# Algorithms at Scale (Week 1)



#### 5 x 4: Use 7 lines.

Can you end at the same place you began?

7 x 7: Use 12 lines

# Sublinear Time / Sampling Algorithms

#### **Basic question:**

# What can we do when we look at only a small part of our input data?

#### Examples:

- Given a graph... only look at a *constant* number of nodes/edges.
- Is the graph connected?
- How many connected components does it have?
- What is the weight of the MST?
- What is the average degree of the graph?
- What is the diameter of the graph?
- What is the best matching on the graph?

	Classical	Sublinear Approximation
Is the graph connected?		
# connected components?		
Weight of MST?		

n : number of nodesm : number of edgesd : degree of graph

W : weight of MST

∈ : error / approximation parameter

	Classical	Sublinear Approximation
Is the graph connected?	O(n+m)	
# connected components?	O(n+m)	
Weight of MST?	$O\left(m\log n ight)$	

n : number of nodesm : number of edgesd : degree of graph

W : weight of MST

∈ : error / approximation parameter

	Classical	Sublinear Approximation		
Is the graph connected?	O(n+m)	$O\left(\frac{1}{\epsilon^2 d}\right)$		
# connected components?	O(n+m)	$O\left(\frac{d}{\epsilon^3}\right)$		
Weight of MST?	$O\left(m\log n ight)$	$O\left(\frac{dW^4\log W}{\epsilon^3}\right)$		

n : number of nodesm : number of edgesd : degree of graph

W : weight of MST∈ : error / approximation parameter

Comparison			does not matter how big the graph is!
	Classical	Sublinear Ap	proximation
Is the graph connected?	O(n+m)	$O\left(\frac{1}{\epsilon^2 a}\right)$	
# connected components?	O(n+m)	$O\left(\frac{d}{\epsilon^3}\right)$	
Weight of MST?	$O\left(m\log n\right)$	$O\left(rac{dW^4\mathrm{l}}{\epsilon^3} ight)$	$\left(\frac{\log W}{3}\right)$

n : number of nodesm : number of edgesd : degree of graph

W : weight of MST∈ : error / approximation parameter

	Classical	Sublinear Approximation
Average degree?	O(n+m)	$O\left(\frac{\sqrt{n}}{\epsilon^{9/2}}\right)$
Diameter?	$O\left(n(m+n)\right)$	$O\left(\frac{1}{\epsilon^3}\right)$
Maximum matching?	$O\left(mn^2 ight)$	$O\left(\frac{d^4}{\epsilon^2}\right)$

n : number of nodesm : number of edgesd : degree of graph

W : weight of MST∈ : error / approximation parameter

# Algorithms

#### Today

Toy example 1: array all 0's?

 Gap-style question: All 0's or far from all 0's?

#### Toy example 2: Faction of 1's?

- Additive ±  $\varepsilon$  approximation
- Hoeffding Bound

#### Is the graph connected?

- Gap-style question.
- O(1) time algorithm.
- Correct with probability 2/3.

### Next week:

Number of connected components in a graph.

 Additive approximation algorithm.

#### Weight of MST

• *Multiplicative* approximation algorithm.

# Trade-off: speed vs. accuracy

Approximate solutions:

Example: relative error

MST(G) = weight of MST

ALG(G) = weight of spanning tree returned by algorithm

 $MST(G)(1-\epsilon) \le ALG(G) \le MST(G)(1+\epsilon)$ 

Example: absolute error

$$MST(G) - \epsilon \leq ALG(G) \leq MST(G) + \epsilon$$

Example: gap error

- If G is connected, then return TRUE.
- If G is  $\in$ -far from connected, then return FALSE.
- Otherwise, don't care.

## Warm-Up Problem: All Zeros



#### Assumptions:

#### Given n element array

- Each element is 0 or 1.
- 0 = good test.
- 1 = failed test.

#### Output: Is the array all 0?



Running time: O(n)

Can we do any better?







