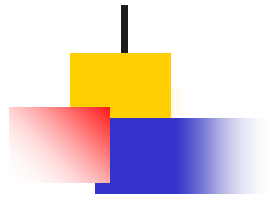


Fundamentals of Machine Learning

UNSUPERVISED LEARNING - CLUSTERING

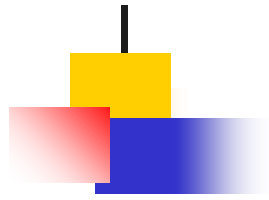
Amit K Roy-Chowdhury

Acknowledgments: Adapted from slides at <https://probml.github.io/pml-book/teaching1.html> by Prof. Saw Shier Nee
Modified by Amit Roy-Chowdhury, UCR.



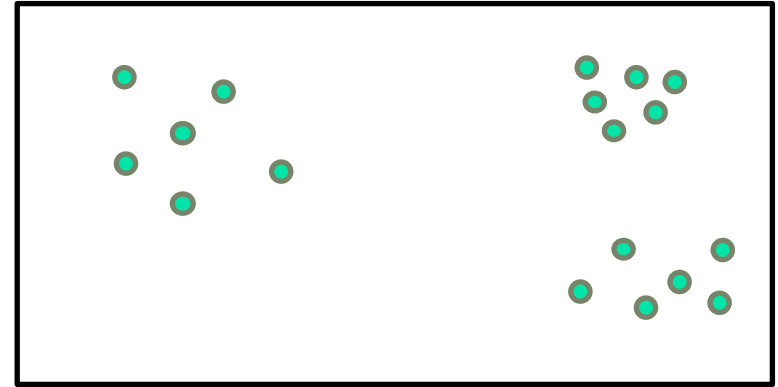
Outline

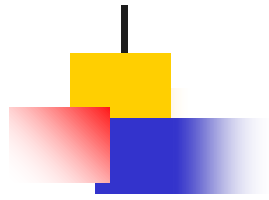
- Dimensionality Reduction methods (PCA) – already covered
- Clustering algorithms – K-means
- Gaussian Mixture Models – EM Algorithm



Clustering

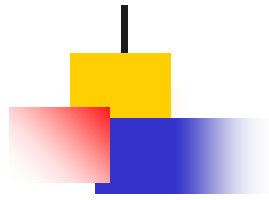
- Unsupervised Learning
- No Ground Truth
- Investigate data structure by grouping them into distinct groups.
- Advantage: useful when don't know what to look for
- Disadvantage: Subjective





Applications

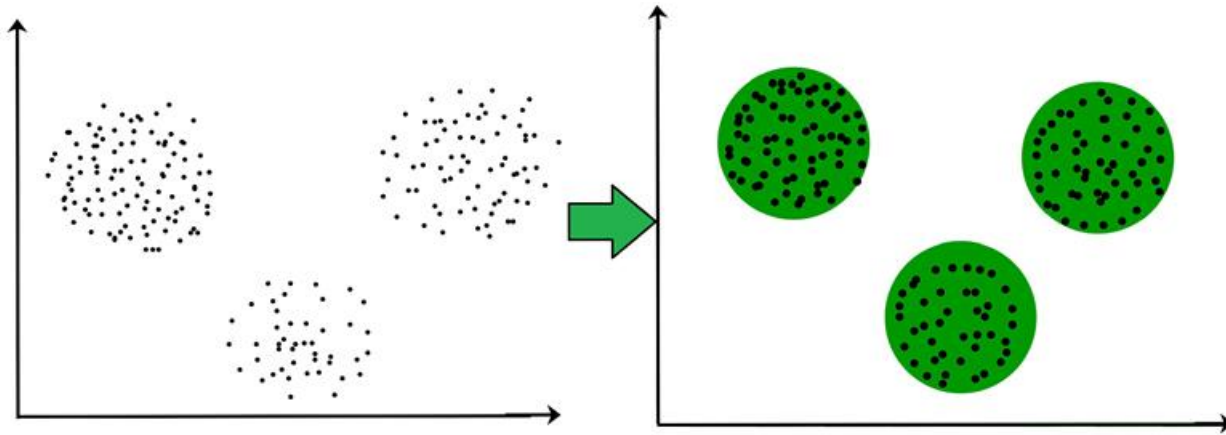
- Find set of representative samples that provide coverage of the data
 - Identify fake news based on the content
 - Group emails
- Exploratory Insight
 - Customer shopping patterns
 - Groups of genes/proteins with similar function

