# Algorithms
# Chapter 9
# Medians and Order Statistics

Assistant Professor: Ching-Chi Lin

林清池 助理教授

chingchi.lin@gmail.com

Department of Computer Science and Engineering
National Taiwan Ocean University

# Outline

▸ **Minimum and maximum**

▸ Selection in expected linear time

▸ Selection in worst-case linear time

# Order statistics

▸ The $i$th **order statistic** of a set of $n$ elements is the $i$th smallest element.

▸ The **minimum** is the first order statistic ($i = 1$).

▸ The **maximum** is the $n$th order statistic ($i = n$).

▸ A **median** is the "halfway point" of the set.

▸ When $n$ is odd, the median is unique, at $i = (n + 1)/2$.

▸ When $n$ is even, there are two medians:

  ▸ The **lower median**: $i = \lfloor (n + 1)/2 \rfloor$, and
  ▸ The **upper median**: $i = \lceil (n + 1)/2 \rceil$.
  ▸ We mean lower median when we use the phrase "the median".

# The selection problem

▸ How can we find the $i$th order statistic of a set and what is the running time?

▸ **Input:** A set $A$ of $n$ (distinct) number and a number $i$, with $1 \le i \le n$.

▸ **Output:** The element $x \in A$ that is larger than exactly $i-1$ other elements of $A$.

▸ The **selection problem** can be solved in ***O(n lg n)*** time.

  ▸ Sort the numbers using an $O(n \lg n)$-time algorithm, such as heapsort or merge sort.

  ▸ Then return the $i$th element in the sorted array.

▸ Are there faster algorithms?

  ▸ An $O(n)$-time algorithm would be presented in this chapter.

# Finding minimum

▸ We can easily obtain an upper bound of $n-1$ comparisons for finding the minimum of a set of $n$ elements.

> ▸ Examine each element in turn and keep track of the smallest one.
>
> ▸ The algorithm is optimal, because each element, except the minimum, must be compared to a smaller element at least once.

MINIMUM($A$)

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

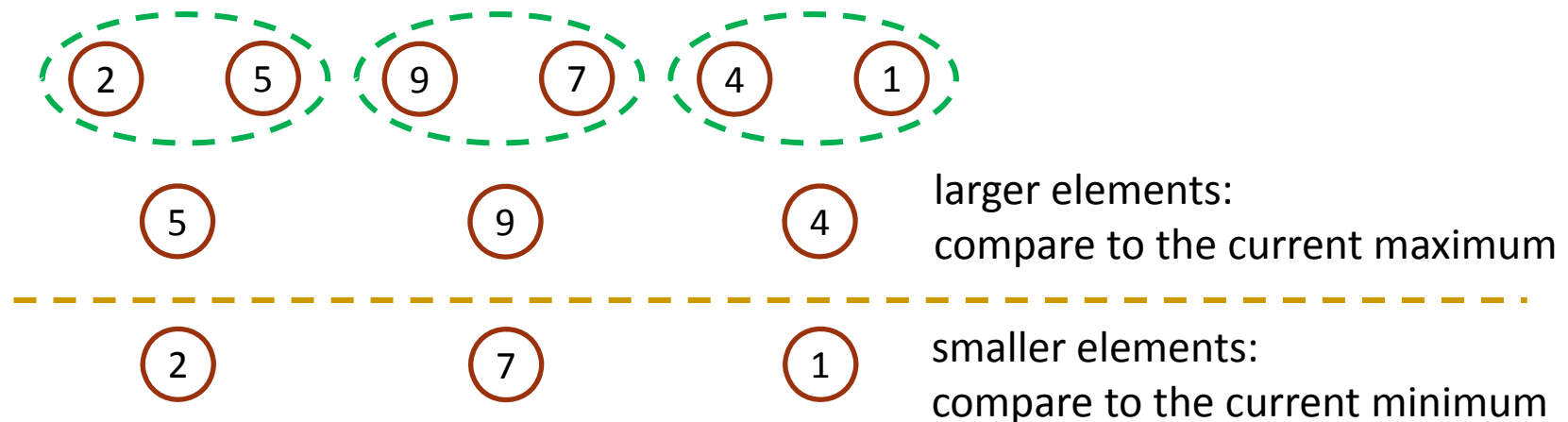▸ The maximum can be found in exactly the same way by replacing the > with < in the above algorithm.

# Simultaneous minimum and maximum$_{1/3}$

- Some applications need both the minimum and maximum.
  - Find the minimum and maximum independently, using $n-1$ comparisons for each, for a total of $2n-2$ comparisons.

