

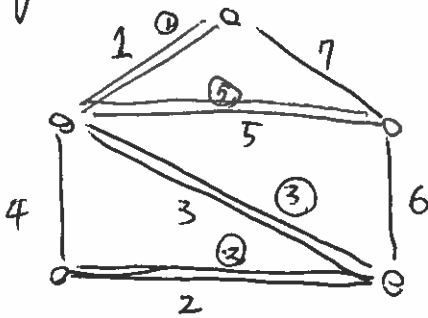
# 10/14/2018 Kruskal's MST Algorithm ①

## (Union-find)

→ Solving MST problem in another way, the another potent greedy algorithm.

How Kruskal's algorithm work?

eg.



- ① Start with the cheapest edge.
- ② Get the next cheapest edge.
- ③ Get the next cheapest edge.
- ④ Skip the next cheapest edge because it will cause the closed loop.
- ⑤ Get the next cheapest edge.
- ⑥ Skip rest of the nodes and done.

Pseudo code

- Sort edges in order of increasing cost.
- $T = \text{None}$  [Rename edges 1, 2, 3, ..., m so that  $C_1 < C_2 < \dots < C_m$ ]
- for  $i = 1$  to  $m$ 
  - if  $T \cup \{i\}$  has no cycles
  - add  $i$  to  $T$
- return  $T$

What is the way to implement the Kruskal?



What will be the straight forward running time? (2)

① Sorting part:  $O(m \log n)$

② Main loop:  $O(m)$

③ Check cycle:  $O(n)$

Check if there is another path to that node.  
[Of course DFS, BFS which is  $O(n)$  Algorithm]

all in all:  $O(m \log n) + O(m \cdot n) \Rightarrow O(m \cdot n)$

Can we do cycle check faster?

$\Rightarrow$  New data structure "Union-Find" will reduce the cycle check to  $O(1)$ -time.

then the sorting will be the bottle neck.

$O(m \log n) + O(n) \Rightarrow O(m \log n)$

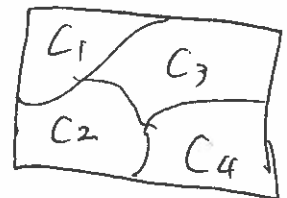
What is Union-Find data structure?

Raison d'être of a union-find:

$\Rightarrow$  Maintain the partition of a set of object.

Operation:

- Find(x): Return the name of group that x belongs to.



- Union( $C_i, C_j$ ): Fuse groups  $C_i$  and  $C_j$  into a single one.

Apply it for Kruskal

- Object = nodes

- Groups = Connected components w.r.t. currently chosen edges  $T$ .

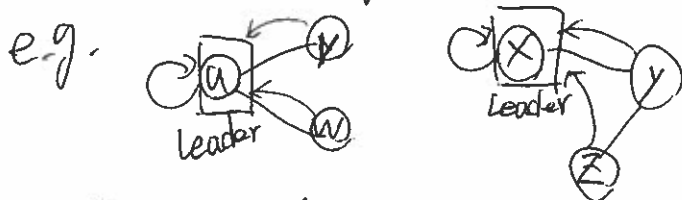
- Adding new edge  $(u, v)$  to  $T_L \Rightarrow$  Fuse connected components of  $u, v$ .

# Union-Find Basics

(3)

- Idea #1 - Maintain one linked structure per connected component of  $(V, T)$
- Each component has an arbitrary leader node.

Invariant: Each node points to the leader of its component.



[ "name" of a component inherited from leader node ]

Key point: Given edge  $(u, v)$ , can check if  $u$  and  $v$  already in same component in  $O(1)$ -time.

[ iff leader pointers of  $u, v$  match ]

i.e.  $\text{Find}(u) == \text{Find}(v)$

$\Rightarrow$  Constant time  $O(1)$

## Maintain the Invariant

Note: When new edge  $(u, v)$  added to  $T$ , connected components of  $u$  and  $v$  merge.

Idea #2: When two components merge, have smaller one inherit the leader of the larger one: (Make update minimum)

$\Rightarrow$  It's still takes  $O(n)$

$\Rightarrow$  But, the number of chances that node ~~with~~ <sup>have</sup> leader update through entire Kruskal algorithm is  $\Theta(\log n)$

Reason: Everytime node's leader pointer gets updated, population of its component at least doubles.

$\Rightarrow$  Can happen  $\leq \log_2 n$  times.

(4)

## Running time Implementing "Union-find"

Sorting:  $O(m \log n)$

Cyclecheck:  $O(m)$  ( $m \times O(1)$ )

Leader node update:  $O(n \log n)$

$\Rightarrow$  all:  $O(m \log n)$  total

$\Rightarrow$  Close to Prim's algorithm

~~Other~~

## Other way to solve MST

$\rightarrow$  Too difficult to prove.

