# 4.7  Clustering [focusing on single linkage]

# Clustering

**Clustering.** Given a set U of n points labeled $p_1$, ..., $p_n$, 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
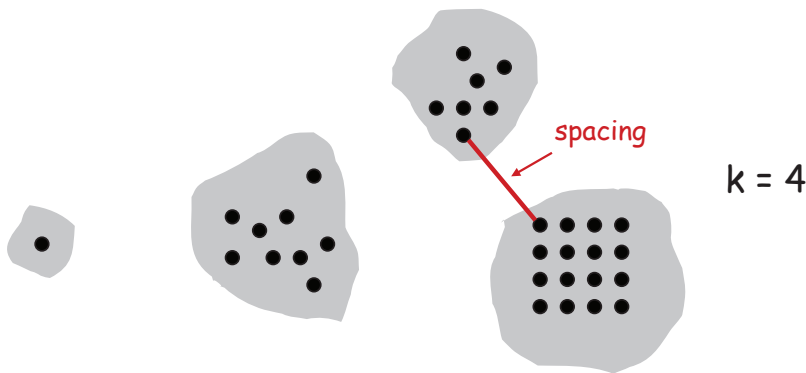
Goal: Divide objects into k non-empty groups (aka clusters): k-clustering

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-linkage 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.
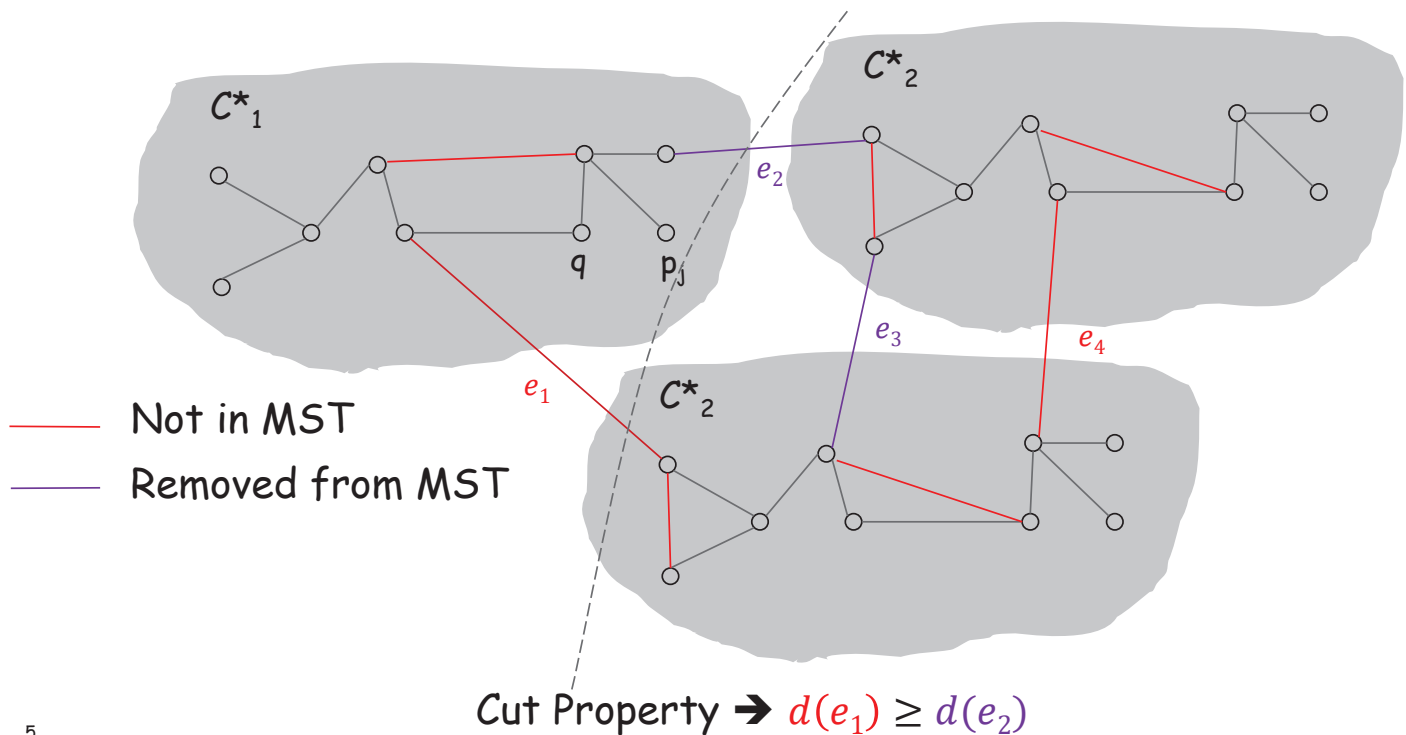Repeat n-k times until there are exactly k clusters.

