## CSE 373: Disjoint sets

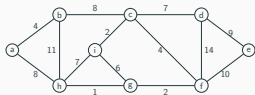Michael Lee
Wednesday, Feb 28, 2018

---

Last time...

▶ **Prim's algorithm:**
  Nearly identical to Dijkstra's, except we use the distance to any already-visited node as the cost.
▶ **Kruskal's algorithm:**
  Loop over edges, from smallest to largest. Use the edge only if it doesn't introduce a cycle.
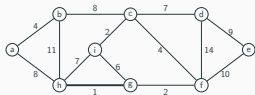
---

**Kruskal's algorithm: example with a weighted graph**

Example of the algorithm:

---

**Kruskal's algorithm: example with a weighted graph**

Example of the algorithm:

---

**Kruskal's algorithm: example with a weighted graph**

Example of the algorithm:

---

**Kruskal's algorithm: example with a weighted graph**
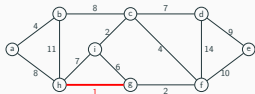
Example of the algorithm:

Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

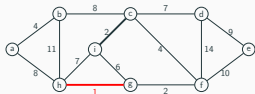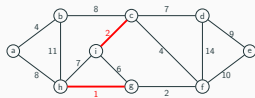Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:

Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

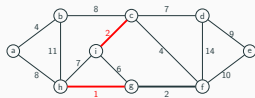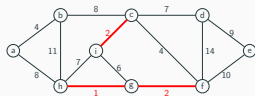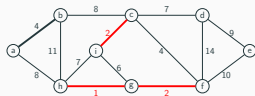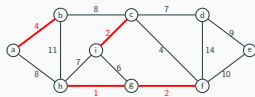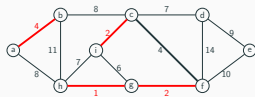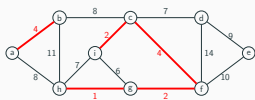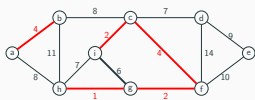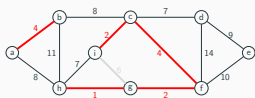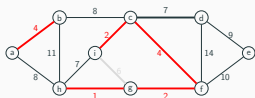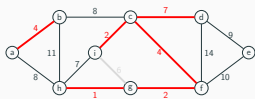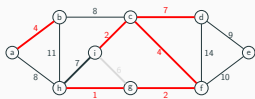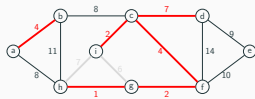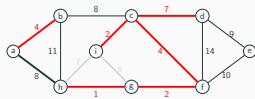Example of the algorithm:



Kruskal's algorithm: example with a weighted graph

Example of the algorithm:

Kruskal's algorithm: example with a weighted graph

Example of the algorithm:

Kruskal's algorithm: example with a weighted graph

Example of the algorithm:

Example of the algorithm:

---

Runtime analysis:

```
def kruskal():
    for (v : vertices):
        makeMST(v)

    sort edges in ascending order by their weight

    mst = new SomeSet<Edge>()
    for (edge : edges):
        if findMST(edge.src) != findMST(edge.dst):
            union(edge.src, edge.dst)
            mst.add(edge)

    return mst
```

Note: assume that...

▶ makeMST(v) takes $\mathcal{O}(t_m)$ time
▶ findMST(v): takes $\mathcal{O}(t_f)$ time
▶ union(u, v): takes $\mathcal{O}(t_u)$ time

---

▶ Making the $|V|$ MSTs takes $\mathcal{O}(|V| \cdot t_m)$ time
▶ Sorting the edges takes $\mathcal{O}(|E| \cdot \log(|E|))$ time, assuming we use a general-purpose comparison sort
▶ The final loop takes $\mathcal{O}(|E| \cdot t_f + |V| \cdot t_u)$ time

Putting it all together:

$$\mathcal{O}(|V| \cdot t_m + |E| \cdot \log(|E|)) + |E| \cdot t_f + |V| \cdot t_u)$$

---

But wait, what exactly is $t_m$, $t_f$, and $t_u$? How exactly do we implement makeMST(v), findMST(v), and union(u, v)?

