

Two Heuristics for Disjoint Set Forests

Union by rank

for each node x , maintain attribute $\text{rank}[x]$ which is an upper bound for $\text{height}(x)$:

$$\text{height}(x) \leq \text{rank}[x]$$

Modify $\text{Union}(\cdot, \cdot)$ so root with smaller rank points to root with higher rank.

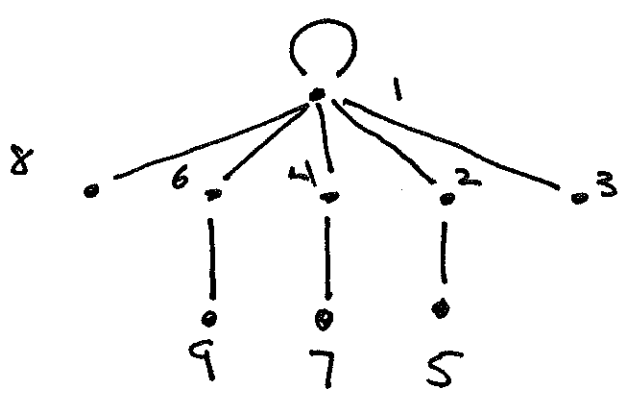
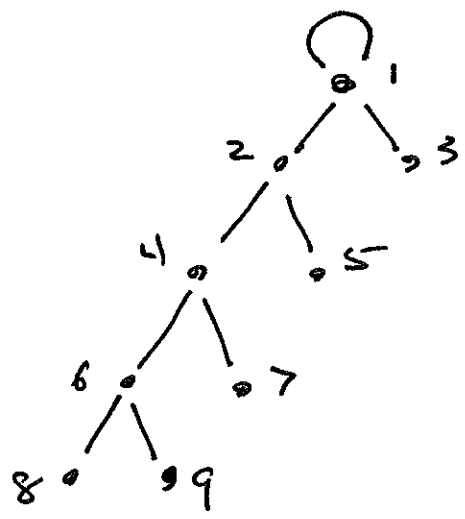
Path Compression

Modify FindSet() to perform

2 Passes :

- up the find-path to find root
- down the find-path setting each node to have root as parent.

Ex.



→
FindSet(8)

MakeSet(x)

1. $P[x] = x$
2. $\text{rank}[x] = 0$

Union(x, y) (union-by-rank)

1. Link (FindSet(x), FindSet(y))

Link(x, y) (Pre: x, y both roots, $x \neq y$)

1. if $\text{rank}[x] > \text{rank}[y]$
2. $P[y] = x$
3. else
4. $P[x] = y$
5. if $\text{rank}[x] == \text{rank}[y]$
6. $\text{rank}[y]++$

Exercise

As long as only MakeSet and Link are called, $\text{rank}[x] = \text{height}(x)$ for all x . Prove by induction on # of link operations.

FindSet(x)

1. if $x \neq p[x]$
2. $p[x] = \text{FindSet}(p[x])$
3. return $p[x]$

Trace on Previous example!

- FindSet(8) !
1. $8 \neq p[8]$
 2. $p[8] = \text{FindSet}(6)$
 3. return 1

FindSet(6): 1. $6 \neq P[6]$
 2. $P[6] = \text{FindSet}(4)$
 3. return 1

FindSet(4): 1. $4 \neq P[4]$
 2. $P[4] = \text{FindSet}(2)$
 3. return 1

FindSet(2): 1. $2 \neq P[2]$
 2. $P[2] = \text{FindSet}(1)$
 3. return 1

FindSet(1): 1. $1 \neq P[1]$ (false)
 2. _____
 3. return 1

✓

Note FindSet(x) may reduce height(x), but it does not alter ranks. Thus

$$\text{height}(x) \leq \text{rank}[x]$$

for all nodes x .

Theorem

using union-by-rank and Path-compression heuristics, any sequence of m MakeSet, Union & FindSet operations (n of which are MakeSets) can be performed in time

$$O(m \cdot \alpha(n))$$

□

Here $\alpha(n)$ is an extremely slow growing function. In any conceivable application $\alpha(n) \leq 4$.

So effectively, $\alpha(n) = \text{const.}$, and runtime $O(m)$.

Amortizing over all m operations, gives run time

$$\frac{O(m)}{m} = O(1)$$

Per operation.

BACK TO SEC. 23.2

RECALL THAT KRUSKAL'S ALGORITHM GROWS SEVERAL TREES SIMULTANEOUSLY UNTIL THEY MERGE INTO A SINGLE MWST. THE IMPLEMENTATION BELOW USES A DISJOINT SET DATA STRUCTURE TO MAINTAIN A COLLECTION OF VERTEX SETS. EACH SET IN THIS COLLECTION IS THE VERTEX SET OF ONE OF THE TREES UNDER CONSTRUCTION.

NOTE THE SIMILARITY OF THIS ALGORITHM TO THE ALGORITHM FOR DETERMINING CONNECTED COMPONENTS FROM 21.1.

KRUSKAL(G, w)

- 1.) $A = \emptyset$
- 2.) for each $v \in V$
- 3.) MakeSet(v)
- 4.) SORT E BY WEIGHT w (NON-DECREASING)
- 5.) for each $(u, v) \in E$ (IN SORTED ORDER)
- 6.) if FindSet(u) \neq FindSet(v)
- 7.) $A = A \cup \{(u, v)\}$
- 8.) Union(u, v)
- 9.) RETURN A

ON EACH ITERATION OF LOOP 5-8 WE PICK THE LIGHTEST EDGE JOINING TWO DISTINCT TREES AND ADD THAT TO A . I.E. AMONGST ALL EDGES WHOSE ADDITION TO A WOULD NOT CREATE A CYCLE, WE PICK ONE OF MINIMUM WEIGHT. THE RESULTING SET $A \subseteq E$ CONSTITUTES THE EDGES OF A MWST IN G .

Run Time

ASSUME WE USE THE DISJOINT SET FOREST IMPLEMENTATION WITH UNION-BY-RANK AND PATH COMPRESSION HEURISTICS, WHICH IS THE ASYMPTOTICALLY FASTEST KNOWN.

INITIALIZATION OF A IN LINE 1 COSTS $O(V)$ TIME. LOOP 2-3 DOES V MAKESET OPERATIONS AND LOOP 5-8 $O(E)$ FINDSET AND UNION OPERATIONS. THUS THE DISJOINT SET OPERATIONS HAVE TOTAL COST $O((V+E) \cdot \alpha(V))$.

NOW SINCE G IS CONNECTED (OTHERWISE THERE IS NO SPANNING TREE) WE HAVE $E \geq V - 1$, WHENCE $V \leq E + 1 = O(E)$, SO THE DISJOINT SET OPERATIONS COST $O(E \alpha(V))$

THE SORT ON LINE 4 COSTS $O(E \lg E)$ TIME (ASSUMING WE USE MERGESORT OR HEAPSORT, OR A SIMILARLY EFFICIENT SORT.)

ALSO NOTE THAT $\alpha(V) = O(\lg V) = O(\lg E)$, SO THE RUN TIME OF KRUSKAL IS

$$O(E \lg E)$$

BUT ALSO SINCE G IS SIMPLE, $E < V^2$ SO $\lg E = O(\lg V)$, AND THE RUN TIME OF KRUSKAL CAN BE EXPRESSED AS

$$O(E \lg V)$$

WHICH IS THE SAME AS PRIM.