

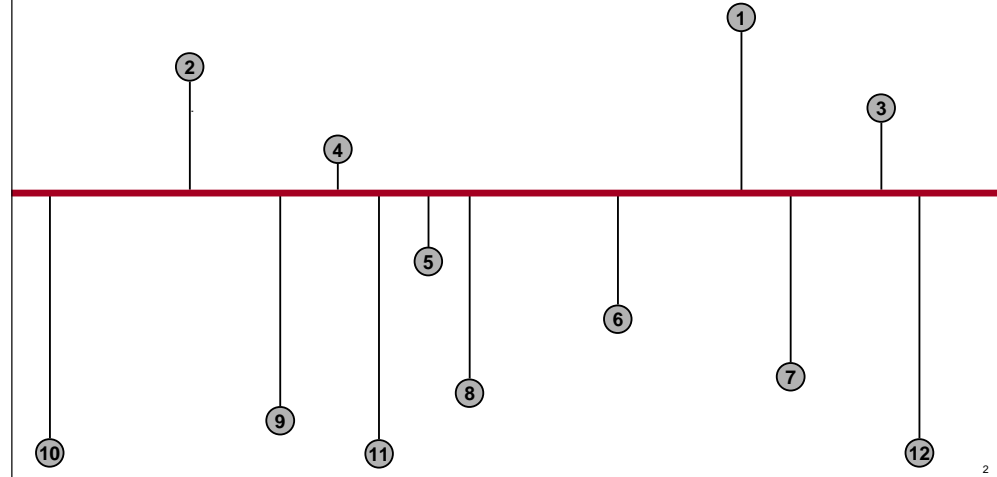
Linear Time Selection



These lecture slides are adapted from CLRS 10.3.

Where to Build a Road?

Given (x,y) coordinates of N houses, where should you build road parallel to x -axis to minimize construction cost of building driveways?

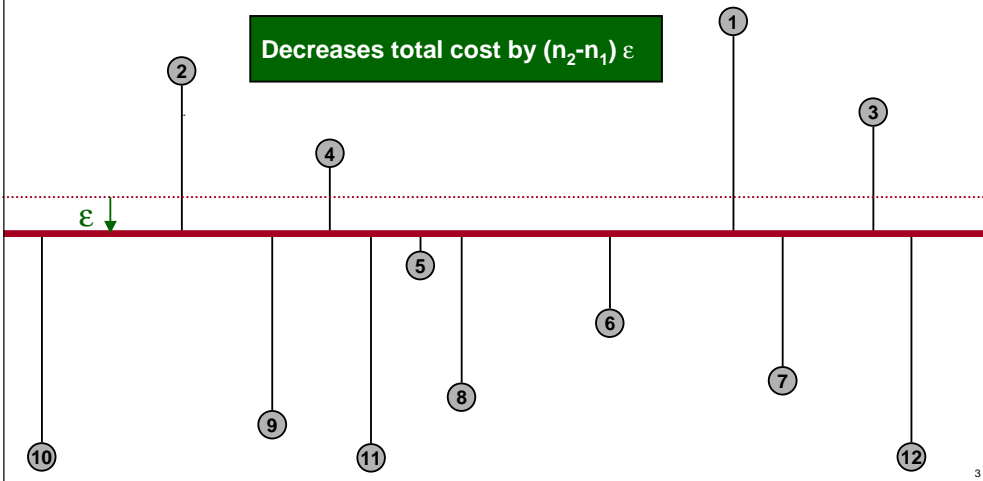


Where to Build a Road?

Given (x,y) coordinates of N houses, where should you build road parallel to x -axis to minimize construction cost of building driveways?

- n_1 = nodes on top.
- n_2 = nodes on bottom.