- In fact, at most $3\lfloor n/2 \rfloor$ comparisons are needed:
  - Maintain the minimum and maximum of elements seen so far.
  - Process elements in pairs.
  - Compare the elements of a pair to each other.
  - Then compare the larger element to the maximum so far, and compare the smaller element to the minimum so far.

- This leads to only 3 comparisons for every 2 elements.

▶ **An observation**

  ▶ If we compare the elements of a pair to each other, the larger can't be the minimum and the smaller can't be the maximum.

  ▶ So we just need to compare the larger to the current maximum and the smaller to the current minimum.

  ▶ It costs 3 comparisons for every 2 elements.

    ▸ The previous method costs 2 comparisons for each element.



larger elements:
compare to the current maximum

smaller elements:
compare to the current minimum

# Simultaneous minimum and maximum

▸ Setting up the initial values for the min and max depends on whether *n* is odd or even.

  ▸ If *n* is even, compare the first two elements and assign the larger to max and the smaller to min.

  ▸ If *n* is odd, set both min and max to the first element.

▸ If *n* is even, # of comparisons $= \dfrac{3(n-2)}{2} + 1 = \dfrac{3n}{2} - 2.$

▸ If *n* is odd, # of comparisons $= \dfrac{3(n-1)}{2} = 3\lfloor n/2 \rfloor.$

▸ In either case, the # of comparisons is $\leq 3\lfloor n/2 \rfloor.$

# Outline

▸ Minimum and maximum

▸ **Selection in expected linear time**
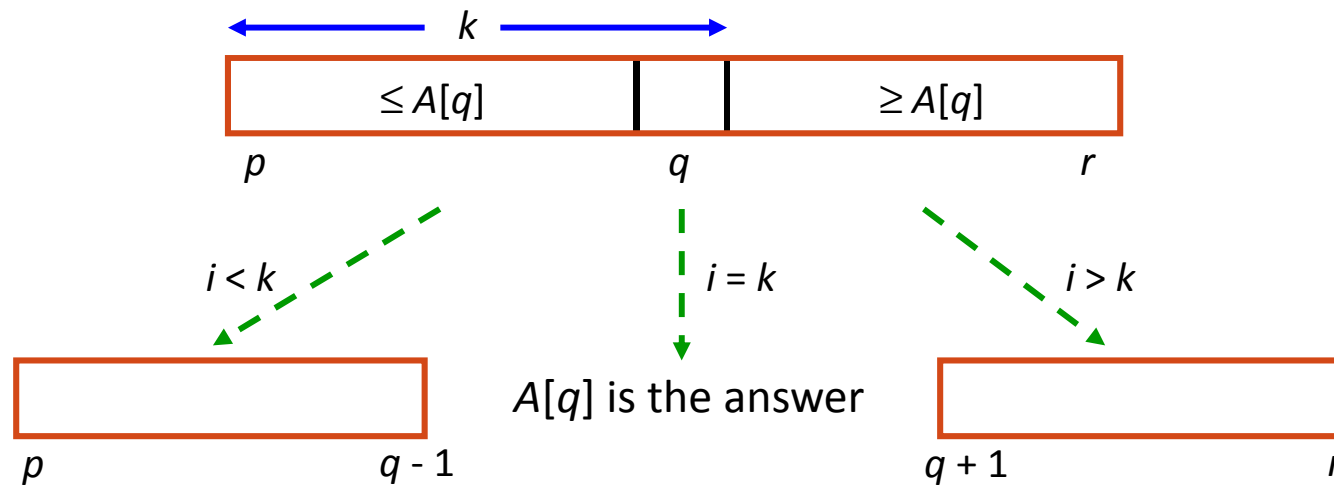
▸ Selection in worst-case linear time

# Selection in expected linear time

▸ In fact, selection of the $i$th smallest element of the array $A$ can be done in $\Theta(n)$ time.

▸ We first present a randomized version in this section and then present a deterministic version in the next section.

▸ The function RANDOMIZED-SELECT:

  ▸ is a divide-and-conquer algorithm,

  ▸ uses RANDOMIZED-PARTITION from the quicksort algorithm in Chapter 7, and

  ▸ recurses on one side of the partition only.

# RANDOMIZED-SELECT procedure$_{1/2}$

1.     RANDOMIZED-SELECT($A$, $p$, $r$, $i$)

2.       **if** $p = r$

3.         **then return** $A[p]$

4.       $q \leftarrow$ RANDOMIZED-PARTITION($A$, $p$, $r$)

5.       $k \leftarrow q - p + 1$

6.       **if** $i = k$      /* the pivot value is the answer */

7.         **then return** $A[q]$

8.       **elseif** $i < k$

9.         **then return** RANDOMIZED-SELECT($A$, $p$, $q - 1$, $i$)

10.       **else return** RANDOMIZED-SELECT($A$, $q$, $r$, $i - k$)

# Algorithm analysis$_{1/4}$

▸ The **worst case:** always recurse on a subarray that is only 1 element smaller than the previous subarray.

