

ASD (2° Mod)

Lesson n. 2 on MST



Time complexity of Prim's Algorithm

THM: $O(n^2)$ with an array; $O(m \log n)$ with a binary heap.

Proof: **DO AS EXERCISE!**

SUGGESTIONS: give answers to :

- How many times a node is explored ?
- How do you represent Q?
- Which operations on Q for every new explored node? How many? How can you implement them?



Implementation. Use a priority queue ala Dijkstra.

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

```
Prim(G, c) {
  foreach (v ∈ V) a[v] ← ∞
  Initialize an empty priority queue Q
  foreach (v ∈ V) insert v onto Q
  Initialize set of explored nodes S ← ∅

  while (Q is not empty) {
    u ← delete min element from Q
    S ← S ∪ { u }
    foreach (edge e = (u, v) incident to u)
      if ((v ∉ S) and (ce < a[v]))
        decrease priority a[v] to ce
  }
```

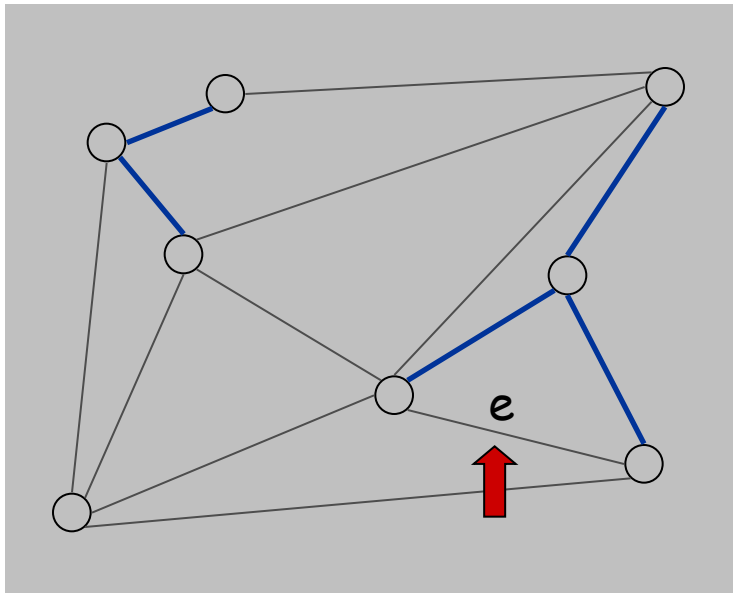
For any visited node $u \in V$, update $O(\text{deg}(u))$ keys in $Q \rightarrow \sum_u \text{deg}(u) = O(m)$
Each update costs $O(\log n)$ (using Heap) \rightarrow Total: $O(m \log n)$



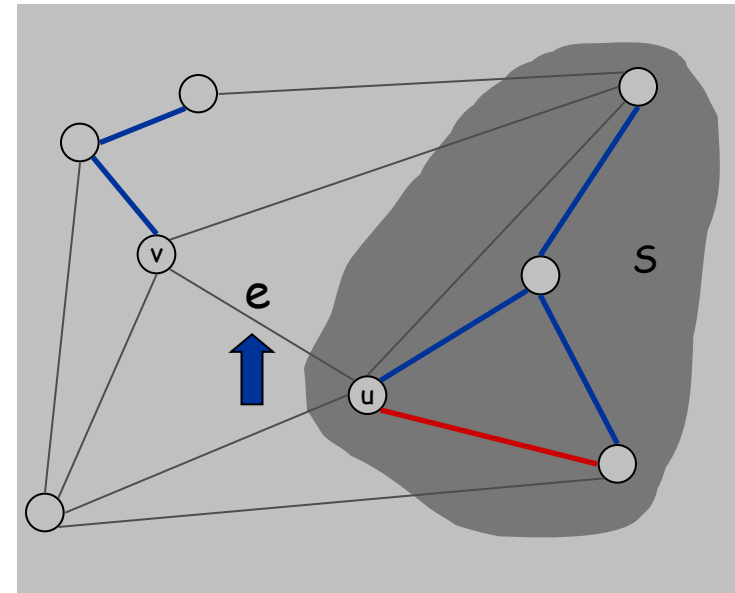
Kruskal's Algorithm: Proof of Correctness

Kruskal's algorithm. [Kruskal, 1956]

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



Case 1



Case 2



Implementation: Kruskal's Algorithm

Implementation. Use the **union-find** data structure.

