#### LEC 18

## CSE 373

# Disjoint Sets I

#### **BEFORE WE START**

Which of the following is true about minimum spanning trees?

- a) Every graph has at least one MST
- b) Every graph has at most one MST
- c) The shortest path from any vertex u to any vertex v is the one that follows the MST
- d) Every MST has a source vertex
- e) Dijkstra's algorithm computes an MST
- f) None of the above

#### pollev.com/uwcse373

Instructor Aaron Johnston

- TAs Timothy Akintilo Brian Chan Joyce Elauria Eric Fan Farrell Fileas
- Melissa Hovik Leona Kazi Keanu Vestil Siddharth Vaidyanathan Howard Xiao

### Announcements

- P3 due TONIGHT at 11:59pm PDT
  - Don't forget to fill out the <u>P3 Feedback Survey</u>!
- EX3 due Friday, 8/07 11:59pm PDT
- P4, the last project of the quarter, released tonight!
  - If you haven't already, please fill out the <u>P4 Partner Form</u> so we can send out partner assignments

## Learning Objectives

After this lecture, you should be able to...

- 1. Describe Kruskal's Algorithm, evaluate why it works, and describe why it needs a new ADT
- 2. Compare and contrast QuickUnion with QuickFind and describe how the two structure optimize for different operations
- 3. Implement WeightedQuickUnion and describe why making the change protects against the worst case find runtime
- 4. Describe path compression at a high level and argue for why it improves runtimes despite not following an invariant

### **Review** Minimum Spanning Trees (MSTs)

- A Minimum Spanning Tree for a graph is a set of that graph's edges that connect all of that graph's vertices (**spanning**) while minimizing the total weight of the set (**minimum**)
  - Note: does NOT necessarily minimize the path from each vertex to every other vertex
  - Any tree with V vertices will have V-1 edges
  - A separate entity from the graph itself! More of an "annotation" applied to the graph, just like a Shortest Paths Tree (SPT)



**Minimum Spanning Tree** 

## **Review** Why do MST Algorithms Work?

- Two useful properties for MST edges. We can think about them from either perspective:
  - Cycle Property: The heaviest edge along a cycle is NEVER part of an MST.
  - **Cut Property:** Split the vertices of the graph into any two sets A and B. The lightest edge between A and B is ALWAYS part of an MST. *(Prim's thinks this way)*
- Whenever you add an edge to a tree you create exactly one cycle. Removing any edge from that cycle gives another tree!
- This observation, combined with the cycle and cut properties form the basis of all of the greedy algorithms for MSTs.
  - greedy algorithm: chooses best known option at each point and commits, rather than waiting for a global view of the graph before deciding

## **Review** Adapting Dijkstra's: Prim's Algorithm

- Normally, Dijkstra's checks for a shorter path from the start.
- But MSTs don't care about individual paths, only the overall weight!
- New condition: "would this be a smaller edge to connect the current known set to the rest of the graph?"



```
prims (G graph, V start)
Map edgeTo, distTo;
initialize distTo with all nodes mapped to \infty, except start to 0
PriorityQueue<V> perimeter; perimeter.add(start);
while (!perimeter.isEmpty()):
  u = perimeter.removeMin()
  known.add(u)
  for each edge (u,v) to unknown v with weight w:
    oldDist = distTo.get(v) // previous smallest edge to v
    newDist = distTo.get(u) + w // is this a smaller edge to v?
    if (newDist < oldDist):</pre>
      distTo.put(u, newDist)
      edgeTo.put(u, v)
      if (perimeter.contains(v)):
        perimeter.changePriority(v, newDist)
      else:
        perimeter.add(v, newDist)
```

## A Different Approach

- Suppose the MST on the right was produced by Prim's
- **Observation**: We basically chose all the smallest edges in the entire graph (1, 2, 3, 4, 6)
  - The only exception was 5. Why shouldn't we add edge 5?
  - Because adding 5 would create a cycle, and to connect A, C, & D we'd rather choose 1 & 4 than 1 & 5 or 4 & 5.
- Prim's thinks "vertex by vertex", but what if you think "edge by edge" instead?
  - Start with the smallest edge in the entire graph and work your way up
  - Add the edge to the MST as long as it connects two new groups (meaning don't add any edges that would create a cycle)



Building an MST "edge by edge" in this graph:

- Add edge 1
- Add edge 2
- Add edge 3
- Add edge 4
- Skip edge 5 (would create a cycle)
- Add edge 6
- Finished: all vertices in the MST!

