

Clustering

Aarti Singh

Machine Learning 10-315

Apr 6, 2022

Some slides courtesy of Eric Xing, Carlos Guestrin



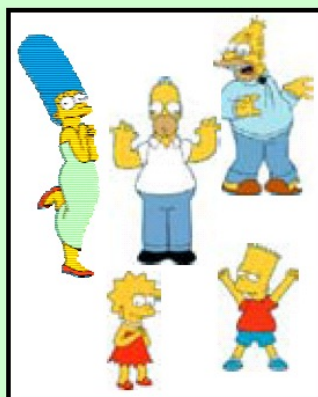
MACHINE LEARNING DEPARTMENT



What is clustering?

- Clustering: the process of grouping a set of objects into classes of similar objects
 - high intra-class similarity
 - low inter-class similarity
 - It is the most common form of **unsupervised learning**

Clustering is subjective



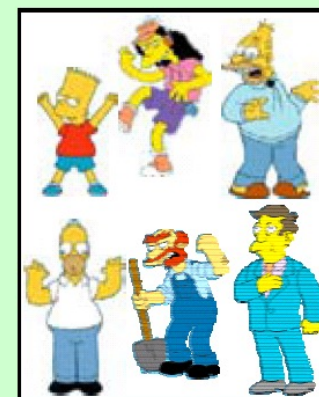
Simpson's Family



School Employees



Females



Males

What is Similarity?

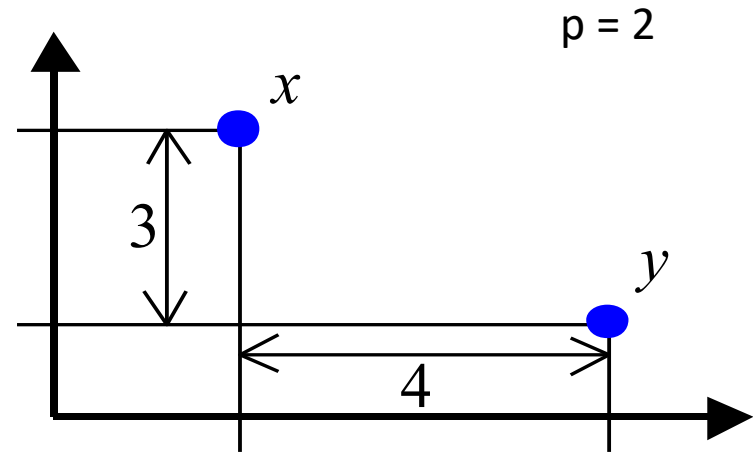


Hard to
define! But *we*
know it when
we see it

- The real meaning of similarity is a philosophical question. We will take a more pragmatic approach - think in terms of a distance (rather than similarity) between vectors or correlations between random variables.

Distance metrics

$$x = (x_1, x_2, \dots, x_p)$$
$$y = (y_1, y_2, \dots, y_p)$$



Euclidean distance

$$d(x, y) = \sqrt{\sum_{i=1}^p |x_i - y_i|^2}$$

✓ 5

Manhattan distance

$$d(x, y) = \sum_{i=1}^p |x_i - y_i|$$

✓ 7

Sup-distance

$$d(x, y) = \max_{1 \leq i \leq p} |x_i - y_i|$$

✓ 4

Correlation coefficient

$$x = (x_1, x_2, \dots, x_p)$$

$$y = (y_1, y_2, \dots, y_p)$$

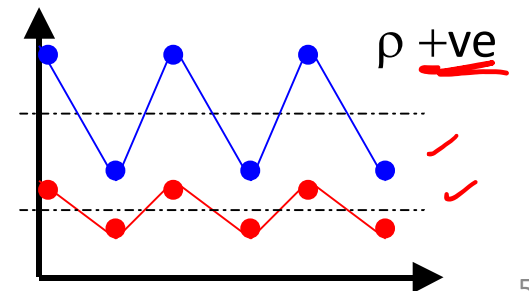
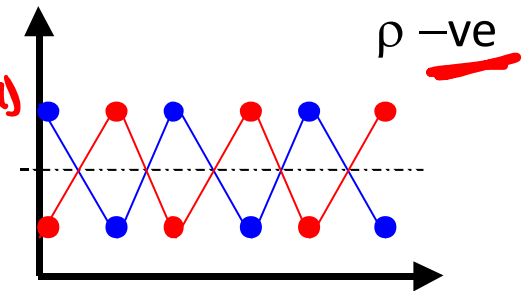
Random vectors (e.g. expression levels of two genes under various drugs)

Pearson correlation coefficient

$$\rho(x, y) = \frac{\sum_{i=1}^p (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^p (x_i - \bar{x})^2 \times \sum_{i=1}^p (y_i - \bar{y})^2}}$$

Handwritten red notes:
 $\text{Cov}(x, y)$ (above the numerator)
 $\sqrt{\text{var}(x) \text{var}(y)}$ (below the denominator)
 Red arrows point from these notes to the corresponding parts of the formula.

