

## Lecture 3:

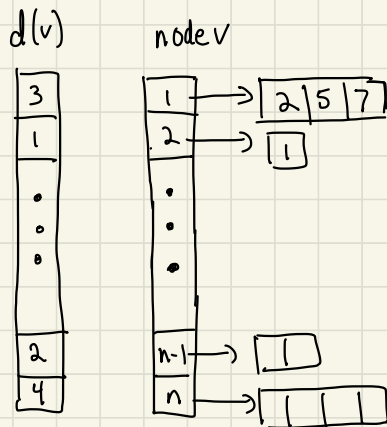
- Estimate average degree
  - recap
  - 2-approximation
  - $1+\epsilon$ -approximation

# Estimating the average degree of a graph

def Average degree  $\bar{d} = \frac{\sum_{u \in V} d(u)}{n}$

Assume:  $G$  simple (no parallel edges, self-loops)  
 $\Omega(n)$  edges (not "ultra-sparse")

Representation via adj list + degrees:



- degree queries: on  $v$  return  $d(v)$
- neighbor queries: on  $(v, j)$  return  $j$ th nbr of  $v$

Naive sampling:

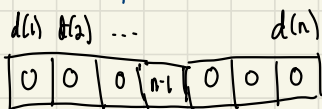
Pick  $O(??)$  sample nodes  $v_1 \dots v_s$

output ave degree of sample:

$$\frac{1}{s} \sum_i d(v_i)$$

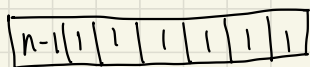
Straight forward Chernoff/Hoeffding needs  $\Omega(n)$  samples

lower bound?



need  $\Omega(n)$  samples to find "needle in haystack"

not a possible degree sequence!!



is possible

Some lower bounds:

"ultrasparse" case:

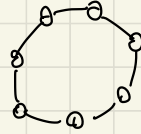
0 edges vs. 1 edge

need  $\Omega(n)$  queries to distinguish

$\Rightarrow$  multiplicative approx needs  $\Omega(n)$

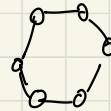
ave deg  $\geq 2$ :

$n$ -cycle  $\bar{d}=2$



vs.

$n - c\sqrt{n}$  cycle  $\bar{d} \approx 2 + c^2$   
+  $c\sqrt{n}$ -clique



need  $\Omega(n^{1/2})$  queries to find  
clique node

Algorithm idea:

group nodes of similar degrees  
estimate average w/in each group

why does this help?

recall Chernoff:

$X_1, \dots, X_r$  iid  $X_i \in [0, 1]$

$$S = \sum_{i=1}^r X_i \quad p = E[X_i] = E[S]/r \quad -\Omega(rp\delta^2)$$

$$\text{Then } \Pr\left[\left|\frac{S}{r} - p\right| \geq \delta p\right] \leq e^{-\Omega(rp\delta^2)}$$

$\Rightarrow$   $r$  needs to be

$$\Omega\left(\frac{1}{p\delta^2}\right)$$

let's assume  $\delta$  is a constant

$X_i$  needs to be in  $[0, 1]$

so if  $X_i \leftarrow \frac{\deg(i)}{n}$

then  $p$  can be as small as  $\frac{1}{n}$

$\Rightarrow$   $r$  needs to be  $\Omega(1/p) = \Omega(n)$

but if  $b \leq \deg(i) \leq (1+\varepsilon)b$

can set  $X_i \leftarrow \frac{\deg(i)}{(1+\varepsilon)b}$

then  $p \geq \frac{1}{1+\varepsilon}$

$\Rightarrow$   $r$  needs to be only  $\Omega(1)$ . Much better!!!

- + each group has bounded variance
- doesn't work for arbitrary  $\epsilon$ 's  
why here?

### Bracketing:

set parameters  $\beta = \frac{\epsilon}{c}$   
 $t = O(\log n / \epsilon)$  #buckets

$$B_i = \{ v \mid (1+\beta)^{i-1} < d(v) \leq (1+\beta)^i \}$$

for  $i \in \{0, \dots, (t-1)\}$

(can add bucket for deg 0 nodes  
 or  
 \* assume none)

note: total degree of nodes in  $B_i$   
 $(1+\beta)^{i-1} |B_i| \leq d_{B_i} \leq (1+\beta)^i |B_i|$

total degree of graph:

$$\sum_i (1+\beta)^{i-1} |B_i| \leq d_{\text{total}} \leq \sum_i (1+\beta)^i |B_i|$$

