



# Algorithms

**MIDTERM EXAM:** Thursday March 9

Lectures 1 to 10

Office Hours today: 1:30-2:30

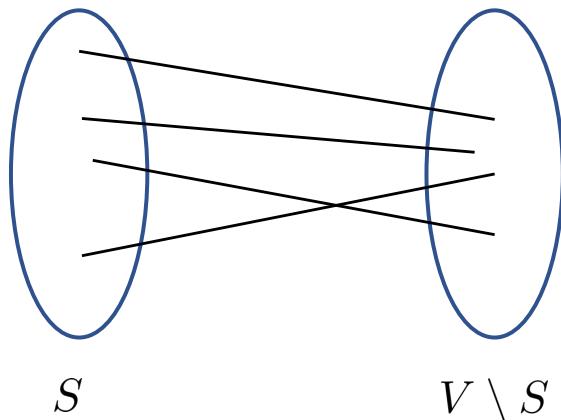
## MINIMUM SPANNING TREE (MST)

*Connected*

Let  $G = (V, E)$  be an input graph.

For any subset  $S \subset V$ , consider the two groups of vertices  $S$  and  $V \setminus S$ .

Such a partition of  $V$  is called a ***cut*** and denoted  $(S, V \setminus S)$ .



The edges going *across* the cut are called ***cut edges***.

Since the graph is connected, there are cut edges.

**Claim.** The lightest cut edge belongs to the MST.      ***Cut property.***

# GREEDY ALGORITHMS FOR MST

## Prim's Algorithm.

Initial tree = any single vertex.

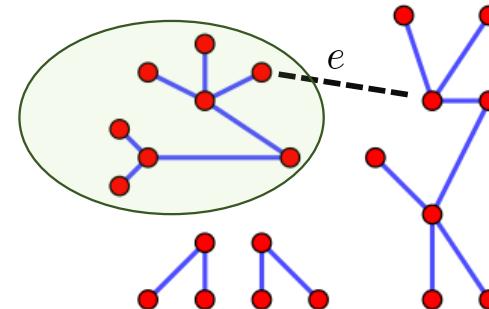
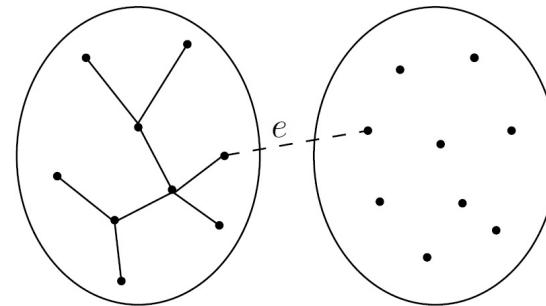
Repeatedly add to the tree the lightest edge connecting a tree vertex to a non-tree vertex until all vertices belong to the tree.

## Kruskal's Algorithm.

At any time, we have a *forest*  
i.e. a collection of trees.

Initially, each vertex is a tree by itself.

We repeatedly add the lightest edge  
that joins two of the trees until we  
have a single tree.



***The correctness of both algorithms follows from the cut property.***

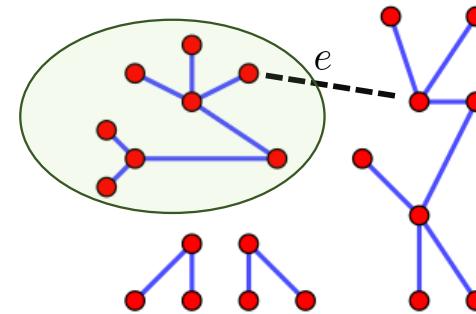
# KRUSKAL'S ALGORITHM

## Kruskal's Algorithm.

At any time, we have a *forest*  
i.e. a collection of trees.

Initially, each vertex is a tree by itself.

We repeatedly add the lightest edge  
that joins two of the trees until we  
have a single tree.



*Implementation?*

## Pseudocode.

Sort the edges in non-decreasing order of their weights.

Let  $e_1, \dots, e_m$  be the sorted order of edges.

for  $i = 1$  to  $m$ :

if the end-points of  $e_i$  are in different trees: *How do we implement this?*

    add  $e_i$  to the MST

## KRUSKAL'S ALGORITHM

This ***disjoint set*** data structure allows us to store a collection of sets and supports the following operations.

|                      |   |
|----------------------|---|
| $\text{makeset}(x)$  | <i>create a new set <math>\{x\}</math></i>  |
| $\text{union}(x, y)$ | <i>combine the set containing <math>x</math> and the set containing <math>y</math> into one set</i> |
| $\text{find}(x)$     | <i>return a unique identifier for the set containing <math>x</math></i>                             |

How do we use this data structure for Kruskal's algorithm?

We think of each tree in the forest as a set (of vertices) in our collection.

