

Week 3: CS5234

Last Week:

- 1) Connected components
- 2) Approx-MST
- 3) Divisibility

Today:

- 1) Diameter
- 2) Avg. degree

Diameter

→ Graph $G = (V, E)$, undirected
n nodes, m edges

One solution: $\forall u \in V$ do: BFS(u)
↳ $O(n(n+m))$ time

Note: different!
 ϵn , not ϵn^2

Defn: G is ϵ -far from G' if must add/delete
 $\lceil \epsilon n \rceil$ edges to G to get G' .

Goal: Given G, D , decide:

- 1) if G has diameter $\leq D$, return **TRUE**
- 2) if G is ϵ -far from graph with diameter $4D+2$, return **FALSE**
- 3) otherwise, either

Surprising! Diameter is a global property

Goal: $O(1/\epsilon^3)$ time

(n, D don't matter!!)

Trick: find a local property to test

Defn: Node $u \in V$ is (K, C) -friendly if there are $\geq K$ nodes within distance C of u .

Key lemma: if $\geq (1 - 1/K)n$ nodes in G are (K, D) -friendly, then G is $2/K$ -close to a graph with diameter $4D + 2$.

Proof by algorithm: Label each node
center
friend
fringe

Repeat:

- 1) Choose unlabelled node u .
- 2) Label u a center
- 3) Label nodes w s.t. $d(u, w) \leq D$ as friend.
- 4) Label unlabelled w s.t. $d(u, w) \leq 2D$ as fringe.
- 5) Delete u and all friends from G

Eventually, all nodes are labelled. How many centers?

For each center, 2 cases:

1) center u is not (K, D) -friendly $\Rightarrow \leq n/K$ centers

2) Center u is (K, D) -friendly.

- a) u is not distance $2D$ from any previously chosen center
(or it would have already been labelled fringe)
- b) initially, u had $\geq K$ nbs at distance $\leq D \Rightarrow$ call them F
(by defn of friendly)
- c) Nodes in F were not previously labelled center/friend.
(if $w \in F$ was friend, then \exists center v s.t. $d(u, v) \leq D$,
so $d(u, w) \leq d(u, v) + d(v, w) \leq 2D$)
- d) All nodes in F labelled friends, deleted
 $\hookrightarrow K$ nodes deleted

Conclusion: $\leq n/K$ non-friendly centers chosen

Total: $\leq 2n/K$ centers chosen

Proof (continued):

1) Choose a center u . Call it C .

2) Add $2^n/k - 1$ edges connecting all centers to u .

Claim: diameter $\leq 4D+2$

→ all nodes labelled \Rightarrow all within $2D$ of a center

$$\rightarrow d(u, v) \leq \underbrace{d(u, C_1)}_{\substack{\uparrow \\ \text{ctr for } u}} + d(C_1, C) + d(C, C_2) + \underbrace{d(C_2, v)}_{\substack{\uparrow \\ \text{ctr for } v}}$$

$$\leq 2D + 1 + 1 + 2D = 4D + 2$$

$\Rightarrow G$ was $2^n/k$ close to a graph (constructed) with diam $\leq 4D+2$

Alg

Diam(G, D, ϵ):

Repeat S times:

Choose u at random.

Do BFS to see if u is $(2^n/\epsilon, D)$ -friendly

If not friendly, return FALSE.

Return TRUE.

1) if diam $\leq D$ and $n \geq 2^n/\epsilon$, then all nodes friendly. \Rightarrow TRUE

2) if $n < 2^n/\epsilon$, then can connect by $< 2^n/\epsilon \leq 2^n/\epsilon$ edges \Rightarrow special case

3) if not ϵ -close to $4D+2$, then $\geq \epsilon^n/2$ nodes not $(2^n/\epsilon, D)$ -friendly
 $k = 2^n/\epsilon$

$$\Pr(u \text{ is not } (2^n/\epsilon, D)\text{-friendly}) \geq (\epsilon^n/2) / n \geq \epsilon/2$$

$$\Pr(\text{no sampled } u \text{ is not } (2^n/\epsilon, D)\text{-friendly}) \leq (1 - \epsilon/2)^S \leq e^{-\epsilon S/2} \leq 1/3$$

\uparrow
independent

Choose $S = 4/\epsilon$

Conclusion: if ϵ -far from $4D+2$, the FALSE w.p. $\geq 2/3$

Cost: $S = 4/\epsilon$ iterations

Cost/iteration: BFS to find $2/\epsilon$ nodes in distance D

$\hookrightarrow \text{Cost} \leq \left\lceil \frac{2}{\epsilon} \right\rceil^2 \leftarrow \begin{array}{l} \leq 2/\epsilon \text{ nodes found} \\ \leq 2/\epsilon \text{ nbrs of each} \\ \text{node visited} \\ \text{unnecessarily} \end{array}$

Total time: $(4/\epsilon) \left\lceil \frac{2}{\epsilon} \right\rceil^2 = O(1/\epsilon^3)$