**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.

# Greedy Clustering Algorithm: Analysis

**Theorem.** Let C* denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the **k-1 most expensive edges** of a MST. Then C* is a k-clustering of max spacing.
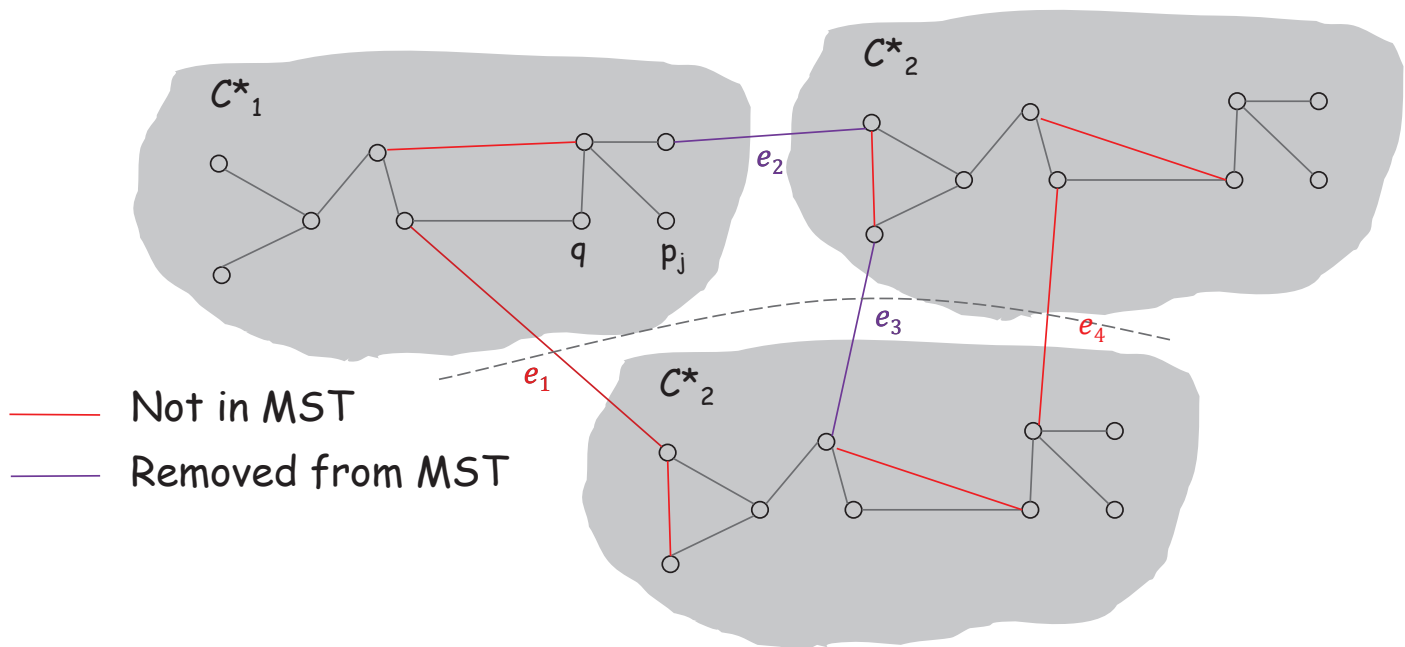
**Observation:** The spacing of C* is the length d* of the $(k-1)^{st}$ most expensive edge (in MST).