- Build set T of edges in the MST.
- Maintain set for each connected component.
- $O(m \log n)$ for sorting and $O(m \alpha(m, n))$ for union-find.

$m \leq n^2 \Rightarrow \log m$ is $O(\log n)$ $\underbrace{\hspace{2cm}}$ essentially a constant

```
Kruskal(G, c) {  
  Sort edges weights so that  $c_1 \leq c_2 \leq \dots \leq c_m$ .  
   $T \leftarrow \phi$   
  
  foreach ( $u \in V$ ) make a set containing singleton  $u$   
  
  for  $i = 1$  to  $m$       are  $u$  and  $v$  in different connected components?  
    ( $u, v$ ) =  $e_i$       ↙  
    if ( $u$  and  $v$  are in different sets) {  
       $T \leftarrow T \cup \{e_i\}$   
      merge the sets containing  $u$  and  $v$   
    }  
    ↖ merge two components  
  return  $T$   
}
```



NEW DATA STRUCTURE!

WE NOW NEED A NEW **DATA STRUCTURE** TO MANAGE
SUBSET OPERATIONS EFFICIENTLY!

THIS PART FOLLOWS THE BOOK *and* IT WILL BE GIVEN IN FEW LESSONS:

*Demetrescu et Al,
Algoritmi e Strutture Dati*



ARBITRARY EDGE COSTS

To remove the assumption that all edge costs are distinct: perturb all edge costs by tiny amounts to break any ties.

Impact. Kruskal and Prim only interact with costs via pairwise comparisons. If perturbations are sufficiently small, MST with perturbed costs is MST with original costs.

↑
e.g., if all edge costs are integers,
perturbing cost of edge e_i by i / n^2

Implementation. Can handle arbitrarily small perturbations implicitly by breaking ties lexicographically, according to index.



Small perturbations of edge costs have no impacts!

Assume $c(e) \geq 1$ for all $e \in E$ and fix $b = 1/n^2$. Now, consider the new instance $I_b = \langle G(V, E), c_b: E \rightarrow \mathbb{R}^+ \rangle$ such that:

$$c_b(e) = c(e) \pm b \text{ and } c(e) \neq c(e') \text{ for any } e \neq e' \text{ (all distinct!)}$$

THM. Let T_b be any MST for I_b then T_b is also a MST for the original instance $I = \langle G(V, E), c: E \rightarrow \mathbb{R}^+ \rangle$

Proof. By contradiction. **Absurd hyp:** exists T^* better than T_b for I ,
We use the following facts:

$$\text{I) } C(T^*, I_b) < C(T^*, I) + 1/n \quad (\text{since } b=1/n^2)$$

$$\text{II) } C(T^*, I) \leq C(T_b, I) - 1 \quad (\text{absurd hyp.})$$

$$\text{III) } C(T_b, I) < C(T_b, I_b) + 1/n \quad (\text{since } b=1/n^2)$$

$$\text{(I)} \leftarrow \text{(II)} \rightarrow \text{IV) } C(T^*, I_b) < C(T_b, I) - 1 + 1/n$$

$$\text{(IV)} \leftarrow \text{(III)} \rightarrow C(T^*, I_b) < C(T_b, I_b) + 1/n - 1 + 1/n \quad (\text{for } n > 2) \\ < C(T_b, I_b) \quad (\text{absurd!})$$



EXERCISE N.1 (MST Properties)

Input: Connected Graph $G(V,E)$; $e \in E$.

