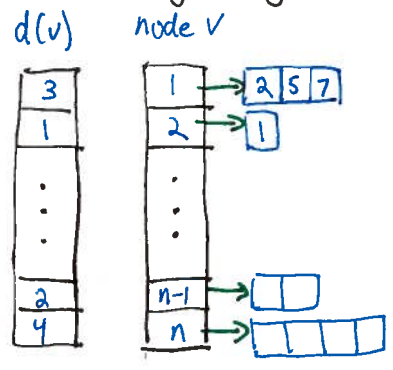


# Approximating Average Degree

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: adjacency list + degrees



- degree queries: on  $v$  return  $d(v)$
- neighbor queries: for  $(v, j)$  return  $j^{\text{th}}$  nbr of  $v$

Naive sampling:

Pick ?? sample nodes  $v_1, \dots, v_s$

output  $\frac{1}{s} \sum_i d(v_i)$  (ave degree of sample)

using straight forward Chernoff/Hoeffding  $\Rightarrow \Omega(\frac{1}{\epsilon^2})$  samples needed

Degree sequences are special?

$(n-1, 0, 0, 0, \dots, 0)$  not possible

$(n-1, 1, 1, 1, \dots, 1)$  is possible

Some lower bounds:

"Ultrasparse case":

need linear time to get any multiplicative approx

graph with 0 edges

ave deg = 0

vs.

graph with 1 edge

ave deg =  $\frac{1}{n}$



need  $\Omega(n)$  queries to distinguish

ave deg  $\geq 2$ :

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



$n - cn^{1/2}$  cycle  $\bar{d} \approx 2 + c^2$   
+  $cn^{1/2}$ -clique



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

Algorithm idea:

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

- + each group has bounded variance
- doesn't work for estimating ave of arbitrary numbers, why should it work here?

Bucketing:

set parameters

$$\beta = \epsilon / c$$

$$t = O(\log n / \epsilon) \quad \# \text{ buckets}$$

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

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

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:

- Take sample  $S$  of nodes
- $S_i \leftarrow S \cap B_i$  (samples that fall in  $i$ th bucket use degree queries to determine this)
- estimate average degree contribution from  $B_i$

using  $S_i$   
ie.  $p_i \leftarrow \frac{|S_i|}{|S|}$

note:  $\forall i$   
 $E[p_i] = E\left[\frac{|S_i|}{|S|}\right]$   
 $= \frac{|B_i|}{n}$

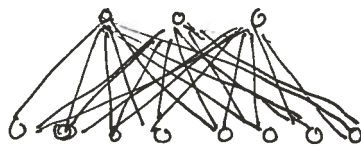
Output  $\sum_i p_i (1+\beta)^{i-1}$

Problem:

$i$  st.  $|S_i|$  is small  
likely come from  $i$  st.  $|B_i|$  small

for these, our estimate of  $|S_i|$  could be terrible

example of problem:



$\leftarrow$  3 nodes, deg  $n-3$

$\leftarrow$   $n-3$  nodes, deg 3

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

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

$|B_a| = n-3$

$|B_b| = 3$

contributes  $(n-3) \cdot 3$  edge

contributes  $3 \cdot (n-3)$  edge

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

Still, maybe good enough for 2-approximation?

Never sampled but contributes  $\frac{1}{2}$  edges !!

Next idea: use "0" for small buckets

Algorithm:

- sample  $S$
- $S_i \leftarrow S \cap B_i$
- For all  $i$

← how big?

if  $|S_i| \geq \sqrt{\frac{\epsilon}{n}} \cdot \frac{|S|}{c \cdot t}$

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

call  $i$  "big"

else

$p_i \leftarrow 0$

call  $i$  "small"

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

← so  $|S| > t \sqrt{\frac{n}{\epsilon}}$

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

$\Rightarrow |S_i| \geq \Omega(\text{polylog } n \times \text{poly } 1/\epsilon)$

Analysis:

1) Output not too large

idealistic (but unrealistic) case

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

$\leq \text{deg of node in } B_i$

realistic case

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

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

bound on sampling error when  $|S_i|$  is big (note that trivial when  $|S_i|$  not big since  $p_i \leftarrow 0$ )



2) Can output be too small?

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

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

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

$\underbrace{\frac{|B_i|}{n} (1+\beta)^i}_{\geq \text{deg of node in } B_i}$

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

For small  $i$  ????

How much undercounting?

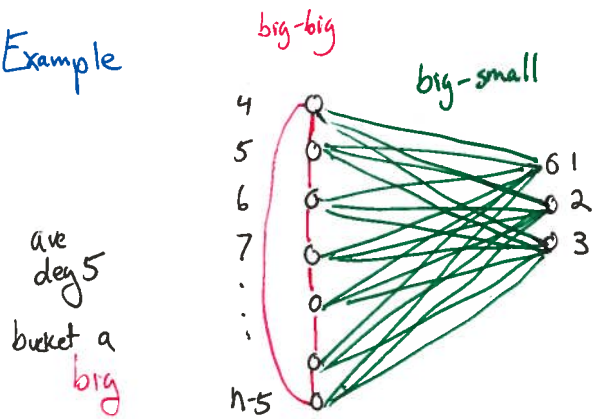
divide edges into 3 types:

