

# Algorithms and Data Structures: Minimum Spanning Trees (Kruskal)

7th Nov, 2014

# Minimum Spanning Tree Problem

Given: *Undirected connected weighted graph*  $(\mathcal{G}, W)$

Output: *An MST of*  $\mathcal{G}$

- ▶ We have already seen the PRIM algorithm, which runs in  $O((m + n) \lg(n))$  time (standard Heap implementation) for graphs with  $n$  vertices and  $m$  edges.
- ▶ In this lecture we will see KRUSKAL's algorithm, a different approach to constructing a MST.

# Kruskal's Algorithm

A **forest** is a graph whose connected components are trees.

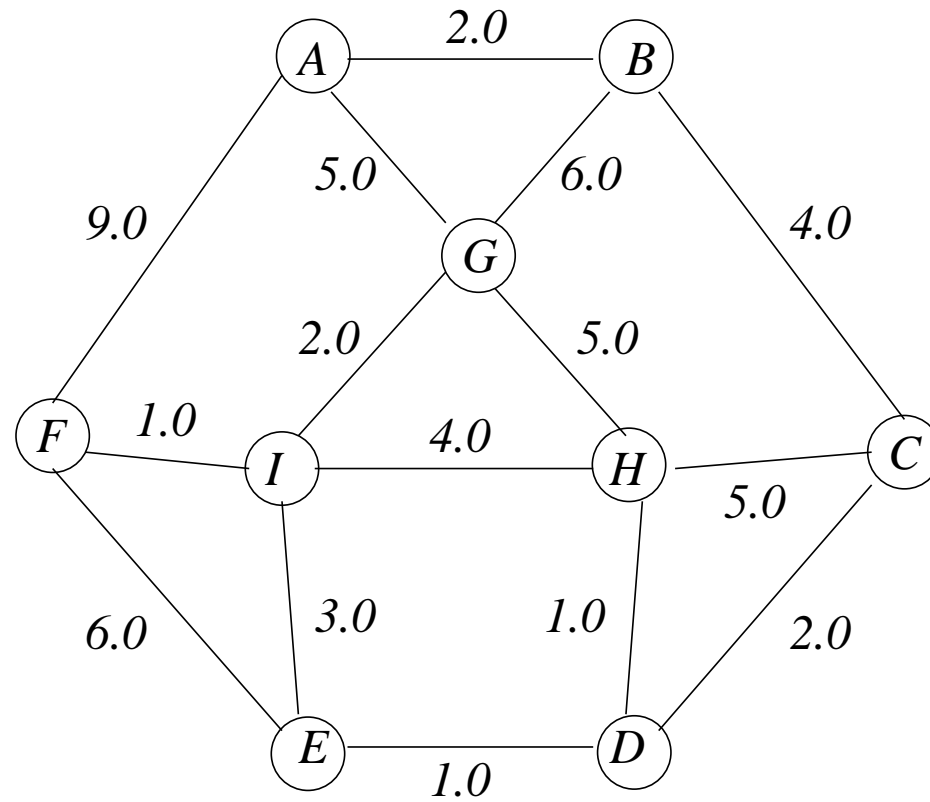
## Idea

Starting from a spanning forest with no edges, repeatedly add edges of minimum weight (never creating a cycle) until the forest becomes a tree.

**Algorithm** KRUSKAL( $\mathcal{G}, W$ )

1.  $F \leftarrow \emptyset$
2. **for all**  $e \in E$  in the order of increasing weight **do**
3.     **if** the endpoints of  $e$  are in different connected components of  $(V, F)$  **then**
4.          $F \leftarrow F \cup \{e\}$
5. **return** tree with edge set  $F$

# Example



## Correctness of Kruskal's algorithm

1. Throughout the execution of  $\text{KRUSKAL}$ ,  $(V, F)$  remains a spanning forest.

*Proof:*  $(V, F)$  is a spanning subgraph because the vertex set is  $V$ . It always remains a forest because edges with endpoints in different connected components never induce a cycle.

2. Eventually,  $(V, F)$  will be connected and thus a spanning tree.

*Proof:* Suppose that after the complete execution of the loop,  $(V, F)$  has a connected component  $(V_1, F_1)$  with  $V_1 \neq V$ . Since  $\mathcal{G}$  is connected, there is an edge  $e \in E$  with exactly one endpoint in  $V_1$ . This edge would have been added to  $F$  when being processed in the loop, so this can never happen.

