

Goal: Estimate # of edges

+ Sample edges almost u.a.r

in  $\tilde{O}\left(\frac{n}{\sqrt{m}}\right)$  queries

$n$  - # of vertices,  $m$  - # of edges

Goal: Estimate # of edges

+ Sample edges almost u.a.r  
in  $\tilde{O}\left(\frac{n}{\sqrt{m}}\right)$  queries

$n$  - # of vertices,  $m$  - # of edges

(You have already seen an algorithm for approximating  
 $m$  up to multiplicative factor in  $\tilde{O}\left(\frac{n}{\sqrt{m}}\right)$  queries  
by Goldreich & Ron

Query Model: Adjacency list :

Query Model: Adjacency list:

- The vertices are labeled arbitrarily  $1..n$ , and Alg knows  $n$ .

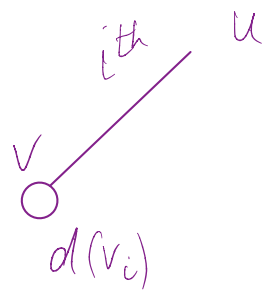
Query Model: Adjacency list:

- The vertices are labeled arbitrarily  $1..n$ , and Alg knows  $n$ .
- Degree queries:  $\text{deg}(v)$  returns  $d_v$ .

$v$   
○  
 $d(v_i)$

## Query Model: Adjacency list:

- The vertices are labeled arbitrarily  $1..n$ , and Alg knows  $n$ .
- Degree queries:  $\text{deg}(v)$  returns  $d_v$ .
- Neighbor queries:  $\text{nbz}(v, i)$  returns the  $i^{\text{th}}$  neighbor of  $v$ .  
if one exists. O.w. returns  $\perp$ .



General approach in sublinear estimation:

General approach in sublinear estimation :

Define some unbiased estimator, s.t.:





General approach in sublinear estimation:

Define some unbiased estimator, s.t.:

1.  $EX[X] = f(m)$  where from  $f(m)$  we can deduce  $m$ .

E.g.:  $EX[X] = m$ , or  $EX[X] = \frac{1}{n} \cdot m$

$$\in \left[ \frac{m}{2n}, \frac{m}{n} \right]$$

2. The (1) maximum value (weaker)

or (2) variance of  $X$  (stronger) are bounded

$$\frac{10 \cdot (H-m)}{n}$$

General approach in sublinear estimation:

Define some unbiased estimator, s.t.:

1.  $EX[X] = f(m)$  where from  $f(m)$  we can deduce  $m$ .

E.g.:  $EX[X] = m$ , or  $EX[X] = \frac{1}{n} \cdot m$

2. The (1) maximum value (weaker)

or (2) variance of  $X$  (stronger) are bounded

Then, by Chernoff inequality (for (1)) and Chebyshev's inequality (for (2))

we can bound the number of sufficient samples. Proof follows.

Cheznoff inequality:

Let  $x_1, \dots, x_k$  be independent random variables, s.t.  $\forall i$   $x_i$  in  $[0, B]$  and let  $x = \sum_{i=1}^k x_i$ . Then:

$$\Pr[|x - EX[x]| > \varepsilon \cdot EX[x]] < 2 \cdot \exp\left(-\frac{\varepsilon^2 \cdot EX[x]}{3 \cdot B}\right)$$

$\Rightarrow$  to get  $\Pr[\text{deviation}] < \delta$ , should set

$$k \text{ to } \frac{3B \cdot \ln(2/\delta)}{\varepsilon^2 \cdot EX[x_i]} \quad (\Rightarrow) \quad \boxed{\frac{\text{Max}}{EX[x_i]}}$$

## Chebyshev's inequality:

- $\Pr[|X - EX[X]| > \varepsilon \cdot EX[X]] < \frac{\text{Var}[X]}{\varepsilon^2 \cdot EX[X]^2}$

- For  $X = \sum_{i=1}^k X_i$ , s.t.  $\forall i, j \quad \text{Var}[X_i] = \text{Var}[X_j]$ ,

$$\Rightarrow \Pr[\dots] < \frac{\sum \text{Var}[X_i]}{\varepsilon^2 \cdot k^2 \cdot EX[X_i]^2} = \frac{k \cdot \text{Var}[X_i]}{\varepsilon^2 \cdot k^2 \cdot EX[X_i]^2} = \frac{\text{Var}[X_i]}{\varepsilon^2 \cdot k \cdot EX[X_i]^2}$$

