

## Lecture 3

Sublinear time approximation  
of average degree

Warning: this is a different  
algorithm than described in  
notes for lecture 2.  
hopefully simpler!!

# Estimating Average Degree

Given  $G = (V, E)$

$\varepsilon \in (0, 1)$  approximation parameter

$\delta \in (0, 1)$  confidence

← lets assume  
 $\delta = 1/4$

Output  $\tilde{d}$  st.  $\Pr[|\tilde{d} - \bar{d}| \leq \varepsilon \bar{d}] \geq 1 - \delta$

where  $\bar{d} = \frac{m}{n}$  (average degree)

Assumption: (1)  $\bar{d} \geq 1$

(2) given access to

• "degree queries":

given  $x$  outputs  $\deg(x)$

• "neighbor queries":

given  $(v, j)$  output  $j^{\text{th}}$  nbr  
of  $v$

Last time:

Question: naive sampling needs  $\Omega(n)$  samples??

Lower bound:

Distinguish

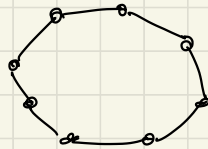
$n$ -cycle  $\bar{d} = 2$



$n - c\sqrt{n}$ -cycle  
+  $c\sqrt{n}$ -clique

$\bar{d} \approx 2 + c^2$

vs.



need  $\Omega(\sqrt{n})$  queries to distinguish?

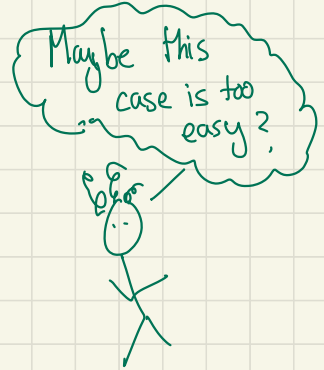
ignore notes for algorithm last time  
today will do simpler algorithm

Today:

Warm up: regular graphs

Assume each node has degree  $d$

Algorithm: output  $d$



Better warmup: almost regular graphs

Assume each node has degree in  $[\Delta, 10\Delta]$

Algorithm:

$$k \leftarrow \frac{50}{\epsilon^2} \ln(2/\delta)$$

For  $i \leftarrow 1$  to  $k$  do

• pick  $v_i \in_u V$

•  $X_i \leftarrow \deg(v_i)$

$$\text{Output } \tilde{d} \leftarrow \frac{1}{k} \sum_{i=1}^k X_i$$

notation:

$x \in_u D$   
means pick  
 $x$  uniformly  
from set  $D$

Runtime:  $O\left(\frac{1}{\epsilon^2} \ln \frac{1}{\delta}\right)$

Behavior:

Claim  $E[\tilde{d}] = \bar{d}$

Pf

$$E[\tilde{d}] = \frac{1}{k} \sum_{i=1}^k E[X_i] = E[X_1]$$

↑  
lin of  
exp

↑  
iid

$$= \sum_{v \in V} \frac{1}{n} \deg(v) = \frac{\sum \deg(v)}{n} = \bar{d} \quad \square$$

Claim  $\Pr[|\bar{d} - \tilde{d}| \leq \epsilon \bar{d}] \geq 1 - \delta$

Pf

Will use following version of Chernoff Bnd:

Thm let  $Y_1 \dots Y_k$  be independent random variables  
st.  $Y_i \in [0, 1]$  &  $Y = \sum_{i=1}^k Y_i$ . For  $b \geq 1$   
 $\Pr[|Y - E[Y]| > b] \leq 2 \cdot \exp(-2b^2/k)$

Note:  $X_i$ 's are not in  $[0,1]$  but are in  $[\Delta, 10\Delta]$  ↓ so can't use Chernoff

$$\text{let } Z_i \leftarrow \frac{X_i}{10\Delta} \text{ then } Z_i \in [0,1]$$

$$Z \leftarrow \sum_{i=1}^k Z_i \quad \tilde{d} = \frac{10\Delta}{k} \cdot Z$$

$$E[Z] = \frac{k}{10\Delta} \cdot E[\tilde{d}] = \frac{k \bar{d}}{10\Delta}$$

$$|\tilde{d} - \bar{d}| \geq \varepsilon \bar{d} \iff \left| \frac{10\Delta}{k} Z - \frac{10\Delta}{k} E[Z] \right| \geq \varepsilon \bar{d}$$