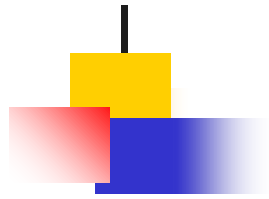


Clustering

Basic idea: group similar instances together

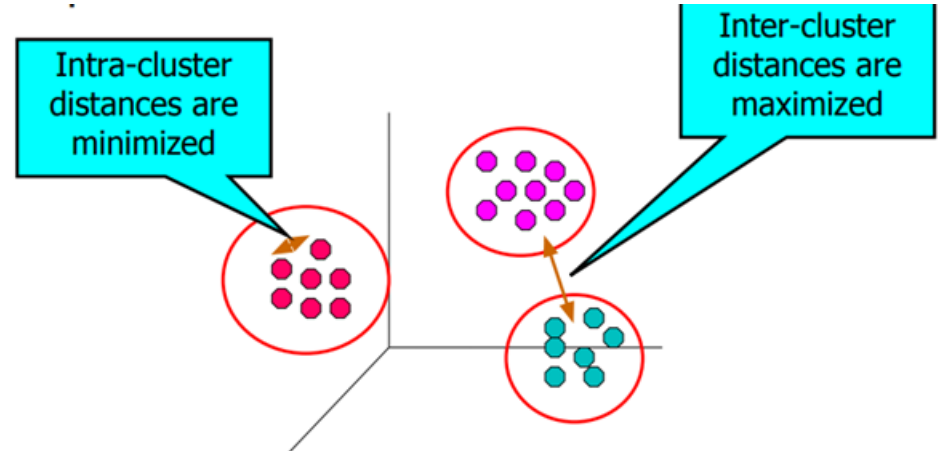
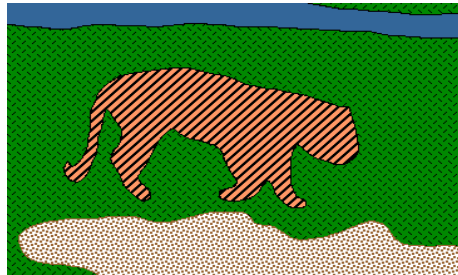
Example: 2D point patterns



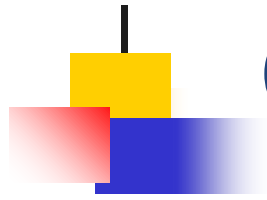


Clustering

- Basic idea: group similar instances together
- What does “similar” mean?
 - Based on Euclidean distance (squared),
or
 - Measure of similarity (or distance) between “points” to be clustered



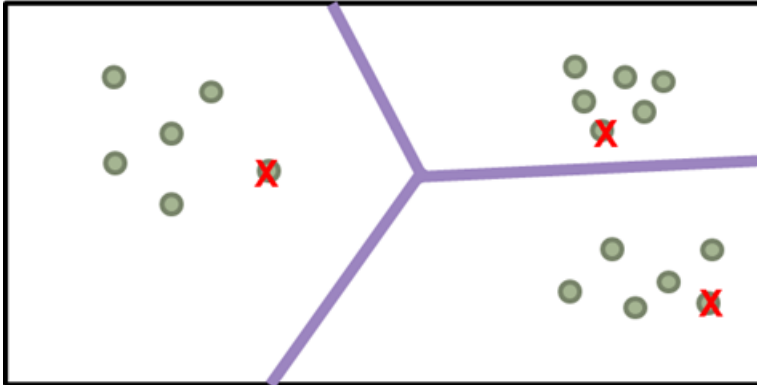
Gestalt Laws seek to formalize this proximity, similarity, continuation, closure, common fate



