

Generic Machine Learning

Clustering and Classification

TCTI-VKAAI-17: Applied Artificial Intelligence

1



K-Means Clustering Algorithm

Also called
'Centroids'

Amount of
clusters

- Input: K , set of points x_1, \dots, x_n
- Place centroids c_1, \dots, c_k at random locations
- Repeat until convergence:

Until nothing
changes

Distance (e.g., Euclidian)
between instance x_i and c_j

- For each point x_i :

- Find the nearest centroid c_j
- Assign the point x_i to cluster j

$$\arg \min_j \overbrace{D(x_i, c_j)}^{\text{Distance (e.g., Euclidian) between instance } x_i \text{ and } c_j}$$

- For each cluster $j = 1, \dots, K$:

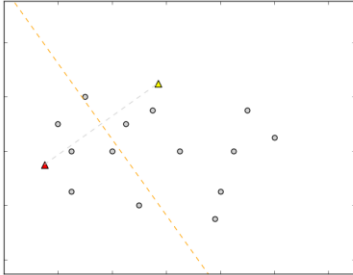
- Calculate new centroid $c_j = \text{mean of all points } x_i \text{ assigned to cluster } j$

$$c_j(a) = \frac{1}{n_j} \sum_{x_i \rightarrow c_j} x_i(a) \quad \text{for } a = 1, \dots, d$$

- Stop when none of the cluster assignments changes

15

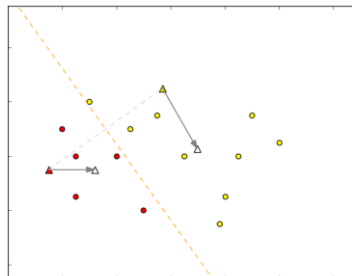
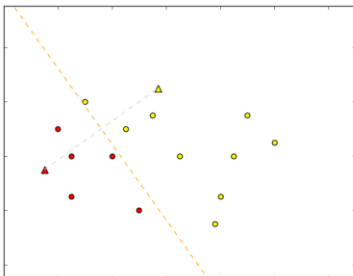
K-Means Example



- Input K , points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

16

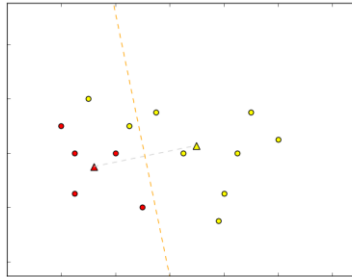
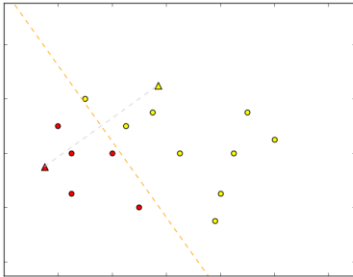
K-Means Example



- Input K , points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

17

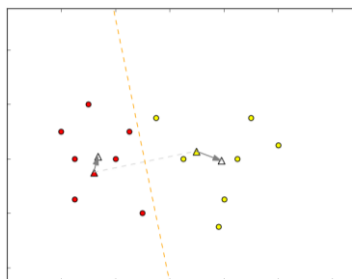
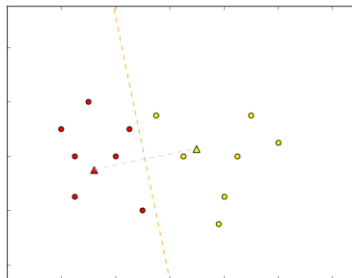
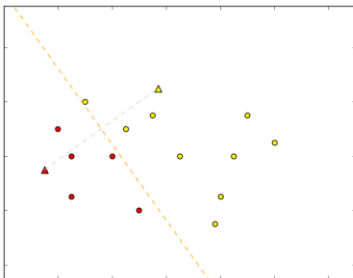
K-Means Example



- Input K points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

18

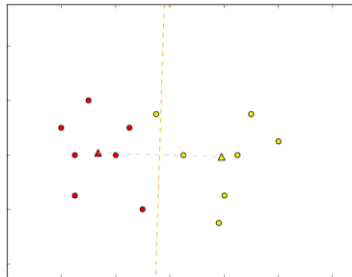
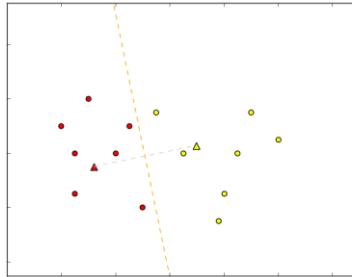
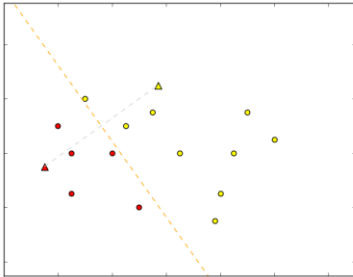
K-Means Example



