

Combinatorial and Graph Algorithms

(Week 2)

Puzzle of the Day:

A bag contains a collection of blue and red balls. Repeat:

- Take two balls from the bag.
- If they are the same color, discard them both and add a blue ball.
- If they are different colors, discard the blue ball and put the red ball back.

What do you know about the color of the final ball?

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: Fraction of 1's?

- Additive $\pm \epsilon$ approximation
- Hoeffding Bound

Is the graph connected?

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

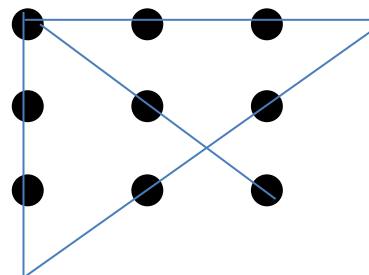
Today:

Number of connected components in a graph.

- *Additive* approximation algorithm.

Weight of MST

- *Multiplicative* approximation algorithm.



9 dots
4 lines

Announcements / Reminders

Problem sets:

Problem Set 1 was due today.

Problem Set 2 will be released tonight.

Today's Problem: Connected Components

Assumptions:

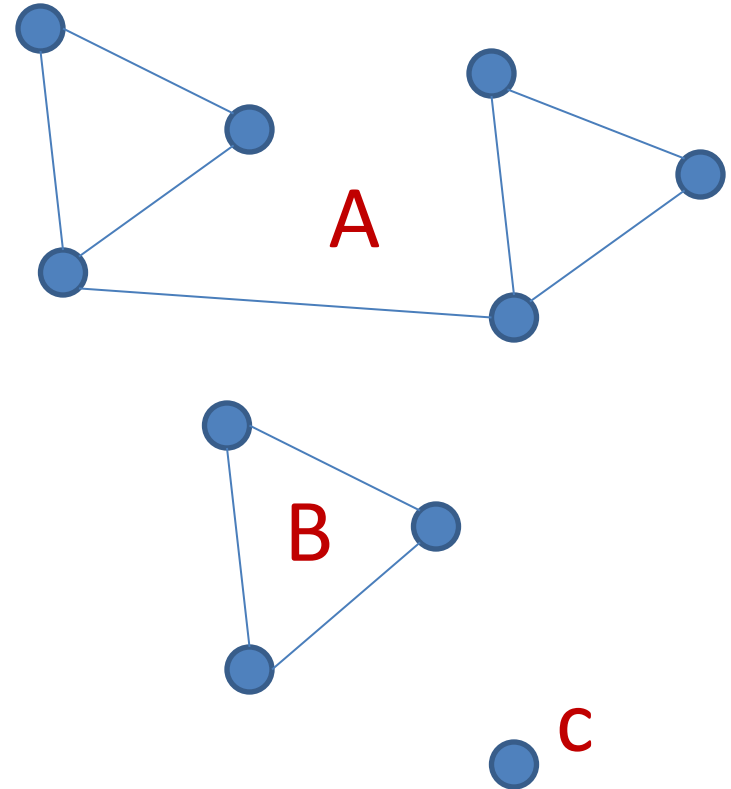
Graph $G = (V, E)$

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

Error term: ε

Output:

Number of connected components.



Example: output 3

Today's Problem: Connected Components

Approximation:

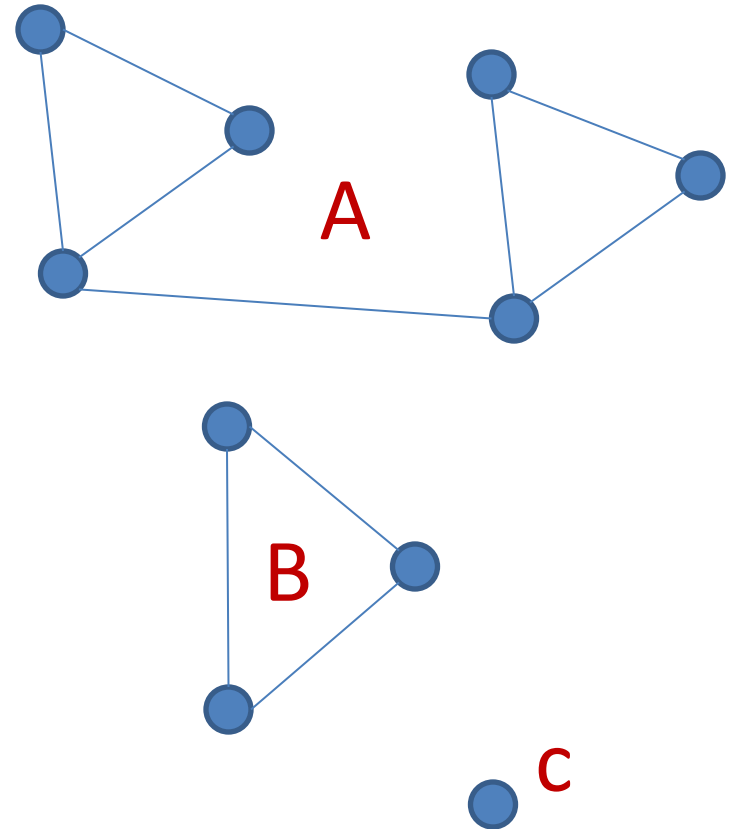
Output **C** such that:

$$CC(G) - \epsilon n \leq C \leq CC(G) + \epsilon n$$

Alternate form:

$$|CC(G) - C| \leq \epsilon n$$

Correct output: **w.p. > 2/3**



Example:

$$\epsilon = 1/10$$

$$\text{Output} \in \{2, 3, 4\}$$

Today's Problem: Connected Components

When is this useful?

What are trivial values of ε ?

What are hard values of ε ?

What sort of applications is this useful for?

Approximate Connected Components

When is this useful?

What are interesting values of ε ?

- What happens when $\varepsilon = 1$?
- What happens when $\varepsilon = 1/(2n)$?

What sort of applications is this useful for?

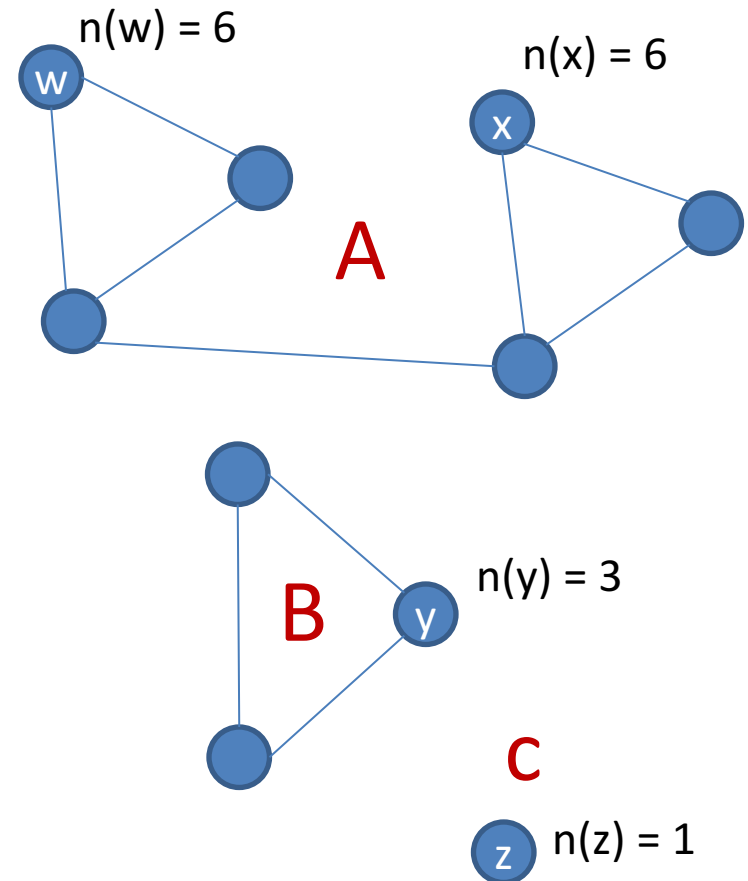
- Large graphs?
- Large social networks?
- The internet?
- Networks with many connected components?
- Number of components follows a heavy tail distribution?

Approximate Connected Components

Key Idea 1:

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .



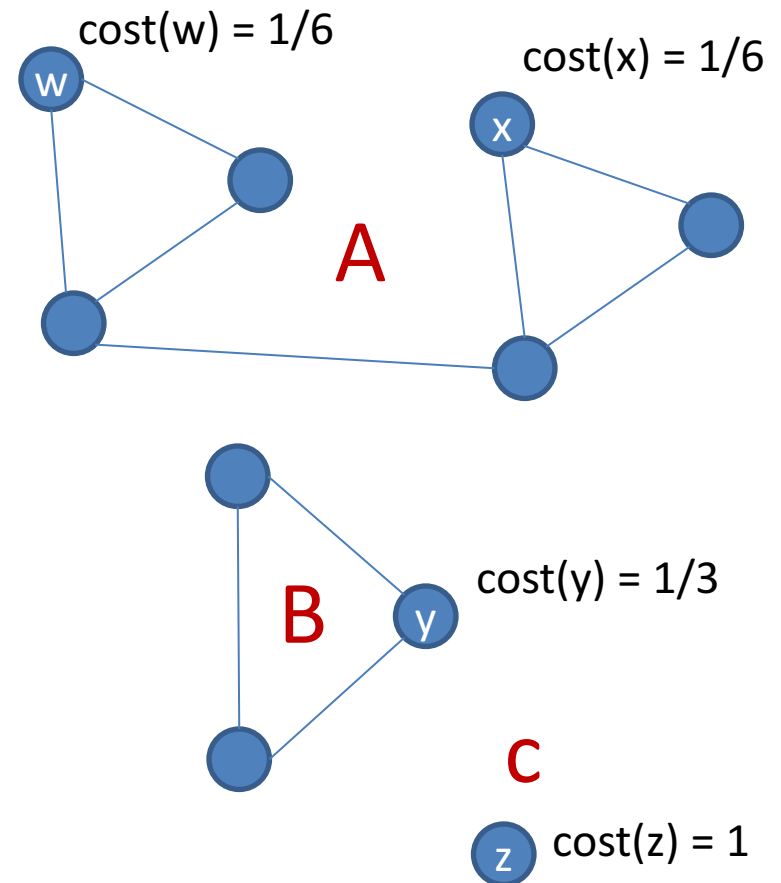
Approximate Connected Components

Key Idea 1:

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\text{cost}(u) = 1/n(u)$.

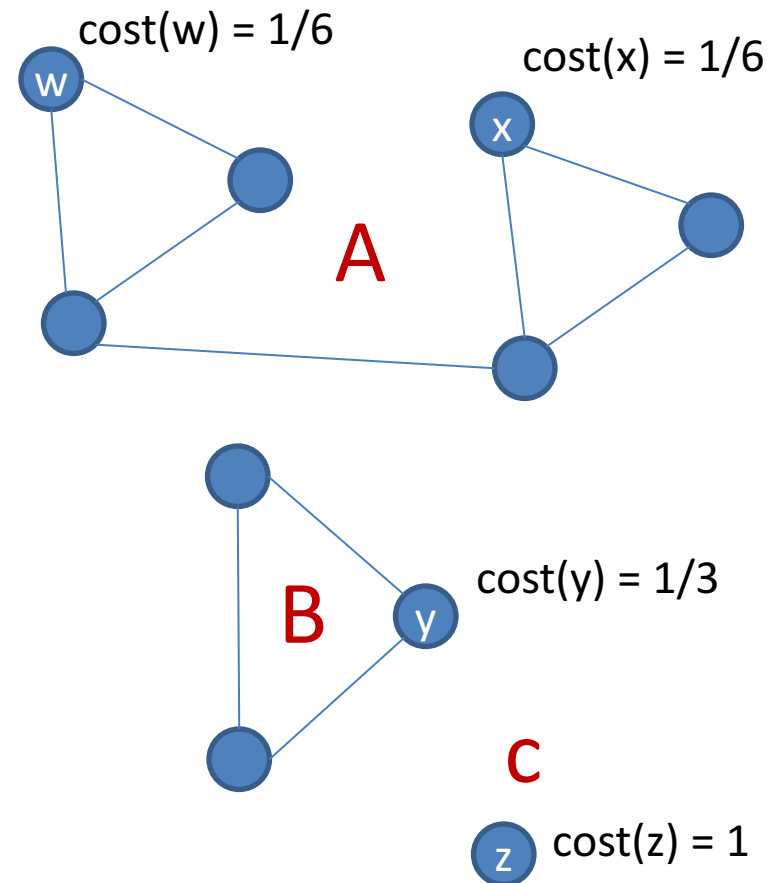


Approximate Connected Components

Key Idea 1:

Why is this useful?

$$\sum_{u \in A} \text{cost}(u) = ??$$

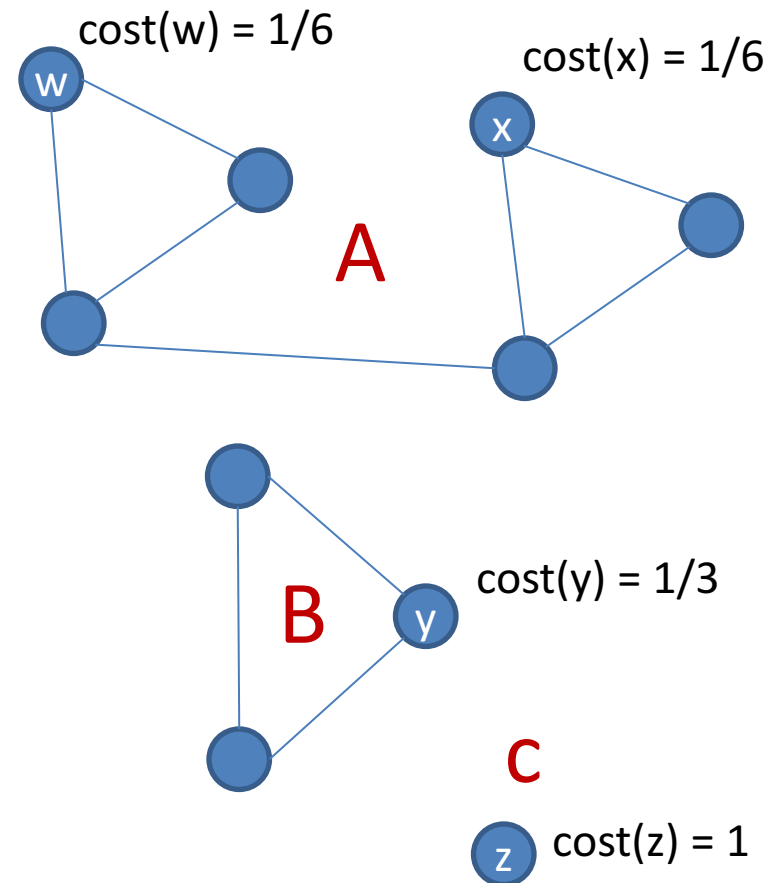


Approximate Connected Components

Key Idea 1:

Why is this useful?

$$\sum_{u \in A} \text{cost}(u) = 1$$



Approximate Connected Components

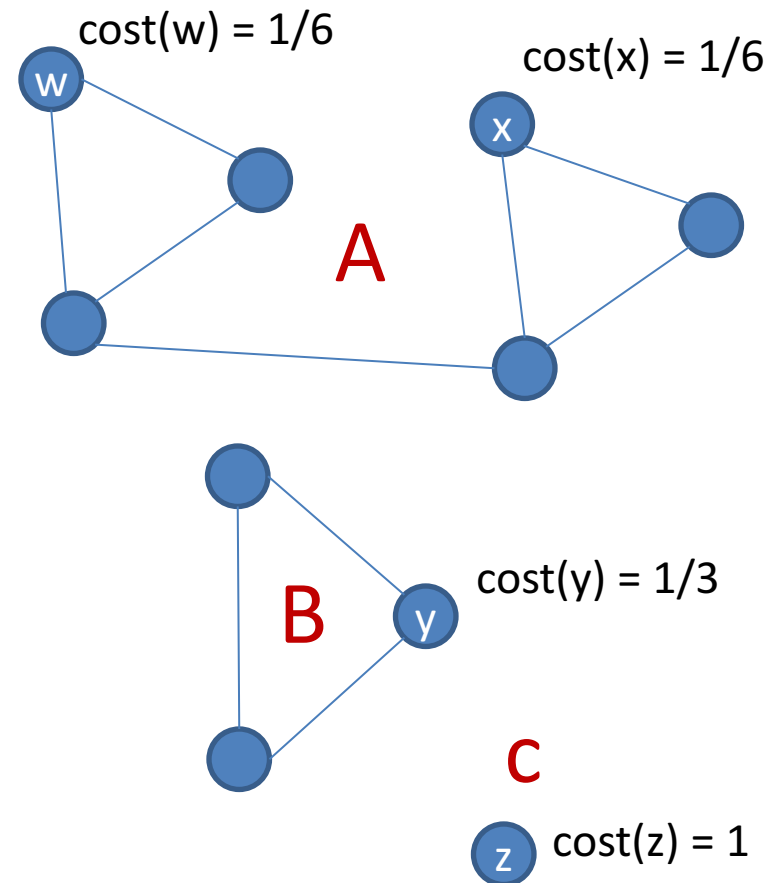
Key Idea 1:

