

# 6.889 Sub-linear Time Algorithms

Prof. Ronitt Rubinfeld

## Lecture 1

### Topics:

- Overview
- Diameter of point set
- #Connected Components
- MST

# Overview

- Big data - how do you access it?  
 too much to view it all?  
 can change by the time you  
 make a decision?

e.g. what is "diameter" of facebook graph?

## • Compromise:

"exactly", "for all", "there exists" statements  
 ⇒ "approximately", "there is a large group" ...

## • Models:

- 1) Random access queries - can access word of input in step  
 - adj matrix vs. adj list  
 - locn i
- 2) Samples - get sample of distribution/input in step

## Course requirements:

- Scribing 25%  
must be in latex
- Psets 35%
- Project 25%
- Class participation 15%  
(includes grading)

## Projects:

- groups of 1-3
- Possible ideas:
  - solve new problem (or try ideas + explain why they fail)
  - read + survey some papers
  - implement an algorithm (or two or three)

# A first Example: Diameter of a point set

Input  $m$  points described by distance matrix  $D$   
st.  $D_{ij}$  is distance from  $i, j$

note input size  $n$  is  $m^2$

- $D$  is: 1) symmetric
- 2) satisfies  $\Delta \neq$   
ie.  $D_{ij} \leq D_{ik} + D_{kj} \quad \forall i, j, k$

## Output

let  $i, j$  be st.  $D_{ij}$  is max  $\Leftarrow D_{ij}$  is "diameter"

Output  $k, l$  st.  $D_{kl} \geq \frac{D_{ij}}{2} \Leftarrow 2$ -approximation

## Algorithm

- Pick  $k$  arbitrarily
- Pick  $l$  to maximize  $D_{kl}$
- Output  $k, l$

runtime:  $O(m) = O(\sqrt{n})$

Why does it work?

$$D_{ij} \leq D_{ik} + D_{kl}$$

$\Delta \nabla$

$$= D_{ki} + D_{kl}$$

symmetry

$$\leq D_{kl} + D_{kl}$$

choice of  $l$

$$\leq 2D_{kl}$$

How many connected components in  $G$ ?

Input  $G = (V, E)$ ,  $\epsilon$   
max degree  $d$

adjacency list representation  
 $|V| = n$   
 $|E| = m \leq d \cdot n$

Output  $y$  st. if  $C = \# \text{ conn comp}$   
then  $C - \epsilon \cdot n \leq y \leq C + \epsilon n$

← "additive approx to w/in  $\epsilon n$ "

A different characterization of # conn components:

notation:  $\forall v$  let  $n_v \equiv \# \text{ nodes in } v\text{'s conn comp.}$

observation:  $\forall$  connected component  $A \subseteq V$   
$$\sum_{u \in A} \frac{1}{n_u} = \sum_{u \in A} \frac{1}{|A|} = 1$$

So new characterization of # conn comp:

$$C = \sum_{u \in V} \frac{1}{n_u}$$

Why is this better?

computing  $\frac{1}{n_u}$  needs  $O(n)$  time?  
sum  $O(n)$  terms?

↔ will estimate!

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

Two ideas:

- 1) estimate  $\frac{1}{n_u}$  quickly
- 2) estimate  $\sum_u \frac{1}{n_u}$  via sampling bounds

Estimating  $\frac{1}{n_u}$ :

def.  $\hat{n}_u \equiv \min \{n_u, 2/\epsilon\}$

$$\hat{C} \equiv \sum_{u \in V} \frac{1}{\hat{n}_u}$$

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

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

← so if can compute  $\hat{C}$  faster,  
it is useful!

Pf of lemma

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

else  $n_u > 2/\epsilon$  so  $\hat{n}_u = 2/\epsilon < n_u$

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

↑  
since  $n_u > 0$

$$\Rightarrow \left| \frac{1}{\hat{n}_u} - \frac{1}{n_u} \right| \leq \epsilon/2 \quad \blacksquare$$

How long 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

$$O(d \cdot 1/\epsilon)$$

↑  
time per step of BFS

How to estimate  $\sum_u \frac{1}{n_u}$ ?

Algorithm estimate  $\hat{c}$

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

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

$\forall u \in U$

compute  $\hat{n}_u$  via above algorithm

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

runtime

$$O\left(\frac{d}{\epsilon} \cdot \frac{1}{\epsilon^3}\right) = O(d/\epsilon^4)$$



Why is it good?