We can do so using a new ADT called the DisjointSet ADT!

---

Review: what is a set?

▶ A set is a "bag" of elements arranged in no particular order.
▶ A set may not contain duplicates.

We implemented a set in project 2: ChainedHashSet

Interesting note: sets come up all the time in math.

---

Properties of a disjoint-set data structure:

▶ A disjoint-set data structure maintains a collection of many different sets.
▶ An item may not be contained within multiple sets. Each set must be disjoint.
▶ Each set is associated with some representative. What is a representative? Any sort of unique "identifier". Examples:
  ▶ We could pick some arbitrary element in the set to be the "representative"
  ▶ We could assign each set some unique integer id.

A disjoint-set has the following core operations:

- **makeSet(x)** – Creates a new set where the only member is $x$. We assign that set a representative.
- **findSet(x)** – Looks up the set containing $x$. Then, returns the representative of that set.
- **union(x, y)** – Looks up the set containing $x$ and the set containing $y$. We combine these two sets together into one. We (arbitrarily) pick one of the two representatives to be the representative of this new set.

9

---

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)
```

10

---

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)
```

Rep: 4
e

Rep: 2     Rep: 3
c          d

Rep: 0     Rep: 1
a          b

10

---

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)
print(findSet(a))
print(findSet(d))
```

Rep: 4
e

Rep: 2     Rep: 3
c          d

Rep: 0     Rep: 1
a          b

10

---

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)
print(findSet(a))
print(findSet(d))
union(a, c)
union(b, d)
print(findSet(a) == findSet(c))
print(findSet(a) == findSet(d))
```

Rep: 4
e

Rep: 2     Rep: 3
c          d

Rep: 0     Rep: 1
a          b

10

---

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)
print(findSet(a))
print(findSet(d))
union(a, c)
union(b, d)
print(findSet(a) == findSet(c))
print(findSet(a) == findSet(d))
```

Rep: 4
e

Rep: 0     Rep: 1
c          d
a          b

10

## The DisjointSet ADT

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)

print(findSet(a))
print(findSet(d))

union(a, c)
union(b, d)
print(findSet(a) == findSet(c))
print(findSet(a) == findSet(d))

union(c, b)
print(findSet(a) == findSet(d))
```

Rep: 4

e

Rep: 0        Rep: 1

c            d

a            b

---

## The DisjointSet ADT

Example:
```
makeSet(a)
makeSet(b)
makeSet(c)
makeSet(d)
makeSet(e)

print(findSet(a))
print(findSet(d))

union(a, c)
union(b, d)
print(findSet(a) == findSet(c))
print(findSet(a) == findSet(d))

union(c, b)
print(findSet(a) == findSet(d))
```

Rep: 4

e

Rep: 0

c            d

a            b

---

## The DisjointSet ADT

What operations does a disjoint-set **NOT** support?

**Answer:** The ability to actually get the entire set.

We can *make* a set, *check* if an item is in a set, and *combine* two sets, but we don't have a built-in way of *getting* the entire set itself.

**Insight:** The few operations we need to support, the more creative our implementation can be.

(If the client really wants the sets, they can get it themselves in $\mathcal{O}(n)$ time – how?)

---

## DisjointSet: implementation

So, how do we implement these?

**Core idea:**
- ▶ We represent each set as a tree
- ▶ The disjoint-set keeps track of a "forest" of trees

**Intuitions:**
- ▶ We want union-ing to be cheap.
  Combining two trees is cheap; we just manipulate pointers.
- ▶ We want a single "representative" per set.
  A tree has a single root!

---

## DisjointSet: implementation

High-level overview:

- ▶ **makeSet(x)**: Adds a new tree (of size 1) to our "forest"
- ▶ **findSet(x)**: Looks up the node, then finds root of tree
- ▶ **union(x, y)**: Combines two trees into one

---

## DisjointSet: implementation

Suppose we call makeSet(...) on 0 through 5.

