

Chapter 9: Medians and Order Statistics

Chi-Yeh Chen

陳奇業

成功大學資訊工程學系



藏行顯光
成就共好

Achieve Securely
Prosper Mutually



國立成功大學 九十週年
90th Anniversary of NCKU

Chapter 9 overview

- The ***i*th order statistic** is the *i*th smallest element of a set of *n* elements
- 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**, at $i = n/2$, and
 - The **upper median**, at $i = n/2 + 1$.
 - We mean lower median when we use the phrase “the median.”

The selection problem

Input: A set A 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 in A .
In other words, the i th smallest element 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 heap sort or merge sort.
- Then return the i th element in the sorted array.

Minimum and Maximum



藏行顯光
成就共好
Achieve Securely
Prosper Mutually
國立成功大學 九十週年
Anniversary of NCKU

Minimum and Maximum

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.
- This is the best we can do, because each element, except the minimum, must be compared to a smaller element at least once.

Minimum and Maximum

The following pseudocode finds the minimum element in array $A[1..n]$:

MINIMUM(A, n)

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

Minimum and Maximum

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

MAXIMUM(A, n)

```
1  $max \leftarrow A[1]$ 
2 for  $i \leftarrow 2$  to  $n$  do
3     if  $max < A[i]$  then
4          $max \leftarrow A[i]$ 
5 return  $max$ 
```

Simultaneous minimum and maximum

- A simple algorithm to find the minimum and maximum: $n-1$ comparisons for the minimum and $n-1$ comparisons for maximum \Rightarrow total $2n-2$ comparisons $\Rightarrow \Theta(n)$ time.
- In fact, at most $3 \lfloor n/2 \rfloor$ comparisons are needed to find both the minimum and maximum !

Simultaneous minimum and maximum (cont.)

- 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.

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. Then process the rest of the elements in pairs.
- If n is odd, set both min and max to the first element. Then process the rest of the elements in pairs.

Analysis of the total number of comparisons

- If n is even, we do 1 initial comparison and then $3(n-2)/2$ more comparisons.

$$\begin{aligned}\text{\# of comparisons} &= \frac{3(n-2)}{2} + 1 \\ &= \frac{3n-6}{2} + 1 \\ &= \frac{3n^2}{2} - 3 + 1 \\ &= \frac{3n^2}{2} - 2\end{aligned}$$

- If n is odd, we do $3\lfloor n/2\rfloor$ comparisons.

In either case, the maximum number of comparisons is $\leq 3\lfloor n/2\rfloor$

Selection in expected linear time



Selection in expected linear time

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

```
1 if  $p = r$  then
2   return  $A[p]$ 
3  $q \leftarrow \text{RANDOMIZED-PARTITION}(A, p, r)$ 
4  $k \leftarrow q - p + 1$ 
5 if  $i = k$  then
6   ▶ pivot value is answer
7   return  $A[q]$ 
8 else if  $i < k$  then
9   return  $\text{RANDOMIZED-SELECT}(A, p, q - 1, i)$ 
10 else
11   return  $\text{RANDOMIZED-SELECT}(A, q + 1, r, i - k)$ 
```



Selection in expected linear time

After the call to RANDOMIZED-PARTITION, the array is partitioned into two subarrays $A[p \dots q - 1]$ and $A[q + 1 \dots r]$, along with a **pivot** element $A[q]$.

- The elements of subarray $A[p \dots q - 1]$ are all $\leq A[q]$.
- The elements of subarray $A[q + 1 \dots r]$ are all $> A[q]$.

Selection in expected linear time

- The pivot element is k th element of subarray $A[p \dots q]$, where $k = q - p + 1$.
- If the pivot element is the i th smallest element (i.e., $i = k$), return $A[q]$.
- Otherwise, recurse on the subarray containing the i th smallest element.
 - If $i < k$, this subarray is $A[p \dots q - 1]$, and we want the i th smallest element.
 - If $i > k$, this subarray is $A[q + 1 \dots r]$ and, since there are k elements in $A[p \dots r]$ that precede $A[q + 1 \dots r]$, we want the $(i - k)$ th smallest element of this subarray.

Analysis

Worst-case running time: $\Theta(n^2)$, because we could be extremely unlucky and always recurse on a subarray that is only 1 element smaller than the previous subarray.

Expected running time: RANDOMIZED-SELECT works well on average.