#### Relaxed (approximate) version:

Output: Is the array *mostly* 0?

if all 0's:return trueif  $\geq \epsilon n$  1's:return falseotherwise:return true or false



#### Relaxed (approximate) version:

Output: Is the array *mostly* 0?

if all 0's:return trueif  $\geq \epsilon n$  1's:return falseotherwise:return true or false



#### Relaxed (approximate) version:

Output: Is the array *mostly* 0?

if all 0's:return trueif  $\geq \epsilon n$  1's:return falseotherwise:return true or false





#### Claim 1: If array is all 0's, then always returns TRUE.



Claim 2: If array has  $\geq \varepsilon n$  1's, then returns FALSE with probability at least 2/3.







Claim 1: If array is all 0's, then always returns TRUE.
Claim 2: If array has ≥ ɛn 1's, then returns FALSE with probability at least 2/3.



What if we want the algorithm to be correct with probability  $\geq 1 - \delta$ ?



What if we want the algorithm to be correct with probability  $\geq 1 - \delta$ ?

Fix s = ?



#### Input: n element array

- Each element is 0 or 1.
- 0 = good test, 1 = failed test.

Output: What fraction of the array is 1's?

Approximation: ± *ε* 

```
Probability correct: 2/3
```





#### **Requirements:**

1) Random variables: Y<sub>1</sub>, Y<sub>2</sub>, ..., Y<sub>s</sub> are random variables.  $Z = \sum_{j=1}^{s} Y_{j}$ 

2) Independent:

Y<sub>1</sub>, Y<sub>2</sub>, ..., Y<sub>s</sub> are *independent*.

3) Bounded:

Each  $Y_j$  is in the range [0,1].

#### **Requirements:**

1) Random variables: Y<sub>1</sub>, Y<sub>2</sub>, ..., Y<sub>s</sub> are random variables.  $Z = \sum_{j=1}^{s} Y_j$ 

2) Independent:

 $Y_1, Y_2, ..., Y_s$  are <u>independent</u>.

3) Bounded: Each Y<sub>j</sub> is in the range [0,1].

#### Example:

- 1) Random variables:
- $Y_j$  = value of array sampled in  $j^{th}$  iteration of the loop.

2) Independent:

Each sample is independent.

3) Bounded:

Each array entry is in [0,1].

If  $Y_1$ ,  $Y_2$ , ...,  $Y_s$  are independent random variables in the range [0,1] and if  $Z = \sum_{j=1}^{s} Y_j$  then:

$$\Pr[|Z - \mathcal{E}[Z]| \ge \delta] \le 2e^{-2\delta^2/s}$$

Conclusion:

The value of Z is within  $\delta$  of its expected value with "good" probability.

Claim:

Imagine you flip a coin n times.

Then you will see:

- at least  $\frac{n}{2} \sqrt{n}$  heads
- at most  $\frac{n}{2} + \sqrt{n}$  heads

with probability at least 2/3.

Claim:

Imagine you flip a coin n times.

Then you will see:

- at least  $\frac{n}{2} \sqrt{n}$  heads
- at most  $\frac{n}{2} + \sqrt{n}$  heads

with probability at least 2/3.

Random variables:

1) Random variables:

 $Y_j = 1$  if flip j is heads.  $Y_j = 0$  if flip j is tails.

$$Z = \sum_{j=1}^{n} Y_j$$

 $\mathbf{E}[Z]=n/2$ 

2) Independent: YES
 3) Bounded: YES

Hoeffding Bound:

 $Y_j = 1$  if flip j is heads.  $Y_j = 0$  if flip j is tails.  $Z = \sum_{j=1}^{n} Y_j$ E[Z] = n/2

 $\Pr[|Z - n/2| \ge \sqrt{n}] \le \Pr[|Z - \mathbb{E}[Z]| \ge \sqrt{n}]$  $\le 2e^{-2(\sqrt{n})^2/n}$  $\le 2e^{-2}$  $\le 1/3$ 

Claim:

Imagine you flip a coin n times.

Then you will see:

- at least  $\frac{n}{2} \sqrt{n}$  heads
- at most  $\frac{n}{2} + \sqrt{n}$  heads

with probability at least 2/3.

