

## Objectives

- Wrap Up Minimum Spanning Tree
- Union-Find data structure
- Clustering

Feb 15, 2013

CSCI211 - Sprenkle

1

## Review

- What is a minimum spanning tree?
- What are three greedy solutions to finding the minimal spanning tree?

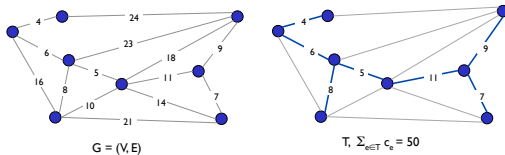
Feb 15, 2013

CSCI211 - Sprenkle

2

## Review: Minimum Spanning Tree

- Spanning tree: spans all nodes in graph
- Given a connected graph  $G = (V, E)$  with positive edge weights  $c_e$ , an **MST** is a subset of the edges  $T \subseteq E$  such that  $T$  is a **spanning tree** whose **sum of edge weights is minimized**



Feb 15, 2013

What were the three algorithms we proposed?

3

## Review: Greedy Algorithms

All three algorithms produce a MST

- **Prim's algorithm.** Start with some root node  $s$  and greedily grow a tree  $T$  from  $s$  outward. At each step, add the cheapest edge  $e$  to  $T$  that has exactly one endpoint in  $T$ .  
 > Similar to Dijkstra's (but simpler)
- **Kruskal's algorithm.** Start with  $T = \emptyset$ . Consider edges in ascending order of cost. Insert edge  $e$  in  $T$  unless doing so would create a cycle.
- **Reverse-Delete algorithm.** Start with  $T = E$ . Consider edges in descending order of cost. Delete edge  $e$  from  $T$  unless doing so would disconnect  $T$ .

What do these algorithms have/do/check in common?

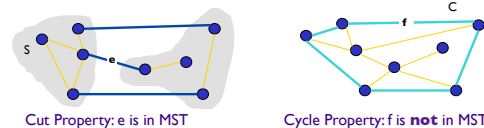
Feb 15, 2013

CSCI211 - Sprenkle

4

## Review: Important Properties

- **Simplifying assumption:** All edge costs  $c_e$  are distinct  
 → MST is unique
- **Cut property.** Let  $S$  be any subset of nodes, and let  $e$  be the min cost edge with exactly one endpoint in  $S$ . Then MST contains  $e$ .
- **Cycle property.** Let  $C$  be any cycle, and let  $f$  be the max cost edge belonging to  $C$ . Then MST does **not** contain  $f$ .



Feb 15, 2013

CSCI211 - Sprenkle

5

## Review: Prim's Algorithm

Similar to Dijkstra's algorithm. Proved optimality with the cut property

- Maintain set of explored nodes  $S$
- For each unexplored node  $v$ , maintain attachment cost  $a[v]$  → cost of cheapest edge  $v$  to a node in  $S$

$O(m \log n)$  with a heap

```

foreach ( $v \in V$ )  $a[v] = \infty$   $O(n)$ 
Initialize an empty priority queue  $Q$ 
foreach ( $v \in V$ ) insert  $v$  onto  $Q$   $O(n \log n)$ 
Initialize set of explored nodes  $S = \emptyset$ 
while ( $Q$  is not empty)  $O(n)$ 
     $u =$  delete min element from  $Q$   $O(\log n)$ 
     $S = S \cup \{u\}$ 
    foreach (edge  $e = (u, v)$  incident to  $u$ )  $O(\deg(u))$ 
        if ( $(v \notin S)$  and ( $c_e < a[v]$ ))  $O(\log n)$ 
            decrease priority  $a[v]$  to  $c_e$ 
  
```

Feb 15, 2013

6

## Kruskal's Algorithm [1956]

- Start with  $T = \phi$
- Consider edges in *ascending order of cost*
- Insert edge  $e$  in  $T$  unless doing so would create a cycle
  - Add edge as long as "compatible"

How can we prove algorithm's correctness?

Feb 15, 2013

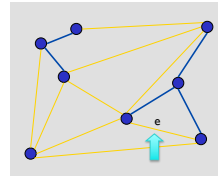
CSCI211 - Sprenkle

7

## Kruskal's Algorithm: Proof of Correctness

What is tricky about implementing Kruskal's algorithm?