Analysis

The running time of RANDOMIZED-SELECT is a random variable $T(n) \Rightarrow$ to obtain $E[T(n)]$ as follows:

- For $k = 1, 2, \dots, n$, define indicator random variable $X_k = I\{\text{subarray } A[p \dots q] \text{ has exactly } k \text{ elements}\}$
- Since $\Pr\{\text{subarray } A[p \dots q] \text{ has exactly } k \text{ elements}\} = 1/n$, Lemma 5.1 says that $E[X_k] = 1/n$.

Analysis

Therefore, we have the recurrence

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

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

Analysis

- Taking expected values gives

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

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

(linearity of expectation)

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

(equation (C.23))

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

Analysis

- We rely on X_k and $T(\max(k-1, n-k))$ being independent random variables in order to apply equation (C.23).
- Looking at the expression $\max(k-1, n-k)$ we have
 - If n is even, each term from $T(\lceil n/2 \rceil)$ up to $T(n-1)$ appears exactly twice in the summation.
 - If n is odd, these terms appear twice and $T(\lfloor n/2 \rfloor)$ appears once.

- Either way,
$$E[T(n)] \leq \frac{2}{n} \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} E[T(k)] + O(n)$$

Analysis

- Solve this recurrence by substitution:
 - Guess that $T(n) \leq cn$ for some constant c that satisfies the initial conditions of the recurrence
 - Assume that $T(n) = O(1)$ for $n < \text{some constant}$
 - Also pick a constant a such that the function described by the $O(n)$ term is bounded from above by an for all $n > 0$

Analysis

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



Analysis

$$\begin{aligned} &= \frac{2c}{n} \left(\frac{n^2 - n}{2} - \frac{n^2/4 - 3n/4 + 2}{2} \right) + an = \frac{c}{n} \left(\frac{3n^2}{4} + \frac{n}{2} - 2 \right) + an \\ &= c \left(\frac{3n}{4} + \frac{1}{2} - \frac{2}{n} \right) + an \leq \frac{3cn}{4} + \frac{c}{2} + an = cn - \left(\frac{cn}{4} - \frac{c}{2} - an \right) \end{aligned}$$



Analysis

To complete this proof, we choose c such that

$$\frac{cn}{4} - \frac{c}{2} - an \geq 0 \quad \frac{cn}{4} - an \geq \frac{c}{2}n \left(\frac{c}{4} - a \right) \geq \frac{c}{2}n \geq \frac{c/2}{c/4 - a}n \geq \frac{2c}{c/4 - 4a}$$



Analysis

- Thus, as long as we assume that $T(n) = O(1)$ for $n < 2c/(c - 4a)$, we have $E[T(n)] = O(n)$.
- Therefore, we can determine any order statistic in linear time on average.

Selection in worst-case linear time



NCKU
National Cheng Kung University

Selection in worst-case linear time

SELECT works on an array of $n > 1$ elements:

1. Divide the n elements into groups of 5. Get $\lceil n/5 \rceil$ groups:

- $\lfloor n/5 \rfloor$ groups with exactly 5 elements and, if 5 does not divide n
- one group with the remaining $n \bmod 5$ elements.

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

- Run insertion sort on each group. Takes $O(1)$ time per group since each group has ≤ 5 elements.
- Then just pick the median from each group, in $O(1)$ time.

Selection in worst-case linear time

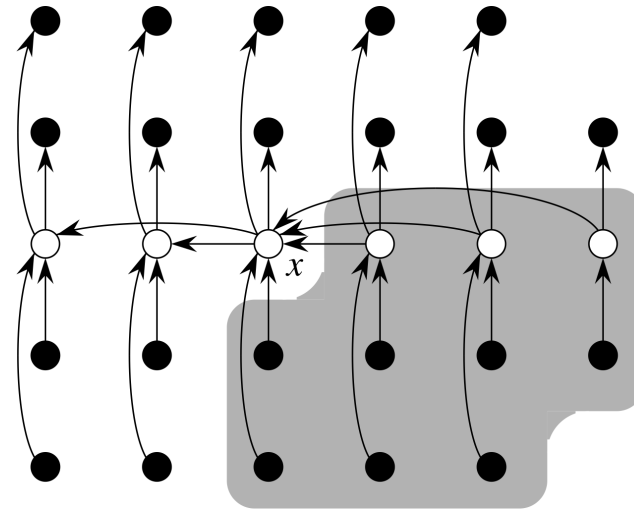
3. Find the median x of the $\lceil n/5 \rceil$ medians by a recursive call to SELECT. (If $\lceil n/5 \rceil$ is even, then follow our convention and find the lower median.)
4. Using the modified version of PARTITION that takes the pivot element as input, partition the input array around x . Let x be the k th element of the array after partitioning, so that there are $k - 1$ elements on the low side of the partition and $n - k$ elements on the high side.