Why is this useful?

$$\sum_{u \in A} \text{cost}(u) = 1$$

$$\sum_{u \in B} \text{cost}(u) = 1$$

$$\sum_{u \in C} \text{cost}(u) = 1$$



Approximate Connected Components

Key Idea 1:

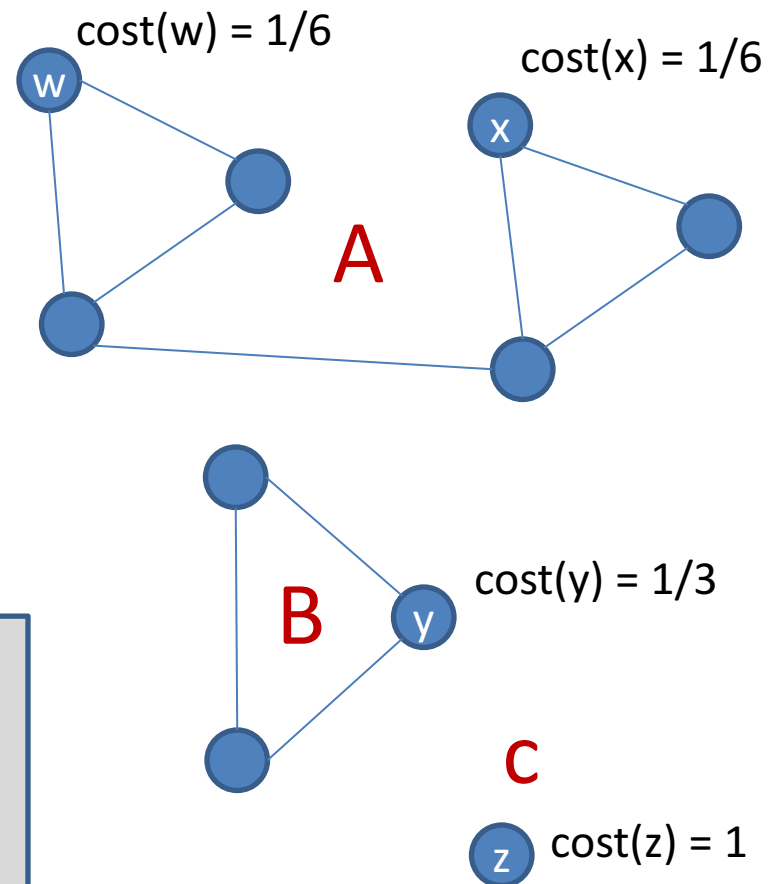
Why is this useful?

$$\sum_{u \in A} \text{cost}(u) = 1$$

$$\sum_{u \in B} \text{cost}(u) = 1$$

$$\sum_{u \in C} \text{cost}(u) = 1$$

$$\sum_{u \in V} \text{cost}(u) = \text{CC}(G)$$

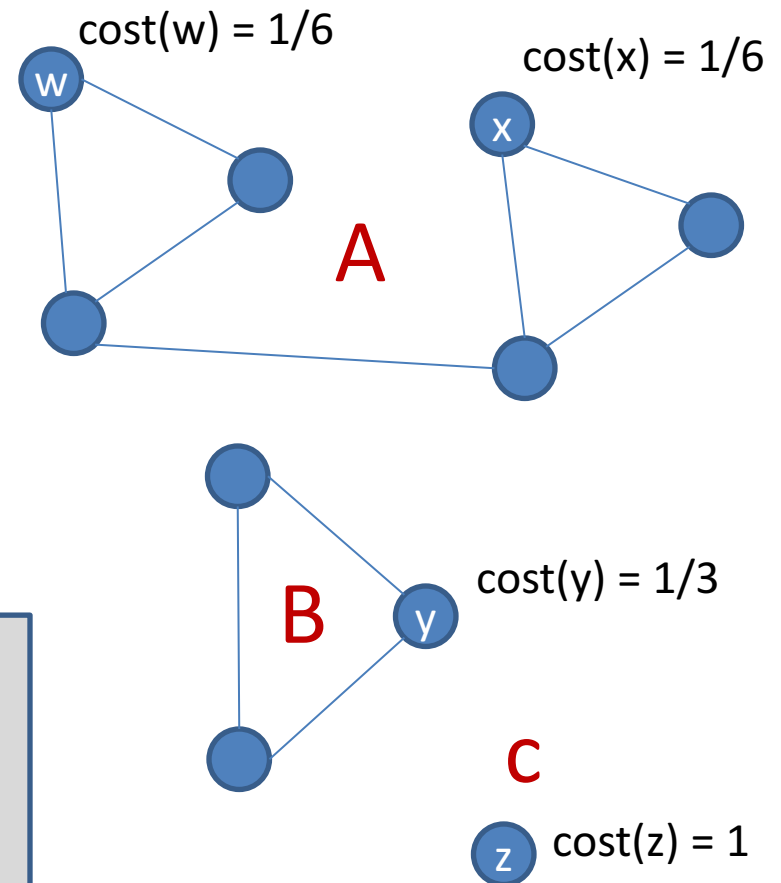


Approximate Connected Components

Algorithm 1

```
sum = 0
for each u in V:
    sum = sum + cost(u)
return sum
```

$$\sum_{u \in V} \text{cost}(u) = \text{CC}(G)$$



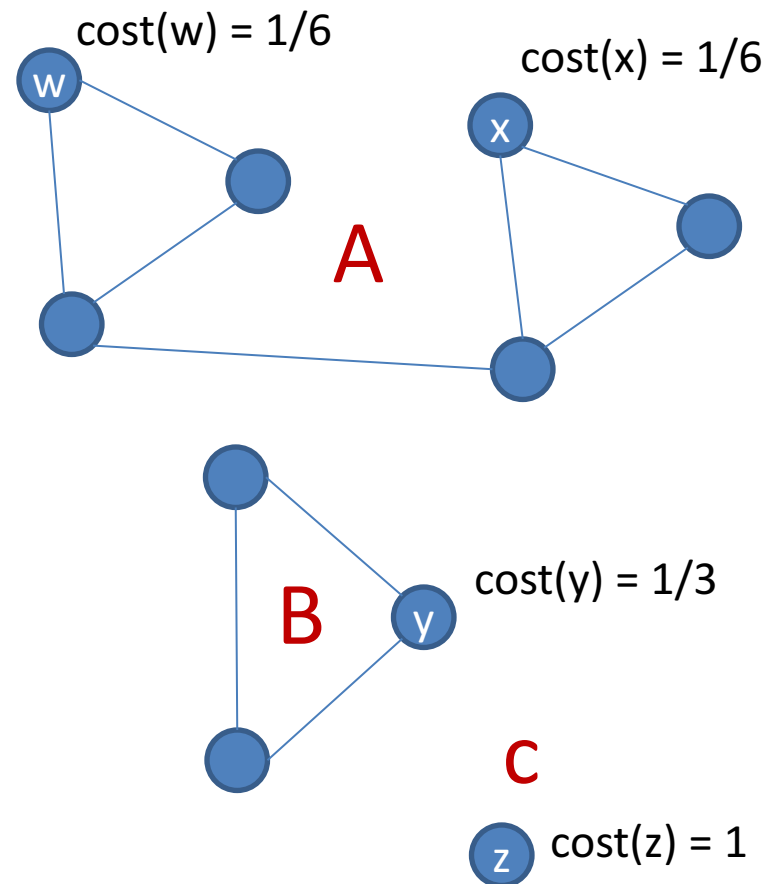
Approximate Connected Components

Algorithm 1

```
sum = 0
for each  $u$  in  $V$ :
    sum = sum + cost( $u$ )
return sum
```

Comments:

- Need a way to *efficiently* compute $\text{cost}(u)$.
- Runs in $O(n)$ time.

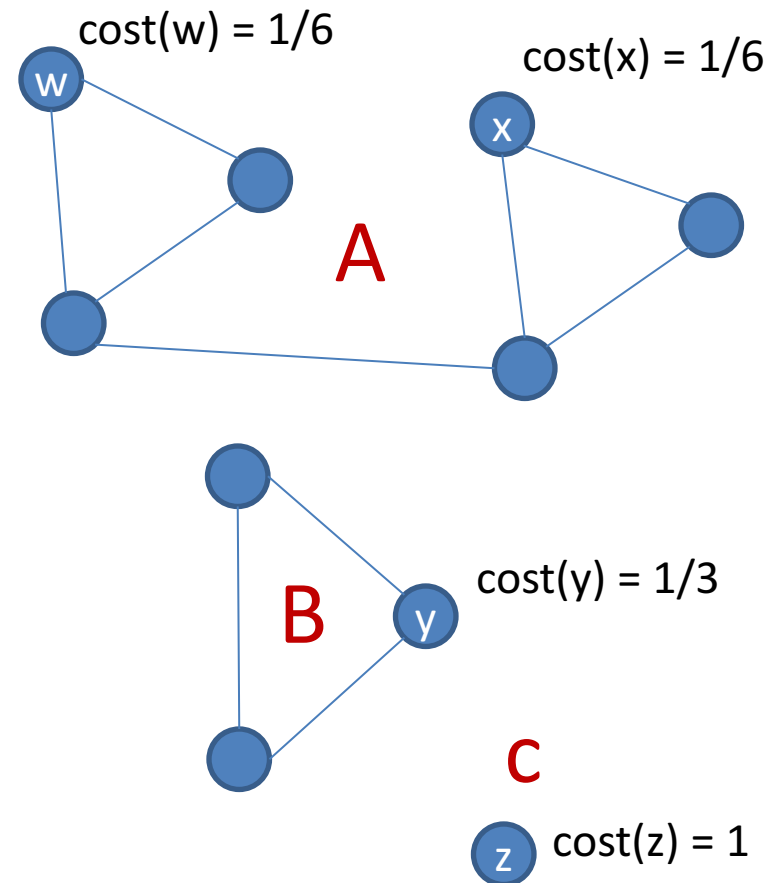


Approximate Connected Components

Key Idea 2: Sampling

Sample

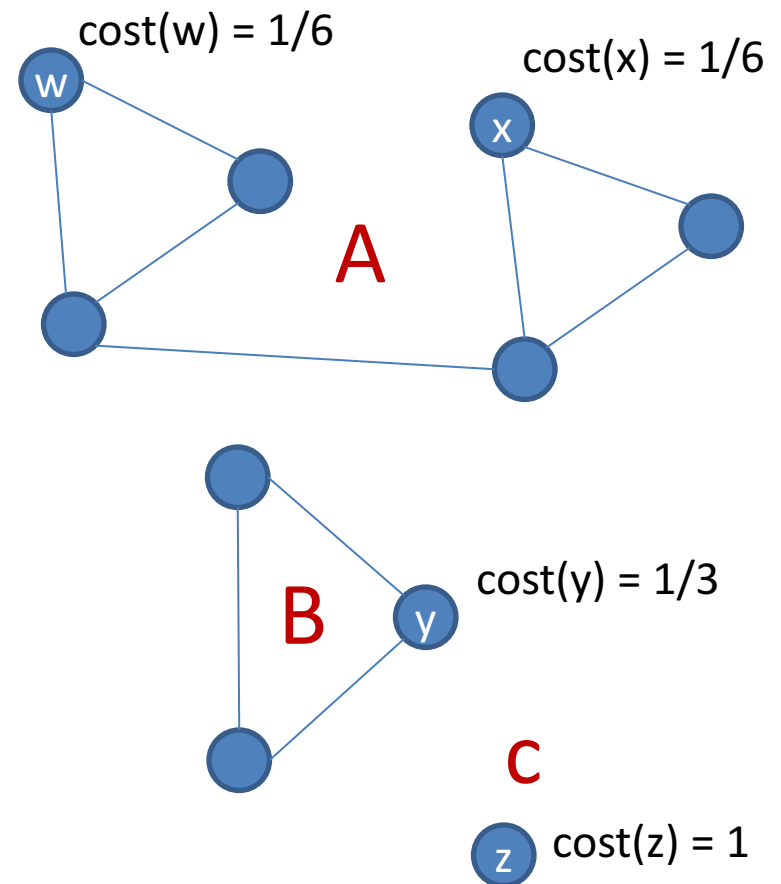
- Choose a small random subset S of V .
- For each node u in S , compute $\text{cost}(u)$.
- Use the sample to estimate the *average* cost of all the nodes.



Approximate Connected Components

Key Idea 2: Sampling

Worries?



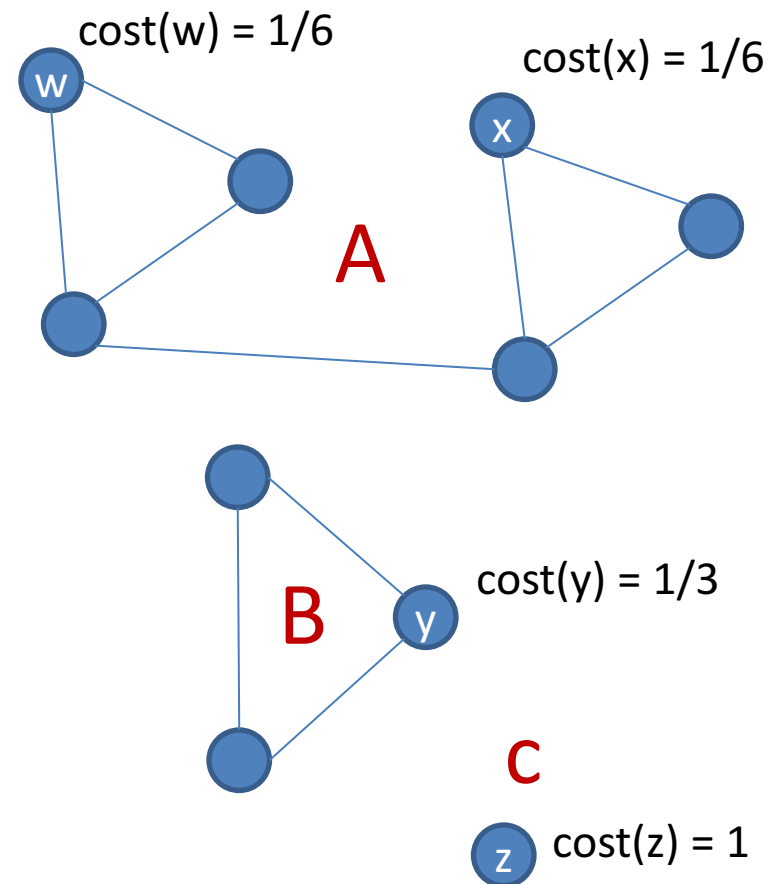
Approximate Connected Components

Key Idea 2: Sampling

Worries?

- Big components are sampled more often than small components?
- Small components may never be sampled?
- Bad examples?

