# Chapter 4

## Greedy Algorithms
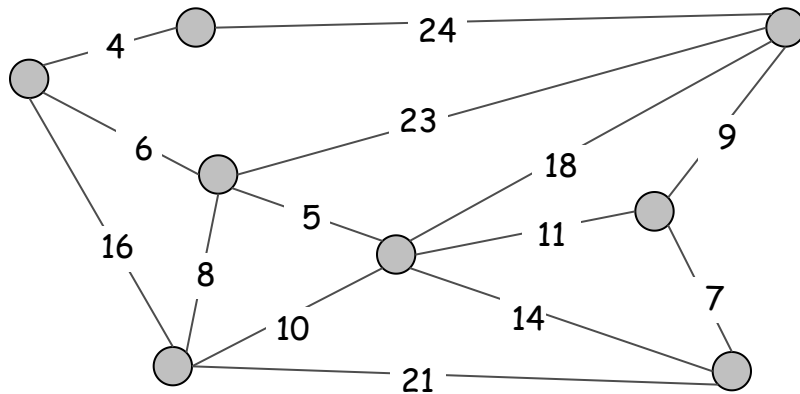
Algorithm Design

**JON KLEINBERG · ÉVA TARDOS**

# 4.5 Minimum Spanning Tree

The minimum spanning tree problem has a long history –
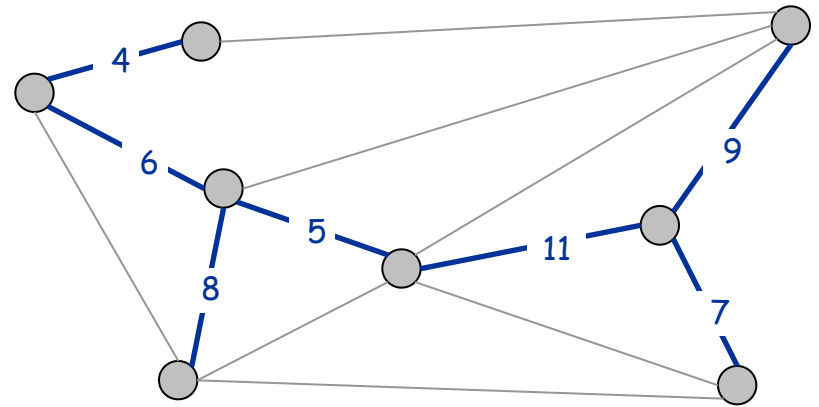the first algorithm dates back to 1926!

# Minimum Spanning Tree

Minimum spanning tree. Given a connected graph $G = (V, E)$ with real-valued 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.

if *(V, T)* is a tree



$G = (V, E)$

$T, \ \Sigma_{e \in T} \ c_e = 50$

# Applications

MST is fundamental problem with diverse applications.

- Network design.
  - telephone, electrical, hydraulic, TV cable, computer, road

- Approximation algorithms for NP-hard problems.
  - traveling salesperson problem, Steiner tree

- Indirect applications.
  - max bottleneck paths
  - LDPC codes for error correction
  - image registration with Renyi entropy
  - learning salient features for real-time face verification
  - reducing data storage in sequencing amino acids in a protein
  - model locality of particle interactions in turbulent fluid flows
  - autoconfig protocol for Ethernet bridging to avoid cycles in a network

- Cluster analysis.