Selection in worst-case linear time

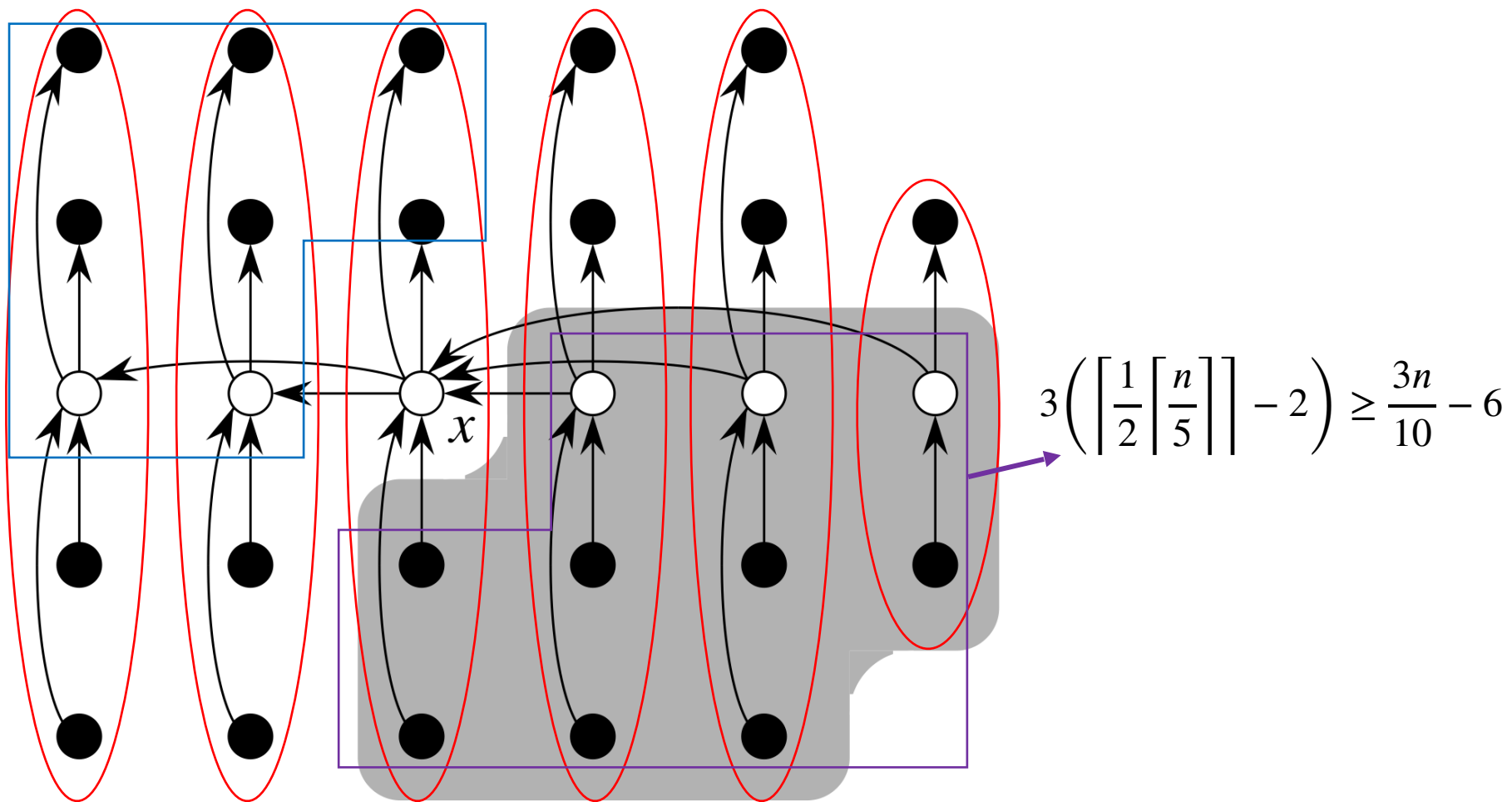
5. Now there are three possibilities:

- If $i = k$, just return x .
- If $i < k$, return the i th smallest element on the low side of the partition by making a recursive call to SELECT.
- If $i > k$, return the $(i - k)$ th smallest element on the high side of the partition by making a recursive call to SELECT.

- Start by getting a lower bound on the number of elements that are greater than the partitioning element x :



[Each group is a column. Each white circle is the median of a group, as found in step 2. Arrows go from larger elements to smaller elements, based on what we know after step 4. Elements in the region on the lower right are known to be greater than x .]



