

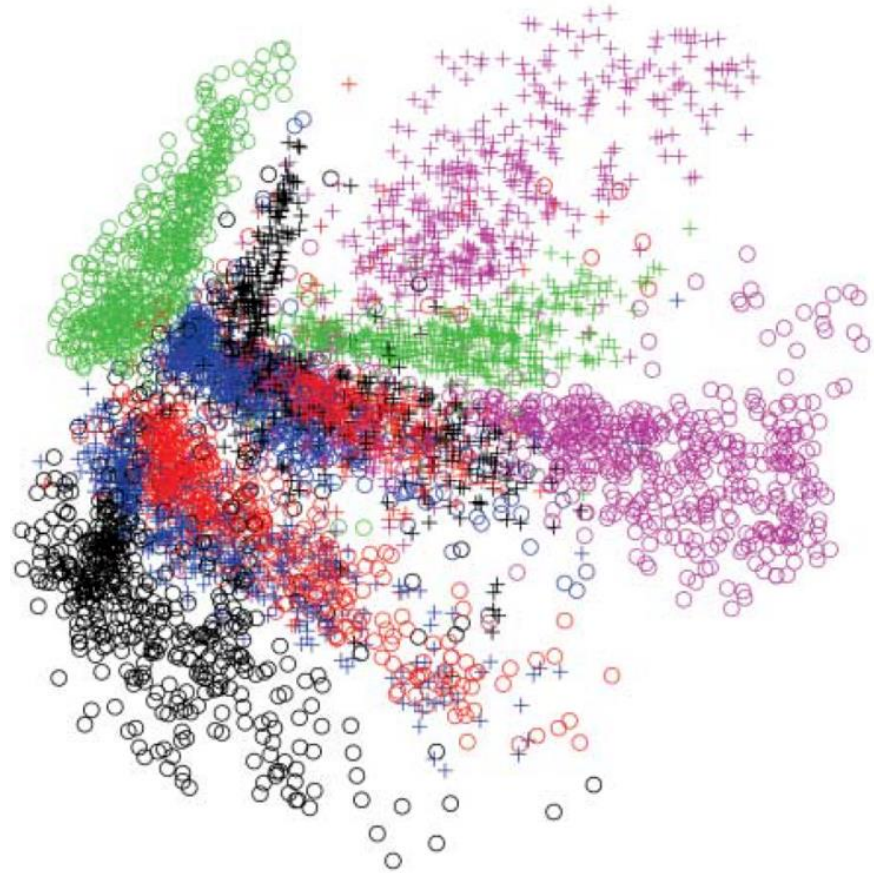
# Announcements

## Assignments

- HW9 (online)
  - Due Thu 4/16, 11:59 pm

# Dimensionality Reduction

## MNIST digit autoencoder



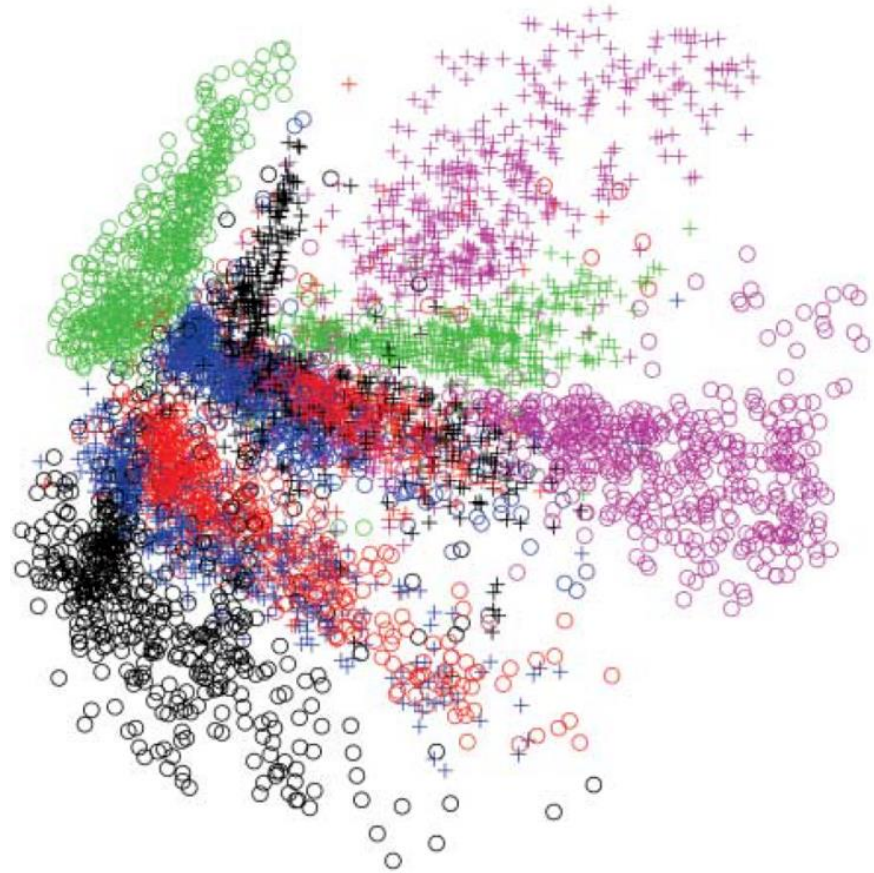
# Piazza Poll 1

Are autoencoders an example of unsupervised learning?

A.

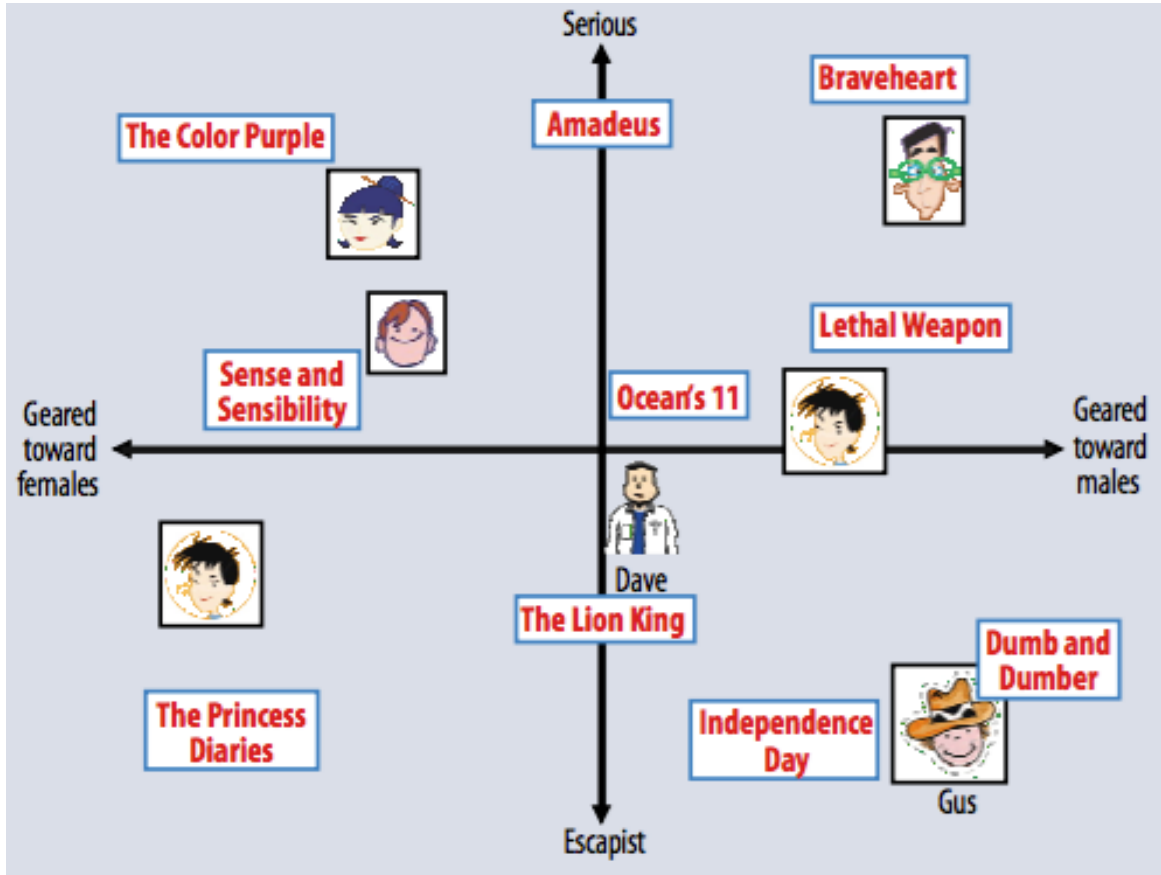
B.

