

Lecture 3:

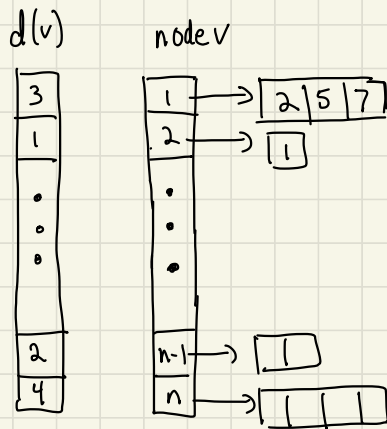
- Estimate average degree
 - recap
 - 2-approximation
 - $1+\epsilon$ -approximation

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:

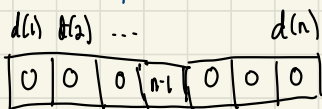
Pick $O(??)$ sample nodes $v_1 \dots v_s$

output ave degree of sample:

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

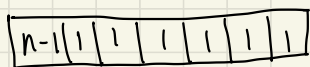
Straight forward Chernoff/Hoeffding needs $\Omega(n)$ samples

lower bound?



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

not a possible degree sequence!!



is possible

Some lower bounds:

"ultrasparse" case:

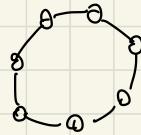
0 edges vs. 1 edge

need $\Omega(n)$ queries to distinguish

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

ave deg ≥ 2 :

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

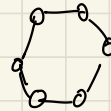


vs.

note:

C can be any const
e.g. 1000

$\left\{ \begin{array}{l} n - C \cdot \sqrt{n} \text{ cycle} \\ + C \sqrt{n} \text{-clique} \end{array} \right. \quad \bar{d} \approx 2 + C^2$



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

\Rightarrow need $\Omega(n^{1/2})$ queries for constant mult approx !!!

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 $\Omega\left(\frac{1}{p\delta^2}\right)$ so p very small is not good!

let's assume δ 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 $X_i \leftarrow \frac{\deg(i)}{(1+\varepsilon)b}$

p is at least a constant \Rightarrow 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 β 's
why here?

Bracketing:

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

$(1+\beta)^t > n$
 when
 $t \geq \frac{\log n}{\log(1+\beta)}$

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

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:

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

First idea for algorithm:

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

• Take sample S of nodes *how many?*

• $S_i \leftarrow S \cap B_i$ ← use degree queries to partition

• estimate $|B_i|$:

define $\delta_j^{(i)} = \begin{cases} 1 & \text{if sample } j \text{ falls in bucket } i \\ 0 & \text{o.w.} \end{cases}$

$$p_i \leftarrow \frac{|S_i|}{|S|} \leftarrow E[p_i] = E\left[\frac{|S_i|}{|S|}\right] = \frac{E\left[\sum_{j=1}^{|S|} \delta_j^{(i)}\right]}{|S|}$$

$$= \cancel{|S|} \cdot \frac{|B_i|}{n} = \frac{|B_i|}{n}$$

• Output $\sum_i p_i (1+\beta)^{i-1}$

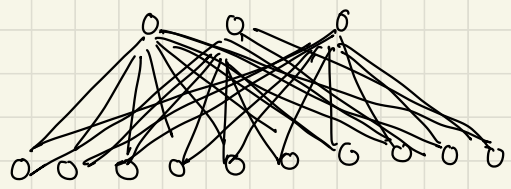
← undercount by defn.

Problem: i s.t. $|S_i|$ is small

$\Rightarrow p_i$ is a bad approx

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

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

B_a contributes $(n-3) \cdot 3$ edges

B_b contributes $3 \cdot (n-3)$ edges

not seen in sample of size $o(n)$

Next idea: use "0" for small buckets

(helps "variance")

Old algorithm:

- Take sample S
- $S_i \leftarrow S \cap B_i$
- estimate $|B_i|$:
$$p_i \leftarrow \frac{|S_i|}{|S|}$$
- Output $\sum_i p_i (1+\beta)^{i-1}$

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

New algorithm:

$t = \#$ buckets

- Take sample S
- $S_i \leftarrow S \cap B_i$
- 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 $\sum_i p_i (1+\beta)^{i-1}$

• how big is S ?