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

▸ The **best case:** always recurse on a subarray that has half of the elements smaller than the previous subarray.

  ▸ $T(n) = T(n/2) + \Theta(n)$
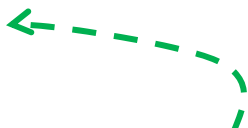     $= \Theta(n)$ (Master Theorem, case 3)

▸ The **average case:**

> ▸ We will show that $T(n) = \Theta(n)$.

> ▸ For $1 \leq k \leq n$, the probability that the subarray $A[p \mathbin{..} q]$ has $k$ elements is $1/n$.

> ▸ To obtain an upper bound, we assume that $T(n)$ is monotonically increasing and that the $i$th smallest element is always in the larger subarray.

> ▸ So, we have

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

$$T(n) \le \frac{1}{n}\sum_{k=1}^{n}(T(\max(k-1,n-k)+O(n)) = \frac{1}{n}\sum_{k=1}^{n}(T(\max(k-1,n-k))) + O(n)$$

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

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

| $k$ | 1 | 2 | ... | $\lceil n/2\rceil$ | $\lceil n/2\rceil+1$ | ... | $n-1$ | $n$ |
|---|---|---|---|---|---|---|---|---|
| max($k-1$, $n-k$) | $n-1$ | $n-2$ | ... | $n-\lceil n/2\rceil$ | $\lceil n/2\rceil$ | ... | $n-2$ | $n-1$ |

▸ If $n$ is even, each term from $T(\lceil n/2\rceil)$ to $T(n-1)$ appears exactly twice.

▸ If $n$ is odd, each term from $T(\lceil n/2\rceil)$ to $T(n-1)$ appears exactly twice and $T(\lfloor n/2\rfloor)$ appears once.

   ▸ Because $k = \lceil n/2\rceil$, $k-1 = n-k = \lfloor n/2\rfloor$.

‣ Solve this recurrence by substitution:

  ‣ Assume $T(n) \leq cn$ for sufficiently large $c$.

  ‣ The function described by the $O(n)$ term is bounded by $an$ for all $n > 0$.

  ‣ Then, we have

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

chose $c$ so that
$$c/4 - a > 0$$

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

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

$$\frac{cn}{4} - \frac{c}{2} - an \geq 0$$

$$= \frac{c}{n}\left(\frac{3n^2}{4} + \frac{n}{2} - 2\right) + an = c\left(\frac{3n}{4} + \frac{1}{2} - \frac{2}{n}\right) + an$$

$$n\left(\frac{c}{4} - a\right) \geq \frac{c}{2}$$

$$\leq \frac{3cn}{4} + \frac{c}{2} + an = cn - \left(\frac{cn}{4} - \frac{c}{2} - an\right).$$

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

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

# Outline

▶ Minimum and maximum

▶ Selection in expected linear time
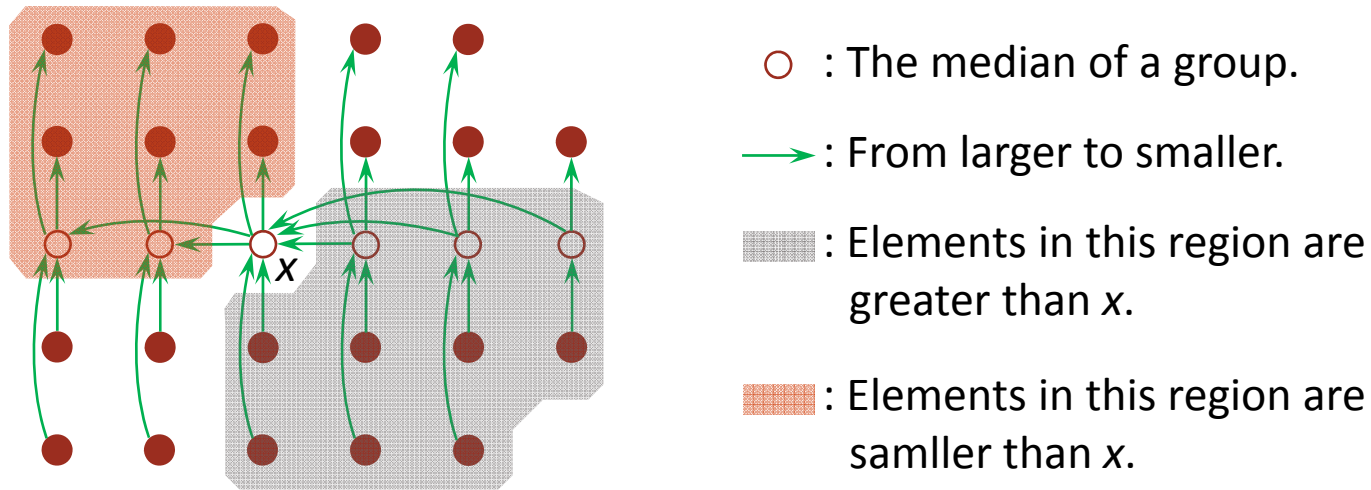
▶ **Selection in worst-case linear time**