$C^*_1$

$C^*_2$

$e_2$

q

p_j

$e_3$

$e_4$

$e_1$

$C^*_2$

—— Not in MST

—— Removed from MST

Cut Property ➡ $d(e_1) \geq d(e_2)$

# Greedy Clustering Algorithm: Analysis

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

**Observation:** The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ most expensive edge (in MST).
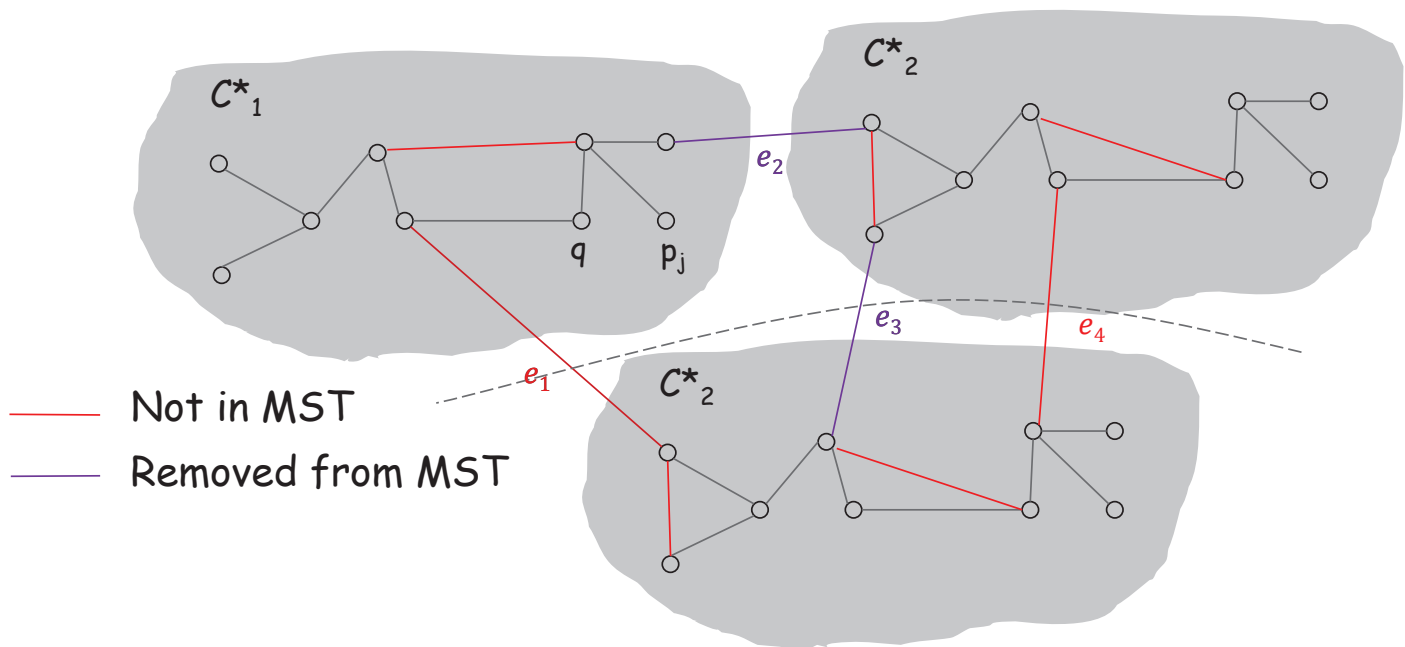


