

Lecture 4

- Estimating the number of connected components in a graph
- Estimating weight of min spanning tree
- Distributed Algorithms vs. sublinear time algorithms

How many connected components in G ?

Input $G = (V, E)$, ε adjacency list representation $|V| = n$
max degree d $|E| = m \leq d \cdot n$

Output y st. if $C = \#$ conn comp of G

then $C - \varepsilon n \leq y \leq C + \varepsilon n$

← additive approx to w in εn

A different way to characterize $\#$ conn. comp.:

notation $\forall v$, let $n_v = \#$ nodes in v 's connected component

Observation \forall connected component $A \subseteq V$

$$\sum_{u \in A} \frac{1}{n_u} = \sum_{u \in A} \frac{1}{|A|} = 1$$

$$\Rightarrow \text{characterization: } C = \sum_{u \in V} \frac{1}{n_u}$$

Why is this better?

Computing $C = \Omega(n^2)$ time? maybe $\Omega(n)$ if careful?

but we will only estimate C

Estimating $c = \sum_{u \in V} \frac{1}{n_u}$:

Two parts:

1) estimating $\frac{1}{n_u}$ quickly

2) estimating $\sum_u \frac{1}{n_u}$ by sampling u 's

↳ estimating the average $\frac{1}{n_u}$ value

Estimating a single $\frac{1}{n_u}$ value:

def let $\hat{n}_u = \min \{n_u, 2/\epsilon\}$

$$\hat{c} = \sum_{u \in V} \frac{1}{\hat{n}_u}$$

Lemma $\forall u, \left| \frac{1}{\hat{n}_u} - \frac{1}{n_u} \right| \leq \epsilon/2$

Pf

if $n_u \leq 2/\epsilon$ then $\hat{n}_u = n_u$ so $\frac{1}{\hat{n}_u} - \frac{1}{n_u} = 0$

if $n_u > 2/\epsilon$ then $\hat{n}_u = 2/\epsilon$

so $0 \leq \frac{1}{n_u} - \frac{1}{\hat{n}_u} = \frac{\epsilon}{2}$

so $\left| \frac{1}{\hat{n}_u} - \frac{1}{n_u} \right| \leq \frac{\epsilon}{2}$

□

Corr. $|C - \hat{C}| \leq \frac{\epsilon n}{2}$

How long does it take to compute \hat{n}_u ?

Algorithm compute \hat{n}_u

Do BFS starting from u until

- visit whole component of u
- or visit $2/\epsilon$ distinct nodes

Output # visited nodes.

runtime $O(d/\epsilon)$ ← time per step of BFS

Estimating sum:

Algorithm estimate c

$$r \leftarrow b/\epsilon^3$$

choose $U = \{u_1, \dots, u_r\}$ random nodes

$\forall u \in U,$

compute \hat{n}_u via above algorithm

Output $\hat{C} = \frac{n}{r} \sum_{u \in U} \frac{1}{\hat{n}_u}$

runtime $O(d/\epsilon^4)$

why is it a good estimate?

Thm $\Pr [|\tilde{c} - c| \leq \epsilon n / 2] \geq 3/4$

Pf.

Chernoff Bnd: X_1, \dots, X_r iid $X_i \in [0, 1]$
 $S = \sum X_i$ $p = E[X_i] = E[S] / r$

Then:
 $\Pr [|S/r - p| \geq \delta p] \leq e^{-\Omega(r p \delta^2)}$

here, $p = E[\frac{1}{n_{u_i}}]$
 $S = \sum_{i=1}^r \frac{1}{n_{u_i}}$
 $\delta = \epsilon / 2$

So $\Pr [|\frac{1}{r} \sum_{i=1}^r \frac{1}{n_{u_i}} - E[\frac{1}{n_{u_i}}]| \geq \frac{\epsilon}{2} E[\frac{1}{n_{u_i}}]] \leq e^{-\Omega(r \cdot \frac{\epsilon^2}{4} \cdot E[\frac{1}{n_{u_i}}])}$

$= \frac{\tilde{c}}{n}$ $= \frac{1}{n} \cdot \sum_{u \in V} \frac{1}{n_u}$ $= \frac{\hat{c}}{n}$

since $\frac{\epsilon}{2} \leq \frac{1}{n_u} \leq 1$
 so $\frac{\epsilon n}{2} \leq \hat{c} \leq n$
 so can make this a large constant
 $-\Omega(\frac{1}{\epsilon} \cdot \frac{\hat{c}}{n})$