(0)   (1)   (2)   (3)   (4)   (5)

Each makeSet(...) adds a new tree to our "forest".
Note that right now, each tree has only one element.

Suppose we call union(3, 5).



We combine those two trees into one.

Assumption: we have an $\mathcal{O}(1)$ way of getting each node.
(E.g. maintain a hashmap of numbers to node objects.)

**Question:** how do we implement findSet(...)?

Once we find a node, move upwards until we're looking at root.

Then, return the root's data field.

15

---

Suppose we call union(5, 4).



**Algorithm:** Find the roots of both trees and add one tree as a subchild of the other.

Which tree becomes the new root? For now, pick randomly.

16

---

Suppose we call union(0, 1), then union(2, 0).



17

---

Now, suppose we call union(2, 4). What happens?



18

---

Now, suppose we call union(2, 4). What happens?



We look up 2 and 3, find their roots, and nest one tree inside the other

18

---

What's the worst-case runtime of our methods?

Better question: are our trees guaranteed to be balanced?

Hint: When union-ing, we pick which tree is nested randomly. Does that guarantee we'll get a balanced tree?

19

## DisjointSet: Analysis

The worst-case scenario:



Possible outcome of calling union(0, 5)

---

## DisjointSet: implementation

So, what are the worst-case runtimes?

- **makeSet(x)**:
  $\mathcal{O}(1)$ – creating the tree takes constant time
- **findSet(x)**:
  $\mathcal{O}(n)$ – if it's a linked list, we need to traverse $n$ elements!
- **union(x, y)**:
  $\mathcal{O}(n)$ – union calls findSet(...) on both elements

...where $n$ is the total number of items added to the disjoint-set.

---

## Improving DisjointSet

How can we improve disjoint sets?

1. **Union-by-rank:**
   Strategy to make sure trees are balanced
2. **Path compression:**
   Hijack findSet(x) and make it do a little extra work to improve overall performance.
3. **Array representation:**
   Takes advantage of cache locality, simplifies implementation, etc.

---

## Union-by-rank

**Problem:** Our trees could be unbalanced

**Solution:**

Let rank(x) be a number representing the upper-bound of the height of x. So, rank(x) $\geq$ height(x).

We then...

1. Keep track of the rank of all trees.
2. When unioning, make the tree with the larger rank the root!
3. If it's a tie, pick one randomly and increase the rank by one.

(Why not keep track of the height? When we look at path compression, keeping track of the height becomes more challenging.)

---

## Union-by-rank

Example: Suppose we call union(1, 5)?

The tree with the root of "6" has the larger rank, so we make it the root.

Note: we're not really "removing" the rank from node 0 – it's just

---

## Union-by-rank

Example: Suppose we call union(1, 5)?

The tree with the root of "6" has the larger rank, so we make it the root.

Note: we're not really "removing" the rank from node 0 – it's just

## Union-by-rank

Example: Suppose we call union(1, 5)?



The tree with the root of "6" has the larger rank, so we make it the root.

Note: we're not really "removing" the rank from node 0 – it's just irrelevant, so we're ignoring it and omitting it from the diagram to save space. We only care about the ranks at the roots.

24

## Union-by-rank

Example: Suppose we call union(5, 11)?



Here, there's a tie. We break the tie arbitrarily, and increment the rank of the new tree by one.

25

## Union-by-rank

Net effect? Our trees stay relatively balanced.

So, what are the worst-case runtimes now?

► makeSet(x):
  $\mathcal{O}(1)$ – still the same
► findSet(x):
  $\mathcal{O}(\log(n))$ – since the tree is balanced
► union(x, y):
  $\mathcal{O}(\log(n))$ – since union calls findSet

26

## Path compression

Consider the following forest:



Suppose we call findSet(3) a few hundred times.

Why do we have to keep finding the root again and again?