$C^*_1$

$C^*_2$

$e_2$

$q$    $p_j$

$e_3$    $e_4$

$e_1$

$C^*_2$

—— Not in MST

—— Removed from MST

Cut Property ➔ $d(e_1), d(e_4) \geq d(e_2)$

Greedy Clustering Algorithm: Analysis

**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.

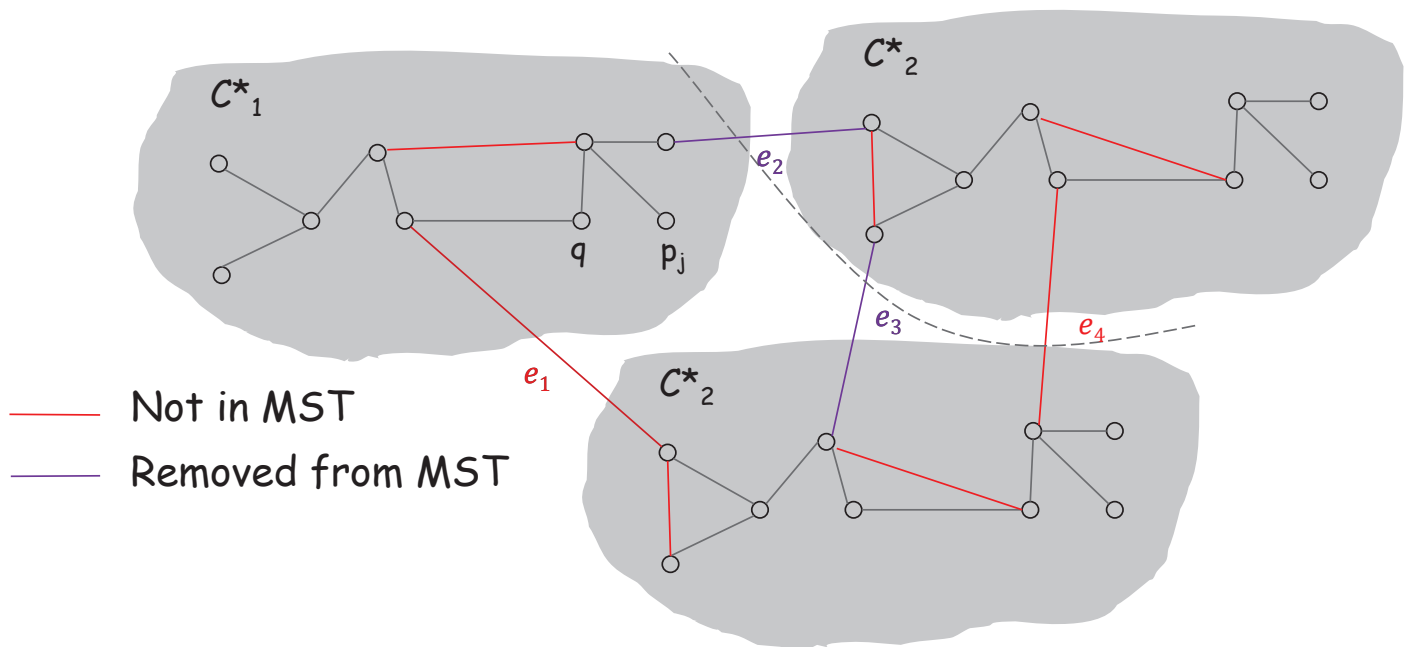Observation: The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ most expensive edge (in MST).



Cut Property ➜ $d(e_1) \geq d(e_3) \rightarrow d(e_1) \geq \max\{d(e_2), d(e_3)\}$