So $\Pr [|\frac{\tilde{c}}{n} - \frac{\hat{c}}{n}| \geq \frac{\epsilon}{2} \frac{\hat{c}}{n}]$
 $= \Pr [|\tilde{c} - \hat{c}| \geq \frac{\epsilon}{2} \hat{c}] \leq e$

$\leq \frac{1}{4}$ by picking good constant in Ω

Corr $\Pr [|c - \tilde{c}| \leq \epsilon n] \geq 3/4$

Pf. $|c - \tilde{c}| \leq |c - \hat{c}| + |\hat{c} - \tilde{c}|$
 always $\leq \epsilon n / 2$ by corr $\leq \epsilon n / 2$ by thm with prob $\geq 3/4$

What is the minimum spanning tree weight?

Input: $G = (V, E)$ adj list representation $n = |V|$
 max degree d
 each edge has weight $w_{uv} \in \{1, \dots, w\} \cup \{\infty\}$
 G is connected (via edges of weight $\leq w$)
 i.e. $w_{uv} \notin E$

Output let $M = \min_{T \text{ spans } G} w(T)$
 \uparrow
 $\sum_{(i,j) \in T} w_{ij}$

output \hat{M} s.t.
 $(1-\epsilon)M \leq \hat{M} \leq (1+\epsilon)M$

Note by assumption of wts, $n-1 \leq w(T) \leq w(n-1)$

A different characterization of M :

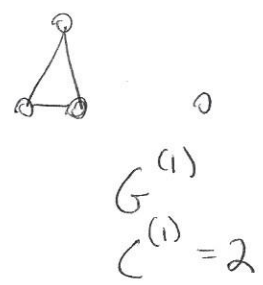
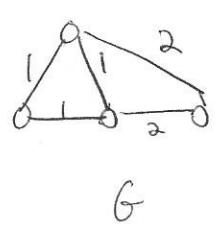
def $G^{(i)} = (V, E^{(i)})$ where $E^{(i)} = \{(u,v) \mid w_{uv} \in \{1, \dots, i\}\}$
 $C^{(i)} = \#$ conn comp of $G^{(i)}$

before characterizing, let's do examples:

1) $w=1$ (only one weight + connected by assumption)

then, $M = n-1$

2) $w=2$



need $\cong C^{(1)} - 1$ edges of wt 2 for MST
and this is enough!

$$\text{so } M = 2(C^{(1)} - 1) + 1 \cdot (n - 1 - (C^{(1)} - 1))$$

$$= n - 2 + C^{(1)}$$

In general:

Claim $M = n - w + \sum_{i=1}^{w-1} C^{(i)}$

Pf

let $\alpha_i = \#$ edges of wt i in any MST of G

all MST's have same value of α_i !
why?

$$\sum_{i>l} \alpha_i = \# \text{ conn comp of } G^{(l)} - 1$$

$$= C^{(l)} - 1$$

where $C^{(0)} = n$

$$M = \sum_{i=1}^w i \cdot \alpha_i$$

$$= \sum_{i=1}^w \alpha_i + \sum_{i=2}^w \alpha_i + \sum_{i=3}^w \alpha_i + \dots + \underbrace{\sum_{i=w}^w \alpha_i}_{= \alpha_w}$$

$$= (n-1) + (C^{(1)} - 1) + (C^{(2)} - 1) + \dots + (C^{(w-1)} - 1)$$

$$= n - w + \sum_{i=1}^{w-1} C^{(i)}$$

□

Approximation Algorithm:

For $i = 1$ to $w-1$

$\hat{C}^{(i)}$ = approx # c.c. of $G^{(i)}$ to within $\frac{\epsilon}{2w}$ (additive error)

Output $\hat{M} = n - w + \sum_{i=1}^{w-1} \hat{C}^{(i)}$

Runtime:

$$O(d / (\epsilon')^4) = O(dw^4 / \epsilon^4) \quad \text{for each call to approximate \# C.C.}$$

↑
we used $\epsilon' = \frac{\epsilon}{2w}$

↑
how do you recompute $G^{(i)}$?
ignore edges of wt $> \epsilon$

Total: $O(dw^5 / \epsilon^4)$

can improve to $O\left(\frac{dw}{\epsilon^2} \log \frac{dw}{\epsilon}\right)$

+ need $\Omega\left(\frac{dw}{\epsilon^2}\right)$

Approximation guarantee:

how?
↓ HWO!