If  $Y_1$ ,  $Y_2$ , ...,  $Y_s$  are independent random variables in the range [0,1] and if  $Z = \sum_{j=1}^{s} Y_j$  then:

$$\Pr[|Z - \mathcal{E}[Z]| \ge \delta] \le 2e^{-2\delta^2/s}$$

Conclusion:

The value of Z is within  $\delta$  of its expected value with "good" probability.




Let f = fraction of 1's in the array.

(The "right" answer for the algorithm is f.)





Unbiased estimator! That's good ...

V = value returned

$$\Pr[|V - f| \le \epsilon] = \Pr[|V - E[V] \le \epsilon]$$

Can we use a Hoeffding Bound here?

 $Y_j = 1$  if sample j is 1.  $Y_j = 0$  if sample j is 0.

 $sum = \sum_{j=1}^{n} Y_j$ E[sum] = s \cdot sum

value returned: sum/sE[value returned] = f

V = value returned

 $\Pr[|V - f| \le \epsilon] = \Pr[|V - \mathbf{E}[V] \le \epsilon]$ 

Can we use a Hoeffding Bound here?

No, not yet. V is not a sum of random variables.  $Y_j = 1$  if sample j is 1.  $Y_j = 0$  if sample j is 0.

 $sum = \sum_{j=1}^{n} Y_j$ E[sum] = s \cdot sum

value returned: sum/sE[value returned] = f

V = value returned

$$\Pr[|V - f| \le \epsilon] = \Pr[|V - \mathbb{E}[V] \le \epsilon]$$
$$= \Pr[|sV - sf| \le s \cdot \epsilon]$$

 $Y_j = 1$  if sample j is 1.  $Y_j = 0$  if sample j is 0.

 $sum = \sum_{j=1}^{n} Y_j$ E[sum] = s \cdot sum

value returned: sum/sE[value returned] = f

Multiple everything by s.

V = value returned

 $Y_j = 1$  if sample j is 1.  $Y_j = 0$  if sample j is 0.

 $sum = \sum_{j=1}^{n} Y_j$ E[sum] = s \cdot sum

 $\cdot \epsilon$ 

value returned: sum/sE[value returned] = f

$$\Pr[|V - f| \le \epsilon] = \Pr[|V - E[V] \le \epsilon]$$
$$= \Pr[|sV - sf| \le s \cdot \epsilon]$$
$$= \Pr[|sum - E[sum]| \le s$$

Can we use a Hoeffding Bound here?

Yes. The random variable *sum* is the sum of independent 0/1 random variables.

Use a Hoeffding Bound!
$$Y_j = 1$$
 if sample  $j$  is 1.  
 $Y_j = 0$  if sample  $j$  is 0. $Y =$  value returned $sum = \sum_{j=1}^{n} Y_j$   
 $E[sum] = s \cdot sum$   
value returned:  $sum/s$   
 $E[value returned] = f$  $\Pr[|V - f| \le \epsilon]$  $=$   
 $\Pr[|sV - sf| \le s \cdot \epsilon]$   
 $=$   
 $\Pr[|sum - \mathbb{E}[sum]| \le s \cdot \epsilon]$   
 $\le 2e^{-2s^2\epsilon^2/s}$ 

Use a Hoeffding Bound!
$$Y_j = 1$$
 if sample  $j$  is 1.  
 $Y_j = 0$  if sample  $j$  is 0. $V =$  value returned $sum = \sum_{j=1}^{n} Y_j$   
 $E[sum] = s \cdot sum$   
value returned:  $sum/s$   
 $E[value returned] = f$  $\Pr[|V - f| \le \epsilon]$  $=$   
 $\Pr[|sV - sf| \le s \cdot \epsilon]$   
 $=$   
 $\Pr[|sum - E[sum]| \le s \cdot \epsilon]$   
 $\le 2e^{-2s^2\epsilon^2/s}$   
 $\le 2e^{-2s\epsilon^2}$ 