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

one of these comes from $t = O\left(\frac{\log n}{\epsilon}\right)$

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

$$\Rightarrow \overset{\text{big}}{|S_i|} = \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{c \cdot t} \geq \Omega(\text{poly} \log n \times \text{poly} \frac{1}{\epsilon})$$

\Rightarrow by union bnd + Chernoff bnd

assume this

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

Why these settings of S ? (ignore dependence on ϵ 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 " δ " 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

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

deg of any node in B_i

for all i
if $|B_i| \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{c \cdot t}$

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

else $p_i \leftarrow 0$

"small"

• Output $\sum_i p_i (1+\beta)^{i-1}$

realistic case

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

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

(note that we are assuming this for all big i
 \forall for all small i we set $p_i = 0$)

2) Can output be too small?

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

since $(1+\beta)(1-\beta) < 1$

$\geq (1-\beta) \sum_i \frac{|B_i|}{n} (1+\beta)^i$

$\geq (1-\beta) \bar{d}$

\geq deg of nodes in B_i

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

for small i ????

$$t \approx O(\log n / \epsilon)$$

$$\frac{|S_i|}{|S|} \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{1}{Ct}$$

big
or small

How much undercounting?

divide edges into 3 types

1) big-big: both endpoints in big buckets counted twice

2) big-small: one endpoint in big bucket counted once
" " " small "

3) small-small: both endpoints in small buckets never counted

big-big ok

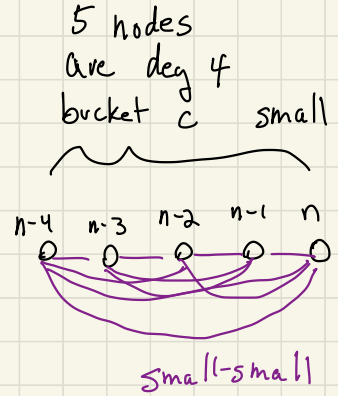
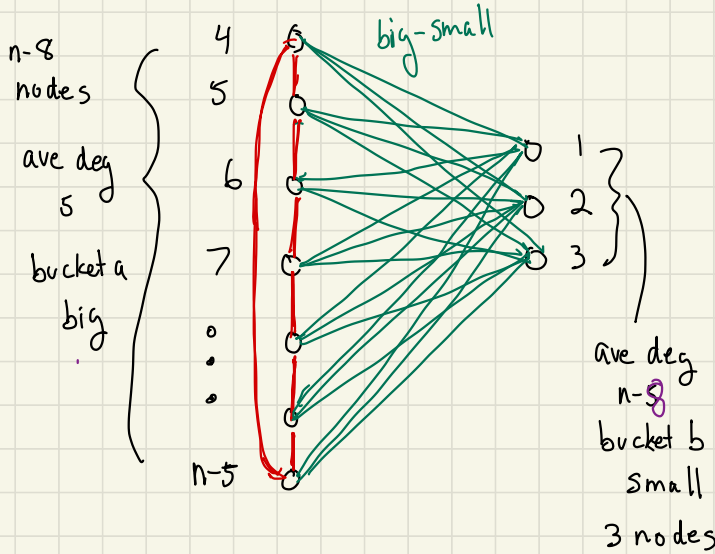
big-small undercounted by $\frac{1}{2}$

small-small not counted at all

} hope for factor 2 approx

Example:

big-big

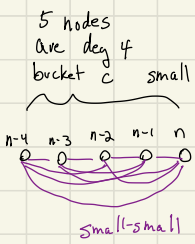
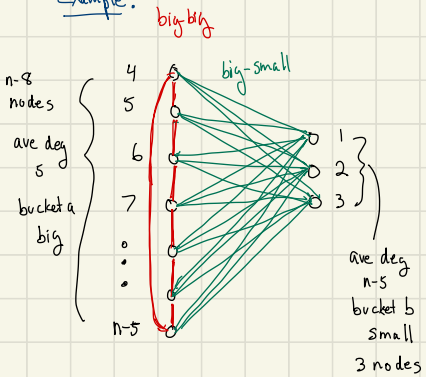


$$\text{Total degree } 5 \cdot (n-8) + (n-8) \cdot 3 + 5 \cdot 4 = 8(n-8) + 20$$

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

algorithm will likely output ≈ 5

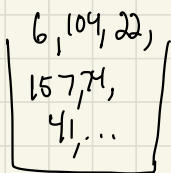
Example:



New algorithm:

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

Samples:



bucket a



bucket b



bucket c

↑
most nodes
in sample

⇒ whp bucket a
big
 $p_a \leftarrow 1$
output = 5

↙ ↘
unlikely to see any sample
whp $p_b = p_c = 0$

Good news:

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

If $|B_i| > \frac{2\sqrt{\varepsilon n}}{c_t}$ then expected size of S_i is

$$\geq |S| \cdot \frac{|B_i|}{n}$$

$$\geq |S| \cdot 2 \cdot \sqrt{\frac{\varepsilon}{n}} \cdot \frac{1}{c_t} = \text{twice the threshold for being big}$$

this is why we set threshold for big as we did

so very likely that algorithm will decide i is big
(sampling Chernoff + union bound)

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

then total # small-small edges is:

$$\leq \left(\frac{2\sqrt{\varepsilon n}}{c_t} \cdot \# \right)^2 = O(\varepsilon n)$$

recall we assumed $\bar{d} \geq 1$

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

First Claim:

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

so far, all we have used are
degree queries!

big-small error
small smaller error

Improving further:

need to improve on "big-small" edges

Main idea:

double count from the big side!

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)

1st use
of
nbr
queries!!



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

sample nodes until fall into bucket i
random nbr query from 1st node
that falls in i

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

deg queries
↓

Output $\hat{d}_i =$ average a_j

Analysis:

easy case: all nodes in B_i have same degree

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

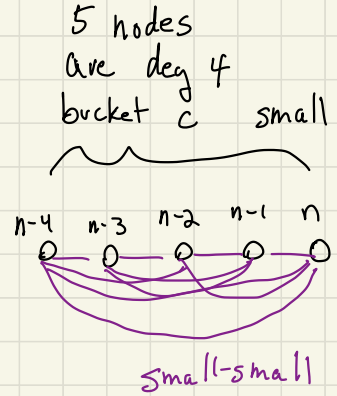
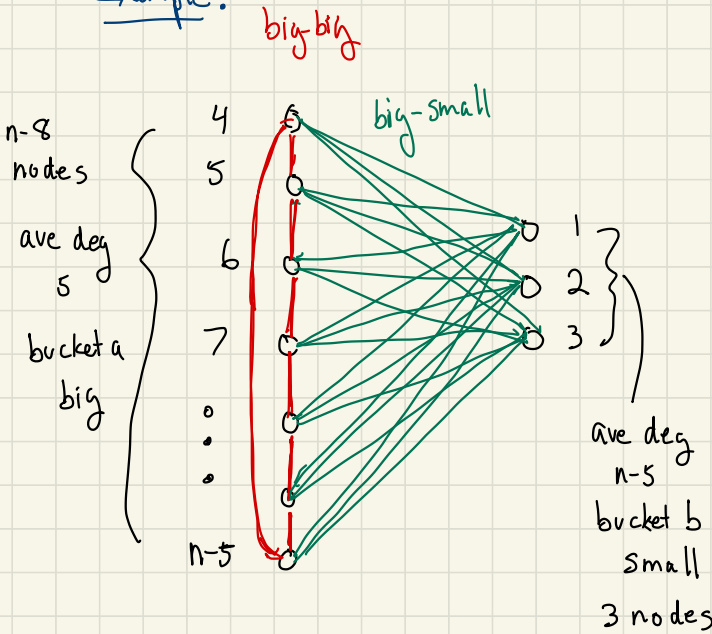
$\Pr[\text{"specific big-small edge } e_{(u,v)} \text{ in } B_i \text{ chosen}] = \frac{1}{|B_i|} \cdot \frac{1}{d}$

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

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

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

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 ≈ 5

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

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

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

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

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

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

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

est fraction
of big-small
edges hanging
off bucket 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?

estimate p_i 's } mult $1 \pm \epsilon$
 α_i 's }

\neq small-small edges } additive