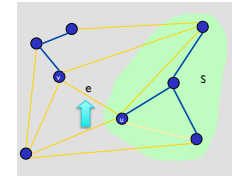
- Consider edges in ascending order of weight
- Case 1: If adding  $e$  to  $T$  creates a cycle, discard  $e$  according to *cycle property* ( $e$  must be max weight)
- Case 2: Otherwise, insert  $e = (u, v)$  into  $T$  according to *cut property* where  $S$  = set of nodes in  $u$ 's *connected component*



Feb 15, 2013

Case 1

CSCI211 - Sprenkle



Case 2

8

## Implementing Kruskal's Algorithm

What is tricky about implementing Kruskal's algorithm?

How do we know when adding an edge will create a cycle?

- What are the properties of a graph/its nodes when adding an edge will create a cycle?

Feb 15, 2013

CSCI211 - Sprenkle

9

## UNION-FIND DATA STRUCTURE

Feb 15, 2013

CSCI211 - Sprenkle

10

## Union-Find Data Structure

- Keeps track of a graph as edges are added
  - Cannot handle when edges are deleted
- Maintains disjoint sets
  - E.g., graph's connected components
- Operations:
  - **Find(u)**: returns name of set containing  $u$ 
    - How utilized to see if two nodes are in the same set?
    - Goal implementation:  $O(\log n)$
  - **Union(A, B)**: merge sets  $A$  and  $B$  into one set
    - Goal implementation:  $O(\log n)$

Feb 15, 2013

Best darn Union-Find Data Structure

11

## Implementing Kruskal's Algorithm

- Using the **union-find** data structure
  - Build set  $T$  of edges in the MST
  - Maintain set for each connected component

### Costs?

```
Sort edge weights so that  $C_1 \leq C_2 \leq \dots \leq C_m$ 
 $T = \{\}$ 
foreach ( $u \in V$ ) make a set containing singleton  $u$ 
for  $i = 1$  to  $m$ 
   $(u, v) = e_i$ 
  if ( $u$  and  $v$  are in different sets)
     $T = T \cup \{e_i\}$ 
    merge the sets containing  $u$  and  $v$ 
return  $T$ 
```

Feb 15, 2013

CSCI211 - Sprenkle

12

## Implementing Kruskal's Algorithm

- Using best implementation of **union-find**
  - Sorting:  $O(m \log n)$   $\leftarrow m \leq n^2 \Rightarrow \log m$  is  $O(\log n)$
  - Union-find:  $O(m \alpha(m, n))$
  - $O(m \log n)$  essentially a constant

```

Sort edges weights so that  $c_1 \leq c_2 \leq \dots \leq c_m$ 
 $T = \{\}$ 
foreach  $(u \in V)$  make a set containing singleton  $u$ 

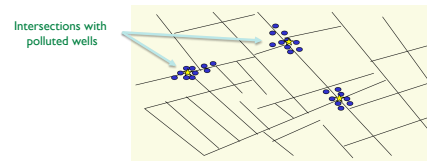
for  $i = 1$  to  $m$ 
     $(u, v) = e_i$ 
    if  $(u$  and  $v$  are in different sets)
         $T = T \cup \{e_i\}$ 
        merge the sets containing  $u$  and  $v$ 
return  $T$ 

```

Feb 15, 2013

CSCI211 - Sprenkle

13



Outbreak of cholera deaths in London in 1850s.  
Reference: Nina Mishra, HP Labs

## CLUSTERING

Feb 15, 2013

CSCI211 - Sprenkle

14

## Clustering

- Given a set  $U$  of  $n$  objects (or points) labeled  $p_1, \dots, p_n$ , classify into coherent groups
  - Problem: Divide objects into clusters so that points in different clusters are far apart
    - Requires quantification of distance
- Applications
  - Routing in mobile ad hoc networks
  - Identify patterns in gene expression
  - Identifying patterns in web application use cases
    - Sets of URLs
  - Similarity searching in medical image databases

Feb 15, 2013

CSCI211 - Sprenkle

15

## Clustering: Distance Function

- Numeric value specifying "closeness" of two objects
- Assume distance function satisfies several natural properties
  - $d(p_i, p_j) = 0$  iff  $p_i = p_j$  (identity of indiscernibles)
  - $d(p_i, p_j) \geq 0$  (nonnegativity)
  - $d(p_i, p_j) = d(p_j, p_i)$  (symmetry)

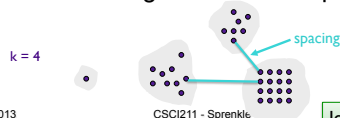
Feb 15, 2013

CSCI211 - Sprenkle

16

## Our Problem: k-Clustering of Maximum Spacing

- k-clustering**. Divide objects into  $k$  non-empty groups
- Spacing**. Min distance between any pair of points in different clusters
- k-clustering of maximum spacing**. Given an integer  $k$ , find a  $k$ -clustering of maximum spacing



Feb 15, 2013

CSCI211 - Sprenkle

Ideas about solving?

