

(21.1) DISJOINT SETS } Chap 21
21.1, 21.2, 21.3

A DISJOINT SET ADT MAINTAINS A COLLECTION OF DYNAMIC SETS

$$\mathcal{S} = \{S_1, S_2, \dots, S_n\}$$

WHICH ARE PAIRWISE DISJOINT: $S_i \cap S_j = \emptyset$ FOR $i \neq j$. EACH SET IN THE COLLECTION IS IDENTIFIED BY A DISTINGUISHED MEMBER CALLED ITS REPRESENTATIVE.

EX.

$$\mathcal{S} = \{ \{a, b, c\}, \{d, e\}, \{f\}, \{g, h\} \}$$

IN MOST APPLICATIONS WE DON'T CARE WHICH MEMBER OF A SET IS USED AS ITS REPRESENTATIVE.

WE WRITE S_x FOR THE (UNIQUE) MEMBER OF \mathcal{S} CONTAINING x . NOTE x NEEDS NOT BE THE REPRESENTATIVE OF S_x . IN THE ABOVE EXAMPLE

$$\mathcal{S} = \{S_a, S_d, S_f, S_g\}$$

A DISJOINT SET ADT THE FOLLOWING OPERATIONS.

SUPPORTS

MakeSet(x)

CREATES A NEW SET $\{x\}$ IN \mathcal{S} . x MAY NOT BELONG TO ANY EXISTING MEMBER OF \mathcal{S} SINCE THE COLLECTION IS TO BE DISJOINT.

$$\mathcal{S} = \mathcal{S} \cup \{\{x\}\}$$

Union(x,y)

UNITE SETS S_x AND S_y . PRECONDITION: $S_x \neq S_y$. THIS OPERATION REDUCES THE NUMBER OF SETS IN \mathcal{S} BY 1.

$$\mathcal{S} = \mathcal{S} - \{S_x\} - \{S_y\} \cup \{S_x \cup S_y\}$$

FindSet(x)

RETURNS (A POINTER TO) THE REPRESENTATIVE OF THE SET S_x .

IN WHAT FOLLOWS WE WILL ANALYZE RUN TIMES OF SEQUENCES OF THESE OPERATIONS. WE ADOPT THE CONVENTIONS THAT

$$n = \# \text{MakeSet}$$

AND

$$m = (\# \text{MakeSet}) + (\# \text{Union}) + (\# \text{FindSet})$$

OPERATIONS. ASSUME THAT SINCE UNION REDUCES $|\mathcal{S}|$ BY 1, WE HAVE

$$(\# \text{union}) \leq n - 1$$

WE ANALYZE SEQUENCES OF OPERATIONS SINCE, AS WE SHALL SEE, THE RUN TIME OF THESE OPERATIONS WILL DEPEND ON THE PARTICULAR HISTORY OF THE DATA STRUCTURE.

TWO QUESTIONS ARISE: (1) OF WHAT USE IS THIS NEW DATA STRUCTURE?, AND (2) HOW CAN IT BE IMPLEMENTED?

TO ANSWER (1) WE PRESENT THE FOLLOWING ALGORITHM WHICH DETERMINES THE CONNECTED COMPONENTS OF AN UNSIRECTED GRAPH $G = (V, E)$.

Components (G)

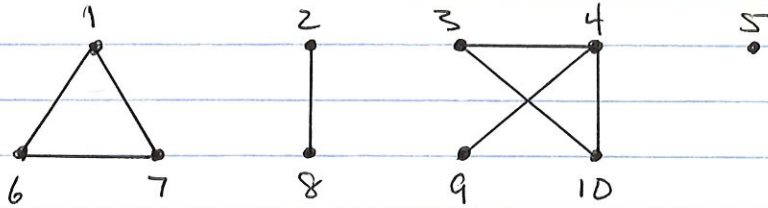
- 1.) for each $v \in V$
- 2.) MakeSet(v)
- 3.) for each $(u, v) \in E$
- 4.) if FindSet(u) \neq FindSet(v)
- 5.) Union(u, v)

SameComponent(u, v)

- 1.) if FindSet(u) == FindSet(v)
- 2.) Return true
- 3.) return false

ANOTHER USE OF THE DISJOINT SET DATA STRUCTURE IS TO DETERMINE A MWST, AS WE'LL SEE WHEN WE RETURN TO CHAP. 23.