Use a Hoeffding Bound!
$$Y_j = 1$$
 if sample  $j$  is 1. $V =$  value returned $Y_j = 0$  if sample  $j$  is 0. $V =$  value returned $sum = \sum_{j=1}^{n} Y_j$  $Pr[|V - f| \le \epsilon] = Pr[|V - E[V] \le \epsilon]$  $value returned: sum/s$  $E[value returned] = f$  $Pr[|sV - sf| \le s \cdot \epsilon]$  $= Pr[|sum - E[sum]| \le s \cdot \epsilon]$  $\le 2e^{-2s^2\epsilon^2/s}$  $\le 2e^{-2s\epsilon^2}$  $\le 2e^{-2s\epsilon^2}$  $\le 2e^{-2s\epsilon^2/\epsilon^2}$  $\le 2e^{-2s\epsilon^2/\epsilon^2}$  $\le 2e^{-2s\epsilon^2/\epsilon^2}$  $\le 2e^{-2s\epsilon^2/\epsilon^2}$ 

Use a Hoeffding Bound!  

$$V = \text{value returned}$$

$$Pr[|V - f| \le \epsilon] = Pr[|V - E[V] \le \epsilon]$$

$$= Pr[|sV - sf| \le s \cdot \epsilon]$$

$$= Pr[|sum - E[sum]| \le s \cdot \epsilon]$$

$$\le 2e^{-2s^2\epsilon^2/s}$$

$$\le 2e^{-2s\epsilon^2}$$

$$\le 2e^{-2\epsilon^2/\epsilon^2}$$

$$\le 2e^{-2}$$

$$\le 1/3$$

Conclusion: value returned V is equal to  $f \pm \varepsilon$  w.p.  $\ge 2/3$ .



Conclusion: sum/s is equal to  $f \pm \varepsilon$  w.p.  $\ge 2/3$ .

# Algorithms

#### Today

Toy example 1: array all 0's?

 Gap-style question: All 0's or far from all 0's?

#### Toy example 2: Faction of 1's?

- Additive ±  $\varepsilon$  approximation
- Hoeffding Bound

#### Is the graph connected?

- Gap-style question.
- O(1) time algorithm.
- Correct with probability 2/3.

### Next week:

Number of connected components in a graph.

 Additive approximation algorithm.

#### Weight of MST

• *Multiplicative* approximation algorithm.

### Assumptions:

#### Graph G = (V,E)

- Undirected
- n nodes
- m edges
- maximum degree d

### Output: Is the graph connected?



### Assumptions:

#### Graph G = (V,E)

- Undirected
- n nodes
- m edges
- maximum degree d

### Format:

- Graph is given as an adjacency list.
- Query: nbr(u, i) returns i<sup>th</sup> neighbor of node u.



#### Exact answer:

BFS solves in: O(m+n)

#### Cannot do faster than: $\Omega(m+n)$



# Challenge #2:

Prove it.

### Assumptions:

#### Graph G = (V,E)

- Undirected
- n nodes
- m edges
- maximum degree d

#### Output:

Is the graph *close to*" connected or *"far from*" connected?



Gap Approximation: If **G** is connected: **Return TRUE** If G is *ɛ*-far from connected: **Return FALSE** Otherwise: Don't care.



### Definition:

G is ε-close to connected if you can add/modify at most εnd entries in the adjacency list to create a new graph G' that is connected.



### Definition:

G is ε-close to connected if you can add/modify at most εnd entries in the adjacency list to create a new graph G' that is connected.



### Definition:

G is ε-close to connected if you can add/modify at most εnd entries in the adjacency list to create a new

graph G' that is connected.

Example: **n** = 10, **d** = 3



### **Definition:**

G is ε-close to connected if you can add/modify at most εnd entries in the adjacency list to create a new

graph G' that is connected.

Example: n = 10, d = 3Add 3 edges to connect graph.

Modify 6 entries in adjacency list.



### **Definition:**

G is  $\varepsilon$ -close to connected if you can add/modify at most  $\varepsilon$ nd entries in the adjacency list to create a new

graph G' that is connected.

Example: n = 10, d = 3 Add 3 edges to connect graph. Modify 6 entries in adjacency list. nd = 30 G is 1/5-close to connected.