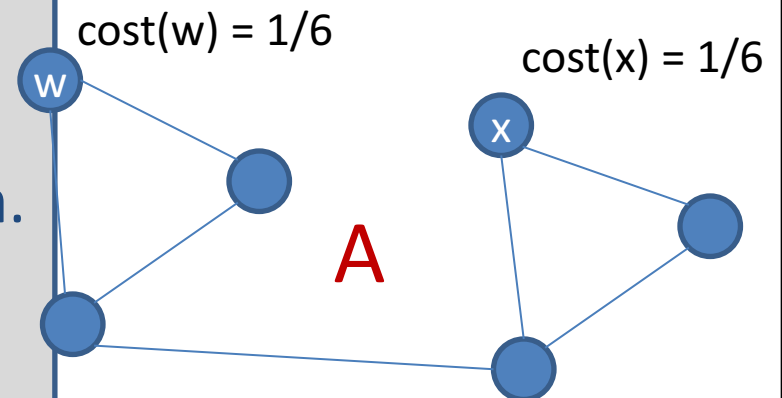
1 component of size 90,
10 components of size 1



Approximate Connected Components

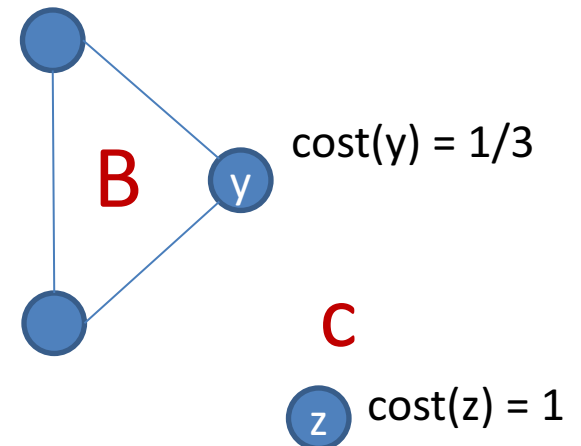
Algorithm 2

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n · (sum/s)
```



Comments:

- (sum/s) is average cost of sample.
- Efficiently compute $\text{cost}(u)$?
- Runs in $O(s)$ time.



Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Define random variables: Y_1, Y_2, \dots, Y_s

u_j = node chosen in j^{th} iteration

Y_j = $\text{cost}(u_j)$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E}[Y_j] = \sum_{i=1}^n \frac{1}{n} \text{cost}(u_i)$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E}[Y_j] = \sum_{i=1}^n \frac{1}{n} \text{cost}(u_i) = \frac{1}{n} \sum_{i=1}^n \text{cost}(u_i)$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

$$Y_j = \text{cost}(u_j)$$

$$\begin{aligned} \mathbb{E}[Y_j] &= \sum_{i=1}^n \frac{1}{n} \text{cost}(u_i) = \frac{1}{n} \sum_{i=1}^n \text{cost}(u_i) \\ &= \frac{1}{n} \text{CC}(G) \end{aligned}$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E}[Y_j] = \frac{1}{n} \text{CC}(G)$$

$$\begin{aligned} \mathbb{E} \left[\sum_{j=1}^s Y_j \right] &= s \mathbb{E}[Y_j] \\ &= \frac{s}{n} \text{CC}(G) \end{aligned}$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E}[Y_j] = \frac{1}{n} \text{CC}(G)$$

$$\mathbb{E}\left[\sum_{j=1}^s Y_j\right] = \frac{s}{n} \text{CC}(G)$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Notice:

Output of algorithm is: $\frac{n}{s} \sum_{j=1}^s Y_j$

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E}[Y_j] = \frac{1}{n} \text{CC}(G)$$

$$\mathbb{E} \left[\sum_{j=1}^s Y_j \right] = \frac{s}{n} \text{CC}(G)$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Notice:

Expected output of algorithm is:

$$\mathbb{E} [n \cdot (sum/s)] = \frac{n}{s} \left(\frac{s}{n} \text{CC}(G) \right) = \text{CC}(G)$$

$$Y_j = \text{cost}(u_j)$$

$$\mathbb{E} [Y_j] = \frac{1}{n} \text{CC}(G)$$

$$\mathbb{E} \left[\sum_{j=1}^s Y_j \right] = \frac{s}{n} \text{CC}(G)$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Important step:

Expected out is number of connected components!

(Algorithm is an unbiased estimator.)

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Notice:

Goal:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n \right\} \leq 1/3$$

Approximate Connected Components

Algorithm 2 Analysis

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n·(sum/s)
```

Notice:

Goal:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} \leq 1/3$$

Approximate Connected Components

Reminder: Hoeffding Bound

Given: independent random variables Y_1, Y_2, \dots, Y_s

Assume: each $Y_j \in [0,1]$

Then:

$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{j=1}^s Y_j \right] - \sum_{j=1}^s Y_j \right| > t \right\} \leq 2e^{-2t^2/s}$$

Approximate Connected Components

Reminder: Hoeffding Bound

Given: independent random variables Y_1, Y_2, \dots, Y_s

Assume: each $Y_j \in [0,1]$

Then:

$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{j=1}^s Y_j \right] - \sum_{j=1}^s Y_j \right| > t \right\} \leq 2e^{-2t^2/s}$$

Goal:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n/2 \right\} \leq 1/3$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

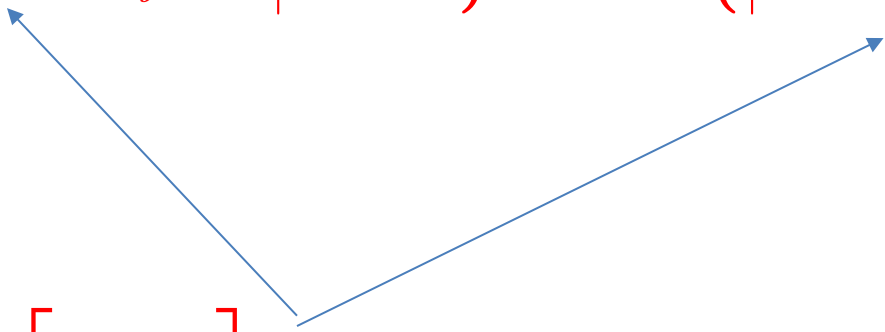
$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} =$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} = \Pr \left\{ \left| \mathbb{E} \left[\frac{n}{s} \sum_{i=1}^s Y_i \right] - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\}$$

$$\mathbb{E} \left[\sum_{j=1}^s Y_j \right] = \frac{s}{n} \text{CC}(G)$$


Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\begin{aligned}\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} &= \Pr \left\{ \left| \mathbb{E} \left[\frac{n}{s} \sum_{i=1}^s Y_i \right] - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} \\ &= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{i=1}^s Y_i \right| > \frac{s}{n} \epsilon n / 2 \right\}\end{aligned}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\begin{aligned}\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n/2 \right\} &= \Pr \left\{ \left| \mathbb{E} \left[\frac{n}{s} \sum_{i=1}^s Y_i \right] - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n/2 \right\} \\ &= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \frac{s}{n} \epsilon n/2 \right\} \\ &= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s/2 \right\}\end{aligned}$$

Approximate Connected Components

$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{j=1}^s Y_j \right] - \sum_{j=1}^s Y_j \right| > t \right\} \leq 2e^{-2t^2/s}$$

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n/2 \right\} = \Pr \left\{ \left| \mathbb{E} \left[\frac{n}{s} \sum_{i=1}^s Y_i \right] - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n/2 \right\}$$

$$= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \frac{s}{n} \epsilon n/2 \right\}$$

$$= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s/2 \right\}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} =$$
$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} \leq 2e^{-2(\epsilon s / 2)^2 / s}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\begin{aligned} \Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} &= \\ \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} &\leq 2e^{-2(\epsilon s / 2)^2 / s} \\ &\leq 2e^{-2\epsilon^2 s / 4} \end{aligned}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} =$$
$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} \leq 2e^{-2(\epsilon s / 2)^2 / s}$$
$$\leq 2e^{-2\epsilon^2 s / 4}$$

$$s = \frac{4}{\epsilon^2}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

$$\begin{aligned} \Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} &= \\ \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} &\leq 2e^{-2(\epsilon s / 2)^2 / s} \\ &\leq 2e^{-2\epsilon^2 s / 4} \\ &\leq 2e^{-\epsilon^2 (4 / \epsilon^2) / 2} \end{aligned}$$

$$s = \frac{4}{\epsilon^2}$$

Approximate Connected Components

Algorithm 2 Analysis

Derivation:

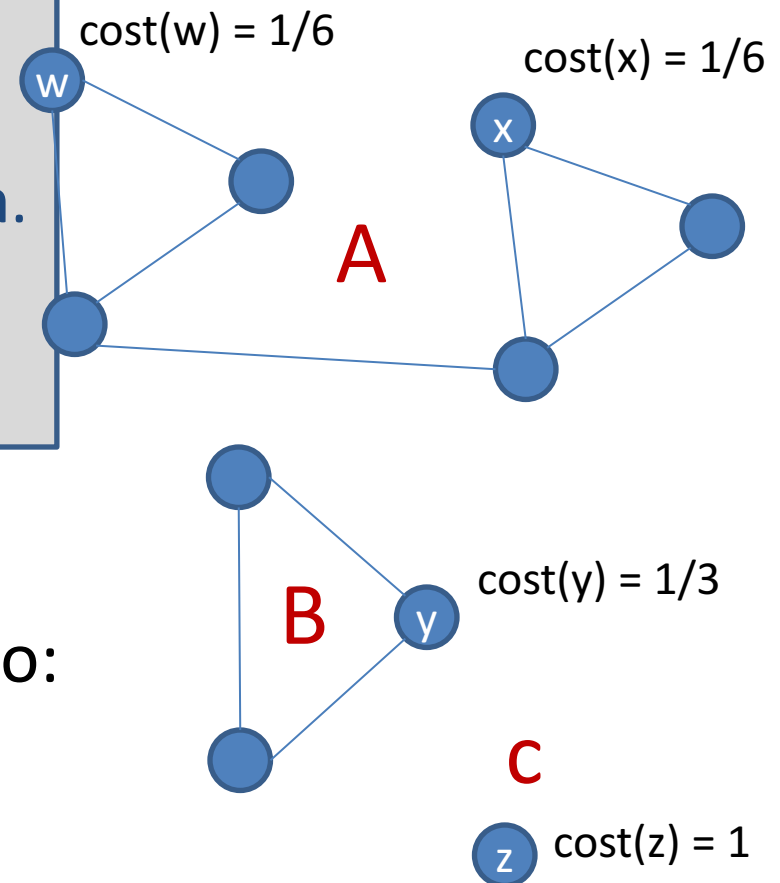
$$\begin{aligned} \Pr \left\{ \left| \text{CC}(G) - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} &= \\ \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} &\leq 2e^{-2(\epsilon s / 2)^2 / s} \\ &\leq 2e^{-2\epsilon^2 s / 4} \\ &\leq 2e^{-\epsilon^2 (4 / \epsilon^2) / 2} \\ &\leq 2e^{-2} \\ &< 1/3 \end{aligned}$$

$$s = \frac{4}{\epsilon^2}$$

Approximate Connected Components

Algorithm 2

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n · (sum/s)
```



We have shown:

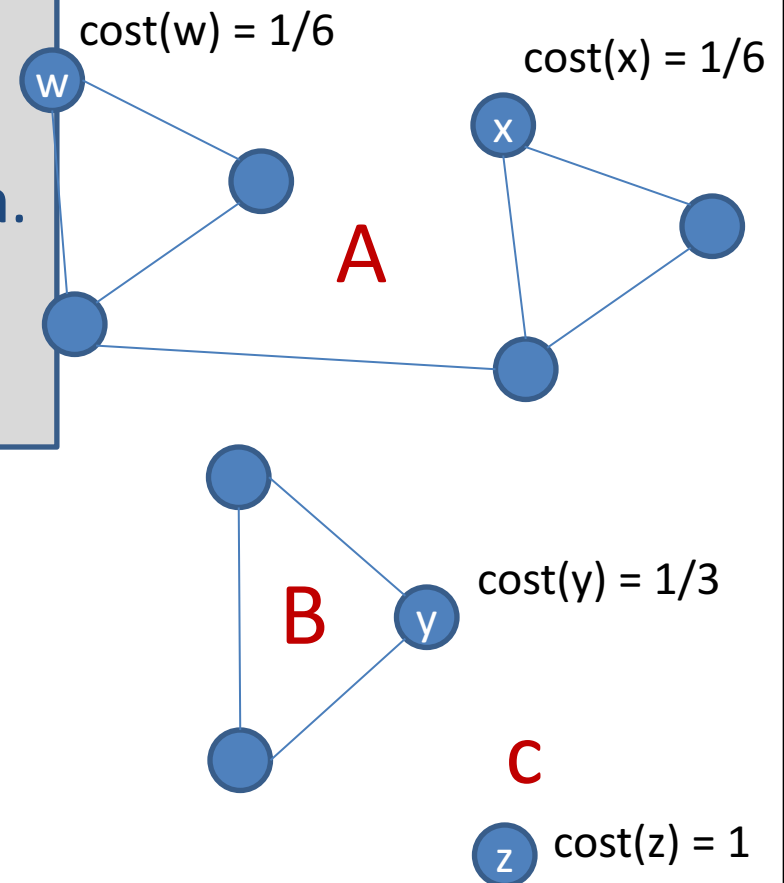
W.p. $> 2/3$, output is equal to:

$$\text{CC}(G) \pm \epsilon n$$

Approximate Connected Components

Algorithm 2

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n · (sum/s)
```



We have shown:

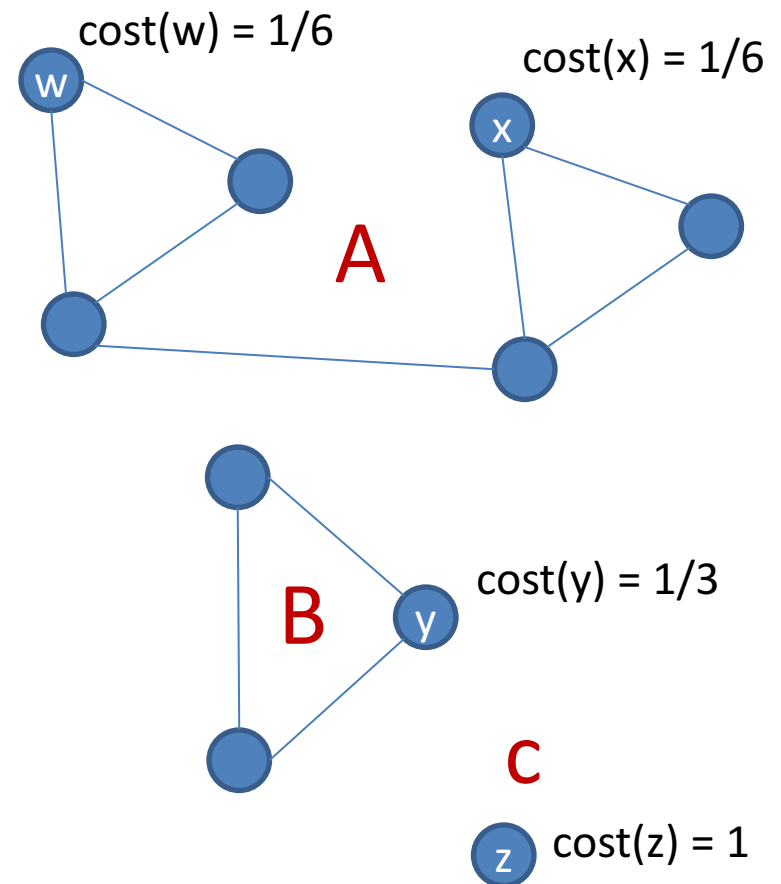
Time: $O(1/\varepsilon^2)$

Approximate Connected Components

Key Idea 2: Sampling

Key problem:

How to efficiently compute $\text{cost}(u)$.



Approximate Connected Components

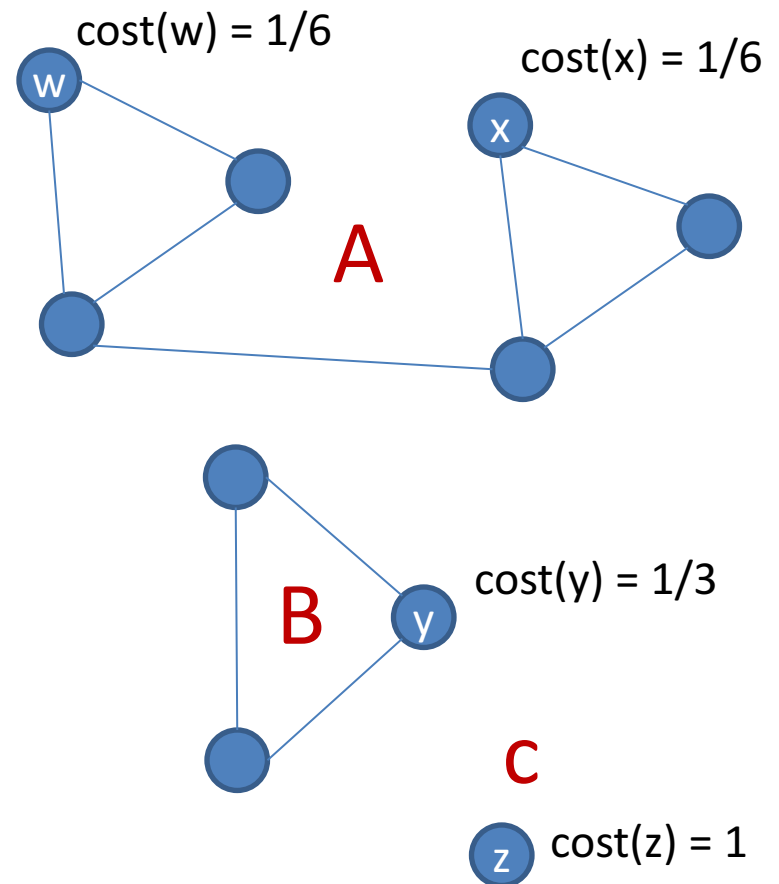
Key Idea 2: Sampling

Key problem:

How to efficiently compute $\text{cost}(u)$.

Key idea 3:

Approximate $\text{cost}(u)$.



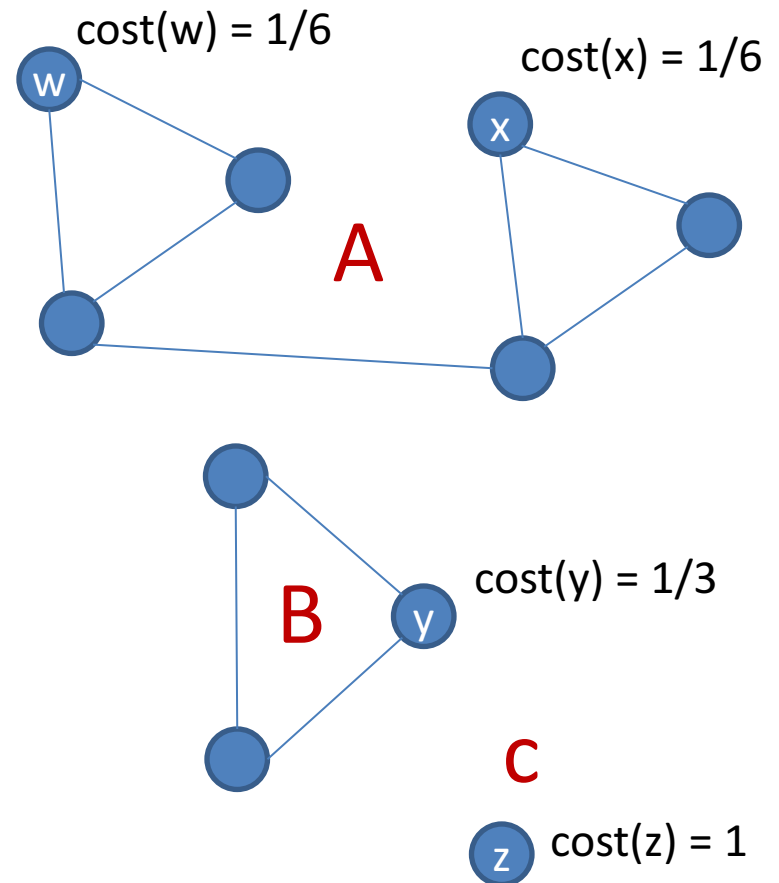
Approximate Connected Components

Key Idea 3: Approximate Cost

Approximate low cost components:

If $\text{cost}(u)$ is small, round up.