First idea for algorithm:

$$G_j = \begin{cases} 1 & \text{if sample } j \\ & \text{falls in bucket } i \\ 0 & \text{o.w.} \end{cases}$$

• Take sample  $S$  of nodes

$$S_i \leftarrow S \cap B_i$$

(samples that fall in  $i$ th bucket  
use degree queries to  
determine)

• estimate  $|B_i|$ :

$$p_i \leftarrow \frac{|S_i|}{|S|}$$

$$\text{note: } E[p_i] = E\left[\frac{|S_i|}{|S|}\right] = \frac{E\left[\sum_{j=1}^{|S|} G_j^{(i)}\right]}{|S|} \\ = \frac{|S|}{|S|} \cdot \frac{|B_i|}{n}$$

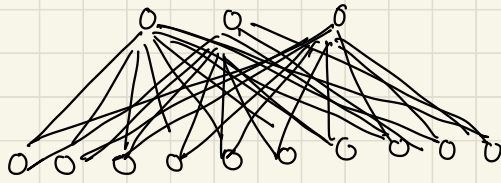
• Output  $\sum_i p_i (1+\beta)^{i-1}$  ← undercounting

Problem:

if  $i$  is st.  $|S_i|$  small, will need lots of  
samples to approx

these likely come from  $B_i$  st.  $|B_i|$  is small

example:



← 3 nodes  
each deg  $n-3$

←  $n-3$  nodes  
each deg 3

$$a \leftarrow i \text{ st. } (1+\beta)^{i-1} \leq 3 \leq (1+\beta)^i$$

$$b \leftarrow i \text{ st. } (1+\beta)^{i-1} \leq n-3 \leq (1+\beta)^i$$

$$\forall c \neq a, b \quad |B_c| = 0$$

$$|B_a| = n-3$$

$$|B_b| = 3$$



both contribute  
 $(n-3) \cdot 3$  edges

but these are not likely to be sampled

still, maybe good enough  
for 2-approximation?

Next idea:

Use "0" for small buckets



Old algorithm:

• Take sample  $S$

•  $S_i \leftarrow S \cap B_i$

• estimate  $|B_i|$ :

$$p_i \leftarrow \frac{|S_i|}{|S|}$$

• Output  $\sum_i p_i (1+\beta)^{i-1}$

New algorithm:

• Take sample  $S$  (how big?)

•  $S_i \leftarrow S \cap B_i$

• estimate  $|B_i|$ :

for all  $i$

if  $|S_i| \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{C \cdot t}$  "big"

use  $p_i \leftarrow \frac{|S_i|}{|S|}$

else  $p_i \leftarrow 0$  "small"

• Output  $\sum_i p_i (1+\beta)^{i-1}$

why  $\sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{C \cdot t}$ ?

let  $|S| = \Theta(\sqrt{n} \text{ polylog } n \cdot \text{poly}(\frac{1}{\epsilon}))$

then  $|S_i| \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{C \cdot t} \Rightarrow |S_i| \geq \Omega(\text{polylog } n \times \text{poly}(\frac{1}{\epsilon}))$

$\Rightarrow$   
union bnd  
+  
Chernoff

$$\forall i \quad (1-\delta) \frac{|B_i|}{n} \leq p_i \leq (1+\delta) \frac{|B_i|}{n}$$

for  $\delta \sim \Theta(\epsilon)$

Why these settings of  $S$ ? (ignore dependence on  $\epsilon$  for now)

\* each bucket that has at least  $\approx \frac{1}{\sqrt{n}}$  fraction of nodes should have enough samples to be able to estimate the fraction.

\* why  $\approx \frac{1}{\sqrt{n}}$ ?

- we will want to argue that "small" buckets represent a very small fraction of the edges so it is ok to zero them out

- remember the clique lower bound example? if we set the "small" threshold to bigger than  $\frac{1}{\sqrt{n}}$  we might miss lots of edges (e.g. a clique on  $\sqrt{n}$  nodes will have  $\Theta(n)$  edges & shouldn't be missed, but represents only  $\frac{\sqrt{n}}{n} = \frac{1}{\sqrt{n}}$  fraction of nodes)

- why is  $\frac{1}{\sqrt{n}}$  small enough?

See later!

\* what is "enough" samples for each bucket?