Output: Decide whether an MST T exists s.t. $e \in T$

Provide an algorithm working in $O(m+n)$ time.

Hint: Combine the CUT Property and the CYCLE one to decide whether e is a MINIMAL BRIDGE.



How detect which is the case for input $G(V,E);c:E\rightarrow R^+; e=(u,w)$ in E ?

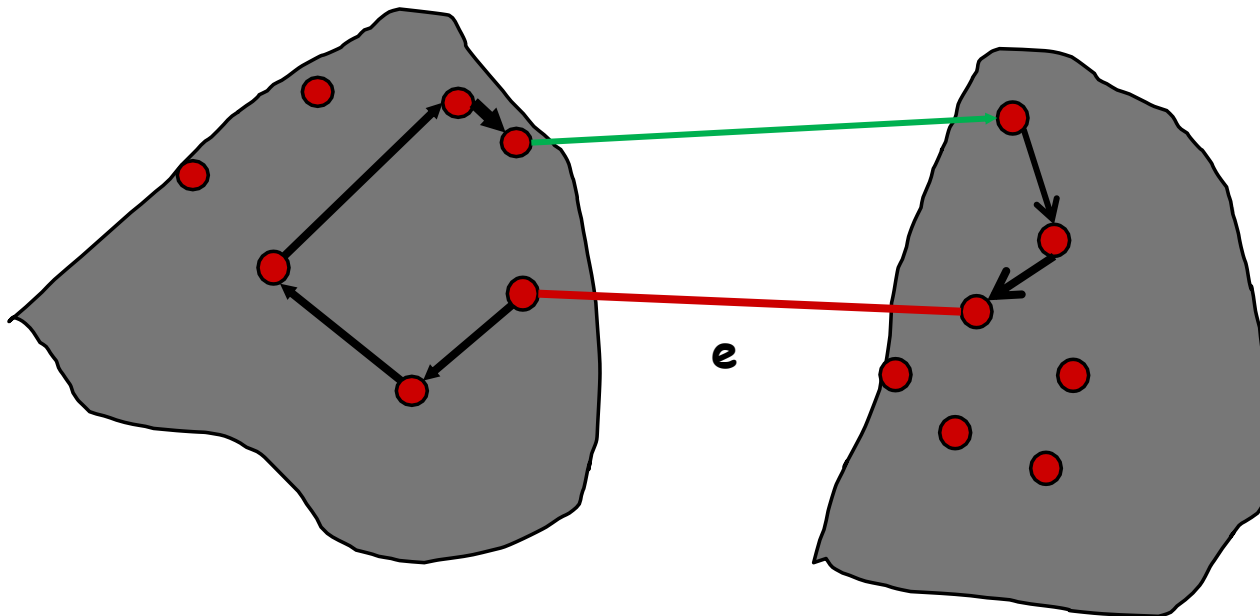
Simple Idea:

- 1) **Remove** all edges from E which are more expensive than e and remove also e . Set $G'(V,E')$ as this new graph.
- 2) **Start** a *BFS* (or a *DFS*) from u (or from w) over $G'(V,E')$
- 3) **If** the computed Tree $T(u)$ contains w then **return**: e does not belong to any MST, otherwise **return**: it does!