How small is small enough?

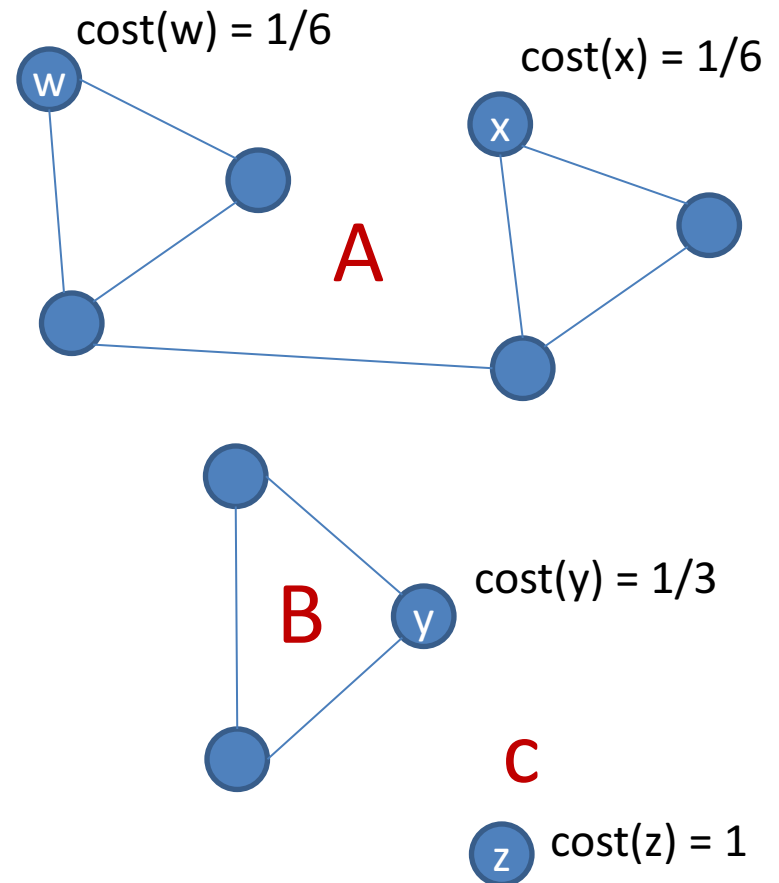


Approximate Connected Components

Key Idea 3: Approximate Cost

Approximate low cost components:

If $\text{cost}(u) < \varepsilon/2$, round up.



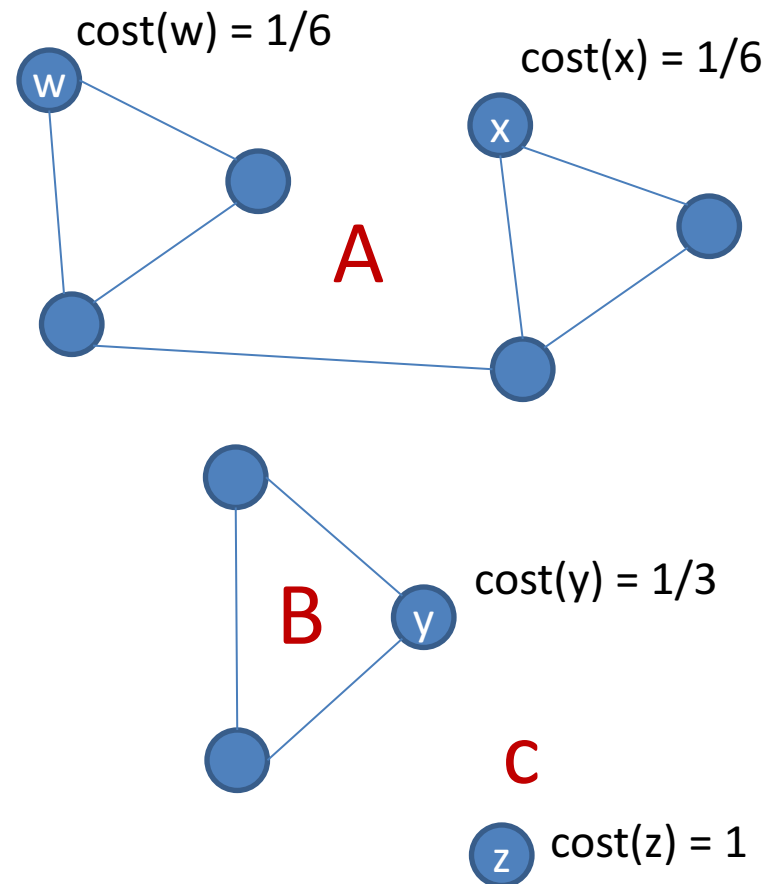
Approximate Connected Components

Key Idea 3: Approximate Cost

Ignore low cost components:

If $\text{cost}(u) < \varepsilon/2$, round up.

Total added cost $\leq \varepsilon n/2$.



Approximate Connected Components

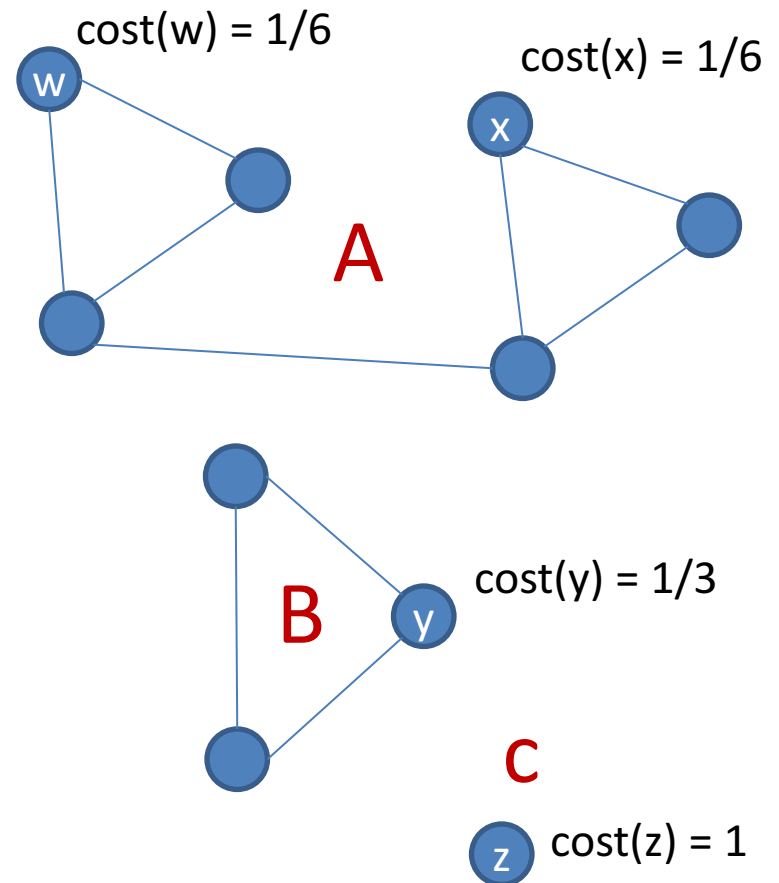
Key Idea 3: Approximate Cost

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\tilde{n}(u) = \min(n(u), 2/\epsilon)$.

Let $\text{cost}(u) = \max(1/n(u), \epsilon/2)$.
 $= 1/\tilde{n}(u)$.



Approximate Connected Components

Key Idea 3: Approximate Cost

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\tilde{n}(u) = \min(n(u), 2/\varepsilon)$.

Let $\text{cost}(u) = \max(1/n(u), \varepsilon/2)$.
 $= 1/\tilde{n}(u)$.

Define:

$$\bar{C} = \sum_{u \in V} \text{cost}(u)$$

Note:

$$\begin{aligned} n(u) &\geq \tilde{n}(u) \\ 1/n(u) &\leq 1/\tilde{n}(u) \end{aligned}$$

Approximate Connected Components

Key Idea 3: Approximate Cost

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\tilde{n}(u) = \min(n(u), 2/\varepsilon)$.

Let $\text{cost}(u) = \max(1/n(u), \varepsilon/2)$.
 $= 1/\tilde{n}(u)$.

Define:

$$\bar{C} = \sum_{u \in V} \text{cost}(u)$$

Note:

$$\begin{aligned} n(u) &\geq \tilde{n}(u) \\ 1/n(u) &\leq 1/\tilde{n}(u) \end{aligned}$$

Approximate Connected Components

Close enough approximation:

$$|\text{CC}(G) - \overline{C}| = \overline{C} - \text{CC}(G)$$

$$\begin{aligned} n(u) &\geq \overline{n}(u) \\ 1/n(u) &\leq 1/\overline{n}(u) \end{aligned}$$

Intuition:

By rounding $\text{cost}(u)$ up to $\varepsilon/2$, we increase the error at most $\varepsilon n/2$.

Approximate Connected Components

Close enough approximation:

$$\begin{aligned} |\text{CC}(G) - \overline{C}| &= \overline{C} - \text{CC}(G) \\ &= \sum_{j=1}^n 1/\overline{n}(u) - \sum_{j=1}^n 1/n(u) \end{aligned}$$

Intuition:

By rounding $\text{cost}(u)$ up to $\varepsilon/2$, we increase the error at most $\varepsilon n/2$.

Approximate Connected Components

Close enough approximation:

$$\begin{aligned} |\text{CC}(G) - \bar{C}| &= \bar{C} - \text{CC}(G) \\ &= \sum_{j=1}^n 1/\bar{n}(u) - \sum_{j=1}^n 1/n(u) \\ &= \sum_{j=1}^n (1/\bar{n}(j) - 1/n(j)) \end{aligned}$$

Intuition:

By rounding $\text{cost}(u)$ up to $\varepsilon/2$, we increase the error at most $\varepsilon n/2$.

Approximate Connected Components

Close enough approximation:

$$\begin{aligned} |\text{CC}(G) - \overline{C}| &= \overline{C} - \text{CC}(G) \\ &= \sum_{j=1}^n 1/\overline{n}(u) - \sum_{j=1}^n 1/n(u) \\ &= \sum_{j=1}^n (1/\overline{n}(j) - 1/n(j)) \\ &\leq \sum_{j=1}^n 1/\overline{n}(j) \end{aligned}$$

Intuition:

By rounding $\text{cost}(u)$ up to $\varepsilon/2$, we increase the error at most $\varepsilon n/2$.

Approximate Connected Components

Close enough approximation:

$$\begin{aligned} |\text{CC}(G) - \overline{C}| &= \overline{C} - \text{CC}(G) \\ &= \sum_{j=1}^n 1/\overline{n}(u) - \sum_{j=1}^n 1/n(u) \\ &= \sum_{j=1}^n (1/\overline{n}(j) - 1/n(j)) \\ &\leq \sum_{j=1}^n 1/\overline{n}(j) \\ &\leq \sum_{j=1}^n \epsilon/2 \end{aligned}$$

Intuition:

By rounding $\text{cost}(u)$ up to $\epsilon/2$, we increase the error at most $\epsilon n/2$.

Approximate Connected Components

Close enough approximation:

$$\begin{aligned} |\text{CC}(G) - \overline{C}| &= \overline{C} - \text{CC}(G) \\ &= \sum_{j=1}^n 1/\overline{n}(u) - \sum_{j=1}^n 1/n(u) \\ &= \sum_{j=1}^n (1/\overline{n}(j) - 1/n(j)) \\ &\leq \sum_{j=1}^n 1/\overline{n}(j) \\ &\leq \sum_{j=1}^n \epsilon/2 \\ &\leq \epsilon n/2 \end{aligned}$$

Intuition:

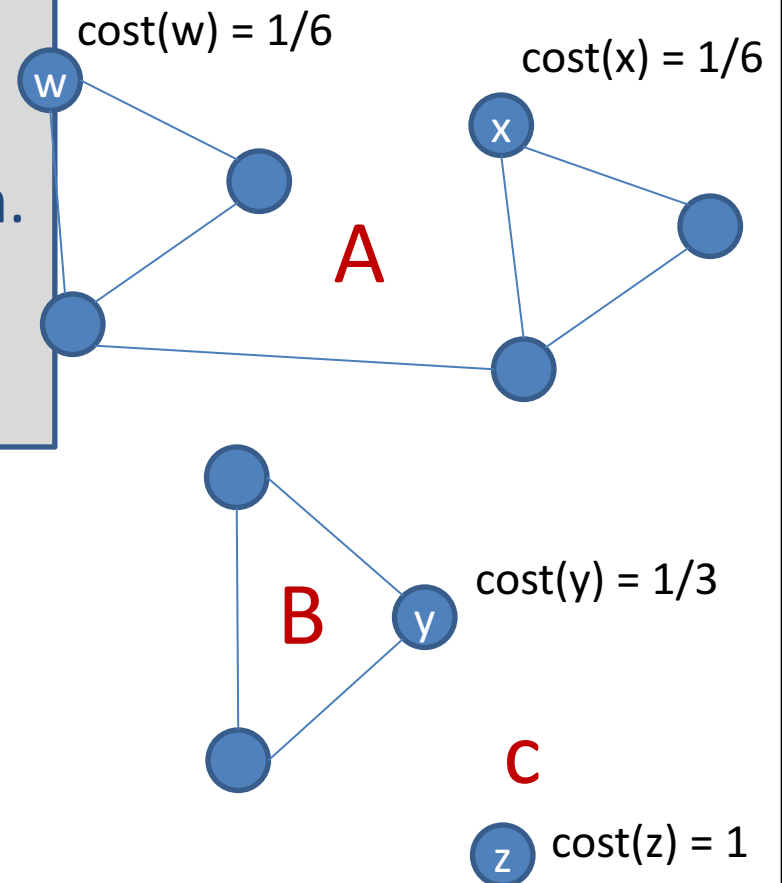
By rounding $\text{cost}(u)$ up to $\epsilon/2$, we increase the error at most $\epsilon n/2$.