So to get  $\Pr[\dots] < \delta$ , we need to set

$$k = \frac{\text{Var}[X_i]}{\varepsilon^2 \cdot EX[X_i]^2 \cdot \delta}$$

## Chebyshev's inequality:

- $\Pr[|X - EX[X]| > \varepsilon \cdot EX[X]] < \frac{\text{Var}[X]}{\varepsilon^2 \cdot EX[X]^2}$

- For  $X = \sum_{i=1}^k X_i$ , s.t.  $\forall i, j \quad \text{Var}[X_i] = \text{Var}[X_j]$ ,

$$\Rightarrow \Pr[\dots] < \frac{\sum \text{Var}[X_i]}{\varepsilon^2 \cdot k^2 \cdot EX[X_i]} = \frac{k \cdot \text{Var}[X_i]}{\varepsilon^2 \cdot k^2 \cdot EX[X_i]} = \frac{\text{Var}[X_i]}{\varepsilon^2 \cdot k \cdot EX[X_i]}$$

So to get  $\Pr[\dots] < \delta$ , we need to set

$$k = \frac{\text{Var}[X_i]}{\varepsilon^2 EX[X_i]^2 \delta}$$

## Back to estimating # edges:

First, naive, attempt:

For every  $v \in V$ , let

$$X_v = d(v)$$

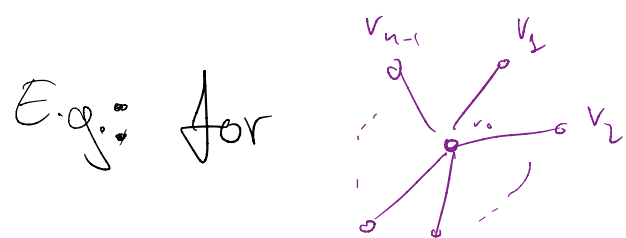
$$\text{Then, } \mathbb{E}[X_v] = \sum_{v \in V} \frac{1}{n} \cdot d(v) = \frac{2m}{n}$$

$\Rightarrow$  unbiased estimator!

But, variance might be too high.

But, *variance* might be too high.

$\exists$  graphs s.t.  $\text{var}_v[X_v] = n \cdot \mathbb{E}[X_v]$



$$\mathbb{E}_v[X_v] = \frac{1}{n} \sum d(v) = \Theta(1)$$

$$\text{Var}[X_v] = \frac{1}{n} \sum x_v^2 = \frac{1}{n} \sum d_v^2 = \frac{1}{n} \cdot (n-1) + \frac{1}{n} = \Omega(n)$$

So, by Chebyshev's we'll have to set  $k$

$$k = \frac{\text{Var}[X_v]}{\varepsilon^2 \mathbb{E}[X_v]^2} = \Omega\left(\frac{n}{\varepsilon^2 \cdot 1}\right) = \Omega(n) \quad \ddot{\smile}$$

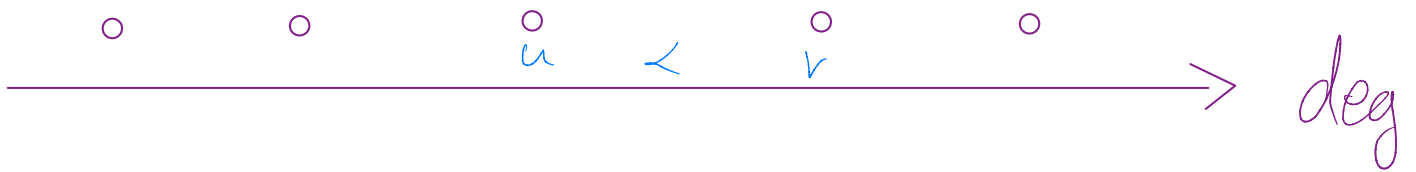


Second attempt :

