

# Kruskal's alg. for MST and disjoint sets

---

Assume we have a data structure of disjoint sets

$$S = \{S_i\} \text{ such that } S_i \cap S_j = \emptyset$$

and supports the following operations.

- Make-Set ( $x$ ) :  $S \leftarrow S \cup \{\{x\}\}$

- Union ( $x, y$ ) :  $S \leftarrow S - \{S_1, S_2\} \cup \{S_1 \cup S_2\}$   
where  $x \in S_1$  and  $y \in S_2$ .

- Find-Set ( $x$ ) : returns a unique object that represents  $S$   
where  $x \in S$ .

Kruskal's alg. for MST.

$T \leftarrow \emptyset$

for each  $v \in V$

do Make-Set( $v$ )

sort  $E$  by increasing edge weight

for each  $(u, v) \in E$  (in sorted order)

do if Find-Set( $u$ )  $\neq$  Find-Set( $v$ )

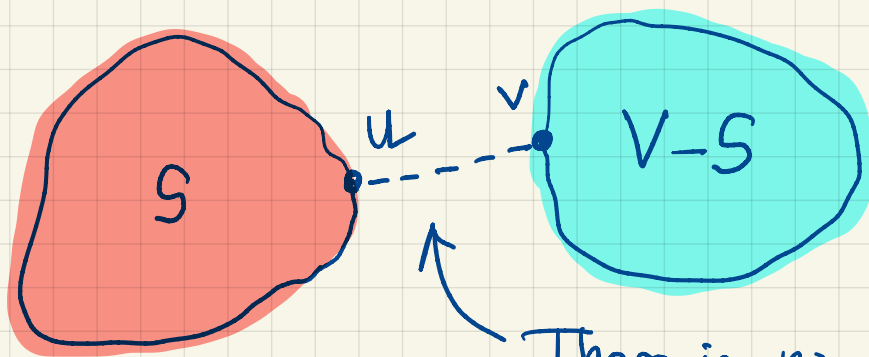
$\triangleright (u, v)$  make no cycle

then  $T \leftarrow T \cup \{(u, v)\}$

Union( $u, v$ )

$\triangleright u$  &  $v$  now in same component

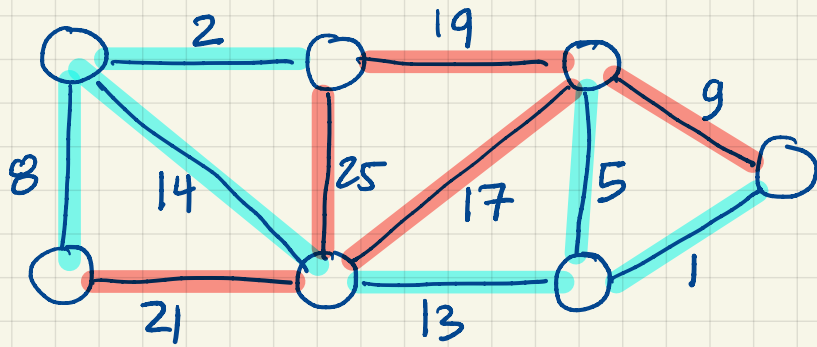
Why does it work? (previous proof)



$$S = \{x \mid \text{Find-Set}(u) = \text{Find-Set}(x)\}$$

There is no smaller (due to sorted order)  
weight edge that crosses the cut

Example:



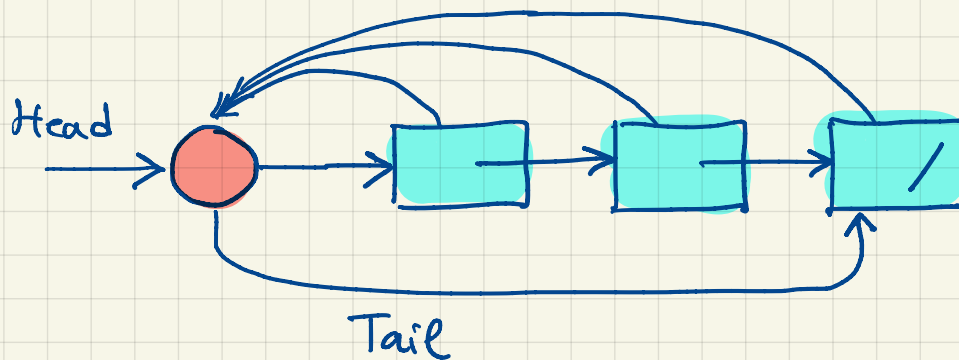
Time:

- sorting:  $\Theta(E \log E) = \Theta(E \log V)$  since graph is connected  
 $= O(E \log V)$  in general since  $E = O(V^2)$   
( $E = \Omega(V)$  if graph is connected)
- $\Theta(V)$  Make-Sets
- $\Theta(E)$  Find-Sets
- $O(V)$  Unions (exactly  $|V|-1$ )

We can implement  $m$  set operations on  $n$  elements in  $\tilde{O}(m \cdot \alpha(m, n))$  time, where  $\alpha$  grows very slowly.  $\alpha(m, n) \leq 4$  if  $m, n = 10^{80}$ .

# Disjoint-Set Implementations

- Each set is a linked list of its elements. Each element points back to "Head" of list, which is the "representative" of the list.



- Using this implementation,  $\text{Make-Set}(x)$  and  $\text{Find-Set}(x)$  can be done in  $O(1)$  time.
- $\text{Union}(x, y)$  takes more time: "Copy" set of  $x$  into that of  $y$ , and update pointers. Using  $y$ 's tail we can quickly identify where to append. This takes time proportional to size of  $x$ 's set.

Worst Case Scenario:

	<u># pointer updates</u>
Union $(x_1, x_2)$ :	1
Union $(x_2, x_3)$ :	2
Union $(x_3, x_4)$ :	3
:	
Union $(x_{n-1}, x_n)$ :	$\frac{n-1}{\theta(n^2)}$ time.

- Improvement: Append smaller set to larger one  
(store size in set header)

A single Union operation can still take  $\theta(n)$  time; e.g. if both sets have  $\frac{n}{2}$  elements.

- But  $m$  set operations, in which  $n$  are Make-Sets, will take  $O(m + n \lg n)$  time. When an elem's back pointer is updated, its set at least doubles in size. This can happen at most  $\lg n$  times per element.

Amortized analysis: Average time, but no probability

$$n \leq m$$

$m$  operations take  $O(m + n \lg n) = O(m \lg n)$

So  $O(\lg n)$  per operation.

We say each operation takes  $O(\lg n)$  amortized time

Note: Average, but no bad sequences. Some operations will take more than  $O(\lg n)$ , but every sequence of  $m$  operations takes at most  $O(m \lg n)$  time.

## Amortized Analysis

No probability

Obtain average performance in the worst-case

No bad sequences

## Average-Case Analysis

May involve prob.

Average performance

Possible bad sequences

## Techniques for Amortized Analysis

- Aggregate analysis
- Accounting method
- Potential method.

What we did in disjoint sets was aggregate analysis.

Example: Stack with an added operation.

Push( $S, x$ ) : Push  $x$  onto stack

Pop( $S$ ) : Pop top of stack and return popped object

Multipop( $S, k$ ) : Remove  $\min(|S|, k)$  objects from top.

Multipop( $S, k$ )

while not stack-empty( $S$ ) and  $k \neq 0$

do Pop( $S$ )

$k \leftarrow k - 1$

Running time of Multipop is  $O(\min(|S|, k))$

which means in the worst-case it's  $O(n)$

Therefore, a sequence of  $n$  operations takes  $O(n^2)$  time.



## Aggregate Analysis

Get a better bound by considering the entire sequence of  $n$  operations.

Claim: Any sequence of  $n$  stack operations take  $O(n)$  time.

- Any object can be popped at most once after it's pushed.

The # times a pop is called on a non-empty stack

(including those in a multipop) is at most equal to # pushes.

- Given  $n$  operations that result in  $m$  pops from within multipop, the running time is  $O(n+m)$

But  $m \leq n$ , so  $O(2n) = O(n)$ .

Each operation runs in  $O(1)$  amortized time

Aggregate: "Treat every operation the same way in terms of time"