- $O(m)$  randomized algorithm exist [Karger-Hein 1995]
- $O(m \alpha(n))$  deterministic [Chazelle 2000]

$\log^* n$  - Optimal algorithm by Pettie 2002's precised asymptotic running time is unknown. (Between  $\Theta(m)$  and  $\Theta(m \alpha(n))$ )

## Why we learn MST?

- Problem ~~itself~~ itself is interesting.
- Good way to learn data structures.
- "Clustering" is useful technique.

# How clustering is useful?

⑤

Informal Goal: Given  $n$  "points" [Webpages, images, etc.] classify them into "coherent groups".

Assumption: ① As input, given the similarity measure (a distance  $d(p, q)$  between each point pair)  
② Symmetric [i.e.  $d(p, q) = d(q, p)$ ]

Examples: Euclidean distance, genome similarity, etc.

Goal: Same cluster  $\Leftrightarrow$  "nearby"

## Max-spacing $k$ -clusterings

Assume: We know  $k = \#$  of clusters desired.

$\Rightarrow$  In practice, can experiment with a range]

Call points  $p$  and  $q$  separated if they're assigned to different clusters

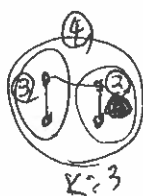
Definition: The spacing of a  $k$ -clustering is  $\min_{\text{separated } p, q} d(p, q)$ . (the bigger, the better)

Problem statement: Given a distance measuring  $d$  and  $k$ , compute the  $k$ -clustering with maximum spacing.

# A Greedy algorithm to solve k-clustering

(6)

e.g.



~~1. Con~~

- ① Set  $k=3$  (Want 3 clusters)
- ② Connect the closest points.  $K=5$
- ③ Connect the closest points.  $K=4$
- ④ Connect the closest points.  $K=3$
- ⑤  $K=3$ , so stop the loop and done.

## Pseudo code

- Initially, each point in a separate cluster.
- Repeat until only  $k$  clusters:
  - let  $p, q$  = closest pair of separated points (Determines the current spacing)
  - merge the cluster containing  $p$  and  $q$  into a single cluster.

⇒ Just like Kruskal's MST algorithm, but stopped early.

Vertices ⇒ Points

Edge cost ⇒ Distances

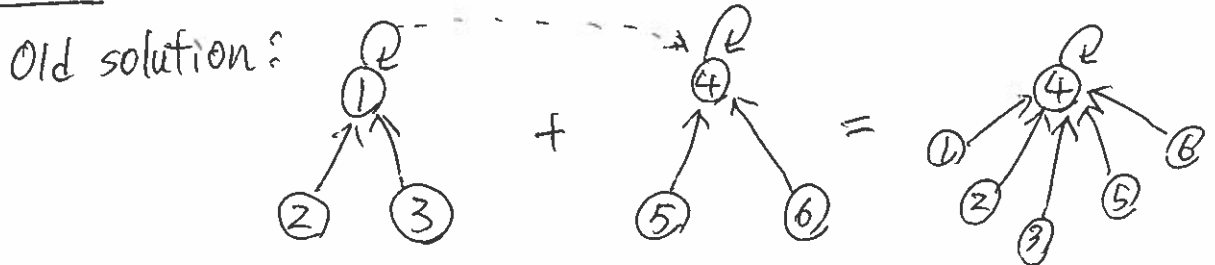
⇒ Called "Single link clustering"

# Lazy Union (Other way to update the leader)

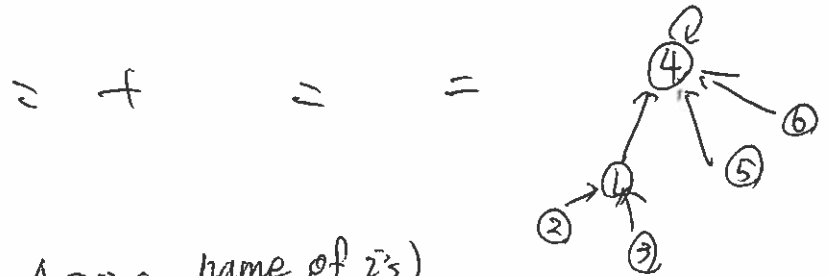
(17)

New idea: Update only one pointer each merge!

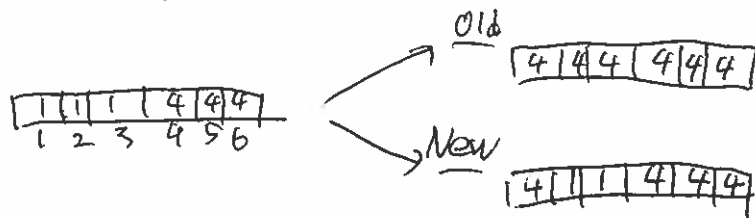
How? :



new solution:



In array rep. (Where  $A[i] \leftrightarrow$  name of  $i$ 's parent)



Pros: Reduction of update time

Cons: Recovering leader is ~~no~~ more than  $O(1)$  time.

