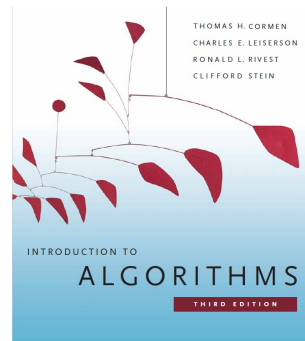




Order Statistics

Prepared by Shmuel Wimer
Courtesy of Prof. Dror Rawitz



Aug 2024

Algorithms and DS II: Order Statistics

1

1



Order Statistics and Selection Problem

The **i th order statistics** of n elements is the i th smallest element.

Given a set A of n distinct elements, the i th order statistics **selection problem** is to find $x \in A$ s.t. x is larger than exactly $i - 1$ elements.

Sorting A in $O(n \log n)$ time enables to solve in $O(n)$ time, but faster solution in $O(n)$ time (no sorting) is possible.

Minimum $i = 1$ and **maximum** $i = n$ problems complexity is $\Theta(n)$.

$\Theta(n)$ is possible for any i th order statistics.

Aug 2024

Algorithms and DS II: Order Statistics

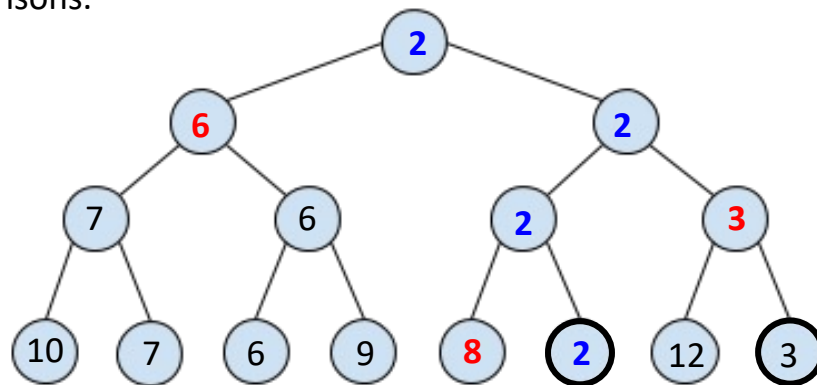
2

2



Finding the 2nd Smallest

Can be found by knockout tournament with $n + \lceil \log n \rceil - 2$ comparisons.



Aug 2024

Algorithms and DS II: Order Statistics

3

3



1. Divide input into $k = \lceil n/2 \rceil$ element pairs.
2. Get the set of k smaller of every pair.
3. If $k == 1$ stop, smallest found. Else set $n = k$ and go back to 1.

Knockout tournament has $n - 1$ comparisons.

Smallest wins after $\lceil \log n \rceil$ comparisons, implying binary tree.

At some comparison **the smallest must beat** the 2nd smallest.

Traversal of the $\lceil \log n \rceil - 1$ beats of smallest discover the 2nd.

Aug 2024

Algorithms and DS II: Order Statistics

4

4



Selection of k th Order Statistics in Expected $\Theta(n)$

A divide-and-conquer algorithm similar to Quicksort, except that it proceeds with only one of the partitions rather than both.

RANDOMIZED-SELECT seeks the i th order statistics of A .

It partitions the array recursively.

It chooses a pivot randomly with RANDOMIZED-PARTITION but proceed with only one side of the partition where i is surely located.

Aug 2024

Algorithms and DS II: Order Statistics

5

5



```

RANDOMIZED-SELECT( $A, p, r, i$ ) //  $i$ th order statistics
  if  $p == r$  return  $A[p]$  //  $A[p]$   $i$ th smallest element
  // partition  $A[p..r]$  into  $A[p..q-1]$  and  $A[q+1..r]$ 
  // around pivot  $A[q]$ 
   $q = \text{RANDOMIZED-PARTITION}(A, p, r)$ 
   $k = q - p + 1$  // number of elements in  $A[p..q]$ 
  if  $i == k$  return  $A[q]$  //  $A[q]$  is  $i$ th smallest element
  else if  $i < k$  //  $i$ th is in the lower part
    return RANDOMIZED-SELECT( $A, p, q - 1, i$ )
  else //  $i$ th is in the higher part
    return RANDOMIZED-SELECT( $A, q + 1, r, i - k$ )
  
```

Aug 2024

Algorithms and DS II: Order Statistics

6

6



Run time analysis

Worst-case run time is $\Theta(n^2)$. When it happens? (HW)

Run time $T(n)$ on $A[p..r]$ of n elements is random variable.