Approximate Connected Components

Algorithm 3

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    sum = sum + cost(u)
return n · (sum/s)
```

We have shown:
Sufficient to approximate
 $\text{cost}(u)$ by rounding up.



Approximate Connected Components

Algorithm 3

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\tilde{n}(u) = \min(n(u), 2/\varepsilon)$.

Let $\text{cost}(u) = \max(1/n(u), \varepsilon/2)$.
 $= 1/\tilde{n}(u)$.

How to efficiently compute $\text{cost}(u)$?

Approximate Connected Components

Algorithm 3

Define: per-node cost

Let $n(u)$ = number of nodes in the connected component containing node u .

Let $\tilde{n}(u) = \min(n(u), 2/\varepsilon)$.

Let $\text{cost}(u) = \max(1/n(u), \varepsilon/2)$.
 $= 1/\tilde{n}(u)$.

How to efficiently compute $\text{cost}(u)$?

Approximate Connected Components

Algorithm 3

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    Perform a BFS from u; stop after seeing  $2/\epsilon$  nodes.
    if BFS found  $> 2/\epsilon$  nodes:
        sum = sum +  $\epsilon/2$ 
    else if BFS found  $n(u)$  nodes:
        sum = sum +  $1/n(u)$ 
return  $n \cdot (\text{sum}/s)$ 
```

Approximate Connected Components

Analysis

Goal:

$$\left| \frac{n}{s} \cdot \text{sum} - \overline{C} \right| \leq \epsilon n / 2$$

Approximate Connected Components

Analysis

Goal:

$$\left| \frac{n}{s} \cdot \text{sum} - \overline{C} \right| \leq \epsilon n / 2$$

Implies:

$$\begin{aligned} \left| \frac{n}{s} \cdot \text{sum} - \text{CC}(G) \right| &\leq \left| \frac{n}{s} \cdot \text{sum} - \overline{C} \right| + |\overline{C} - \text{CC}(G)| \\ &\leq \epsilon n / 2 + \epsilon n / 2 \\ &\leq \epsilon n \end{aligned}$$

Approximate Connected Components

Algorithm 3 Analysis

Define random variables: Y_1, Y_2, \dots, Y_s

u_j = node chosen in j^{th} iteration

Y_j = $\text{cost}(u_j)$

Rounded up cost



Approximate Connected Components

Algorithm 3 Analysis

Define random variables: Y_1, Y_2, \dots, Y_s

$$\begin{aligned} E[Y_j] &= \sum_{i=1}^n \frac{1}{n} \text{cost}(u_i) = \frac{1}{n} \sum_{i=1}^n \text{cost}(u_i) \\ &= \frac{1}{n} \overline{C} \end{aligned}$$

Approximate Connected Components

Algorithm 3 Analysis

Unbiased estimator:

$$\begin{aligned} \mathbb{E} \left[\sum_{j=1}^s Y_j \right] &= s \mathbb{E} [Y_j] \\ &= \frac{s}{n} \overline{C} \end{aligned}$$

Approximate Connected Components

Algorithm 3 Analysis

Notice:

Expected output of algorithm is:

$$\mathbb{E} [n \cdot (sum/s)] = \frac{n}{s} \left(\frac{s}{n} \overline{C} \right) = \overline{C}$$

Approximate Connected Components

Algorithm 3 Analysis

Goal:

$$\Pr \left\{ \left| \overline{C} - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} \leq 1/3$$

Approximate Connected Components

Algorithm 3 Analysis

Derivation:

$$\begin{aligned}\Pr \left\{ \left| \overline{C} - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} &= \Pr \left\{ \left| \mathbb{E} \left[\frac{n}{s} \sum_{i=1}^s Y_i \right] - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} \\ &= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \frac{s}{n} \epsilon n / 2 \right\} \\ &= \Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\}\end{aligned}$$

Approximate Connected Components

Algorithm 3 Analysis

Derivation:

$$\Pr \left\{ \left| \bar{C} - \frac{n}{s} \sum_{j=1}^s Y_j \right| > \epsilon n / 2 \right\} =$$

$$\Pr \left\{ \left| \mathbb{E} \left[\sum_{i=1}^s Y_i \right] - \sum_{j=1}^s Y_j \right| > \epsilon s / 2 \right\} \leq 2e^{-2(\epsilon s / 2)^2 / s}$$

$$\leq 2e^{-2\epsilon^2 s / 4}$$

$$\leq 2e^{-\epsilon^2 (4/\epsilon^2) / 2}$$

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

$$< 1/3$$

$$s = \frac{4}{\epsilon^2}$$

Approximate Connected Components

Analysis

Goal:

$$\left| \frac{n}{s} \cdot \text{sum} - \overline{C} \right| \leq \epsilon n / 2$$

Implies:

$$\begin{aligned} \left| \frac{n}{s} \cdot \text{sum} - \text{CC}(G) \right| &\leq \left| \frac{n}{s} \cdot \text{sum} - \overline{C} \right| + |\overline{C} - \text{CC}(G)| \\ &\leq \epsilon n / 2 + \epsilon n / 2 \\ &\leq \epsilon n \end{aligned}$$

Approximate Connected Components

Algorithm 3

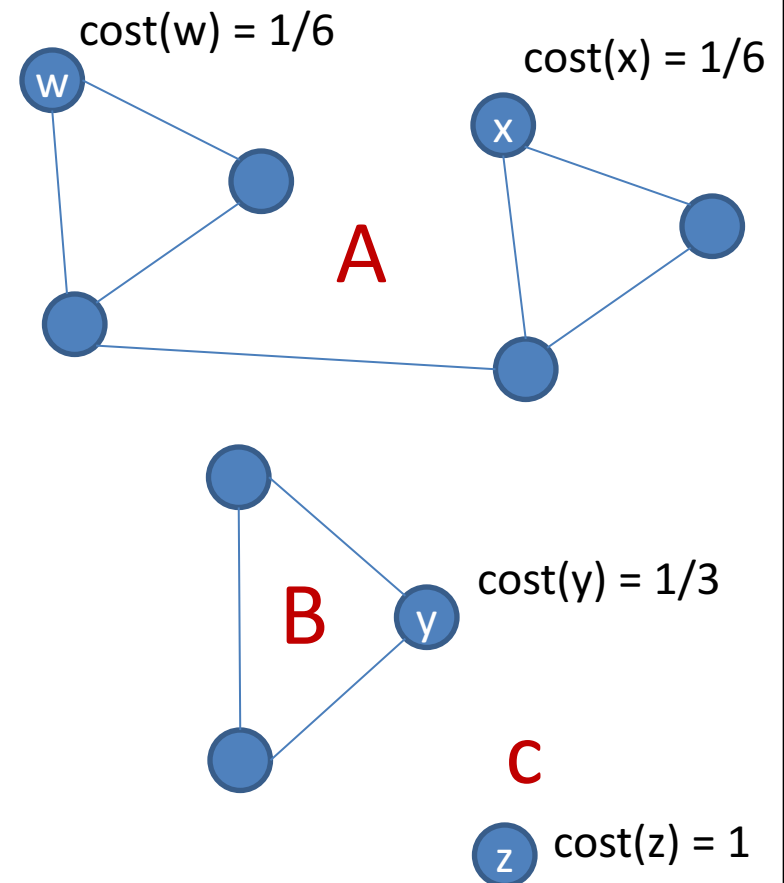
```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    Perform a BFS from u; stop after seeing  $2/\epsilon$  nodes.
    if BFS found  $> 2/\epsilon$  nodes:
        sum = sum +  $\epsilon/2$ 
    else if BFS found  $n(u)$  nodes:
        sum = sum +  $1/n(u)$ 
return  $n \cdot (\text{sum}/s)$ 
```

Approximate Connected Components

Algorithm 3

We have shown:

With probability $> 2/3$,
output is equal to:
 $CC(G) \pm \epsilon n$



Approximate Connected Components

Algorithm 3

Cost of BFS: $O((2/\epsilon) \cdot d)$

sum = 0

for $j = 1$ to s :

Choose u uniformly at random.

Perform a BFS from u ; stop after seeing $2/\epsilon$ nodes.

if BFS found $> 2/\epsilon$ nodes:

sum = sum + $\epsilon/2$

else if BFS found $n(u)$ nodes:

sum = sum + $1/n(u)$

return $n \cdot (\text{sum}/s)$

Approximate Connected Components

Algorithm 3

```
sum = 0
for j = 1 to s:
    Choose u uniformly at random.
    Perform a BFS from u; stop after seeing 2/ε nodes.
    if BFS found > 2/ε nodes:
        sum = sum + ε/2
    else if BFS found n(u) nodes:
        sum = sum + 1/n(u)
return n·(sum/s)
```

Cost of BFS: $O((2/\epsilon) \cdot d)$



Total cost:

$$O(s(2/\epsilon) \cdot d) =$$

$$O((1/\epsilon^2)(2/\epsilon)d) =$$

$$O(d/\epsilon^3)$$

Approximate Connected Components

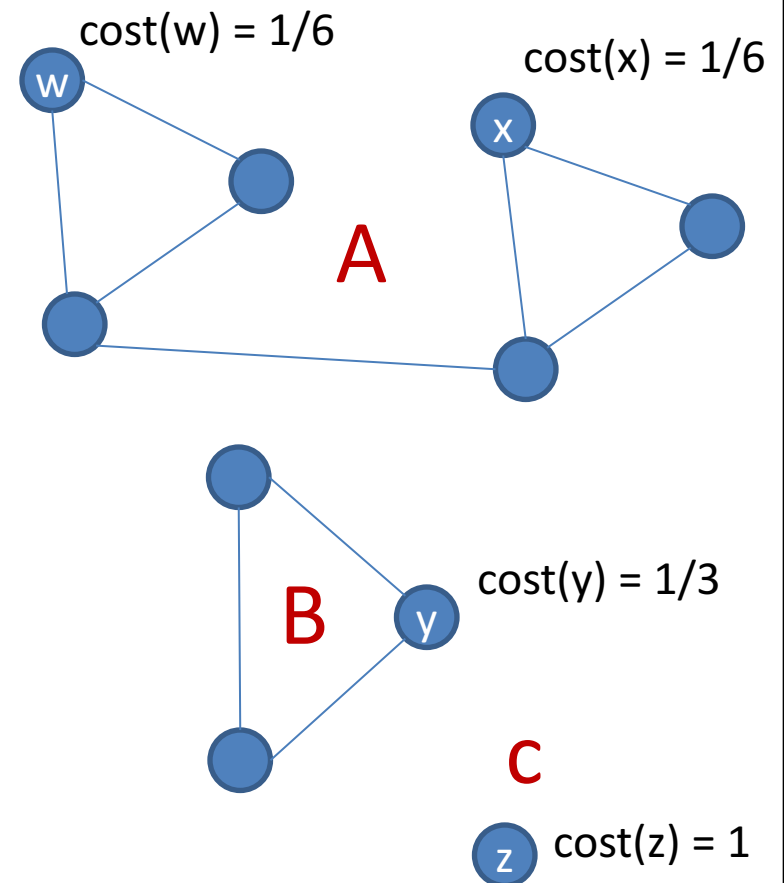
Algorithm 3

We have shown:

With probability $> 2/3$,
output is equal to:

$$CC(G) \pm \epsilon n$$

Running time: $O\left(\frac{d}{\epsilon^3}\right)$



Approximate Connected Components

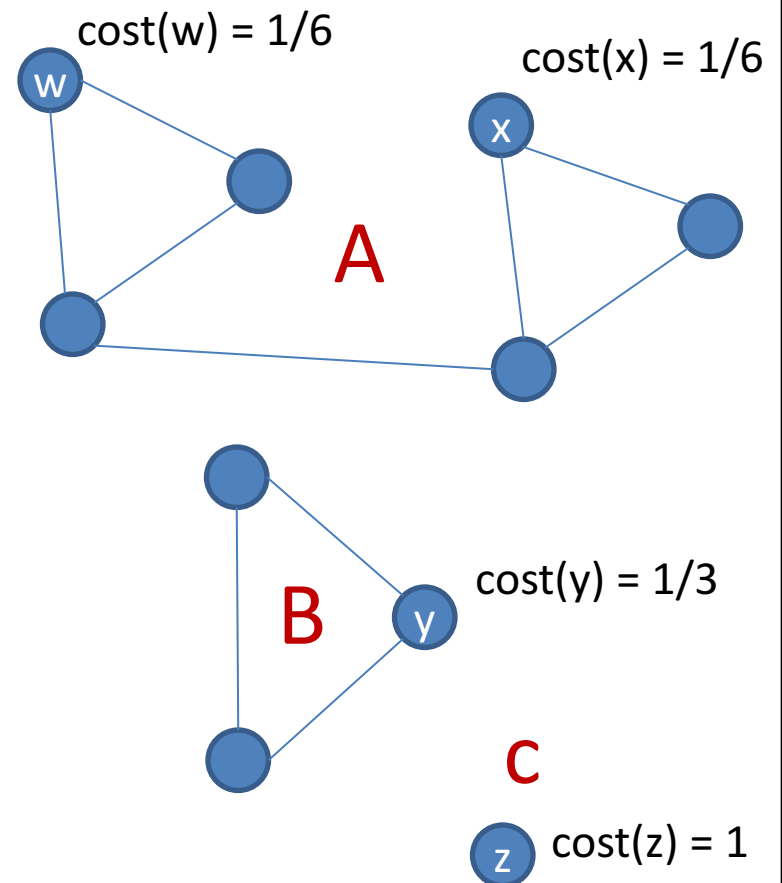
Algorithm 3

We have shown:

With probability $> 1 - 1/\delta$,
output is equal to:

$$CC(G) \pm \epsilon n$$

Running time: $O\left(\frac{d \ln \delta}{\epsilon^3}\right)$



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: Fraction of 1's?

- Additive $\pm \varepsilon$ approximation
- Hoeffding Bound

Is the graph connected?

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

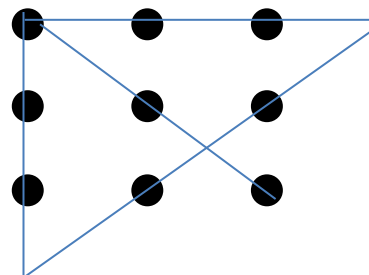
Today:

Number of connected components in a graph.

- Approximation algorithm.

Weight of MST

- Approximation algorithm.



9 dots
4 lines

Today's Problem: Minimum Spanning Tree

Assumptions:

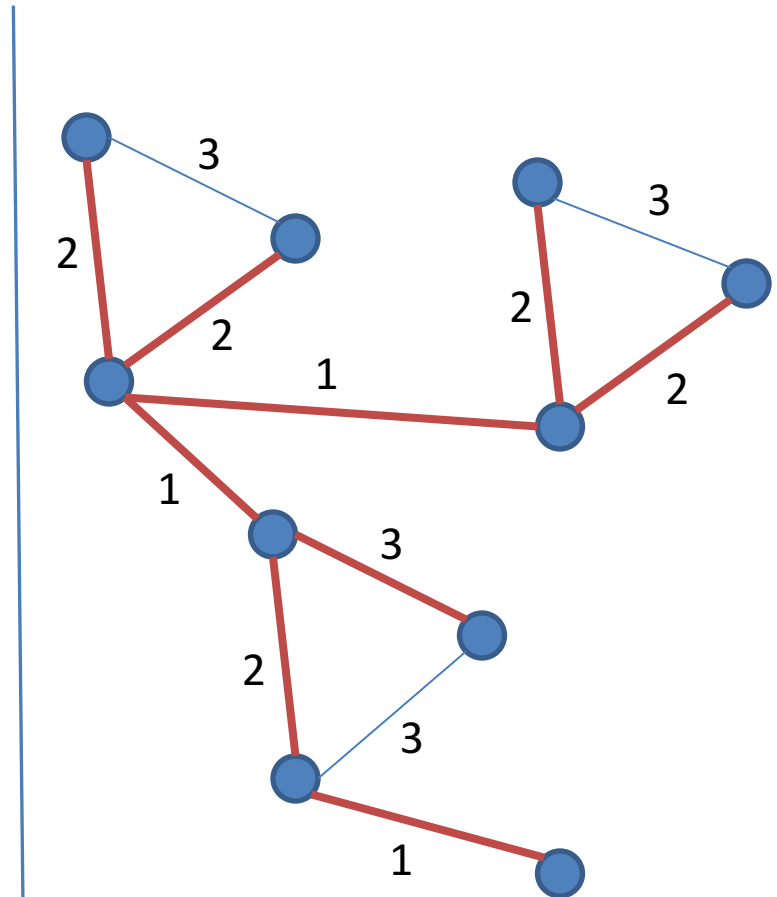
Graph $G = (V, E)$

- Undirected
- Weighted, max weight W
- Connected
- n nodes
- m edges
- maximum degree d

Error term: $\varepsilon < 1/2$

Output:

Weight of MST.



Example: output 16

Today's Problem: Minimum Spanning Tree

Approximation:

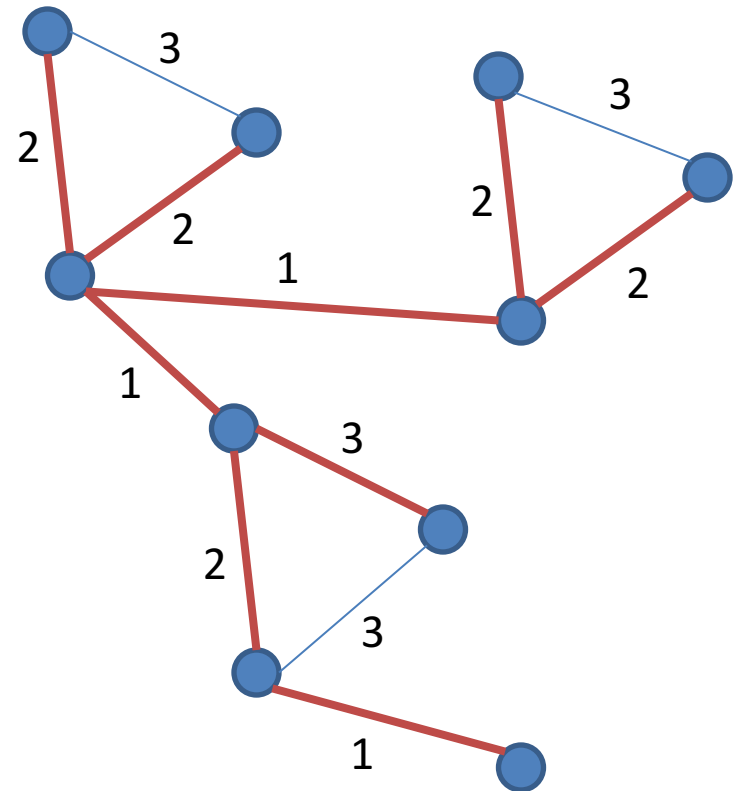
Output **M** such that:

$$\text{MST}(G)(1 - \epsilon) \leq M \leq \text{MST}(1 + \epsilon)$$

Alternate form:

$$M = \text{MST}(G)(1 \pm \epsilon)$$

Correct output: **w.p. > 2/3**



Example:

$$\epsilon = 1/4$$

$$\text{Output} \in [12, 20]$$

Today's Problem: Minimum Spanning Tree

When is this useful?

What are trivial values of ε ?

What are hard values of ε ?

What sort of applications is this useful for?

Why multiplicative approximation for MST and additive approximation for connected components?