Degree estimation

$G = (V, E)$, undirected, connected

Alg

AvgDeg(G, ϵ)

$\min = n$

repeat K times

total = 0

Sampling loop

repeat S times:

choose $u \in V$ at random

total = total + degree(u)

if total/ $S < \min$ then

$\min = \text{total}/S$

return \min

Seems unlikely to work:

$A = \{n, n, n, n, 0, 0, \dots, 0\}$

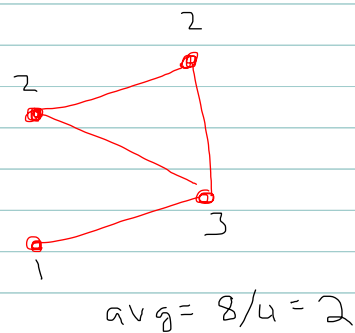
avg = 4

Sampling $\Rightarrow 0$

Sampling does not
find average

Key insight: a graph is not an arbitrary array
 it has structure Cannot have degrees $[n, n, n, n, 1, 1, 1, \dots, 1]$

Setup n nodes
 m edges
 $\text{avg degree} = \frac{1}{n} \sum_u \deg(u) = \frac{1}{n} [2m]$



Fix ϵ

Goal: $(1/2 - \epsilon)$ -approx in
 $O\left(\frac{\sqrt{n}}{\epsilon^{3/2}}\right)$ time

Step 1: Show answer $\leq (1 + \epsilon)d$

Markov's Inequality!!

For one sampling loop: let $X_j = \deg$ sample j

$$X = \sum X_j = \text{total}, \quad E[X_j] = \frac{1}{n} \sum_u \deg(u) = \frac{2m}{n} = d$$

$$E[X] = sd$$

$$\Pr\left[\frac{X}{s} \geq (1 + \epsilon)d\right] = \Pr[X \geq (1 + \epsilon)E[X]] \leq \frac{E[X]}{(1 + \epsilon)E[X]}$$

$$\leq \left(\frac{1}{1 + \epsilon}\right) \leq 1 - \epsilon/2$$

$$\text{Recall: if } X < 1: (1 + X)^{-1} \geq 1 - X$$

$$\leq 1 - X/2$$

$$\text{Taylor: } (1 + X)^{-1} = 1 - X + X^2/2 - \dots$$

$$\text{Conclusion: } \Pr\left[\frac{X}{s} > (1 + \epsilon)d\right] \leq 1 - \epsilon/2$$

Too big? [cont]

$$\Pr(\text{all } K \text{ sample loops } > (1+\varepsilon)d) \leq (1-\varepsilon/2)^K$$

↑
independent

choose $K = 8/\varepsilon$

$$\leq (1-\varepsilon/2)^{8/\varepsilon}$$

$$\leq e^{-4} \leq 1/8$$

Conclusion: w.p. $\geq 7/8$, answer is $\leq (1+\varepsilon)d$

What about showing answer is not $\leq (1-\varepsilon)d$??

How do you find heavy nodes?

If miss heavy nodes, then answer is too small!

Key problem: heavy (large deg) values have big impact on final result but are hard to find.

Solution: "Ignore" heavy nodes

Defn

$H = \sqrt{\varepsilon n}$ nodes with largest deg

L = all other nodes

M = max degree in L .

Question: if we only sample L , is that good enough?

(How to sample only L ?? Postponed...)

Defn $\deg(H) = \sum_{u \in H} \deg(u)$

$$\deg(L) = \sum_{u \in L} \deg(u)$$

$$\deg(G) = \sum_u \deg(u) = 2m = \deg(H) + \deg(L)$$

Claim $\deg(L) \geq m(1-\epsilon)$

$$\deg(H) \leq 2 \left(\sqrt{\epsilon n} \right) + m \leq 2 \left\lceil \sqrt{\epsilon n} \right\rceil + m \leq \epsilon n + m \leq \epsilon m + m \leq (1+\epsilon)m$$

\uparrow edges between L and H

$n \leq m$

$$\deg(L) = \deg(G) - \deg(H) \leq 2m - m(1+\epsilon) = m - \epsilon m = m(1-\epsilon)$$

Conclusion: if we estimate avg deg L , we are within $(\frac{1}{2}-\epsilon)$

of avg deg G [i.e., $\frac{m}{n}(1-\epsilon)$ is $\geq \frac{2m}{n}(\frac{1}{2}-\epsilon)$]

How to estimate $\deg(L)$?

(Don't know which are in L .)