RANDOMIZED-PARTITION(A, p, r) likely returns any element as pivot $A[q]$. $\Rightarrow \Pr[|[p..q]| = k] = 1/n, 1 \leq k \leq n$.

Define random variable $X_k = \begin{cases} 1 & |[p..q]| = k \\ 0 & \text{otherwise} \end{cases} \Rightarrow E[X_k] = \frac{1}{n}$.

It is unknown a priori whether selection proceeds with $A[p..q-1]$, $A[q]$ or $A[q+1..r]$, so larger interval is safely assumed.

Aug 2024

Algorithms and DS II: Order Statistics

7

7



When $X_k = 1$, the implied subarrays have size $k-1$ and $n-k$. \Rightarrow

$$T[n] \leq \sum_{k=1}^n X_k \cdot (T(\max(k-1, n-k)) + O(n)).$$

Taking expectation and applying its linearity, there is

$$E[T[n]] \leq \sum_{k=1}^n E[X_k \cdot T(\max(k-1, n-k))] + O(n)$$

Independence: $\Pr[X_k = 0 / 1 | \max(k-1, n-k)] = \Pr[X_k = 0 / 1]$

$$= \sum_{k=1}^n E[X_k] \cdot E[T(\max(k-1, n-k))] + O(n)$$

$$= \frac{1}{n} \sum_{k=1}^n E[T(\max(k-1, n-k))] + O(n).$$

Aug 2024

Algorithms and DS II: Order Statistics

8

8



There is $\max(k - 1, n - k) = \begin{cases} k - 1 & \text{if } k > \lceil n/2 \rceil \\ n - k & \text{if } k \leq \lceil n/2 \rceil \end{cases}$.

In $\sum_{k=1}^n E[\cdot]$, for even n each term from $T(\lceil n/2 \rceil)$ to $T(n - 1)$ appears twice. Same for odd n plus $T(\lceil n/2 \rceil)$ which appears once. \Rightarrow

$$E[T[n]] \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} E[T(k)] + O(n).$$

Above recursion is solved by substitution, yielding $E[T(n)] = O(n)$.
(HW, see CLRS 9.2).

Aug 2024

Algorithms and DS II: Order Statistics

9

9



Selection of k th Order Statistics in $O(n)$ Worst Case

SELECT(A, k)

