# Median Finding Algorithms

There are a few approaches to finding the median of a list. One approach is to sort the elements of the list and then compute the position of the median. That approach takes $\Theta(n \log n)$ time.
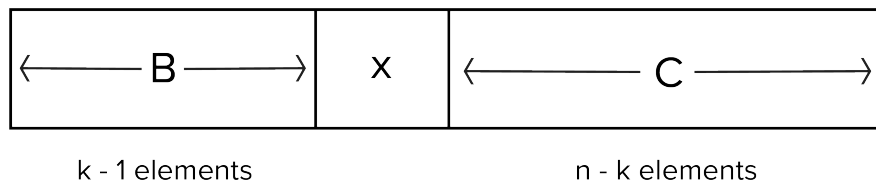
Another approach inspired by quicksort is to find a pivot, and recurse on a subarray based on the rank of the pivot. The corresponding pseudocode is presented below.

```
1    SELECT(S, i):
2      Pick x in S cleverly
3      Compute k = rank(x)
4      B = Set of elements in S, such that elements < x
5      C = Set of elements in S, such that elements > x
6      if k == i:
7        return x
8      else if k > i:
9        return SELECT(B, i)
10     else if k < i:
11       return SELECT(C, i - k)
```

The intuition here is if we can eliminate either B or C cleverly, and prove that the subarray we recurse on is always a fraction of the previous array, the algorithm will run in $\Theta(n)$ time.



k - 1 elements          n - k elements

Without a clever strategy, if we're selecting the $n - 1$ element and $k = 1$ every time, then we're eliminating only one element on every iteration. There would be $n$ iterations, each taking $\Theta(n)$ time to partition the subarray. Therefore, the algorithm would run in $\Theta(n^2)$.

**Clever Partitioning**

1. Arrange $S$ into columns of size 5 ($\lceil n/5 \rceil$ columns).

2. Sort each column in linear time (top down)

3. Find the medians of medians as $x$.

Half of the $\lceil n/5 \rceil$ groups contribute at least 3 elements $> x$ except for 1 trailing group which could contain fewer than 5 elements and 1 group containing $x$. Therefore, there are at least $3(\lceil n/10 \rceil - 2)$ elements $> x$ and at least $3(\lceil n/10 \rceil - 2)$ elements $< x$.
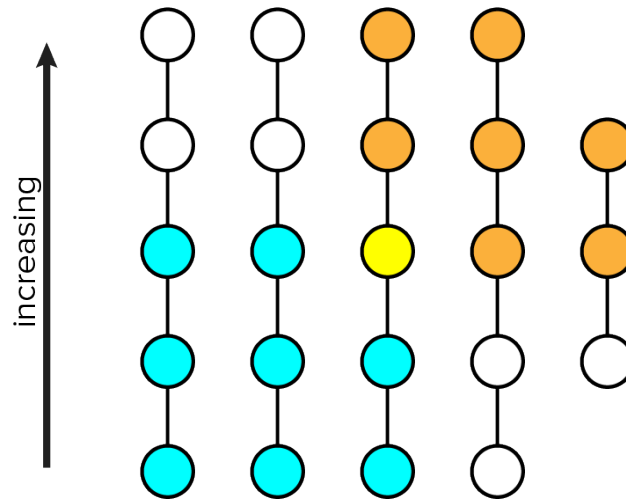
1

Figure 1: The cyan circles are the elements lower than $x$ (denoted by the yellow circle) and the orange circles are the elements greater than $x$

The recurrence can now be written as

$$f(n) = \begin{cases} \Theta(1), & \text{for } n \leq 140 \\ T(\lceil n/5 \rceil) + T(\frac{7n}{10} + 6) + \Theta(n), & \text{for } n > 140 \end{cases}$$

**Proof**

Let's prove that $T(n) < cn$ via induction. The base case is when $n \leq 140$, in which case, $T(n) < cn$ for some arbitrarily large $c$.

Let's now consider the case when $n > 140$. First, we can substitute the base case into our equation. Then we can add an additional $c$ term.

$$T(n) \leq c\lceil n/5 \rceil + c(\frac{7n}{10} + 6) + an$$
$$\leq \frac{cn}{5} + c + \frac{7nc}{10} + 6c + an$$
$$= cn + (-\frac{cn}{10} + 7c + an)$$

It's important to now check for when $\left(-\frac{cn}{10} + 7c + an\right) = 0$.

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

The expression is always true for $n \geq 140$ and $c \geq 20a$.

## Code

```python
1  import math
2  import random
3
4  const = 20
5  sub_size = 5
6  def SELECT(S, i):
7    if len(S) < const:
8      return sorted(S)[i] # assuming rank is zero indexed
9    num_medians = int(math.ceil(len(S) / sub_size))
10   medians = []
11   for curr_index in range(0, len(S), sub_size):
12     sub = S[curr_index : min(curr_index + sub_size, len(S))]
13     mid_index = int(math.floor(len(sub) / 2))
14     medians.append(sorted(sub)[mid_index])
15
16   median = SELECT(medians, int(math.floor(len(medians) / 2)))
17   B = []
18   C = []
19   for element in S:
20     if element < median:
21       B.append(element)
22     elif element > median:
23       C.append(element)
24   if len(B) == i:
25     return median
26   elif len(B) > i:
27     return SELECT(B, i)
28   else:
29     return SELECT(C, i - len(B) - 1)
30
31 arr = range(40)
32 random.shuffle(arr)
```

```
33  print(arr)
34  print(SELECT(arr, 15))
```