

Selection

CS 4231, Fall 2020

Mihalis Yannakakis

Selection (Order Statistics)

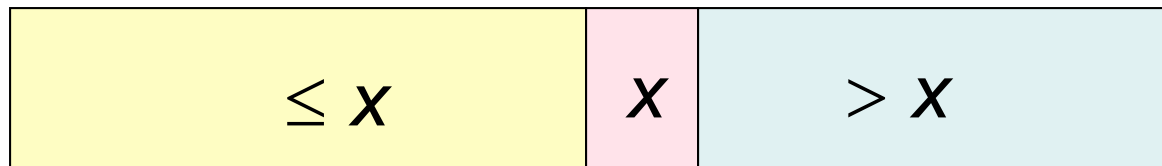
- **Input:** Set A of n numbers (or more generally, elements from an ordered domain), number i , $1 \leq i \leq n$
- **Output:** i -th smallest (rank i) element of A
- $i=1$: minimum
- $i=n$: maximum
- $i=(n+1)/2$: median

Straightforward solution

- Sort A
- Output i-th element
- $\Theta(n \log n)$
- Can we do better?
- Min, Max: $\Theta(n)$ easy
- General rank i ?

Divide and Conquer Selection of i-th smallest element

- **Divide:** Partition the input array A of elements with respect to a **pivot** element x into two parts:



- **Conquer & Combine:**
Check if $|\text{left part}| \geq i$, $= i-1$, or $< i-1$,
and accordingly recurse on left part,
return x , or recurse on right part

Example

Select 6th smallest element in array

3	8	4	2	9	7	5
---	---	---	---	---	---	---

Partition the array w.r.t pivot 5

3	4	2	5	9	7	8
---	---	---	---	---	---	---

Select recursively the 2nd smallest element in

right part:

9	7	8
---	---	---

Randomized Selection

- RANDOMIZED-SELECT(A, p, r, i)

Input: array A , indices $p \leq r$, number $i \leq r - p + 1$

if $p = r$ **then return** $A[p]$

$q = \text{RANDOMIZED-PARTITION}(A, p, r)$

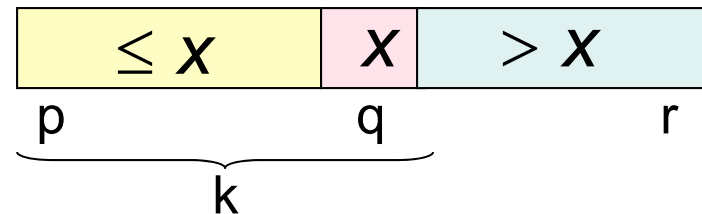
$k = q - p + 1$

if $i = k$ **then return** $A[q]$

else if $i < k$

then return $\text{RANDOMIZED-SELECT}(A, p, q - 1, i)$

else return $\text{RANDOMIZED-SELECT}(A, q + 1, r, i - k)$



Main Call: $\text{RANDOMIZED-SELECT}(A, 1, n, i)$

Analysis

- **Worst Case:** Unbalanced partition

$$T(n) = T(n-1) + \Theta(n) \Rightarrow T(n) = \Theta(n^2)$$

- **Lucky Case:** Balanced partition $\frac{1}{2} : \frac{1}{2}$

$$T(n) = T(n/2) + \Theta(n)$$

Master theorem: $a=1$, $b=2$, $f(n)=n$, $n^{\log_b a} = 1$

$$\Rightarrow T(n) = \Theta(n)$$

Expected Time Analysis

- Assume distinct elements
- Let $T(n)$ =expected running time
- Partition $(k, n-k-1)$ with probability $1/n$ for each $k=0, \dots, n-1$
- At worst we'll recurse on the larger part

$$\begin{aligned} T(n) &\leq \sum_{k=0}^{n-1} \frac{1}{n} T(\max(k, n-k-1)) + \Theta(n) \\ &= \frac{1}{n} \sum_{k=0}^{\lfloor n/2 \rfloor - 1} T(n-k-1) + \frac{1}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} T(k) + \Theta(n) \\ &\leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} T(k) + \Theta(n) \end{aligned}$$