1. Order all vertices according to degree

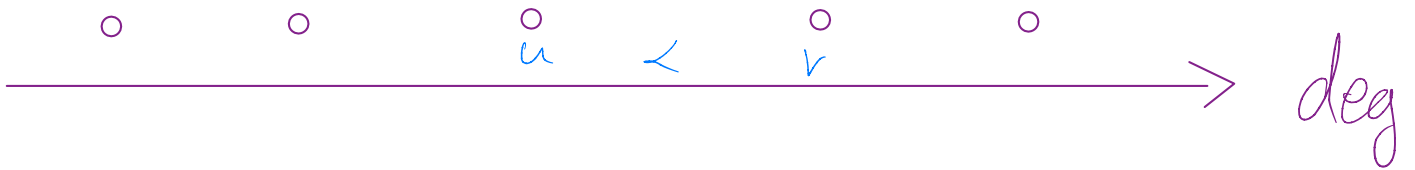
Second attempt :

1. Order all vertices according to degree



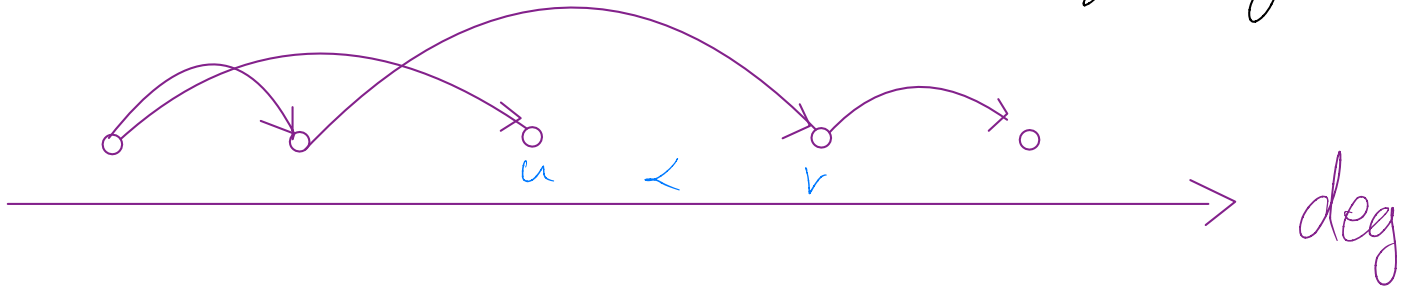
Second attempt :

1. Order all vertices according to degree
2. Orient edges from low deg to high deg vertices



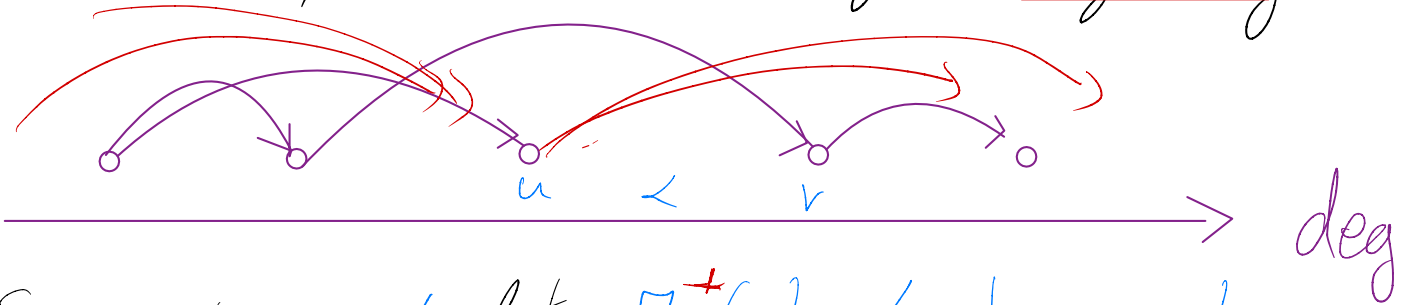
Second attempt :

1. Order all vertices according to degree
2. Orient edges from low deg to high deg vertices



Second attempt :