Let $e := (u, w)$ and $c(e) := c$

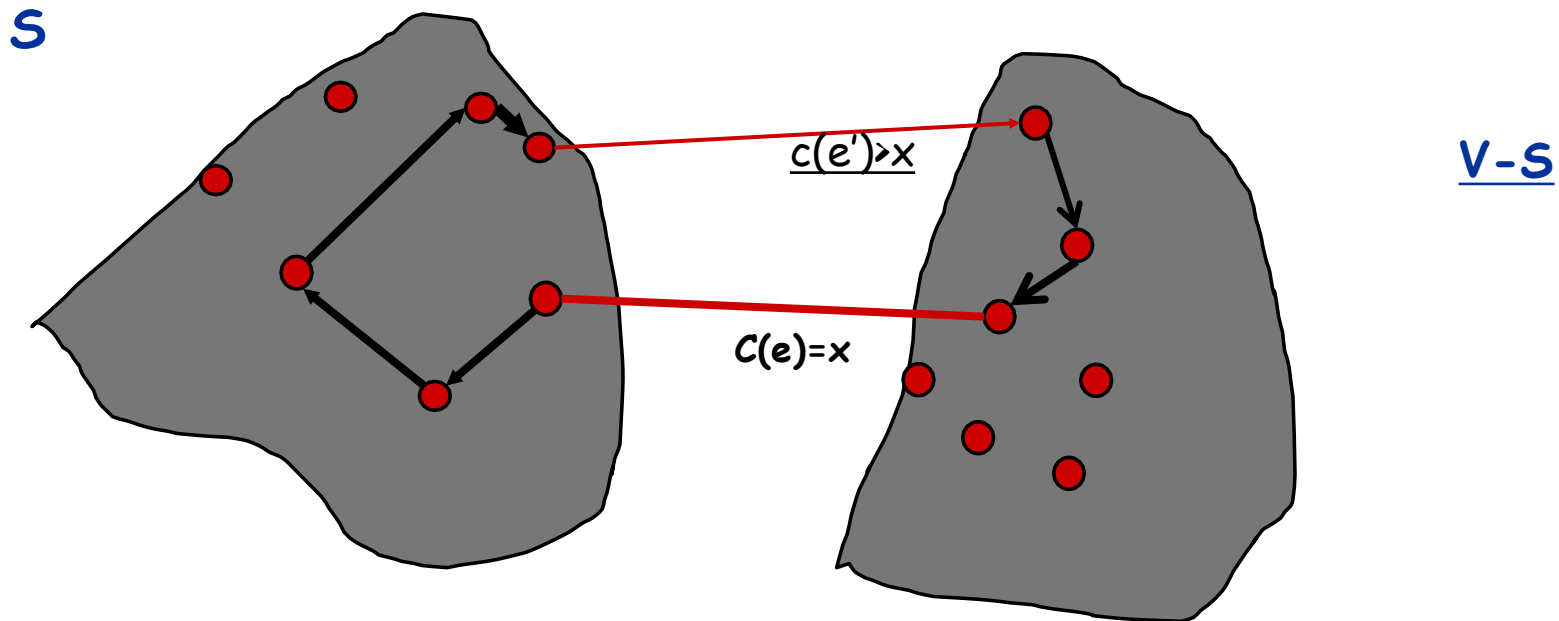
I) case: there is a Path (forming a cycle with e) where e is the most expensive!



Then, by the cycle property $\rightarrow e$ cannot belong to any MST



II case: the Set **S** of all nodes reachable from **u** with edges cheaper than **x** does not contain **w** and, so, the set **V-S** contains **w** and all nodes reachable from **w** with edges cheaper than **x**



Then, by the cut property $\rightarrow e$ belongs to any MST



Excercise n.2

Prove or Confute the following Statements:

a) Let $\langle G(V,E) w \rangle$ be s.t. G is connected and all edges have distinct weights. Let e^* be the edge of minimal weight. Does e^* always belong to an MST ?

b) Let T be an MST for $\langle G(V,E),w \rangle$ and consider the NEW instance $\langle G(V,E),w^2 \rangle$ where

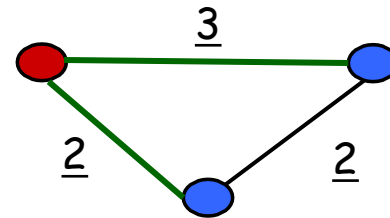
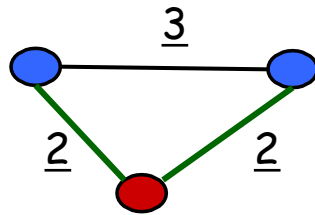
$$\text{for any } e \in E : w^2(e) = (w(e))^2$$

Is T an MST for $\langle G(V,E),w^2 \rangle$ as well?