### **Definition:**

G is ε-close to connected if you can add/modify at most εnd entries in the adjacency list to create a new graph G' that is connected.

G is *ε*-far from connected if it is not *ε*-close to connected.



### Assumptions:

#### Graph G = (V,E)

- Undirected
- n nodes
- m edges
- maximum degree d

#### Output:

If G is connected: return TRUE. If G is  $\varepsilon$ -far from connected: return FALSE. Else: don't care.





Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/4 connected components.



Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\ge \varepsilon dn/4$  connected components.

#### Proof:

Assume G has  $\leq \epsilon dn/4$  connected components.

Then add  $\epsilon$ dn/4-1 edges to build connected graph G'. That requires modifying  $\leq \epsilon$ dn/2 entries in adjacency list.

G is *ɛ*-close to connected.

ted, then

#### Oops!

Cannot always add an edge without increasing the degree of the graph.

#### Proot:

Assume G has  $\leq \epsilon dn/4$  connected components.

Then add  $\epsilon$ dn/4-1 edges to build connected graph G'. That requires modifying  $\leq \epsilon$ dn/2 entries in adjacency list.

G is *ɛ*-close to connected.

Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\ge \varepsilon dn/4$  connected components.

Proof:

Assume G has  $\leq \epsilon dn/4$  connected components.

For each connected component, if every node has degree d:



Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\ge \varepsilon dn/4$  connected components.

Proof: Assume G has ≤ ɛdn/4 connected components.

For each connected component, if every node has degree d:

If it has k nodes, find a spanning tree with k-1 edges. Remove any one edge not in spanning tree.





Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\ge \varepsilon dn/4$  connected components.

#### Proof:

Assume G has  $\leq \epsilon dn/4$  connected components.

Delete  $\leq \epsilon dn/4$  edges so each connected component has at least one node with degree < d.

Then add  $\leq \epsilon dn/4-1$  edges to build connected graph G'.



### Lemma: If G is $\varepsilon$ -far from connected, then it has $\ge \varepsilon dn/4$ connected components.

#### Proof:

Assume G has  $\leq \epsilon dn/4$  connected components.

Delete  $\leq \epsilon dn/4$  edges so each connected component has at least one node with degree < d.

Then add  $\leq \epsilon dn/4-1$  edges to build connected graph G'.



Modifies ≤ *ε*dn entries in adjacency list.

G is *ɛ*-close to connected.

Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/4 connected components.



Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.



Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.

Proof: Counting argument.


Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.

### Proof:

Counting argument. Assume not.

There are at least *convected* components.

At most half can be twice the average size...



Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.

## Proof:

Counting argument.

Assume not.

- There are at least *convected* components.
- At most  $\varepsilon dn/8$  are of size  $\leq 8/\varepsilon d$ .
- At least *ɛ*dn/8 are of size > 8/*ɛ*d.



#### Lemma:

If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.

### Proof:

Counting argument. Assume not.

- There are at least *ɛdn/4* connected components.
- At most  $\varepsilon dn/8$  are of size  $\leq 8/\varepsilon d$ .
- At least *ɛdn/8* are of size > 8/*ɛd*.

 $\left(\frac{\epsilon dn}{8}\right) \left\lfloor \frac{8}{\epsilon d} + 1 \right\rfloor > n \text{ $\longrightarrow$ CONTRADICTION}$ 

Lemma: If G is  $\varepsilon$ -far from connected, then it has  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.

Proof: Counting argument.



## Connected(G, n, d, $\varepsilon$ )

### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

#### Return **TRUE**



## **Connected(G, n, d, ε)**

#### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

Return **TRUE** 

Claim: Each BFS takes time at most  $d(8/\varepsilon d)$ .

## Connected(G, n, d, $\varepsilon$ )

#### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

Return **TRUE** 



Claim: Each BFS takes time at most  $d(8/\epsilon d) = 8/\epsilon$ .

Proof: Explore at most  $(8/\varepsilon d)$  nodes of degree at most d.