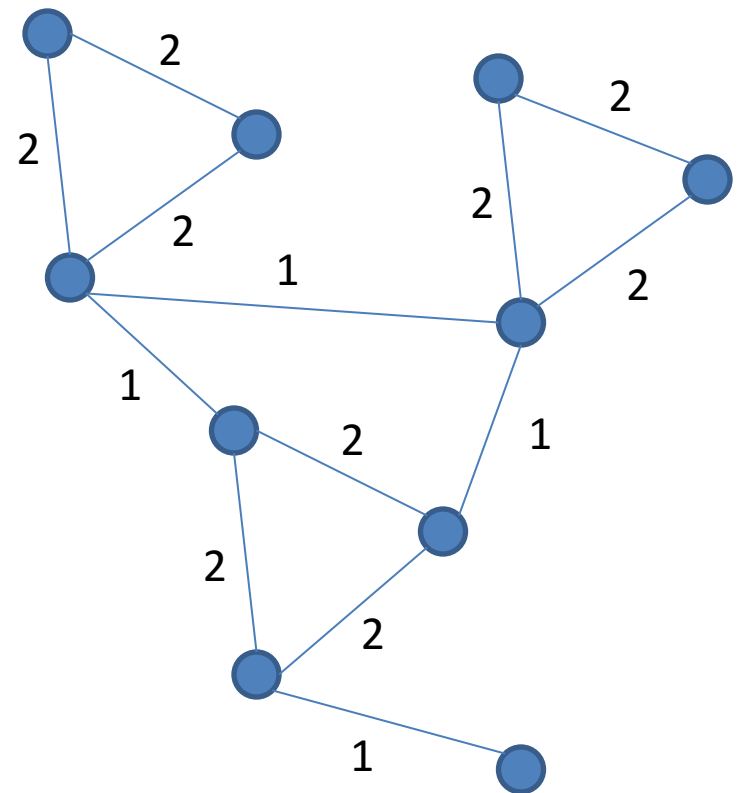
Simple Minimum Spanning Tree

Assume all weights 1 or 2

Which edges must be in MST?

How many weight-2 edges in MST?

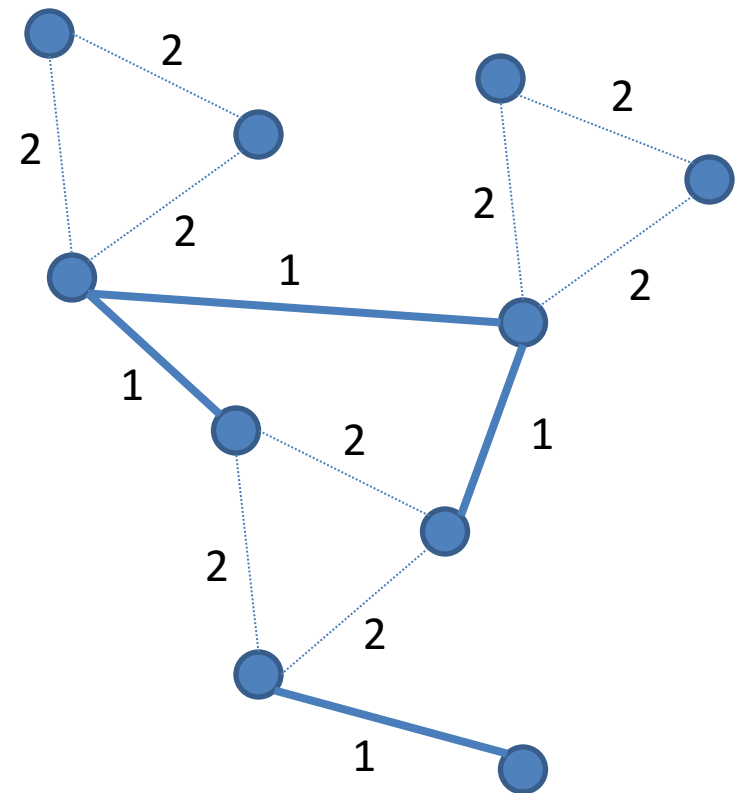
Best (exact) algorithm?



Simple Minimum Spanning Tree

Assume all weights 1 or 2

Let G_1 = graph containing only edges of weight 1.

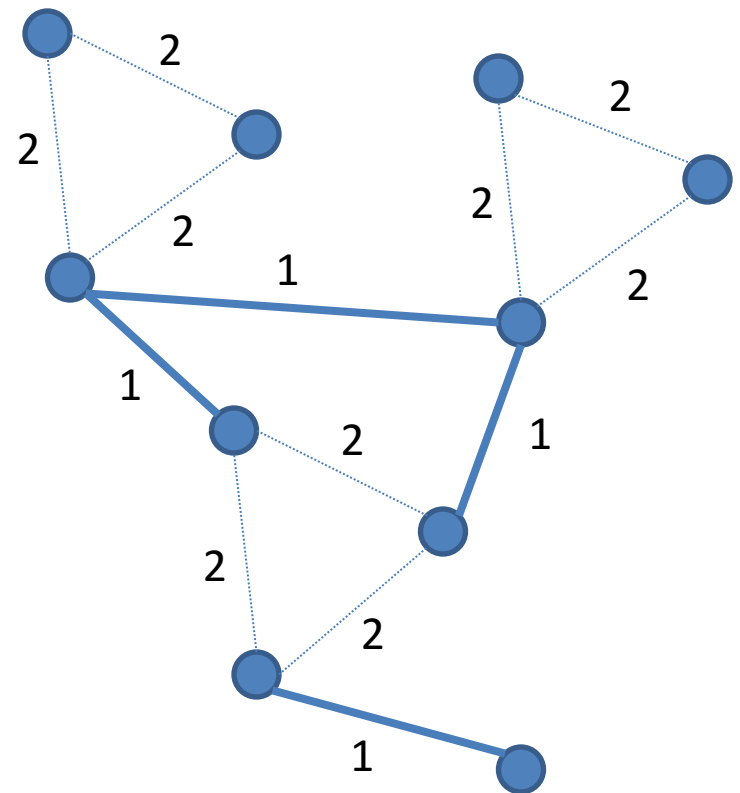


Simple Minimum Spanning Tree

Assume all weights 1 or 2

Let G_1 = graph containing only edges of weight 1.

Let C_1 = number of connected components in G_1 .



Ex: $C_1 = 6$

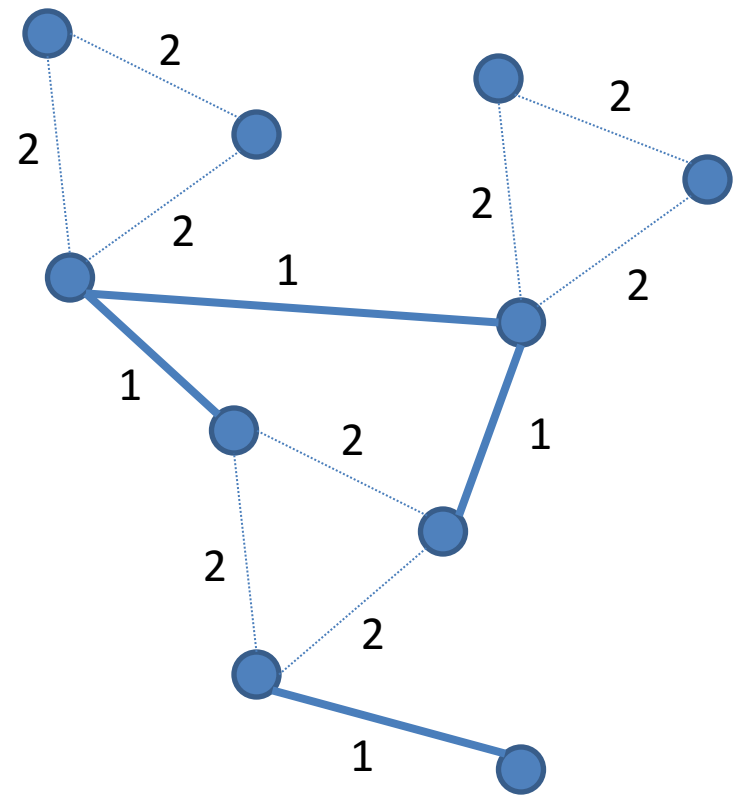
Simple Minimum Spanning Tree

Assume all weights 1 or 2

Let G_1 = graph containing only edges of weight 1.

Let C_1 = number of connected components in G_1 .

Claim: MST contains example $C_1 - 1$ edges of weight 2.



Ex: $C_1 = 6$

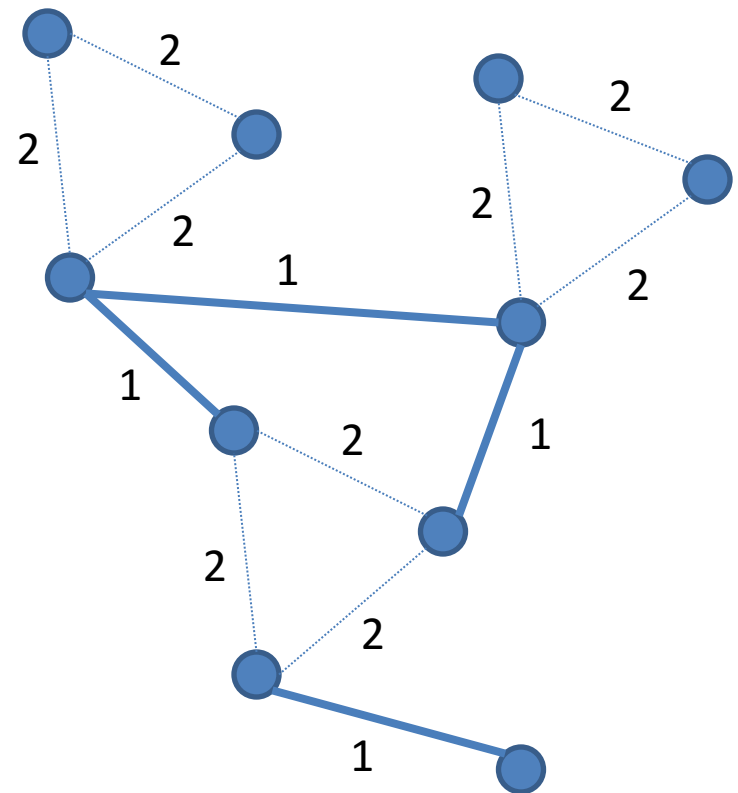
Simple Minimum Spanning Tree

Assume all weights 1 or 2

Claim: MST contains example
 $C_1 - 1$ edges of weight 2.

Basic MST Property:

For any cut, minimum weight edge across cut is in MST.



Ex: $C_1 = 6$

Simple Minimum Spanning Tree

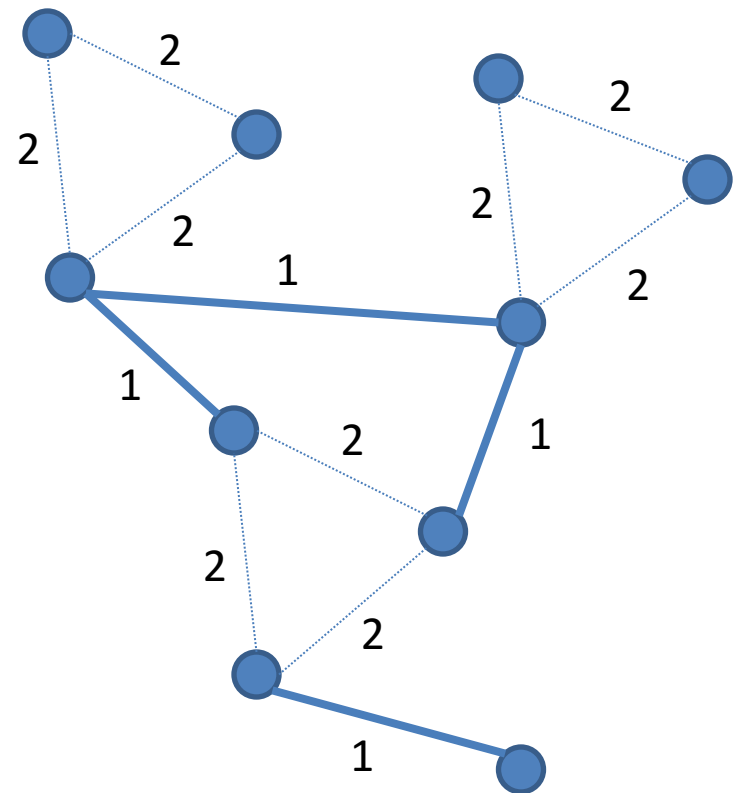
Assume all weights 1 or 2

Claim: MST contains example
 $C_1 - 1$ edges of weight 2.

Algorithm:

For any connected component,
add minimum weight outgoing
edge.

Here all the edges have weight 2,
so add $C_1 - 1$ edges of weight 2.



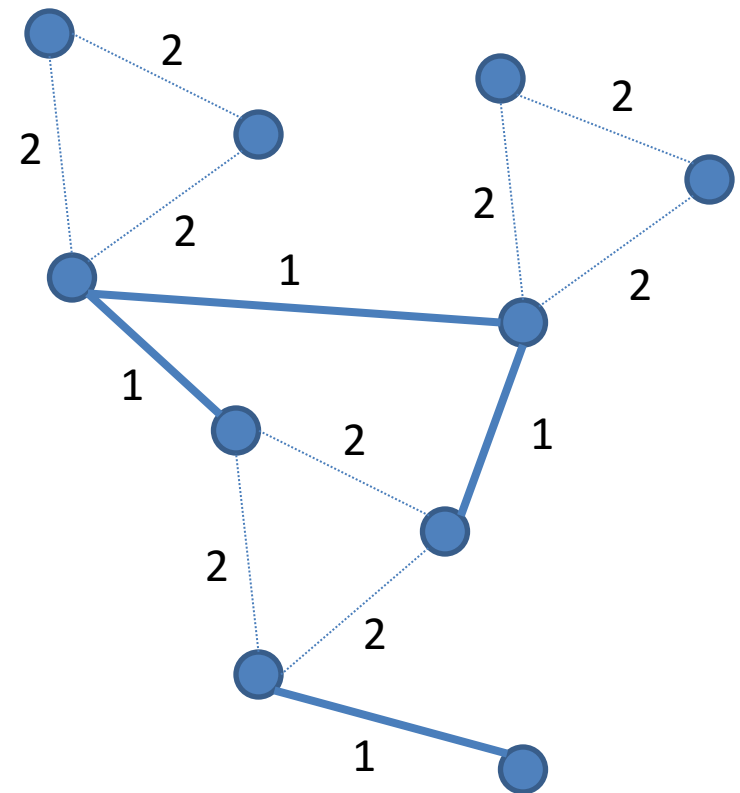
Ex: $C_1 = 6$

Simple Minimum Spanning Tree

Assume all weights 1 or 2

Claim: MST contains example
 $C_1 - 1$ edges of weight 2.

Weight of MST?



Ex: $C_1 = 6$

Simple Minimum Spanning Tree

Assume all weights 1 or 2

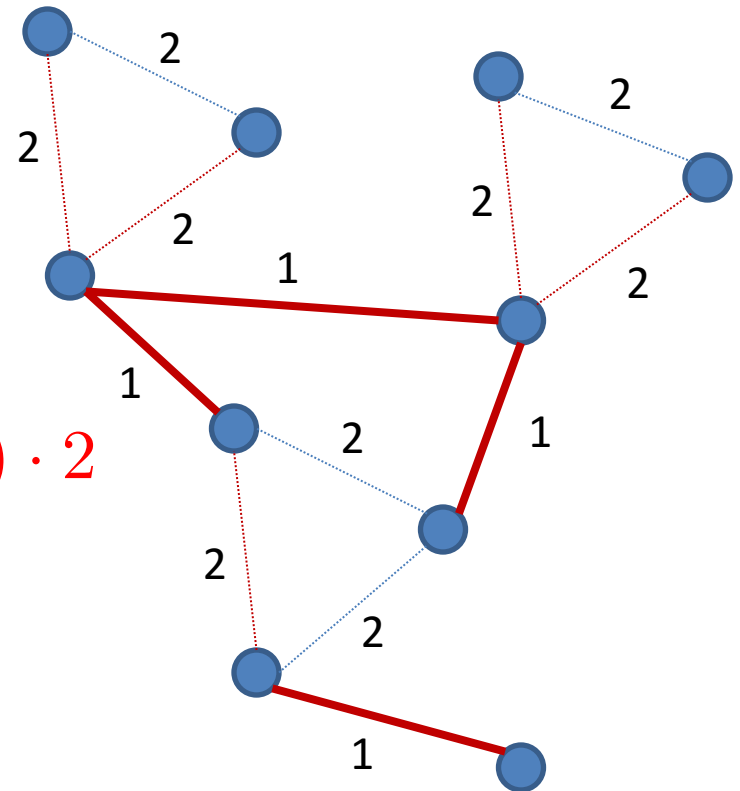
Claim: MST contains example
 $C_1 - 1$ edges of weight 2.