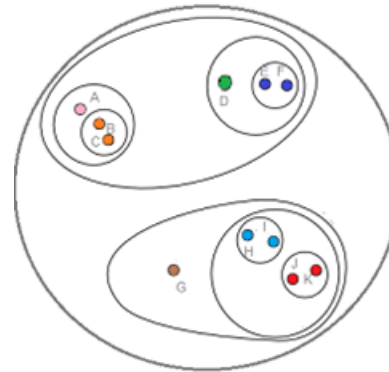
Clustering

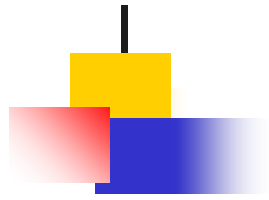
- Partitional - data in non-overlapping subsets. One data object is in one subset
- Hierarchy - data are in nested clusters, organized in a hierarchical tree

Partitional



Hierarchy





Kmeans Algorithm

Also known as Lloyd's algorithm

Input: x_1, x_2, \dots, x_n , given $x \in \mathbb{R}^d$

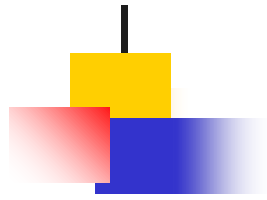
Output: 'Centers', $\mu_1, \mu_2, \dots, \mu_k$, given $\mu \in \mathbb{R}^d$

Assign each data point to its closest centers: $z_n^* = \arg \min_k \|x_n - \mu_k\|_2^2$

Compute the cluster centers:

$$\mu_k = \frac{1}{N_k} \sum_{n: z_n = k} x_n$$

Iterate between them till convergence – closely related to Expectation Maximization



Summary - Kmeans

A particularly simple method for clustering is K-means, which is identical to the generalized Lloyd algorithm we know from vector quantization, just applied to clustered data.

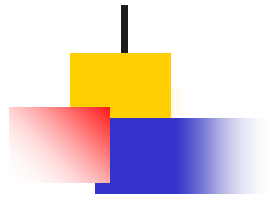
The idea is to represent each cluster k by a center point c_k and assign each data point x_n to one of the clusters k , which can be written in terms of index sets C_k .

The center points and the assignment are then chosen such that the mean squared distance between data points and center points is minimized

$$\text{objective function} \leftarrow J = \sum_{j=1}^k \sum_{i=1}^n \|x_i^{(j)} - c_j\|^2$$

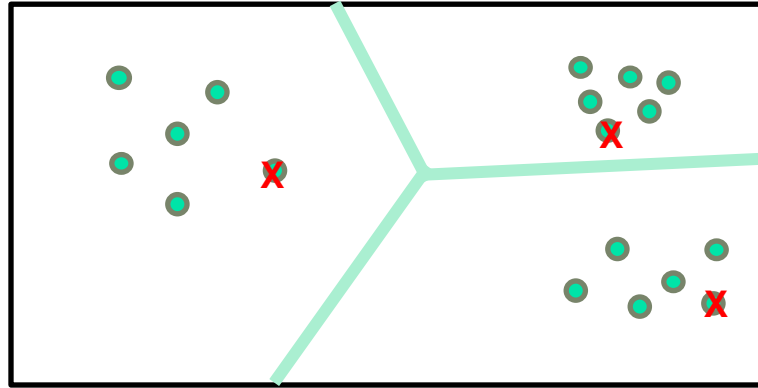
Diagram illustrating the K-means objective function J :

- k : number of clusters
- n : number of cases
- i : case i
- j : centroid for cluster j
- $\|x_i^{(j)} - c_j\|^2$: Distance function (squared distance between data point $x_i^{(j)}$ and centroid c_j)

