

Given a set S of n numbers,
we want to find an element in S
that has rank = k ($1 \leq k \leq n$)

Selection

5 20 9 15 8

$n = 5$

$k = 2$

5 8 9 15 20
 $k=2$

Do we need to sort to find the rank k
element?

Can we do selection in $O(n \log n)$
comparisons

Median finding : $k = \lfloor \frac{n}{2} \rfloor$ or $\lfloor \frac{n}{2} \rfloor$

Considered important for even splitting for
divide-and-conquer algorithms like sorting

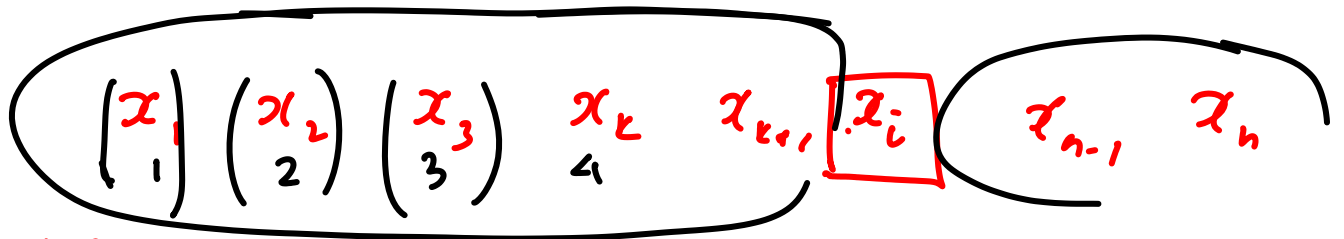
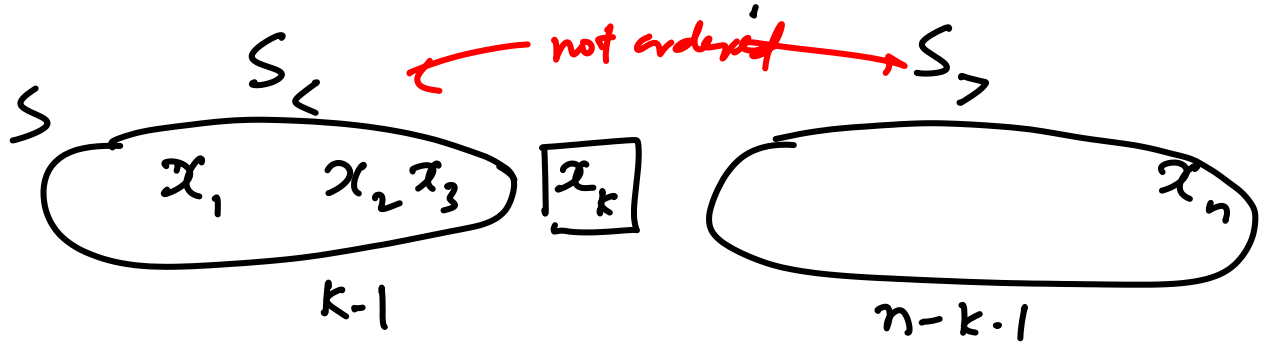
Goal : Design an $O(n)$ time for selection

SpL case $k = 1$ / $k = n$

$k = 2$. $2n$? (two scans)

$\Rightarrow O(kn)$ for k scans

Why can't we find minimum faster than $n-1$ comparisons?



if $i > k$ - then $x_k < x_i$

$i < k$ - then $x_k > x_i$

$S_{< x_i}$ or $S_{> x_i}$ can be generated by comparing all elements with x_i

At least one of the subsets $S_{< x_i}$ or $S_{> x_i}$ can be discarded

[We repeat (recursively) the same in the remaining subset with $k = k / k - i$ until $k = i$