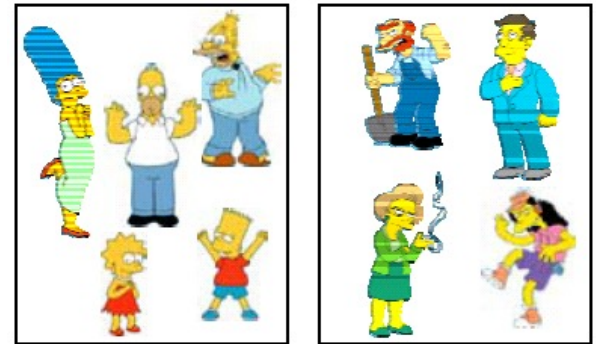
where $\bar{x} = \frac{1}{p} \sum_{i=1}^p x_i$ and $\bar{y} = \frac{1}{p} \sum_{i=1}^p y_i$.



Clustering Algorithms

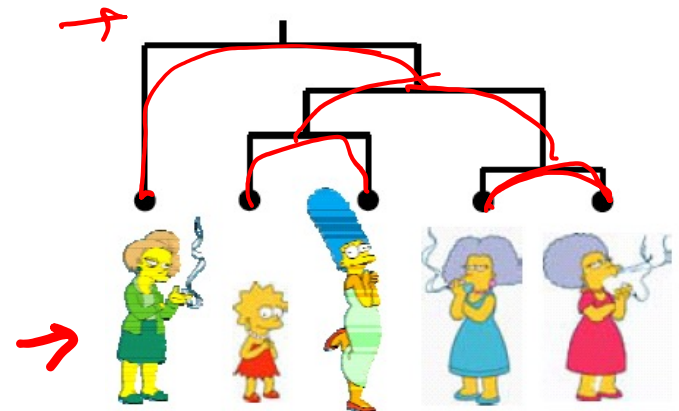
- **Partition algorithms**

- K means clustering ✓
- Mixture-Model based clustering ✓



- **Hierarchical algorithms**

- Single-linkage ✓
- Average-linkage ✓
- Complete-linkage ✓
- Centroid-based ✓



Partitioning Algorithms

- Partitioning method: Construct a partition of n objects into a set of K clusters
- Given: a set of objects and the number K
- Find: a partition of K clusters that optimizes the chosen partitioning criterion
 - Globally optimal: exhaustively enumerate all partitions ↗
 - Effective heuristic method: K-means algorithm ↗

K-Means

Algorithm

Input – Desired number of clusters, k ✓

Initialize – the k cluster centers (randomly if necessary)

Iterate –

1. Assign points to the nearest cluster centers ✓
2. Re-estimate the k cluster centers (aka the **centroid** or **mean**), by assuming the memberships found above are correct.

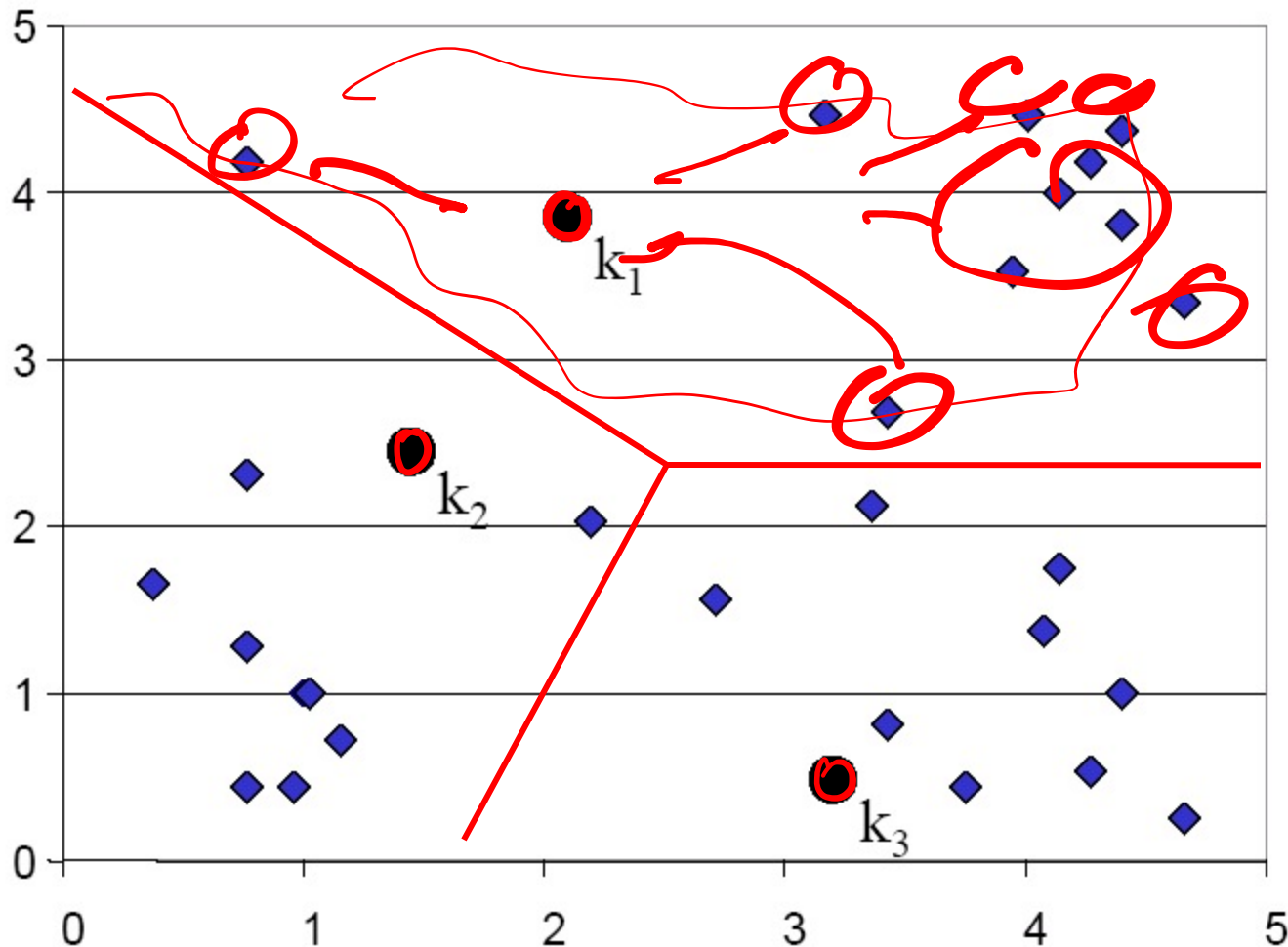
$$\underline{\vec{\mu}_k} = \frac{1}{\underline{c_k}} \sum_{\underline{i \in \mathcal{C}_k}} \vec{x}_i \quad \checkmark$$