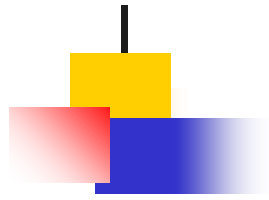


K-means - Intuition

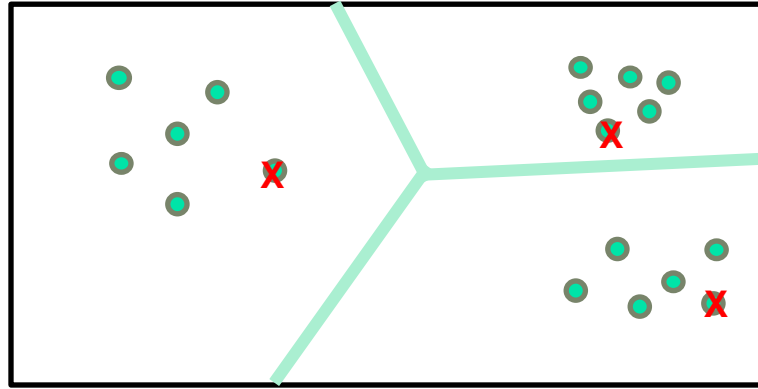
Say we want three clusters, first initialize random centre



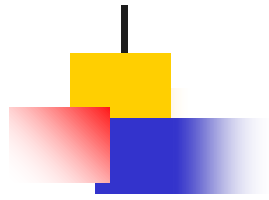
Generate optimal partition



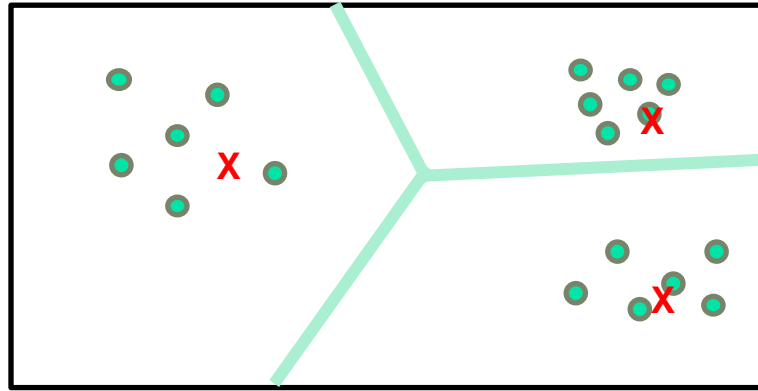
K-means - Intuition



Update the centre by computing the mean of coordinates in their respective region.



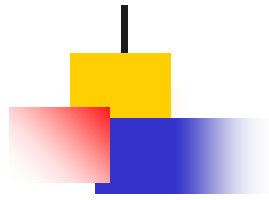
K-means - Intuition



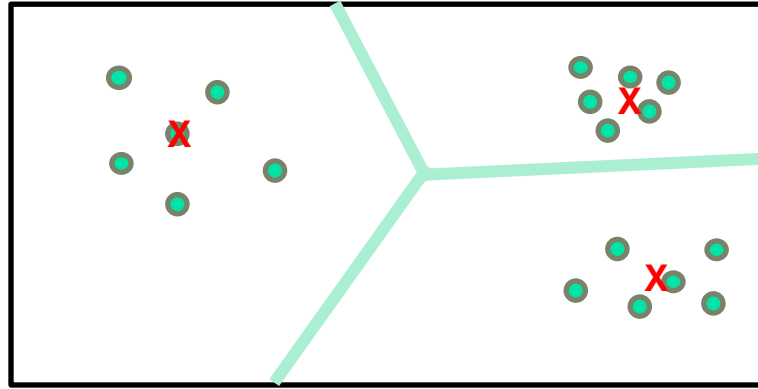
Updated centre

Compute Euclidean distance between points to the closest centres.

Identify the three regions



K-means - Intuition



Centre no longer changes (converge). Now, we have identify the three clusters.