C.



# Recommender Systems

## Matrix Factorization



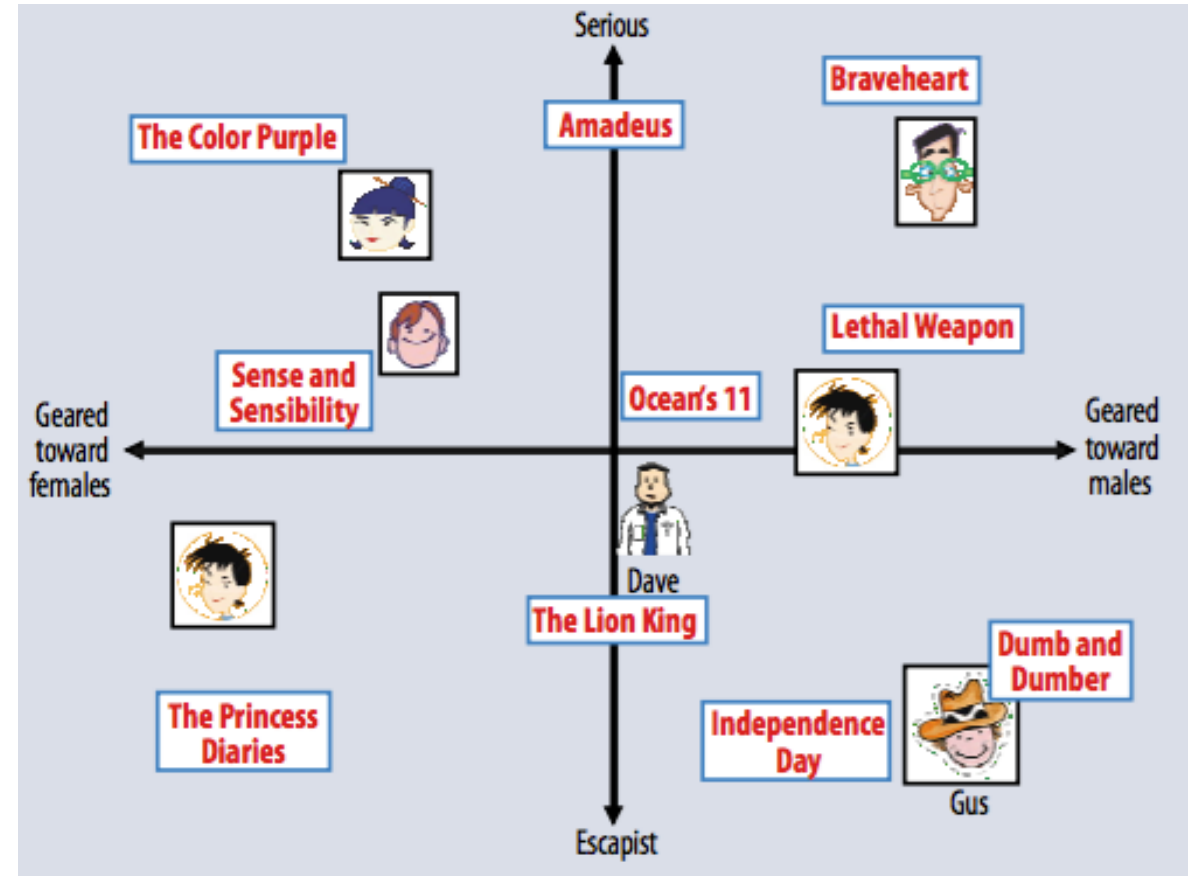
# Piazza Poll 2

Are recommender systems an example of unsupervised learning?

A.

B.

C.



# Plan

## Last week

- Dimensionality reduction
  - PCA
  - Autoencoders
- Recommender systems

## This week

- Clustering

## Next week

- Learning Theory

# Introduction to Machine Learning

## Clustering

Instructor: Pat Virtue

# Clustering, Informal Goals

**Goal:** Automatically partition **unlabeled** data into groups of similar datapoints.

**Question:** When and why would we want to do this?

**Useful for:**

- Automatically organizing data.
- Understanding hidden structure in data.
- Preprocessing for further analysis.
  - Representing high-dimensional data in a low-dimensional space (e.g., for visualization purposes).

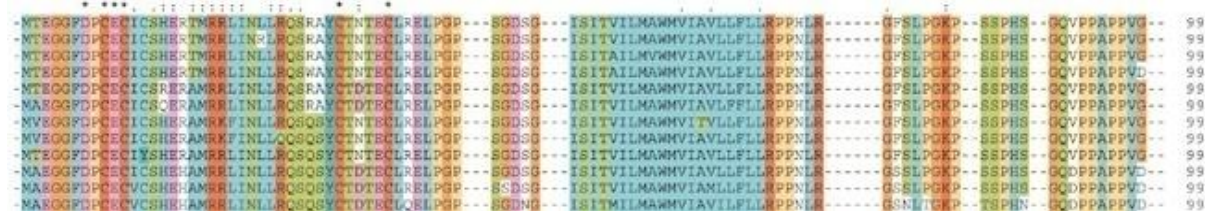


# Applications (Clustering comes up everywhere...)

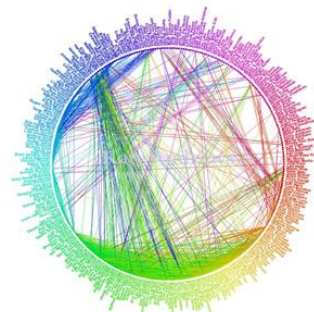
Cluster news articles or web pages or search results by topic.



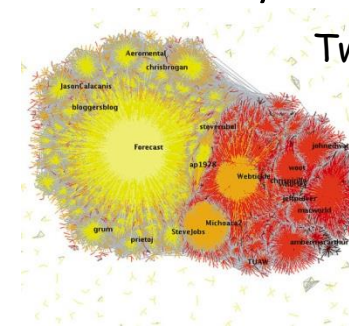
- Cluster protein sequences by function or genes according to expression profile.



- Cluster users of social networks by interest (community detection).



Facebook network



Twitter Network

# Applications (Clustering comes up everywhere...)

Cluster customers according to purchase history.



- Cluster galaxies or nearby stars (e.g. Sloan Digital Sky Survey)