# Greedy Algorithms

**Kruskal's algorithm.**  Start with T = $\phi$. 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.

**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.
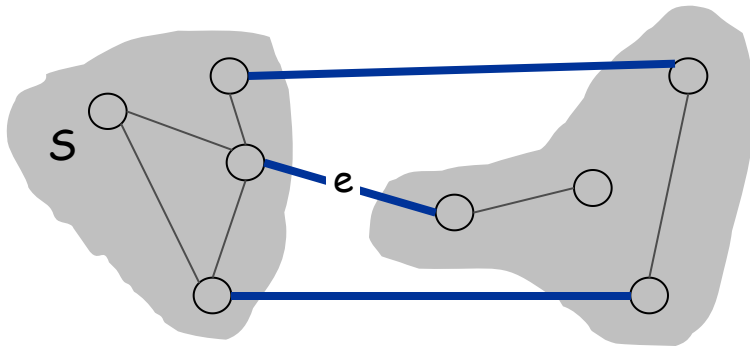
**Remark.**  All three algorithms produce an MST.
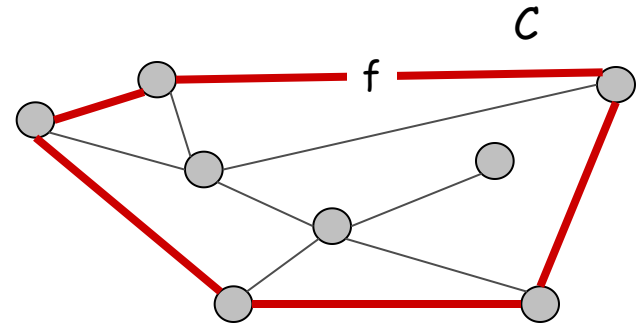
# Greedy Algorithms

Simplifying assumption.  All edge costs $c_e$ are distinct.

Cut property.  Let S be any subset of nodes, and let e be the min cost edge with exactly one endpoint in S.  Then the MST contains e.

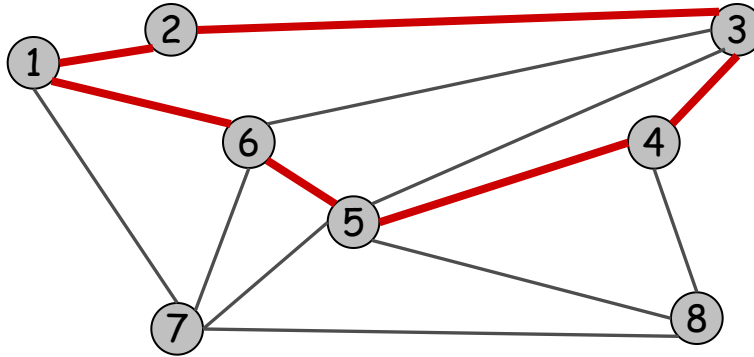Cycle property.  Let C be any cycle, and let f be the max cost edge belonging to C.  Then the MST does not contain f.
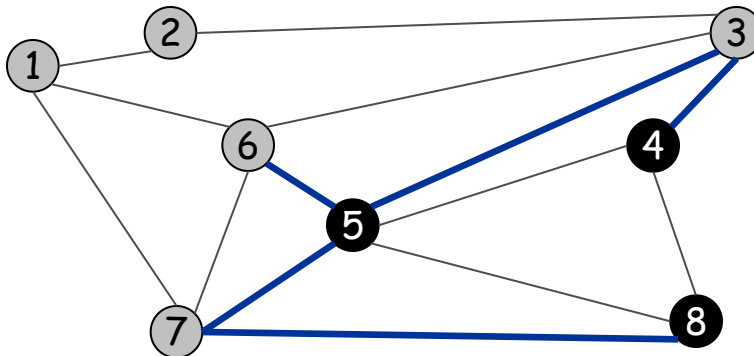
e is in the MST

f is not in the MST

# Cycles and Cuts

Cycle.  Set of edges that form a-b, b-c, c-d, …, y-z, z-a.



Cycle C  =  1-2, 2-3, 3-4, 4-5, 5-6, 6-1
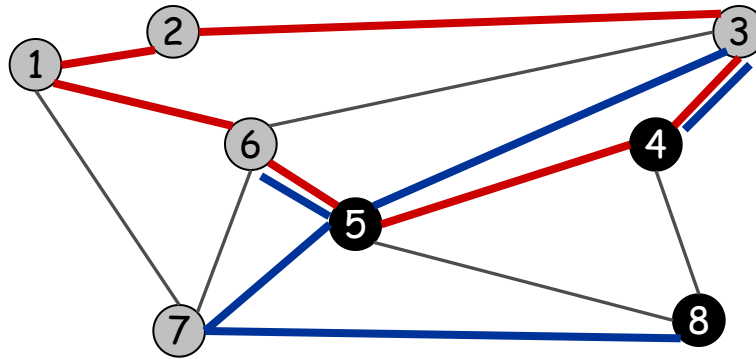
Cutset.  A cut S is a subset of nodes.  The corresponding cutset D is the subset of edges with exactly one endpoint in S.