EX



EDGE	
	{1}, {2}, {3}, {4}, {5}, {6}, {7}, {8}, {9}, {10}
(1,7)	{1,7}, {2}, {3}, {4}, {5}, {6}, {8}, {9}, {10}
(2,8)	{1,7}, {2,8}, {3}, {4}, {5}, {6}, {9}, {10}
(3,4)	{1,7}, {2,8}, {3,4}, {5}, {6}, {9}, {10}
(6,7)	{1,6,7}, {2,8}, {3,4}, {5}, {9}, {10}
(9,4)	{1,6,7}, {2,8}, {3,4,9}, {5}, {10}
(4,10)	{1,6,7}, {2,8}, {3,4,9,10}, {5}
(1,6)	" " " "
(3,10)	" " " "

SECTIONS 21.2 & 21.3 PRESENT TWO METHODS OF IMPLEMENTING A DISJOINT SET DATA STRUCTURE, WHICH ANSWERS QUESTION (2).

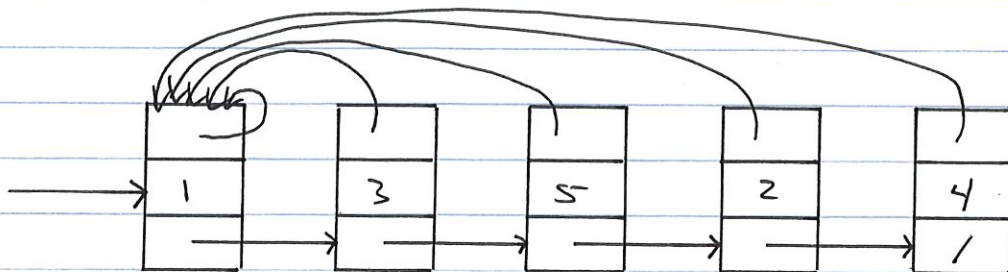
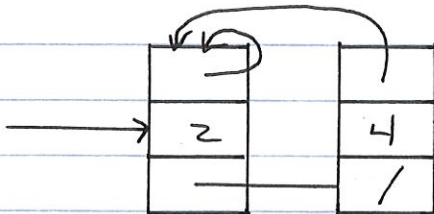
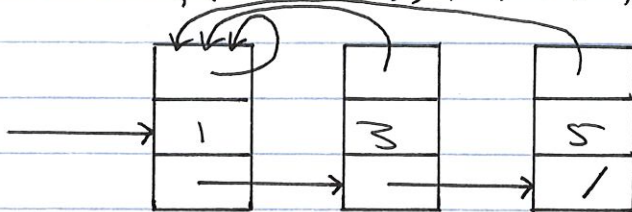
(21.2) LINKED LIST REPRESENTATION

THE SETS IN \mathcal{S} CAN BE REPRESENTED BY LINKED LISTS WHERE THE FIRST OBJECT IN A LIST SERVES AS THAT SET'S REPRESENTATIVE

EACH LIST NODE MAINTAINS A KEY FIELD AS WELL AS POINTERS TO THE NEXT NODE AND THE REPRESENTATIVE NODE

	Rep
NODE	Key
	next

Ex. $\mathcal{S} = \{\{1, 3, 5\}, \{2, 4\}\}$



RESULT OF Union(3,4)

MakeSet(x)

1.) CREATE NEW list with some object x

FindSet(x)

1.) return Rep[x]

Union(x, y)

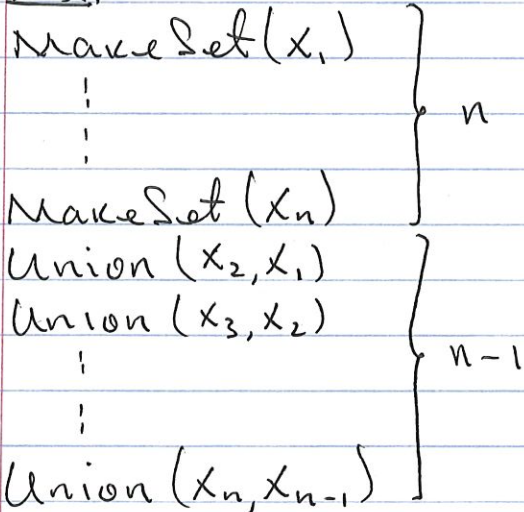
1.) CONCATENATE lists S_x AND S_y

2.) UPDATE Rep POINTERS in S_y TO POINT TO Rep[x]

MakeSet AND FindSet RUN in $\Theta(1)$ TIME
WHILE Union TAKES $\Theta(|S_y|)$ TIME.

