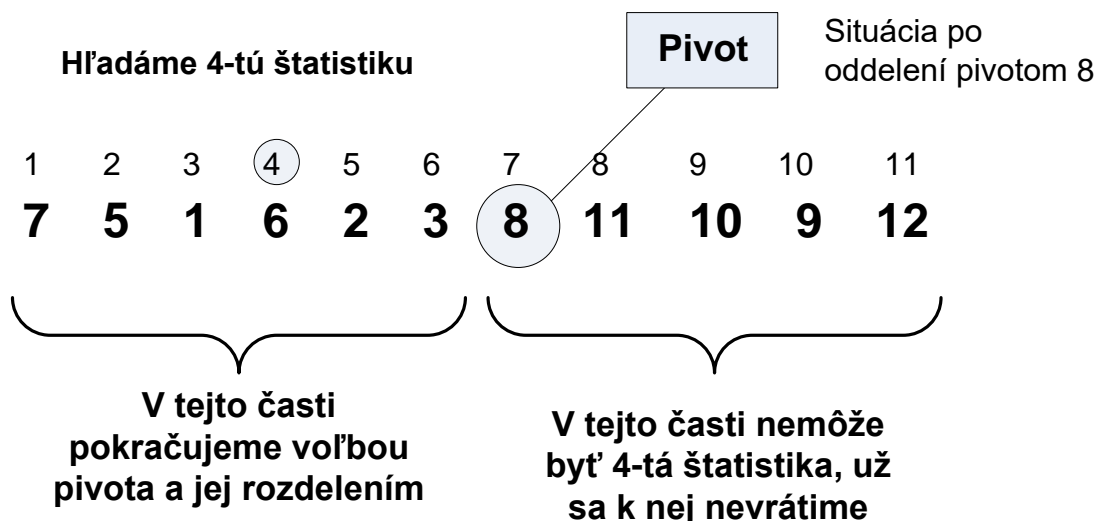
Medián: $i = (n + 1) / 2$ ak n je nepárne,
 $i = n / 2$ a $i = n / 2 + 1$ (dva mediány) ak n je párne.

Ukážeme si štyri algoritmy pre nájdenie i - tej poradovej štatistiky z n prvkov (prvé tri algoritmy predpokladajú uloženie prvkov v poli, štvrtý predpokladá uloženie prvkov vo vyváženom binárnom vyhľadávacom strome):

- I. Poradová štatistika v očakávanom (priemerná zložitosť) aj najhoršom čase (najhoršia zložitosť) $O(n \cdot \log_2 n)$
 1. utried' prvky algoritmom s najhoršou zložitosťou $O(n \cdot \log_2 n)$ (napr. Heapsort)
 2. v i - tom prvku je požadovaná štatistika
- II. Poradová štatistika v očakávanom čase $O(n)$, najhoršom $O(n^2)$

Algoritmus je založený na tej istej myšlienke ako QuickSort (partitioning – oddeľovanie). Po oddelení sa však nevenuje obom častiam (ako QuickSort) ale iba tej, v ktorej sa nachádza index i .

Príklad:



Nepriaznivá najhoršia zložitosť je spôsobená náhodným výberom pivota (prípady, že za pivota sa vyberie vždy prvok blízky maximu alebo minimu v spracováwanej časti), tak ako je to aj v QuickSorte. Rovnako ako pri QuickSorte existuje viacero možných stratégií ako vybrať pivota. Jednou z nich výber mediánu z prvého, posledného a prostredného prvku. Inou možnosťou je náhodný výber pivota.

III. Poradová štatistika v očakávanom aj najhoršom čase $O(n)$: Princíp rozdeľovania zostáva podobný, ako v algoritme II ale deliaci prvok (pivot) sa vyberá deterministicky. Pivot sa vyberá ako medián z mediánov menších častí (o dĺžke 5 alebo iné nepárne číslo) a využíva sa skutočnosť, že ak tento medián mediánov častí postupnosti použijeme ako pivot pre rozdelenie zaručí nám to najhoršiu zložitosť $O(n)$.

1. Rozdeľ n prvkov uložených v poli na $\lfloor n/5 \rfloor$ skupín po 5 prvkov a najviac jednu obsahujúcu zvyšných $n \bmod 5$ prvkov.
2. Nájdi medián každej z $\lfloor n/5 \rfloor$ častí napr. tak, že sa utriedi InsertSortom a zoberie sa prostredný prvok (v prípade párneho počtu menší z nich).
3. Nasad' rekurzívne tento algoritmus na nájdenie mediánu x z $\lfloor n/5 \rfloor$ mediánov nájdených v kroku 2.
4. Teraz, po výstupe z rekurzie x síce nie je mediánom všetkých prvkov, ale je dokázané, že ak ho použijeme ako pivot, žiadna z častí (partition) nemôže obsahovať menej ako $(3n/10)-6$ prvkov (pre veľké n takmer 30%).

