

## Minimum Spanning Tree (MST)

Input: (undirected connected)  $G = (V, E, c)$

$c_e$  is the cost of the edge  $e$ ,

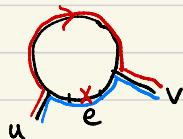
$c_e > 0$  and w.l.o.g all  $c_e$ 's are distinct.

Output: a subset  $T \subseteq E$  so that  $G = \langle V, T \rangle$  is connected and the total cost  $\sum_{e \in T} c_e$  is minimized.

Claim (4.1b):  $T$  is a tree.

Proof: suppose  $T$  contains a cycle  $C$ ,

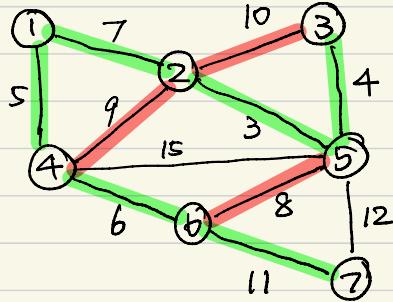
then if we any edge  $e \in C$ , the graph  $(V, T - \{e\})$  is still connected.



## Kruskal's algorithm:

1. Start with  $T = \emptyset$ .
2. Sort edges in  $E$  by cost in ascending order.
3. Iterate over all edges, add  $e$  to  $T$  as long as it does not create a cycle.

Example:



—  $\in T$   
—  $\notin T$

Output:

$T = \text{all green edges}$

$$\begin{aligned} \text{cost}(T) &= 3 + 4 + 5 + 6 + 7 + 11 \\ &= 36. \end{aligned}$$

## Correctness:

Theorem (4.18): Kruskal's algorithm produces an MST.  
(page 146)

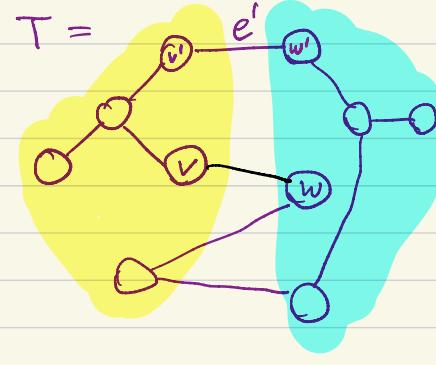
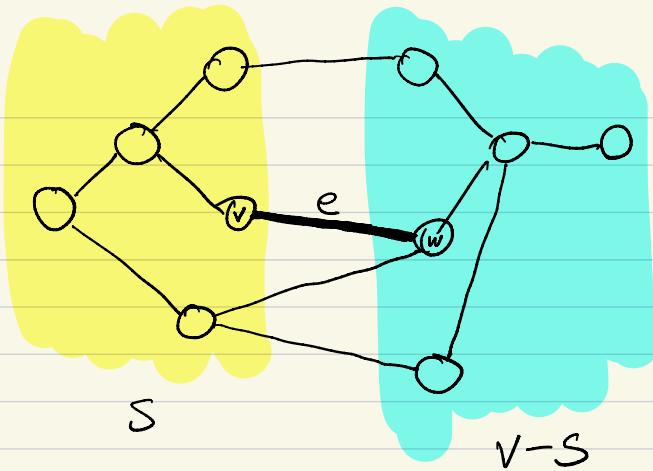
Lemma (4.17): Assume  $c_e$ 's are distinct.

Fix  $S \subseteq V$ .

Let  $e$  be the minimum cost edge with one end in  $S$  and one end in  $V-S$ .

Then  $e$  must appear in every MST.

Proof.



Suppose  $e = (v, w)$  is the cheapest edge between  $S$  and  $V-S$ .

Suppose there is an MST  $T$  s.t.  $e \notin T$ .

Because  $T$  is a spanning tree,

there is a path  $P$  from  $v$  to  $w$  in  $T$ .

Since  $v \in S$  and  $w \notin S$ ,

$P$  must leave  $S$  at some point.

Let  $e'$  be the first edge in  $P$  that leaves

$(v', w')$

$S$ .

$T + \{e\} - \{e'\}$  is a spanning tree with smaller cost.

- the cost is smaller because  $C_e < C_{e'}$ .

- $T + \{e\}$  contains a unique cycle.  $C$

and moreover  $e' \in C$ .

$\Rightarrow T + \{e\} - \{e'\}$  is connected.

" $T + \{e\} - \{e'\}$ " has  $(n-1)$  edges.

Proof (correctness of Kruskal):

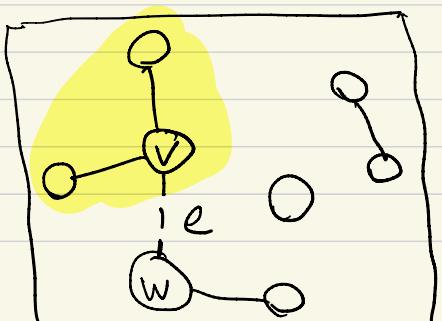
Consider an edge  $e = (v, w)$  added by Kruskal's algorithm.

Let  $S$  be  $v$ 's connected component just before adding  $e$ .