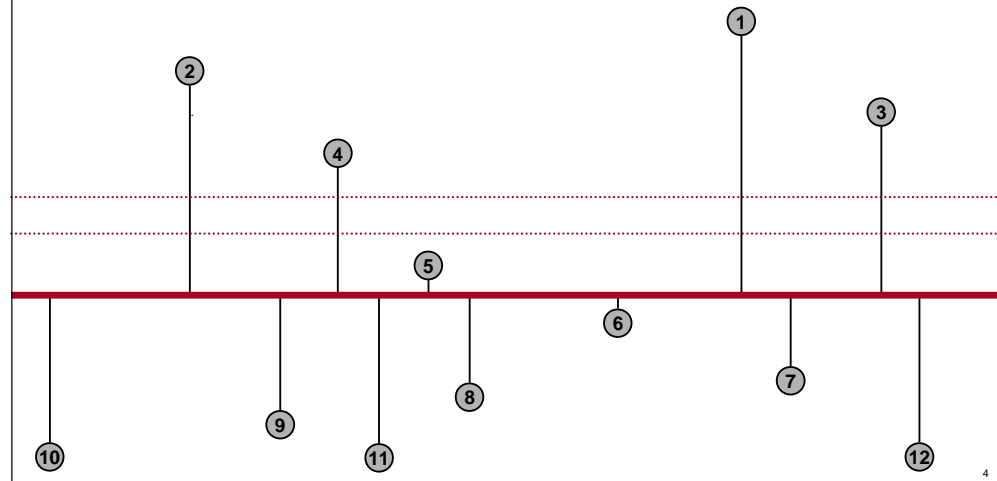
Decreases total cost by $(n_2 - n_1) \epsilon$



Where to Build a Road?

Given (x,y) coordinates of N houses, where should you build road parallel to x -axis to minimize construction cost of building driveways?

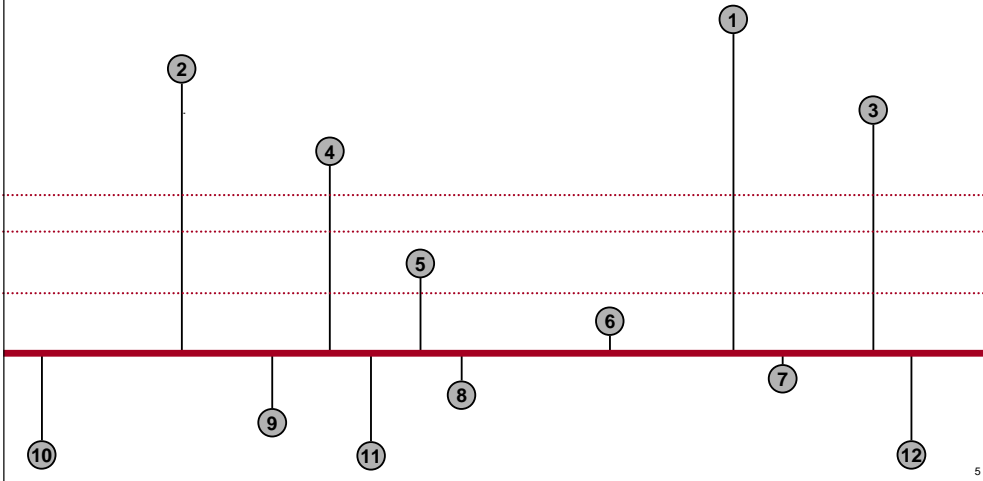
- n_1 = nodes on top.
- n_2 = nodes on bottom.



Where to Build a Road?

Given (x,y) coordinates of N houses, where should you build road parallel to x -axis to minimize construction cost of building driveways?

- n_1 = nodes on top.
- n_2 = nodes on bottom.

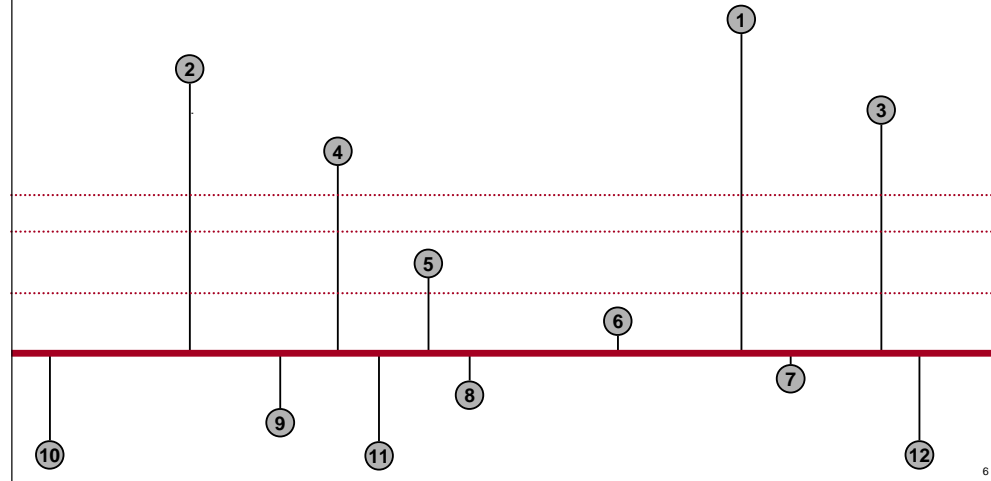


5

Where to Build a Road?