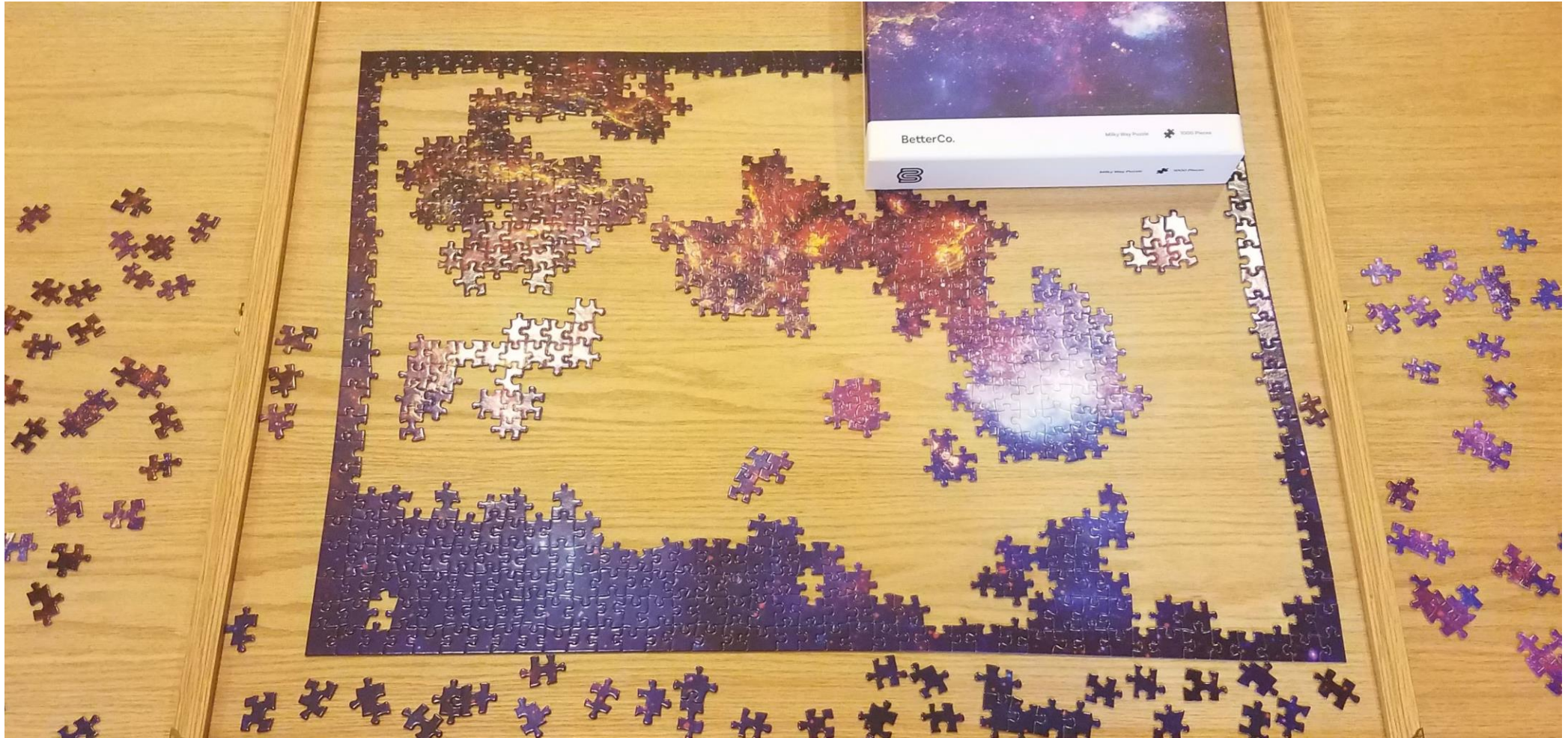


- And many many more applications....



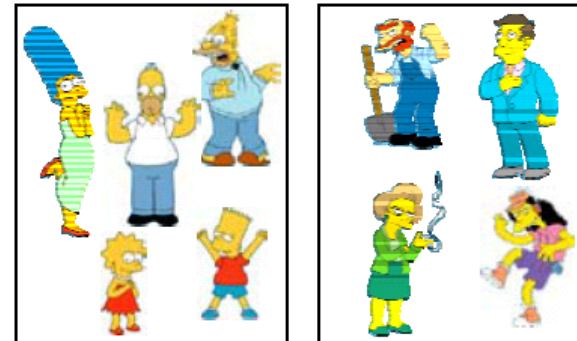
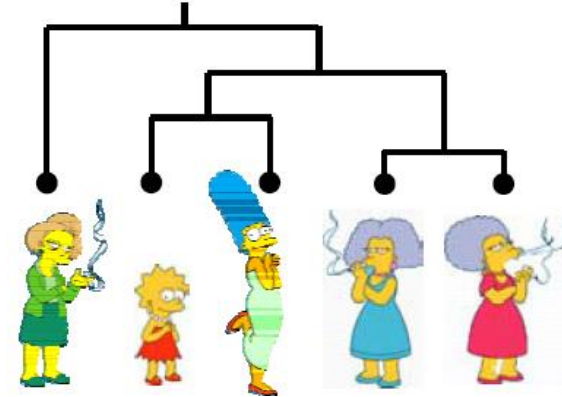
# Clustering Applications

Jigsaw puzzles!



# Clustering Algorithms

- Hierarchical algorithms
  - Bottom-up: Agglomerative Clustering
  - Top-down: Divisive
- Partition algorithms
  - K means clustering
  - Mixture-Model based clustering



# Hierarchical Clustering

- Bottom-Up Agglomerative Clustering

Starts with each object in a separate cluster, and repeat:

- Joins the most similar pair of clusters,
- Update the similarity of the new cluster to others until there is only one cluster.

Greedy - less accurate but simple to implement

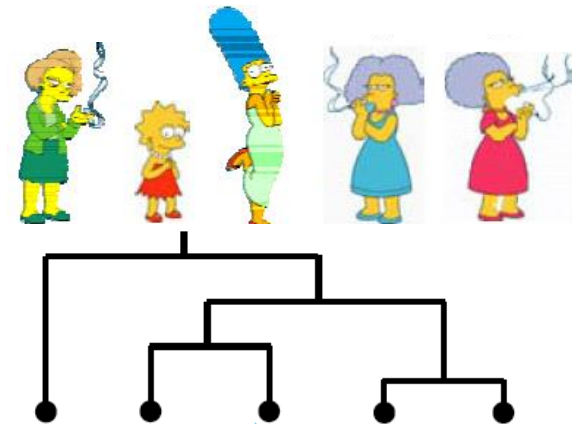


- Top-Down divisive

Starts with all the data in a single cluster, and repeat:

- Split each cluster into two using a partition algorithm
- Until each object is a separate cluster.

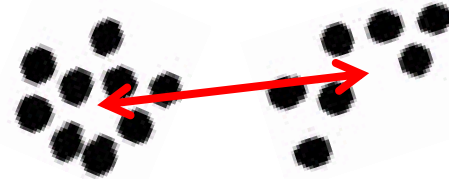
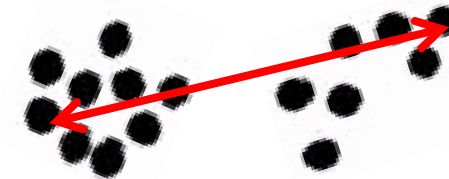
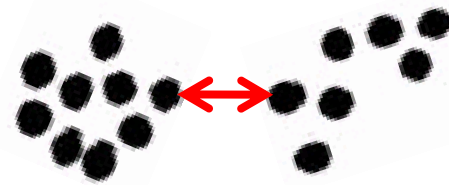
More accurate but complex to implement