↑  
 $E[\tilde{d}]$

$$\iff |Z - E[Z]| \geq \frac{k}{10 \cdot \Delta} \cdot \varepsilon \bar{d}$$

Use Chernoff on  $Z$ 's

$$\text{with } b = \frac{k}{10\Delta} \varepsilon \bar{d}$$

$$\Pr\left[|Z - E[Z]| \geq \frac{k}{10\Delta} \varepsilon \bar{d}\right] \leq 2 \cdot e^{-\left(\frac{2k^2 \varepsilon^2 \bar{d}^2}{100 \Delta^2 \cdot k}\right)}$$

$$= 2 e^{-\frac{1}{50} \cdot \frac{k \varepsilon^2 \bar{d}^2}{\Delta^2}}$$

$$\leq 2 e^{-\frac{k \varepsilon^2}{50}}$$

$$= 2 e^{-\frac{(50/\varepsilon^2)(\ln 2/\delta) \cdot \varepsilon^2}{50}} = \delta$$

$\bar{d} \geq \Delta$  by  
assumption on all  
degrees  $\geq \Delta$



## General Case:

by Markov's  $\neq$ ,  $\leq \frac{1}{2}$  nodes have degree

$\geq c \cdot \bar{d}$ . Can we use that?

- So most nodes satisfy warmup case!

the rest of the nodes can have huge degree!!

- what about the rest?



define total order " $\prec$ " on nodes:

assume distinct IDs

def.  $u \prec v$  if

- $\deg(u) < \deg(v)$

or •  $\deg(u) = \deg(v)$

+  $ID(u) < ID(v)$

$$\deg^+(u) = \# \text{ nbrs of } u \text{ st. } u \prec v$$



orienting edges from small to large,  $\deg^+(u)$  counts "out-edges"

Observation  $\sum_{u \in V} \deg^+(u) = m = \frac{n}{2} \cdot \bar{d}$

(since each edge only counted once instead of twice as in  $\sum_u \deg(u)$ )

idea estimate average  $\left( \deg^+(u) \right)_u$

problem? we can query  $\deg(u)$   
not  $\deg^+(u)$

benefit:

Lemma  $\forall v \in V \quad \deg^+(v) \leq \sqrt{m}$

Proof

define  $H \subseteq V$  to be  $\sqrt{m}$  nodes  
with highest rank (degree) wrt.  $\alpha$

$\forall v \in H, \deg^+(v) \leq \sqrt{m}$  since  
edges "leaving"  $v$  go to bigger nodes



(which must also be in  $H$ )

$\forall v \in V \setminus H, \deg^+(v) \leq \deg(v) \leq \sqrt{2m}$ :

Why? if not,  $\deg(v) > \sqrt{2m}$  ← assume for contradiction

but all  $w$  in  $H$  have

$$\deg(w) \geq \deg(v) > \sqrt{2m}$$

so total degree

$$> \underbrace{|H| \cdot \sqrt{2m}}_{\text{contribution from } H} + \underbrace{\text{something positive}}_{\text{contribution from } V \setminus H}$$

$$> \sqrt{2m} \cdot \sqrt{2m} = 2 \cdot m$$

but sum of degrees =  $2 \cdot m$

Algorithm:

$$K \leftarrow \frac{16}{\epsilon^2} \sqrt{n}$$

for  $i=1$  to  $K$

pick  $v_i \in_r V$  (1)


pick  $u_i \in_r N(v_i)$  (2)

if  $u_i \not\sim v_i$  then  $X_i \leftarrow 2 \deg(v_i)$

else  $X_i \leftarrow 0$

$$\text{return } \tilde{d} = \frac{1}{K} \sum_{i=1}^K X_i$$

→ ←

Symbol for "contradiction" 

Question to think about:  
why the "2"?