# SELECT algorithm$_{1/2}$

▸ **Idea:** Guarantee a good split when the array is partitioned.

▸ The SELECT algorithm:

1. Divide $n$ elements into groups of 5 elements.

2. Find median of each of the $\lceil n/5 \rceil$ groups.

   ▸ Run insertion sort on each group.

   ▸ Then just pick the median from each group.

3. Use SELECT recursively to find median $x$ of the $\lceil n/5 \rceil$ medians.

4. Partition the $n$ elements around $x$.

   ▸ Let $x$ be the $k$th element of the array after partitioning.

   ▸ There are $k-1$ elements on the low side of the partition and $n-k$ elements on the high side.

5. Now there are three possibilities:

▶ If $i = k$, then return $x$.

▶ If $i < k$, then use SELECT recursively to find $i$th smallest element on the low side.

▶ If $i > k$, then use SELECT recursively to find $(i-k)$th smallest element on the high side.

○ : The median of a group.

⟶ : From larger to smaller.

: Elements in this region are greater than $x$.

: Elements in this region are samller than $x$.

# Time complexity$_{1/4}$

▸ At least half of the medians are $\geq x$.

  ▸ Precisely, at least $\lceil \lceil n/5 \rceil / 2 \rceil$ medians $\geq x$.

▸ These group contributes 3 elements that are $> x$, except for 2 of the groups:

  ▸ the group containing $x$, and

  ▸ the group with $< 5$ elements.

▸ The number of elements greater than $x$ is at least:

  ▸ $3 \left( \lceil \lceil n/5 \rceil / 2 \rceil - 2 \right) \geq$ **$3n/10 - 6$**.

▸ Similarly, at least $3n/10 - 6$ elements are less than $x$.

▸ Thus, SELECT is called recursively on $\leq$ **$7n/10 + 6$** elements in step 5.

# Time complexity<sub>3/4</sub>

▶ The SELECT algorithm:

1. Divide $n$ elements into groups of 5 elements. $O(n)$

2. Find median of each of the $\lceil n/5 \rceil$ groups. $O(n)$

   ▶ Run insertion sort on each group.

   ▶ Then just pick the median from each group.

3. Use SELECT recursively to find median $x$ of the $\lceil n/5 \rceil$ medians. $T(\lceil n/5 \rceil)$

4. Partition the $n$ elements around $x$. $O(n)$

   ▶ Let $x$ be the $k$th element of the array after partitioning.

   ▶ There are $k - 1$ elements on the low side of the partition and $n - k$ elements on the high side.

5. Now there are three possibilities: $T(7n/10 + 6)$

   ▶ If $i = k$, then return $x$.

   ▶ If $i < k$, then use SELECT recursively to find $i$ th smallest element on the low side.

   ▶ If $i > k$, then use SELECT recursively to find $(i - k)$th smallest element on the high side.

▶ Time complexity: $T(n) \leq T(\lceil n/5 \rceil) + T(7n/10 + 6) + O(n)$.

# Time complexity$_{4/4}$

‣ Solve this recurrence by substitution:

  ‣ Assume $T(n) \leq cn$ for sufficiently large $c$.

  ‣ The function described by the $O(n)$ term is bounded by $an$ for all $n > 0$.

  ‣ Then, we have

$$
\begin{aligned}
T(n) &\leq c\lceil n/5 \rceil + c(7n/10 + 6) + an \\
&\leq cn/5 + c + 7cn/10 + 6c + an \\
&= 9cn/10 + 7c + an \\
&= cn + (-cn/10 + 7c + an)
\end{aligned}
$$

  ‣ This last quantity is $\leq cn$ if we choose $c \geq 20a$.

$$
\begin{aligned}
-cn/10 + 7c + an &\leq 0 \\
cn/10 - 7c &\geq an \qquad \text{Notice: } (n/n\text{-}70) \leq 2 \text{ for } n \geq 140. \\
c(n - 70) &\geq 10an \\
c &\geq 10a(n/(n - 70))
\end{aligned}
$$