Initially, we create a singleton for each vertex using *makeset*.

We process the edges in non-decreasing the order of their weights.

We add an edge  $(u, v)$  to the MST iff  $\text{find}(u) \neq \text{find}(v)$ .

If an edge  $(u, v)$  is added, we execute *union* $(u, v)$ .

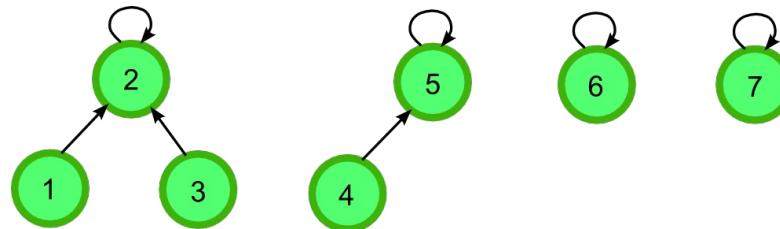
# DISJOINT SET DATA STRUCTURE

How do we implement the data structure?

We represent each set in our collection using a tree.

*This tree is internal to the data structure and has nothing to do with the MST we're building.*

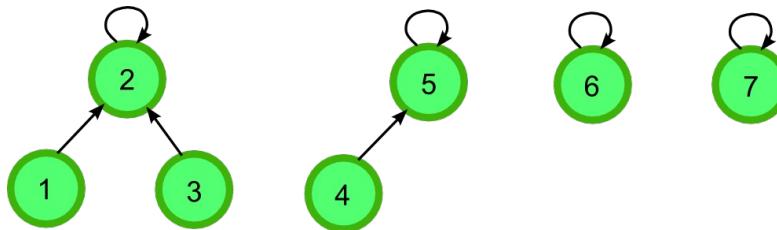
The unique identifier for a set is the root of the tree representing it.



Every node has a pointer to its parent (**not** to the children).

If the node is the root of a tree, its parent pointer points to itself.

# DISJOINT SET DATA STRUCTURE



Each node also has a *rank*.

*makeset(x)* : Creates a single node containing  $x$ .

Its parent pointer points to itself.

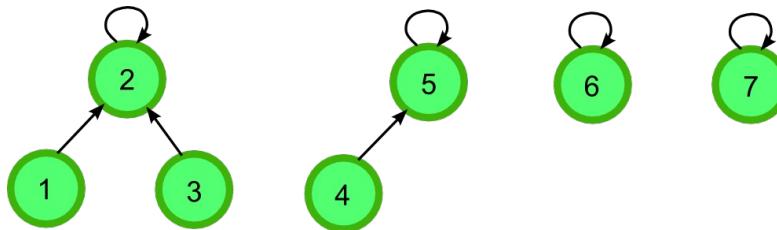


Its rank is 0.

*find(x)* : Starting at the node containing  $x$ , we follow the parent pointer until we reach a root node.

We return a pointer to this node.

# DISJOINT SET DATA STRUCTURE



Each node also has a *rank*.

$\text{union}(x, y) :$  Let  $r_1 = \text{find}(x)$  and  $r_2 = \text{find}(y)$ .

If  $r_1 \neq r_2$ :

If  $\text{rank}(r_1) \leq \text{rank}(r_2)$ :

Set  $\text{parent}(r_1) = r_2$ .

Else:

Set  $\text{parent}(r_2) = r_1$ .

If the ranks are equal, we increment the rank of the parent after merging.



If  $\text{rank}(r_1) == \text{rank}(r_2)$ : increment  $\text{rank}(r_2)$

This procedure ensures that the rank of a node is strictly smaller than its parent.

# DISJOINT SET DATA STRUCTURE

Each node also has a **rank**.

$union(x, y)$  : Let  $r_1 = find(x)$  and  $r_2 = find(y)$ .

If  $r_1 \neq r_2$ :

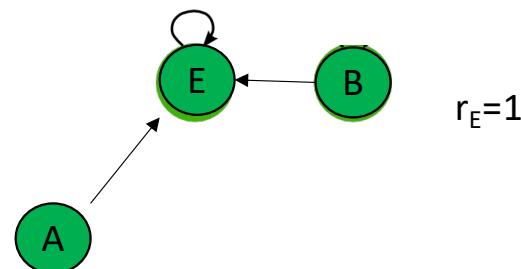
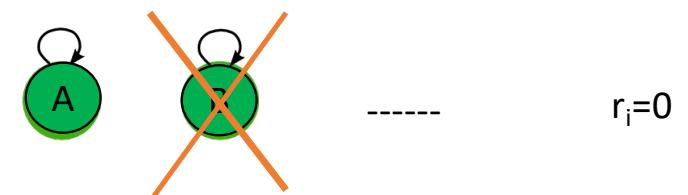
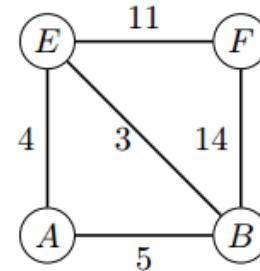
If  $rank(r_1) \leq rank(r_2)$ :

Set  $parent(r_1) = r_2$ .