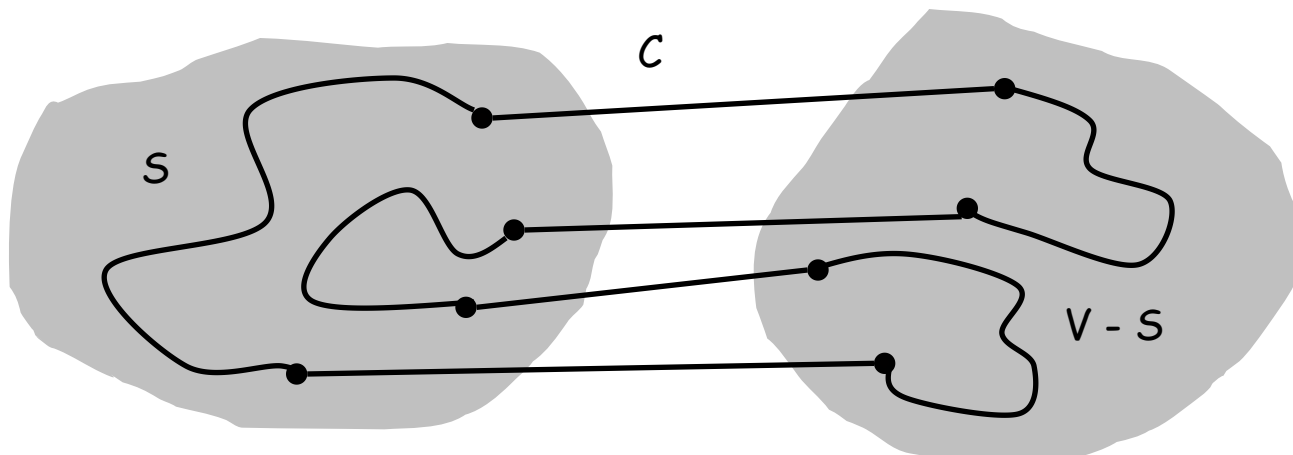Cut S      = { 4, 5, 8 }
Cutset  D = 5-6, 5-7, 3-4, 3-5, 7-8

# Cycle-Cut Intersection

Claim. A cycle and a cutset intersect in an even number of edges.



Cycle  C = 1-2, 2-3, 3-4, 4-5, 5-6, 6-1
Cutset D = 3-4, 3-5, 5-6, 5-7, 7-8
Intersection = 3-4, 5-6
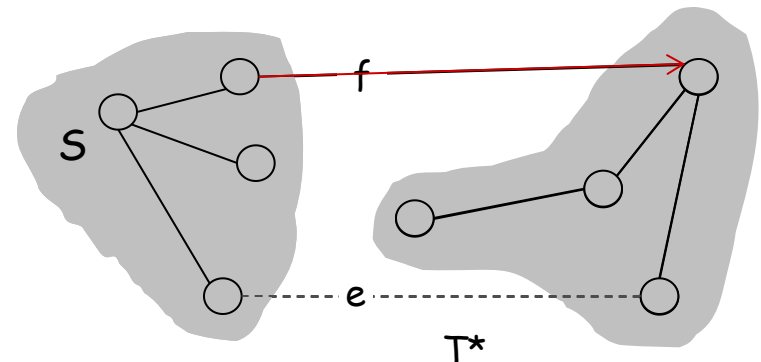Cut S = {4, 5, 8}

Pf. (by picture)

# Greedy Algorithms

**Simplifying assumption.** All edge costs $c_e$ are distinct.

**Cut property.** Let S be any subset of nodes, and let e be the min cost edge with exactly one endpoint in S. Then the MST T* contains e.