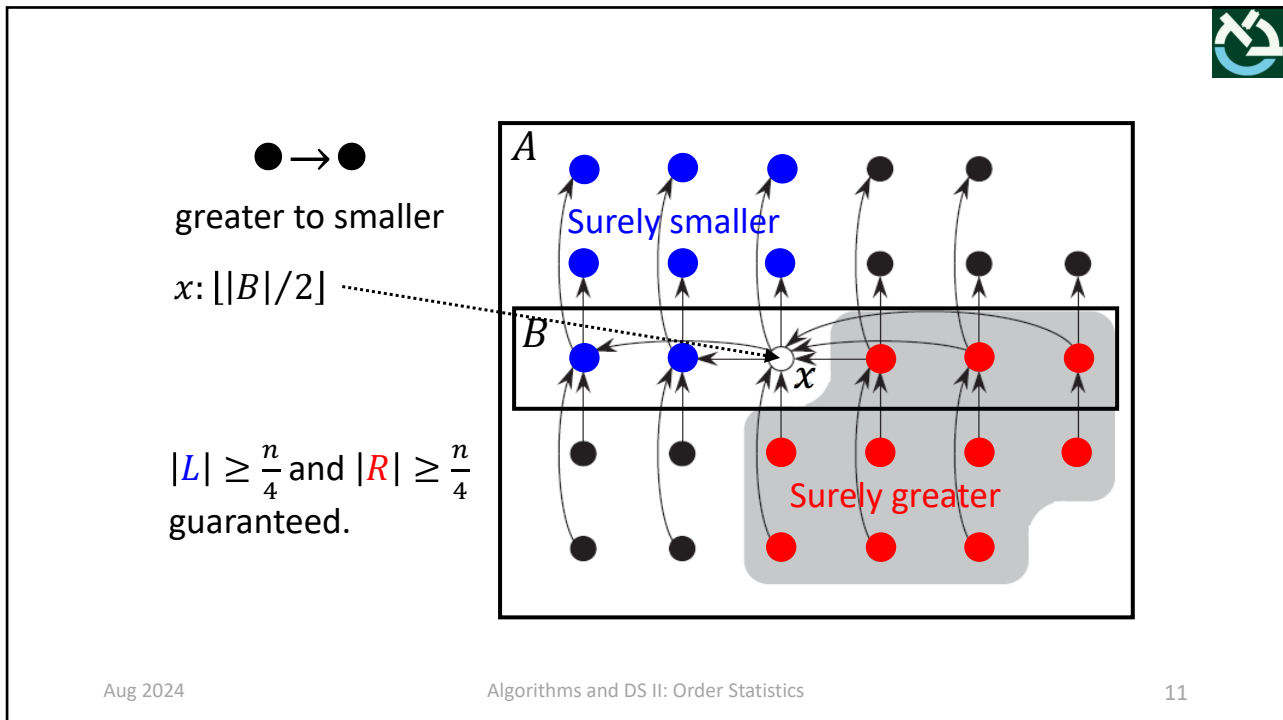
1. Divide input into $\lceil \frac{n}{5} \rceil$ groups of 5 elements, one possibly smaller.
2. Find the median of each group. Let B be the set of all medians.
3. Let $x = \text{SELECT}\left(B, \left\lceil \frac{|B|}{2} \right\rceil\right)$ be B 's median.
4. Divide A around x : $L = \{a \in A \mid a < x\}$, $R = \{a \in A \mid a > x\}$.
5. If $k = |L| + 1$ return x . $k > |L| + 1$
6. If $k \leq |L|$ then SELECT(L, k) else SELECT($R, k - (|L| + 1)$).