Else:

Set  $parent(r_2) = r_1$ .

If  $rank(r_1) == rank(r_2)$ : increment  $rank(r_2)$



Complexity?

And so on ....

## DISJOINT SET DATA STRUCTURE

**Claim.** If the root of tree has rank  $r$ , the tree contains  $\geq 2^r$  nodes.

*Proof.* By induction. Base case:  $r = 0$ . The tree has 1 node. ✓

Inductive Hypothesis: the statement is true for  $r = k$ .

We need to show that the statement is true for  $r = k + 1$ .

When does the rank of a node increase from  $k$  to  $k + 1$ ?

When we combine two trees whose roots have rank are  $k$ .

By our inductive hypothesis, both trees have  $\geq 2^k$  nodes.

Thus, after merging, the tree has  $2^{k+1}$  nodes. ■

If  $n$  is the total number of elements in the data structure,  
what is the maximum rank any node can have?  $\lfloor \log_2 n \rfloor$

# DISJOINT SET DATA STRUCTURE

How much time does  $\text{makeset}(\cdot)$  take?  $O(1)$

|                      |  |
|----------------------|--|
| $\text{makeset}(x)$  | create a new set $\{x\}$   |
| $\text{union}(x, y)$ | combine the set containing $x$ and the set containing $y$ into one set |
| $\text{find}(x)$     | return a unique identifier for the set containing $x$                  |

How much time does  $\text{find}(\cdot)$  take?

$O(\log n)$  since rank increases monotonically on the path from a node to the root of its tree, and the maximum rank is  $\leq \log_2 n$ .

How much time does  $\text{union}(\cdot, \cdot)$  take?

$O(\log n)$  since, apart from two  $\text{find}$  operations, it takes  $O(1)$  time.

What is the running time of Kruskal's algorithm using this data structure?

$O(m \log m)$  time for sorting the edges by weight.

How many of each data structure operation do we need?

$n$  makesets,  $2m$  finds,  $n - 1$  unions :  $O(m \log m + n + 2m \log n + n-1) = O(n + m \log m)$

Note:  $\log m \leq 2 \log n$  since  $m \leq n^2$ . So final complexity is:  $O(n + m \log n)$  (same as Prim)

```

Let  $r_1 = \text{find}(x)$  and  $r_2 = \text{find}(y)$ .
If  $r_1 \neq r_2$ :
    If  $\text{rank}(r_1) \leq \text{rank}(r_2)$ :
        Set  $\text{parent}(r_1) = r_2$ .
    Else:
        Set  $\text{parent}(r_2) = r_1$ .
    If  $\text{rank}(r_1) == \text{rank}(r_2)$ : increment  $\text{rank}(r_2)$ 

```

```

Sort the edges in non-decreasing order of their weights.
Let  $e_1, \dots, e_m$  be the sorted order of edges.
for i = 1 to m:
    if the end-points of  $e_i$  are in different trees:
        add  $e_i$  to the MST

```

## Exercise

## MINIMUM SPANNING TREE (MST)

**Fact.** If all the edge weights are distinct, then there is a unique MST.

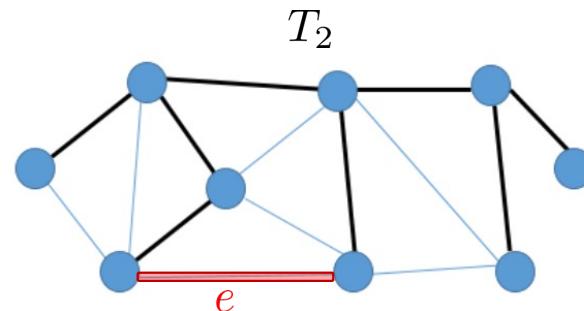
*Proof.* Suppose that each edge of a graph  $G$  has distinct weight but there are two distinct MST's  $T_1$  and  $T_2$ .

Consider the *lightest* edge  $e$  that is present in one MST but not the other.

W.L.O.G. assume that  $e$  is in  $T_1$  but not  $T_2$ .

If we add  $e$  to  $T_2$ , we will get a cycle  $C$ .

Can  $e$  be the *heaviest* edge in this cycle?



No, since that would mean that all edges in the cycle were present in  $T_1$ .

Then, replacing the heaviest edge in the cycle by  $e$  yields a lighter MST!



## MINIMUM SPANNING TREE (MST)

There are faster algorithms for Computing the MST.

There is a randomized algorithm that takes  $O(m + n)$  time *in expectation*.

