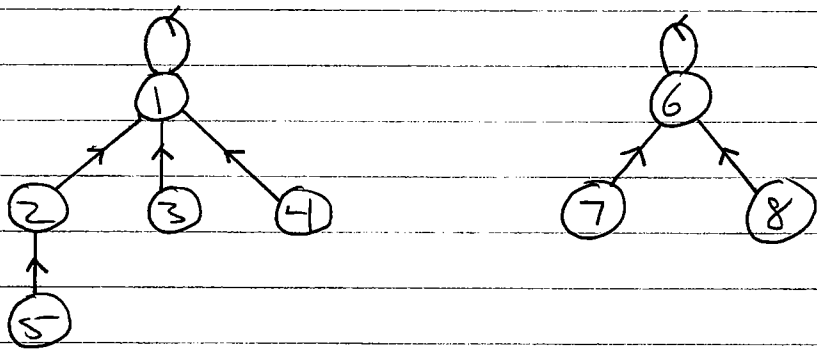


(21.3) DISJOINT SET FORESTS

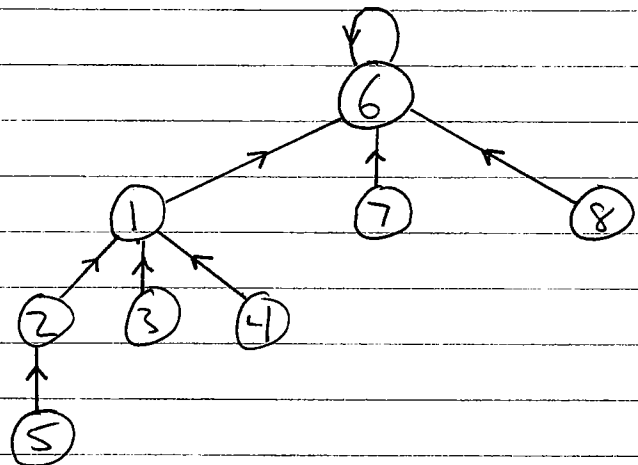
ANOTHER IMPLEMENTATION STRATEGY IS TO REPRESENT EACH SET IN \mathcal{S} BY A ROOTED TREE, i.e. \mathcal{S} IS A DISJOINT SET FOREST.

EACH NODE POINTS ONLY TO ITS PARENT. THE ROOT OF EACH TREE IS ITS OWN PARENT, AND IS ITS SET'S REPRESENTATIVE.

Ex



RESULT OF Union(3, 8) :



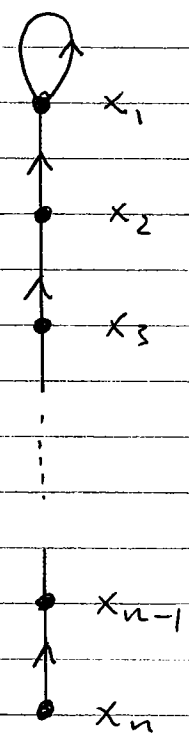
CONVENTION: Union(x, y) CAUSES REP[x] TO POINT TO REP[y].

MakeSet CREATES A TREE WITH JUST ONE NODE WHICH POINTS TO ITSELF. Union CAUSES THE ROOT OF ONE TREE TO POINT TO THE ROOT OF ANOTHER. BOTH OPERATIONS HAVE COST $\Theta(1)$.

FindSet FOLLOWS PARENT POINTERS UNTIL THE ROOT IS REACHED. THE NODES VISITED ON THIS PATH TOWARD THE ROOT CONSTITUTE A FIND PATH. THE COST OF FindSet IS $\Theta(\text{Length-of-FindPath})$.

IT IS POSSIBLE USING THESE OPERATIONS TO CREATE A TREE WHICH IS SIMPLY A LINEAR LINKED LIST.

EX
 MakeSet(x_1)
 ⋮
 MakeSet(x_n)
 Union(x_n, x_{n-1})
 Union(x_{n-1}, x_{n-2})
 ⋮
 Union(x_3, x_2)
 Union(x_2, x_1)



THIS SEQUENCE IS NOT BAD IN ITSELF (RUN TIME $\Theta(2n-1) = \Theta(m)$), BUT IF CREATED A TREE IN WHICH FindSet will run slowly (i.e. in $\Theta(n)$ TIME.)

TWO HEURISTICS ARE USED TO IMPROVE THE RUN TIME OF SEQUENCES WHICH INCLUDE MANY FindSet OPERATIONS.

Union By RANK

FOR EACH NODE x , MAINTAIN A VALUE $\text{rank}[x]$, WHICH IS AN UPPER BOUND ON THE HEIGHT OF x .