Aug 2024

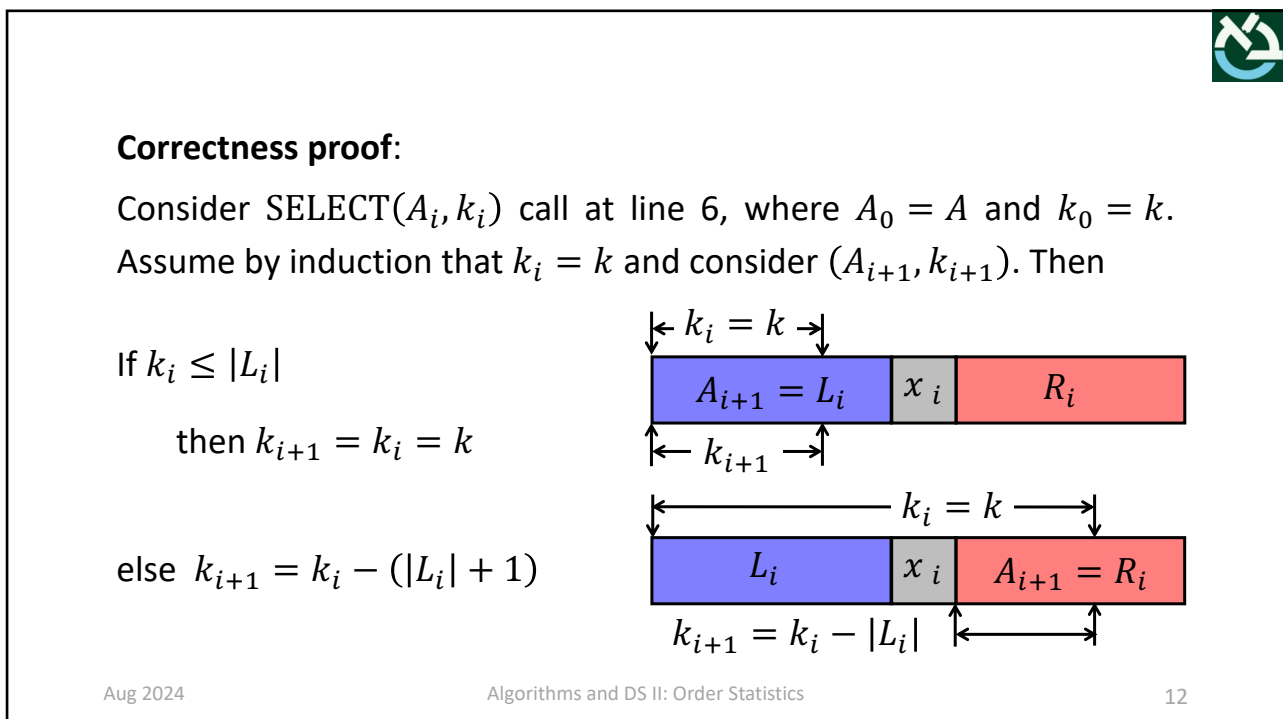
Algorithms and DS II: Order Statistics

10

10



11



12



Run time analysis: Assume w.l.o.g that all the elements are distinct.

Then at least half of the medians in step 2 are not smaller than x .

\Rightarrow At least half of the $\lceil \frac{n}{5} \rceil$ groups contribute 3 elements $> x$, maybe except one smaller group and the group containing x .

$\Rightarrow |R| \geq 3 \left(\frac{1}{2} \lceil \frac{n}{5} \rceil - 2 \right) \geq \frac{3n}{10} - 6$. Similarly, $|L| \geq \frac{3n}{10} - 6$.

Consequently, $\frac{3n}{10} - 6 \leq |L|, |R| \leq \frac{7n}{10} + 6$.

Aug 2024

Algorithms and DS II: Order Statistics

13

13



Let $T(n)$ be the worst-case run time. Steps 1, 2 and 4 take $O(n)$ time.

Steps 3 takes $T(\lceil n/5 \rceil)$ time.

Step 6 takes at most $T(7n/10 + 6)$ time.

The following recurrence is in order

$$T(n) = T(\lceil n/5 \rceil) + T(7n/10 + 6) + O(n).$$

We show that $T(n) = O(n)$.

Assume $T(n) \leq cn$ for sufficiently large n and a constant c . Then

Aug 2024

Algorithms and DS II: Order Statistics

14

14



$$T(n) \leq c(\lceil n/5 \rceil) + c(7n/10 + 6) + an \quad O(n) = an$$

$$\leq c(n/5 + 1) + c(7n/10 + 6) + an$$

$$= 9cn/10 + 7c + an \quad \text{by assumption}$$

$$= cn + (-cn/10 + 7c + an) \leq cn$$

$$\Rightarrow -cn/10 + 7c + an \leq 0$$

$$\Rightarrow (1) \quad c \geq 10a(n/(n - 70)) \rightarrow 10a \text{ as } n \rightarrow \infty.$$

Choosing $c \geq 20a$ will satisfy (1). ■

Aug 2024

Algorithms and DS II: Order Statistics

15

15



What happens if 5 is replaced by 3? 7? (HW).

SELECT is a deterministic algorithm. We can use probabilistic pivot choice in line 3 (as in Quicksort) and then proceed in one side of the partitions.

Show that its expected run time is linear (HW).

Quicksort can use SELECT for balanced partition, yielding worst-case $O(n \log n)$ time deterministic Quicksort.

Impractical because of the large constants in SELECT.

Aug 2024

Algorithms and DS II: Order Statistics

16

16