1. Order all vertices according to degree
2. Orient edges from low deg to high deg vertices

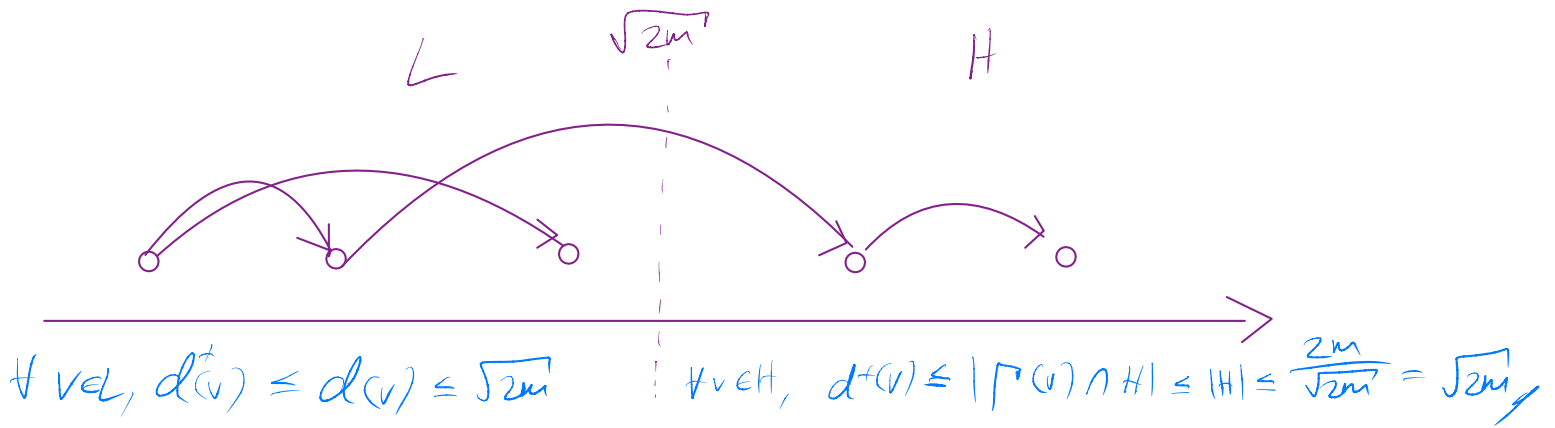


3. For each  $u \in V$  let  $\Gamma^+(u) = \{v \mid u < v\}$   
and let  $d^+(u) = |\Gamma^+(u)|$ .

Observe that  $\sum_{v \in V} d^+(v) = 2m$

(since every edge is counted from exactly one endpoint)

Claim:  $\forall v, d^+(v) \leq \sqrt{2m}$



2<sup>nd</sup> (wishful) attempt:

$$\sum_{v \in V} d^+(v) = m$$

$$\text{Let } \chi_v = d^+(v)$$

$$\text{Then } \mathbb{E} \chi_v = \frac{1}{n} \sum_{v \in V} d^+(v) = \frac{m}{n}$$

→ unbiased estimator!

Also, by previous analysis,

$$\max_{v \in V} \chi_v \leq \sqrt{2m} \Rightarrow \text{bounded max. value!}$$

Hence, by Chebyshev's / Chernoff,

$$\text{Sufficient to take } k = \frac{\max}{\mathbb{E} \chi} \cdot \frac{3 \log(2/\delta)}{\epsilon^2} = \tilde{O}\left(\frac{\sqrt{m}}{m/n}\right) = \tilde{O}\left(\frac{n}{\sqrt{m}}\right)$$

Problem: We don't have access to  $d^+(v)$  queries.

Solution: We'll set a random variable so  
that in Expectation,  $E[X_u] = d^+(u)$



Problem: We don't have access to  $d^+(v)$  queries.

Solution: We'll set a random variable so that in Expectation,  $E[X_u] = d^+(u)$

Given  $u \in V$ :

1. Query  $d(u)$  &

2. Query a uniform nbr of  $u$ , denoted  $v$   $d(u)$

3. Let

$$X_u = \begin{cases} d(u) & \text{if } u < v \\ 0 & \text{o.w.} \end{cases}$$



$$E_v[X_u | u] = \frac{d^+(u)}{d(u)} \cdot d(u) = d^+(u)$$

Problem: We don't have access to  $d^+(v)$  queries.