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
```

```
vMST = msts.find(v)
if (uMST != vMST):
  finalMST.add(edge (u, v))
  msts.union(uMST, vMST)
```

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
for each edge (u,v) in ascending order:
```

```
uMST = msts.find(u)
vMST = msts.find(v)
if (uMST != vMST):
  finalMST.add(edge (u, v))
  msts.union(uMST, vMST)
```

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
for each edge (u, v) in according order;
```

```
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
       finalMST.add(edge (u, v))
       msts.union(uMST, vMST)
```

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
```

```
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
       finalMST.add(edge (u, v))
       msts.union(uMST, vMST)
```

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
```

```
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
       finalMST.add(edge (u, v))
       msts.union(uMST, vMST)
```

• This "edge by edge" approach is how Kruskal's Algorithm works!

- Key Intuition: Kruskal's keeps track of isolated "islands" of vertices (each is a sub-MST)
  - Start with each vertex as its own "island"
  - If an edge connects two vertices within the same "island", it forms a cycle! Discard it.
  - If an edge connects two vertices in different "islands", add it to the MST! Now those "islands" need to be combined.



"islands"

```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
```

```
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
       finalMST.add(edge (u, v))
       msts.union(uMST, vMST)
```

## **Prim's Demos and Visualizations**

- Dijkstra's Visualization
  - <u>https://www.youtube.com/watch?v=1oiQ0hrVwJk</u>
  - Dijkstra's proceeds radially from its source, because it chooses nearby edges by *path length from source*
- Prim's Visualization
  - <u>https://www.youtube.com/watch?v=6uq0cQZOyoY</u>
  - Prim's jumps around the perimeter, because it chooses nearby edges by edge weight (there's no source)
- Kruskal's Visualization
  - <u>https://www.youtube.com/watch?v=ggLyKfBTABo</u>
  - Kruskal's jumps around the entire graph, because it chooses from all edges purely by edge weight (while preventing cycles)

## Selecting an ADT

- Kruskal's needs to find what MST a vertex belongs to, and union those MSTs together
  - Our existing ADTs don't lend themselves well to "unioning" two sets...
  - Let's define a new one!

```
kruskalMST(G graph)
Set(?) msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
        finalMST.add(edge (u, v))
        msts.union(uMST, vMST)
```



## Disjoint Sets ADT (aka "Union-Find")

- Kruskal's will use a Disjoint Sets ADT under the hood
  - Conceptually, a single instance of this ADT contains a "family" of sets that are disjoint (no element belongs to multiple sets)