Termination –

If none of the objects changed membership in the last iteration, exit.
Otherwise go to 1.

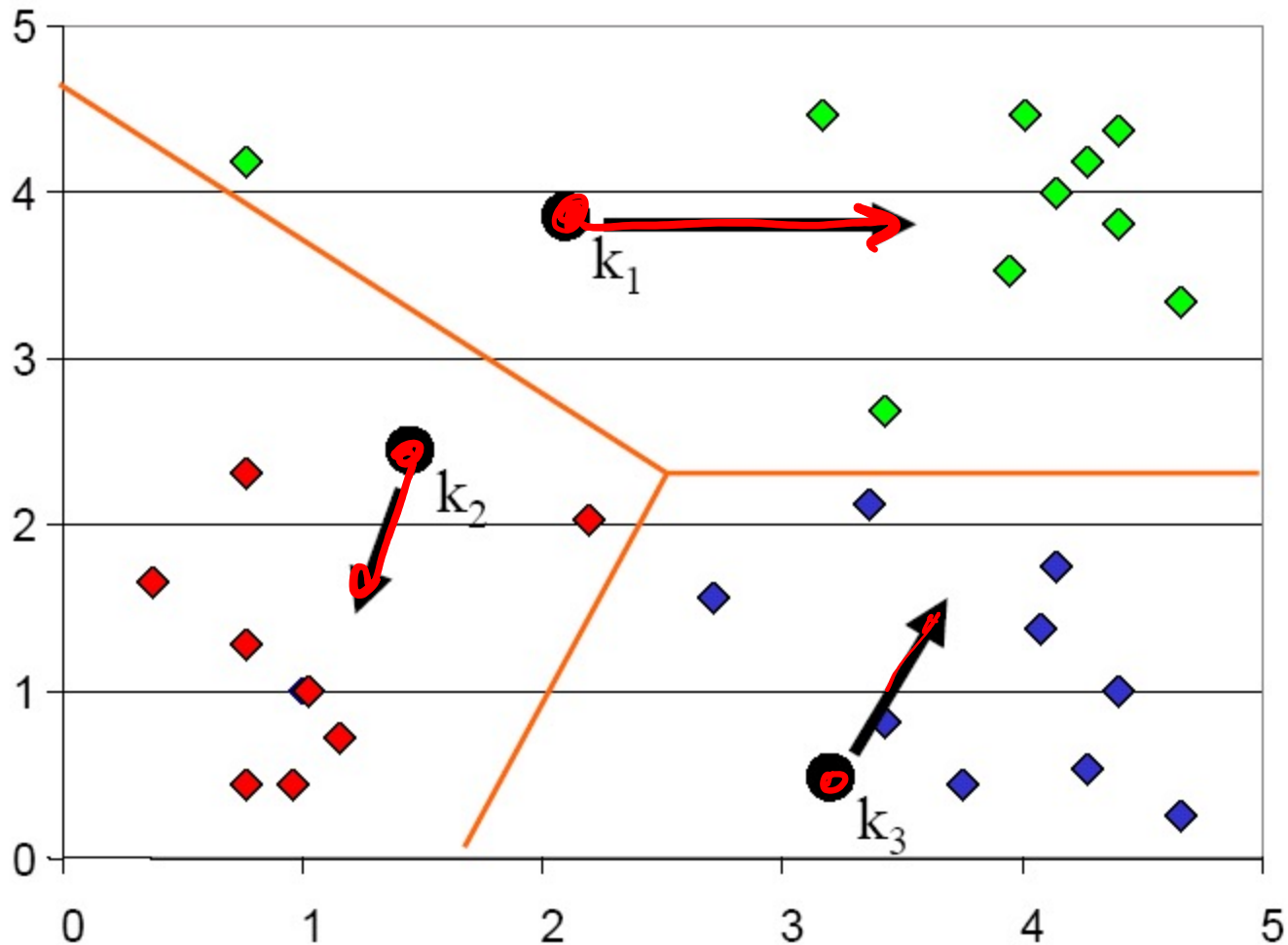
K-means Clustering: Step 1

¹/₃

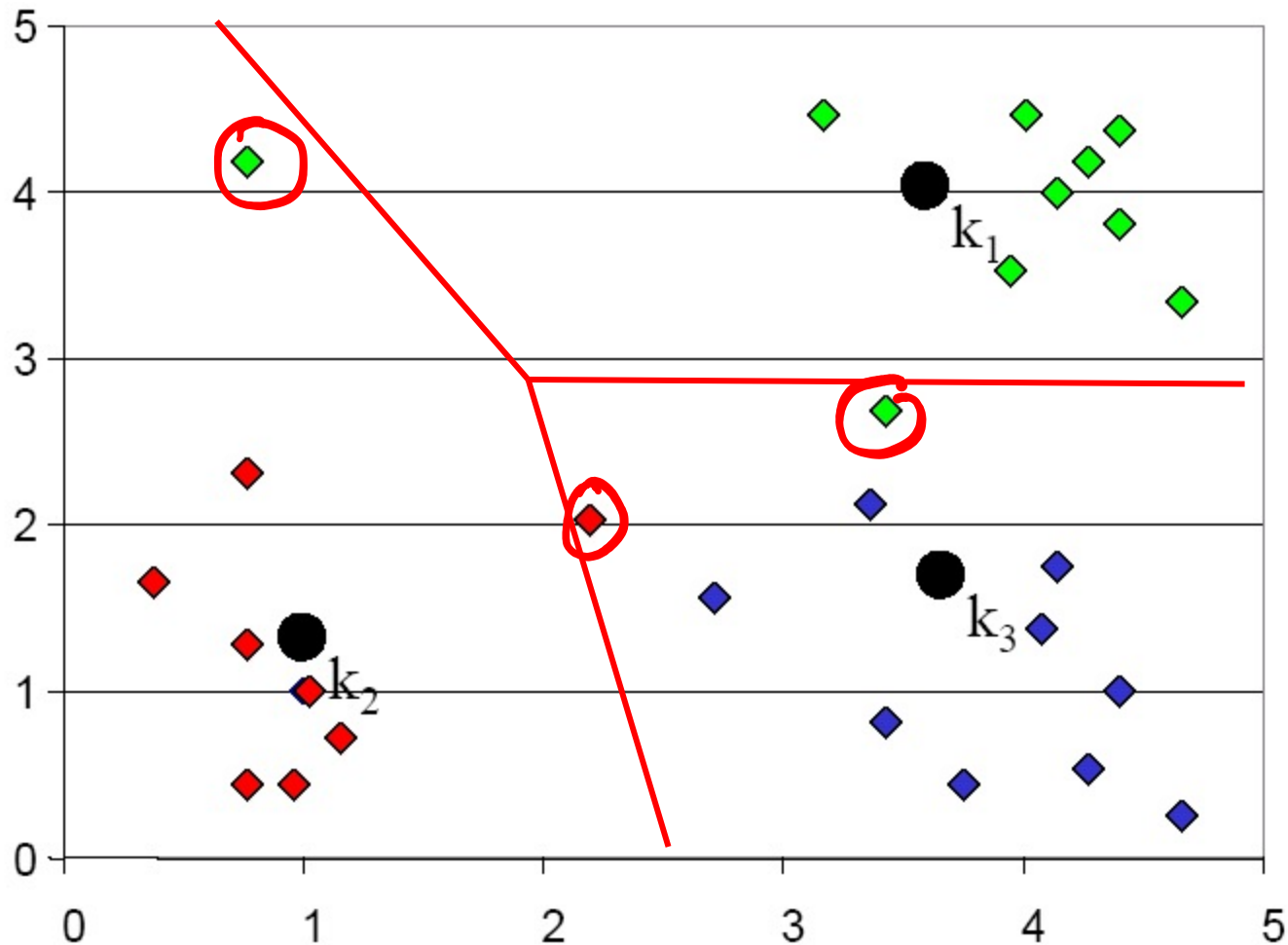


Voronoi
diagram

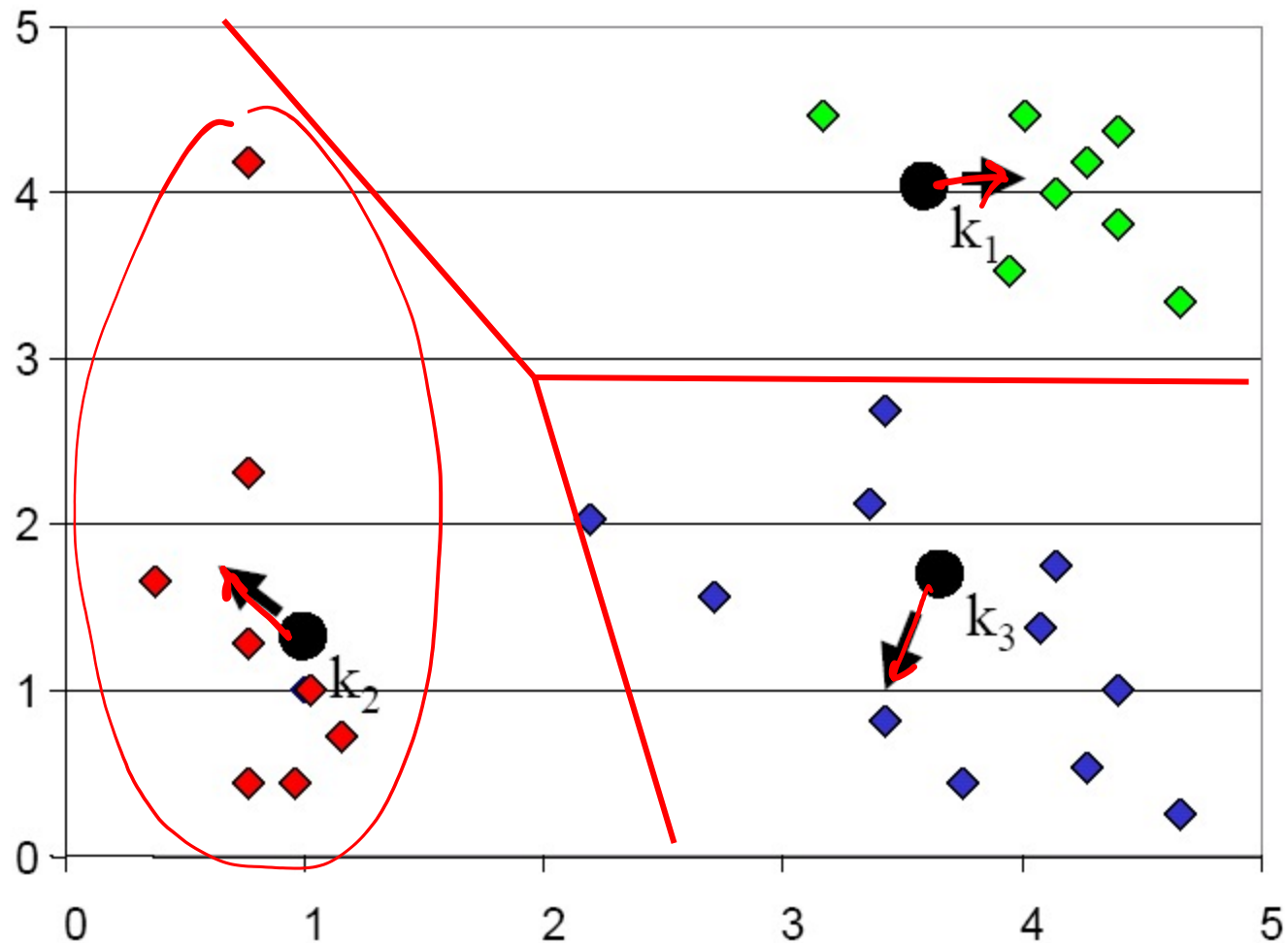
K-means Clustering: Step 2



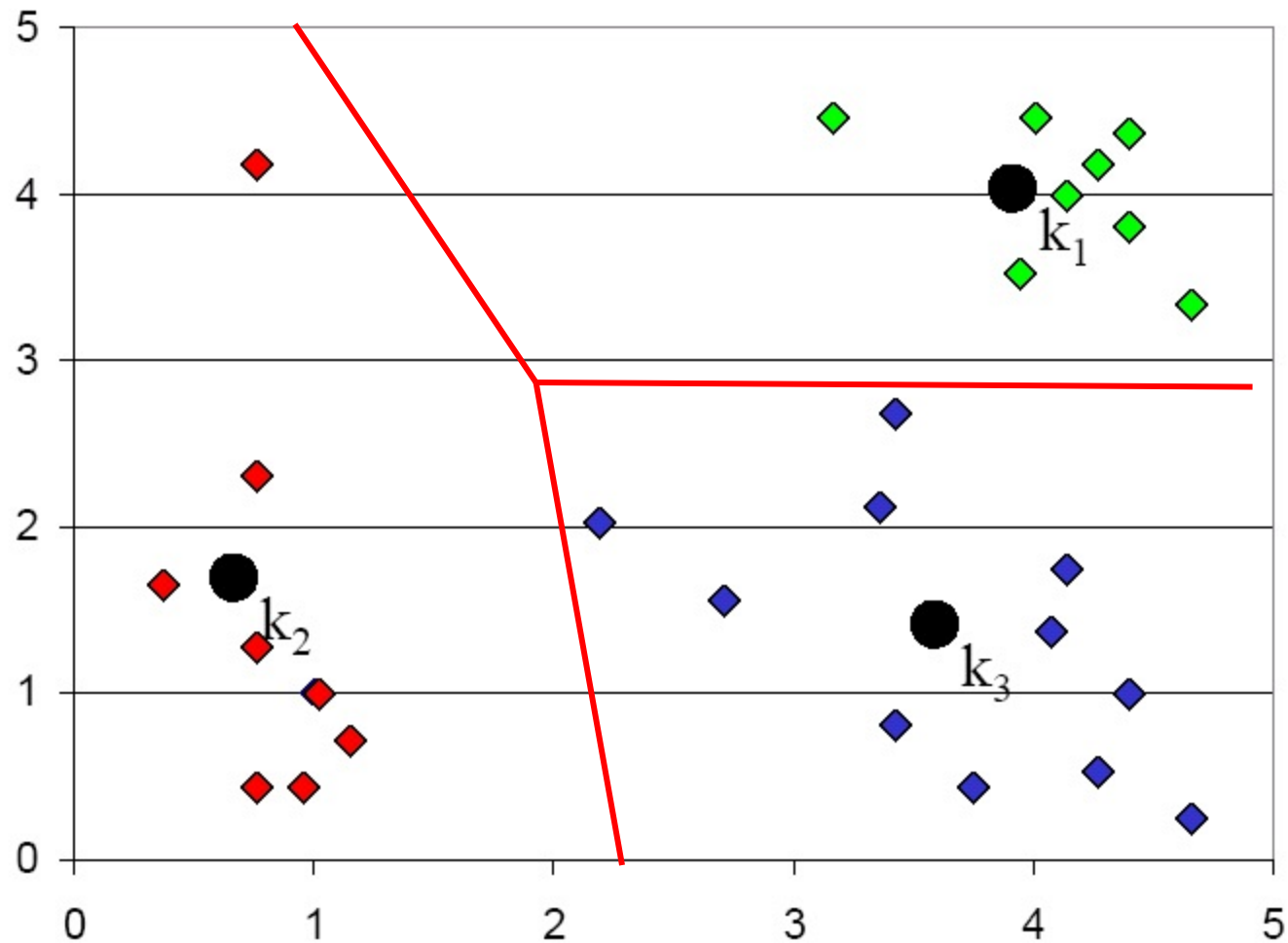
K-means Clustering: Step 3



K-means Clustering: Step 4



K-means Clustering: Step 5



K-means Recap ...

- Randomly initialize k centers
 - $\mu^{(0)} = \underbrace{\mu_1^{(0)}}_{\text{red}}, \dots, \underbrace{\mu_k^{(0)}}_{\text{red}}$

