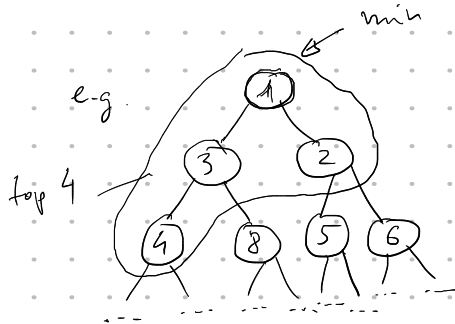


Soft heap applications

Structured selection

Historical note: optimal algorithms were given by [Fredrickson] [Fredrickson-Johnson] these were very complicated, Soft heap makes it easier!

binary heap



size: n (possibly infinite)

keys n : min-heap order

task: find k^{th} smallest key

(or find k smallest keys)

($k \leq n$)

Note: there can be missing children in heap, just replace them by ∞

Problem: selection in a heap.

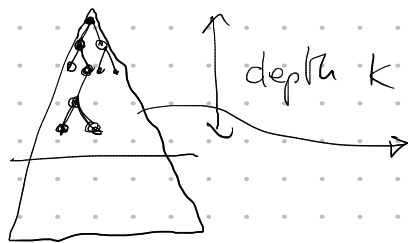
Very naive solution: extract-min k times

time:

$$O(k \cdot \log n)$$

naive solution:

only need to look at top k levels



$$\text{size } 2^k - 1$$

time:

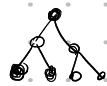
$$O(k^2)$$

Wait $O(k)$

General selection

x_1, \dots, x_n , no info about ordering (selection $\Theta(n)$ time)

our problem: Selection in heaps



(partial order)

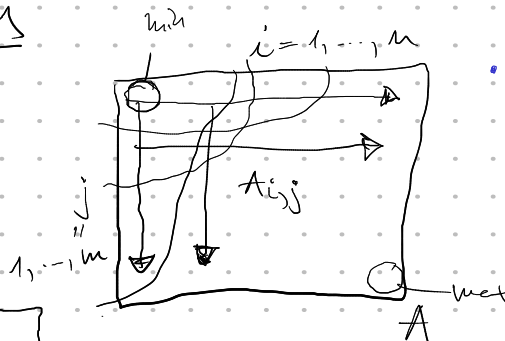
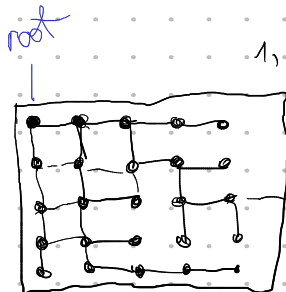
"semi sorted" how to select efficiently?

Sorted input

$x_1 < x_2 < \dots < x_n$ (selection trivial)

Example application

(1) Sorted matrix



every row, every column is sorted

Which $A_{i,j}$ is the k^{th} smallest

$O(k)$ instead of $O(mn)$

view as heap

(some children missing, we can decide which edges to add)

② Selection from $X + Y$

Given $X = \{x_1, \dots, x_n\}$
 $Y = \{y_1, \dots, y_m\}$

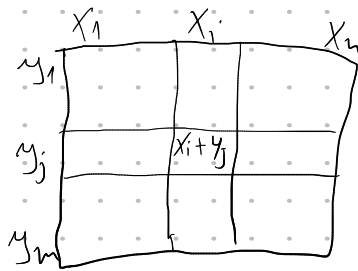
$A_{i,j} = x_i + y_j \quad \forall i, j$
 Task: find k^{th} smallest sum in $X+Y$

$X+Y = \{A_{i,j} \mid i \in [n], j \in [m]\}$ $n \cdot m$ sums

If X, Y sorted

$x_1 < x_2 < \dots < x_n$
 $y_1 < y_2 < \dots < y_m$

use approach ①
 $O(k)$ time



If X, Y not sorted, first sort in time $O(n \log n + m \log m)$
 → sorted Matrix

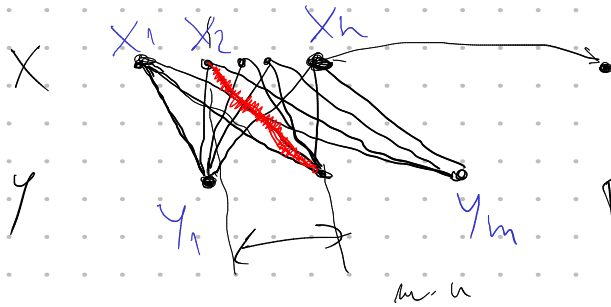
(can solve more efficiently even if X, Y are not sorted.)

see exercise

Example application.