# Greedy Clustering Algorithm:  Analysis

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

**Observation:** The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ most expensive edge (in MST).



Cut Property ➔ $d(e_4) \geq \max\{d(e_2), d(e_3)\}$

Greedy Clustering Algorithm: Analysis

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

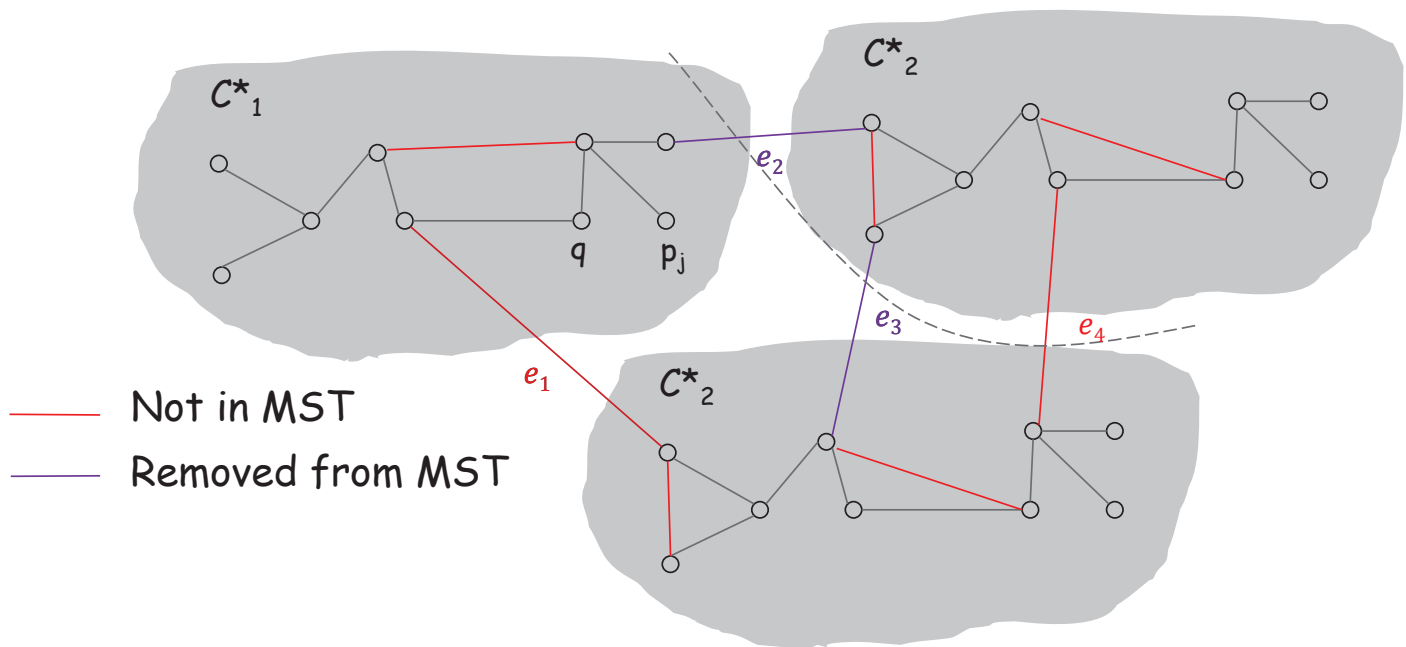**Observation:** The spacing of $C*$ is the length $d*$ of the $(k-1)^{st}$ most expensive edge (in MST).

$C*_1$

$C*_2$

$e_2$

$q$

$p_j$

$e_3$

$e_4$

$e_1$

$C*_2$

—— Not in MST

—— Removed from MST

Spacing = $\max\{d(e_2), d(e_3)\}$

9

Greedy Clustering Algorithm:  Analysis

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

Proof?

Greedy Clustering Algorithm:  Analysis

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

Pf.  Let C denote some other clustering $C_1, ..., C_k$.
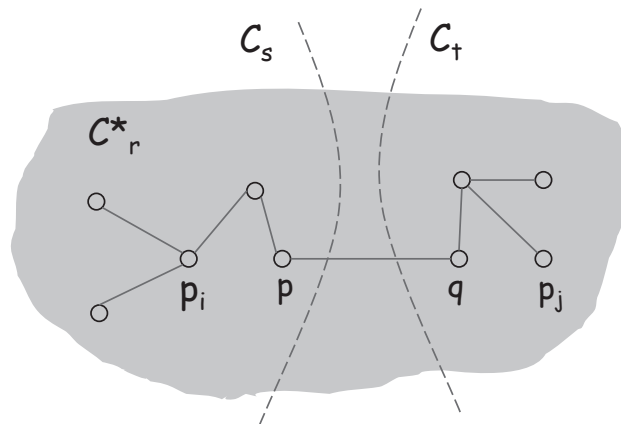The spacing of C* is the length d* of the $(k-1)^{st}$ most expensive edge (in MST).