K-means Algorithm

K-means clustering algorithm

1. Randomly initialize the cluster centers, c_1, \dots, c_K
2. Given cluster centers, determine points in each cluster
 - For each point p , find the closest c_i . Put p into cluster i
3. Given points in each cluster, solve for c_i
 - Set c_i to be the mean of points in cluster i
4. If c_i have changed, repeat Step 2

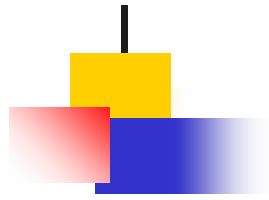
Properties

Will always converge to *some* solution

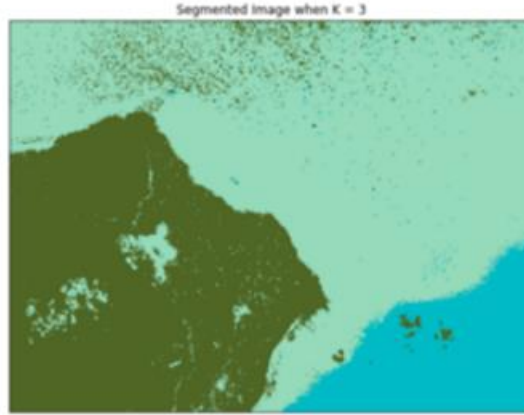
Can be a “local minimum”

- does not always find the global minimum of objective function:

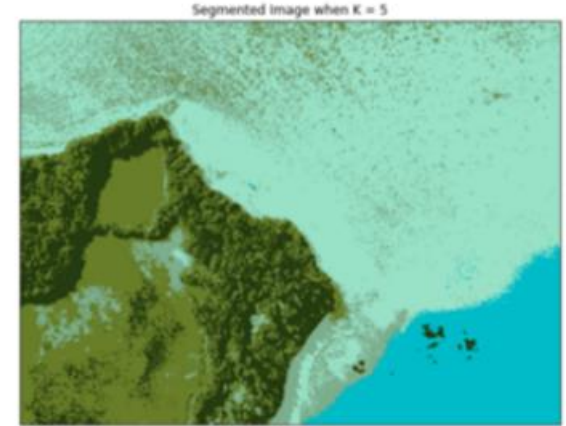
$$\sum_{\text{clusters } i} \sum_{\text{points } p \text{ in cluster } i} \|p - c_i\|^2$$



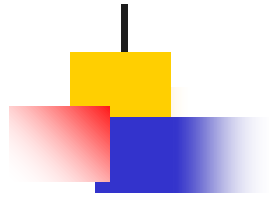
Kmeans for Segmentation



K=3

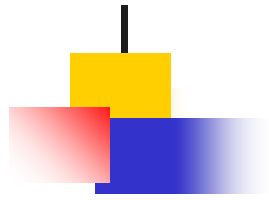


K=5

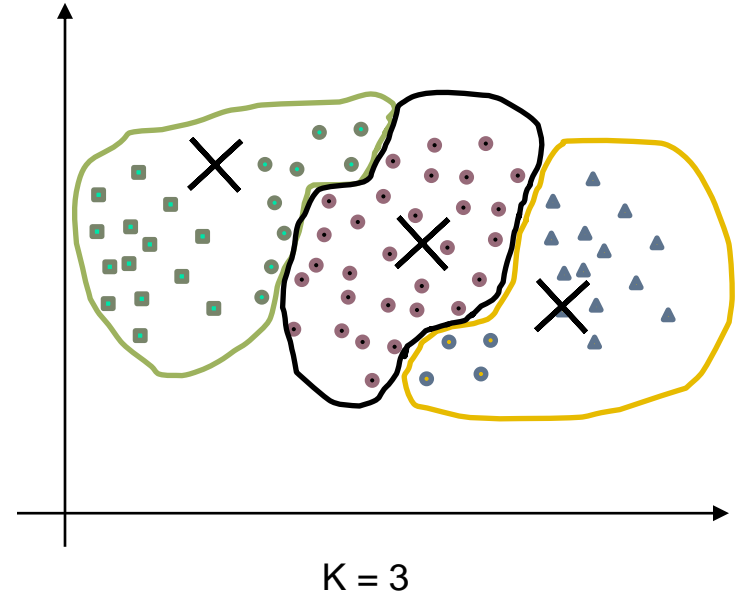
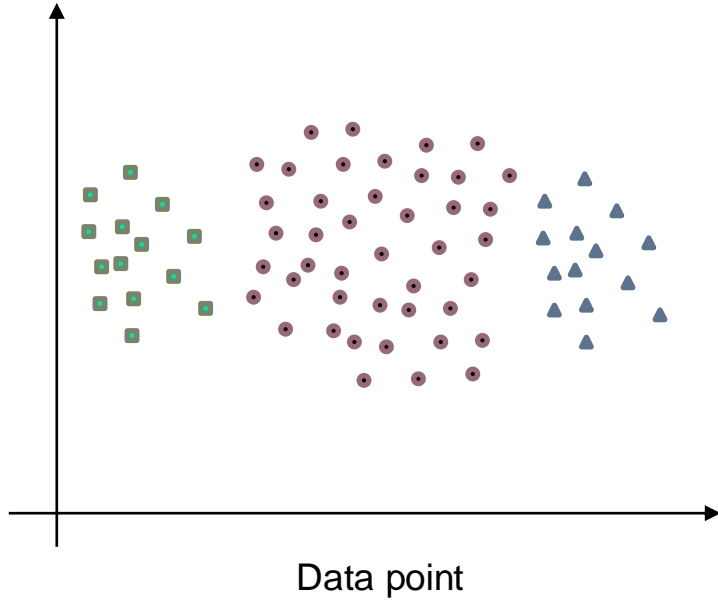