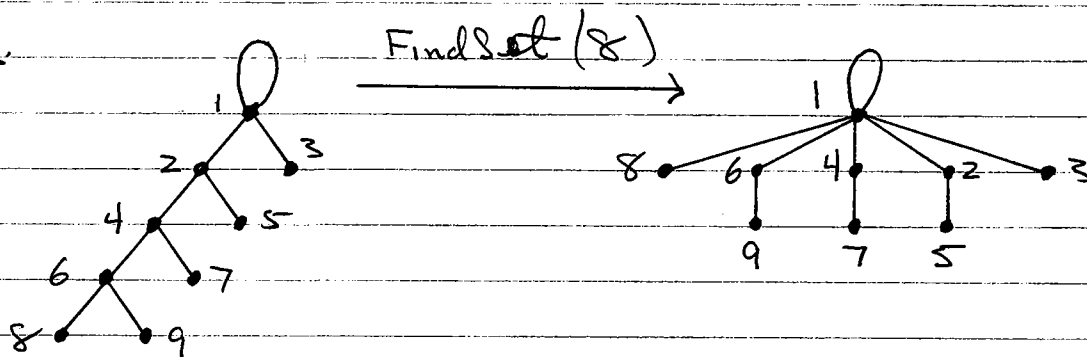
$$\text{height}(x) \leq \text{rank}[x].$$

MODIFY Union SO THAT THE ROOT WITH SMALLER RANK IS MADE TO POINT TO THE ROOT WITH LARGER RANK.

PATH COMPRESSION

THIS IS A MODIFICATION OF FindSet WHICH CAN REDUCE THE HEIGHT OF A TREE. IT INVOLVES TWO PASSES: PERFORM ONE PASS UP THE FIND PATH TO FIND THE ROOT, THEN MAKE ANOTHER PASS DOWN THE FIND PATH CAUSING EACH NODE ALONG THAT PATH TO POINT TO THE ROOT.

Ex.



FindSet(8) ALONG THE TREE AS ABOVE, THEN RETURNS 1.

MakeSet(x)

- 1.) $P[x] \leftarrow x$
- 2.) $rank[x] \leftarrow 0$

Union(x, y) (using UNION-BY-RANK)

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

Link(x, y) (Pre: x, y BOTH ROOT NODES)

- 1.) if $rank[x] > rank[y]$
- 2.) $P[y] \leftarrow x$
- 3.) else
- 4.) $P[x] \leftarrow y$
- 5.) if $rank[x] = rank[y]$
- 6.) $rank[y] \leftarrow rank[y] + 1$

NOTE: AS LONG AS ONLY MakeSet AND Link OPERATIONS ARE PERFORMED, $rank[x]$ IS ACTUALLY EQUAL TO height(x). THIS CAN BE PROVED BY INDUCTION ON THE NUMBER OF Link OPERATIONS PERFORMED.

FindSet(x) (using PATH COMPRESSION)

- 1.) if $x \neq P[x]$
- 2.) $P[x] \leftarrow \text{FindSet}(P[x])$
- 3.) return $P[x]$.

WE TRACE FindSet ON THE PRECEDING EXAMPLE:

- FindSet(8):
- 1.) $8 \neq P[8]$
 - 2.) $P[8] \leftarrow \text{FindSet}(6)$
 - 3.) return 1

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

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

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

- FindSet(1)
- 1.) $1 = P[1]$
 - 2.)
 - 3.) return 1

NOTE THAT $\text{FindSet}(x)$ MAY REDUCE THE HEIGHT OF THE TREE CONTAINING x , WHILE IT DOES NOT ALTER ANY RANKS. THUS FindSet OPERATIONS MAINTAIN THE INEQUALITY

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

FOR ALL NODES x .

THE UNION-BY-RANK AND PATH COMPRESSION HEURISTICS WORK TOGETHER TO KEEP THE HEIGHTS OF TREES SMALL, WHICH LOWERS THE COST OF FindSet .

THEOREM

USING THE UNION-BY-RANK AND PATH COMPRESSION HEURISTICS, ANY SEQUENCE OF m MAKESET, UNION, AND 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 (DEFINED IN 21.4), IN ANY CONCEIVABLE APPLICATION $\alpha(n) \leq 4$, SO THE ABOVE RUN TIME IS EFFECTIVELY LINEAR, I.E. $O(m)$.

AMMORTIZING THE COST OVER ALL OPERATIONS GIVES $O(1)$ TIME PER OPERATION ON AVERAGE, WHICH IS VERY EFFICIENT.

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.) $A \leftarrow \emptyset$
- 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.) $A \leftarrow A \cup \{(u, v)\}$
- 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.