- Input K points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

19

K-Means Example

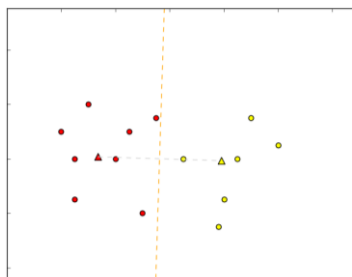
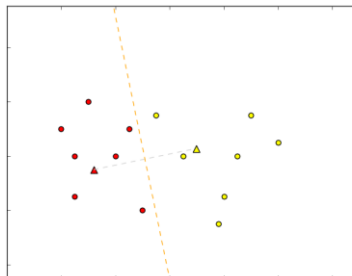
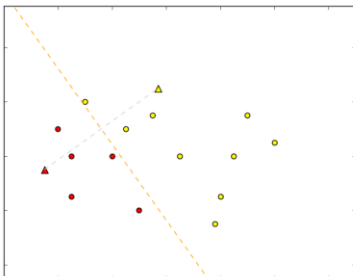


- Input K points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

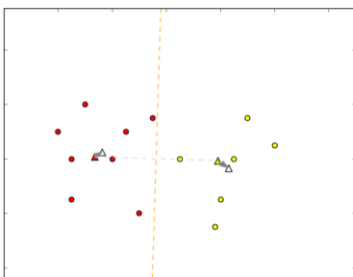


20

K-Means Example



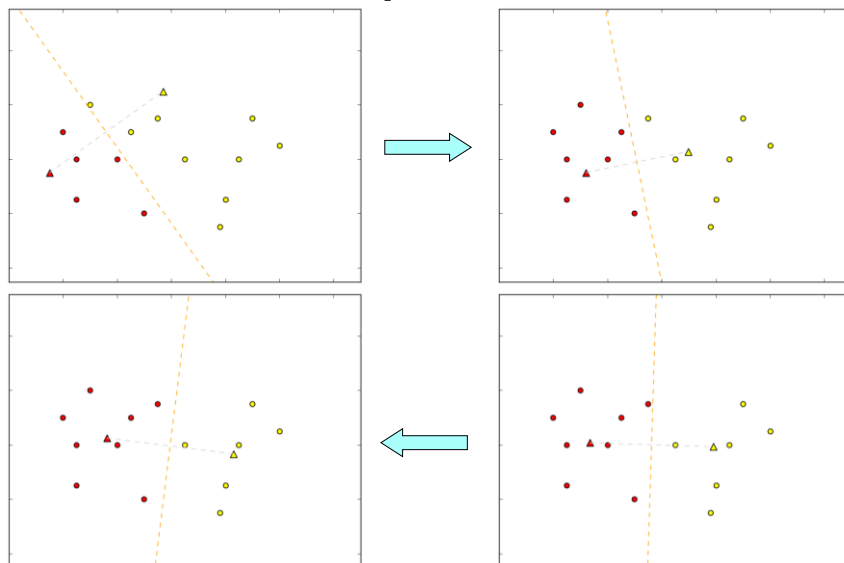
- Input K points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j



21



K-Means Example



- Input K points x_1, \dots, x_n
- Place centroids randomly
- Repeat until convergence:
 - For each point x_i :
 - Find nearest centroid c_j
 - Assign x_i to cluster j
 - For each cluster j :
 - Compute new centroid c_j

22



K-Means Properties

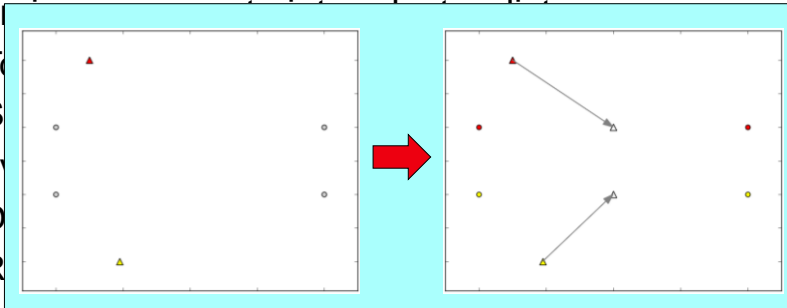
- Minimizes aggregate intra-cluster distance
 - Total squared distance from point to centre of its cluster
 - Same as variance if Euclidian distance is used
- Converges to local minimum
 - Different starting points \rightarrow very different results
 - Run several times with random starting points
 - Pick clustering that yields smallest aggregate distance
- Nearby points might not end up in the same cluster
 - The following clustering is a stable local minimum

