

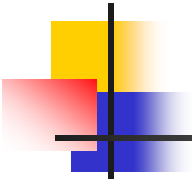


## 9. Medians and order statistics

---

Hsu, Lih-Hsing

Computer Theory Lab.

- 
- The *i*th **order statistic** of a set of  $n$  element is the *i*th smallest.
  - The **selection problem** can be specified formally as follows:
    - **Input:** A set of  $n$  (distinct) numbers and a number  $i$ , with  $1 \leq i \leq n$ .
    - **Output:** the element  $x \in A$  that is larger than exactly  $i - 1$  other elements of  $A$ .



## 9.1 Minimum and Maximum

---

MINIMUM( $A$ )

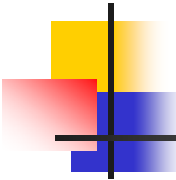
```

1   $min \leftarrow A[1]$ 
2  for  $i \leftarrow 2$  to  $\text{length}[A]$ 
3    do if  $min > A[i]$ 
4      then  $min \leftarrow A[i]$ 
5  return  $min$ 

```

Analysis:  $*O(n)$

- The expected number of times that line 4 is executed is  $O(\log n)$ .

- 
- 
- The second smallest of  $n$  elements can be found with  $n + \lceil \log n \rceil - 2$  comparisons in the worst case.
  - $\left\lceil \frac{3n}{2} \right\rceil - 2$  comparisons are necessary in the worst case to find both the maximum and the minimum of  $n$  elements.

## 9.2 Selection in expected linear time

```

RANDOMIZED_SELECT( $A, p, r, i$ )
1  if  $p = r$ 
2      then return  $A[p]$ 
3   $q \leftarrow \text{RANDOMIZED\_PARTITION}(A, p, r)$ 
4   $k \leftarrow q - p + 1$ 
5  if  $i = k$       ► the pivot value is the answer
6      then return  $A[q]$ 
7  elseif  $i < k$ 
8      then return  $\text{RANDOMIZED\_SELECT}(A, p, q-1, i)$ 
9  else return  $\text{RANDOMIZED\_SELECT}(A, q+1, r, i-k)$ 

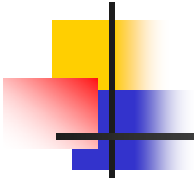
```

## Analysis

For  $k = 1, 2, \dots, n$ , we define indicator random variables  $X_k$  where  $X_k = I \{ \text{the subarray } A[p..q] \text{ has exactly } k \text{ elements} \}$ , and so we have

$$E[X_k] = 1/n .$$

$$\begin{aligned}
 T(n) &\leq \sum_{k=1}^n X_k \cdot (T(\max(k-1, n-k)) + O(n)) \\
 &= \sum_{k=1}^n (X_k \cdot T(\max(k-1, n-k)) + O(n))
 \end{aligned}$$



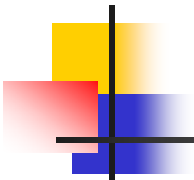
Taking expected values, we have

$$\begin{aligned}
 & E[T(n)] \\
 & \leq E\left[\sum_{k=1}^n X_k \cdot T(\max(k-1)) + O(n)\right] \\
 & = \sum_{k=1}^n E[X_k \cdot T(\max(k-1, n-k))] + O(n) \\
 & = \sum_{k=1}^n E[X_k] \cdot E[T(\max(k-1, n-k))] + O(n) \\
 & = \sum_{k=1}^n \frac{1}{n} \cdot E[T(\max(k-1, n-k))] + O(n) \\
 & \max(k-1, n-k) = \begin{cases} k-1 & \text{if } k > \lceil n/2 \rceil \\ n-k & \text{if } k \leq \lceil n/2 \rceil \end{cases} \\
 & E[T(n)] \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} E[T(k)] + O(n)
 \end{aligned}$$

Chapter 10

P.7

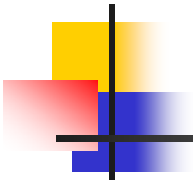
Computer Theory Lab.



$$\begin{aligned}
 E[T(n)] & \leq \frac{2}{n} \sum_{k=\lceil n/2 \rceil}^{n-1} ck + an \\
 & = \frac{2c}{n} \left( \sum_{k=1}^{n-1} k - \sum_{k=1}^{\lfloor n/2 \rfloor} k \right) + an \\
 & = \frac{2c}{n} \left( \frac{(n-1)n}{2} - \frac{(\lfloor n/2 \rfloor - 1)\lfloor n/2 \rfloor}{2} \right) + an \\
 & \leq \frac{2c}{n} \left( \frac{(n-1)n}{2} - \frac{(n/2 - 2)(n/2 - 1)}{2} \right) + an \\
 & = \frac{2c}{n} \left( \frac{n^2 - n}{2} - \frac{n^2/4 - 3n/2 + 2}{2} \right) + an \\
 & = \frac{c}{n} \left( \frac{3n^2}{4} + \frac{n}{2} - 2 \right) + an \\
 & = c \left( \frac{3n}{4} + \frac{1}{2} - \frac{2}{n} \right) + an \\
 & \leq \frac{3cn}{4} + \frac{c}{2} + an \\
 & = cn - \left( \frac{cn}{4} - \frac{c}{2} - an \right).
 \end{aligned}$$

Chapter 10

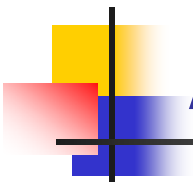
P.8



We choose the constant  $c$  so that  $c/4 - a > 0$ , i.d.,  $c > 4a$ , we can Divide both sides by  $c/4 - a$ , giving

$$n \geq \frac{c/2}{c/4 - a} = \frac{2c}{c - 4a}.$$

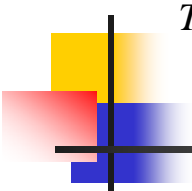
Thus, if we assume that  $T(n) = O(1)$  for  $n < 2c/(c-4a)$ , we have  $T(n) = O(n)$ .



