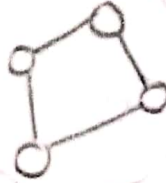


# Day 40: MST and Disjoint sets.



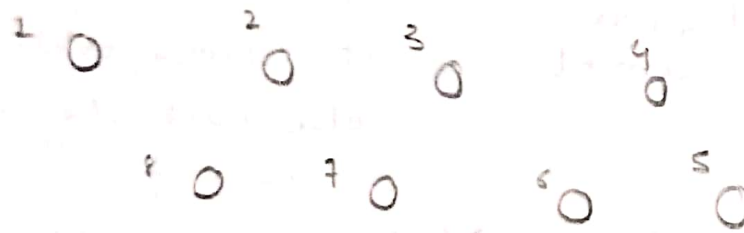
Undirected graph

Disjoint sets can be used to segregate different connected components.

Disjoint set union

can be used when edges are encountered one by one and not at once. (online)

Disjoint sets is an online data structure



- i) edge  $NA NB$
- ii) query :

one way could be to calculate the number of connected components at each query using BFS or DFS, which can be a little inefficient, where disjoint set union can help

number of connected components / clusters

Edge :  $NA - NB$

$c_1 = \text{findCluster}(NA)$

$c_2 = \text{findCluster}(NB)$

if  $(c_1 \neq c_2)$  :

union  $(c_1, c_2)$  ;

when all nodes are init, to 1.

starting at n

by just maintaining a count of clusters

Disjoint set union functions in the way that it finds the cluster of each candidate node and if they belong to different clusters, merge or join them and make them one cluster.

(in other language a connected component)

A unique identifier for each cluster can be ~~is not~~ one of the nodes which can be regarded as root of the tree. (which can act as the identity of the cluster)

↓

Each cluster with N nodes can be thought to have N-1 edges (since when we get nodes in the same cluster, we don't operate on them.)  
(tree like structure)

(int) findCluster (node A) {

if (parent[A] == -1) return A;

the root node would have no parent

return findCluster (parent[A]);

}

won't work for nodes other than the roots.

void union (int c1, int c2) {

findCluster() ← parent[c2] = c1;

return;

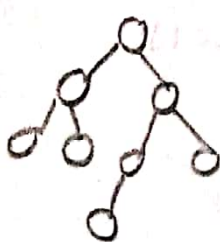
}



make c1 parent of c2

findCluster

Time complexity:  $O(n) * (n^2 \text{ max queries}) \rightarrow O(n^3)$



max possible edges in a connected graph

need to be optimized.

Such a cluster is also possible. In this case findCluster will become  $O(n)$ .

space complexity: parent array  $\rightarrow O(n)$ .

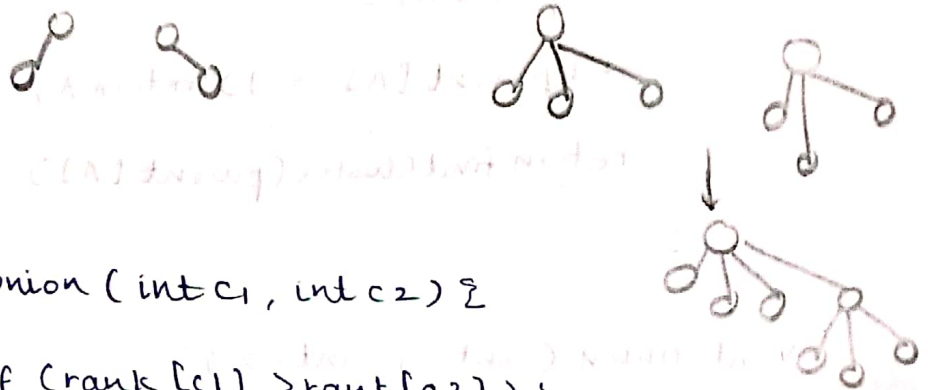
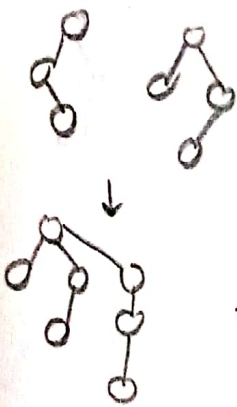
path compression: the parent for each node in the cluster can be the root of the cluster. that will reduce the complexity for find cluster to  $O(1)$ .

```

int findCluster(int A) {
    if (parent[A] == -1) return A;
    int root = findCluster(parent[A]);
    parent[A] = root;
    return root;
}

