# Minimum Spanning Trees

## Application to Clustering

Algorithms: Design and Analysis, Part II

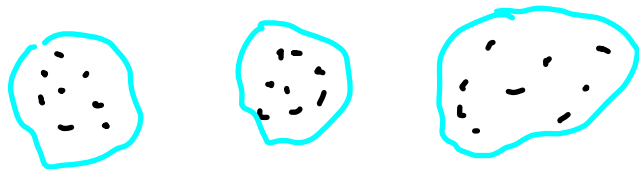# Clustering

Informal goal!: given n "points" [Web pages, images, genome fragments, etc.] classify into "coherent groups".

Assumptions: ① As input, given a (dis)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 $\Longleftrightarrow$ "nearby"

Tim Roughgarden

# Max-Spacing k-Clusterings

<span style="color:red">**Assume**</span>: we know $k := \#$ of clusters desired.

<span style="color:blue">[in practice, can experiment with a range of values]</span>

Call points $p \, \xi \, q$ <span style="color:green">separated</span> if they're assigned to different clusters.

<span style="color:red">**Definition**</span>: the <span style="color:green">spacing</span> of a k-clustering is

$$\min_{\text{separated } p, q} d(p, q).$$  <span style="color:cyan">(the bigger, the better]</span>

<span style="color:red">**Problem statement**</span>: given a distance measure $d$ and $k$, compute the k-clustering with maximum spacing.

Tim Roughgarden

# A Greedy Algorithm

- 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 clusters containing p & q into a single cluster

Note: just like kruskal's MST algorithm, but stopped early.
- points ⟷ vertices; distances ⟷ edge costs; point pairs ⟷ edges

⟹ called single-link clustering

(k=3)

Tim Roughgarden