## Greedy Clustering Algorithm

- Single-link k-clustering algorithm**
  - Form a graph on the vertex set  $U$ , corresponding to  $n$  clusters
  - Find the closest pair of objects such that *each object is in a different cluster* and add an edge between them
  - Repeat  $n-k$  times until there are exactly  $k$  clusters

How is this related to the MST?

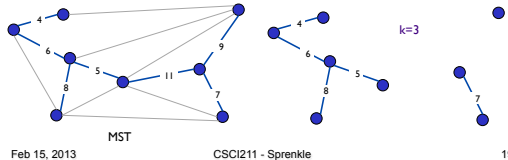
Feb 15, 2013

CSCI211 - Sprenkle

18

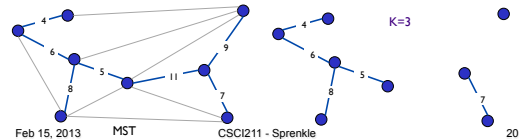
## Greedy Clustering Algorithm

- **Key observation:** Same as Kruskal's algorithm
  - Except we stop when there are  $k$  connected components
- **Remark.** Equivalent to finding MST and deleting the  $k-1$  most expensive edges



## Greedy Clustering Algorithm: Analysis

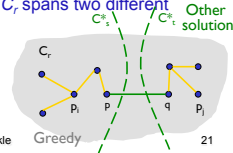
- **Theorem.** Let  $C$  denote the clustering  $C_1, \dots, C_k$  formed by deleting the  $k-1$  most expensive edges of a MST.  $C$  is a  $k$ -clustering of *max spacing*.
- **Pf Intuition:**
  - What can we say about  $C$ 's spacing?
    - Within clusters and between clusters
  - What if  $C$  isn't optimal?
    - What does that mean about  $C$ 's clusters vs (optimal)  $C^*$ 's clusters?



## Greedy Clustering Algorithm: Analysis

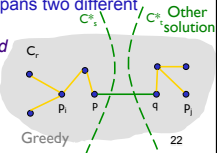
- **Theorem.** Let  $C$  denote the clustering  $C_1, \dots, C_k$  formed by deleting the  $k-1$  most expensive edges of a MST.  $C$  is a  $k$ -clustering of *maximum spacing*.
- **Pf Sketch.** Let  $C^*$  denote some other clustering  $C^*_1, \dots, C^*_k$ .  $C^*$  and  $C$  must be different; otherwise we're done.
  - The spacing of  $C$  is length  $d$  of  $(k-1)^{\text{st}}$  most expensive edge
  - Let  $p_i, p_j$  be in the same cluster in Greedy solution  $C$  (say  $C_r$ ) but different clusters in other solution  $C^*$ , say  $C^*_s$  and  $C^*_t$
  - Some edge  $(p, q)$  on  $p_i-p_j$  path in  $C$ , spans two different clusters in  $C^*$

What do we know about  $(p, q)$ ?



## Greedy Clustering Algorithm: Analysis

- **Theorem.** Let  $C$  denote the clustering  $C_1, \dots, C_k$  formed by deleting the  $k-1$  most expensive edges of a MST.  $C$  is a  $k$ -clustering of *maximum spacing*.
- **Pf.** Let  $C^*$  denote some other clustering  $C^*_1, \dots, C^*_k$ .  $C^*$  and  $C$  must be different; otherwise we're done.
  - The spacing of  $C$  is length  $d$  of  $(k-1)^{\text{st}}$  most expensive edge
  - Let  $p_i, p_j$  be in the same cluster in  $C$  (say  $C_r$ ) but different clusters in  $C^*$ , say  $C^*_s$  and  $C^*_t$
  - Some edge  $(p, q)$  on  $p_i-p_j$  path in  $C$ , spans two different clusters in  $C^*$
  - All edges on  $p_i-p_j$  path have length  $\leq d$  since Kruskal chose them
  - Spacing of  $C^*$  is at most  $\leq d$  since  $p$  and  $q$  are in different clusters



## Looking ahead

- Wiki: Chapter 4, Section 2, 4-6 (skipping section 3)
  - Due Tues midnight after break
- PS 5 due Friday after break