

def median(L):

L.sort()

return L[len(L)//2]

) can we find medians
more efficiently?

Why?

- Find good pivots in divide + conquer algorithms such as quicksort (either median of whole input or of sample)
- Statistics - find central/representative point for cloud of (1-d) data points



if points are true value + normally distributed noise (Bell curve) - use average

if a small number of points could be very far from main cloud of points (outliers) - median more robust will still be accurate as long as #outliers < 50%
"breakdown point" of median = 0.5

- Data whose ordering is meaningful
but whose numerical value is not

Movie ratings by # of stars

Meaningful: you gave one movie more stars than
another

Not meaningful: this movie is one star better
than this other one

(diff. $4\star - 3\star \neq 2\star - 1\star$)

Not meaningful: average of star ratings
Meaningful: median

- Facility location

- Where along this road should I place my store to minimize average customer driving time?

(ignore price, zoning, ...)

- given n points on a line x_i
median = point y
that minimizes $\frac{1}{n} \sum_{i=0}^{n-1} |x_i - y|$

Proof: $\frac{d}{dx} \left(\frac{1}{n} \sum |x_i - y| \right)$ ~~$\frac{d}{dy}$~~

\uparrow each x_i contributes $\pm \frac{1}{n}$

$$= \left(\text{fraction of points } x_i > y \right) - \left(\text{fraction of points } x_i < y \right)$$

$O(n)$ - time median computation

more general problem: select k th smallest

($k=0 \Rightarrow \min$ $k=n-1 \Rightarrow \max$ $k=n/2 \Rightarrow \text{median}$)

Deterministic - complicated, tricky analysis, impractical

Randomized, based on quicksort -
simple and practical

Basic idea -

interleave quicksort
recursion with
selection problem
instead of doing
whole quicksort
algorithm first -

faster \leftarrow [Skipping some
irrelevant recursive
calls]

def quickselect(L, k)
if $\text{len}(L) \leq 1$: return $L[0]$
choose random pivot p
partition L into
 X : elements $< p$
 p itself
 Y : elements $> p$
if $k < \text{len}(X)$:
 return quickselect(X, k)
else if $k == \text{len}(X)$:
 return p
return quickselect($Y, k - \text{len}(X) - 1$)

adjust k to account for removing X, p

after checking
that $k == 0$

Analysis of quickselect

Option 1: like quicksort,

$$E[\text{\#compares}] = \sum_{\text{pair } x, y} \text{Pr}[x, y \text{ are compared}]$$

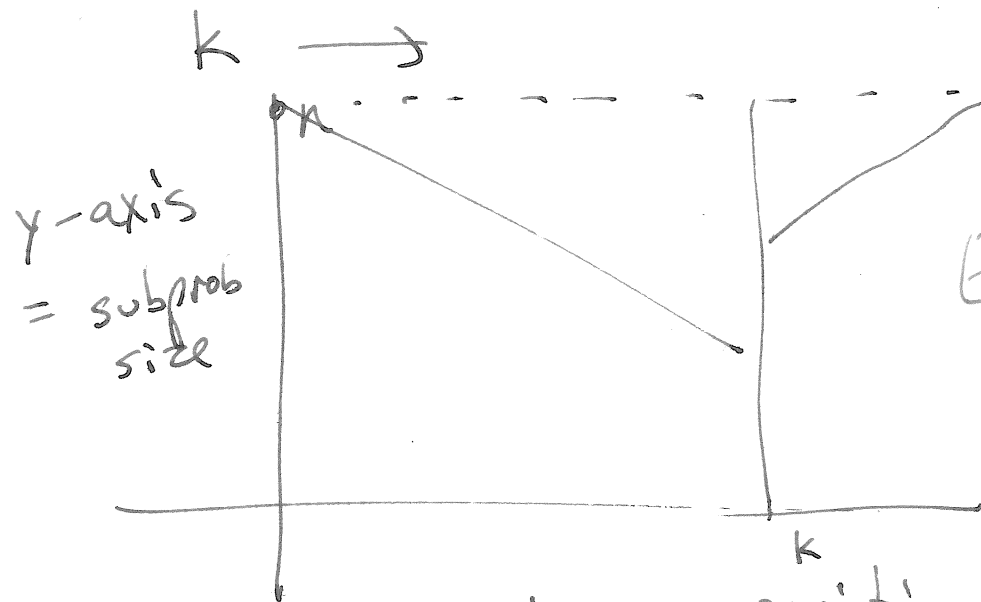
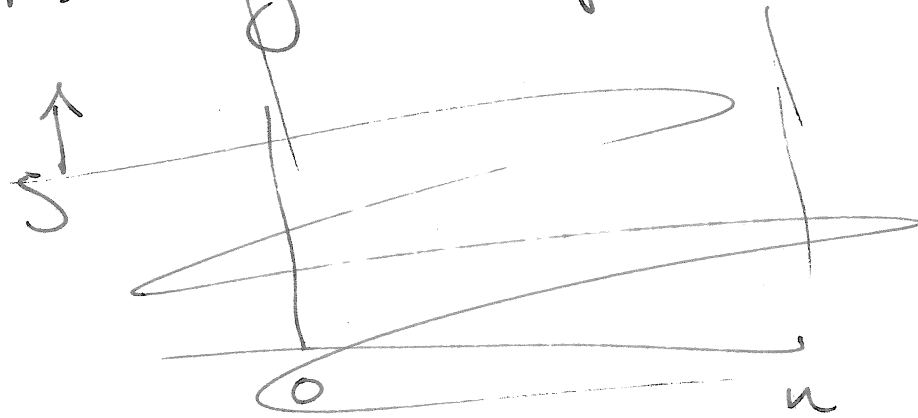
Option 2: recurrence

$$E[\text{Time}(n)] = \underset{\text{at top level} \rightarrow}{O(n)} + \sum_{\substack{\text{subproblem} \\ \text{size } s}} \text{Pr}(s) \times \text{Time}(s)$$

$$\begin{array}{l} \text{"="} \\ \text{(not really)} \end{array} \quad O(n) + T(\text{average subprob. size})$$

swapping \sum and $T()$
only works if T already
known to be linear

Average subproblem size?



Claim:
for worst case
choice of k ,
 $E[\text{subprob size}]$
 $\leq \frac{3n}{4}$

x-axis = position of pivot
in sorted list

$$T(n) = O(n) + \frac{3}{4}T(n)$$

$$= O(n)$$

plug into induction proof

$$T(n') \quad n' \leq n$$

linear by induct. hyp.

\Rightarrow justifies Σ/T swap