# Greedy Clustering Algorithm:  Analysis

**Theorem.** Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the $k-1$ most expensive edges of a MST. Then $C^*$ 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 (in MST).
Let $p_i$, $p_j$ be in the same cluster in $C^*$, say $C^*_r$, but different clusters in $C$, say $C_s$ and $C_t$.
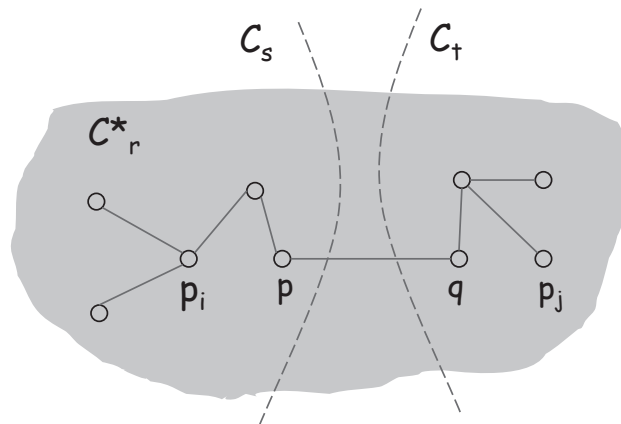
## Greedy Clustering Algorithm: Analysis

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

**Pf.** Let $C$ denote some other clustering $C_1, \ldots, C_k$.
  The spacing of $C^*$ is the length $d^*$ of the $(k-1)^{st}$ most expensive edge (in MST).
  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 $C^*_r$ spans two different clusters in $C$.
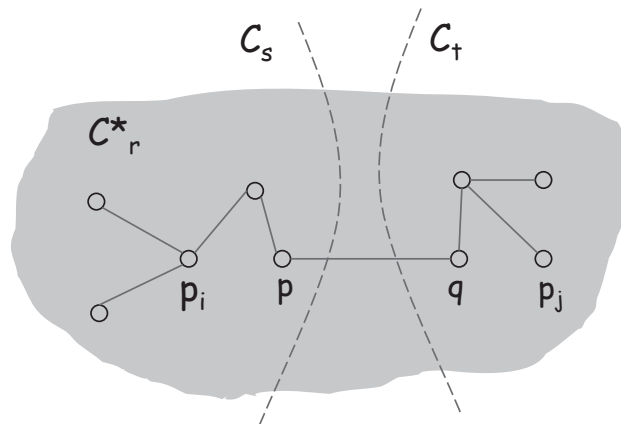
## Greedy Clustering Algorithm: Analysis

**Theorem.** Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the k-1 most expensive edges of a MST. Then $C^*$ 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 (in MST).

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 $C^*_r$ spans two different clusters in $C$.

All edges on $p_i$-$p_j$ path have length $\leq d^*$ since Kruskal chose them (and we did not delete them from our MST).

# Greedy Clustering Algorithm:  Analysis

**Theorem.** Let $C^*$ denote the clustering $C^*_1, ..., C^*_k$ formed by deleting the k-1 most expensive edges of a MST. Then $C^*$ 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 (in MST).

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 $C^*_r$ spans two different clusters in $C$.

All edges on $p_i$-$p_j$ path have length $\leq d^*$ since Kruskal chose them (and we did not delete them from our MST).

Spacing of $C$ is $\leq d^*$ since $p$ and $q$ are in different clusters. ∎