

A first Example: Diameter of a point set

Input  $m$  points described by distance matrix  $D$

st.  $D_{ij}$  is distance from  $i, j$

+  $D$  is: 1) symmetric

2) satisfies  $\Delta \neq$

ie.  $D_{ij} \leq D_{ik} + D_{kj}$

$\forall i, j, k$

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

Output

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

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

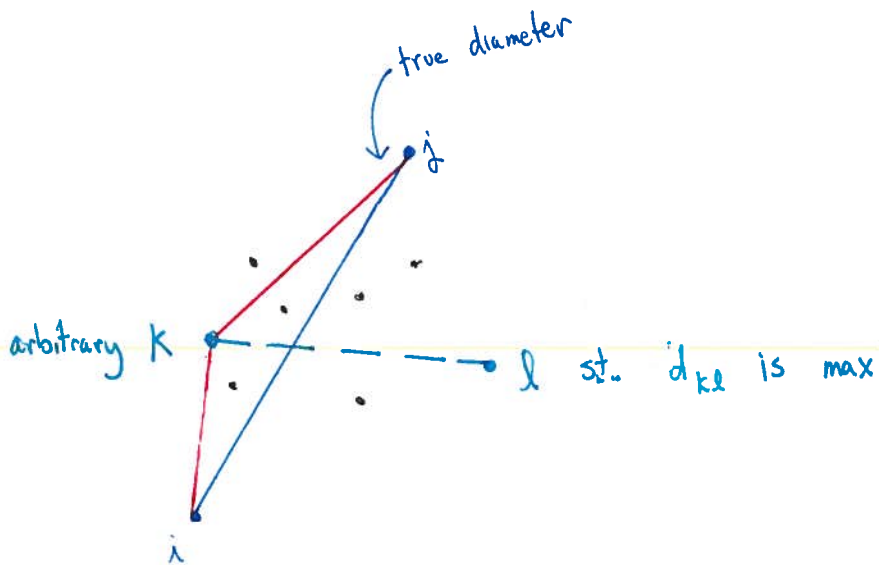
Algorithm

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

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

Why does it work?

$$\begin{aligned}
 D_{ij} &\leq D_{ik} + D_{kj} && \Delta \neq \\
 &= D_{ki} + D_{kj} && \text{symmetry} \\
 &\leq D_{kl} + D_{kl} && \text{choice of } l \\
 &\leq 2D_{kl}
 \end{aligned}$$



## A second example:

How many connected components in  $G$ ?

NOTE:  
← can solve via  
DFS/BFS in  
linear time

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  $\# \text{ conn components}$ :

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  $\# \text{ conn comp}$ :

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

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  $\leftarrow$  additive estimate

2) estimate  $\sum_u \frac{1}{n_u}$  via sampling bounds  $\leftarrow$  additive estimate but will use to get multiplicative error bound

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

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

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

idea  $n_u$  could be really big, so hard to compute exactly  
but if  $n_u$  really big then  $\frac{1}{n_u}$  is really small so can approx  $\frac{1}{n_u}$  by 0 or  $\frac{\epsilon}{2}$

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

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

$\leftarrow$  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  $\hat{n}_u = 2/\epsilon < n_u$

$$\Rightarrow \begin{array}{c} 0 \leq \frac{1}{n_u} - \frac{1}{\hat{n}_u} = \frac{\epsilon}{2} \\ \uparrow \\ \text{since } n_u > 0 \end{array}$$

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

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

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

use average value estimate  
to estimate sum

runtime  $O((d/\epsilon) \cdot 1/\epsilon^3) = O(d/\epsilon^4)$

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 should pick r to be  $\sim \frac{1}{\epsilon^2} \delta^2$   
 $\uparrow p \delta^2$   
 $\uparrow \Omega(\frac{1}{\epsilon^2})$   
 $\uparrow \Omega(\frac{1}{\epsilon^2})$   
 another reason why we need  $\frac{1}{n} \geq \frac{\epsilon}{2}$

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  $\geq 2$

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  $b \geq 16$

so that probability  $\leq e^{-2} \leq \frac{1}{4}$

Now we are done:

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

Pf.  $|c - \tilde{c}| \leq |c - \hat{c}| + |\hat{c} - \tilde{c}|$  by  $\Delta \neq$

↑

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

↑

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

---