

CS5234

Week 1

August 16

Part One: Sublinear time

Ex Graph is 1TB

Disk scan: 200 Mb/s \Rightarrow 83 min

Disk seek: 1 Mb/s \Rightarrow 11.5 days

BFS takes $>$ 1 week!!!

What can we do without scanning the whole graph?

Problem

Alg

Is G connected?

BFS

connected components?

$O(n+m)$

$\hookrightarrow O(1/\epsilon^2 d) / O(d/\epsilon^3)$

Weight of MST?

Prim's

$O(m \log n)$

$\hookrightarrow O(dW^4 \log(W) / \epsilon^3)$

Average degree?

$$\hookrightarrow O(\sqrt{n} \cdot \epsilon^{-9/2})$$

Scan
 $O(n)$

Diameter?

$$\hookrightarrow O(1/\epsilon^3)$$

BFS
 $O(n(m+n))$

Max matching

$$\hookrightarrow O(d^4/\epsilon^2)$$

Edmunds
Blossom Alg

Trade-off: approximate solution

Eg:

$$\text{MST}(G)[1-\epsilon] \leq \text{ALG}(G) \leq \text{MST}(G)[1+\epsilon]$$

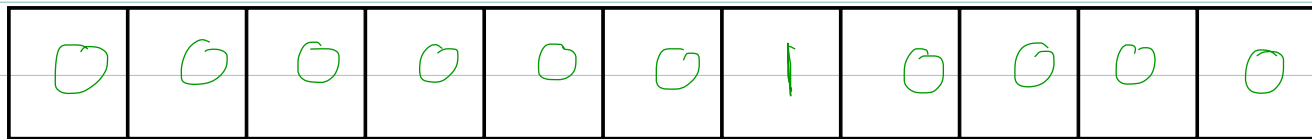
or

{ if G is connected \Rightarrow true

{ if G is " ϵ -far" from connected \Rightarrow false

Warm up:

n element array



0 = good
1 = error

Is array all 0's?

Alg: check all cells

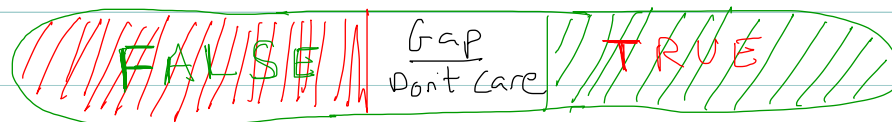
Time: $O(n)$

Impossible to do better: $\Omega(n)$

Is array mostly 0's??

{ IF all 0's: return true
IF $> \epsilon n$ 1's: return false
Else: don't care

Gap:



All-Zeros (A, ϵ)

Repeat S times:

Choose random $i \in [1, n]$

if $A[i] = 1$ then return false

Return true

Claim: If A is all 0, then alg returns true.

Fix $S = 2/\epsilon$

time: $O(1/\epsilon)$

Claim: If A has $\geq \epsilon n$ 1's, then
alg returns false.

Proof: For a sample i , $\Pr(A[i] = 1) \geq \frac{\epsilon n}{n} \geq \epsilon$

$$\begin{aligned} \Pr(\text{all samples are 0}) &\leq (1 - \epsilon)^S \\ &\leq (1 - \epsilon)^{2/\epsilon} \\ &\leq e^{-2} \\ &\leq 1/3 \end{aligned}$$

Don't forget:

$$\begin{aligned} e^{-x} &\geq 1 - x \\ (1 - \frac{1}{x})^x &\leq e^{-1} \end{aligned}$$

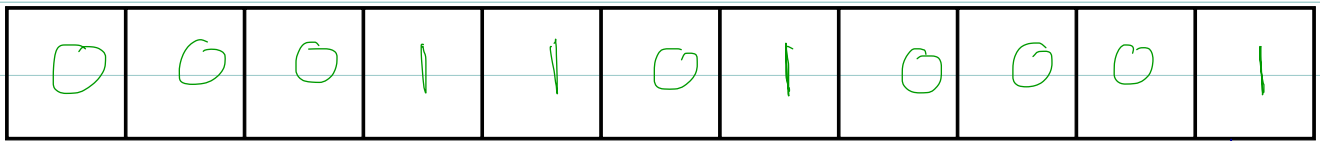
With probability $\geq 2/3$, correct

What if
you want
error $\leq \delta$?

$$\text{Set } S = \frac{1}{\epsilon} \ln \frac{1}{\delta}$$

Warm up #2:

n element array



0 = good
1 = error

What fraction of the array is 1's?

e.g., $4/11$

Tolerate ϵ error

e.g., 0.4 ± 0.05

Probability of error $\leq 1/3$

PercentZeros (A, ϵ)

sum = 0

Repeat S times:

Choose random $i \in [1, n]$

sum = sum + $A[i]$

Return (sum / S)

Will set $S = 1/\epsilon^2 \Rightarrow O(1/\epsilon^2)$ time

Hoeffding Bound

$Y_1, Y_2, \dots, Y_s = \text{iid random variables}$

$$Y_i \in [0, 1]$$

$$Y = \sum Y_i$$

$$\Pr[|Y - E[Y]| \geq \delta] \leq 2e^{-2\delta^2/s}$$

Set $Y_i = 1$ if i^{th} sample is 1
0 if i^{th} sample is 0

$$Y = \sum Y_i$$

Alg returns (Y/s)