Pf. (exchange argument)
- Suppose e does not belong to T*, and let's see what happens.
- Adding e to T* creates a cycle C in T*.
- Edge e is both in the cycle C and in the cutset D corresponding to S
  $\Rightarrow$ there exists another edge, say f, that is in both C and D.
- T' = T* $\cup$ { e } - { f } is also a spanning tree.
- Since $c_e < c_f$, cost(T') < cost(T*).
- This is a contradiction. ∎

Previous slide: an even number of edges



S

f
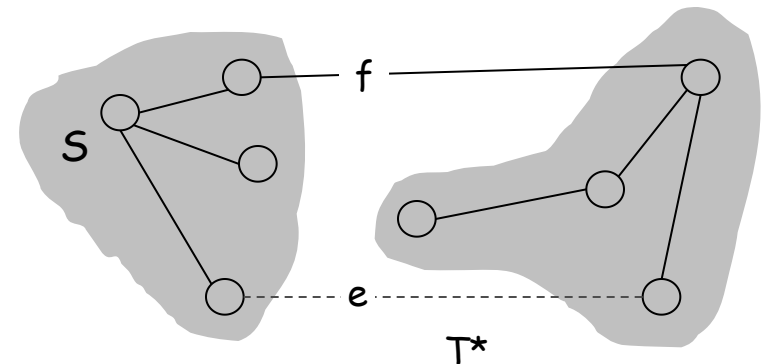
e

T*

# Greedy Algorithms

Simplifying assumption. All edge costs $c_e$ are distinct.

Cycle property. Let C be any cycle in G, and let f be the max cost edge belonging to C. Then the MST T* does not contain f.

Pf. (exchange argument)
- Suppose f belongs to T*, and let's see what happens.
- Deleting f from T* creates a cut S in T*.
- Edge f is both in the cycle C and in the cutset D corresponding to S
  $\Rightarrow$ there exists another edge, say e, that is in both C and D.
- T' = T* $\cup$ { e } - { f } is also a spanning tree.
- Since $c_e$ < $c_f$, cost(T') < cost(T*).
- This is a contradiction. ■

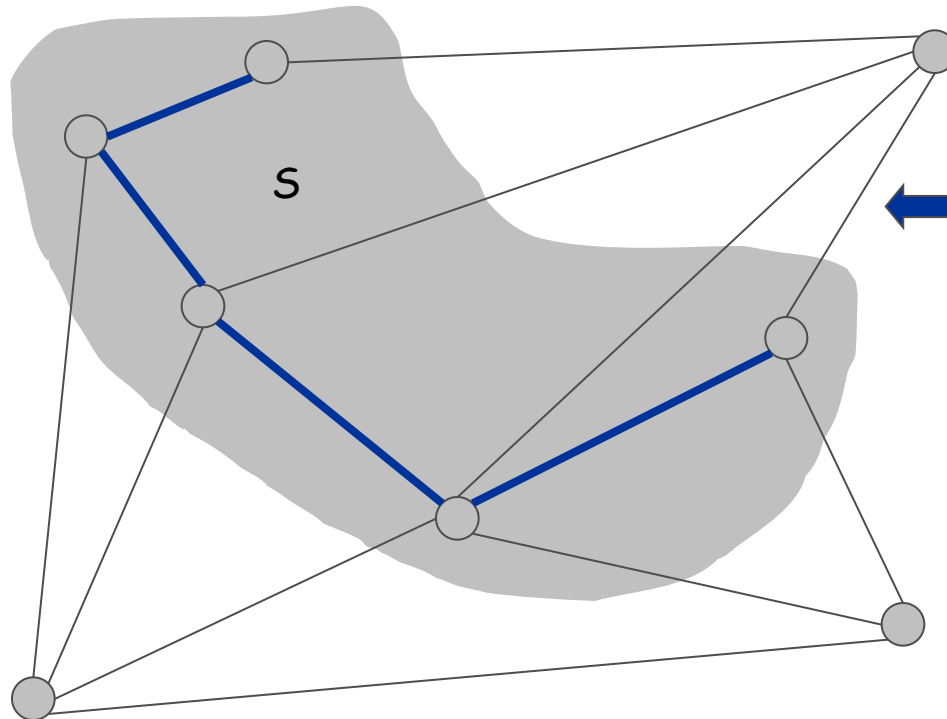Previous 2 slides: an even number of edges