23



K-Means Properties

- Minimize the sum of squared distances from each point to its assigned cluster center
 - To find the best clustering, you need to try many different initial cluster centers
 - Since there are many different initial cluster centers, you need to try many different initial cluster centers
- Compute the sum of squared distances from each point to its assigned cluster center
 - Distance from each point to its assigned cluster center
 - Repeat the process until the sum of squared distances is minimized
 - Pick clustering that yields smallest aggregate distance
- Nearby points might not end up in the same cluster
 - The following clustering is a stable local minimum

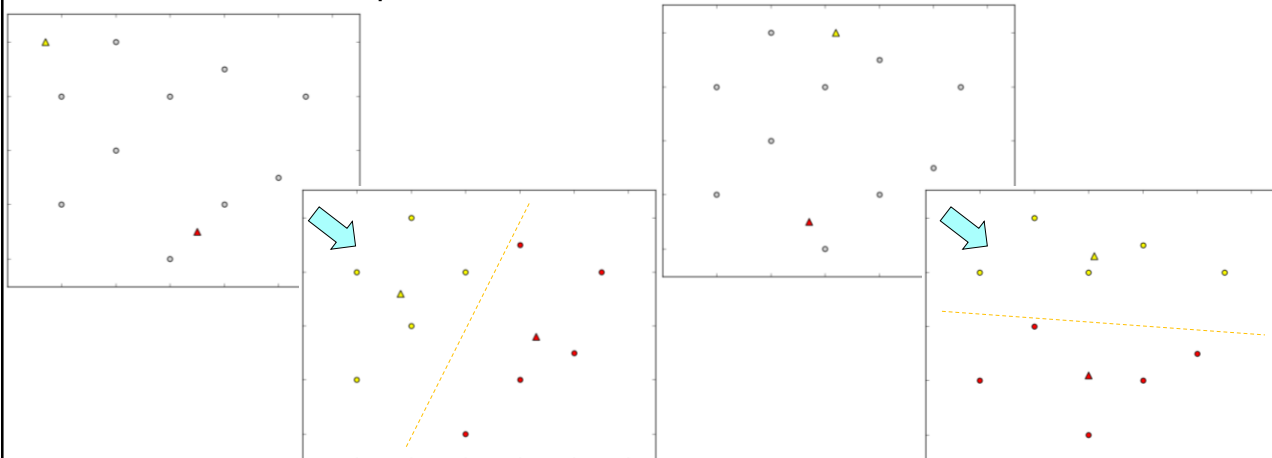


24



Convergence

- K-Means converges to a **local minimum**
 - Solution depends on initial values of K's





How many clusters?

- How many clusters are there in the data?

- Class labels may suggest the value of K (e.g., digits 0...9)

- Optimize distance V : for $K = 2, 3, \dots$

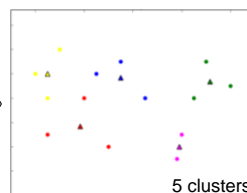
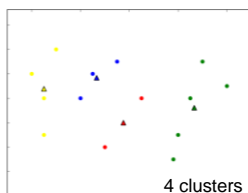
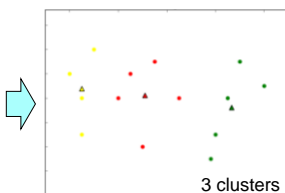
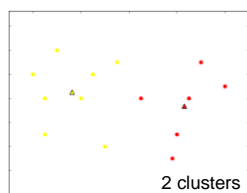
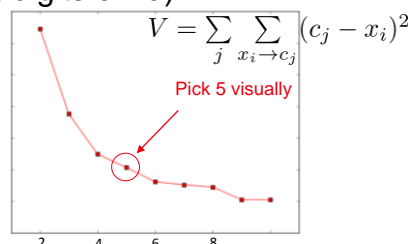
- Run K-Means, record distance

- Problem: V minimized when $K = n$

- What if we use a validation set?

- Visually from scree plot

- Point where 'mountain' ends, 'rubble' begins



26

Conclusion



- Classification and Clustering appear similar

- Classification: supervised categorization of observations

- Clustering: unsupervised determination of (potential) categories

- k Nearest Neighbors simple and effective

- But very costly to run effectively

- K-Means also rather simple and effective (and fast!)

- Some unwanted results due to randomness of chosen centroids

- Meaning one has to calculate several times for a data set → more costly

27