Weight of MST?

$$(n - (C_1 - 1) - 1) \cdot 1 + (C_1 - 1) \cdot 2$$

$$= n + C_1 - 2$$

Ex: $10 + 6 - 2 = 14$



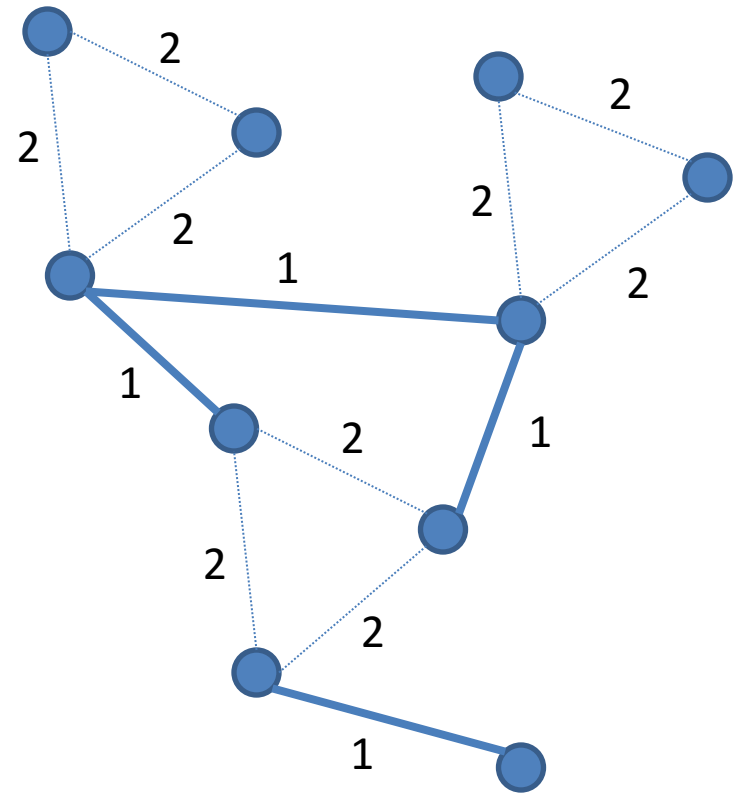
Ex: $C_1 = 6$

Simple Minimum Spanning Tree

Assume all weights 1 or 2

Weight of MST: $n + C_1 - 2$

Algorithm idea?



Ex: $C_1 = 6$

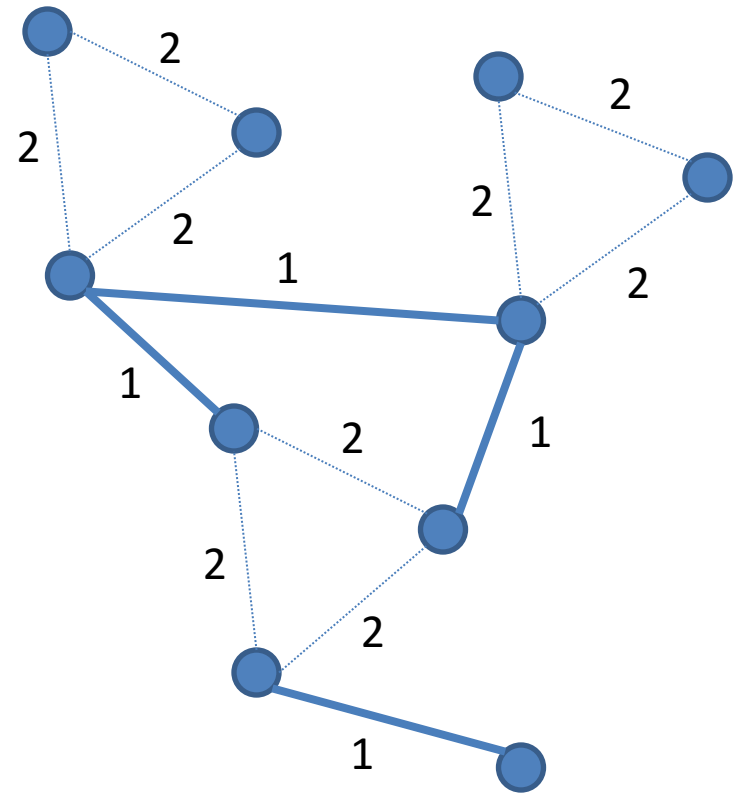
Simple Minimum Spanning Tree

Assume all weights 1 or 2

Weight of MST: $n + C_1 - 2$

Algorithm idea:

Approximate connected components of G_1 .



Ex: $C_1 = 6$

Approximate Minimum Spanning Tree

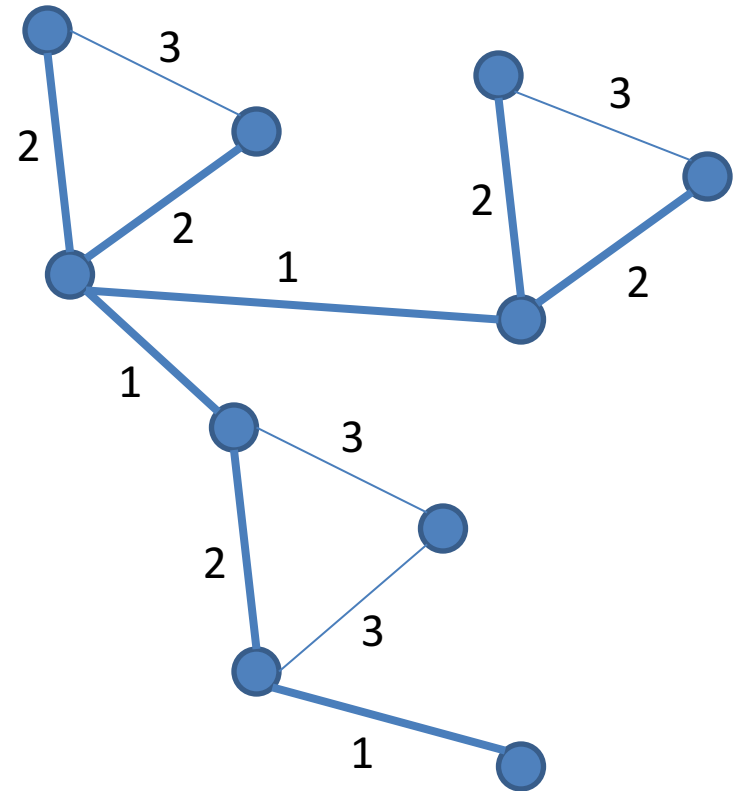
Weights $\{1, 2, \dots, W\}$

Let G_1 = graph containing only edges of weight 1.

Let G_2 = graph containing only edges of weight $\{1, 2\}$.

...

Let G_j = graph containing only edges of weights $\{1, 2, \dots, j\}$.



Ex: G_2

Approximate Minimum Spanning Tree

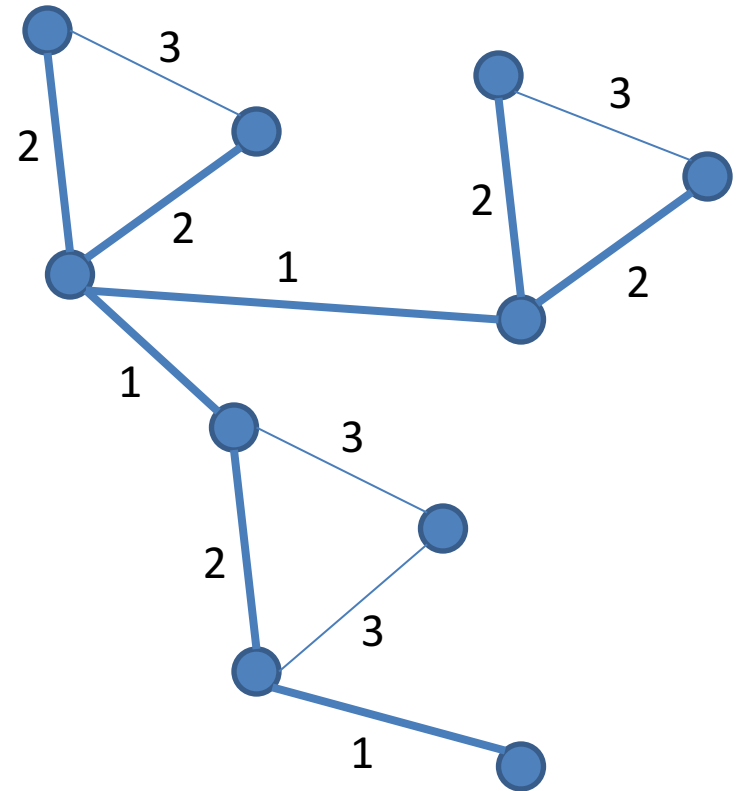
Weights $\{1, 2, \dots, W\}$

Let C_1 = number CC in G_1 .

Let C_2 = number CC in G_2 .

...

Let C_j = number CC in G_j .



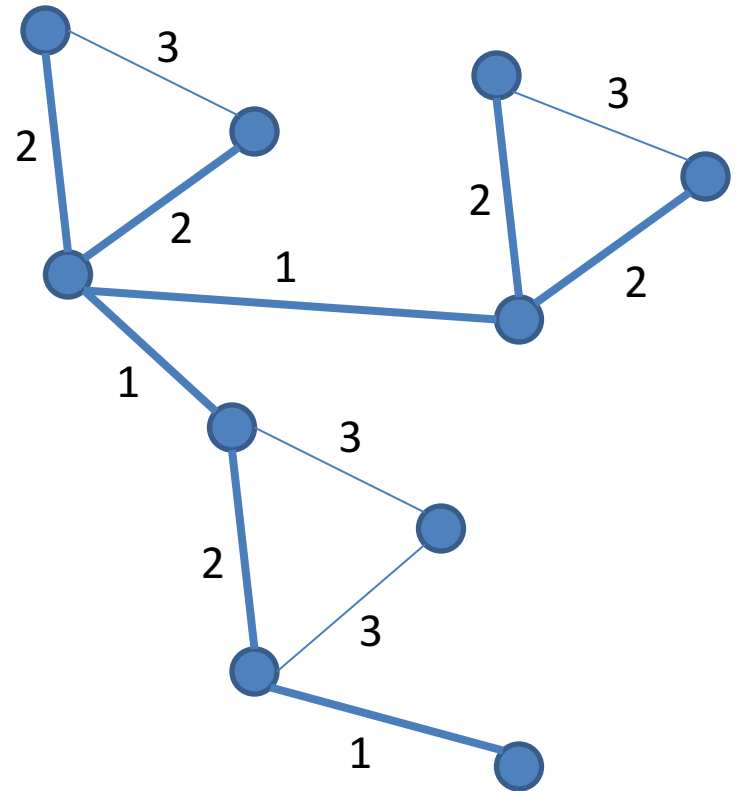
Ex: G_2

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Claim:

MST(G) contains $C_j - 1$ edges
of weight $> j$.



Ex: G_2

Approximate Minimum Spanning Tree

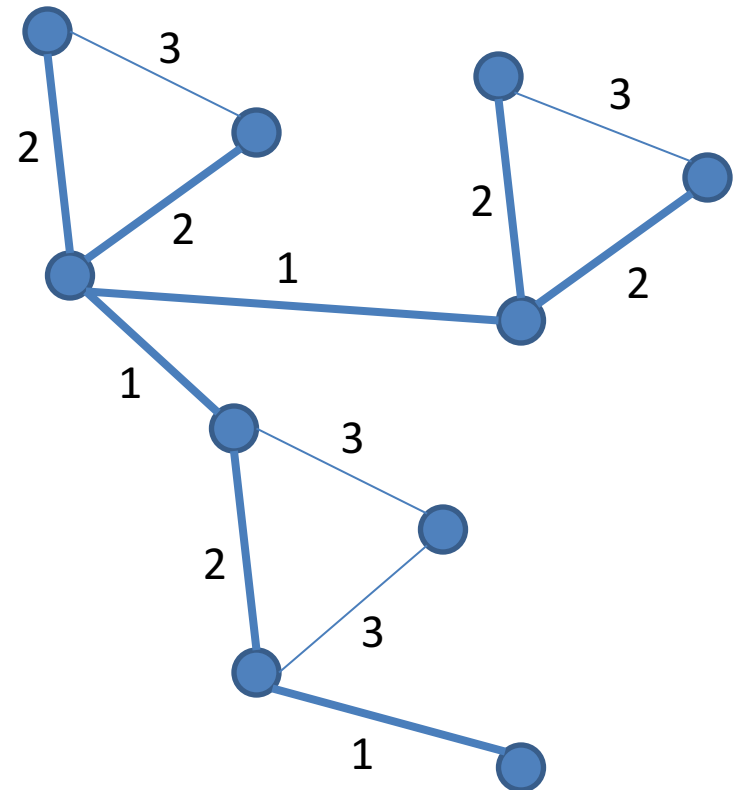
Weights $\{1, 2, \dots, W\}$

Claim:

MST(G) contains $C_j - 1$ edges of weight $> j$.

Why?

There are C_j connected components in G_j . There must be $C_j - 1$ edges connecting them, and each must have weight $> j$.



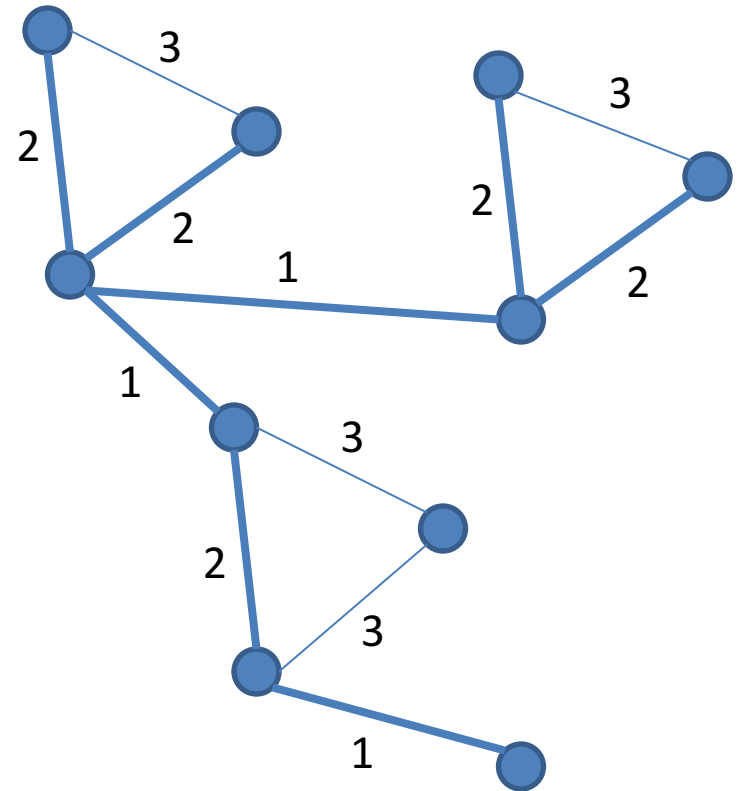
Ex: G_2

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Lemma:

$$\text{MST}(G) = n - W + \sum_{j=1}^{W-1} C_j$$



Ex: G_2

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Edges of weight 1:

$n - 1$ edges total in MST

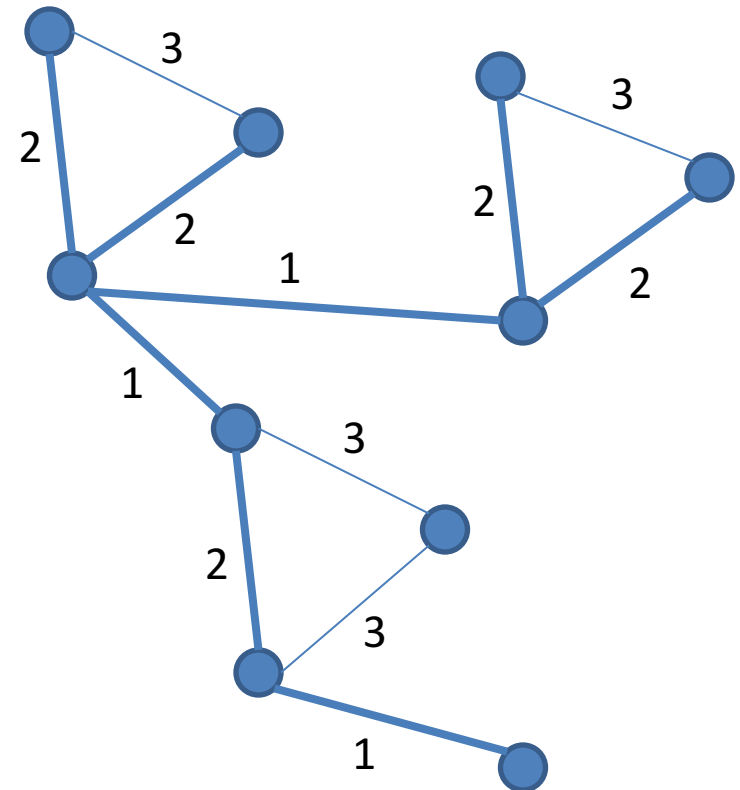
$C_1 - 1$ edges of weight > 1



$(n - 1) - (C_1 - 1)$ edges of weight 1.



$(n - C_1)$ edges of weight 1.