```
kruskalMST(G graph)
DisjointSets<V> msts; Set finalMST;
initialize msts with each vertex as single-element MST
sort all edges by weight (smallest to largest)
for each edge (u,v) in ascending order:
    uMST = msts.find(u)
    vMST = msts.find(v)
    if (uMST != vMST):
        finalMST.add(edge (u, v))
        msts.union(uMST, vMST);
```



#### DISJOINT SETS ADT

#### State

Family of Sets

- disjoint: no shared elements
- each set has a representative (either a member or a unique ID)

#### Behavior

makeSet(value) - new set with value as only member (and representative) find(value) - return representative of the set containing value union(x, y) - combine sets containing x and y into one set with all elements, choose single new representative

## Project 4: Mazes!

- You find yourself trapped in the Labyrinth of Greek legend – bummer!
- How do you solve a maze?
  - If we could model a maze as a graph, we'd just need an algorithm to find a path from s to t... Maybe even the shortest path?

• How do you *generate* a maze?

 We'd love an algorithm that is guaranteed to connect s to t (spanning), but only produces one path from s to t (tree)...



## **Project 4: Mazes**

- It turns out that *randomizing* the weights of a graph and then computing the MST is a **fantastic way to generate** mazes!
- In P4, you'll do both: Implement Dijkstra's to solve an arbitrary maze, then implement Kruskal's (and a Disjoint Set) to generate those mazes
- This project is really application-heavy!
  - Graphical User Interface (GUI) for viewing mazes and solving them
  - Significantly more starter code than past projects, to give you practice integrating with an existing codebase
  - A major part of the challenge in P4 is reading through the starter code to understand what you need to interface with! Don't underestimate the time that takes.
- 2 week project, and 2 weeks worth of work. It's never been more important to start early!



### **Lecture Outline**



## **Case Study: Disjoint Sets**

- Today's lecture on the data structures which implement the Disjoint Sets ADT is an interesting case study in data structure design and iterative design improvements
  - Good chance to dust off your metacognitive skills!
  - In particular, try to identify what observations we make in each data structure that inspire improvements in the next data structure. How could you apply a similar skill to your own data structures?



Metacognitive Opportunities Ahead!

## Can we use an existing data structure?

aka "Can we just throw maps at the problem?"

### Maps to Sets (baseline): map from representative ID to set of values



find(value): scan
through every set unde
every representative

union(x, y): copy all elements from set pointed to by x into set pointed to by y

DISJOINT SETS ADT
<pre>State Family of Sets • disjoint: no shared elements • each set has a representative (either     a member or a unique ID)</pre>
<pre>Behavior makeSet(value) - new set with value as only member (and representative) find(value) - return representative of the set containing value union(x, y) - combine sets containing x and y into one set with all elements, choose single new representative</pre>

	Maps to Sets
<pre>makeSet(value)</pre>	Θ(1)
<pre>find(value)</pre>	$\Theta(n)$
union(x, y)	$\Theta(n)$

### **QuickFind Implementation**

**QuickFind**: map from value to representative ID



<pre>find(Farrell)</pre>	$\rightarrow$	1
find(Eric) $\rightarrow$	2	
<pre>find(Farrell)</pre>	! =	find(Eric)
<pre>find(Farrell)</pre>	==	<pre>find(Keanu)</pre>

- If we store values as the keys, we can take advantage of fast lookup to make find fast!
- But what about union?

**find(value):** lookup element and return corresponding rep.

Maps to SetsQuickFindmakeSet(value) $\Theta(1)$  $\Theta(1)$ find(value) $\Theta(n)$  $\Theta(1)$ union(x, y) $\Theta(n)$  $\Theta(n)$ 

union(x, y): scan through all elements to find those in same set, update to new rep.

#### DISJOINT SETS ADT

#### State

Family of Sets

- disjoint: no shared elements
- each set has a representative (either a member or a unique ID)

#### Behavior

makeSet(value) - new set with value as only member (and representative) find(value) - return representative of the set containing value union(x, y) - combine sets containing x and y into one set with all elements, choose single new representative

### **Lecture Outline**



## **QuickUnion Data Structure**

- Fundamental idea:
  - QuickFind tracks each element's ID
  - QuickUnion tracks each element's *parent*. Only the root has an ID!
    - Each set becomes tree-like, but something slightly different called an **up-tree**: store pointers from children to parents!



## QuickUnion: Find

#### find(Eric):

jump to Eric node travel upward until root return ID

- Key idea: can travel upward from any node to find its representative ID
- How do we jump to a node quickly?
  - *Also* store a map from value to its node (Omitted in future slides)

find(Farrell) → 1
find(Eric) → 2
find(Farrell) != find(Eric)
find(Farrell) == find(Keanu)



Keanu (1)

Farrell

Howard (3)

Joyce (2)

Melissa

Eric

**Brian** 

## **QuickUnion: Union**

- Key idea: easy to simply rearrange pointers to union entire trees together!
- Which of these implementations would you prefer?
  - Vote in the Participants Panel!



## QuickUnion: Union



- We prefer the right implementation because by changing just the root, we effectively pull the entire tree into the new set!
  - If we change the first node instead, we have to do more work for the rest of the old tree
  - A rare example of constant time work manipulating a factor of n elements

### **QuickUnion: Why bother with the second root?**



- Key idea: will help minimize runtime for future find() calls if we keep the height of the tree short!
  - Pointing directly to the second element would make the tree taller

## **QuickUnion: Checking in on those runtimes**

	Maps to Sets	QuickFind	QuickUnion
<pre>makeSet(value)</pre>	Θ(1)	Θ(1)	Θ(1)
<pre>findSet(value)</pre>	$\Theta(n)$	Θ(1)	$\Theta(n)$
union(x, y)	$\Theta(n)$	$\Theta(n)$	Θ(1) 🜟

\*Only if we discount the runtime from union's calls to find! Otherwise,  $\Theta(n)$ .

 However, for Kruskal's, not a bad assumption: we only ever call union with roots anyway!