# Bottom-up Agglomerative clustering

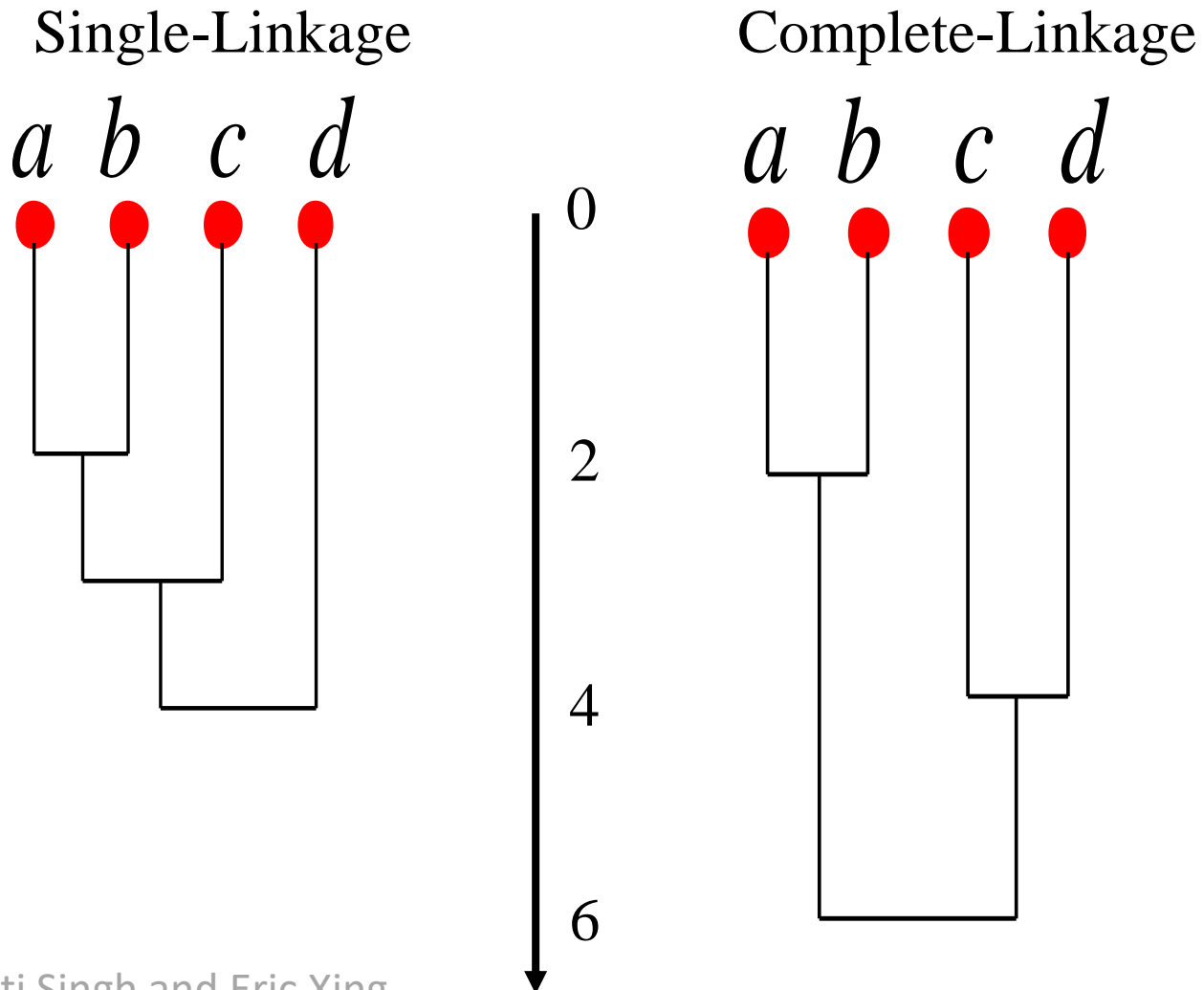
Different algorithms differ in how the similarities are defined (and hence updated) between two clusters

- Single-Linkage
  - Nearest Neighbor: similarity between their closest members.
- Complete-Linkage
  - Furthest Neighbor: similarity between their furthest members.
- Centroid
  - Similarity between the centers of gravity
- Average-Linkage
  - Average similarity of all cross-cluster pairs.



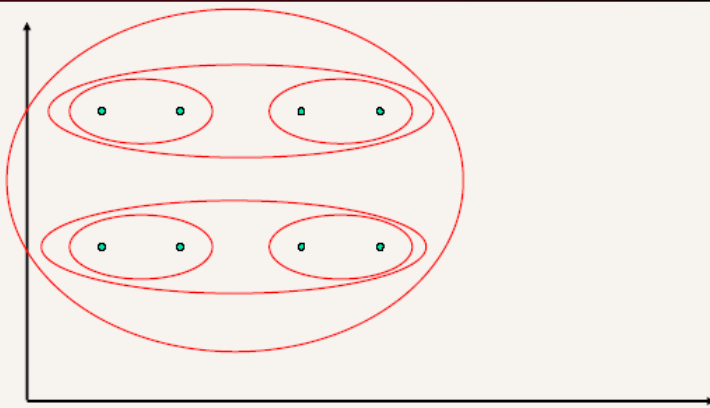


# Dendrograms

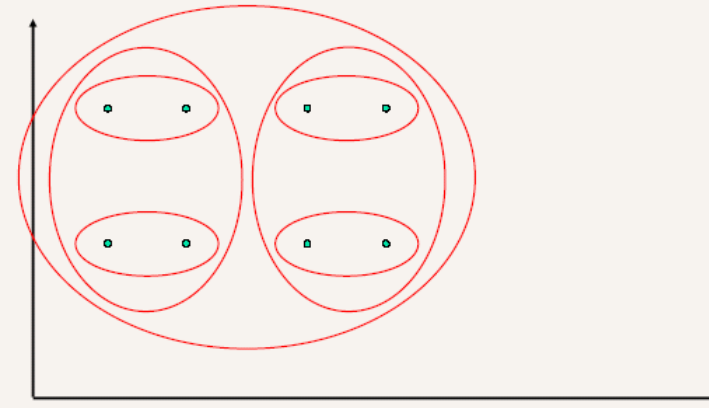


# Another Example

Single Link Example



Complete Link Example





# Single vs. Complete Linkage

## Shape of clusters

Single-linkage

allows anisotropic and non-convex shapes

Complete-linkage

assumes isotropic, convex shapes



# 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** – Data,  $\mathbf{x}^{(i)}$ , Desired number of clusters,  $K$

**Initialize** – the  $K$  cluster centers (randomly if necessary)

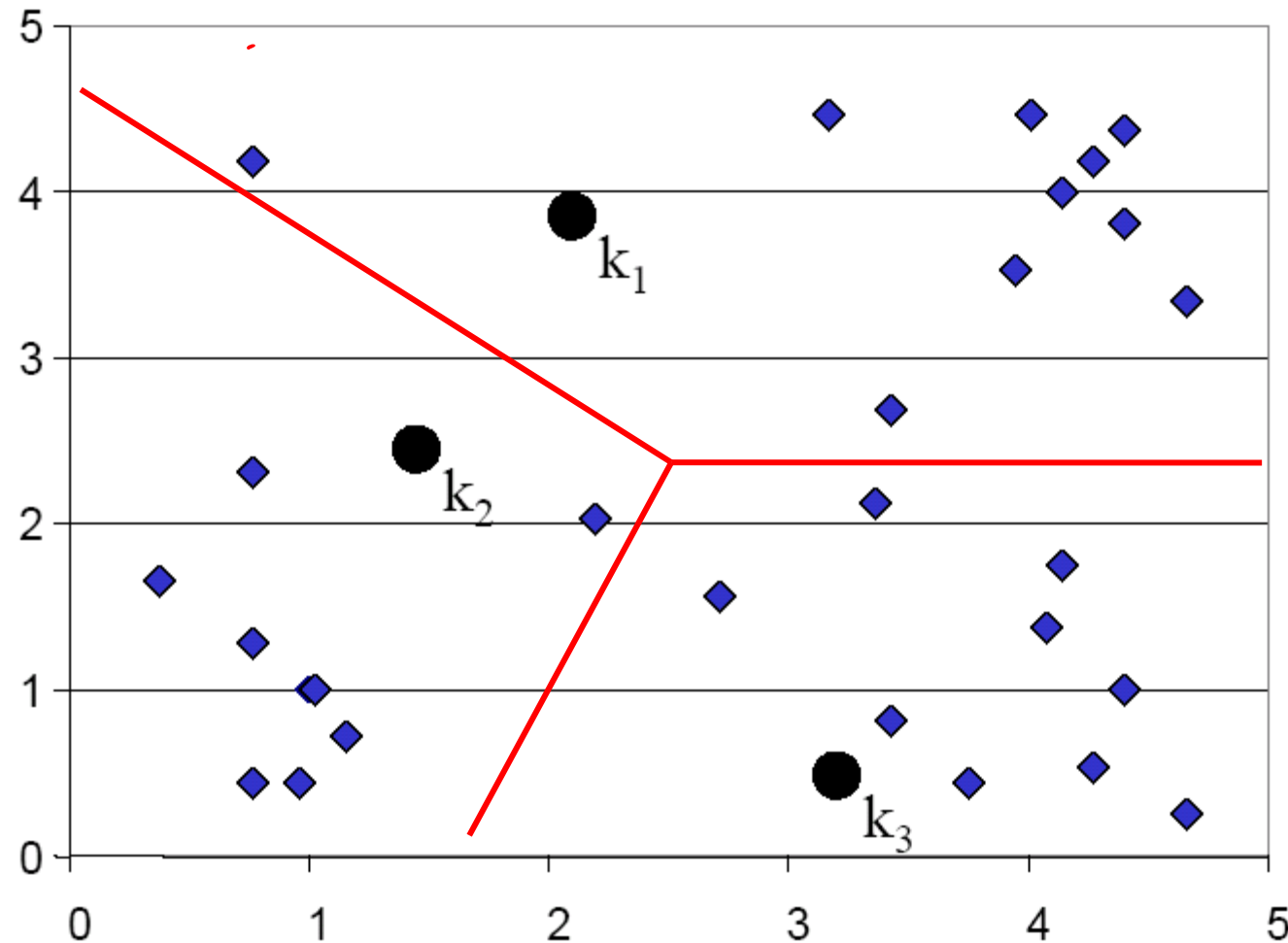
**Iterate** –

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.