S

f

e

T*

Prim's algorithm.  [Jarník 1930, Dijkstra 1957, Prim 1959]

- Initialize S = any node.
- Apply cut property to S.
- Add min cost edge in cutset corresponding to S to T, and add one new explored node u to S.

# Implementation: Prim's Algorithm

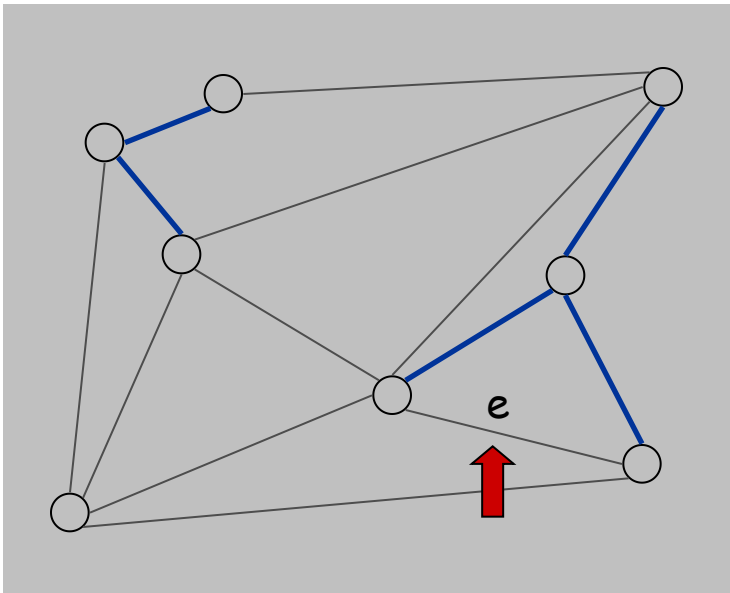Implementation. Use a priority queue ala (just like) Dijkstra.
- Maintain set of explored nodes S.
- For each unexplored node v, maintain attachment cost a[v] = cost of cheapest edge v to a node in S.
- $O(n^2)$ with an array; $O(m \log n)$ with a binary heap.

```
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 ← φ
    Pick any node s and a[s] ← 0

    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 (c_e < a[v]))
                decrease priority a[v] to c_e
    }
}
```
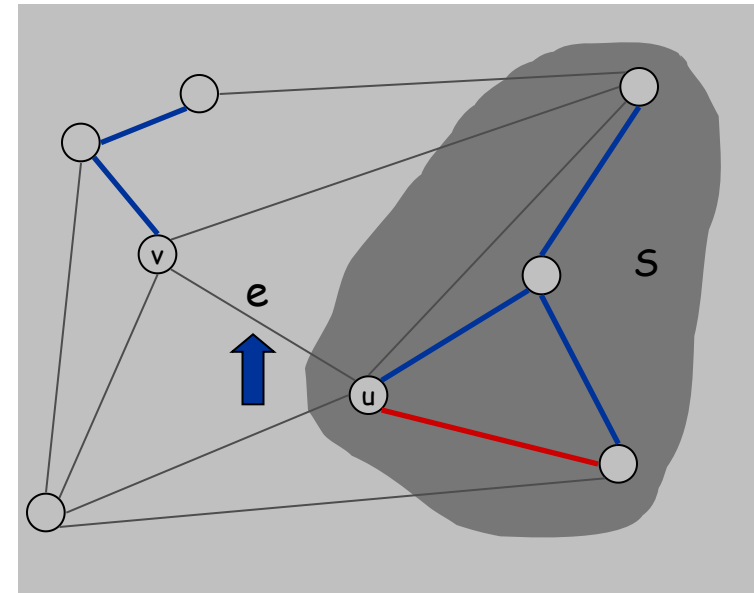
# Kruskal's Algorithm:  Proof of Correctness

