➤



$X = k_1, k_2 ... k_n$

$\text{Rank}(q, s) = j$

Let $\text{Rank}(q, X) = r_j$

$$Prob\left(\left|r_j - j\frac{n}{s}\right| > \sqrt{3\alpha}\frac{n}{\sqrt{s}}\sqrt{\log n}\right) \le n^{-\alpha}$$

➤ To begin with, all the keys are alive, N=n;

Repeat

1. In one pass through the data pick a random sample S

    Each alive key will be in S with a prob. of $\frac{M}{2N}$

    $$E[|S|] = \frac{M}{2}$$

2. Pick $l_1$ & $l_2$ from S such that,

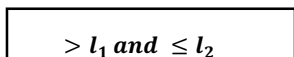    $\text{Rank}(l_1, S) = i \cdot \frac{s}{N} - \sqrt{4\alpha s \cdot \log N}$

    and $\text{Rank}(l_2, S) = i \cdot \frac{s}{N} + \sqrt{4\alpha s \cdot \log N}$

3. Do one more pass through the data and store Y = {q: q is alive and

    $l_1 < q \le l_2$} in the disk;

| | alive keys |
|---|---|

| $> l_1 \text{ and } \le l_2$ | Y |
|---|---|

let $n_1 = |\{q : q \text{ is alive and } q < l_1\}|;$

$n_2 = |\{q : q \text{ is alive and } q \leq l_2\}|;$

4. If $i < n_1$ or $i > n_2$

   or if $|Y| > \dfrac{N}{M^{0.4}}$ then start all over;

   else $i = i - n_1$ and $N = |Y|;$

   only the elements in Y are alive;

Until $N \leq M;$

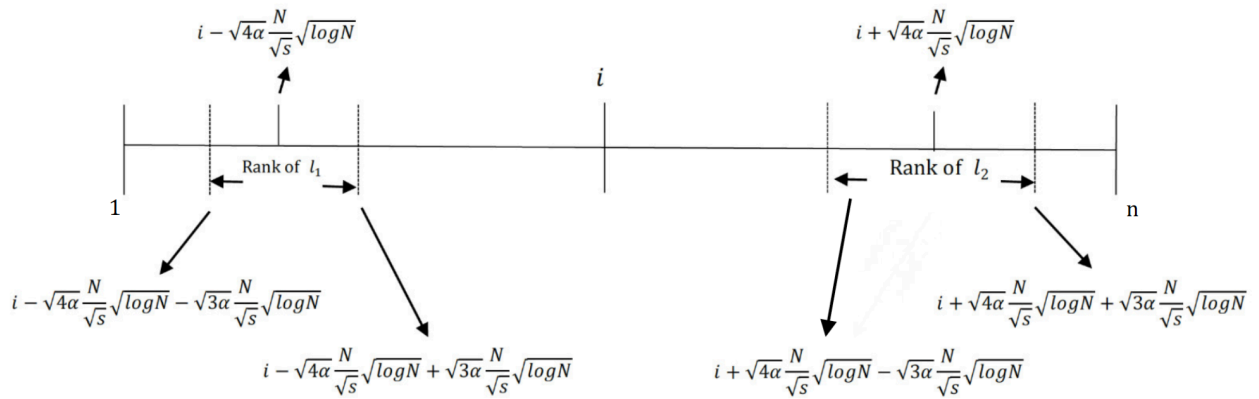Identify and output the $i^{th}$ smallest from the alive keys.

➢ *Analysis:*

In step 1, the # of sample keys in S is Binomial with parameters N and $\dfrac{M}{2N}$

Using Chernoff bounds, $|s| = \tilde{\theta}(M)$

Expected Rank of $l_1$ in $X = [i \cdot \dfrac{s}{N} - \sqrt{4\alpha \cdot s \cdot \log N}] \dfrac{N}{s}$

$$= i - \sqrt{4\alpha} \dfrac{N}{\sqrt{s}} \sqrt{\log N}$$

Expected Rank of $l_2$ in $X = i + \sqrt{4\alpha} \dfrac{N}{\sqrt{s}} \sqrt{\log N}$

$$|Y| \leq i + \left(\sqrt{4\alpha} + \sqrt{3\alpha}\right)\frac{N}{\sqrt{s}}\sqrt{\log N} - \left[i - \left(\sqrt{4\alpha} + \sqrt{3\alpha}\right)\frac{N}{\sqrt{s}}\sqrt{\log N}\right]$$

$$\leq 2\left(\sqrt{4\alpha} + \sqrt{3\alpha}\right)\frac{N}{\sqrt{s}}\sqrt{\log N} \quad \text{with a prob. of } \geq (1 - N^{-\alpha})$$