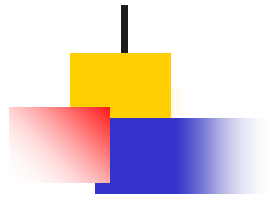
Limitation

- K-means has problems when clusters are of different
 - Sizes
 - Densities
 - Non-globular shapes
- K-means has problems when the data contains outliers.

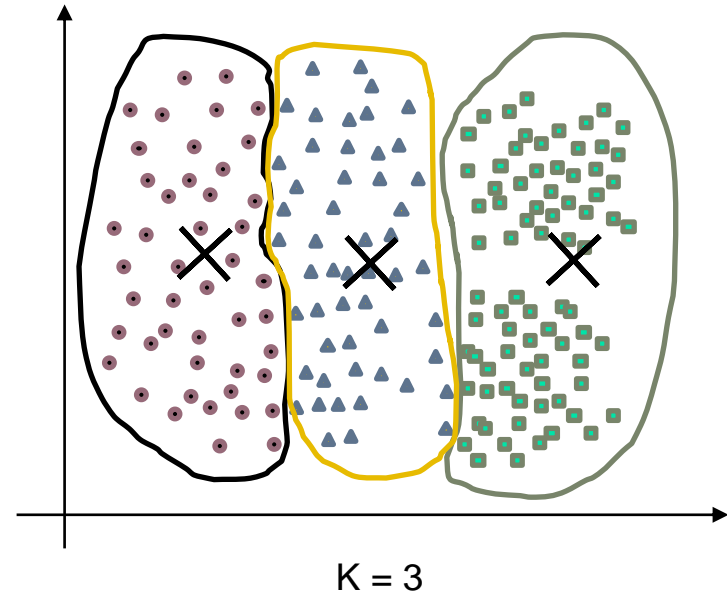
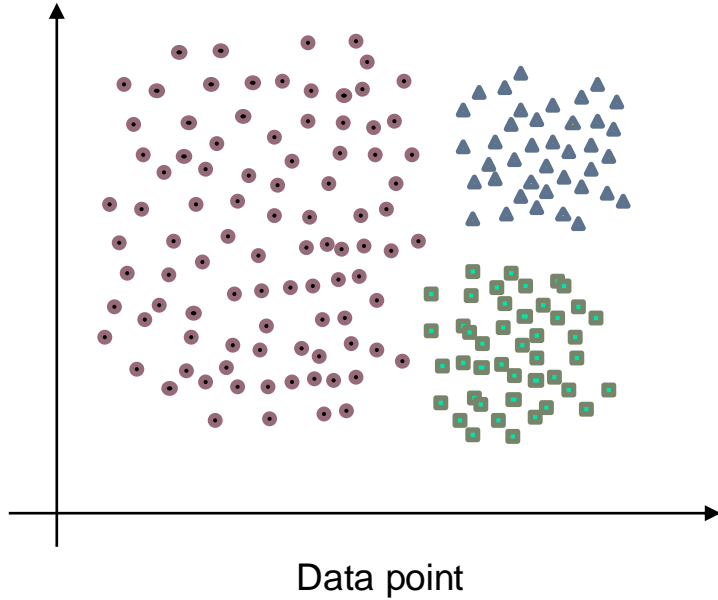


