# Lecture 2

Topics:

- Sublinear time approximation of average degree

- Estimate number of connected components
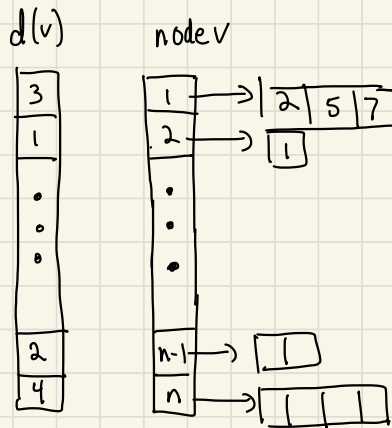
## Estimating the average degree of a graph

<u>def</u>  Average degree   $\bar{d} = \dfrac{\sum_{u \in V} \deg(u)}{n} = \dfrac{2m}{n}$

Assume:   $G$ simple   (no parallel edges, self-loops)

$\qquad \Omega(n)$ edges   (not "ultra-sparse")

Representation via adj list + degrees:



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

# Estimating Average Degree

Given   $G = (V, E)$

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

$\delta \in (0,1)$ confidence          $\leftarrow$ lets assume

$\delta = \frac{1}{4}$

Output   $\tilde{d}$  s.t.  $\Pr\left[ |\tilde{d} - \bar{d}| \leq \varepsilon\,\bar{d} \right] \geq 1 - \delta$

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

Last time we saw that "naive sampling" ie.

Pick $O(??)$ sample nodes $v_1 .. v_s$

Output ave degree of sample:

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

Does not work so well, although we did prove that if <u>all</u> $deg(v)$ are in $[\Delta, 10\Delta)$ then constantly many samples are sufficient, for the naive sampling algorithm

In general, we saw a handwavy argument that

- $\Omega(n)$ time is need to give a multiplicative estimate for average degree (this used "ultrasparse" graphs)

- $\Omega(\sqrt{n})$ time is needed for estimating average degree, even when the average degree is $> 1$.

Today we will see the general case, & a different algorithm

<u>General    Case</u> :   "Order" edges   to   control   outdegree

Our plan:
    define   total order  "$\prec$"   on   nodes:
                                     $\curvearrowright$ assume  distinct
                                              ID's

      <u>def</u>. $u \prec v$    if

$$\circ \quad \deg (u) < \deg (v)$$
$$\text{or} \quad \bullet \quad \deg (u) = \deg (v)$$
$$+ \quad ID (u) < ID (v)$$

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



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

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

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

<u>idea</u> estimate $\underset{u}{average} \left( deg^+(u) \right)$ $= \frac{1}{2} \cdot \bar{d}$
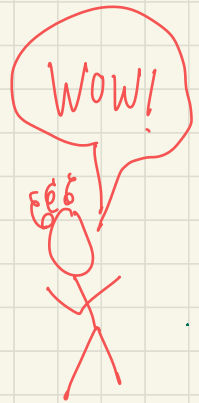
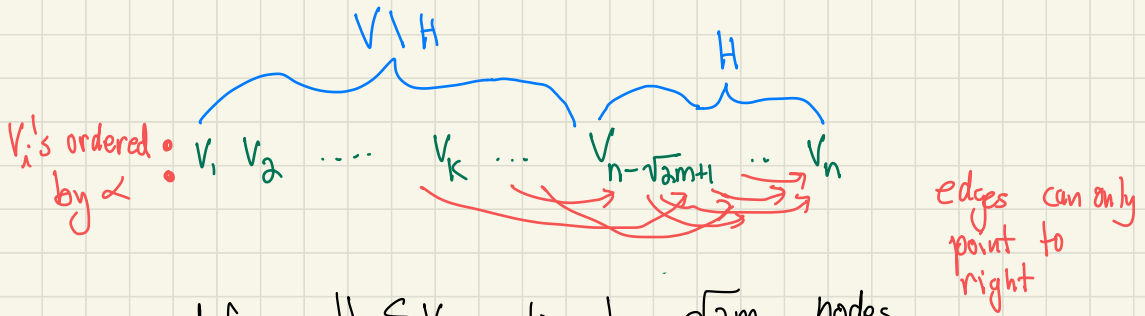problem? we can query $deg(u)$
not $deg^+(u)$

benefit:

<u>Lemma</u> $\forall v \in V$ $deg^+(v) \leq \sqrt{2m}$

<u>Proof</u>

Consider order of $v$'s by $\alpha$:

$V \setminus H$     $H$

$V_i$'s ordered by $\prec$ • $V_1$ $V_2$ .... $V_k$ ... $V_{n-\sqrt{2m}+1}$ .. $V_n$

edges can only point to right

define $H \subseteq V$ to be $\sqrt{2m}$ nodes
with highest rank (degree) wrt. $\prec$

**heavy nodes:**
$\forall v \in H, \deg^+(v) \le \sqrt{2m}$ since
edges "leaving" $v$ go to bigger nodes
(which must also be in $H$)

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

Why?

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

but all $w$ in $H$ have
$\deg(w) \ge \deg(v) > \sqrt{2m}$

so total degree $\sum\limits_{v} d(v)$

$> |H| \cdot \sqrt{2m} + $ something positive

$\underbrace{\qquad\qquad}$ contribution from H  $\underbrace{\qquad\qquad}$ contribution from $V \backslash H$

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

but sum of degrees $= 2 \cdot m$

$\rightarrow \leftarrow$

$\underbrace{\qquad}$ Symbol for "contradiction"

Algorithm:

$K \leftarrow \dfrac{16}{\varepsilon^2} \sqrt{n}$

for $i = 1$ to $K$

    pick $v_i \in_r V$       (1)

    pick $u_i \in_r N(v_i)$    ← neighbor of $v$     (2)

    if $v_i < u_i$ then $X_i \leftarrow 2 \deg(v_i)$

           else $X_i \leftarrow 0$

return $\tilde{d} = \dfrac{1}{K} \sum\limits_{i=1}^{K} X_i$

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

Pf

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

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

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

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

$= \frac{1}{n} \cdot \sum_{v \in V} \sum_{\substack{u \in N(v) \\ \& v \alpha u}} \frac{1}{\deg(v)} \cdot 2 \cdot \deg(v)$   ⏜ if $v \alpha u$
then
$X_i = 2 \deg(v)$
else $X_i = 0$

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

But how many samples do we need to
assure that we are close to

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

**Pf**   $\text{Var}[X_i] \equiv E[X_i^2] - E[X_i]^2 \leq E[X_i^2]$   $\Big\}$ as above

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

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

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

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

2 useful facts about variance!

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

  then   $\text{Var}[Y] = \frac{1}{K} \text{Var}[X]$   $\Big\}$ important

  *so can reduce variance by sampling + averaging more!*

  but pairwise independence is good enough

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

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

Pf

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

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

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

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

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

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

$\underset{\parallel}{\phantom{=}}$

$\dfrac{2m}{n}$

$= \dfrac{4\sqrt{2m}\cdot n}{\varepsilon^2\cdot2m\cdot k} = \dfrac{4n}{\varepsilon^2\cdot\sqrt{2m}\cdot k}$

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

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

$\leq \dfrac{1}{4}$    since $\sqrt{\dfrac{n}{2m}} = \sqrt{\dfrac{1}{\bar{d}}}$

       $\leq 1$   since
       we assumed $\bar{d}\geq1$

$\Rightarrow$ good estimate with prob $\geq 3/4$

How do we improve probability of success?

      See    HW 0!