Given (x,y) coordinates of N houses, where should you build road parallel to x -axis to minimize construction cost of building driveways?

Solution: put street at median of y coordinates.



6

Order Statistics

Given N linearly ordered elements, find i^{th} smallest element.

- **Min:** $i = 1$
- **Max:** $i = N$
- **Median:** $i = \lfloor (N+1) / 2 \rfloor$ and $\lceil (N+1) / 2 \rceil$
- $O(N)$ for min or max.
- $O(N \log N)$ comparisons by sorting.
- $O(N \log i)$ with heaps.

Can we do in worst-case $O(N)$ comparisons?

- Surprisingly, yes. (Blum, Floyd, Pratt, Rivest, Tarjan, 1973)

Assumption to make presentation cleaner.

- All items have distinct values.

7

Select

Similar to quicksort, but throw away useless "half" at each iteration.

- Select i^{th} smallest element from a_1, a_2, \dots, a_N .

```

Select ( $i^{\text{th}}, N, a_1, a_2, \dots, a_N$ )
x ← Partition( $N, a_1, a_2, \dots, a_N$ ) ← x = partition element
k ← rank(x)

if (i == k)
    return x
else if (i < k)
    b[] ← all items of a[] less than x
    return Select( $i^{\text{th}}, k-1, b_1, b_2, \dots, b_{k-1}$ )
else if (i > k)
    c[] ← all items of a[] greater than x
    return Select( $(i-k)^{\text{th}}, N-k, c_1, c_2, \dots, c_{N-k}$ )
    
```

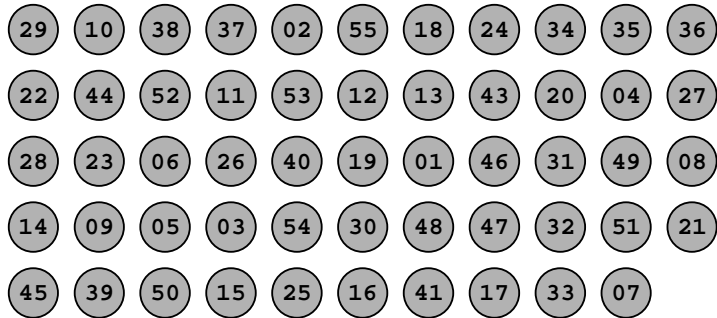
Want to choose x so that x is (roughly) the i^{th} largest.

8

Partition

Partition().

- ➔ Divide N elements into $\lfloor N/5 \rfloor$ groups of 5 elements each, plus extra.



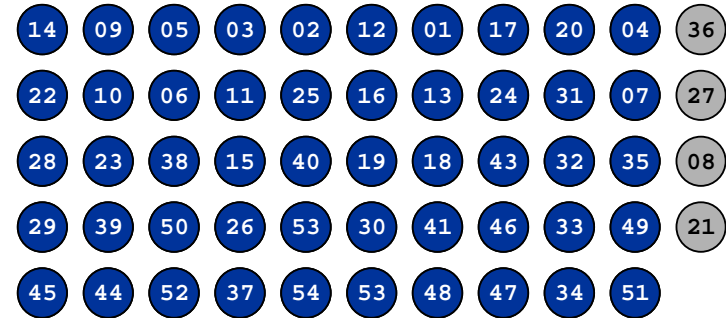
N = 54

9

Partition

Partition().

- Divide N elements into $\lfloor N/5 \rfloor$ groups of 5 elements each, plus extra.
- ➔ Brute force sort each of the 5-element groups.