Expected Time Analysis

$$T(n) \leq \frac{2}{n} \sum_{k=\lfloor n/2 \rfloor}^{n-1} T(k) + an$$

Solution: $T(n) \leq cn$

Proof:
(Induction step)

$$\begin{aligned} T(n) &\leq \frac{2c}{n} \cdot \sum_{k=\lfloor n/2 \rfloor}^{n-1} k + an \\ &= \frac{2c}{n} \cdot \frac{(\lfloor n/2 \rfloor + n - 1)(n - \lfloor n/2 \rfloor)}{2} + an \\ &\leq \frac{2c}{n} \cdot \frac{(3n/2)(n/2)}{2} + an \\ &\leq (3c/4 + a)n \\ &\leq cn \quad \text{provided } 3c/4 + a \leq c, \text{ i.e., } 4a \leq c \end{aligned}$$

$$T(n) = \Theta(n)$$

Randomized Selection algorithm

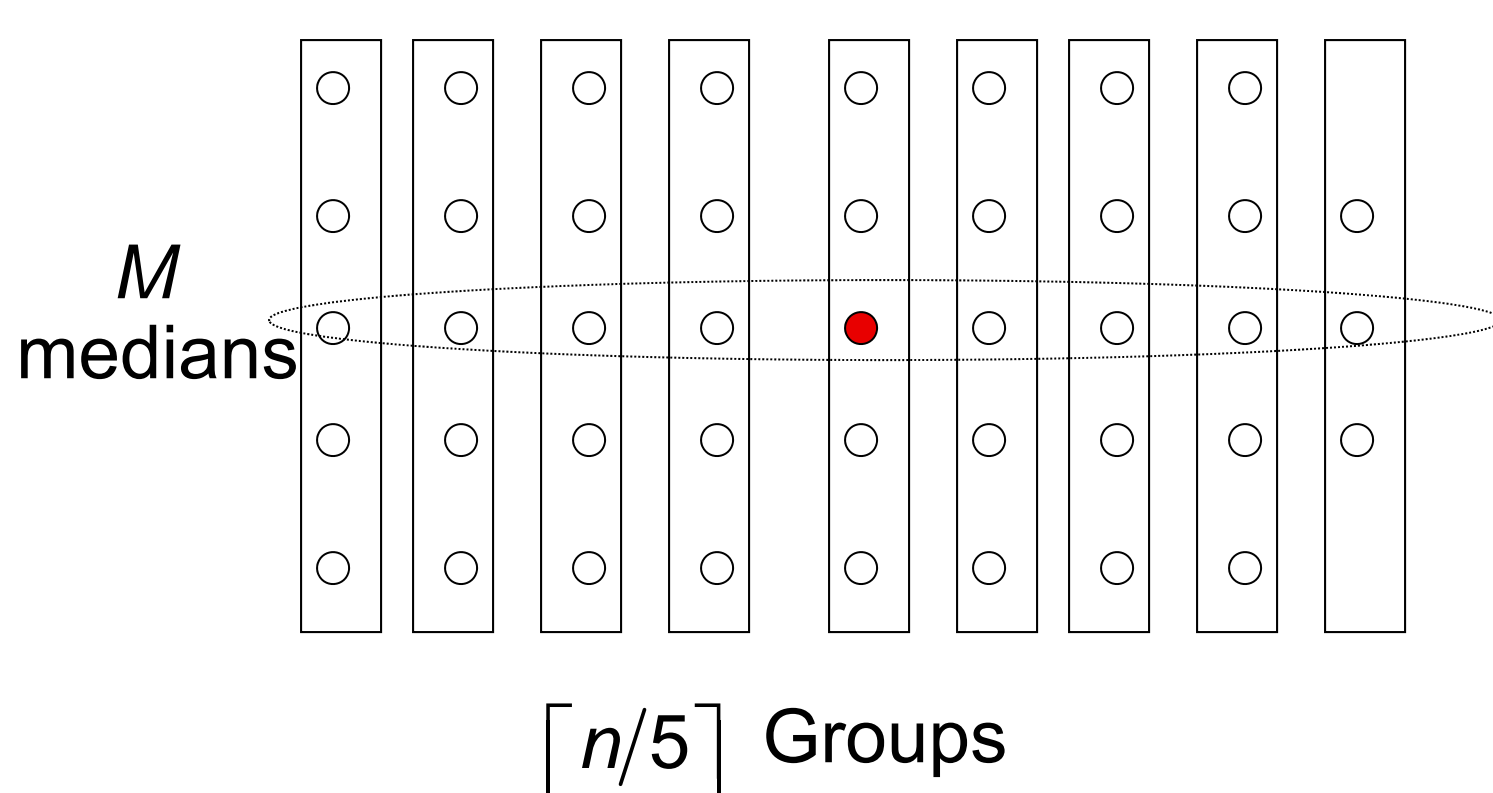
- Practical, fast, in-place
 - Linear expected time
(if equal elements use 3-way partition)
 - Quadratic worst case time
-
- Is there a deterministic algorithm that runs in worst case linear time?
 - Yes