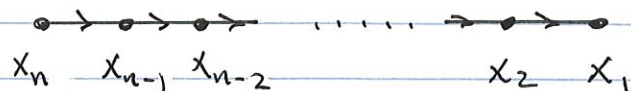
THE $\Theta(|S_y|)$ RUN TIME FOR Union CAUSES CERTAIN SEQUENCES OF OPERATIONS TO HAVE POOR TOTAL RUN TIME.

Ex.



$$m = n + (n-1) = 2n - 1$$

THIS SEQUENCE CREATES A SINGLE LINKED LIST



THE n MAKESET OPERATIONS TAKE $\Theta(n)$ TIME.
 THE $n-1$ UNION OPERATIONS PERFORM A TOTAL
 OF

$$1 + 2 + \dots + (n-1) = \frac{n(n-1)}{2}$$

REPRESENTATIVE UPDATES TAKING $\Theta(n^2)$ TIME.
 THE TOTAL RUN TIME OF THIS SEQUENCE
 IS THEREFORE $\Theta(n + n^2) = \Theta(n^2)$. IN TERMS
 OF $m = 2n - 1$ THIS IS $\Theta(m^2)$ TIME TO PERFORM
 A TOTAL OF m OPERATIONS. IN OTHER WORDS
 $\Theta(m)$ TIME PER OPERATION, ON AVERAGE.

WHY ARE WE INTERESTED IN THE RUN TIME
 OF SEQUENCES OF OPERATIONS? BECAUSE
 IN A DISJOINT SET DATA STRUCTURE, THE RUN
 TIME OF AN OPERATION (E.G. UNION) DEPENDS
 ON THE PARTICULAR HISTORY OF OPERATIONS
 PERFORMED.

ANALYZING THE RUN TIME OF A SEQUENCE,
 THEN AVERAGING OVER ALL OPERATIONS PERFORMED,
 IS CALLED AMORTIZED ANALYSIS. THIS IS
 OFTEN DONE IN SITUATIONS WHERE HISTORY
 MATTERS.

GOING BACK TO THE PREVIOUS EXAMPLE, WE SEE
 THE PROBLEM WITH THAT SEQUENCE IS THAT
 WE ALWAYS APPEND THE LONGER LIST TO THE
 SHORTER ONE, WHICH INCREASES THE NUMBER
 OF REPRESENTATIVE UPDATES.

WE CAN IMPROVE PERFORMANCE BY SORTING THE LENGTH OF EACH LIST, THEN ALWAYS APPEND THE SHORTER LIST TO THE LONGER, CAUSING FEWER REPRESENTATIVE UPDATES. WE CALL THIS THE WEIGHTED UNION HEURISTIC.

THEOREM

USING THE WEIGHTED UNION HEURISTIC WITH LINKED LIST REPRESENTATION OF DISJOINT SETS, ANY SEQUENCE OF m MAKESET, UNION, AND FINDSET OPERATIONS, n OF WHICH ARE MAKESETS, TAKES TIME $O(m + n \lg n)$.

PROOF:

P. 566

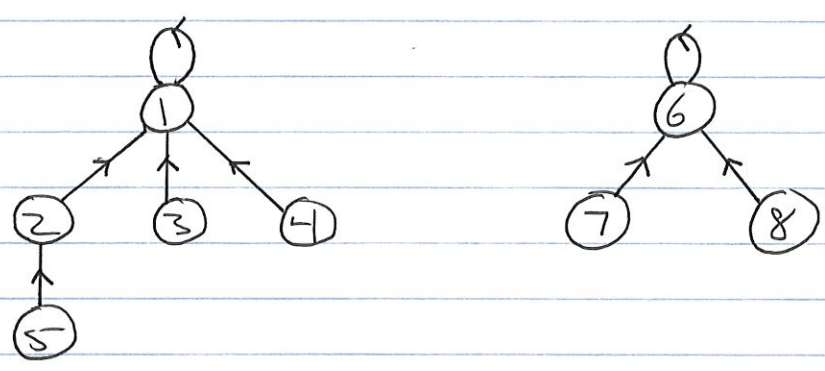
NOTE THAT THE PREVIOUS EXAMPLE TAKES TIME $\Theta(n) = \Theta(m)$ WITH WEIGHTED UNION HEURISTIC, WHICH IS SOMEWHAT BETTER THAN THAT GUARANTEED BY THE ABOVE THEOREM. AVERAGING OVER ALL m OPERATIONS IN THIS EXAMPLE GIVES $\Theta(1)$ AMORTIZED TIME.