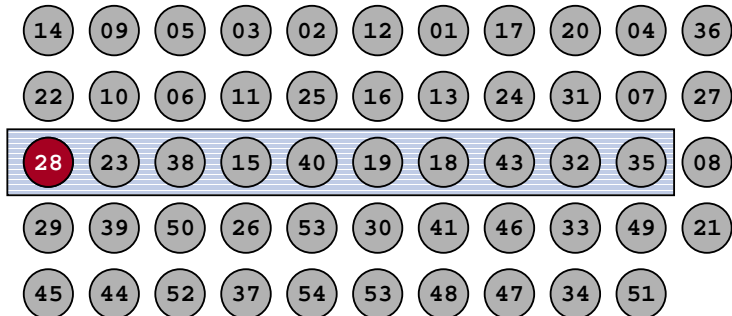


10

Partition

Partition().

- Divide N elements into $\lfloor N/5 \rfloor$ groups of 5 elements each, plus extra.
- Brute force sort each of the 5-element groups.
- ➔ Find x = "median of medians" by select() on $\lfloor N/5 \rfloor$ medians.



11

Select

Select (i^{th} , N, a_1, a_2, \dots, a_N)

if (N is small) use mergesort

Divide $a[]$ into groups of 5, and let $m_1, m_2, \dots, m_{\lfloor N/5 \rfloor}$ be list of medians.

$x \leftarrow \text{Select}(N/10, m_1, m_2, \dots, m_{\lfloor N/5 \rfloor})$ ← median of medians

$k \leftarrow \text{rank}(x)$

if ($i == k$) // Case 1
return x

else if ($i < k$) // Case 2
 $b[] \leftarrow$ all items of $a[]$ less than x
return Select(i^{th} , $k-1$, b_1, b_2, \dots, b_{k-1})

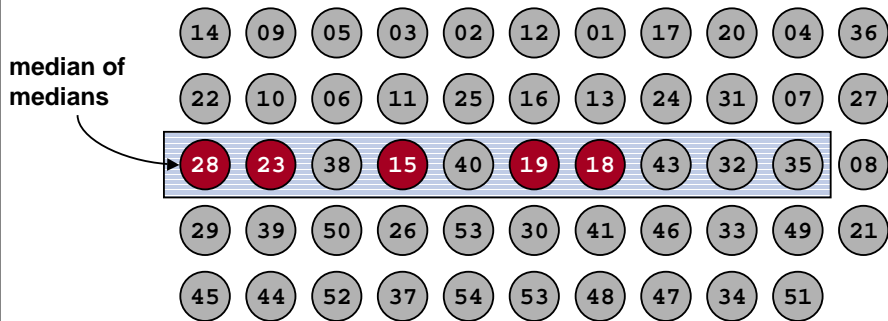
else if ($i > k$) // Case 3
 $c[] \leftarrow$ all items of $a[]$ greater than x
return Select($(i-k)^{\text{th}}$, $N-k$, c_1, c_2, \dots, c_{N-k})

12

Selection Analysis

Crux of proof: delete roughly 30% of elements by partitioning.

- At least 1/2 of 5 element medians $\leq x$
 - at least $\lfloor \lfloor N/5 \rfloor / 2 \rfloor = \lfloor N/10 \rfloor$ medians $\leq x$

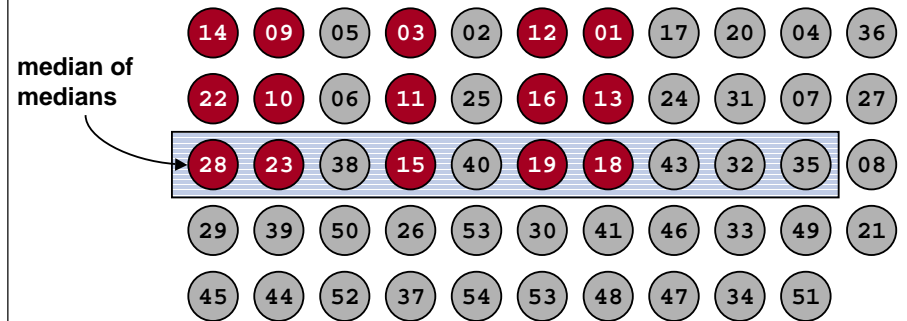


13

Selection Analysis

Crux of proof: delete roughly 30% of elements by partitioning.