Analysis

- At least half of the median found in step 2 are $\geq x$.
- Look at the groups containing these medians that are $\geq x$. All of them contribute 3 elements that are $> x$ (the median of the group and the 2 elements in the group greater than the group's median), except for 2 of the groups: the group containing x (which has only 2 elements $> x$) and the group with < 5 elements.
- Forget about these 2 groups. That leaves $\geq \left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2$ groups with 3 elements known to be $> x$.
- Thus, we know that at least
$$3 \left(\left\lceil \frac{1}{2} \left\lceil \frac{n}{5} \right\rceil \right\rceil - 2 \right) \geq \frac{3n}{10} - 6$$
elements are $> x$.

Analysis

Symmetrically, the number of elements that are $< x$ is $\geq 3n/10 - 6$.

Therefore, **when we call SELECT recursively in step 5, it's on $\leq 7n/10 + 6$ elements.**

Develop a recurrence for the worst-case running time of SELECT:

- Steps 1, 2, and 4 each take $O(n)$ time:
 - Step 1: making groups of 5 elements takes $O(n)$ time.
 - Step 2: sorting $\lceil n/5 \rceil$ groups in $O(1)$ time each.
 - Step 4: partitioning the n -element array around x takes $O(n)$ time.
- Step 3 takes time $T(\lceil n/5 \rceil)$.
- Step 5 takes time $\leq T(7n/10 + 6)$, assuming that $T(n)$ is monotonically increasing.

Analysis

- Assume that $T(n) = O(1)$ for small enough n . We'll use $n \leq 140$ as “small enough.”
- Thus, we get the recurrence

$$T(n) \leq \begin{cases} O(1) & \text{if } n < 140 \\ T(\lceil n/5 \rceil) + T(7n/10 + 6) + O(n) & \text{if } n \geq 140 \end{cases}$$

Solve this recurrence by substitution:

- **Inductive hypothesis:** $T(n) \leq cn$ for some constant c and all $n > 0$.
- Assume that c is large enough that $T(n) \leq cn$ for all $n < 140$. So we are concerned only with the case $n \geq 140$.
- Pick a constant a such that the function described by the $O(n)$ term in the recurrence is $\leq an$ for all $n > 0$.

Analysis

Substitute the inductive hypothesis in the right-hand side of the recurrence:

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



Analysis

This last quantity is $\leq cn$ if

$$-\frac{cn}{10} + 7c + an \leq 0 \quad \frac{cn}{10} - 7c \geq an \quad cn - 70c \geq 10an \quad (n - 70) \geq 10a$$

$$c \geq 10a \left(\frac{n}{n - 70} \right)$$

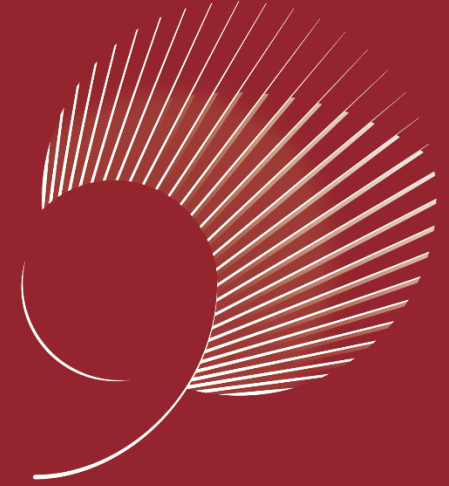
Analysis

- Because we assumed that $n \geq 140$, we have $n/(n - 70) \leq 2$.
- Thus, $20a \geq 10an/(n - 70)$, so choosing $c \geq 20a$ gives $c \geq 10an/(n - 70)$, which in turn gives us the condition we need to show that $T(n) \leq cn$.
- We conclude that $T(n) = O(n)$, so that SELECT runs in linear time in all cases.

Analysis

Notice that SELECT and RANDOMIZED-SELECT determine information about the relative order of elements only by comparing elements.

- Sorting requires $\Omega(n \lg n)$ time in the comparison model.
- Sorting algorithms that run in linear time need to make assumptions about their input.
- Linear-time selection algorithms do not require any assumptions about their input.
- Linear-time selection algorithms solve the selection problem without sorting and therefore are not subject to the $\Omega(n \lg n)$ lower bound.



藏行顯光 成就共好

Achieve Securely
Prosper Mutually



國立成功大學 九十週年
90th Anniversary of NCKU



國立成功大學
National Cheng Kung University