- we will need to argue that we are getting good estimates of  $\frac{|B_i|}{n}$  for each big bucket

Chernoff bound  
union bound over  $\log n$  buckets

so need prob of having bad estimate " $\delta$ " set to  $\leq \frac{1}{\log n}$  per bucket

Chernoff will also depend on accuracy parameter  $\beta = \frac{\epsilon}{c}$

So if we set  $S \approx \sqrt{n} \cdot \text{poly}(\frac{1}{\epsilon}) \cdot \text{poly}(\log n)$

we should be more than ok

to get buckets with  $\frac{1}{\sqrt{n}}$  fraction of nodes  
this comes in everywhere  
to satisfy Chernoff & union bounds

# Analysis:

1) Output not too large:

idealistic case

$$\text{Suppose } \forall i \quad p_i = \frac{|B_i|}{n},$$

$$\text{then } \sum_i p_i (1+\beta)^{i-1} = \sum_i \frac{|B_i|}{n} (1+\beta)^{i-1} \leq \bar{d}$$

$\leq$  deg of nodes in  $B_i$

realistic case

$$\text{Suppose } \forall i \quad p_i \leq \frac{|B_i|}{n} (1+\gamma)$$

$$\Rightarrow \sum_i p_i (1+\beta)^{i-1} \leq \bar{d} (1+\gamma)$$

e.g. when  $i$  is big

2) Can output be too small?

$$\text{if } \forall i \quad p_i = \frac{|B_i|}{n} \quad \text{then } \sum_i p_i (1+\beta)^{i-1} = \sum_i \frac{|B_i|}{n} (1+\beta)^{i-1}$$

$$\begin{aligned} &\xrightarrow{\text{multiply by } (1+\beta)(1-\beta) < 1} \geq (1-\beta) \sum_i \frac{|B_i|}{n} (1+\beta)^i \\ &\geq (1-\beta) \bar{d} \end{aligned}$$

by sampling, for big  $i$ ,  $p_i \approx \frac{|B_i|}{n} (1-\epsilon)$

for small  $i$  ????

How much undercounting?

divide edges into 3 types

- type determined by run of algorithm
- 1) big-big: both endpoints in big buckets counted twice
  - 2) big-small: one endpoint in big bucket, " " "small" counted once
  - 3) small-small: both endpoints in small buckets "never" counted

note: small-small can be a big problem

big-small only undercounted by a factor of 2

Example:

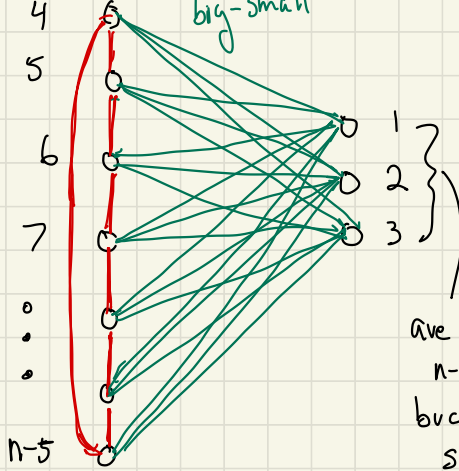
big-big

big-small

$n-8$   
nodes

ave deg  
5

bucket a  
big

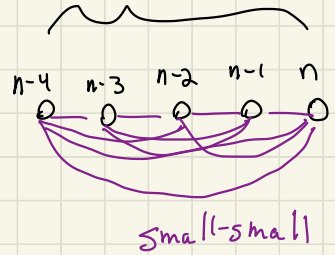


ave deg  
 $n-5$   
bucket b  
small  
3 nodes

5 nodes

ave deg 4

bucket c small



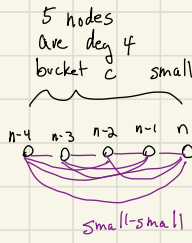
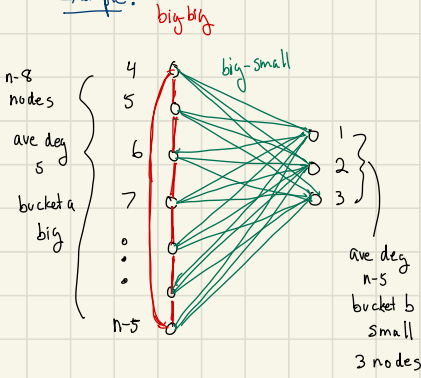
small-small

$$\text{Total degree: } 5 \cdot (n-8) + (n-8) \cdot 3 + 4 \cdot 5 = 8(n-8) + 20$$

ave degree  $\approx 8$

algorithm will likely output  $\approx 5$

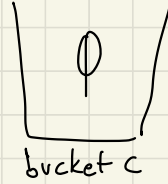
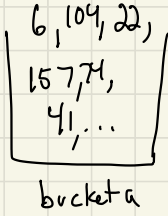
Example:



New algorithm:

- Take sample  $S$  (how big?)
- $S_i \leftarrow S \cap B_i$
- estimate  $|B_i|$ :
  - for all  $i$ 
    - if  $|S_i| \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{c \cdot t}$  <sup>"big"</sup>
    - use  $p_i \leftarrow \frac{|S_i|}{|S|}$
    - else  $p_i \leftarrow 0$  <sup>"small"</sup>
- Output  $\sum_i p_i (1+\beta)^{i-1}$

Samples:



↑  
most nodes here

whp  $p_a \ll 1$

output  $\approx 1.5$

↗  
few, if any, in these buckets  
⇒ whp  $b+c$  are small so likely  
that  $p_b = p_c = 0$

Good news:

Small buckets can't have many nodes  
 $\Rightarrow$  bound on total # small-small edges

$$\begin{aligned} \text{if } |B_i| > \frac{2\sqrt{\epsilon n}}{ct} \text{ then expected size of } S_i \\ \text{is } &\geq |S_i| \cdot \frac{|B_i|}{n} \\ &\geq |S_i| \cdot 2\sqrt{\frac{\epsilon}{n}} \cdot \frac{1}{ct} \end{aligned} \quad \left. \begin{array}{l} \text{twice} \\ \text{threshold} \\ \text{for} \\ \text{"big"} \end{array} \right\}$$

so likely algorithm will  
decide that  $i$  "big"

Assume for all  $i$  "small" that  $|B_i| \leq \frac{2\sqrt{\epsilon n}}{ct}$

then total # small-small edges

$$\leq \left( \underbrace{\frac{2\sqrt{\epsilon n}}{ct}}_{\substack{\# \text{ nodes} \\ \text{per small} \\ \text{bucket}}} \cdot \underbrace{t}_{\# \text{ buckets}} \right)^2 = O\left(\frac{\epsilon n}{c^2}\right) = O(\epsilon n)$$



if ignore small-small edges,  
they affect approx of  $\bar{d}$   
by  $\leq \frac{\epsilon n}{n} = \epsilon$  additive factor  
 $\leq (1+\epsilon)$  multiplicative factor

assume  $\bar{d} \geq 1$

First Claim:

Algorithm gives factor  $(2+\epsilon)$ -mult approx

large-small underestimated by factor  $\leq 2$

small-small error

Improving further:

need to improve on "big-small" edges

Can we estimate fraction of them

+ correct for them?

e.g. by sampling random edges?

New queries:

random neighbor query ( $v$ ):

given  $v$ , return random nbr of  $v$

implementation:

1. degree query to  $v$

2. pick random  $i \in [1..deg(v)]$

3. neighbor query ( $v, i$ )

pick (almost) random edge in (big) bucket  $i$ :

pick random edge by sampling nodes

until one falls in bucket  $i$

return random nbr query from that node

Estimate fraction big-small in  $B_i$  (big):

repeat  $O(1/\epsilon)$  times

pick random node  $u \in B_i$

$e \leftarrow$  random nbr of  $u$

set  $a_j$  to be  $\begin{cases} 1 & \text{if } e \text{ "big-small"} \\ 0 & \text{o.w.} \end{cases}$   
( $e$  is "big-big")

Output  $\alpha_i =$  average  $a_j$

Analysis:

easy case: all nodes in  $B_i$  have same degree  $d$

$T_i \leftarrow$  # big-small edges in  $B_i$

$$\Pr[\text{"big-small" edge } e \text{ in } B_i \text{ chosen}] = \frac{1}{|B_i|} \cdot \frac{1}{d}$$

$\parallel$   
 $(u,v)$   
only one of  $u,v$  big  
wlog assume  $u$  big

$\underbrace{\quad}_{\text{prob } u \text{ is chosen}} \cdot \underbrace{\quad}_{\text{prob } (u,v) \text{ output given } u \text{ chosen}}$