3. Throughout the execution of  $\text{KRUSKAL}$ ,  $(V, F)$  is contained in some MST of  $\mathcal{G}$ .

*Proof:* Similar to the proof of the corresponding statement for Prim's algorithm. **Will prove in week 9 Tutorial.**

# Data Structures for Disjoint Sets

- ▶ A **disjoint set** data structure maintains a collection  $\mathcal{S} = \{S_1, \dots, S_k\}$  of **disjoint sets**.
- ▶ The sets are *dynamic*, i.e., they may change over time.
- ▶ Each set  $S_i$  is identified by some *representative*, which is some member of that set.

## Operations:

- ▶ **MAKE-SET( $x$ )**: Creates new set whose only member is  $x$ . The representative is  $x$ .
- ▶ **UNION( $x, y$ )**: Unites set  $S_x$  containing  $x$  and set  $S_y$  containing  $y$  into a new set  $S$  and removes  $S_x$  and  $S_y$  from the collection.
- ▶ **FIND-SET( $x$ )**: Returns representative of the set holding  $x$ .

# Implementation of Kruskal's Algorithm

**Algorithm** KRUSKAL( $\mathcal{G}, W$ )

1.  $F \leftarrow 0$
2. **for all** vertices  $v$  of  $\mathcal{G}$  **do**
3.     MAKE-SET( $v$ )
4. sort edges of  $\mathcal{G}$  into non-decreasing order by weight
5. **for all** edges  $(u, v)$  of  $\mathcal{G}$  in non-decreasing order by weight **do**
6.     **if** FIND-SET( $u$ )  $\neq$  FIND-SET( $v$ ) **then**
7.          $F \leftarrow F \cup \{(u, v)\}$
8.         UNION( $u, v$ )
9. **return**  $F$

# Analysis of KRUSKAL

Let  $n$  be the number of vertices and  $m$  the number of edges of the input graph

- ▶ Line 1:  $\Theta(1)$
- ▶ Loop in Lines 2–3:  $\Theta(n \cdot T_{\text{MAKE-SET}}(n))$
- ▶ Line 4:  $\Theta(m \lg m)$
- ▶ Loop in Lines 5–8:  $\Theta\left(m \cdot T_{\text{FIND-SET}}(n) + n \cdot T_{\text{UNION}}(n)\right)$ .
- ▶ Line 9:  $\Theta(1)$

Overall:

$$\Theta\left(n\left(T_{\text{MAKE-SET}}(n) + T_{\text{UNION}}(n)\right) + m\left(\lg m + T_{\text{FIND-SET}}(n)\right)\right)$$



## Analysis of KRUSKAL (overview)

$$T(n, m) = \Theta\left(n(T_{\text{MAKE-SET}}(n) + T_{\text{UNION}}(n)) + m(\lg m + T_{\text{FIND-SET}}(n))\right)$$

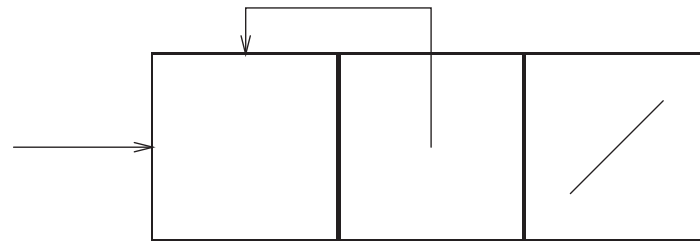
With standard efficient implementations of disjoint sets this amounts to

$$T(n, m) = \Theta(m \lg(m)).$$

- ▶ *NOT* better than the standard Heap implementation of PRIM for typical implementations of disjoint sets.
- ▶ Always have to sort the weights when using KRUSKAL:
  - ▶  $\Theta(m \lg(m))$  if the weights are arbitrarily large.

# Linked List Implementation of Disjoint Sets

Each element represented by a pointer to a cell:



Use a linked list for each set.

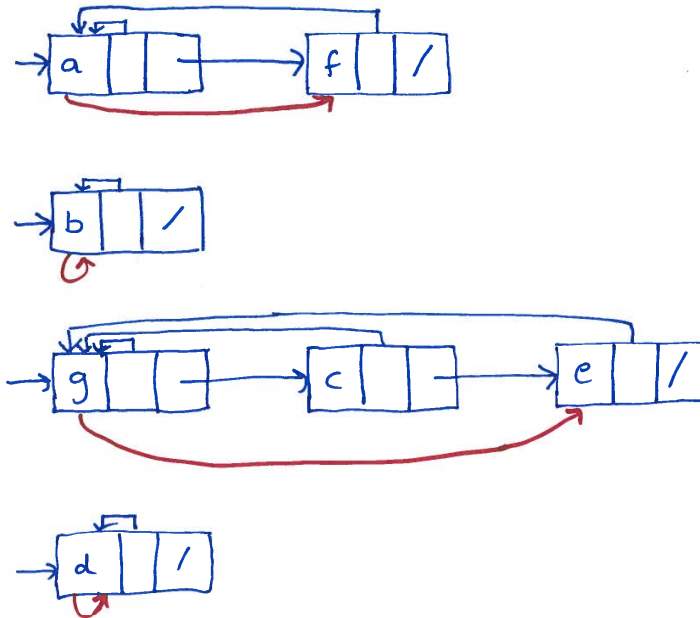
Representative of the set is at the head of the list.

