CSE 4502/5717
BIG DATA ANALYTICS

LECTURE ON 9-27-22

# RANDOMIZED SELECTION: FLOYD & RIVEST (1977)

① PICK A RANDOM SAMPLE $S$
   with $|S| = 8$.

② PICK $l_1$ and $l_2$ FROM $S$ s.t.:

   ⓐ The $i^{th}$ smallest element $8$
      $x \in [l_1, l_2]$.

   ⓑ $\left| \left\{ l \in X : l_1 \leq l \leq l_2 \right\} \right|$ is "SMALL"

③ SCAN through X and keep only

$$Y = \{ z \in X : \ell_1 \le z \le \ell_2 \}.$$

Make sure that Conditions 2a and 2b are MET.

Let $X_1 = \{ z \in X : z < \ell_1 \}$.

Let $n_1 = |X_1|$. Let $n_2 = |Y|$.

If $c > n_1$ and $c \le (n_1 + n_2)$ then 2a IS MET.

Let Rank $(q, S) = \dot{j}$, FOR some element $q \in S$. Rank $(q, X) = \gamma_{\dot{j}}$

$$E\left[\gamma_{\dot{j}}\right] = \dot{j}\,\frac{n}{s}.$$

LEMMA:

$$\text{Prob}\left[\left|\gamma_{\dot{j}} - \dot{j}\,\frac{n}{s}\right| > \sqrt{3\alpha}\,\frac{n}{\sqrt{s}}\sqrt{\log n}\right] \le n^{-\alpha}.$$

Let $\boxed{d_r = \sqrt{32} \, \frac{n}{\sqrt{8}} \sqrt{\log n}}$ ———①

$$E[\hat{r}_j] = \hat{j}\frac{n}{8}.$$

The Lemma Says that $r_j$ lies in this interval with a high prob.

Interval endpoints: $\hat{j}\frac{n}{8} - d_r$ and $\hat{j}\frac{n}{8} + d_r$

OUT_OF_CORE ALGORITHMS TO BEGIN
$\qquad$ with $N = n$.

REPEAT UNTIL

① Do ONE SCAN through the input
& keep every d. in S with a
Prob. $\dfrac{M}{2N}$.

② Let $l_1 \in S$ be Such that

$$Rank(d_1, S) = \boxed{i \dfrac{8}{n} - \sqrt{4\alpha \log n}}$$

Let $l_2 \in S$ be Such that

$$Rank(l_2, S) = \boxed{i \dfrac{8}{n} + \sqrt{4\alpha \log n}}$$

③ SCAN through the input to identify $Y$ & write it in the Disk.

———————————————|X
|————————————————|X

|——————————|Y

④ If any of the two conditions $2a$ and $2b$ is not met, Start all OVER; $i = i - n_1$;

⑤ $N = |Y|$ ;

UNTIL $N \leq M$

PERFORM AN APPROPRIATE Selection on the remaining elements and Output.

---

ANALYSIS: Use Chernoff bounds to

Show that $S \leq \frac{3}{4} M$ w.h.p.