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

Pf

Chernoff Bnd:  $X_1 \dots X_r$  iid  $X_i \in [0,1]$   
 $S = \sum_{i=1}^r X_i$   $p = E[X_i] = E[S]/r$   
 Then:  $\Pr [ |\frac{S}{r} - p| \geq \delta p ] \leq e^{-\Omega(rp\delta^2)}$

here  $X_i = \frac{1}{\hat{n}_{u_i}}$

$$p = E\left[\frac{1}{\hat{n}_{u_i}}\right] = \frac{1}{n} \cdot \sum_{u \in V} \frac{1}{\hat{n}_{u_i}} = \frac{\hat{c}}{n}$$

$$\delta = \frac{\epsilon}{2}$$

$$\frac{S}{r} = \frac{1}{r} \sum_{i=1}^r \frac{1}{\hat{n}_{u_i}} \approx \frac{\hat{c}}{n}$$

so  $\Pr [ |\frac{\tilde{c}}{n} - \frac{\hat{c}}{n}| \geq \frac{\epsilon}{2} \cdot \frac{\hat{c}}{n} ] = \Pr [ |\tilde{c} - \hat{c}| \geq \frac{\epsilon}{2} \hat{c} ]$

$$\leq e^{-\left(\frac{b}{\epsilon^2} \cdot \frac{\hat{c}}{n} \cdot \frac{\epsilon^2}{4}\right)}$$

} want this to be  $\geq 2$  so that probability  $\leq e^{-2} \leq \frac{1}{4}$

Since  $\frac{\epsilon}{2} \leq \frac{1}{\hat{n}_u} \leq 1$   
 summing over  $u$ :  $\frac{\epsilon n}{2} \leq \hat{c} \leq n$   
 so  $\frac{\epsilon}{2} \leq \frac{\hat{c}}{n} \leq 1$

$$\leq e^{-\left(\frac{b}{\epsilon^2} \cdot \frac{\epsilon}{2} \cdot \frac{1}{4}\right)}$$

pick  $b \geq 16$

Now we are done:

$$\underline{\text{Corr.}} \quad \Pr[|C - \tilde{C}| \leq \epsilon n] \geq 3/4$$

$$\underline{\text{Pf.}} \quad |C - \tilde{C}| \leq |C - \hat{C}| + |\hat{C} - \tilde{C}| \quad \text{by } \Delta \neq$$

↑  
always  $\leq \frac{\epsilon n}{2}$   
by corr

↑  $\leq \frac{\epsilon n}{2}$  by thm  
with prob  $\geq 3/4$

# Approximating Min Spanning Tree (MST)

Input  $G = (V, E)$  adj list representation  $n = |V|$   
 max degree  $d$   
 each edge has weight  $w_{uv} \in \{1..W\} \cup \{\infty\}$   
 $\epsilon$   $G$  connected ie.  $w_{uv} \notin E$

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

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

Assumption on wts  $\Rightarrow n-1 \leq w(T) \leq w(n-1)$

A different characterization of MST:

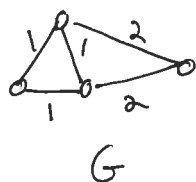
def  $G^{(i)} = (V, E^{(i)})$  where  $E^{(i)} = \{(u,v) \mid w_{uv} \in \{1..i\}\}$

$C^{(i)} = \#$  conn comp of  $G^{(i)}$

Some examples before characterization:

1)  $w=1$  only size 1 weights + connected by assumption  
here  $M=n-1$

2)  $w=2$  weights  $\in \{1, 2\}$



$G^{(1)}$   
 $C^{(1)}=2$

idea of Kruskal:

use as many wt 1 edges as you can  
only need wt 2 edges to connect the components

$\Rightarrow$  need  $C^{(1)} - 1$  wt 2 edges in MST

(recall that total # of edges is  $n-1$ )

Total wt of MST:

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

$\uparrow$   
1 for each edge

$\uparrow$  additional 1  
for wt 2 edges

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

Pf.

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

↑  
Kruskal's tells us that all MST's have same value of  $\alpha_i$

why?

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

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

where  $C^{(0)} = n$  (no edges in  $G^{(0)}$ )

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

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

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

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



### Approximation Algorithm :

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

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

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

↑  
adds poly( $\log \frac{w}{\epsilon}$ )  
factor to runtime

### Runtime :

$\tilde{O}(d/(\epsilon')^4) = \tilde{O}(dw^4/\epsilon^4)$  for each call to approx # c.c.

Total  $\tilde{O}(dw^5/\epsilon^4)$

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

(Can improve to  $O(\frac{dw}{\epsilon^2} \log \frac{dw}{\epsilon})$ )

+ need  $\Omega(dw/\epsilon^2)$

### Approximation guarantee:

"error" if approx error of approx #cc is too big ( $> \epsilon'$ )

how?  
Hw!

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

Pr [all calls to approx #cc give output that is  $\epsilon'$  additive approx]  $\geq 1 - \frac{w}{4w}$  ← union bound

If happens:  $|M - \hat{M}| \leq w \cdot \frac{\epsilon n}{2w} = \frac{\epsilon n}{2}$  ← small additive error = 3/4

→ since  $M \geq n-1 \geq n/2$ ,  $|M - \hat{M}| \leq \epsilon M$  ← small multiplicative error

this is where we use lower bound on edge wts

**Conclusion:** runtime depends only on  $d, w, 1/\epsilon$  gives additive/multiplicative error