

Median and Order Statistics

Andres Mendez-Vazquez

September 30, 2018

Contents

1	Introduction	2
2	Selection Problem	2
3	On the way to a better solution	2
4	Selection in expected linear time	3
5	Analysis of the Worst Case Selection	4

1 Introduction

One of the most important tasks when analyzing a collection of numbers is to find the i th smallest element. Examples of these elements are

- The minimum, $i=1$.
- The maximum, $i = n$.
- The Median for n odd, $i = \frac{n+1}{2}$.

Then, it is necessary to find fast solutions to find this statistics.

2 Selection Problem

There is a clever way to see the finding of the i th statistics as a selection problem (Look at the slides).

- Input: A set A of n (distinct) numbers and an integer i , with $1 \leq i \leq n$.
- Output: The element $x \in A$ that is larger than exactly $i-1$ other elements of A .

Here, a classic way to solve the problem is to sort the elements with merge sort. Then, find the statistics. However, we can do much better.

3 On the way to a better solution

In order to find the minimum alone, using simple comparisons in a sequence of n numbers, the naive algorithm can find that element in n . However, if we try to find the maximum and minimum at the same time, we can do much better:

- Take two elements at the same time.
- Compare them between them to get the min and max in the tuple.
- Compare the smallest with the actual minimum do the same with the biggest one.

This algorithm takes 3 comparisons per sets of two elements the we are bounded by $3 \lfloor \frac{n}{2} \rfloor$. For example, if n is even we have one initial comparison followed by

$$3 \left(\frac{n-2}{2} \right) + 1$$

This is giving us that the total number of comparisons is at most $3 \lfloor \frac{n}{2} \rfloor$

4 Selection in expected linear time

The analysis of the Randomized-selection can be done assuming the following:

- $X_k = I \{\text{the subarray } A[p \dots q] \text{ has exactly } k \text{ elements}\}$ with $E[X_k] = \frac{1}{n}$
(Assuming that the elements are distinct)

Now, we need to bound the recursive function $T(n)$ describing the Randomized-select algorithm. This can be done, if we assume that

- The i th element is always in the largest partition size.

Therefore,

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

Then, we take the expected value:

$$\begin{aligned} E[T(n)] &\leq E \left[\sum_{k=1}^n X_k \times (T(\max(k-1, n-k)) + O(n)) \right] \\ &= \sum_{k=1}^n E[X_k \times (T(\max(k-1, n-k)))] + O(n) \\ &= \sum_{k=1}^n E[X_k] E[T(\max(k-1, n-k))] + O(n) \\ &= \sum_{k=1}^n \frac{1}{n} \times E[T(\max(k-1, n-k))] + O(n) \end{aligned}$$

Now if we take in account that

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

This means that

- if n is even each term $T(\lceil \frac{n}{2} \rceil)$ to $T(n-1)$ appears exactly twice.
- if n is odd each term appears twice and $T(\lfloor \frac{n}{2} \rfloor)$ appears once.

Then, we have the following recursion:

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

Thus, if we assume that $E[T(n)] \leq cn$

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

Now, we need that $\frac{cn}{4} - \frac{c}{2} - an \geq 0$ or $n \left(\frac{c}{4} - a \right) \geq \frac{c}{2} > 0$. Thus, $\frac{c}{4} - a > 0$ i.e. we choose $c > 4a$. Therefore, $n \geq \frac{\frac{c}{2}}{\frac{c}{4} - a} = \frac{2c}{c-4a}$. Thus, if we assume that $T(n) = O(1)$, we have that $E[T(n)] = O(n)$.

5 Analysis of the Worst Case Selection

For this, we look at the following image

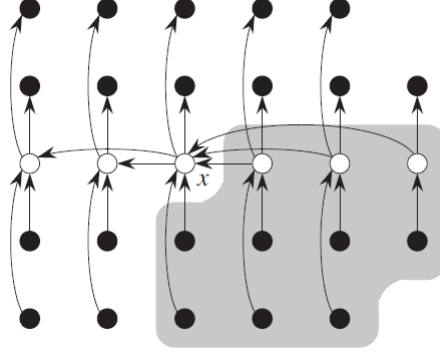


Figure 1: Look at this

1. At least half of the medians are greater than or equal to the median-of-medians x .
2. At least half of the $\lceil \frac{n}{5} \rceil$ groups contribute with at least 3 elements greater than x , except for two groups: one that has less than 5 elements and the one containing x .
3. Then, the number of elements greater than x is at least $3 \left(\lceil \frac{1}{2} \lceil \frac{n}{5} \rceil \rceil - 2 \right) \geq \frac{3n}{10} - 6$.
4. Similarly, at least $\frac{3n}{10} - 6$ elements are less than x .
5. We can then in the worst case we have select is called recursively in at most $\frac{3n}{10} - 6 + \frac{3n}{10} - 6 < \frac{7n}{10} + 6$.
6. Steps 1, 2 and 4 take $O(n)$ time.

Then, if we assume that

- $a < b \Rightarrow T(a) < T(b)$ (Monotonically increasing).

Thus, we have the following recurrence:

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