## Another analysis

$$\begin{aligned} T(n) &\leq \frac{1}{n} [T(\max(1, n-1)) + \sum_{k=1}^{n-1} T(\max(k, n-k))] + O(n) \\ &\leq \frac{1}{n} [T(n-1) + 2 \sum_{k=1}^{n-1} T(k)] + O(n) \\ &= \frac{2}{n} \sum_{k=\lceil \frac{n}{2} \rceil}^n T(k) + O(n) \end{aligned}$$

GUESS  $T(n) \leq cn \Rightarrow T(n) = O(n)$



$$T(n) \leq \frac{2}{n} \sum_{k=\lceil \frac{n}{2} \rceil}^{n-1} ck + O(n) \leq \frac{2c}{n} \left( \sum_{k=1}^{n-1} k - \sum_{k=1}^{\lceil \frac{n}{2} \rceil - 1} k \right) + O(n)$$

$$\leq \frac{2c}{n} \left( \frac{(n-1)n}{2} - \frac{\left( \lceil \frac{n}{2} \rceil - 1 \right) \lceil \frac{n}{2} \rceil}{2} \right) + O(n)$$

$$\leq c(n-1) - \frac{c}{n} \left( \frac{n}{2} - 1 \right) \left( \frac{n}{2} \right) + O(n)$$

$$= c \left( \frac{3}{4}n - \frac{1}{2} \right) + O(n) \leq cn$$

Pick  $c$  large enough so that  $c(\frac{n}{4} + \frac{1}{2})$  dominates the  $O(n)$  term.



## 9.3 Selection in worst-case linear time

### Algorithm

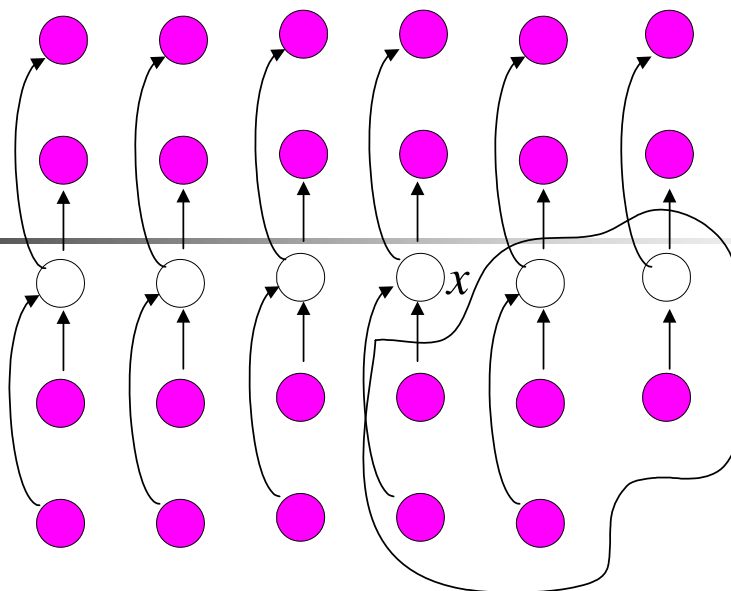
1. Divide the  $n$  elements of the input array into  $\lfloor n/5 \rfloor$  groups of 5 elements each and at most one group made up of the remaining  $n \bmod 5$  elements.
2. Find the **median** of each of the  $\lceil n/5 \rceil$  groups by first insertion sorting the element of each group and picking its median.
3. Use **SELECT** recursively to find the median  $x$  of the  $\lceil n/5 \rceil$  medians found in Step 2.

4. Partition the input array around the *median-of-medians*  $x$  using a modified version of partition. Let  $k$  be the number of elements on the lower side of the partition, so that  $x$  is the  $k$ th smallest element and there are  $n - k$  elements on the high side of the partition.
5. If  $i = k$ , then return  $x$ . Otherwise, use SELECT recursively to find the  $i$ th smallest element on the low side if  $i < k$ , or  $(i-k)$ th smallest element on the high side if  $i > k$ .

Chapter 10

P.13

Computer Theory Lab.



at least

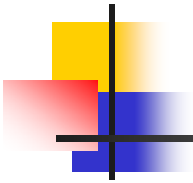
$$3\left(\left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2\right) \geq \frac{3n}{10} - 6$$

nodes

$$T(n) \leq \begin{cases} \Theta(1) & \text{if } n \leq 140 \\ T\left(\frac{n}{5}\right) + T\left(\frac{7n}{10} + 6\right) + O(n) & \text{if } n > 140 \end{cases}$$

Chapter 10

P.14



Assume that  $T(n) \leq cn$  for some  $c$  and small  $n \leq 140$ .

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

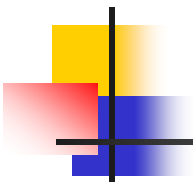
$$\leq cn/5 + c + 7cn/10 + 6c + an$$

$$\leq 9cn/10 + 7c + an$$

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

which is at most  $cn$  if

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



Pick  $c$  large enough so that  $c(\frac{n}{10} - 7)$  is larger than the function described by  $O(n)$  term for all  $n > 140$ .

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

- QUICKSORT can be improved into  $O(n \log n)$ .