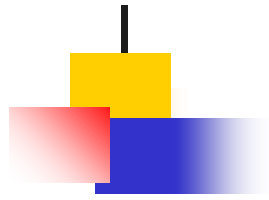
Limitation – Different Size



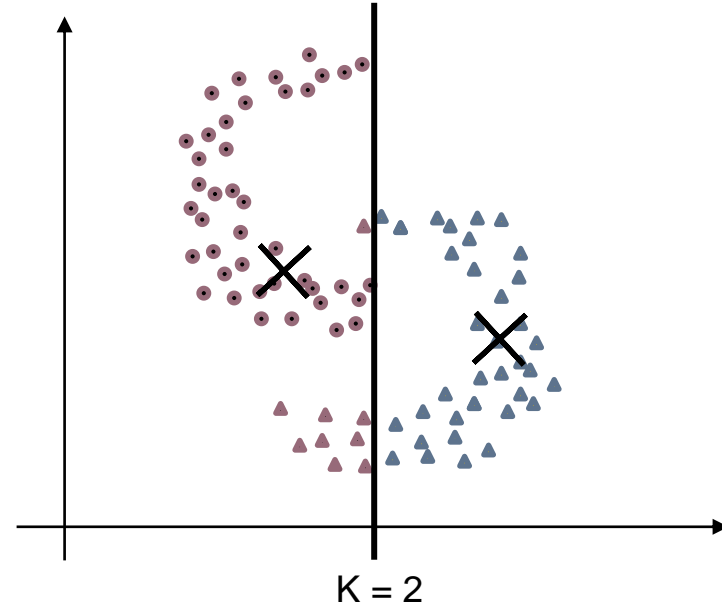
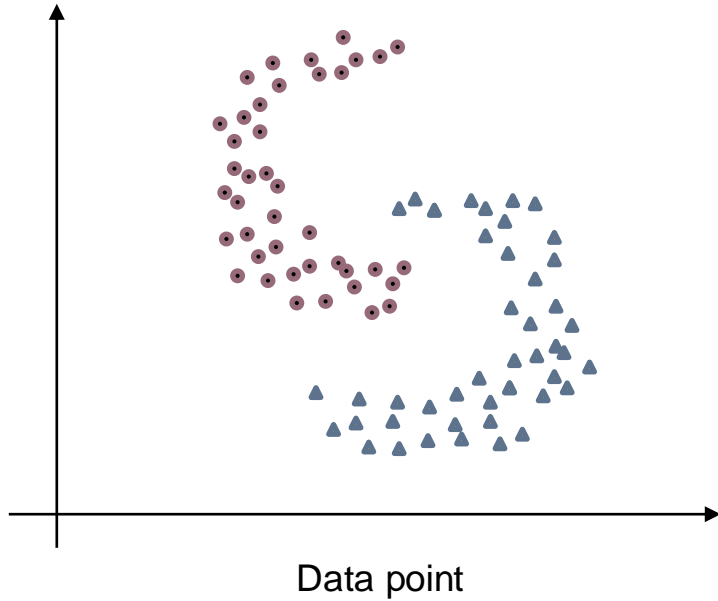