Ex: G_2

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Edges of weight $j+1$:

$C_j - 1$ edges of weight $> j$

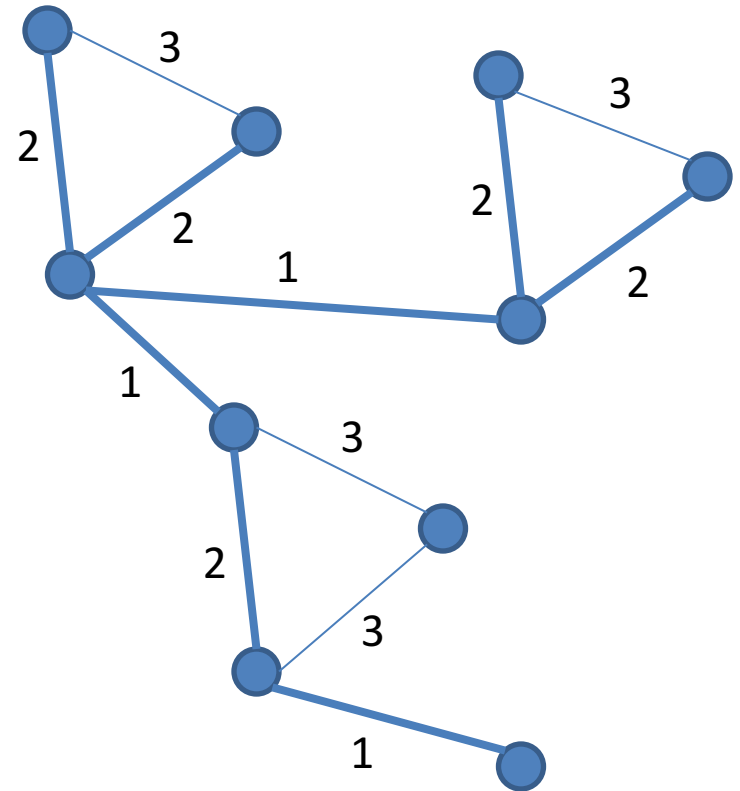
$C_{j+1} - 1$ edges of weight $> j+1$



$(C_j - 1) - (C_{j+1} - 1)$ edges of weight $j+1$.



$(C_j - C_{j+1})$ edges of weight $j+1$.



Note: $C_j \geq C_{j+1}$

Ex: G_2

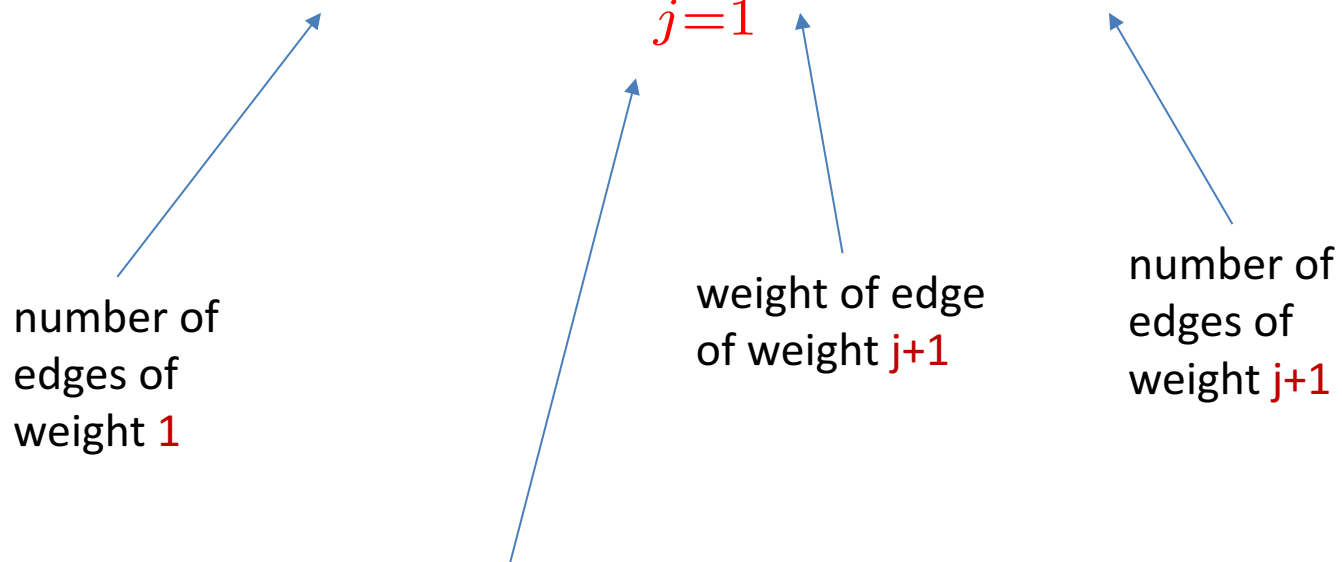
Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Sum the weights:

$$\text{MST}(G) = (n - C_1) + \sum_{j=1}^{W-1} (j+1)(C_j - C_{j+1})$$

number of
edges of
weight 1



weight of edge
of weight $j+1$

number of
edges of
weight $j+1$

Note: sum is from $j = 1$ to $W-1$.

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Sum the weights:

$$\begin{aligned} \text{MST}(G) &= (n - C_1) + \sum_{j=1}^{W-1} (j+1)(C_j - C_{j+1}) \\ &= (n - C_1) + (2C_1 - 2C_2) + (3C_2 - 3C_3) \\ &\quad + (4C_3 - 4C_4) + \dots \\ &\quad + (WC_{W_1} - WC_W) \end{aligned}$$

Approximate Minimum Spanning Tree

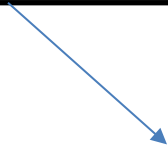
Weights $\{1, 2, \dots, W\}$

Sum the weights:

$$\begin{aligned}\text{MST}(G) &= (n - C_1) + \sum_{j=1}^{W-1} (j+1)(C_j - C_{j+1}) \\ &= (n - C_1) + (2C_1 - 2C_2) + (3C_2 - 3C_3) \\ &\quad + (4C_3 - 4C_4) + \dots \\ &\quad + (WC_{W-1} - WC_W) \\ &= n + C_1 + C_2 + \dots + C_{W-1} - WC_W\end{aligned}$$

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

$$C_W = 1$$


Sum the weights:

$$\text{MST}(G) = n + C_1 + C_2 + \dots + C_{W-1} - WC_W$$

$$= n + C_1 + C_2 + \dots + C_{W-1} - W$$

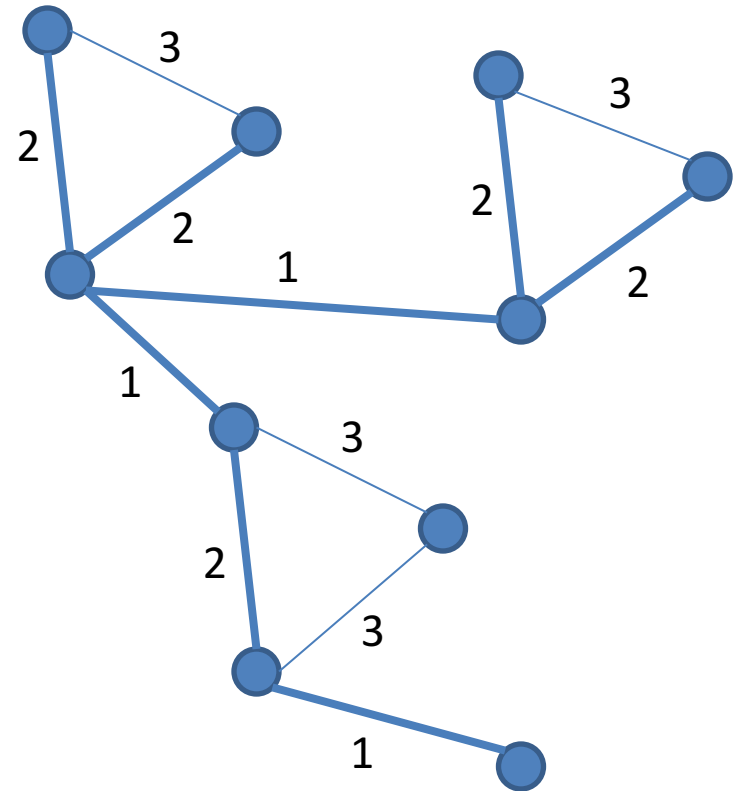
$$= n - W + \sum_{j=1}^{W-1} C_j$$

Approximate Minimum Spanning Tree

Weights $\{1, 2, \dots, W\}$

Lemma:

$$\text{MST}(G) = n - W + \sum_{j=1}^{W-1} C_j$$



Ex: G_2

Approximate Minimum Spanning Tree

Algorithm ApproxMST

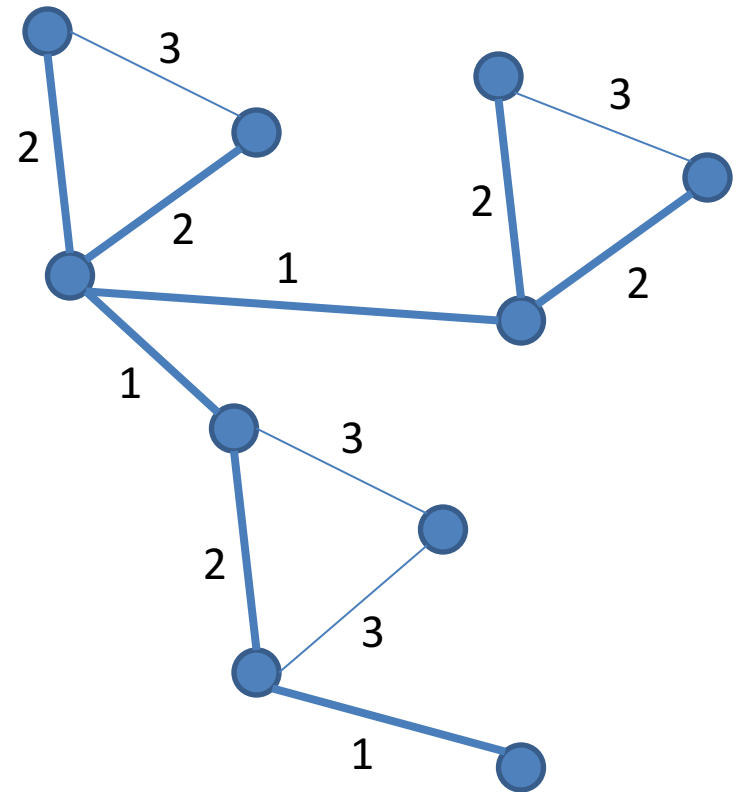
$\text{sum} = n - W$

for $j = 1$ to $W - 1$:

$X_j = \text{AproxCC}(G_j, d, \varepsilon', \delta)$

$\text{sum} = \text{sum} + X_j$

return sum



Ex: G_2

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
     $X_j = \text{AproxCC}(G_j, d, \varepsilon', \delta)$ 
    sum = sum +  $X_j$ 
return sum
```

Set: $\varepsilon' = \varepsilon/W$

Sum of errors: $\leq W(\varepsilon n/W) \leq \varepsilon n$

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
     $X_j = \text{AproxCC}(G_j, d, \epsilon', \delta)$ 
    sum = sum +  $X_j$ 
return sum
```

Guarantee for each AproxCC:

$$\Pr \{ |X_j - C_j| > \epsilon n / W \} < 1/3$$

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
    Xj = AproxCC(Gj, d, ε', δ)
    sum = sum + Xj
return sum
```

Guarantee for each AproxCC:

$$\Pr \{ |X_j - C_j| > \epsilon n / W \} < 1/3$$

Not good enough: $\Pr\{\text{all correct}\} \cong (2/3)^W$

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
    Xj = AproxCC(Gj, d, ε', δ)
    sum = sum + Xj
return sum
```

Set $\varepsilon' = \varepsilon/W$, $\delta = 1/(3W)$

Error probability: $\Pr \{\text{any fails}\} \leq \sum_{j=1}^{W-1} \frac{1}{3W}$

$$\leq \frac{W-1}{3W} < 1/3$$

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
    Xj = AproxCC(Gj, d, ε', δ)
    sum = sum + Xj
return sum
```

Set $\varepsilon' = \varepsilon/W$, $\delta = 1/(3W)$

Guarantee for each AproxCC:

$$\Pr \{ |X_j - C_j| > \varepsilon n / W \} < \frac{1}{3W}$$

Approximate Minimum Spanning Tree

Error Calculation

```
sum = n - W
for j = 1 to W - 1:
     $X_j = \text{AproxCC}(G_j, d, \epsilon', \delta)$ 
    sum = sum +  $X_j$ 
return sum
```

Set: $\epsilon' = \epsilon/W$, $\delta = 1/(3W)$

Sum of errors: $\leq W(\epsilon n/W) \leq \epsilon n$

→ $\text{MST}(G) - \epsilon n \leq \text{sum} \leq \text{MST}(G) + \epsilon n$

Approximate Minimum Spanning Tree

Error Calculation

$$\text{MST}(G) \geq n - 1 \geq n/2$$

Approximate Minimum Spanning Tree

Error Calculation

$$\text{MST}(G) \geq n - 1 \geq n/2$$

$$\text{MST}(G) - \epsilon n \leq \text{sum} \leq \text{MST}(G) + \epsilon n$$

Approximate Minimum Spanning Tree

Error Calculation

$$\text{MST}(G) \geq n - 1 \geq n/2$$

$$\text{MST}(G) - \epsilon n \leq \text{sum} \leq \text{MST}(G) + \epsilon n$$

$$\begin{aligned} \text{MST}(G) + \epsilon n &\leq \text{MST}(G) + \epsilon(2\text{MST}(G)) \\ &\leq \text{MST}(G)(1 + 2\epsilon) \end{aligned}$$

Approximate Minimum Spanning Tree

Error Calculation

$$\text{MST}(G) \geq n - 1 \geq n/2$$

$$\text{MST}(G) - \epsilon n \leq \text{sum} \leq \text{MST}(G) + \epsilon n$$

$$\begin{aligned} \text{MST}(G) + \epsilon n &\leq \text{MST}(G) + \epsilon(2\text{MST}(G)) \\ &\leq \text{MST}(G)(1 + 2\epsilon) \end{aligned}$$

$$\begin{aligned} \text{MST}(G) - \epsilon n &\geq \text{MST}(G) - \epsilon(2\text{MST}(G)) \\ &\geq \text{MST}(G)(1 - 2\epsilon) \end{aligned}$$

Approximate Minimum Spanning Tree

Error Calculation

$$\text{MST}(G) \geq n - 1 \geq n/2$$

$$\text{MST}(G) - \epsilon n \leq \text{sum} \leq \text{MST}(G) + \epsilon n$$

$$\begin{aligned} \text{MST}(G) + \epsilon n &\leq \text{MST}(G) + \epsilon(2\text{MST}(G)) \\ &\leq \text{MST}(G)(1 + 2\epsilon) \end{aligned}$$

$$\begin{aligned} \text{MST}(G) - \epsilon n &\geq \text{MST}(G) - \epsilon(2\text{MST}(G)) \\ &\geq \text{MST}(G)(1 - 2\epsilon) \end{aligned}$$

$$\text{MST}(G)(1 - 2\epsilon) \leq \text{MST}(G) \leq \text{MST}(G)(1 + 2\epsilon)$$

Approximate Minimum Spanning Tree

Running time

```
sum = n - W
for j = 1 to W - 1:
    Xj = AproxCC(Gj, d, ε', δ)
    sum = sum + Xj
return sum
```

Set $\epsilon' = \epsilon/W$, $\delta = 1/(3W)$

Running time: $O\left(W \cdot \frac{d \ln(1/(1/3W))}{(\epsilon/W)^3}\right)$

Approximate Minimum Spanning Tree

Running Time

```
sum = n - W
for j = 1 to W - 1:
    Xj = AproxCC(Gj, d, ε', δ)
    sum = sum + Xj
return sum
```

Set $\varepsilon' = \varepsilon/W$, $\delta = 1/(3W)$

Running time: $O\left(W \cdot \frac{d \ln(1/(1/3W))}{(\varepsilon/W)^3}\right) = O\left(\frac{dW^4 \log W}{\varepsilon^3}\right)$

Approximate MST

Summary

We have shown:

With probability $> 2/3$, output is equal to:
 $\text{MST}(G)(1 \pm \epsilon n)$

Running time:

$$O\left(\frac{dW^4 \log W}{\epsilon^3}\right)$$

Approximate MST

Summary

Note:

See: Chazelle, Rubinfeld, Trevisan

Impossible to do better than:

$$\Omega\left(\frac{dW}{\epsilon^2}\right)$$

Best known:

$$O\left(\frac{dW}{\epsilon^2} \log \frac{dW}{\epsilon}\right)$$

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: Fraction of 1's?

- Additive $\pm \varepsilon$ approximation
- Hoeffding Bound

Is the graph connected?

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

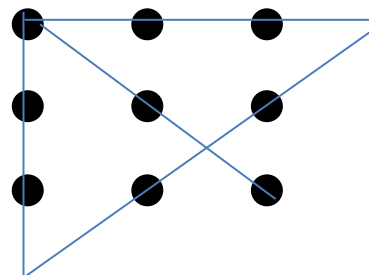
Today:

Number of connected components in a graph.

- Approximation algorithm.

Weight of MST

- Approximation algorithm.



9 dots
4 lines