- $w \notin S$  (o.w. adding  $e$  creates a cycle).
- no edge from  $S$  to  $V-S$  has been considered (any such edge would have been added).

Thus,  $e$  must be the cheapest edge that leaves  $S$ .  
so it is safe to include  $e$ .

Moreover, Kruskal's output must be connected.



## Prim's algorithm.

1. Pick a root  $s \in V$ .
2.  $S = \{s\}$ . (All nodes in  $S$  are connected.)
3. Repeat the following until  $S = V$ :

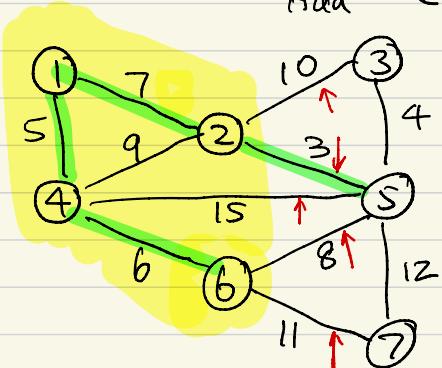
Add  $v$  to  $S$  where  $v \notin S$  minimizes the "attachment cost".

(breaking ties arbitrarily)

min  
 $e = (u, v), u \in S$   
 $e = (u, v)$  to  $T$ .

4. Return  $T$ .

Example:



## Correctness:

- Every time the algorithm adds an edge  $e$ ,  $e$  is the cheapest edge between  $S$  and  $V - S$ .
- $T$  is a spanning tree at the end.

Runtime of Prim :  $O(m \log n)$ .

Similar to Dijkstra's algorithm.

We will use a min-heap  $H$ .

Overall.  $\begin{cases} O(n) & H.\text{min}() \quad \text{each call takes } O(1) \text{ time} \\ O(m) & H.\text{changekey}() \quad \text{each call takes } O(\log n) \text{ time} \end{cases}$

$$\text{Overall runtime} = O(n \cdot 1 + m \cdot \log n) \\ = O(m \log n).$$

## Runtime of Kruskal:

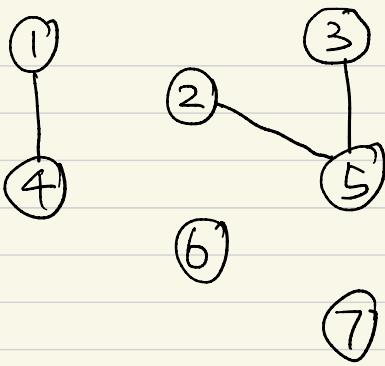
1. Sort all edges:  $O(m \log m) = O(m \log n)$  ( $m \leq n^2$ )
2. Check (in total  $O(m)$  times) whether adding  $e$  would cause a cycle. **How?** ( $O(n)$  via DFS)

Can we do it faster?

## The Union-Find data structure. (page 152)

- $\text{Init}()$ : all elements are in separate sets.
- $\text{Union}(x, y)$ : merge the set containing  $x$  with the set containing  $y$ .
- $\text{Find}(x)$ : return the "name" of the set that  $x$  is in.

We use this data structure to maintain the set of connected components (which are changing over time) in Kruskal's algorithm.



representative element

When we consider any edge  $e = (u, v)$

If ( $\text{find}(u) \neq \text{find}(v)$ ) then

Add  $e$  to  $T$

$\text{union}(u, v)$ .

End if.

### Union-Find First Attempt.

arrays  $\text{name}[i]$ : name of the set that contains  $i$ .  
 $\text{size}[j]$ : size of the set with name  $j$ .

functions  $\text{Init}()$ :  $\text{name}[i] = i$ ,  $\text{size}[i] = 1$ ;

$\text{Find}(x)$ : return  $\text{name}[x]$ ;

$\text{Union}(x, y)$ :  $x = \text{name}[x]$ ;  
 $y = \text{name}[y]$ ;

If ( $\text{size}[x] < \text{size}[y]$ ) swap( $x, y$ );  
(so w.l.o.g  $\text{size}[x] \geq \text{size}[y]$ )

For every element  $i$  in set "y":

$\text{name}[i] = x$ ;

End For

$\text{size}[x] = \text{size}[x] + \text{size}[y]$ ;

Example: Before

	1	2	3	4	5	6	7	(index)
name :	1	2	2	1	2	6	7	
size :	2	3				1	1	

$\text{Union}(4, 5)$        $x=4$ ,  $y=5 \Rightarrow x=1$      $y=2$ .

$x=2$      $y=1$      $\Downarrow$     (swapped so that  $\text{size}[x] \geq \text{size}[y]$ )

for  $i \in \{1, 4\}$ .

End for     $\text{name}[i] = 2$

$\text{size}[2] = \text{size}[2] + \text{size}[1]$ .

After.

	1	2	3	4	5	6	7	(index)
name	2	2	2	2	2	6	7	
size	5					1	1	

We always merge the smaller set into the bigger one!

Runtime: Init:  $O(n)$ . Find :  $O(1)$

Claim: The first  $K$  union operations run in total time  $O(K \log K)$ .