$X = \{x_1, \dots, x_n\}$ we moments.
 $Y = \{y_1, \dots, y_m\}$

Compare Result = single number



Hodges-Lehmann estimator

pick line with median slope

pick median from $Y - X$
 $(+X) + Y$

"Robust estimator"

③ Given $x_1 < x_2 < \dots < x_n$

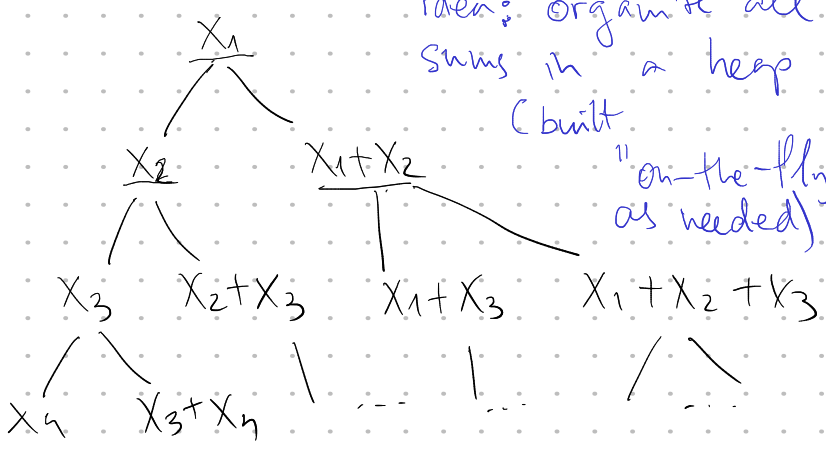
Task: Find k^{th} smallest sum on a subset of $\{x_1, \dots, x_n\}$

$\emptyset, x_1, x_2, x_3, \dots, x_1+x_2, \dots, x_1+x_3+x_5+x_6, \dots, x_99+x_{100}, \dots$
 1st 2nd 3rd 4th

possible subsets 2^n Wait time $O(k)$

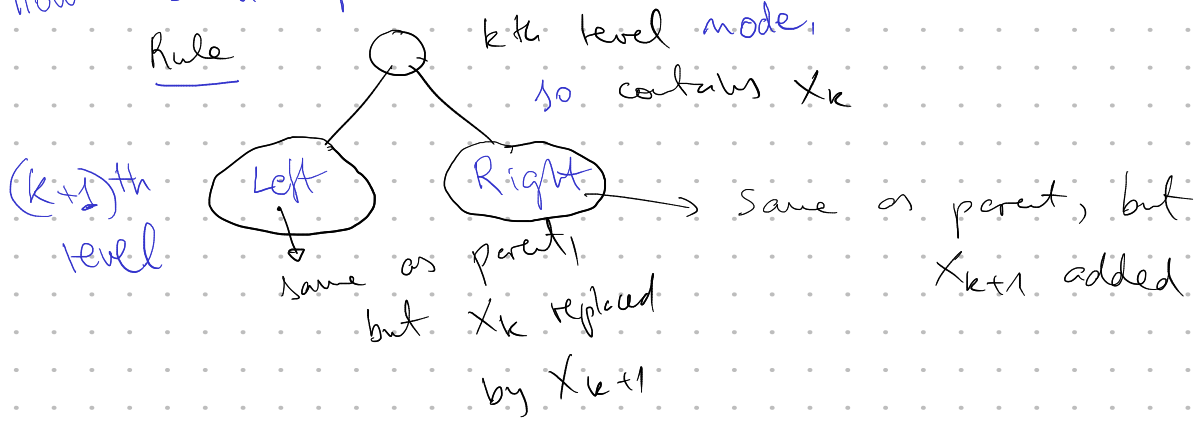
Idea: organize all sums in a heap (built "on-the-fly" as needed)

- top k levels: all sums with elements in $\{x_1, \dots, x_k\}$
- in k th level: all sums containing x_k with elements in $\{x_1, \dots, x_k\}$



How to build heap?