```
union(A, B):
  rootA = find(A)
  rootB = find(B)
  set rootA to point to rootB
```



pollev.com/uwcse373

## QuickUnion: Let's Build a Worst Case

Even with the "use-the-roots" implementation of union, try to come up with a series of calls to union that would create a worst-case runtime for find on these Disjoint Sets:

Тор

QuickUnion: Let's Build a Worst Case



find(A):
 jump to A node
 travel upward until root
 return ID

union(A, B):
 rootA = find(A)
 rootB = find(B)
 set rootA to point to rootB

## **I** Poll Everywhere

pollev.com/uwcse373

## QuickUnion: Let's Build a Worst Case

Even with the "use-the-roots" implementation of union, try to come up with a series of calls to union that would create a worst-case runtime for find on these Disjoint Sets:





find(A):
 jump to A node
 travel upward until root
 return ID

union(A, B):
 rootA = find(A)
 rootB = find(B)
 set rootA to point to rootB

## Analyzing the QuickUnion Worst Case

- How did we get a degenerate tree?
  - Even though pointing a root to a root usually helps with this, we can still get a degenerate tree **if we put the root of a large tree under the root of a small tree.**
  - In QuickUnion, rootA always goes under rootB
    - But what if we could ensure the smaller tree went under the larger tree?



### **Lecture Outline**



## WeightedQuickUnion

- Goal: Always pick the smaller tree to go under the larger tree
- Implementation: Store the number of nodes (or "weight") of each tree in the root
  - Constant-time lookup instead of having to traverse the entire tree to count

union(A, B):
 rootA = find(A)
 rootB = find(B)
 put lighter root under heavier root







Perfect! Best runtime we can get.

- union()'s runtime is still dependent on find()'s runtime, which is a function of the tree's height
- What's the worst-case height for WeightedQuickUnion?

```
union(A, B):
  rootA = find(A)
  rootB = find(B)
  put lighter root under heavier root
```























Ν	Н
1	0
2	1
4	2
8	3

- Consider the worst case where the tree height grows as fast as possible
- Worst case tree height is Θ(log N)



## Why Weights Instead of Heights?

- We used the number of items in a tree to decide upon the root
- Why not use the height of the tree?
  - HeightedQuickUnion's runtime is asymptotically the same: Θ(log(n))
  - It's easier to track weights than heights, even though WeightedQuickUnion can lead to some suboptimal structures like this one:



### WeightedQuickUnion Runtime

	Maps to Sets	QuickFind	QuickUnion	WeightedQuickUnion
<pre>makeSet(value)</pre>	Θ(1)	Θ(1)	Θ(1)	Θ(1)
<pre>find(value)</pre>	$\Theta(n)$	Θ(1)	$\Theta(n)$	$\Theta(\log n)$
<pre>union(x, y) assuming root args</pre>	$\Theta(n)$	$\Theta(n)$	Θ(1)	Θ(1)
union(x, y)	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(\log n)$

• This is pretty good! But there's one final optimization we can make: **path compression** 

### **Lecture Outline**



## **Modifying Data Structures for Future Gains**

- Thus far, the modifications we've studied are designed to *preserve invariants* 
  - E.g. Performing rotations to preserve the AVL invariant
  - We rely on those invariants always being true so every call is fast
- Path compression is entirely different: we are modifying the tree structure to *improve future performance* 
  - Not adhering to a specific invariant
  - The first call may be slow, but will optimize so future calls can be fast

### Path Compression: Idea

• This is the worst-case topology if we use WeightedQuickUnion



• Idea: When we do find(15), move all visited nodes under the root

- Additional cost is insignificant (we already have to visit those nodes, just constant time work to point to root too)

### Path Compression: Idea

• This is the worst-case topology if we use WeightedQuickUnion



- Idea: When we do find(15), move all visited nodes under the root
  - Additional cost is insignificant (we already have to visit those nodes, just constant time work to point to root too)
- Perform Path Compression on every find(), so future calls to find() are faster!