### **Connected(G, n, d, ε)**

### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

#### Return **TRUE**

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

**Proof:** 
$$\left(\frac{16}{\epsilon d}\right)\left(\frac{8}{\epsilon}\right) = O\left(\frac{1}{\epsilon^2 d}\right)$$

## Connected(G, n, d, $\varepsilon$ )

#### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

Return **TRUE** 

Claim: If G is connected, returns TRUE.

Proof: Immediate. No component has  $\leq 8/\varepsilon d$  nodes.

## **Connected(G, n, d, ε)**

#### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

Return **TRUE** 

Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .



Proof: ...

Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .

#### Proof:

If G is  $\varepsilon$ -far from connected, then it has at least  $\varepsilon$ dn/8 connected components of size  $\leq 8/\varepsilon$ d.



Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .

#### Proof:

If G is  $\varepsilon$ -far from connected, then it has at least  $\varepsilon$ dn/8 connected components of size  $\ge 1$  and  $\le 8/\varepsilon$ d.

Each iteration / sample has probability at least  $\geq \frac{\frac{\epsilon dn}{8}}{n} \geq \frac{\epsilon d}{8}$  of finding small

connected component and returning **FALSE**.



Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .

#### Proof:

Each iteration / sample has probability at least  $\geq \frac{\frac{\epsilon dn}{8}}{n} \geq \frac{\epsilon d}{8}$  of finding small

connected component and returning FALSE.

Pr[no iteration finds a small connected component]

$$\leq \left(1 - \frac{\epsilon d}{8}\right)^{s}$$

$$\leq \left(1 - \frac{\epsilon d}{8}\right)^{\frac{16}{\epsilon d}}$$

$$\leq e^{-2}$$

$$\leq 1/3$$

Death Bed Fact: (0 < x < 1)

 $e^{-x} = 1 - x + x^2/2 - \dots$ 

Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .

#### Proof:

Each iteration / sample has probability at least  $\geq \frac{\frac{\epsilon dn}{8}}{n} \geq \frac{\epsilon d}{8}$  of finding small

connected component and returning false.

 $\Pr[\text{no iteration finds a small connected component}] \leq$ 

→ Some iteration finds a small CC and returns FALSE with probability at least 2/3.

$$\leq \left(1 - \frac{\epsilon d}{8}\right)^{s}$$

$$\leq \left(1 - \frac{\epsilon d}{8}\right)^{\frac{16}{\epsilon d}}$$

$$\leq e^{-2}$$

$$\leq 1/3$$

Death Bed Fact: (0 < x < 1)

 $e^{-x} = 1 - x + x^2/2 - \dots$ 

## **Connected(G, n, d, ε)**

#### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

#### Return **TRUE**

Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .



## **Connected(G, n, d, ε)**

### Repeat 16/*ɛ*d times:

- Choose random node u.
- Do a BFS from u, stopping after 8/ɛd nodes are found.
- If CC of u has ≤ 8/εd nodes, return FALSE.

#### Return **TRUE**

- Claim: Total time is  $O(1/\varepsilon^2 d)$ .
- Claim: If G is connected, returns TRUE.
- Claim: If G is  $\varepsilon$ -far from connected, then returns FALSE with probability  $\ge 2/3$ .



### General idea:

- Use sampling and local approximation to understand global graph properties.
- For what other interesting properties can you do this?

### Questions to think about:

- Is gap approximation useful?
- Is there a better notion of "close to connected"?
- For what values of  $\varepsilon$  and d is this actually fast?
- What happens in dense graphs?
- Can you find a faster algorithm? In theory? In practice?

## Announcements / Reminders

### Problem sets:

Problem Set 1 will be released tomorrow.

Problem Set 1 will be due next week.

# Summary

### Last Week:

Toy example 1: array all 0's?

 Gap-style question: All 0's or far from all 0's?

#### Toy example 2: Faction of 1's?

- Additive ±  $\varepsilon$  approximation
- Hoeffding Bound

#### Is the graph connected?

- Gap-style question.
- O(1) time algorithm.
- Correct with probability 2/3.

# Today:

Number of connected components in a graph.

• Approximation algorithm.

#### Weight of MST

• Approximation algorithm.