- |  |   |   |               |
|--|---|---|---------------|
| type is determined by run of algorithm | } | 1) big-big - both endpts in big buckets                 | counted twice |
|  |   | 2) big-small - one endpt in big bucket<br>" " " small " | counted once  |
|  |   | 3) small-small - both endpts in small buckets           | never counted |

[see example]

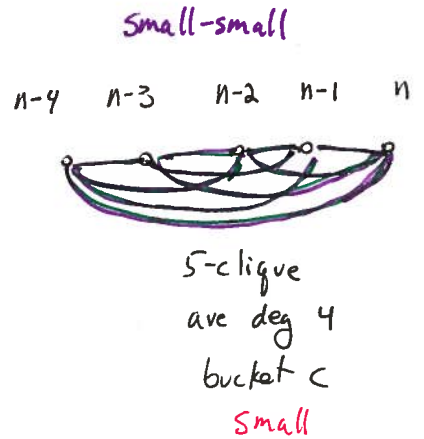
note: big-big + big-small get counted (off by factor of two)  
but small-small can be a real problem

Example



ave deg 5  
bucket a  
big

ave deg  $n-5$   
bucket b  
Small

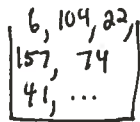


Total degree

$$5 \cdot (n-5) + (n-5)(3) + 4.5 = 8(n-5) + 20$$

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

Samples



bucket a



bucket b



bucket c

↑  
most nodes here

⇒ (whp) bucket a is big, in fact,

whp  $P_a \leftarrow 1$

output  $\approx 5$



very few nodes in these buckets  
so unlikely to see any samples

⇒ (whp) b + c are small

$P_b \leftarrow 0$     $P_c \leftarrow 0$

# big-small edges:  $3 \cdot (n-5)$

Fraction of big-big + big-small:  $\frac{3(n-5)}{5(n-5)} = 3/5$

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

$$\text{Output } 1 \cdot \underbrace{\left(1 + \frac{3}{5}\right) \left(1 + \beta\right)^4}_{\approx 5} \approx 8$$

Good news: Small buckets can't have too many nodes  
 $\Rightarrow$  can bound total # small-small edges

if  $|B_i| > \frac{2\sqrt{\epsilon n}}{c\epsilon}$  then Expected size of  $S_i$  is  $\geq |S_i| \cdot \frac{|B_i|}{n}$   
 $\geq |S_i| \cdot 2\sqrt{\frac{\epsilon}{n}} \cdot \frac{1}{c\epsilon}$

*(Red arrow pointing to  $\frac{2\sqrt{\epsilon n}}{c\epsilon}$ ):  $O(\frac{\log n}{\epsilon})$*

*(Green arrow pointing to  $\frac{1}{c\epsilon}$ ): twice the threshold for being "big"*

So very likely algorithm will decide via Chernoff bnds that  $i$  is "big"

So assume  $|B_i| \leq \frac{2\sqrt{\epsilon n}}{c\epsilon}$  for all  $i$  "small"

then total # small-small edges

$$\leq \left( \frac{2\sqrt{\epsilon n}}{c\epsilon} \cdot t \right)^2 = O\left(\frac{\epsilon n}{c^2}\right) = O(\epsilon n)$$

if we ignore them, they affect approx of  $\bar{A}$  by  $\leq (1+\epsilon)$  multiplicative factor

*(Green arrow pointing to  $\frac{2\sqrt{\epsilon n}}{c\epsilon} \cdot t$ ): # nodes / small bucket*

*(Green arrow pointing to  $t$ ): # buckets*

$\leq \epsilon n$  additive factor

*(Purple arrow pointing to  $(1+\epsilon)$ ): here we assume graph has ave degree  $\geq 1$*

First Claim:

Algorithm almost gives factor 2 mult approx

since large-small underestimated by  $\leq$  factor  $\frac{1}{2}$

we get  $(2+\epsilon)$ -multiplicative approx

*(Purple arrow pointing to  $(2+\epsilon)$ ): large-small error*

*(Purple arrow pointing to  $\epsilon$ ): small-small error*



Improving further:

need to do better on "big-small" edges ...

can we estimate the fraction of them + correct for them?

can do via sampling if we can pick a "random" edge

New queries:

random neighbor query ( $v$ ):

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

- implementation:
1. degree query to  $v$
  2. pick random  $i \in [1.. \text{deg}(v)]$
  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  $\begin{cases} 1 \\ 0 \end{cases}$

if  $e$  is "big-small"  
 o.w. ( $e$  is "big-big")

Output  $\alpha_i =$  average  $a_j$

Analysis :

Easy case : All nodes in  $B_i$  have same degree

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

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

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

$\uparrow$   
 $e=(u,v)$  only one of  $u,v$  is big since  $e$  is "big-small"

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

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

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

$\uparrow$   
estimate to  $(1+\epsilon)$ -mult factor to get

$(1+\epsilon)(1+\beta)$  estimate of  $\frac{T_i}{n}$  via  $\alpha_i \rho_i (1+\beta)^{i-1}$   
undercount of edges in  $B_i$

### Final Algorithm :

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

- $S_i \leftarrow S \cap B_i$

- For all  $i$

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

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

For all  $v \in S_i$

- Pick random nbr  $u$  of  $v$

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

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

else use  $p_i \leftarrow 0$

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

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

includes big-big  
+ one side of big-small

other side of big-small  
correction