Claim  $E[X_i] = \bar{d}$

PF

$$E[X_i] = \sum_{v \in V} \Pr[v \text{ chosen in (1)}] \cdot E[X_i \mid v \text{ chosen in (1)}]$$

$$= \sum_{v \in V} \frac{1}{n} \cdot E[X_i \mid v \text{ chosen in (1)}]$$

$$= \frac{1}{n} \sum_{v \in V} \sum_{u \in N(v)} \Pr[u \text{ chosen in (2)} \mid v \text{ chosen in (1)}]$$

$$\times E[X_i \mid u \text{ chosen in (2)} + v \text{ chosen in (1)}]$$

$$= \frac{1}{n} \cdot \sum_{v \in V} \sum_{\substack{u \in N(v) \\ v \neq u}} \frac{1}{\deg(v)} \cdot 2 \cdot \deg(v)$$

if  $v \neq u$   
then  
 $2 \deg(v)$   
else 0

$$= \frac{2}{n} \cdot \sum_{v \in V} \deg^+(v) = \frac{2m}{n} = \bar{d}$$

▀

But how many samples do we need to assure that we are close to expectation? Here is where we use graph properties!

Claim  $\text{Var}[X_i] \leq 4\sqrt{2m} \bar{d}$

Pf  $\text{Var}[X_i] = E[X_i^2] - E[X_i]^2 \leq E[X_i^2]$  ↘ as above

$$= \frac{1}{n} \sum_{v \in V} \sum_{\substack{u \in N(v) \\ v < u}} \frac{1}{\deg(v)} \underbrace{(2 \deg(v))^2}_{X_i^2}$$

$$= \frac{4}{n} \sum_{v \in V} \underbrace{\deg^+(v)}_{\leq \sqrt{2m}} \cdot \deg(v)$$

← key insight

$$\leq \frac{4}{n} \cdot \sqrt{2m} \sum_{v \in V} \deg(v)$$

$$\leq 4 \cdot \sqrt{2m} \cdot \bar{d}$$



2 useful facts about variance!

• Lemma let  $Y = \frac{1}{k} \sum_{i=1}^k X_i$  where  $X_i$ 's are iid

then  $\text{Var}[Y] = \frac{1}{k} \text{Var}[X]$

so can reduce variance by sampling averaging more!

important but pairwise independence is good enough

• Chebyshev's  $\neq$ :  $\Pr[|X - E[X]| \geq b] \leq \frac{\text{Var}[X]}{b^2}$

Lemma  $\Pr [ |\tilde{d} - \bar{d}| \leq \varepsilon \bar{d} ] \geq 3/4$

Pf

$$E[\tilde{d}] = \bar{d} \quad \text{by lin of expectation}$$

$$\text{Var}[\tilde{d}] \leq \frac{4 \cdot \sqrt{2m}}{k} \cdot \bar{d}$$

since  $\bar{d} = E[\tilde{d}]$

$$\Pr [ |\tilde{d} - \bar{d}| \geq \varepsilon \bar{d} ] = \Pr [ |\tilde{d} - E[\tilde{d}]| \geq \varepsilon \bar{d} ]$$

$$\leq \frac{\text{Var}[\tilde{d}]}{(\varepsilon \bar{d})^2}$$

$$\leq \frac{\frac{4 \sqrt{2m}}{k} \cdot \bar{d}}{\varepsilon^2 \bar{d}^2} = \frac{4 \sqrt{2m}}{\varepsilon^2 \bar{d} \cdot k}$$

$$= \frac{4 \sqrt{2m} \cdot n}{\varepsilon^2 \cdot 2m \cdot k} = \frac{4n}{\varepsilon^2 \sqrt{2m} \cdot k}$$

$$= \frac{\sqrt{n}}{4 \cdot \sqrt{2m}}$$

pick  $k = \frac{16}{\varepsilon^2} \sqrt{n}$

$$\leq \frac{1}{4}$$

since  $\sqrt{\frac{n}{2m}} = \sqrt{\frac{1}{\bar{d}}}$   
 $\leq 1$  since  
we assumed  $\bar{d} \geq 1$

How do we improve probability of success?

See HW 0!