```

$O(1) \leftarrow$



```

void union(int c1, int c2) {

```

```

    if (rank[c1] > rank[c2]) {
        parent[c2] = c1;
        rank[c1] += rank[c2];
    }

```

① maintain a rank array which indicates the no. of nodes in the cluster headed by this node.

```

    else {

```

```

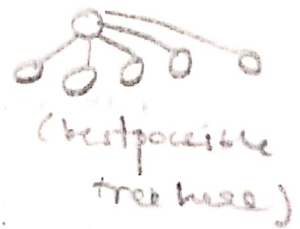
        parent[c1] = c2;

```

```

        rank[c2] += rank[c1];
    }

```



② rank array can be the measure of

height of the cluster headed by that vertex

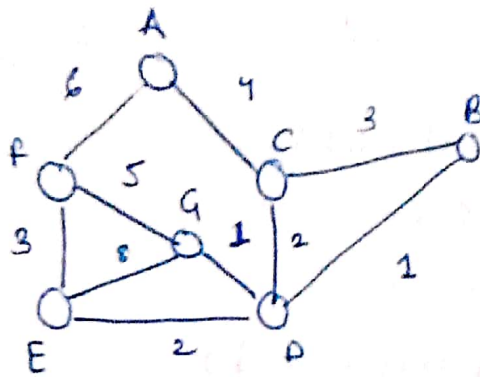
(height will increase only when the height of the cluster is same)

amortized  $O(1)$  per update on edge  $\rightarrow n^2$  operations

space:  $2n$  space  $\rightarrow O(n)$  (for 2 arrays utilized)



d.



Minimum spanning tree.

↓  
all vertices connected,  
 $n-1$  edges, with minimum  
sum weight.

i) sort edges as per weights

ii) for each edge till compiled number of edges  $= n-1$

$O(E \log E + E)$

check if both points are in  
different clusters or same. If  
different, take union. otherwise  
ignore.

(Kruskal's algorithm)

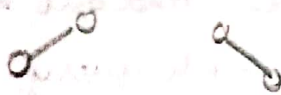
↓  
if edges are already

sorted, then  $O(E)$

since find cluster &

union are amortized  $O(1)$ .

Ex



4 → 2



$n \rightarrow$  union  $\rightarrow O(n \log n)$

$n^2 - n$

$n(n-1) \rightarrow$  no union ( $O(1)$ )

8 → 3



16 → 4



max  
height  $\Rightarrow (\log_2 n)$   
of any  
cluster

\* union based on height as rank

```
void union(int c1, int c2) {
```

```
    c1 = find cluster(c1);
```

```
    c2 = find cluster(c2);
```

```
    if (rank[c1] > rank[c2])
```

```
        parent[c2] = c1;
```

```
    else if (rank[c2] > rank[c1])
```

```
        parent[c1] = c2;
```

```
    else {
```

```
        parent[c2] = c1;
```

```
        rank[c1] += 1;
```

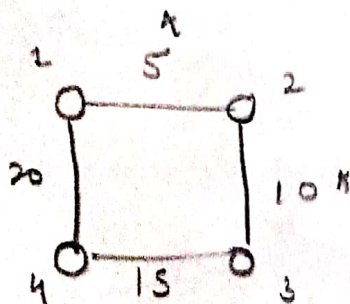
```
    }
```

```
}
```

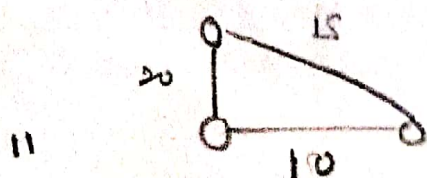
limits the  
height to log.  
for each  
cluster

Q. Given a graph (bush), may or maynot be connected.  
Given a starting node and a weight  $w$ . we can traverse  
an edge if our weight is less than the edge's  
weight. Given a bunch of queries, find out how  
many nodes can be visited for each query

starting node } for each query i/p.  
weight



$10^5$  nodes & edges  
 $10^5$  queries.



\* (queries are  
offline)

merge the query weights & edge weights in descending order.



for any query, we only want to consider the edges with weight more than query weight



→ rank can be size of the cluster

edge: union



query: cluster which start node is part of, since I know that the clusters formed are using weights more than the query weight.



return the size of the cluster of the start node