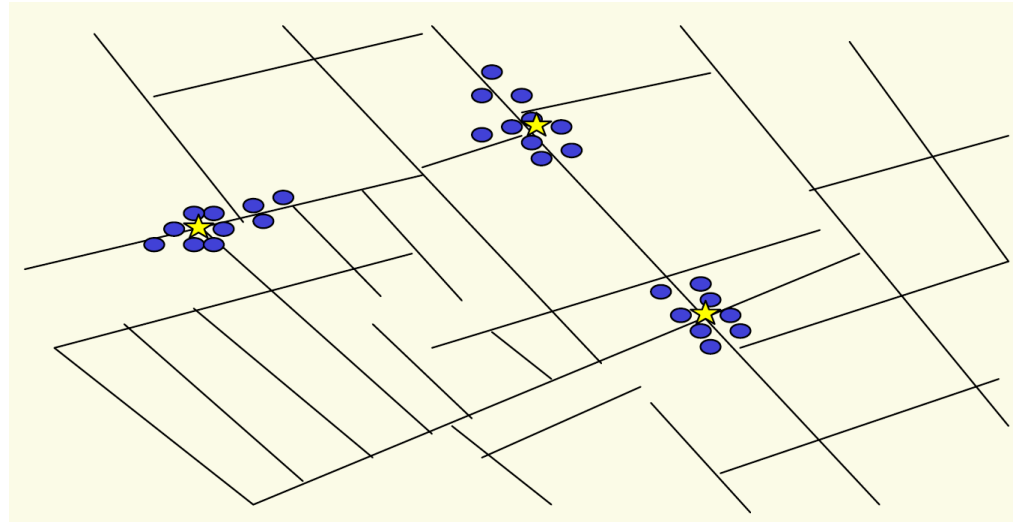


Shortest Path Tree vs Minimum Spanning Tree

- Consider an instance of the **MST** problem: $\langle G(V,E);c: E \rightarrow \mathbb{R}^+ \rangle$
- FACT 1:** The **MST** $T(s)$ computed by Prim's Algorithm may depend by the source node s . But $T(s)$ is a **feasible** and **optimal** solution for instance $\langle G(V,E);c: E \rightarrow \mathbb{R}^+ \rangle$, for **any choice of s'** in V .
- Consider an instance of the **SPT** problem: $\langle G(V,E);c: E \rightarrow \mathbb{R}^+;s \rangle$
- FACT 2:** The **SPT** $T(s)$ (and its **cost**) computed by any correct algorithm may depend on the choice of s



4.7 Clustering



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



Clustering

Clustering. Given a set U of n objects labeled p_1, \dots, p_n , classify into coherent groups.

↑
photos, documents, micro-organisms

Distance function. Numeric value specifying "closeness" of two objects:

$\text{distance}(p_i, p_j)$

↑
number of corresponding pixels whose intensities differ by some threshold

Fundamental problem. Divide into clusters so that points in different clusters **are far apart**.

- Routing in mobile ad hoc networks.
- Identify patterns in gene expression.
- Document categorization for web search.
- Similarity searching in medical image databases
- Skycat: cluster 10^9 sky objects into stars, quasars, galaxies.