**Termination** –

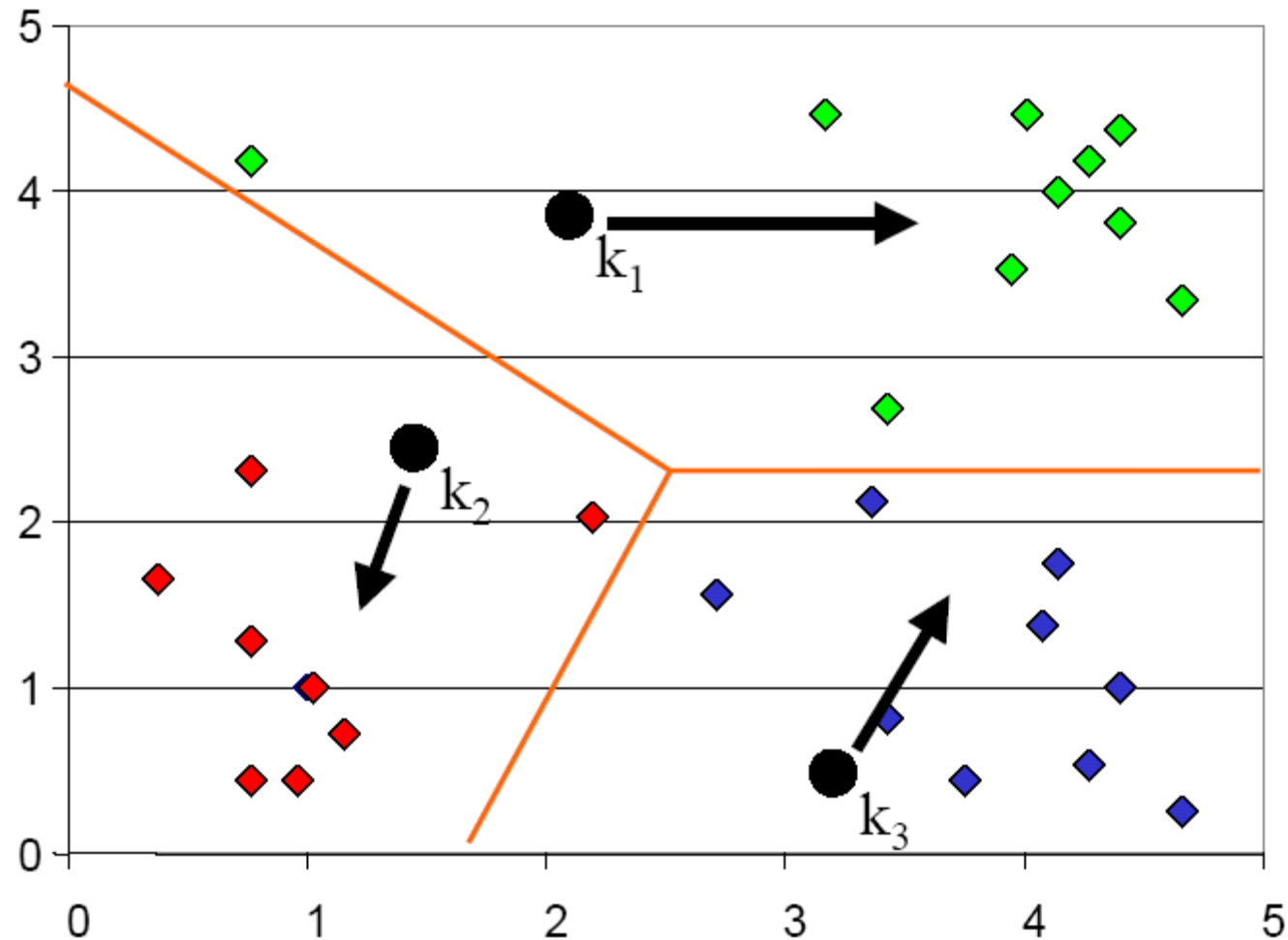
If none of the objects changed membership in the last iteration, exit.  
Otherwise go to 1.