K-means Recap ...

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

Iterate $t = 0, 1, 2, \dots$

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $\underline{C^{(t)}}(\underline{j}) \leftarrow \arg \min_{\underline{i=1, \dots, k}} \|\underline{\mu_i^{(t)}} - \underline{x_j}\|^2$

*cluster assignment
for point j*

K-means Recap ...

- Randomly initialize k centers

- $\mu^{(0)} = \mu_1^{(0)}, \dots, \mu_k^{(0)}$

Iterate $t = 0, 1, 2, \dots$

- **Classify:** Assign each point $j \in \{1, \dots, m\}$ to nearest center:

- $C^{(t)}(j) \leftarrow \arg \min_{i=1, \dots, k} \|\mu_i^{(t)} - x_j\|^2$

- **Recenter:** μ_i becomes centroid of its points:

- $\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2 \quad i \in \{1, \dots, k\}$

- Equivalent to $\mu_i \leftarrow$ average of its points!

What is K-means optimizing?

- Potential function $F(\mu, C)$ of centers μ and point allocations C :

$$F(\mu, C) = \sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2$$

$C(j)$ = cluster to which j is assigned

$$= \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2$$

- Optimal K-means:

$$\square \min_{\mu} \min_C F(\mu, C)$$

$f(x_1, x_2)$

➤ Is the K-means objective convex?

K-means algorithm

- Optimize potential function:

$$\min_{\mu} \min_C F(\mu, C) = \min_{\mu} \min_C \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2 \quad \leftarrow$$

- K-means algorithm:** (coordinate descent on F)

(1) Fix μ , optimize C

Expected cluster assignment ✓

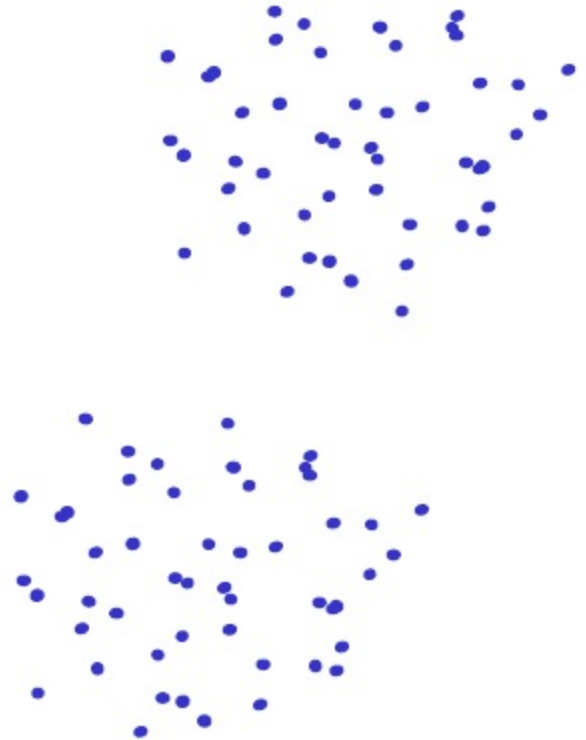
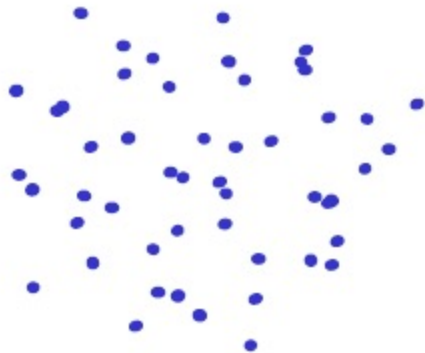
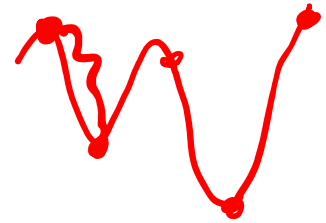
(2) Fix C , optimize μ