Rule



obs all possible sums appear in heap.

Induction

x_k : largest element of sum

$$S = x_{i_1} + x_{i_2} + \dots + x_{i_j} + \underline{x_k}$$

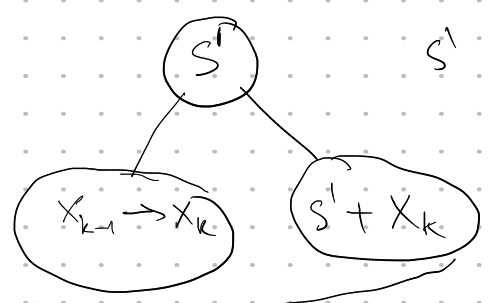
$$i_1 < i_2 < \dots < i_j < k$$

(consider an arbitrary sum S over a subset of $\{x_1, \dots, x_n\}$ and show that it appears in the heap)

Claim: S appears in heap at level k .

$$S' = x_{i_1} + x_{i_2} + \dots + x_{i_j} + x_{k-1} \quad \text{if } i_j < k-1$$

$$= x_{i_1} + x_{i_2} + \dots + x_{i_j} \quad \text{if } i_j = k-1$$



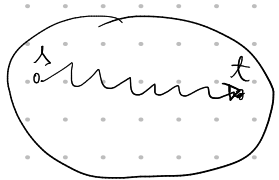
S' in the heap at level $k-1$ (inductive claim)

one of these exactly \underline{S}

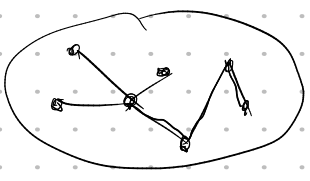
Selecting k th sum \rightarrow using selection in a heap, will take $O(k)$ time

4) Selection in heaps can be used for structured selective problems, e.g. in graphs.

- shortest path from s to t .

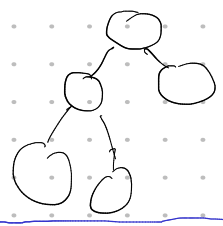


- MST



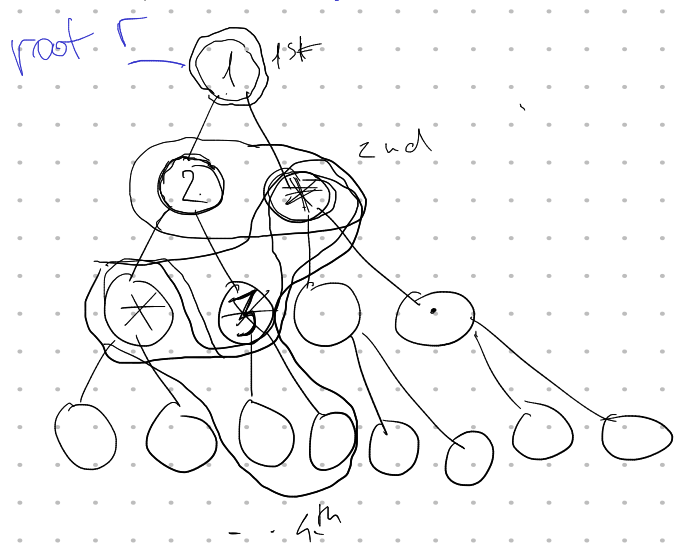
- k^{th} shortest path $s \rightarrow t$
- top k shortest paths
- top k MST

- can be reduced to selection in heaps



Algorithms for Selection from heaps

First approach: $O(k \log k)$ time algorithm



Idea: maintain a set of candidate k^{th} elements.

When pop minimum x , add to set two children of x .