# K-means Clustering: Step 1

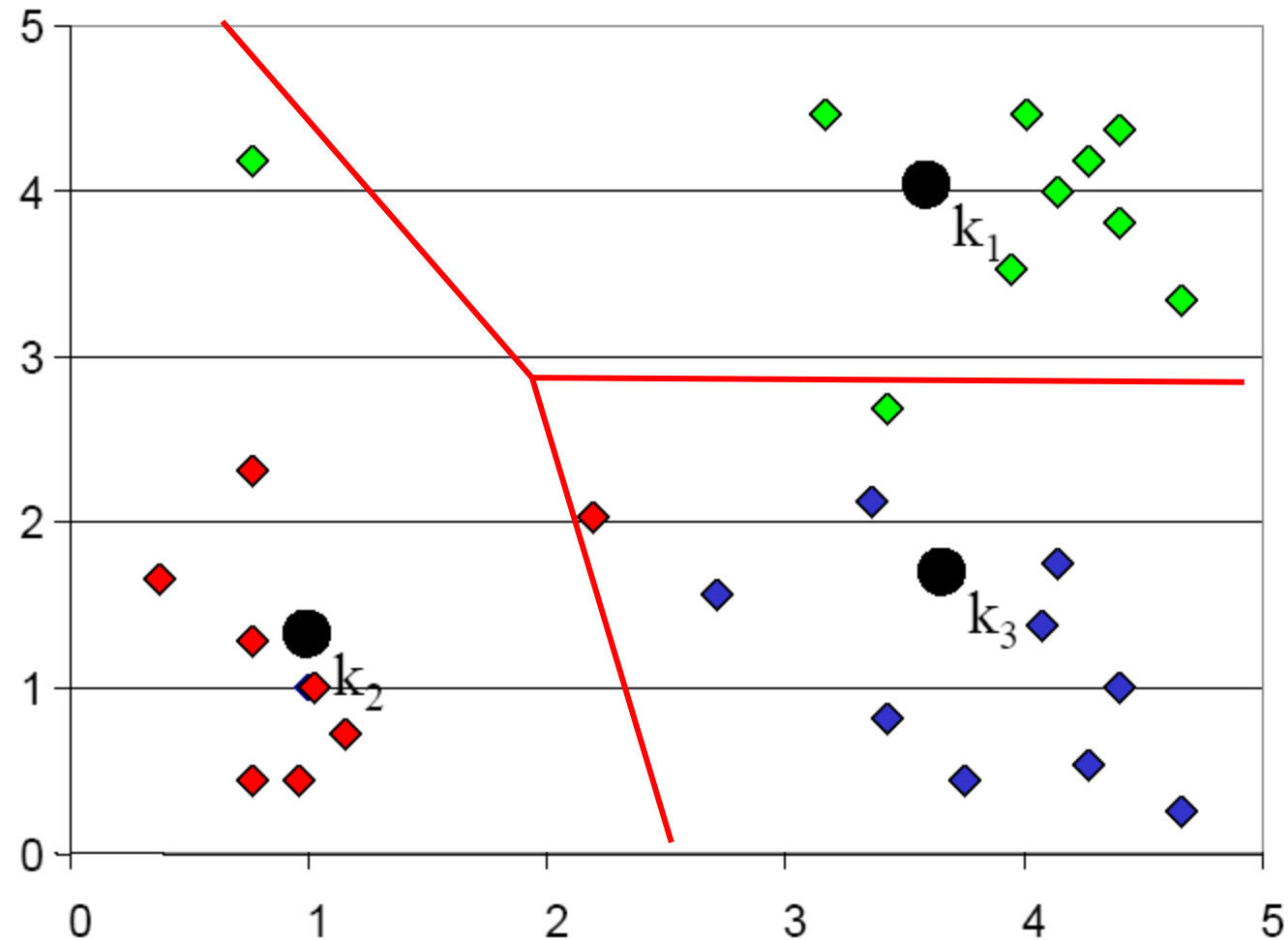


**Voronoi  
diagram**

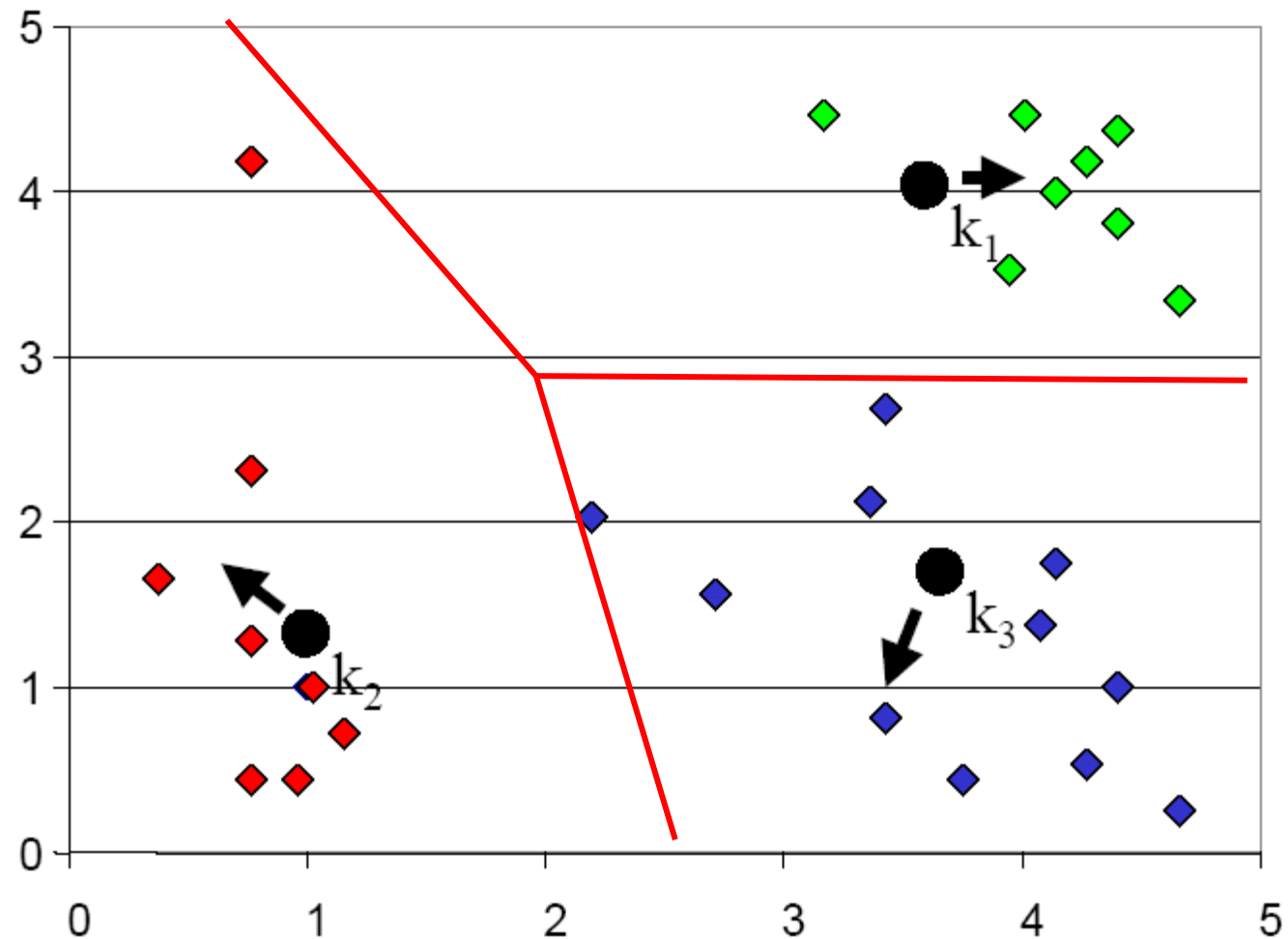
# K-means Clustering: Step 2



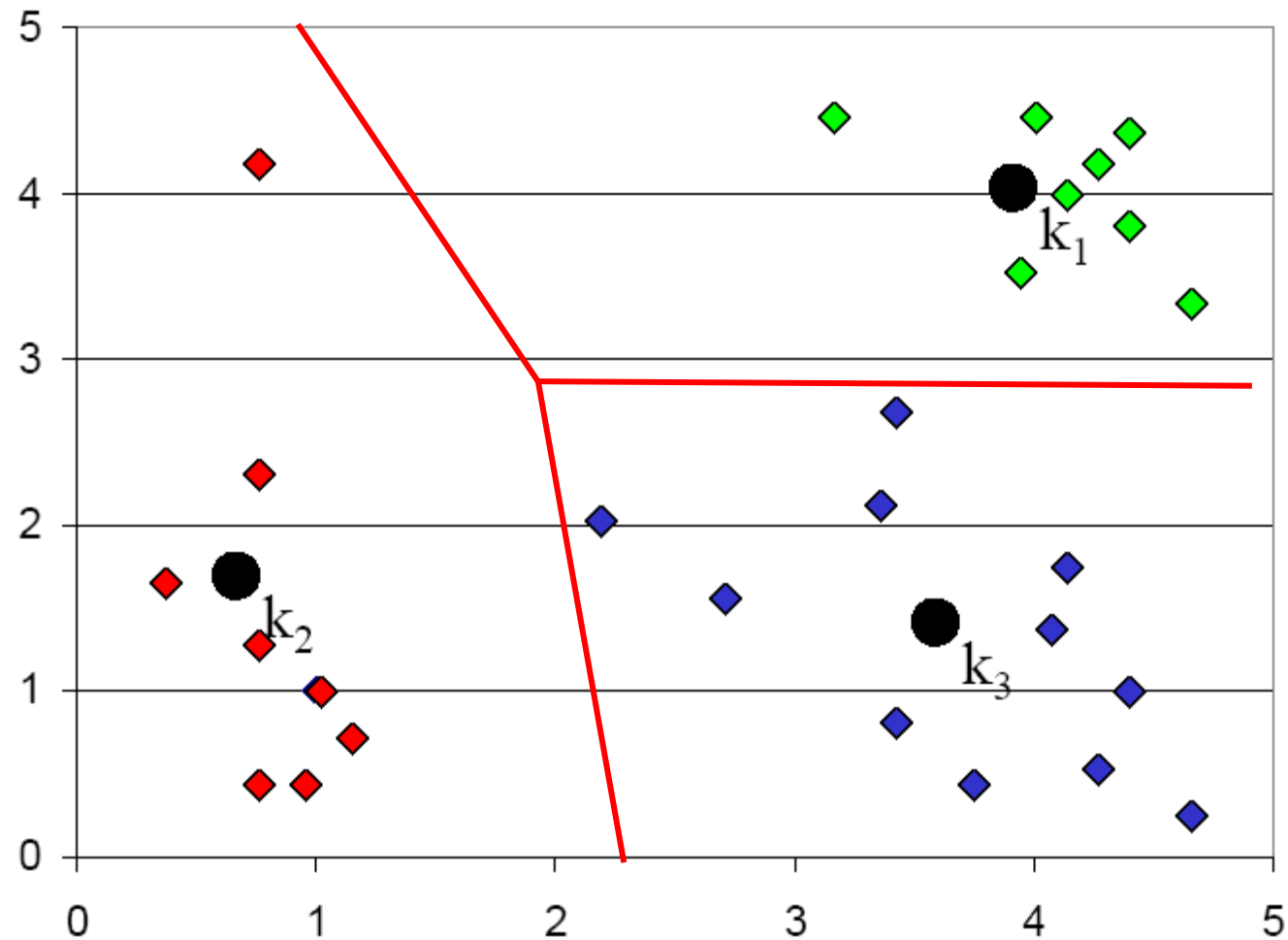
# K-means Clustering: Step 3



# K-means Clustering: Step 4



# K-means Clustering: Step 5



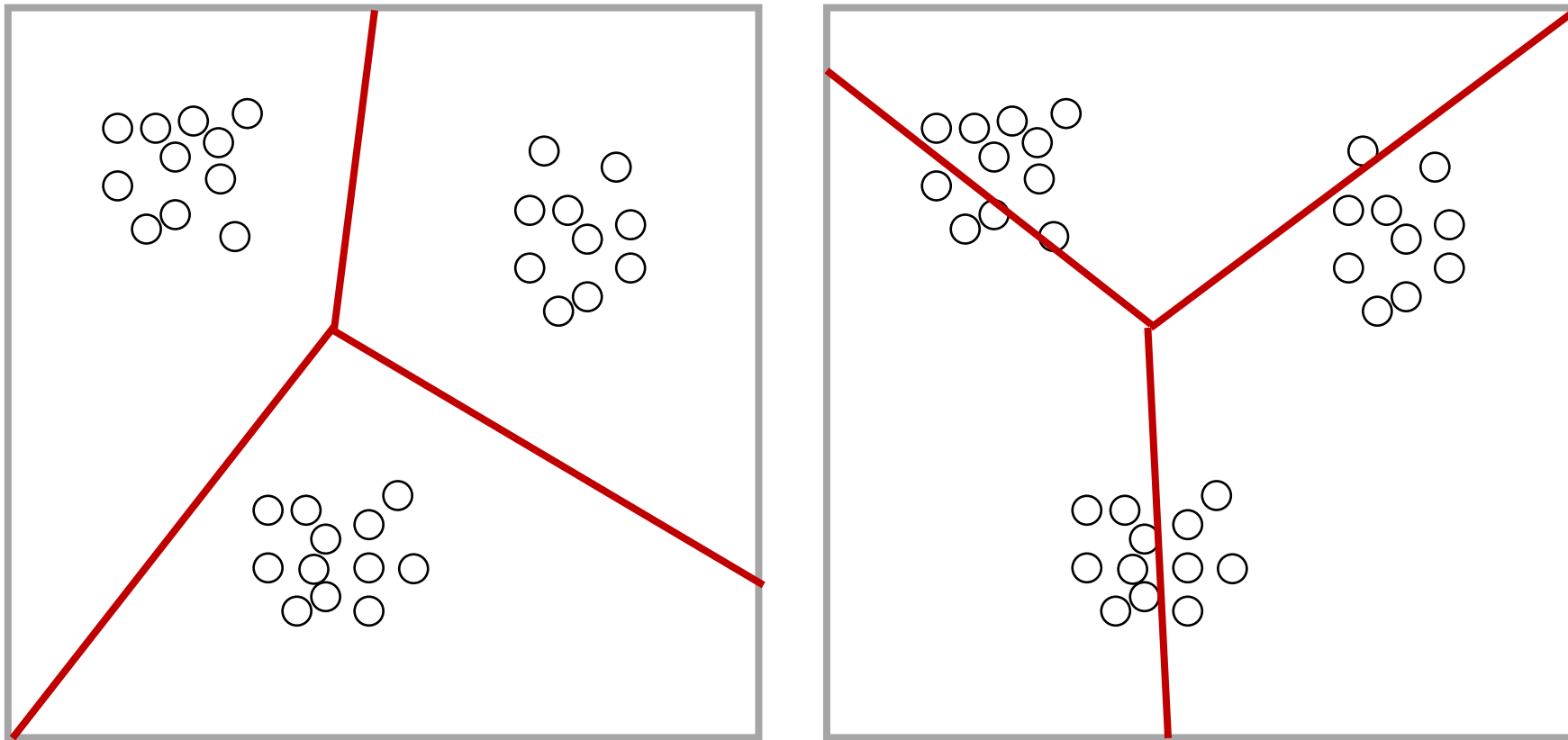


# K-means Optimization

## Optimization recipe

# K-means Optimization

Question: Which of these partitions is “better”?



# K-means Optimization

# K-means Optimization

Computational complexity

$$\mathbf{C}, \mathbf{z} = \operatorname{argmin}_{\mathbf{C}, \mathbf{z}} \sum_{i=1}^N \|\mathbf{x}^{(i)} - \mathbf{c}_{\mathbf{z}(i)}\|_2^2$$

# K-means Optimization

Alternating minimization

# K-means Optimization

Alternating minimization

# Alternating minimization

Where have we seen this before?

Recommender Systems

$$\min_{\mathbf{U}, \mathbf{V}} J(\mathbf{U}, \mathbf{V}) \quad J(\mathbf{U}, \mathbf{V}) = \sum_{i,j \in \mathcal{S}} \left( R_{ij} - \mathbf{u}^{(i)T} \mathbf{v}^{(j)} \right)^2$$