Maximum likelihood for center ✓

Generalization: EM (Expectation-Maximization) algorithm

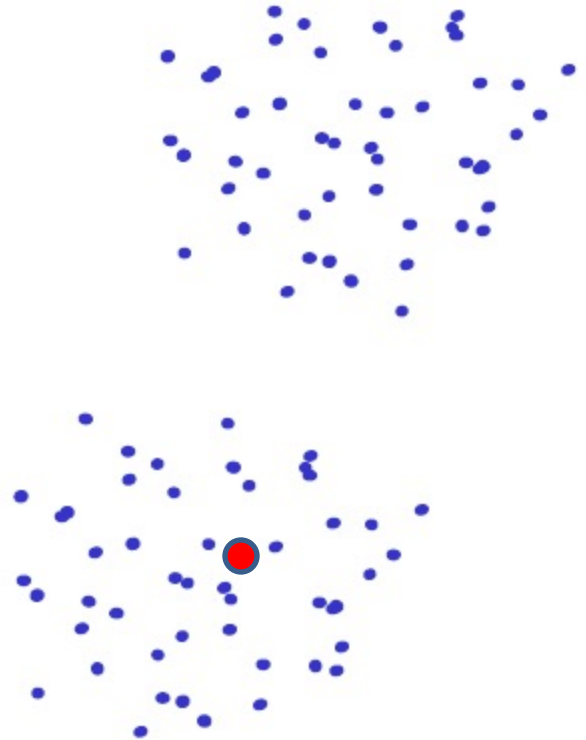
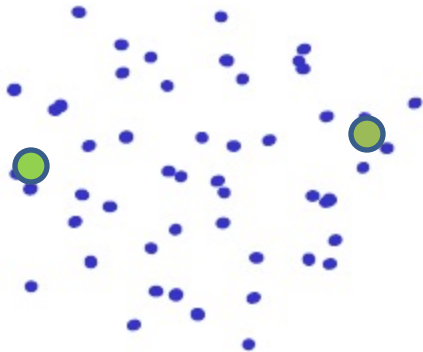
Seed Choice

- Results are quite sensitive to seed selection.



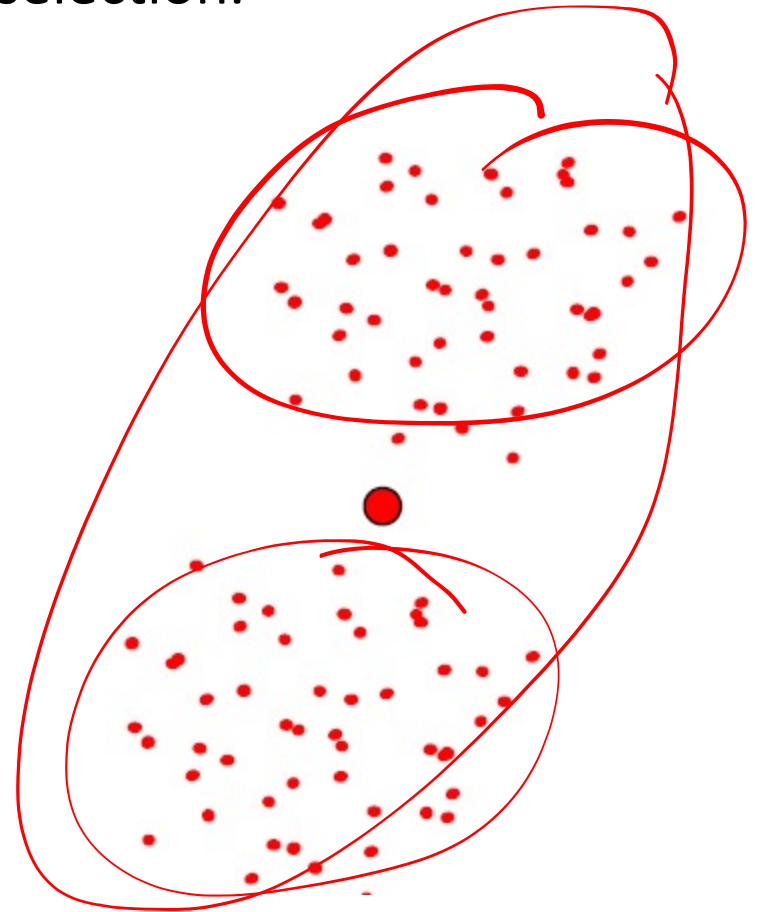
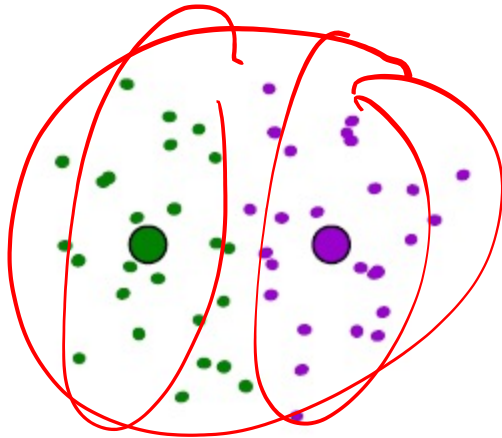
Seed Choice