Trick: define new random vars.

$$y_i = \min(X_i, M)$$

\leftarrow max value in L

if $x_j \in H$, $y_j = M$

if $x_j \in L$, $y_j = x_j$

$$x_j \geq y_j$$

New goal: bound Y

$$X \geq Y = \sum y_j$$

$$\begin{aligned}
E[y_j] &= \frac{1}{n} \sum_j \min(\deg(j), M) \\
&= \frac{1}{n} \sum_{u \in H} \min(\deg(u), M) + \frac{1}{n} \sum_{u \in L} \min(\deg(u), M) \\
&\geq \frac{1}{n} \sum_{u \in L} \min(\deg(u), M) \\
&\geq \frac{1}{n} \sum_{u \in L} \deg(u) \\
&\geq \frac{1}{n} \deg(L) \\
&\geq \frac{m}{n} (1-\varepsilon) \\
&\geq (d/2) (1-\varepsilon)
\end{aligned}$$

$$\begin{aligned}
E[y_j] &\geq \frac{1}{n} \sum_{u \in H} \min(\deg(u), M) \\
&\geq \frac{1}{n} \cdot |H| \cdot M \\
&\geq \frac{1}{n} (\sqrt{n\varepsilon}) M \geq \frac{M}{\sqrt{n/\varepsilon}}
\end{aligned}$$

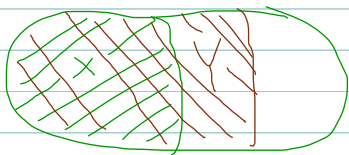
Show: $\Pr[Y \leq (1-\delta) E[Y]]$ is small, i.e., $\leq 1/8K$

$$\Rightarrow \Pr[X \leq S(1-\delta)(d/2)(1-\varepsilon)] \leq \Pr[Y \leq S(1-\delta)(d/2)(1-\varepsilon)]$$

$$\uparrow Y \leq X$$

$$\leq \Pr[Y \leq (1-\delta) E[Y]] \leq 1/8K$$

$$\uparrow E[Y] \geq (d/2)(1-\varepsilon) S$$



$$\Pr[X/S \leq (1-\delta)(d/2)(1-\varepsilon)] \leq 1/8K$$

Break: Chernoff Bounds

z_1, z_2, \dots, z_n are r.v. $\in [0, 1]$

$$Z = \sum z_i$$

$$E[Z] = M$$

$$0 \leq \varepsilon \leq 1$$

$$1) \Pr(Z \geq (1+\varepsilon)M) \leq e^{-\varepsilon^2 M/3}$$

$$2) \Pr(Z \leq (1-\varepsilon)M) \leq e^{-\varepsilon^2 M/3}$$

$$3) \Pr(|Z-M| \leq \varepsilon M) \leq 2e^{-\varepsilon^2 M/3}$$

Wait!! What if $z_i \in [1, K]$?

$$1) \Pr(Z \geq (1+\varepsilon)M) \leq e^{-\varepsilon^2 M/3K}$$

$$2) \Pr(Z \leq (1-\varepsilon)M) \leq e^{-\varepsilon^2 M/3K}$$

$$3) \Pr(|Z-M| \leq \varepsilon M) \leq 2e^{-\varepsilon^2 M/3K}$$

$$\Pr[Y \leq (1-\delta)E[Y]] \leq e^{-E[Y]\delta^2/2M}$$
$$\leq e^{-sM/\sqrt{M}\delta^2} \quad \leftarrow E[Y] \geq sM/\sqrt{M}\delta$$

$$\text{Choose } s = \frac{8\sqrt{M}\delta^2 \ln K}{\delta^2}$$

$$\leq e^{-4\ln K} \leq 1/8K$$

$$\Rightarrow \Pr\left[\frac{X}{s} \leq (1-\delta)(1-\varepsilon)\left(\frac{d}{2}\right)\right] \leq 1/8K$$

$$\Pr[\text{all iterations good}] \geq \left(1 - \frac{1}{8K}\right)^K \geq e^{-1/8} \geq 7/8$$

$$\Pr(\text{some iteration bad OR answer too big}) \leq 1/8 + 1/8 \leq 1/4$$

$$\implies \text{correct v.p.} \geq 3/4$$

$$\text{Set } \delta = \varepsilon: \text{ correct} = \quad 1) \text{ answer} \leq (1 + \varepsilon) d$$

$$2) \text{ answer} \geq (1 - \delta)(1 - \varepsilon) \left(d/2\right)$$

$$\geq (1 - \varepsilon)^2 \left(d/2\right) \quad (1 - \varepsilon)^2 = (1 - 2\varepsilon + \varepsilon^2) \geq 1 - 2\varepsilon$$

$$\geq (1 - 2\varepsilon) \left(d/2\right)$$

$$\geq d \left[1/2 - \varepsilon\right]$$

$$d \left[1/2 - \varepsilon\right] \leq \text{answer} \leq d(1 + \varepsilon)$$

$$\text{Time: } O(KS)$$

$$K = \frac{8}{\varepsilon} \quad S = \frac{8 \sqrt{n/\varepsilon} \ln(8/\varepsilon)}{\varepsilon^2}$$

$$O(KS) \leq O\left(\sqrt{n/\varepsilon} \cdot \frac{1}{\varepsilon^4}\right) = O\left(\sqrt{n} \varepsilon^{-9/2}\right)$$