$$\textcircled{1} \min_{\mathbf{u}} J(\mathbf{u}, \underline{\mathbf{v}})$$

$$\mathbf{u}^{t+1} \leftarrow \mathbf{u}^t - \alpha \nabla_{\mathbf{u}} J$$

$$\textcircled{2} \min_{\mathbf{v}} J(\underline{\mathbf{u}}, \mathbf{v})$$

$$\mathbf{v} \leftarrow \mathbf{v} - \alpha \nabla_{\mathbf{v}} J$$

# Alternating minimization

Two different approaches

$$\min_{\alpha, \beta} J(\alpha, \beta)$$



# Alternating minimization

Two different approaches

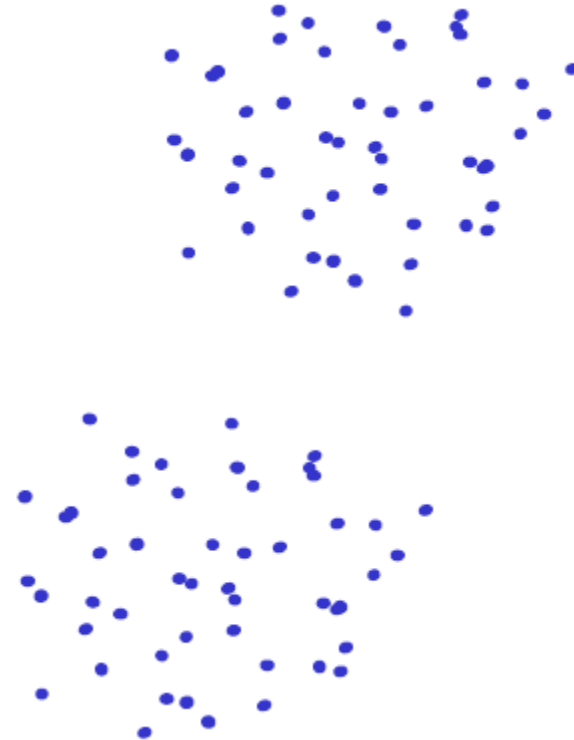
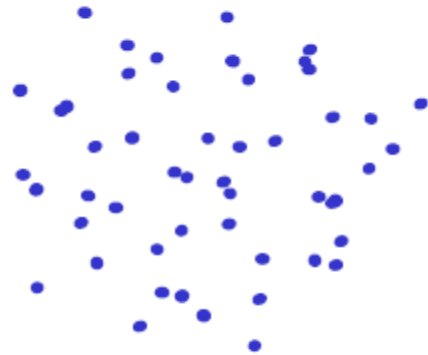
$$\min_{\theta_1, \theta_2} J(\theta_1, \theta_2)$$

# Computational Complexity

- At each iteration,
  - Computing cluster centers: Each object gets added once to some cluster:  $O(N)$
  - Computing distance between each of the  $N$  objects and the  $K$  cluster centers is  $O(KN)$
- Assume these two steps are each done once for  $l$  iterations:  $O(lKN)$