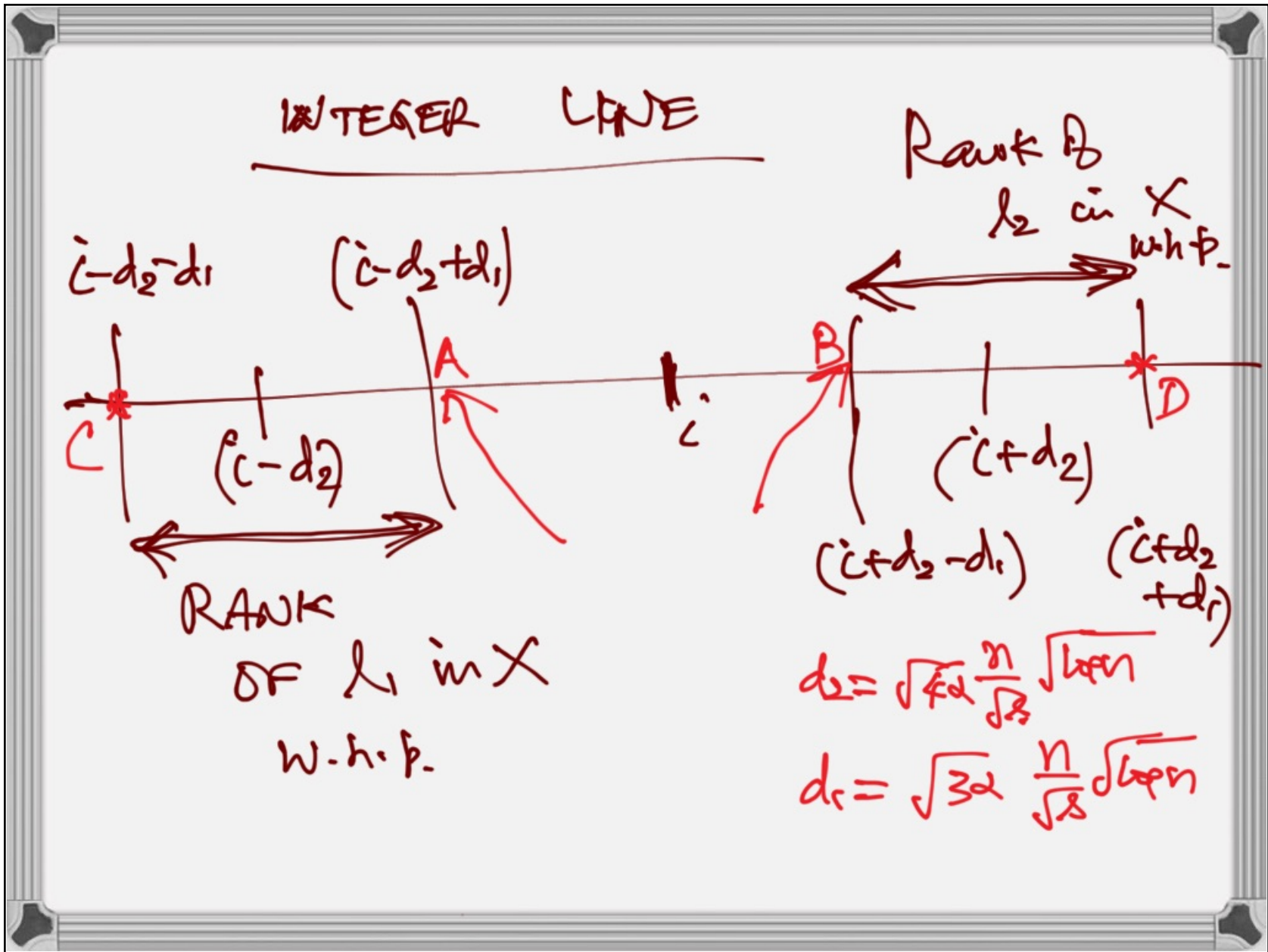Step 1 takes $\frac{N}{B}$ I/O Operations.

Step 2 takes no I/O.

Step 3 takes $\frac{N}{B}$ I/o Operations.

$$E\left[\text{Rank}(l_1, x)\right] = i - \sqrt{4\alpha}\, \frac{N}{\sqrt{8}}\sqrt{\log n}$$

$$E\left[\text{Rank}(l_2, x)\right] = i + \sqrt{4\alpha}\, \frac{N}{\sqrt{8}}\cdot\sqrt{\log n}$$

$$d_2 = \sqrt{4\alpha}\, \frac{N}{\sqrt{8}}\sqrt{\log N} \qquad\longrightarrow ②$$

$$d_1 = \sqrt{3\alpha}\, \frac{N}{\sqrt{8}}\sqrt{\log N} \qquad\longrightarrow ①$$

# INTEGER LINE

$i - d_2 - d_1$

$(i - d_2 + d_1)$

Rank B

$l_2$ in X

w.h.p.

A

C

$(c - d_2)$

$l_1$

B

D

$(c + d_2)$

RANK OF $l_1$ in X

w.h.p.