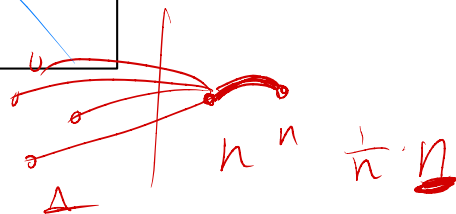
Solution: We'll set a random variable so that in Expectation,  $E[X_u] = d^+(u)$

Given  $u \in V$ :

1. Query  $d(u)$  & a uniform nbr of  $u$ , denoted  $v$
2. Query  $d(v)$
3. Let

$$X_u = \begin{cases} \underline{d(u)} & \text{if } u < v \\ 0 & \text{otherwise} \end{cases}$$

$$E[X_u] = \frac{d^+(u)}{d(u)} \cdot d(u) = \underline{d^+(u)}$$



Problem: We don't have access to  $d^+(v)$  queries.

Solution: We'll set a random variable so that in Expectation,  $E[X_u] = d^+(u)$

Given  $u \in V$ :

1. Query  $d(u)$  & a uniform nbr of  $u$ , denoted  $v$

2. Query  $d(v)$

3. Let

$$X_u = \begin{cases} d(u) & \text{if } u < v \\ 0 & \text{otherwise} \end{cases}$$

$$E[X_u] = \frac{d^+(u)}{d(u)} \cdot d(u) = d^+(u)$$

Problem: max value  $d(u)$  can be as high as  $n$ .

Problem: We don't have access to  $d^+(v)$  queries.

Solution: We'll set a random variable so that in Expectation,  $E[X_u] = d^+(u)$

Given  $u \in V$ :

1. Query  $d(u)$  & a uniform nbr of  $u$ , denoted  $v$
2. Query  $d(v)$
3. Let

$$X_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2m/\epsilon} \\ 0 & \text{otherwise} \end{cases}$$

$$E[X_u] = \frac{d^+(u)}{d(u)} \cdot d(u) = d^+(u)$$

~~Problem: max value  $d(u)$  can be as high as  $n$ .~~

Analysis:

$$\chi_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2m/\epsilon} \\ 0 & \text{otherwise} \end{cases}$$

Analysis:

$$x_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2n/\epsilon} \\ 0 & \end{cases}$$

Observe that there are two sources of randomness:

The choice of  $u$  \& the choice of  $u$ 's nb2,  $v$ .

Analysis:

$$X_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2m/\epsilon} \\ 0 & \text{otherwise} \end{cases}$$

Observe that there are two sources of randomness:

The choice of  $u$  & the choice of  $u$ 's nbz,  $v$ .

Fix some  $u \in V$ .

If  $d(u) \leq \sqrt{2m/\epsilon}$ ,  $\mathbb{E}_{v \in \mathcal{N}(u)} [X_u | u] = \underline{d^+(u)}$

Analysis:

$$X_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2n/\epsilon} \\ 0 & \text{if } u > v \text{ \& } d(u) > \sqrt{2n/\epsilon} \end{cases}$$

Observe that there are two sources of randomness:

The choice of  $u$  & the choice of  $u$ 's nbz,  $v$ .

Fix some  $u \in V$ .

$$\text{If } d(u) \leq \sqrt{2n/\epsilon}, \quad \underset{v \in \mathcal{N}(u)}{\text{EX}} [X_u | u] = d^+(u)$$

$$\text{If } \underline{d(u) > \sqrt{2n/\epsilon}} \quad \text{then} \quad \underline{\text{EX}[X_u | u] = 0}$$



Analysis:

$$X_u = \begin{cases} d(u) & \text{if } d(u) \leq \sqrt{2m/\epsilon} \text{ \& } u < v \\ 0 & \text{otherwise} \end{cases}$$

Observe that there are two sources of randomness:

The choice of  $u$  & the choice of  $u$ 's nbz,  $v$ .