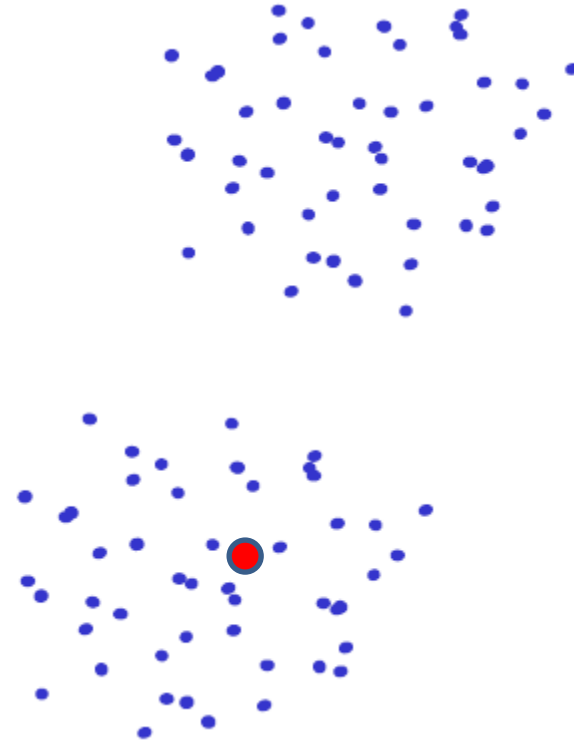
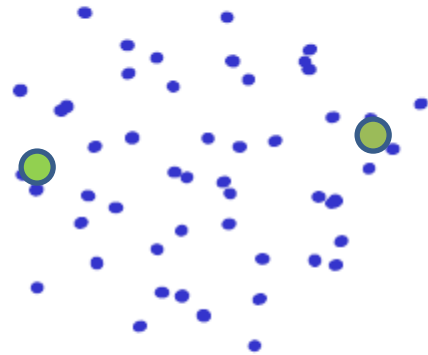
# Issues: Seed Choice

- Results are quite sensitive to seed selection.



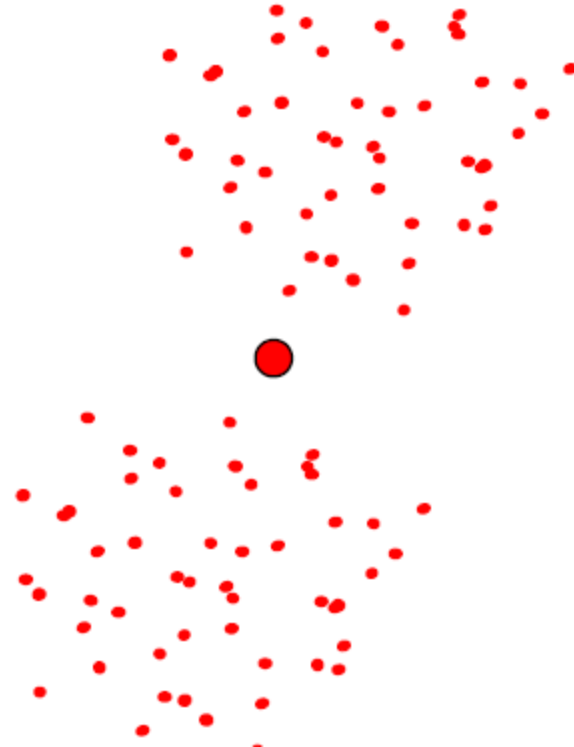
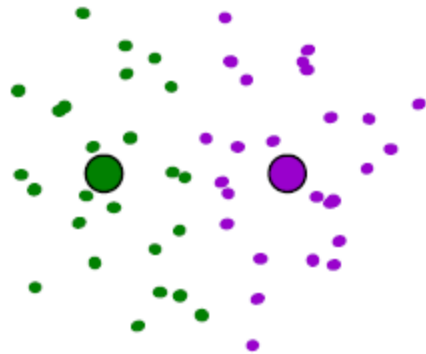
# Issues: Seed Choice

- Results are quite sensitive to seed selection.

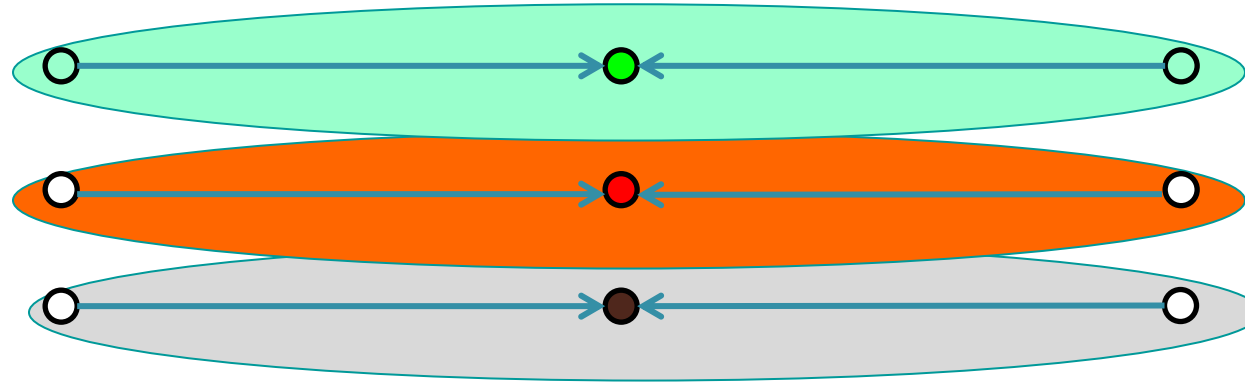


# Issues: Seed Choice

- Results are quite sensitive to seed selection.



# Issues: Seed Choice



K-means always converges, but it may converge at a local optimum that is different from the global optimum, and in fact could be arbitrarily worse in terms of its objective.

# Issues: 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

- Shape of clusters
  - Assumes isotropic, equal variance, convex clusters
- Sensitive to Outliers
  - use K-medoids





# K-medoids

Use actual training point as cluster center (medoid) rather than mean

- More robust to outliers pulling the mean away
- Better interpretability, can “visualize” the medoid
- More work to compute medoid than mean

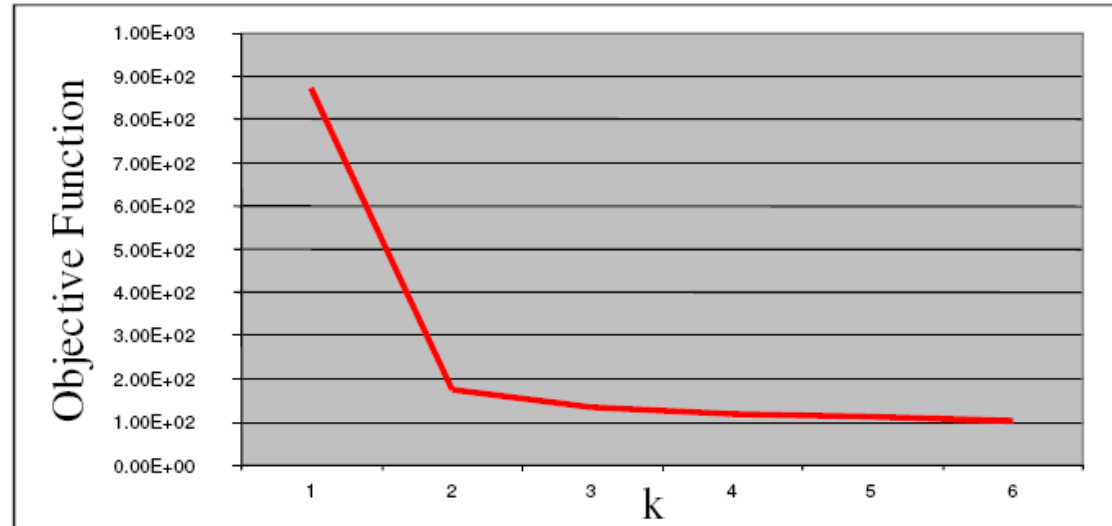
# Other Issues

- Number of clusters K

- Objective function

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

- Look for “Knee” in objective function



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

# K-means algorithm

- Optimize potential function:

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

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

**(1)** Fix  $\mu$ , optimize C      **Expected** cluster assignment

**(2)** Fix C, optimize  $\mu$       **Maximum** likelihood for center

Next lecture, we will see a generalization of this approach:

**EM algorithm**