$(c + d_2 - d_1)$

$(c + d_2 + d_1)$

$$d_2 = \sqrt{4\alpha} \, \frac{n}{\sqrt{8}} \sqrt{\log n}$$

$$d_1 = \sqrt{3\alpha} \, \frac{n}{\sqrt{8}} \sqrt{\log n}$$

$$|Y| \leq D - c \quad \text{w.h.p}$$

$$= (i + d_2 + d_1) - (i - d_1 - d_2)$$

$$= 2(d_1 + d_2)$$

$$= 2 \frac{N}{\sqrt{S}} \sqrt{\log N} \left( \sqrt{4\alpha} + \sqrt{3\alpha} \right).$$

$\Rightarrow$ CONDITION 2a holds w.h.p.

Also, 2b holds. $|Y| = \tilde{O}\left( \dfrac{N}{M^{0.4}} \right)$

Total # of I/o Operations

$$= 2\frac{n}{B} + 2\frac{n}{M^{0.4}B} + 2\frac{n}{M^{0.8}B} + \cdots$$

$$= (2+\epsilon)\frac{n}{B} \quad \text{For some } \epsilon < 1.$$

w.h.p.

PROBLEM: GRAPH SEARCH:-

INPUT:- AN UNDIRECTED GRAPH $G(V,E)$.

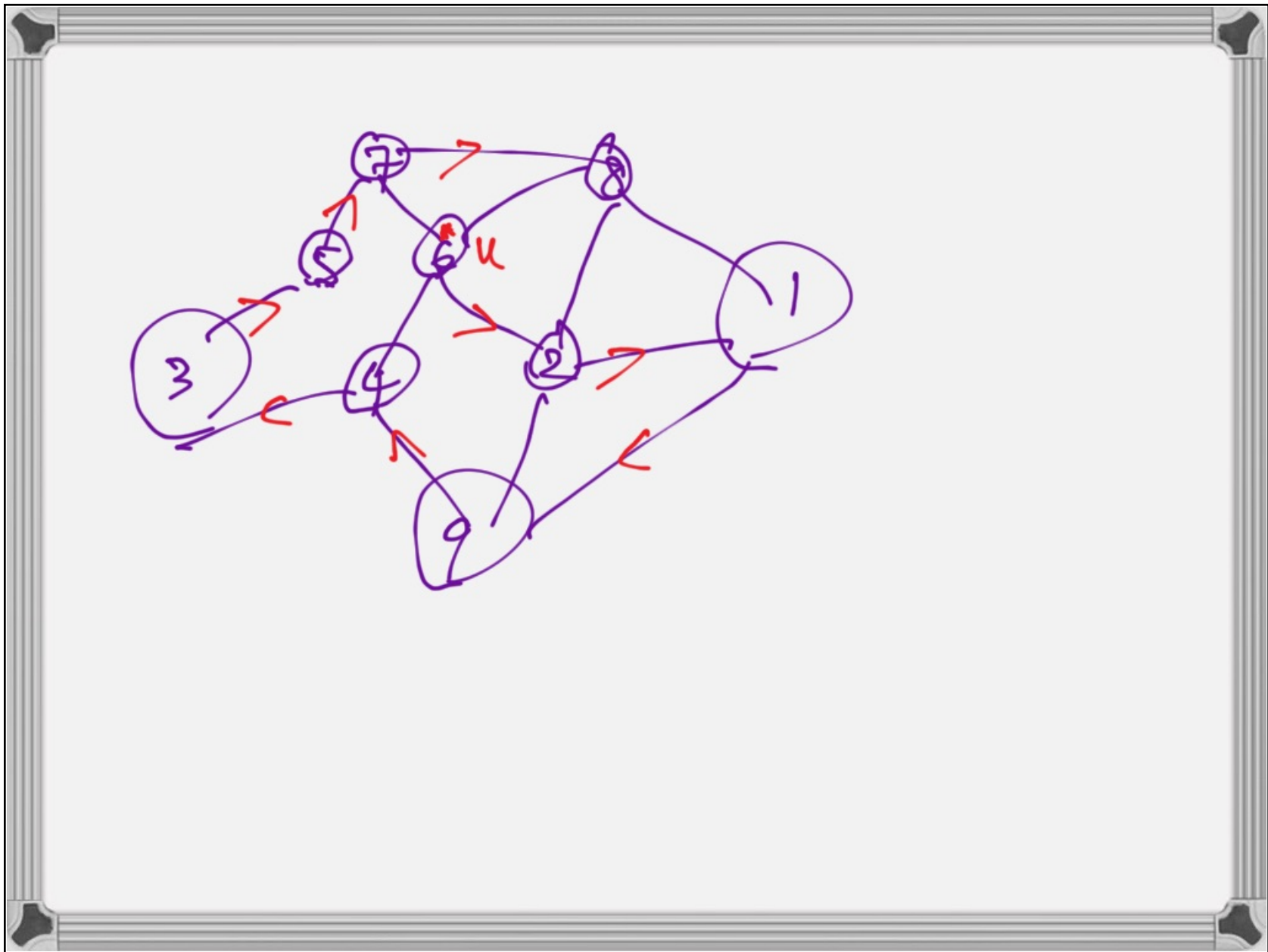Goal: SEARCH the graph.