Fix some  $u \in V$ .

$$\text{If } d(u) \leq \sqrt{2m/\epsilon}, \quad \underset{v \in \mathcal{N}(u)}{\text{EX}} [X_u | u] = d^+(u)$$

$$\text{If } d(u) > \sqrt{2m/\epsilon} \text{ then } \text{EX}[X_u | u] = 0$$

$$\text{Therefore, } \text{EX}_{u,v} [X_u] = \frac{1}{n} \sum_{\substack{u \in V, \\ d(u) \leq \sqrt{2m/\epsilon}}} d^+(u)$$

$$EX_{u,v} [X_u] = \frac{1}{n} \sum_{\substack{u \in V \\ d(u) \leq \sqrt{2m/\epsilon}}} d^+(u)$$

---

$$X_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2m/\epsilon} \\ 0.6 & 0 \end{cases}$$

$$EX_{u,v} [X_u] = \frac{1}{n} \sum_{\substack{u \in V \\ d(u) \leq \sqrt{2m/\epsilon}}} d^+(u)$$

$$X_u = \begin{cases} d(u) & \text{if } u < v \text{ \& } d(u) \leq \sqrt{2m/\epsilon} \\ 0 & \text{otherwise} \end{cases}$$

Claim:  $\sum_{u \in V, d(u) \leq \sqrt{2m/\epsilon}} d^+(u) \in [(1-\epsilon)m, m]$

Proof: There are at most  $\frac{2m}{\sqrt{2m/\epsilon}} = \sqrt{\epsilon 2m}$

vertices whose  $\text{deg.} > \sqrt{2m/\epsilon}$ . Denote their set by  $H$

$$\sum_{v \in H} d^+(v) \leq \binom{|H|}{2} \leq \epsilon m$$

$$\sum_{u \in V, d(u) \leq \sqrt{2m/\epsilon}} d^+(u) = \sum_{u \in V} d^+(u) - \sum_{u \in H} d^+(u) \geq (1-\epsilon)m$$

and also  $\sum_{u \in V, d(u) \leq \sqrt{2m/\epsilon}} d^+(u) \leq \sum_{u \in V} d^+(u) \leq m$   $\square$

Hence, we have:

$$X_u = \begin{cases} d(u) & \text{if } u < v \\ 0 & \text{otherwise} \end{cases} \quad \text{and } d(u) \leq \sqrt{2m/\epsilon}$$

$$EX_{u,v} [X_u] \in \frac{1}{n} \cdot [(1-\epsilon)m, m]$$

$$\wedge \max_{u,v} \{X_u\} \leq \sqrt{2m/\epsilon}$$

$\Rightarrow$  By Chernoff, sufficient to set

$$k = \frac{\max \{X_u\} \cdot 3 \log(1/\delta)}{EX[X_u]} \leq \frac{n}{\sqrt{m}} \cdot \frac{6 \log(2/\delta)}{\epsilon^{2.5}}$$

Hence, we have:

$$EX_{u,v} [X_u] \in \frac{1}{n} \cdot [(1-\varepsilon)m, m]$$

$$\Delta \max_{u,v} \{X_u\} \leq \sqrt{2m/\varepsilon}$$

$\Rightarrow$  By Chernoff, sufficient to set

$$k = \frac{\max \{X_u\}}{EX[X_u]} \cdot \frac{3 \log(1/\delta)}{\varepsilon^2} \leq \boxed{\frac{n}{\sqrt{m}}} \cdot \frac{6 \log(2/\delta)}{\varepsilon^{2.5}}$$

Final alg: Input:

- $\tilde{m} \in [m, 2m]$  const. approx. of  $m$ .
- $\epsilon$  - approximation parameter
- $\delta$  - confidence parameter

1. Repeat  $k = \frac{n}{\tilde{m}} \cdot \frac{\log(V\delta)}{\epsilon^2\delta}$  times:

a. Sample a vertex  $u_i \in V$  u.o.r. & query  $d(u)$

b. Sample a uniform nbr of  $u$ , denote it  $v_i$  and query  $d(v_i)$

c. Let  $X_{u_i} = \begin{cases} d(u_i) & \text{if } u \leq v \text{ \& } d(u_i) \leq \sqrt{2m/\epsilon} \\ 0 & \text{o.w.} \end{cases}$

2. Let  $X = \frac{1}{k} \sum_{i=1}^k X_{u_i}$

3. Return  $\tilde{m} = n \cdot X$  as the estimate.