**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.
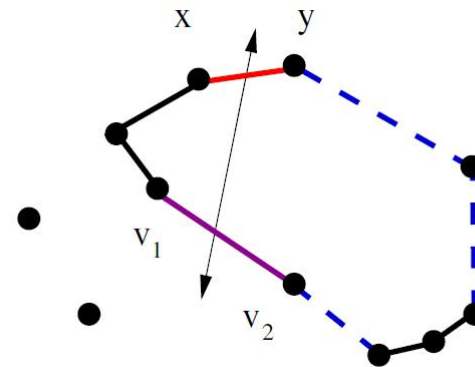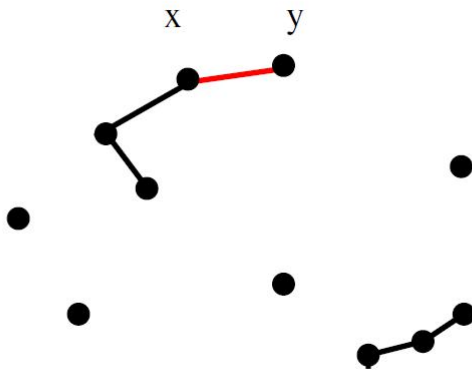


Case 1

Case 2

# Kruskal's Algorithm: Proof of Correctness

Pf. (Proof by Contradiction)

✓ If Kruskal's algorithm is not correct, these must be some graph G where it does not give the minimum cost spanning tree.

✓ If so, there must be a first edge (x; y) Kruskal adds such that the set of edges cannot be extended into a minimum spanning tree.



✓ When we added (x; y) there is no path between x and y, or it would have created a cycle. Thus adding (x; y) to the optimal tree must create a cycle.

✓ But at least one edge in this cycle must have been added after (x; y), so it must be heavier because when Kruskal added (x,y) it did not created a cycle.

✓ Deleting this heavier edge leaves a better MST than the optimal tree, yielding a contradiction!
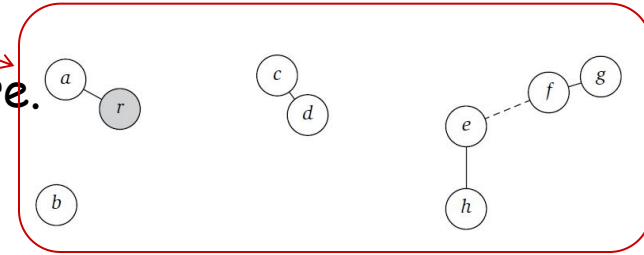
# Implementation: Kruskal's Algorithm

Union(Find(u),Find(v))

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.

$m \le n^2 \Rightarrow \log m$ is $O(\log n)$        essentially a constant

```
Kruskal(G, c) {
    Sort edges weights so that c₁ ≤ c₂ ≤ ... ≤ cₘ.
    T ← φ

    foreach (u ∈ V) make a set containing singleton u

    for i = 1 to m        are u and v in different connected components?
        (u,v) = eᵢ
        if (u and v are in different sets) {
            T ← T ∪ {eᵢ}
            merge the sets containing u and v
        }
    return T
}
```

