

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?
Six Degrees of Separation

- 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

What are we really studying?

Randomized Algorithms?

Approximation Algorithms?

Parallel Algorithms?

Communication Complexity?

Statistics?

Learning?

YES!!

What kinds of datasets?

Graphs : e.g., diameter, # connected components, MST, bipartiteness, clusterability...

Functions : e.g., monotonicity, convexity, linearity, juntas...

Distributions : e.g., uniformity, independence, entropy, monotonicity, ...

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)

Grading psets: { not turned in, \checkmark -, \checkmark , \checkmark + }

Comments
solutions + typical errors

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

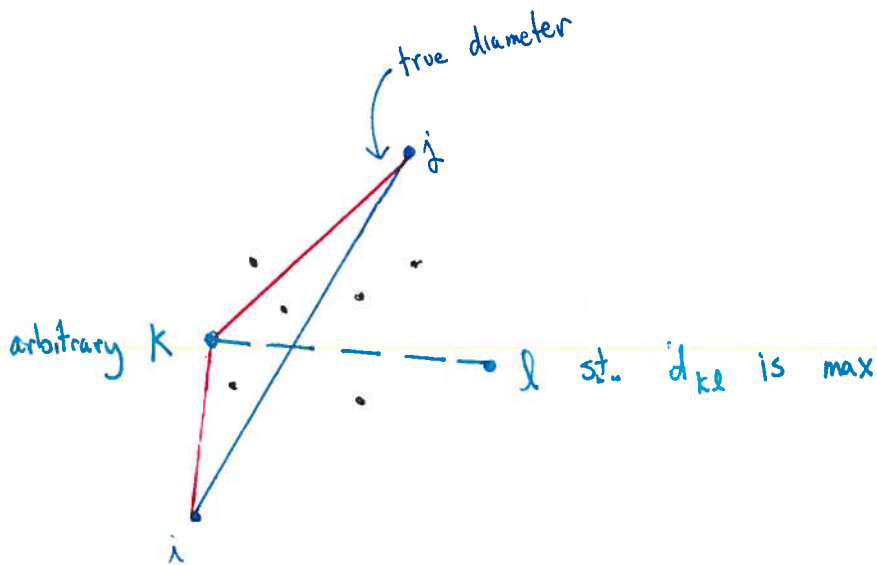
Algorithm

- 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 $G = (V, E)$, ϵ
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 ϵn "

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

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

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

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

so should pick r to be $\sim \frac{1}{\epsilon^2}$
 $\uparrow p \delta^2$
 $\uparrow \Omega(\frac{1}{\epsilon^2})$
 \uparrow 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 ≥ 2

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

summing over u: $\frac{\epsilon n}{2} \leq \hat{c} \leq n$

so $\frac{\epsilon}{2} \leq \frac{\hat{c}}{n} \leq 1$

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 $\geq 3/4$

Approximating Min Spanning Tree (MST)

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

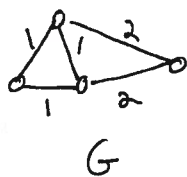
def $G^{(i)} = (V, E^{(i)})$ where $E^{(i)} = \{(u,v) \mid w_{uv} \in \{1..i\}\}$

$C^{(i)} = \#$ 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 $\in \{1, 2\}$



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

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

\uparrow
Kruskal's tells us that all MST's have
same value of α_i
why?

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

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

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

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

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

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



Approximation Algorithm :

For $i = 1$ to $w-1$
 $\hat{C}^{(i)} = \text{approx \# cc of } G^{(i)} \text{ to within } \frac{\epsilon'}{2w} \cdot n \text{ (additive error)}$
 Output $\hat{M} = n-w + \sum_{i=1}^{w-1} \hat{C}^{(i)}$
 ↑ adds poly($\frac{w}{\epsilon}$) factor to runtime

Runtime :

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

Call approx # cc with error probability $\leq \frac{1}{4w}$
 "error" if approx error of approx # cc is too big ($> \epsilon'$)
 how? Hwo!
 Pr [all calls to approx # cc give output that is ϵ' additive approx] $\geq 1 - \frac{w}{4w}$ ← union bound
 = $3/4$
 If happens: $|M - \hat{M}| \leq w \cdot \frac{\epsilon n}{2w} = \frac{\epsilon n}{2}$ ← small additive error
 → 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