27

## Path compression

**Observation:** To find root, we must also traverse these nodes:



What if, next time, we could just jump straight to the root?
Same for the other nodes we visited

28

## Path compression

**Observation:** To find root, we must also traverse these nodes:



What if, next time, we could just jump straight to the root?
Same for the other nodes we visited

28

**Slide (page 28):** Path compression

**Observation:** To find root, we must also traverse these nodes:

What if, next time, we could just jump straight to the root?
Same for the other nodes we visited

28

**Slide (page 29):** Path compression

So, let's do it!

Now what happens if we try calling findSet(3)?

29

**Slide (page 29):** Path compression

So, let's do it!

Now what happens if we try calling findSet(3)?

29

**Slide (page 29):** Path compression

So, let's do it!

Now what happens if we try calling findSet(3)?

29

**Slide (page 29):** Path compression

So, let's do it!

Now what happens if we try calling findSet(3)?

29

**Slide (page 30):** Path compression

One additional note: path compression changes the heights of our trees.

This means it could be the case that rank ≠ height.

Is this a problem?

**Answer:** No; proof is beyond the scope of this class

30

## Path compression: runtime

Now, what are the worst-case and best-case runtime of the following?

- `makeSet(x)`:
  $\mathcal{O}(1)$ – still the same
- `findSet(x)`:
  In the best case, $\mathcal{O}(1)$, in the worst case $\mathcal{O}(\log(n))$
- `union(x, y)`:
  In the best case, $\mathcal{O}(1)$, in the worst case $\mathcal{O}(\log(n))$

## Back to Kruskal's

Why are we doing this? To help us implement Kruskal's algorithm:

```
def kruskal():
    for (v : vertices):
        makeMST(v)

    sort edges in ascending order by their weight

    mst = new SomeSet<Edge>()
    for (edge : edges):
        if findMST(edge.arc) != findMST(edge.dst):
            union(edge.arc, edge.dst)
            mst.add(edge)

    return mst
```

- `makeMST(v)` takes $\mathcal{O}(t_m)$ time
- `findMST(v)`: takes $\mathcal{O}(t_f)$ time
- `union(u, v)`: takes $\mathcal{O}(t_u)$ time

## Back to Kruskal's

We concluded that the runtime is:

$$\mathcal{O}\left(\underbrace{|V| \cdot t_m}_{\text{setup}} + \underbrace{|E| \cdot \log(|E|)}_{\text{sorting edges}} + \underbrace{|E| \cdot t_f + |V| \cdot t_u}_{\text{core loop}}\right)$$

Well, we just said that in the worst case:

- $t_m \in \mathcal{O}(1)$
- $t_f \in \mathcal{O}(\log(|V|))$
- $t_u \in \mathcal{O}(\log(|V|))$

So the worst-case overall runtime of Kruskal's is:

$$\mathcal{O}(|V| + |E| \cdot \log(|E|) + (|E| + |V|) \cdot \log(|V|))$$

## Back to Kruskal's

Our worst-case runtime:

$$\mathcal{O}(|V| + |E| \cdot \log(|E|) + (|E| + |V|) \cdot \log(|V|))$$

One minor improvement: since our edge weights are numbers, we can likely use a *linear sort* and improve the runtime to:

$$\mathcal{O}(|V| + |E| + (|E| + |V|) \cdot \log(|V|))$$

We can drop the $|V| + |E|$, since they're dominated by the last term:

$$\mathcal{O}((|E| + |V|) \cdot \log(|V|))$$

...and we're left with something that's basically the same as Prim's algorithm.

## Disjoint-sets, amortized analysis

...or are we?

**Observation:** each call to `findSet(x)` improves all future calls. How much of a difference does that make?

Interesting result:

It turns out `union` and `find` are *amortized* $\log^*(n)$.

## Disjoint-sets, amortized analysis

**Iterated log**

The expression $\log^*(n)$ is equivalent to the number of times you need to compute $\log(x)$ to bring the value down to at most 1

Example:

- $\log^*(2) = \log(2) = 1$
- $\log^*(4) = \log(\log(4)) = 2$
- $\log^*(8) = \log(\log(\log(8))) = 3$
- $\log^*(65536) = \log^*(2^{2^{2^2}}) = 4$
- $\log^*(2^{65536}) = \ldots = 5$

What is $2^{65536}$?