Clustering of Maximum Spacing

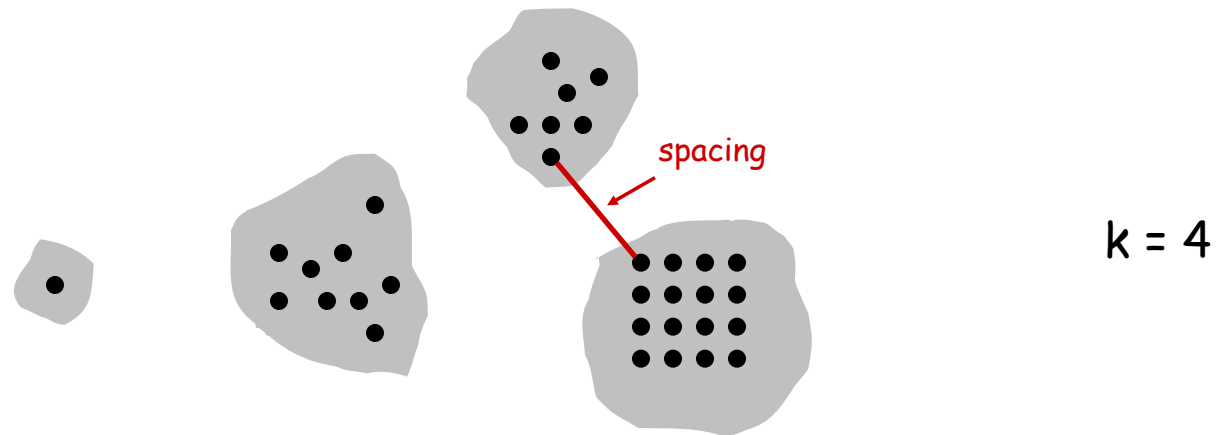
k-Clustering. Divide objects into k non-empty groups.

Distance function. Assume it 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)

Spacing. Min distance between any pair of points in different clusters.

Clustering of maximum spacing. Given an integer k , find a k -clustering of **maximum spacing**.



K-Clustering

EXERCISE: Provide formal definition of the K-Clustering Problem:

- INSTANCE: $I = \langle \dots \rangle$??????
- FEASIBLE SOLUTIONS: $Y = \dots$????
- COST OF A FEASIBLE SOLUTION: $c(Y)$
- GOAL: MINIMIZE/MAXIMIZE ???

- Determine the main parameters of the size of the Instance



PROBLEM MODEL

MAIN IDEA:

USE WEIGHTED GRAPHS TO REPRESENT THE INSTANCE !

$G(V,E)$, $d: E \rightarrow \mathbb{R}^+$, where

- $V = U = \{ p_1, \dots, p_n \}$
- $E = \{ \text{ALL non-ordered pairs in } V \}$ (i.e. complete graph)
- $\text{Cost}(e) = d(p_i, p_j)$ for any $e = \{ p_i, p_j \}$



Greedy Clustering Algorithm

Single-link k -Clustering algorithm.

- Form a graph on the vertex set U , corresponding to n clusters.
- Find the **closest** pair (edge) of objects (p,p') such that p & p' are not in the same cluster, and add **an edge** between them: so merging 2 clusters.
- Repeat $n-k$ times until there are exactly k clusters.

Key Obs. 1. This procedure is precisely Kruskal's algorithm (except we stop when there are k connected components).

Key Obs. 2. Equivalent to finding an MST T and deleting the $k-1$ most expensive edges from T (thus forming k connected components).

Proofs of 1 and 2: Exercises



Greedy Clustering Algorithm: Analysis

Theorem. Let \mathcal{C}^* denote the clustering $\mathcal{C}^*_1, \dots, \mathcal{C}^*_k$ formed by deleting the $k-1$ most expensive edges of an **MST**. Then, \mathcal{C}^* is a k -clustering of maximal spacing.

Pf. Let \mathcal{C} denote some other clustering $\mathcal{C}_1, \dots, \mathcal{C}_k$.

- The spacing of \mathcal{C}^* is the length d^* of the $(k-1)^{\text{st}}$ most expensive edge.
- Let p, p' be in the same cluster in \mathcal{C}^* , say \mathcal{C}^*_r , but different clusters in \mathcal{C} , say \mathcal{C}_s and \mathcal{C}_t .
- Some edge (q, q') on $p \rightarrow p'$ path in \mathcal{C}^*_r spans two diff. clusters in \mathcal{C} .
- All edges on $p \rightarrow p'$ path have length $\leq d^*$ since Kruskal chooses them.
- Spacing of \mathcal{C} is $\leq d^*$ since q and q' are in different clusters of \mathcal{C} .