- ## Príklad:

70 52 01 46 92 13 88 11 10 09 12 54 21 22 71 45 66 32 19 44 72 06 59 35 42 97 15 16 58 62 49 40 18

01 46 52 70 92 09 10 11 13 88 12 21 22 54 71 19 32 44 45 66 06 35 42 59 72 15 16 58 62 97 18 40 49

52	11	22	44	42	58	40
-----------	-----------	-----------	-----------	-----------	-----------	-----------

Diagram illustrating a Huffman tree structure for the sequence 11 22 42 44 52 40 58. The root node is 110, which splits into 11 and 10. Node 10 splits into 1 and 0. Node 1 splits into 11 and 0. Node 0 splits into 1 and 0. Node 11 splits into 11 and 0. Node 0 splits into 1 and 0. Node 11 splits into 11 and 0. Node 0 splits into 1 and 0. The leaf nodes are 11, 22, 42, 44, 52, 40, and 58.

Diagram illustrating the selection of a pivot element for partitioning:

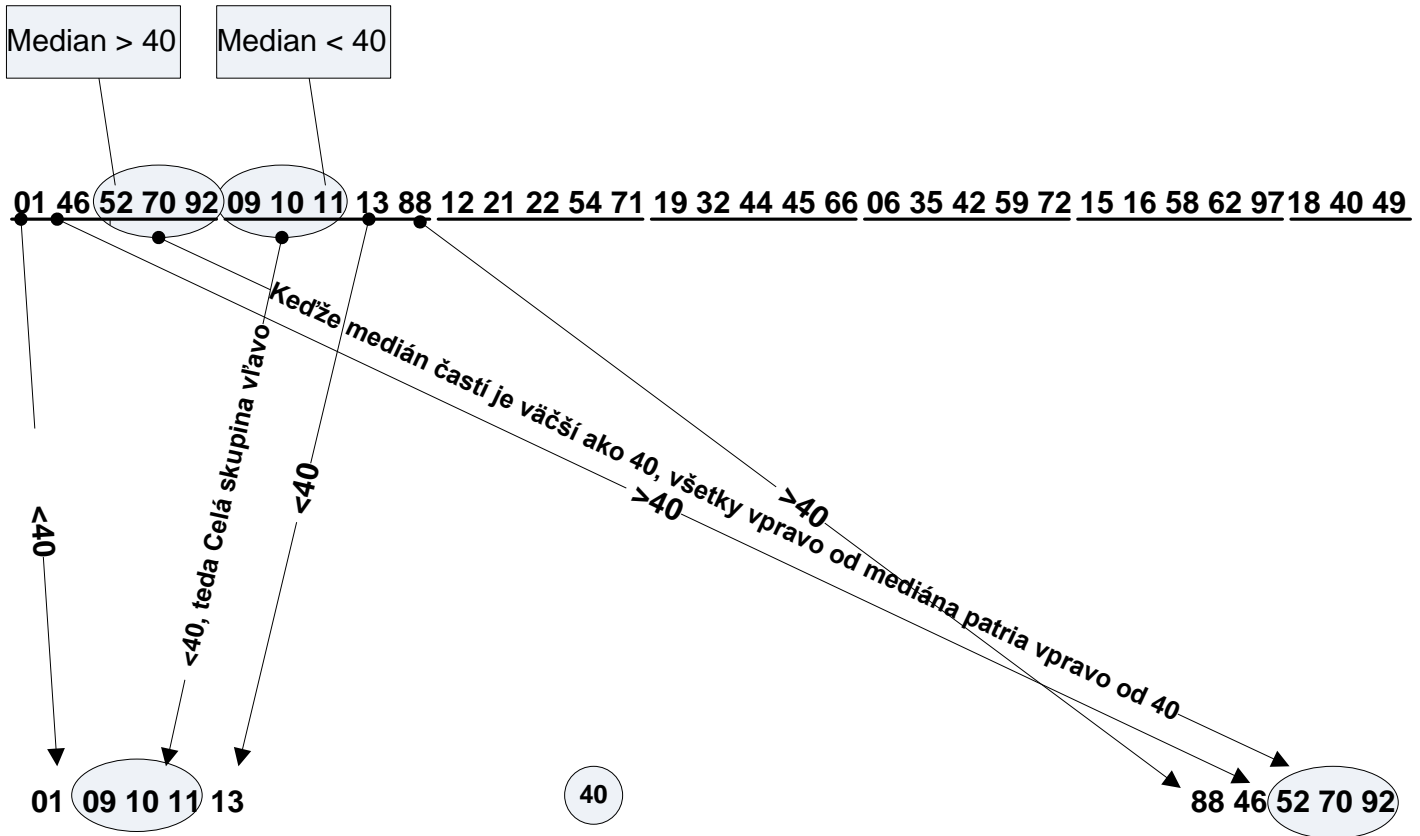
The array is: 01 13 11 10 09 12 21 22 32 19 06 35 15 16 23 18 40 70 52 46 92 88 54 71 45 66 44 72 59 42 97 58 62 49

The pivot element is 40, indicated by the label "pre rozdelenie" (before partitioning) and an arrow pointing to it.