$2^{65536} =$

2003529930406846464649790723515602557504478254755697514419
2650169737108940595563114530895061308809333481010382343429072
6318182294938211881266886950636476154702916016504187191635158796
6347219442930927982084309104855900570159318959639524863372367
2030029166695921561087649488925450985089114570376752085002063
71563702366126359747144807111748158809141357427209671901 51 83
6282560618091458852699826141425030123391082736038437678764 49
0432059603791244900537556031403507616256247603186379312648 47
0374378295497561377098160461441330869211810248595915238019533
1030292162800160586670105651646750568378424963842244845292
5373614425336143737290883037946012747249584148649159306472520
15155669392628188069165079631086433277530726714399815850881129
26289011342377827055674210800700652839633221550778312142885 51
675554073451072131124273990546299821971769155054839205223804 35
70 4584819795639315785351001899200002414196370681355984046403947
2194016069517690156119726982337890017641517190051133466306898
1402193834814354263873065395529696913880241581618595611006403
6211979610185953402871672001226046424923851113934004635162
3867567078745294646709038865477434832178970127644555529409092
0219595857516229733335761595523948852975799540284719435299135

37

4376370598692891375715374000198639433246489005254310662966916
5243419174691389632476560289415199775477703313806478134230950 96
1909606545913008901888875880847336259560654448885014473357060
5881709016210849971452956834406197969056546981363116205357936
9791403236328496233046421066136200220175787851857409162050489
7117818204018728293994344618622432800983732376493181478984 81
1945271300744022076568091037620399920349202390066262644919091 6
7985461515778839060397720759227937885224129430101745808686226 3
3692847258514030396155585643303854506886522131148136384083847
7826379045960718687672850976347127198889068047824323039471865
0525660078150729861141430305816927924971409161059417185352275
8875044775922183011587807019755357222414000195481020056617735
8978149953232520858975346354700778669040642901676380816174055
0405117670093673202804549339027992491867306539931640720492238
4748152806191669009338057321208163507076343516698696250209690
2316285935007187419057916124153689751480826190484794657173660
1005892476655445840838347905441448176842553272073150863493 47
6051374197795251903650321980201087647383686825310251833775339
0886142618480037400808223810407646887847164755294532694766170
0424461063311238021134588694532200116564076327023074292426051

38

6340696503084422585596703927186946115851379338647569974856867
0079823960604394788516649260304945061743412365828352144806
7266768418070837548622114082365798029612000274413244384324023
3125740354501935242877643080235085085860899627744581646808 5
7875115807014743674819676950499916439982843572904153781434 38
8473034842619033884149403136619854257635577105335580206622 1
8557706008255128889333229643628198438613239570671914096385 3
3832374343758830859233722284642879962456054769324289984326 52
6773783731732806321075321123868060467470842805116648870908 47
7029120816110491255559832236624486855665140268464120969498 259
0565519216188104341226838996283071654868525536914850299539675
5039549383718534059009618748947399288043249637316575380367 35
8671017578394818471798498246948060532081996066183434012476 69
6639519778021441199752546704080608499344178256285097726523709
8986515394621930046073645079262129759176982938923670151709920
9153156781443971248475706237804600009182933230680570046559
1458387208088016887445835976929258465124763087148566313528934
1661174906175266714926721761283308457393646024458289257133 88
7783905630048248379983996920292222154861459023734782268252163
9957440801727144146179559226175083889020074169926238300282286

39

...I got tired of copying and pasting, but we're not even a fourth of the way through.

Punchline? $\log^*(n) \leq 5$, for basically any reasonable value of $n$.

Runtime of Kruskal? $\mathcal{O}\left((|E| + |V|)\log^*(|V|)\right) \approx \mathcal{O}\left(|E| + |V|\right)$

40

But wait!

Somebody then came along and proved that find and union are amortized $\mathcal{O}(\alpha(n))$ – the inverse of the Ackermann function.

This grows even more slowly then $\log^*(n)$!

41