Call approx # cc with error probability $\leq \frac{1}{4 \cdot w}$

is. prob of approx error $> \epsilon'$

Pr [all calls to approx # cc give output that is ϵ' additive approx] $\geq 1 - \frac{w}{4w} = 3/4$

If ↑ happens,

$$|M - \hat{M}| \leq n \cdot w \cdot \frac{\epsilon}{2w} = \frac{\epsilon n}{2}$$

← additive error small

since $M \geq n-1 \geq \frac{n}{2}$ (for $n \geq 2$) ← this is where we use l.b. on edge wts

$$|M - \hat{M}| \leq \epsilon M$$

← also multiplicative approx error is small

Distributed Algorithms vs. sublinear time algorithms on SPARSE graphs

↑
max deg $\leq d$

Again, sparse graphs: \circ max degree d
adj list representation

A problem to solve:

Vertex Cover

$V' \subseteq V$ is "Vertex Cover" (VC) if $\forall (u, v) \in E$

either $u \in V'$ or $v \in V'$

VC Question: What is min size of VC?

Note: in $\text{deg} \leq d$ graph, $|VC| \geq \frac{m}{d}$ since each node can cover $\leq d$ edges

(VC is NP-complete, but there is a polytime 2-multiplicative approximation)

Can you approximate V.C. in sublinear time?

multiplicative? no! graph with n edges $|VC| = 0$ \Rightarrow can't distinguish these cases in sublinear time but must answer 0 in first case & > 0 in second.
graph with 1 edge $|VC| = 1$

additive? hard! need some mult error
computationally hard to approx to better than 1.36 factor (maybe even 2)

Combination?

def. \hat{y} is (α, ϵ) -estimate of soln value y for minimization problem if

$$y \leq \hat{y} \leq \alpha y + \epsilon$$



allows mult + additive error

(analogous defn for maximization problems)

Some Background on distributed Algorithms

- Network

- processors \rightarrow max degree d known to all
- links

- Communication round

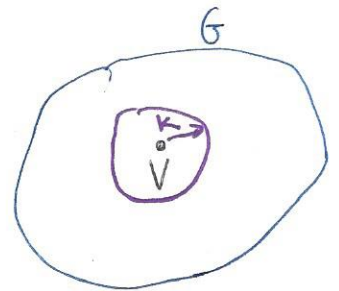
- nodes send messages to neighbors

def. Vertex Cover problem for distributed networks:

- Network graph = input graph (i.e. network computes on itself) ↙ not some other graph
- at end, each node knows if in or out of VC

Main insight on why fast distributed \Leftrightarrow sublinear time:

in k-round algorithm, output of node v only depends on nodes at distance at most k from v . At most d^k of these!



⇒ Can simulate V 's view of distributed computation in $\leq d^k$ time
↓ figure out if v in or out of VC

comment: if algorithm is randomized, v needs to know random bits (or be able to construct) of all d^k nbrs. \leftarrow must be consistent

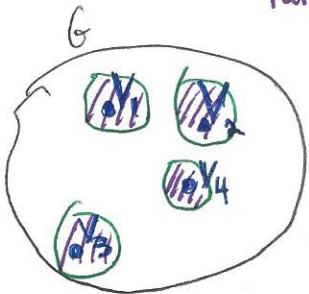
But are there fast VC distributed algorithms?

YES, will see some soon

↑ often called "local distributed algorithms"

How do you use this to approximate VC in sublinear time?

Parnas-Ron framework:



Sample nodes of graph v_1, \dots, v_r

for each v_i , simulated distributed algorithm to see if $v_i \in VC$

Output $\frac{\#v_i\text{'s in VC}}{r}$ on

Runtime $O(r \cdot d^{k/r}) \approx O(\frac{1}{\epsilon^2} \cdot d^{k/r})$

Proof of correctness Chernoff/Hoeffding bnds

fast distributed algorithm for VC:

$i \leftarrow 1$

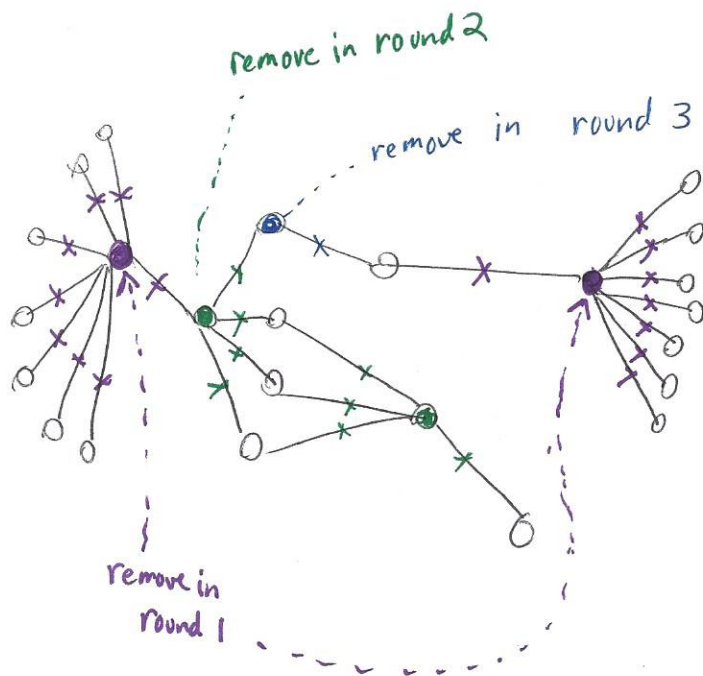
While edges remain:

- remove vertices of degree $> d/2^i$ + adjacent edges
- update degrees of remaining nodes
- increment i

Output all removed nodes as VC

#rounds: $\log d$

example:



is it a VC?

no edges remain at end

all removed along with some adjacent vertex

Is it a good approximation?

Thm let $VC(G) =$ size of min VC of G

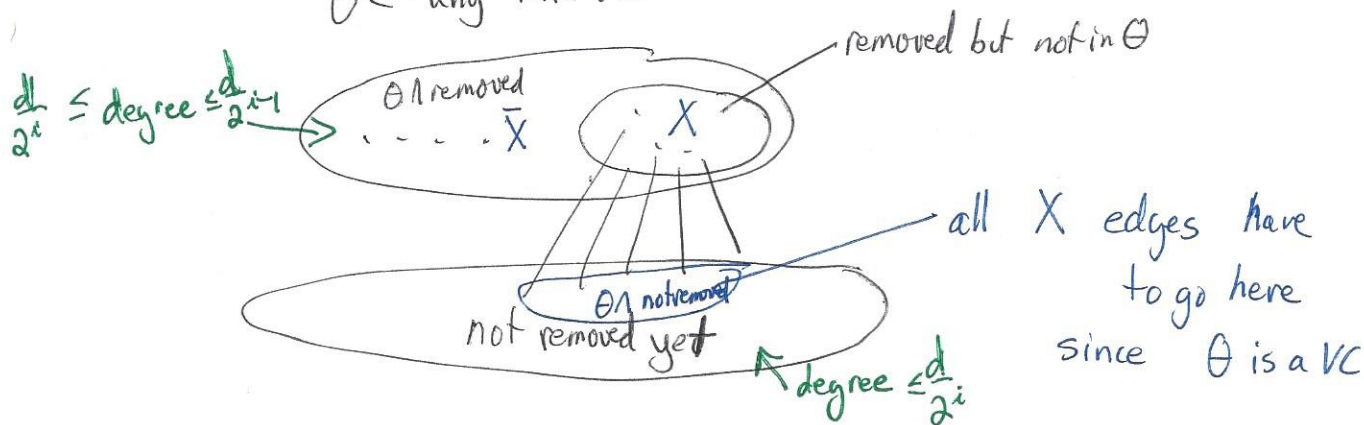
Then, $VC(G) \leq \text{output} \leq (2 \log d + 1) VC(G)$

\uparrow since output is VC \uparrow to prove

Proof.

Claim: in each iteration, add $\leq 2 VC_G$ new vertices

why? all nodes removed have deg bet $\frac{d}{2^{i+1}} + \frac{d}{2^i}$
 $\Theta \leftarrow$ any min VC



so $|X| \cdot \frac{d}{2^i} \leq |\Theta| \cdot \frac{d}{2^{i+1}}$

so $|X| \leq 2|\Theta|$