Deterministic Selection Algorithm

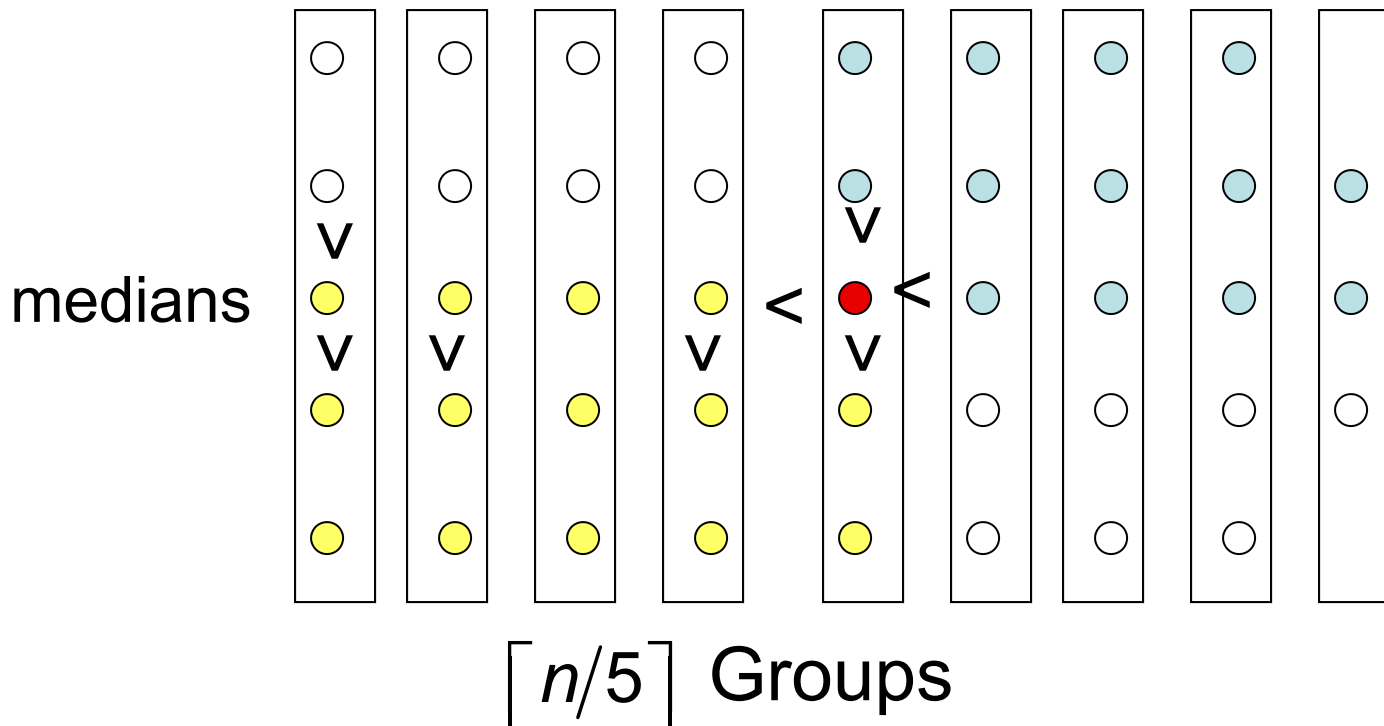
- Divide the given set A of elements into groups of 5
- Find the median in each group directly
→ subset M
- Recursively select the median x of M
- Partition A with respect to pivot x
- Recurse on the appropriate part

Key property: Partition wrt x not too unbalanced

Choice of pivot x



Choice of pivot x



At least $1/2$ groups have median $\leq x$
 from each such group (but one) 3 elements $\leq x$
 \Rightarrow at least $(3n/10)-2$ elements are $\leq x$

Similarly, at least $(3n/10)-2$ elements are $\geq x$

Analysis

$$T(n) \leq T\left(\left\lceil \frac{n}{5} \right\rceil\right) + T\left(\frac{7n}{10} + 2\right) + \Theta(n)$$

Solution: $T(n) = \Theta(n)$

Proof (by substitution): Guess $T(n) \leq cn - b$

$$T(n) \leq T(\lceil n/5 \rceil) + T(7n/10 + 2) + an$$

$$\leq c(n/5 + 1) + c(7n/10 + 2) + an - 2b$$

$$= n[c(1/5 + 7/10) + a] + 3c - 2b$$

$$\leq cn - b$$

provided $c(1/5 + 7/10) + a \leq c$, i.e., $10a \leq c$

and $3c \leq b$

Deterministic vs. Randomized Selection

- Deterministic runs always in $O(n)$ time
- But in practice, randomized algorithm is faster because hidden constant of deterministic algorithm in $O(n)$ is large