$$\text{so } \Pr[a_j = 1] = E[a_j] = \frac{T_i}{d \cdot |B_i|}$$

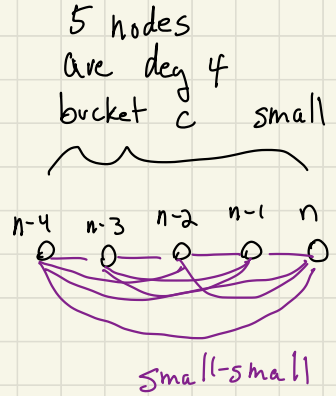
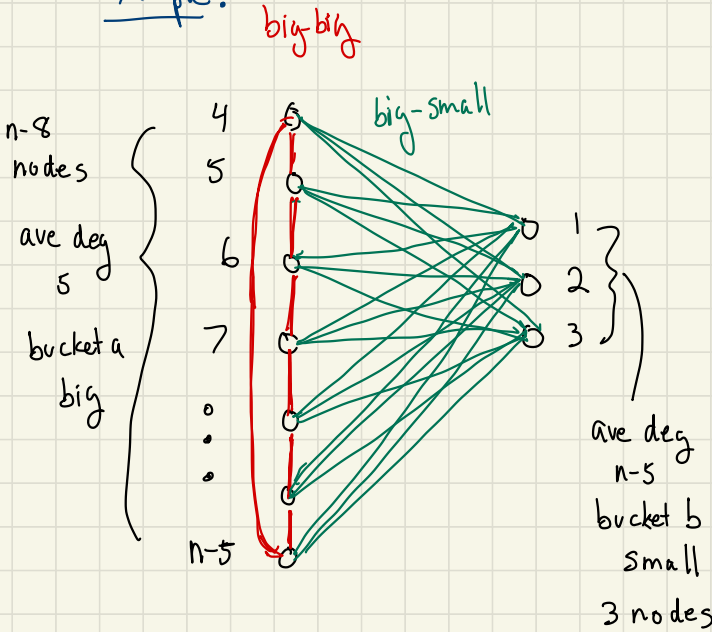
general case: all nodes in  $B_i$  have degrees within  $(1+\beta)$  factor of each other

$$\frac{1}{|B_i| (1+\beta)^i} \leq \Pr[\text{"big-small" edge } e \text{ in } B_i \text{ chosen}] \leq \frac{1}{|B_i| (1+\beta)^{i-1}}$$

$$\frac{T_i}{|B_i| (1+\beta)^i} \leq \underbrace{E[a_j]}_{\substack{\text{estimate to} \\ 1 \pm \varepsilon \text{-mult factor} \\ \text{to get} \\ (1 \pm \varepsilon)(1+\beta) \text{ estimate} \\ \text{of } \frac{T_i}{n} \text{ via } \underbrace{\alpha_i p_i}_{\substack{\text{undercount} \\ \text{of } \# \\ \text{edges in } B_i}} (1+\beta)^{i-1}} \leq \frac{T_i}{|B_i| (1+\beta)^{i-1}} \Rightarrow$$

$$E[a_j] |B_i| (1+\beta)^{i-1} \leq T_i \leq E[a_j] |B_i| (1+\beta)^i$$

Example:



$$\text{Total degree: } 5 \cdot (n-8) + (n-8) \cdot 3 + 4 \cdot 5 = 8(n-8) + 20$$

$$\text{ave degree} \approx 8$$

algorithm will likely output  $\approx 5$

$$\# \text{ big-small edges slots: } 3 \cdot (n-8)$$

$$\text{Fraction of big-big over big-small: } \approx \frac{3(n-8)}{5(n-8)} = \frac{3}{5}$$

$$E[a_j] = \frac{3}{5}$$

$$\text{Output } 1 \cdot \left(1 + \frac{3}{5}\right) \underbrace{\left(1 + \frac{3}{5}\right)^2}_{\approx 5} \approx 8$$

## Final Algorithm:

• sample  $\Theta\left(\frac{\sqrt{n}}{\epsilon} t\right)$  nodes + place in  $S$

•  $S_i \leftarrow S \cap B_i$

• For all  $i$

if  $|S_i| \geq \sqrt{\frac{\epsilon}{n}} \frac{|S|}{ct}$

use  $p_i \leftarrow \frac{|S_i|}{|S|}$

for all  $v \in S_i$

• Pick random nbr  $u$  of  $v$

•  $\chi(v) \leftarrow \begin{cases} 1 & \text{if } u \text{ is small} \\ 0 & \text{o.w.} \end{cases}$

$$\alpha_i \leftarrow \frac{|\{v \in S_i \mid \chi(v) = 1\}|}{|S_i|}$$

else use  $p_i \leftarrow 0$

• Output

$\sum_{\text{large } i}$

$$p_i (1 + \alpha_i) (1 + \beta)^{i-1}$$

big-big +  
one side of big-small

Correction to get  
other side of  
big-small

Where do errors come from?

estimating  $p_i$ 's } mult  $1+\epsilon$ -factor  
estimating  $d_i$ 's }  
Small-small edges } additive  $\epsilon n$  error