- Results are quite sensitive to seed selection.



Seed Choice

- Results are quite sensitive to seed selection.



Seed Choice

- Results can vary based on random seed selection.
 - Some seeds can result in poor convergence rate, or convergence to sub-optimal clustering.
 - Try out multiple starting points (very important!!!) ✓
 - k-means ++ algorithm of Arthur and Vassilvitskii
- key idea: choose centers that are far apart
- (probability of picking a point as cluster center \propto distance from nearest center picked so far)

Other Issues

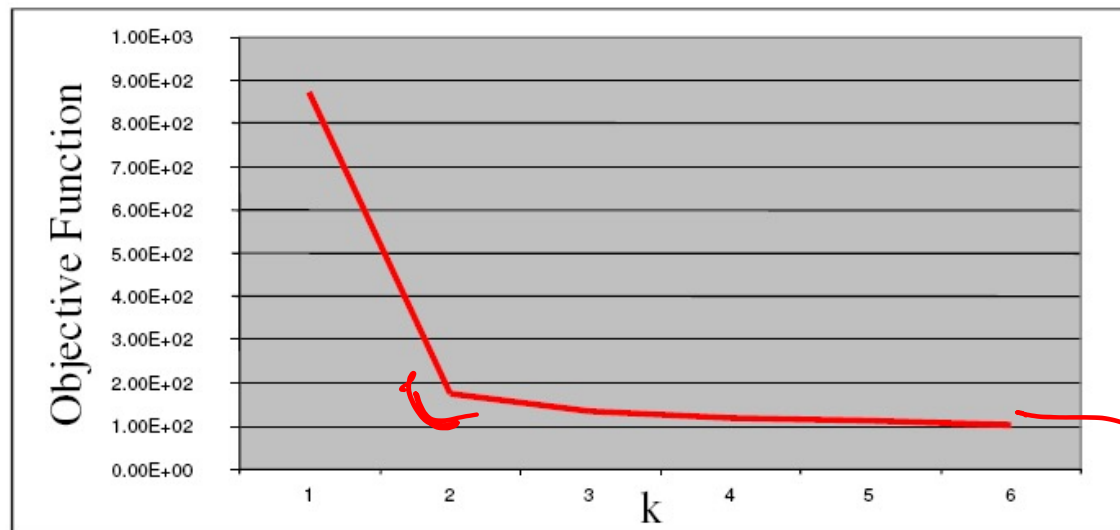
- Number of clusters **K**

- Objective function

$$\sum_{j=1}^m \|\mu_{C(j)} - x_j\|^2 \quad \leftarrow$$

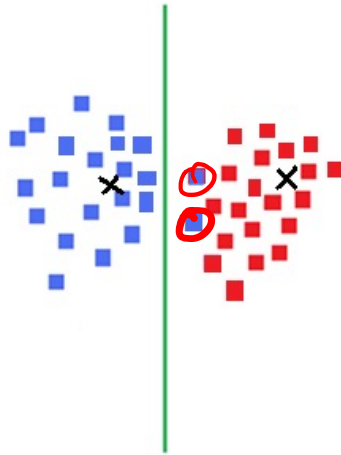
➤ Can you pick K by minimizing the objective over K?

- Look for “Knee” in objective function



Other Issues

- Sensitive to Outliers
 - use K-medoids

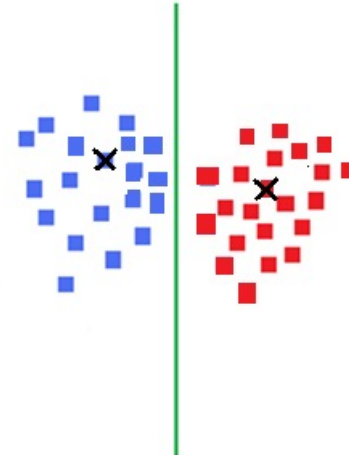


$$d(\mu, x) = (\mu - x)^2$$

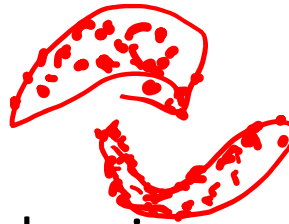
↓

$$d(\mu, x) = |\mu - x|$$

↓
median



- Shape of clusters
 - Assumes isotropic, equal variance, convex clusters



Partitioning Algorithms

- K-means
 - hard assignment: each object belongs to only one cluster
- Mixture modeling
 - soft assignment: probability that an object belongs to a cluster



Generative approach