**V tejto čast nemôže
byť 4-tá štatistika, už
sa k nej nevrátíme**

Zefektívnenie procesu rozdeľovania:

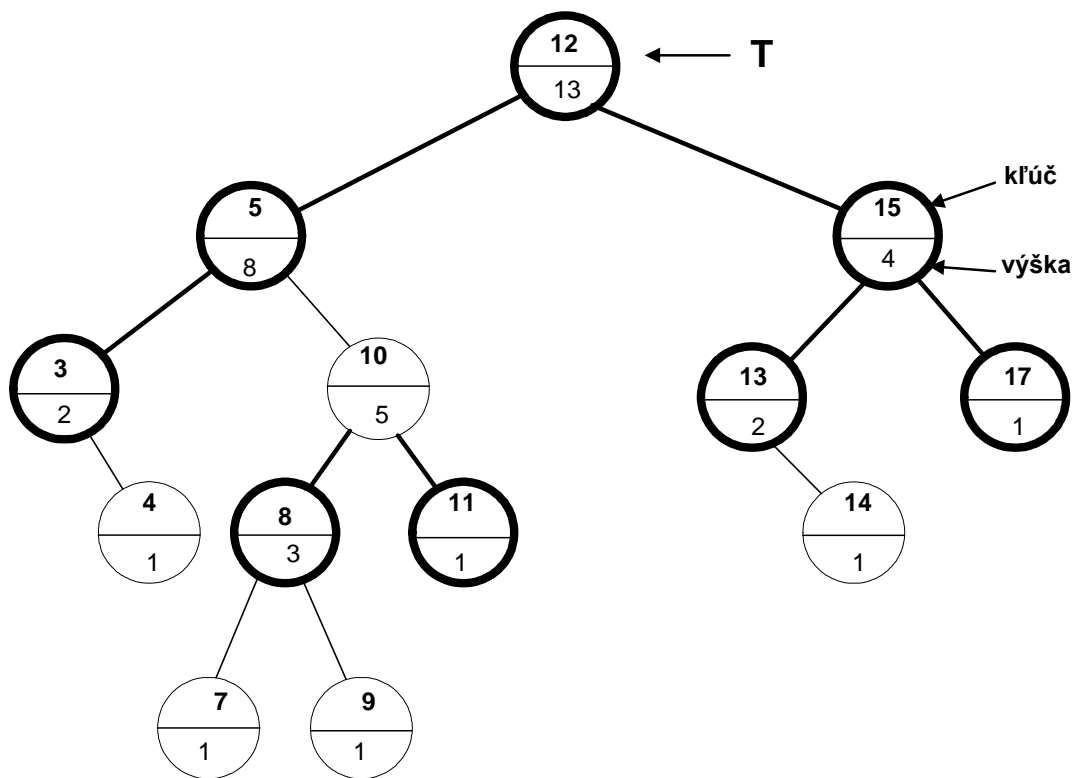
- ak je medián skupiny menší ako medián mediánov, sú menšie aj všetky vľavo od neho v skupine a teda sa môžu presunúť do ľavej časti bez jednotlivého porovnávania
- analogicky, keď je medián skupiny väčší



IV. Poradová štatistika v očakávanom aj najhoršom čase $O(\log_2 n)$

Predpokladá voľbu štruktúry dát vo forme ľubovoľného vyváženého binárneho vyhľadávacieho stromu (napr. RB strom).

Takýto strom doplníme tak, že každý vrchol x bude obsahovať informáciu $size(x)$ o počte vrcholov v podstrome, ktorého je koreňom vrátane koreňa samého.



Ak dohodneme $\text{size}(\text{nil}) = 0$, tak:
 $\text{size}(x) = \text{size}(\text{left}(x)) + \text{size}(\text{right}(x)) + 1$

Ukážeme si rekurzívnu formu algoritmus, ktorý nájde i - tú poradovú štatistiku v strome T :

funkcia $\text{Statistika}(\text{strom } T, \text{prvok } i)$

1. $r = \text{size}(\text{left}(x))$
2. ak $i = r + 1$, vráť x (i – tý prvok nájdený), inak krok 3
3. ak $i < r + 1$, vráť $\text{Statistika}(\text{left}(x), i)$,
 inak vráť $\text{Statistika}(\text{right}(x), i - r - 1)$

Pri reorganizáciách spôsobených operáciami *Vlož*, *Odober* je nutné aktualizovať hodnoty *size*.