merge two components

# Lexicographic Tiebreaking

To remove the assumption that all edge costs are distinct:  perturb all edge costs by tiny amounts to break any ties.
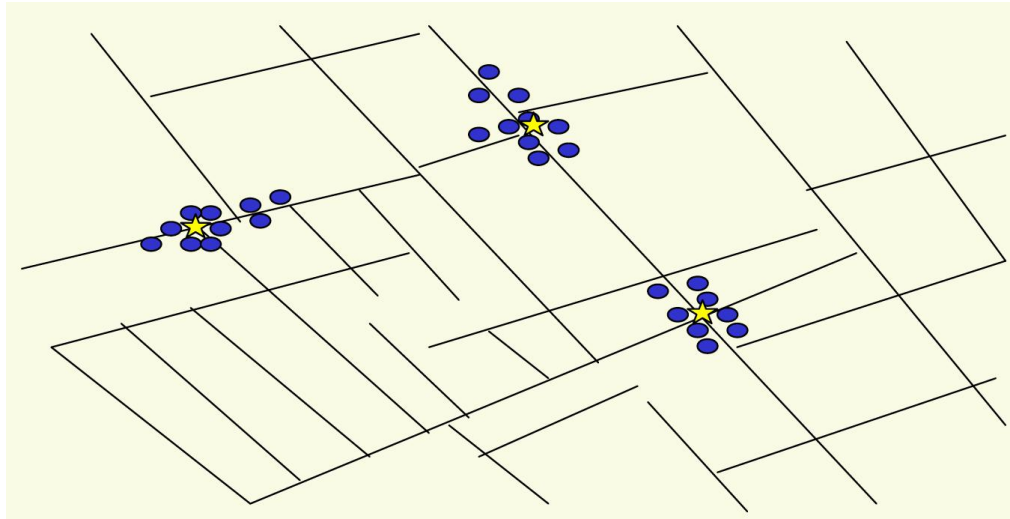
Impact.  Kruskal and Prim only interact with costs via pairwise comparisons.  If perturbations are sufficiently small, MST with perturbed costs is MST with original costs.

↑
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.

```
boolean less(i, j) {
    if       (cost(e_i) < cost(e_j))  return true
    else if  (cost(e_i) > cost(e_j))  return false
    else if  (i < j)                  return true
    else                              return false
}
```

# 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 $p_1, \ldots, p_n$, classify into coherent groups.

↑
photos, documents. micro-organisms

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

↑
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.
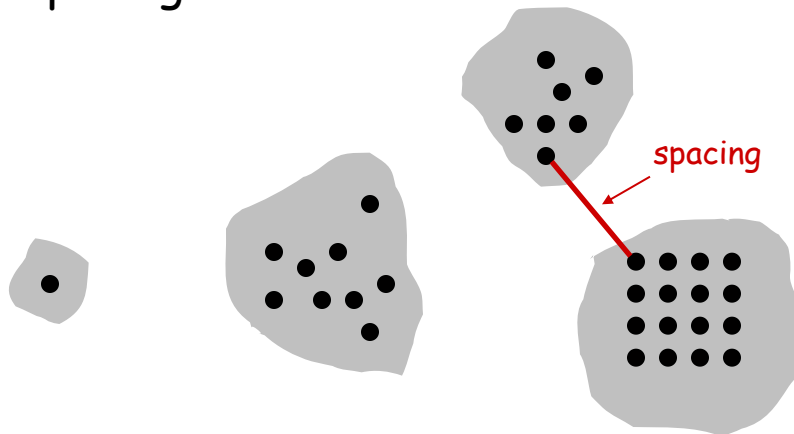
# Clustering of Maximum Spacing

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.



spacing

k = 4

# 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 to merge these two clusters into a new cluster (one cluster less).
- Repeat n-k times until there are exactly k clusters.

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

**Remark.** Equivalent to finding an MST and deleting the k-1 most expensive edges.

**Theorem.** Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the k-1 most expensive edges of a MST. $C^*$ is a k-clustering of max spacing.

**Pf.** Let C denote some other clustering $C_1, ..., C_k$.

- The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ 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 $C_t$ (which must be possible, otherwise same)
- Some edge $(p, q)$ on $p_i$-$p_j$ path in $C^*_r$ spans two different clusters in C.
- All edges on $p_i$-$p_j$ path have length $\leq d^*$
  since Kruskal chose them before the d*.
- Spacing of C is $\leq d^*$ since p and q
  are in different clusters.
- C is not max spacing
  → $C^*$ is a k-clustering of max spacing