Use a heap data structure to maintain candidate set.

```

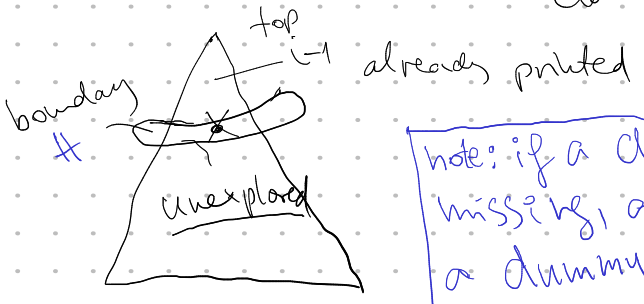
H := make-heap()
insert(r, H)
for i = 1..k
  x = extract-min(H)
  print(x)
  insert(x.left, H)
  insert(x.right, H)
  
```

** e.g. skew heap or binary heap*

$O(\log k)$

Correctness

by induction
 i^{th} print outputs the smallest element

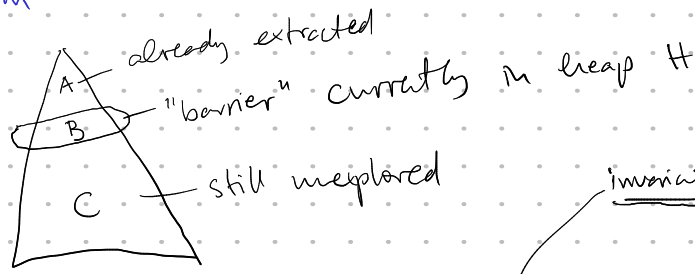


note: if a child missing, add a dummy node

running time

- in each iteration heap size increases by 1.
 - At most $k+1$ elements in heap
- $O(k \log k)$

Review of $O(k \log k)$ time algorithm



invariant: all A is smaller than all in B, C

$|A| = k$, done
return A

induction

• true when $A = \emptyset$

• $\text{extract_min}(H) \rightarrow X$

X: key min in B U C

(every one in C has an ancestor in B)

moving X to A maintains invariant

Observation:

Algorithm returns top k

elements in increasing order



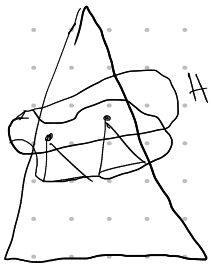
we sorted top k elements



as long as we sort, we cannot do better than $O(k \log k)$

Second approach = $O(k)$ time algorithm

Idea: Same as previous algorithm, but use soft-heap for maintaining H.



$\log_{1/2} \in O(1)$, so running time $O(k)$

modification: add children of X to H, when X is extracted or when X is corrupted.

What we expect from soft heap?

make-heap $\rightarrow \epsilon$
 $O(1)$

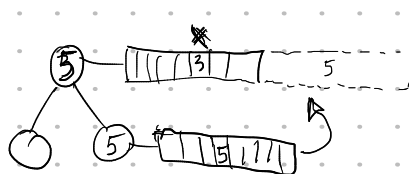
insert } $O(\log^{1/2} \epsilon)$
extract-min }

no corruption

return (X, C)

list of corrupted items

\rightarrow want soft heap to tell us whenever it corrupts some element.



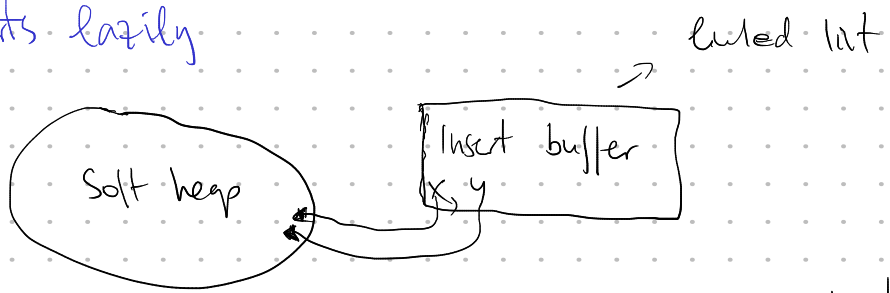
recall: corruption happens during soft-up

for each list move, one item corrupted, charge output cost to soft-up cost.