Each cell has a pointer direct to the representative (head of the list).

# Example

Linked list representation of

$\{a, f\}$ ,  $\{b\}$ ,  $\{g, c, e\}$ ,  $\{d\}$ :



The "representatives" are a, b, g and d respectively

last[ ] pointers are in red

# Analysis of Linked List Implementation

MAKE-SET: constant ( $\Theta(1)$ ) time.

FIND-SET: constant ( $\Theta(1)$ ) time.

UNION: Naive implementation of

UNION( $x, y$ )

appends  $x$ 's list onto end of  $y$ 's list.

**Assumption:** Representative  $y$  of each set has attribute  $\text{last}[y]$ : a pointer to last cell of  $y$ 's list.

**Snag:** have to update “representative pointer” in each cell of  $x$ 's list to point to the representative (head) of  $y$ 's list.

Cost is:

$\Theta(\text{length of } x\text{'s list})$ .

# Notation for Analysis

Express running time in terms of:

$\hat{n}$ : the number of MAKE-SET operations,

$\hat{m}$ : the number of MAKE-SET, UNION and FIND-SET operations overall.

## Note

1. After  $\hat{n} - 1$  UNION operations only one set remains.
2.  $\hat{m} \geq \hat{n}$ .

# Weighted-Union Heuristic

## Idea

Maintain a “length” field for each list. To execute

UNION( $x, y$ )

append shorter list to longer one (breaking ties arbitrarily).

## Theorem 1

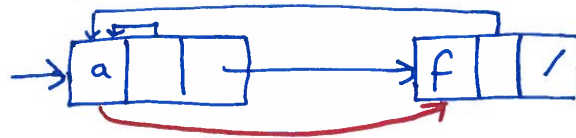
*Using the linked-list representation of disjoint sets and the weighted-union heuristic, a sequence of  $\hat{m}$  MAKE-SET, UNION & FIND-SET operations,  $\hat{n}$  of which are MAKE-SET operations, takes*

$$O(\hat{m} + \hat{n} \lg \hat{n})$$

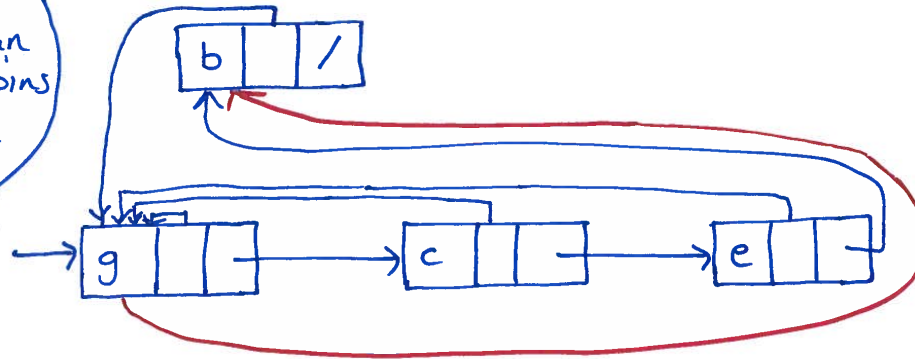
*time.*

**“Proof”:** Each element appears at most  $\lg \hat{n}$  times in the short list of a UNION.

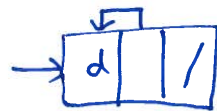
# Example (UNION(g, b))



b's list is SHORTER than g's, so it joins the end of g's list



- b's "representative" pointer changes to point at g-cell
- e's "next" pointer changes to point at b-cell
- g's "last" pointer changes to point at b-cell



result of performing UNION(g, b)

## KRUSKAL with Linked lists (weighted union)

The run-time for KRUSKAL (for  $\mathcal{G} = (V, E)$  with  $|V| = n, |E| = m$ ) is

$$T(n, m) = \Theta\left(n\left(T_{\text{MAKE-SET}}(n) + T_{\text{UNION}}(n)\right) + m\left(\lg m + T_{\text{FIND-SET}}(n)\right)\right)$$

In terms of the collection of “Disjoint-sets” operations, we have  $\hat{m} = 2n + 2m$  operations,  $\hat{n} = n$  which are UNION. So

$$\begin{aligned} T(n, m) &= \Theta(m \lg(m) + (2n + 2m) + n \lg(n)) \\ &= \Theta(m \lg(m)) \end{aligned}$$