[https://en.wikipedia.org/wiki/Expected\\_linear\\_time\\_MST\\_algorithm](https://en.wikipedia.org/wiki/Expected_linear_time_MST_algorithm)

There is a deterministic algorithm that takes  $O(m \alpha(m, n))$  time.

<https://www.cs.princeton.edu/courses/archive/fall05/cos528/handouts/A%20Minimum%20Spanning.pdf>

$\alpha$  is the inverse Ackermann function

## HUFFMAN ENCODING

Each quantized value is encoded in binary.

However, in general some quantized values are more frequent than other.

Consider a simple example with four quantized values represented by the symbols  $A, B, C$  and  $D$  whose frequencies are:

$A$  : 75 million,  $B$  : 2 million,  $C$  : 22 million and  $D$  : 31 million.

One option would be to encode each value using two bits.

How many bits does this require? *260 million*

Can we do better by encoding the values with codes of different lengths?

Yes! Here is a better encoding.  $A$  : 0,  $B$  : 100,  $C$  : 101 and  $D$  : 11.

How many bits does this require? *209 million*

## HUFFMAN ENCODING

We have four quantized values  $A, B, C$  and  $D$  with frequencies:

$A : 75$  million,  $B : 2$  million,  $C : 22$  million and  $D : 31$  million.

The encoding  $A : 0$ ,  $B : 100$ ,  $C : 101$  and  $D : 11$  requires 209 million bits.

Can we do even better? How about the following encoding?

$A : 0$ ,  $B : 01$ ,  $C : 11$  and  $D : 1$ .

This encoding is not good because it is not uniquely decipherable.

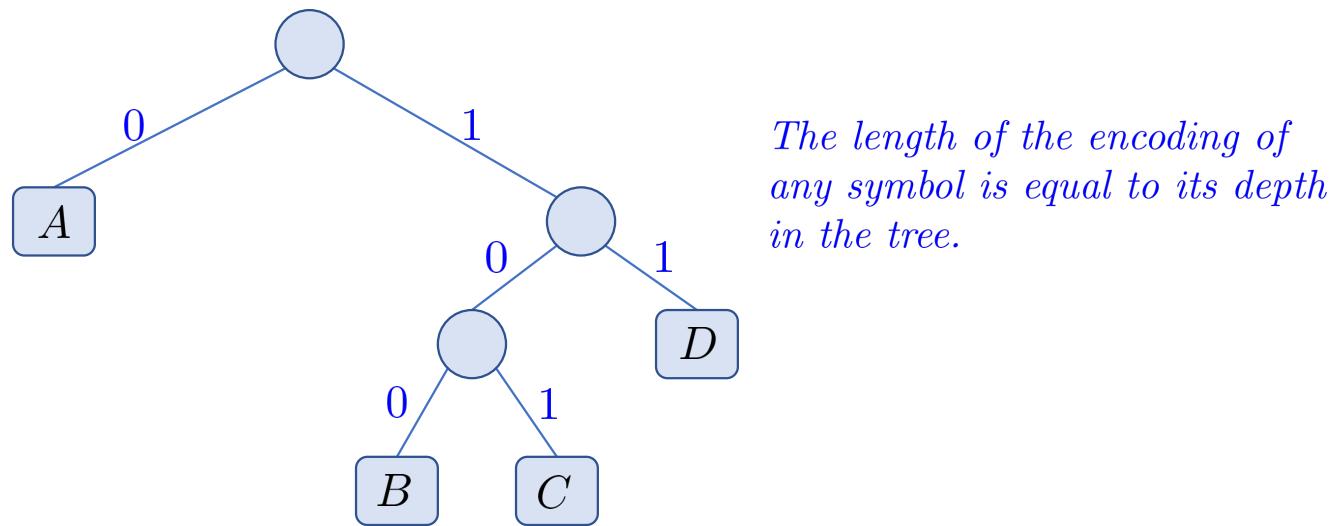
For example 11 can either be  $C$  or  $DD$ .

We will use a *prefix-free code* to avoid this. *No code is a prefix of another.*

## HUFFMAN ENCODING

The encoding  $A : 0$ ,  $B : 100$ ,  $C : 101$  and  $D : 11$  is prefix free.

Any encoding that is prefix free can be represented as a binary tree whose leaves are the symbols.



Let  $f(A), f(B), f(C)$  and  $f(D)$  denote the frequencies of  $A, B, C$  and  $D$ .

Then the length of the encoding is:  $1 \cdot f(A) + 3 \cdot f(B) + 3 \cdot f(C) + 2 \cdot f(D)$ .

## HUFFMAN ENCODING

In general, there are  $n$  symbols whose frequencies are given in an array  $f[1..n]$ .

We want a tree representing a prefix free codes that minimizes

$$\sum_{i=1}^n \text{depth}(i) \cdot f(i).$$