→ Modify soft heap design, so that corruption happens only during extract-min

→ insert should not cause corruption

idea: do inserts lazily



Amortized analysis
 β unchanged

insert(x)
 insert(y) → items put into insert buffer
 no actual restructuring

extract-min
 everything from buffer inserted,
 then do extract-min (insert buffer emptied)

$H \leftarrow \text{make_heap}(\frac{1}{4})$ ----- \rightarrow soft-heap w. error-param $\epsilon = \frac{1}{4}$
 $L \leftarrow \text{make_list}()$
 insert(r, H)

for $i = \overline{1, k}$

$(x, k) \leftarrow \text{extract_min}(H)$

insert(x, L)

insert($x.\text{left}, H$)

insert($x.\text{right}, H$)

for all $c \in K$ *(newly corrupted items)*

insert($c.\text{left}, H$)

insert($c.\text{right}, H$)

// top k elements are in $(L \cup H)$ *size $O(k)$*

Select top k from $L \cup H$

from gk items

Proof of obs.

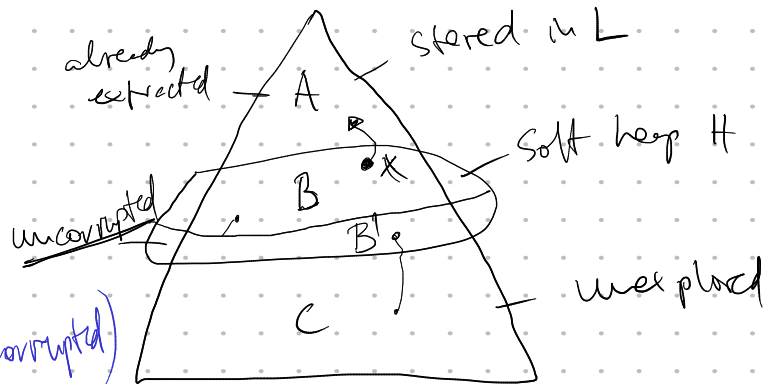
x has smallest current key in B

$x.\text{key} \leq x.\text{current}$

$x.\text{key}$ smaller than all current keys in B

$x.\text{key} \leq y.\text{key}$ for all $y \in B'$

every element in C has an ancestor in B'
 every key in $C > x.\text{key}$



Obs. all in A have smaller key than all in C .

if $|A| = k$

\Rightarrow top k elements cannot be in C

\Rightarrow top k can only be in

$A \cup B$
 $(L \cup H)$

finish the job with a standard selection step

uncomputed

Show that $|A \cup B| \in O(k)$

Let I denote # insertions into H

Let K denote # corruptions

(# corrupted items in heap)
 $\leq \epsilon \cdot I$

$$I \leq 2k + 2K$$

$$K \leq \epsilon \cdot I + k$$

$|B \setminus B'|$ everything extracted already

Observe that

$$I = |A \cup B|$$

(everything in the heap and everything extracted must have been inserted at some point)

(Solve for I)

$$K \leq \epsilon \cdot (2k + 2K) + k$$

$$\Rightarrow k(1 - 2\epsilon) \leq k(1 + 2\epsilon)$$

$$\Rightarrow k \leq k \frac{1 + 2\epsilon}{1 - 2\epsilon}$$

$$\Rightarrow I \leq \frac{2k}{1 - 2\epsilon} + \frac{2k}{1 - 2\epsilon} \frac{1 + 2\epsilon}{1 - 2\epsilon} \in O(k)$$

$$\epsilon = \frac{1}{4}$$

$$\frac{\frac{3}{2}}{\frac{1}{2}} = 3$$

$$I \leq 8k$$

$$\Rightarrow |A \cup B| \leq 8k$$

last step:

Select top k from $\leq 8k$ items (takes $O(k)$ time.)

Overall $O(k)$ time

to select top- k in heap.

Observation

We assumed that heap is binary. We could handle larger constant degree too.

- First approach extends easily to d -ary heap.
- Second approach also extends, but calculation changes

now $I \leq \underline{d} \cdot k + \underline{d} \cdot k$

(we may need to set ϵ smaller, e.g. $\epsilon = \frac{1}{2d}$)