Let $f = \text{fraction 1's}$

$$\Pr(Y_i = 1) = f, \quad E[Y_i] = 1 \cdot \Pr(Y_i = 1) + 0 \cdot \Pr(Y_i = 0)$$

$$E[Y] = sf \quad = f$$

$$E[Y/s] = f$$

$$\Pr[|Y/s - f| \geq \epsilon] = \Pr[|Y - fS| \geq \epsilon S]$$
$$= \Pr[|Y - E[Y]| \geq \epsilon S]$$

$$\text{Set } s = 1/\epsilon^2 \quad \leq 2e^{-2(\epsilon S)^2/s}$$
$$\leq 2e^{-2\epsilon^2 \cdot 1/\epsilon^2} \leq 1/3$$

Graph Connectivity

Graph $G = (V, E)$

n nodes

max degree d

$d \geq 3$

$[m \leq dn]$

Adjacency list

query: $\text{nbr}(u, i)$ returns
 i^{th} neighbor of u

Question 1: is G disconnected?

Requires $\Omega(m+n)$ time

Question 2: is G "very" disconnected?

Defn: Graph G is " Σ -far" from
connected if you need to
modify $\Sigma n d$ entries in
adjacency list to connect it.

Note: adding an edge requires modifying
2 entries: $e = (u, v) \Rightarrow$ update u and v
 \Rightarrow can delete/add $\Sigma n d / 2$ edges

Goal: if G is connected \rightarrow TRUE
 if G is ϵ -far from connected \rightarrow FALSE
 otherwise \rightarrow don't care

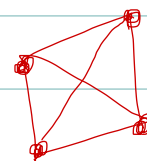
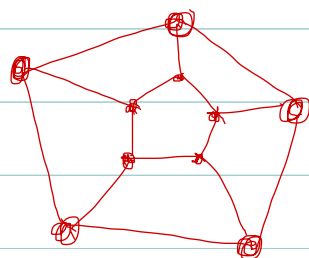
Pr(error) $\leq 1/3$

Preliminary claims

Lemma 1: if G is ϵ -far from connected
 then it has $\geq \epsilon dn/4$ Connected
 components.

Proof: Assume G has $< \epsilon dn/4$ C.C.
 Then can connect G by adding
 $< \epsilon dn/4 - 1$ edges \Rightarrow contradiction

$d=3$



Does not have
 2 nodes with free edge slot

cannot add
 1 edge to connect
 since $d=3$

Full CC with l nodes has $\geq \frac{(l-1)d}{2}$ edges

Spanning tree has $l-1 < \frac{(l-1)d}{2}$ edges

\Rightarrow First delete one edge from full CC.

Lemma 1 (cont.)

Total changes: $< \epsilon dn/4$ edges deleted
(to make root)

$< \epsilon dn/4 - 1$ edges added
(to connect)

$\Rightarrow \leq \epsilon dn$ changes to adjacency
list to connect G

$\Rightarrow G$ not ϵ -far from connected

\Rightarrow contradiction

Lemma 2: if G is ϵ -far from connected
then: $\geq \epsilon dn/8$ connected components
in G are of size $\leq 8/\epsilon d$

Proof: Counting

If $\epsilon dn/8$ have $> 8/\epsilon d \Rightarrow > n$ nodes

$\Rightarrow < \epsilon dn/8$ have $> 8/\epsilon d$

Lemma 1 says $\geq \epsilon dn/4$ in total

$\Rightarrow \geq \epsilon dn/4 - \epsilon dn/8 = \epsilon dn/8$ have $\leq 8/\epsilon d$

Connected(G, n, d, ϵ)

Repeat $\lceil 1/\epsilon d \rceil$ times

Choose u at random

Do a BFS from u , stopping when
 $\lceil 8/\epsilon d \rceil$ nodes are found

IF $|CC(u)| \leq \lceil 8/\epsilon d \rceil$, return FALSE

Return TRUE

Claim: time is $O(1/\epsilon^2 d)$

Proof: Each BFS takes time $\lceil 8/\epsilon d \rceil d$.

Repeat $\lceil 1/\epsilon d \rceil$.

$\Rightarrow \lceil 1/\epsilon d \rceil \lceil 8/\epsilon d \rceil d = O(1/\epsilon^2 d)$

Claim: if G is connected, returns TRUE

Lemma 3: if G is ϵ -far from connected
then returns FALSE

Proof: By Lemma 2, $\exists \epsilon d n / 8$ CC of
size $\leq \lceil 8/\epsilon d \rceil$.

Each has ≥ 1 node

$\Pr(u \text{ is in CC of size } \leq \lceil 8/\epsilon d \rceil) \geq \frac{(\epsilon d n / 8)}{n} \geq \epsilon d / 8$

Lemma 3 (continued):

$$\begin{aligned} \Pr(\text{return TRUE}) &\leq \left(1 - \frac{\epsilon d}{8}\right)^{16/\epsilon d} \\ &\leq e^{-\left(\frac{\epsilon d}{8}\right)\left(\frac{16}{\epsilon d}\right)} \leq e^{-2} \\ &\leq 1/3 \end{aligned}$$

error!

\Rightarrow alg is correct w.p. $\geq 2/3$