VISITED $[i] =$ False; $1 \leq i \leq |V|$.

DFS $(u)$

Process $(u)$; visited $[u] := true$;

For $w \in Adj(u)$ do

IF ( visited $(w)$ then

DFS $(w)$;

## Out-of-Core Algorithm:

Assume that $M = \Theta(|V|)$.

Note that for every node in the graph, we have to access all 8 its neighbors.

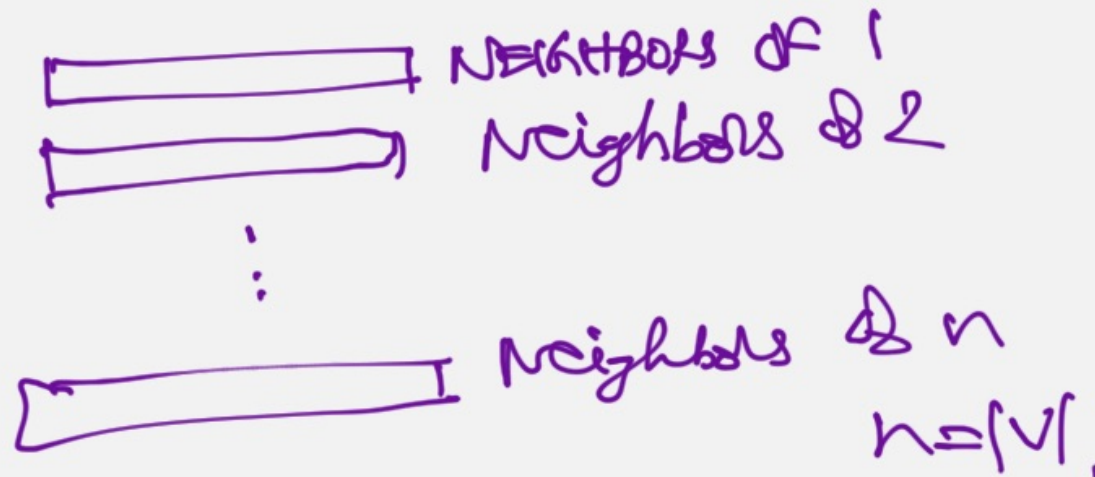Assume that we use Adjacency List. and BREADTH FIRST SEARCH.

## BFS (G):

Start FROM a node $u$.

Visit     nodes at a distance 1;

Visit nodes at a distance 2;

and so on.

```
[_____]  NEIGHBORS OF 1
[_____]  Neighbors of 2
        ⋮
[_____]  Neighbors of n
                    n = |V|.
```

## I/o Complexity:

$$\sum_{u \in V} \left\lceil \frac{d_u}{B} \right\rceil \ ; \qquad d_u \rightarrow \text{DEGREE OF } u$$

$$\text{FOR ANY } u \in V.$$

$$\leq \sum_{u \in V} \left( \frac{d_u}{B} + 1 \right) = \frac{2|E|}{B} + |V|$$

$$= O\left( \frac{|E|}{B} + |V| \right)$$

This is $\cancel{\text{formal}}$ if $|E| > |V|B$.

# PARALLEL DISKS MODEL:



COMPUTER

ONE I/O.

DISKS 1 2 ... D

## SORTING:

The input will be across the disks. The output will also be across the disks.