- At least 1/2 of 5 element medians $\leq x$
 - at least $\lfloor \lfloor N/5 \rfloor / 2 \rfloor = \lfloor N/10 \rfloor$ medians $\leq x$
- At least $3 \lfloor N/10 \rfloor$ elements $\leq x$.

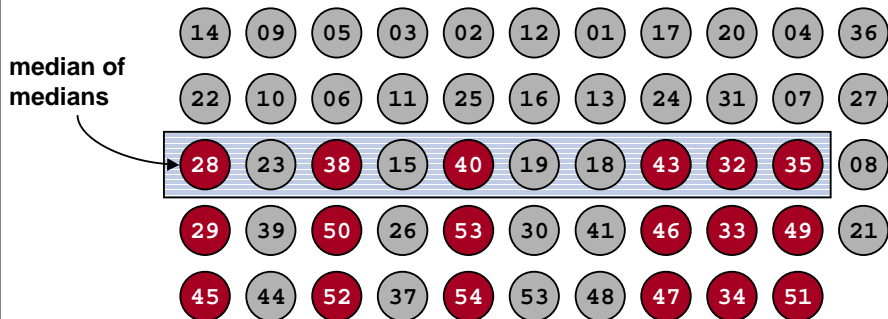


14

Selection Analysis

Crux of proof: delete roughly 30% of elements by partitioning.

- At least 1/2 of 5 element medians $\leq x$
 - at least $\lfloor \lfloor N/5 \rfloor / 2 \rfloor = \lfloor N/10 \rfloor$ medians $\leq x$
- At least $3 \lfloor N/10 \rfloor$ elements $\leq x$.
- At least $3 \lfloor N/10 \rfloor$ elements $\geq x$.



15

Selection Analysis

Crux of proof: delete roughly 30% of elements by partitioning.

- At least 1/2 of 5 element medians $\leq x$
 - at least $\lfloor \lfloor N/5 \rfloor / 2 \rfloor = \lfloor N/10 \rfloor$ medians $\leq x$
- At least $3 \lfloor N/10 \rfloor$ elements $\leq x$.
- At least $3 \lfloor N/10 \rfloor$ elements $\geq x$.
 - \Rightarrow **select()** called recursively (Case 2 or 3) with at most $N - 3 \lfloor N/10 \rfloor$ elements.

$C(N)$ = # comparisons on a file of size N .

$$C(N) \leq \underbrace{C(\lfloor N/5 \rfloor)}_{\text{median of medians}} + \underbrace{C_{\leq}(N - 3 \lfloor N/10 \rfloor)}_{\text{recursive select}} + \underbrace{O(N)}_{\text{insertion sort}}$$

Now, solve recurrence.

- Apply master theorem?
- Assume N is a power of 2?
- Assume $C(N)$ is monotone non-decreasing?

16

Selection Analysis

Analysis of selection recurrence.

- $T(N)$ = # comparisons on a file of size $\leq N$.
- $T(N)$ is monotone, but $C(N)$ is not!

$$T(N) \leq \begin{cases} 20cN & \text{if } N < 50 \\ T(\lfloor N/5 \rfloor) + T(N - 3\lfloor N/10 \rfloor) + cN & \text{otherwise} \end{cases}$$

Claim: $T(N) \leq 20cN$.

- Base case: $N < 50$.
- Inductive hypothesis: assume true for $1, 2, \dots, N-1$.
- Induction step: for $N \geq 50$, we have:

$$\begin{aligned} T(N) &\leq T(\lfloor N/5 \rfloor) + T(N - 3\lfloor N/10 \rfloor) + cN \\ &\leq 20c\lfloor N/5 \rfloor + 20c(N - 3\lfloor N/10 \rfloor) + cN \\ &\leq 20c(N/5) + 20c(N) - 20c(N/4) + cN \\ &= 20cN \end{aligned}$$

$$\text{For } n \geq 50, \\ 3\lfloor N/10 \rfloor \geq N/4.$$

17

Linear Time Selection Postmortem

Practical considerations.

- Constant (currently) too large to be useful.
- Practical variant: choose random partition element.
 - $O(N)$ expected running time ala quicksort.
- Open problem: guaranteed $O(N)$ with better constant.

Quicksort.

- Worst case $O(N \log N)$ if always partition on median.
- Justifies practical variants: median-of-3, median-of-5.

18