➤ _**w.h.p. (with high probability)**_

$$|Y| = O\left[\left(\sqrt{4\alpha} + \sqrt{3\alpha}\right)\frac{N}{\sqrt{M}}\sqrt{\log N}\right] \qquad \text{note: } N = M^c \rightarrow \log N = c \cdot \log M$$

If N is a polynomial in M,

then $|Y| = O\left(\frac{N}{M^{0.4}}\right)$  w.h.p.

➤ _**Example**_

$M = 10^9, N \leq M^4$

I/O complexity:

$$\frac{2n}{B} + \frac{n}{M^{0.4}} \cdot \frac{2}{B} + \frac{2n}{M^{0.8}B} + \cdots \leq (2 + \varepsilon)\frac{n}{B}, \ \ for \ any \ constant \ \varepsilon > 0$$

➤ _**A Graph Problem:**_

Minimum spanning tree(MST)

Problem:

Input:     a weighted, connected and undirected graph, G(V,E)

Output:    A minimum spanning tree for G

➤ _**Prim's algorithm**_

Grow a subtree by adding one edge at a time, starting with the lightest edge.
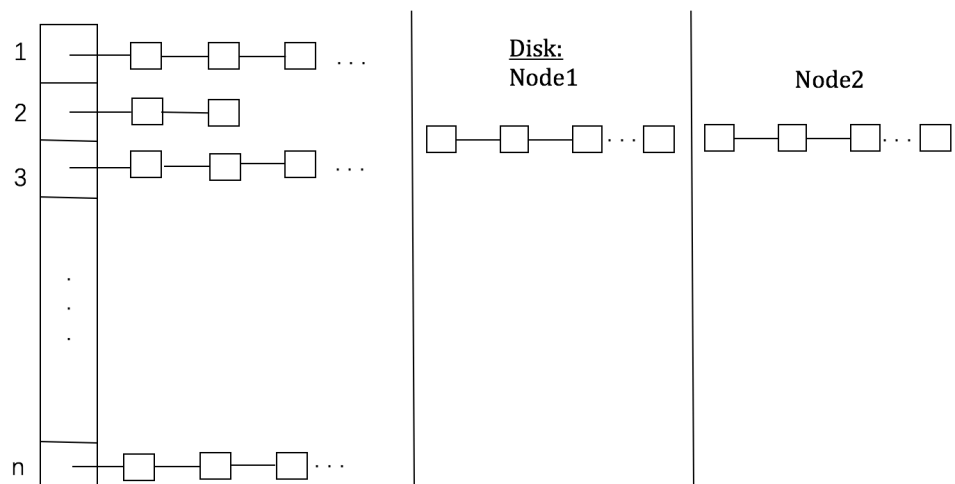
Let x be any node outside the tree

NEAR[x] = the closest tree neighbor of x;

       1. pick the lightest edge e=(a,b) in E;

       2. for every u ∈ V−{a,b} do

           if weight(u,a) < weight(u,b) then

           NEAR[u] = a; else NEAR[u]=b;

       3. for u ∈ V−{a,b} do

insert u into a 2-3 tree Q where the key value is weight(u, NEAR[u]);

4. for i=1 to (n-2) do

find the node u in Q with the least key;

insert (u, NEAR[u]) into T;

for every w ∈ Adj(u) do                    *Adj(u)  →  Adjacent to u

if weight(w, NEAR[w]) > weight(w,u) then

NEAR[w] = u;

➢ Assume that the input is in adjacency lists form



➢ Assume that  $M = \theta(n)$, n = |V|. This means that Q can be stored in core memory;

## Analysis:

- Step1: takes  $\dfrac{|E|}{B}$  I/O operations

- Step2: takes  $\leq \dfrac{2|V|}{B}$  I/O operations

- Step3: No I/O operations; Let  $d_u$ be the degree of $u$ for any $u \in V$.

- Step4: takes  $\leq \sum_{u \in V} \left\lceil \dfrac{d_u}{B} \right\rceil \leq \sum_{u \in V} \left( \dfrac{d_u}{B} + 1 \right) = O\left( \dfrac{|E|}{B} + |V| \right)$  I/O operations.

Note that this algorithm is optimal when  $|E| \geq |V|B$.