

Suppose we are given an array of n distinct elements. The i th **order statistic** is the i th smallest element.

- $i = 1$: minimum element
- $i = n$: maximum element
- $i = \lfloor (n + 1)/2 \rfloor$: median

Trivial algorithm: Sort the input using merge sort. Return the element at index i . Worst case running time is $\Theta(n \log n)$.

Can we develop an algorithm which has a better worst-case running time than the trivial algorithm?

Consider a divide and conquer algorithm for this problem:

1. Use a randomized partitioning scheme similar to the one used in quicksort (choose a pivot and partition the elements into two sets “around” the pivot).
2. Let k be the index of the pivot after partitioning.
3. If $i = k$, then the pivot is the element we are looking for, and we return the pivot.
4. If $i < k$, then we know that the element we are looking for is in the first subarray, and we recursively find the i th smallest element in this subarray.
5. If $i > k$, then we now that the element we are looking for is in the second subarray, and we recursively find the element in position $i - k$ in this subarray.

$i = 7$

6	10	13	5	8	3	2	11
---	----	----	---	---	---	---	----

↑
pivot

2	5	3	6	8	13	10	11
---	---	---	---	---	----	----	----

↑
pivot

Next Subproblem

$$\text{update: } i = i - k \Rightarrow i = 7 - 4 = \textcircled{3}$$

$\text{Rand-Select}(A, p, q, i)$

{

if $p == q$, return $A[p]$

$r = \text{Rand-Partition}(A, p, q)$

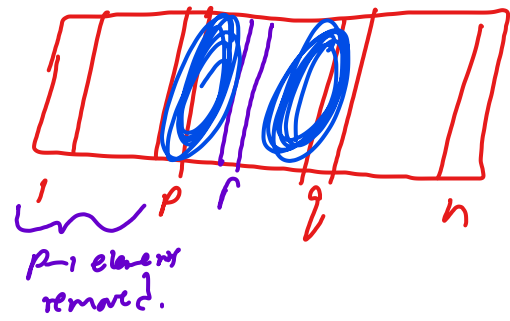
$k = r - (p - 1)$

if $i == k$, return $A[r]$

else if $i < k$, return $\text{Rand-Select}(A, p, r-1, i)$

else, return $\text{Rand-Select}(A, r+1, q, i-k)$

}



Analysis:

If we choose a good pivot whose largest subarray is of size at most $9n/10$:

- $T(n) \leq T(9n/10) + dn$

$$a=1, b=\frac{10}{9}, f(n)=n^1, n^{\log_b a} = n^0$$

$$n^0 \in \Omega(n^{0+\epsilon}) \text{ for } \epsilon=1$$

$$af\left(\frac{1}{b}\right) = \frac{9n}{10} \Rightarrow C = \frac{9}{10} \checkmark$$

$$\text{Case 3 holds, } T(n) \in \Theta(n)$$

If we get unlucky and pick the worst possible pivot:

- $T(n) = T(n-1) + dn$

$$\begin{array}{c} n \\ | \\ n-1 \\ | \\ n-2 \\ | \\ \vdots \\ | \end{array}$$

$$\sum_{i=1}^n i = \Theta(n^2)$$

What is the expected running time of the algorithm?

The running time will depend on the sizes of the subproblems we generate. The two subarrays computed by partition will be of size $(k, n - k - 1)$ for some $k \in \{0, 1, \dots, n - 1\}$.

To obtain an upper bound on the running time, we will assume that the i th element always falls in the larger subarray (i.e. the subproblem size will be $\max(k, n - k - 1)$).

Thus we can express the running time of the algorithm in the following way:

$$T(n) = \begin{cases} T(\max(0, n-1)) + dn & \text{if } 0 \leq n-1 \text{ split} \\ T(\max(1, n-2)) + dn & \text{if } 1 \leq n-2 \text{ split} \\ \vdots \\ T(\max(n-1, 0)) + dn & \text{if } n-1 \leq 0 \text{ split} \end{cases}$$

We can express the running time as a sum, if we introduce the following indicator RV X_k :

$$X_k = \begin{cases} 1 & \text{if } k:n-k-1 \text{ split} \\ 0 & \text{otherwise} \end{cases}$$

Using these indicator RVs, we have that the running time of the algorithm is:

$$\begin{aligned} T(n) &= \sum_{k=0}^{n-1} X_k \cdot (T(\max(k, n-k-1)) + dn) \\ &\leq 2 \cdot \sum_{k=\lceil \frac{n}{2} \rceil}^{n-1} X_k \cdot (T(k) + dn) \end{aligned}$$

Calculating $E[T(n)]$:

$$E[T(n)] \leq E\left[2 \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} X_k \cdot (T(k) + d_n)\right]$$

$$= 2 \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} E[X_k \cdot (T(k) + d_n)]$$

$$= 2 \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} (E[X_k] \cdot E[T(k) + d_n])$$

↑ $\frac{1}{n}$

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

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

$\sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} d_n = d_n$

$$= \left(\frac{2}{n} \sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} E[T(k)] \right) + d_n$$

Note that the subproblem $T(k)$ is a RV (which depends on pivot choices for the future subproblems). We can use induction to show that $E[T(n)] = \underline{O(n)}$.

Inductively Assume: $E[T(k)] \leq C \cdot k$ for some $C > 0$
for all $k < n$.

We will use the following fact: $\sum_{k=\lfloor \frac{n}{2} \rfloor}^{n-1} k \leq \frac{3}{8} n^2$.

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

$$\leq \frac{\cancel{2}C}{\cancel{n}} \frac{3\cancel{n}^2}{\cancel{8}4} + dn$$

$$= \frac{3}{4} Cn + dn$$

$$= \underbrace{Cn}_{\text{desired}} - \underbrace{\frac{1}{4}Cn + dn}_{\text{residual}}$$

$$\leq Cn$$

↑ true when $-\frac{1}{4}Cn + dn \leq 0$
 $\boxed{\frac{d}{C} \leq 1}$

The randomized algorithm described is excellent in practice (linear expected running time); however, the worst case running time ($\Theta(n^2)$) is slower than the trivial algorithm.

Is it possible to obtain an algorithm whose worst-case running time is better than the $\Theta(n \log n)$ running time of merge sort?

Answer is yes [Blum, Floyd, Pratt, Rivest, and Tarjan - 1973].

The idea is to recursively generate a good pivot.

So the running time is $T(n) = T(n/5) + T(7n/10 + 3) + dn$. We will use induction to show that $T(n) = O(n)$.

The worst case running time is $O(n)$, but the hidden constants are quite large. The algorithm is mainly of theoretical interest.

In practice, it is better to use the randomized algorithm than this algorithm.