(21.3) Disjoint Set Forests

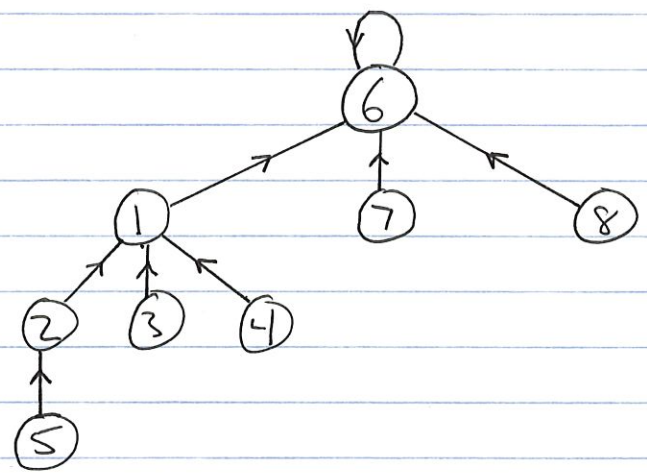
ANOTHER IMPLEMENTATION STRATEGY IS TO REPRESENT EACH SET IN \mathcal{S} BY A ROOTED TREE, i.e. \mathcal{S} IS A DISJOINT SET FOREST.

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].

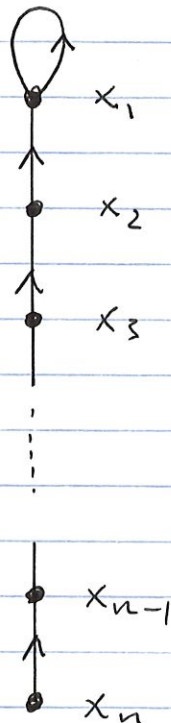
MakeSet CREATES A TREE WITH JUST ONE NODE WHICH POINTS TO ITSELF. Union CAUSES THE ROOT OF ONE TREE TO POINT TO THE ROOT OF ANOTHER. BOTH OPERATIONS HAVE COST $\Theta(1)$.

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 $\Theta(2n-1) = \Theta(n)$), 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 .

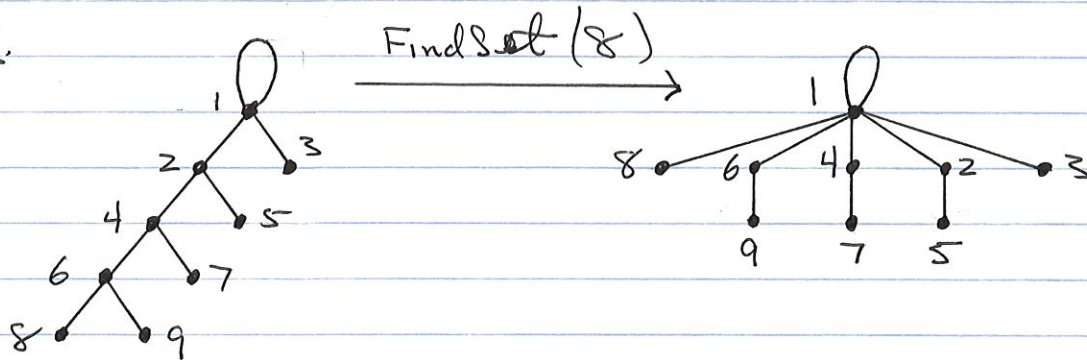
$$\text{height}(x) \leq \text{rank}[x].$$

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.) $P[x] = x$
- 2.) $rank[x] = 0$

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

- 1.) $Link(\text{FindSet}(x), \text{FindSet}(y))$

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

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

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.) $P[x] = \text{FindSet}(P[x])$
- 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 THIS 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.) $A = \emptyset$
- 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.