$\Rightarrow$  Kruskal runs in time:  $O(m \log n + m + n \log n) = O(m \log n)$

### Proof sketch:

After  $k > 1$  union operations, the largest set has size  $O(k)$ .

How many times can  $\text{name}[i]$  change?

Once per  $\text{union}(A, B)$  if  $i \in B$  and  $|B| \leq |A|$ .

At most  $O(\log k)$  times because the size of the set  $\text{name}[i]$  at least doubles each time.

At most  $2k$   $\text{name}[i]$  are changed.

$$O(k \log k).$$

Not very ideal because sometimes  $\text{Union}(x, y)$  can take  $\Omega(n)$  time.

Example:  $\text{name} = \underbrace{11 \dots 11}_{n/2} \quad \underbrace{22 \dots 22}_{n/2}$

### Union-Find Second Attempt



index	1	2	3	4	5	6	7	(- if no parent)
parent	-1	-1	2	1	2	-1	-1	

$\text{Init}()$ :  $\text{parent}[i] = -1 \quad \forall i$ .

$\text{Find}(x)$ : return the root of  $x$ .

$\text{Find}(x)$ :

```
while ( $\text{parent}[x] \neq -1$ )
     $x = \text{parent}[x]$ 
end while
return  $x$ 
```

$\text{Union}(x, y)$ :  $x = \text{Find}(x)$ ,  $y = \text{Find}(y)$ .

merge the smaller set into the bigger one.

If  $\text{size}[x] \geq \text{size}[y]$

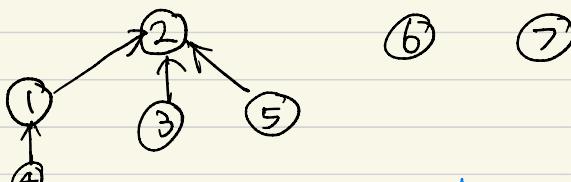
$\text{parent}[y] = x$ :  $\text{size}[x] = \text{size}[x] + \text{size}[y]$

Else

$\text{parent}[x] = y$ ,

$\text{size} \dots$

Example:  $\text{Union}(4, 5)$ .



Runtime of

Kruskal under this Union-Find implementation =  $O(m \log n + m \log n + n)$

Runtime:

$\text{Init}()$ :  $O(n)$

$\text{Find}()$ :  $O(\log n)$

the depth of any tree is at most  $O(\log n)$

$\text{Union}()$ :  $O(1)$

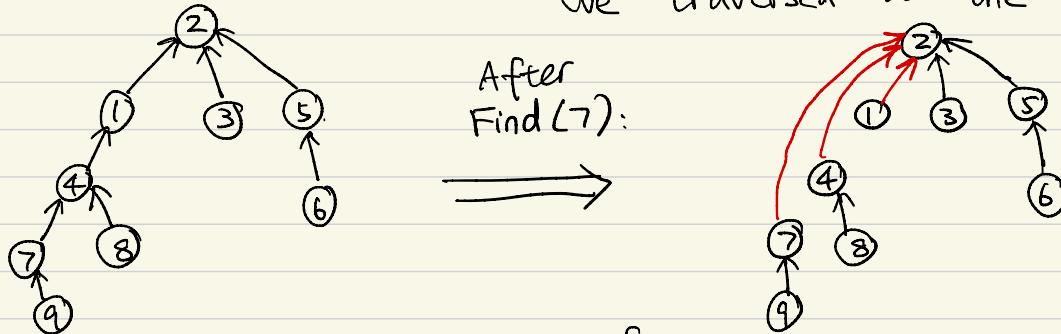
after  $x = \text{find}(x)$  and  $y = \text{find}(y)$ .

sorting find union

## Union-find Third Attempt (the real version)

Building on the second attempt, we add one more optimization.

Path compression: After  $\text{find}(x)$ , point all the elements we traversed to the root.



```
Implementation: int Find(x) {
    If ( $x == \text{parent}[x]$ ) return  $x$ ;
    root = Find( $\text{parent}[x]$ );  $\leftarrow$  recursion.
     $\text{parent}[x] = \text{root}$ ;
    return  $\text{root}$ 
}
```

$\text{find}(7) \rightarrow \text{find}(4) \rightarrow \text{find}(1) \rightarrow \text{find}(2) = 2$   
 $\text{parent}[7]=2 \leftarrow \text{parent}[4]=2 \leftarrow \text{parent}[1]=2 \leftarrow$

Upshot: Next time we call  $\text{find}(7)$  before any union operation takes  $O(1)$  time.

Runtime: First  $K$  Union/Find operations

[Hopcroft and Ullman '73]

A "simple" proof (available on Union-Find Wikipedia page) shows that these  $K$  operations runs in  $O(K \cdot \log^* n)$ .

$\log^* n$ : iterated logarithm: # of times log needs to be applied for  $n$  to become 1.

$$\log_2^*(2^{2^{2^{2^2}}}) = 5$$

??  
19728  
10

$O(m \cdot \alpha(n))$  [Tarjan '75]

$\uparrow$   
inverse Ackermann function (grows even slower than  $\log^* n$ ).

$\Theta(m \cdot \alpha(n))$  [Fredman and Saks '89]