Suppose $i = \frac{n}{2}$

$$n + \frac{n}{2} + \frac{n}{4} + \dots < 2n$$

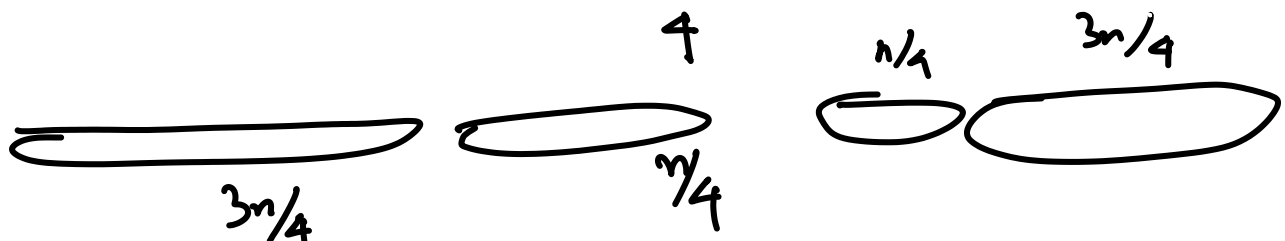
So find k^{th} ranked element can be reduced to finding median efficiently



Very unbalanced partitioning will lead to $\Omega(n-1 + n-2 + \dots) = \Omega(n^2)$

The rank of the partitioning element x_i is key to the success of the approach

Suppose $i = \frac{3n}{4}$



For any rank αn , $0 < \alpha < 1$, we will be doing $n + \alpha n + \alpha^2 n + \dots = O(n)$

As long as the rank of the partitioning element is "approx in the middle"

$$\alpha = \frac{1}{4} \quad \left[\frac{n}{4}, \frac{3n}{4} \right]$$

$x_1, x_2, x_3, \left[\underbrace{x_i}_{x_{n/4} \text{ to } x_{3n/4}} \right], \dots, x_n$
 approximate median

Choose x_i uniformly at random from S

\Rightarrow Rank of the partitioning element is uniformly distributed in $[1, 2, \dots, n]$

Say \nearrow rank of x_i $\rightarrow Y$ is the random variable with a value in $[1, 2, \dots, n]$

$$Pr[Y = j] = \frac{1}{n} \quad 1 \leq j \leq n$$

What is the prob that $\frac{n}{4} \leq Y \leq \frac{3n}{4}$
 $= \frac{1}{2}$

If Y is not in this range we choose another splitter randomly till we succeed

What is the expected no. of times we repeat the above

Let Z be a random variable that has value k with prob p_k

Then Expectation of $Z = E[Z] = \sum k \cdot p_k$

$Z = \# \text{ of trials}$

$Z = 1 \quad 2 \quad 3 \quad 4 \quad \dots$

Prob $\frac{1}{2} \quad \frac{1}{4} \quad \frac{1}{8}$

$$E[Z] = 2 = \sum i \cdot \frac{1}{2^i}$$

For a geometric distribution with parameter p
the Expectation is $\frac{1}{p}$

For ^{any} n v.v. $Z_1, Z_2, Z_3, \dots, Z_t$ (not necessarily independent)

$$E[Z_1 + Z_2 + \dots + Z_t] = \sum_{i=1}^t E[Z_i]$$

Linearity of expectation

What is the expected running time of
the selection algorithm based on
random choice of splitter

Let X_1 be the random variable for the
comparisons in 1st round

X_i : v.v. for # comparisons in i^{th} round

$$\text{Total cost} = X_1 + X_2 + \dots + X_t$$

$$\begin{aligned}
 E[\text{Total cost}] &= E[X_1 + X_2 + \dots + X_i] \\
 &= \sum_{i=1}^P E[X_i]
 \end{aligned}$$

$t = \log_{4/3} n$
 $\therefore O(\log n)$

$$E[X_1] \leq 2n$$

$$E[X_2] \leq 2 \times \frac{3}{4}n$$

$$E[X_i] \leq 2 \times \left(\frac{3}{4}\right)^{i-1} n$$

$$2n \sum \left(\frac{3}{4}\right)^i \leq O(n)$$

What does it mean by Expected
 running time $\leq f(n)$

Markov inequality

$$Pr[X \geq c E[X]] \leq \frac{1}{c}$$

for non-neg rand variables