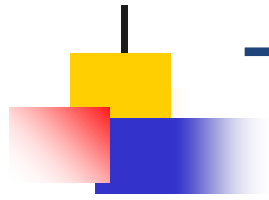
Limitation – Different Density





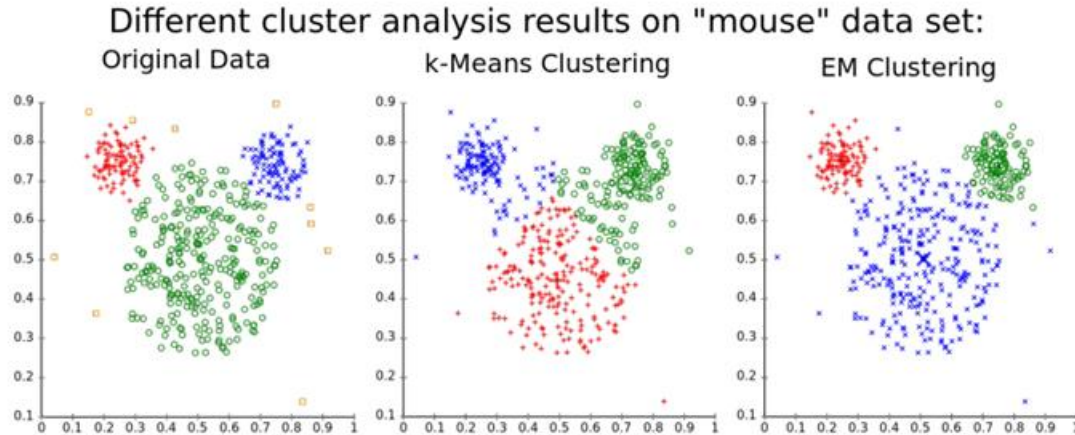
Limitation – non globular shape





To overcome k-means limitation

- Generate more clusters. Then combine similar cluster
- EM clustering using Gaussian Mixture Model
 - Intuition: Assume that the data generating process is a mixture of Multivariate Normals



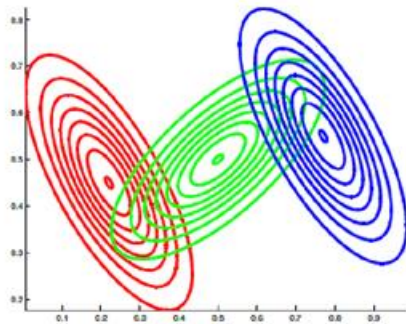
[source: Wikipedia, Public Domain image]

Gaussian Mixture Model/Mixture of Gaussian

$$p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k p_k(\mathbf{y})$$

where p_k is the k 'th mixture component, and π_k are the mixture weights which satisfy $0 \leq \pi_k \leq 1$ and $\sum_{k=1}^K \pi_k = 1$.

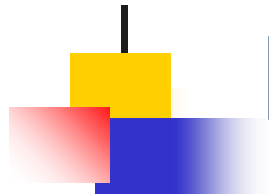
$$p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{y}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$



A mixture of 3 Gaussians in 2d.



A surface plot of the overall density.



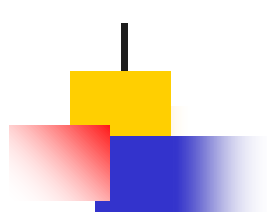
Expectation Maximization

$$\mathcal{D} = \{y_n : n = 1 : N\}$$

Need to find $\text{MLE } \hat{\theta} = \operatorname{argmax} \log p(\mathcal{D}|\theta)$

Approach:

1. E step: given current parameters, compute ownership of each point (hidden variable)
2. M step: given ownership probabilities, update parameters to maximize likelihood function
3. repeat until convergence



Expectation Maximization - Details

E Step: responsibility of cluster k for generating data point n

$$r_{nk}^{(t)} = p^*(z_n = k | \mathbf{y}_n, \boldsymbol{\theta}^{(t)}) = \frac{\pi_k^{(t)} p(\mathbf{y}_n | \boldsymbol{\theta}_k^{(t)})}{\sum_{k'} \pi_{k'}^{(t)} p(\mathbf{y}_n | \boldsymbol{\theta}_{k'}^{(t)})}$$

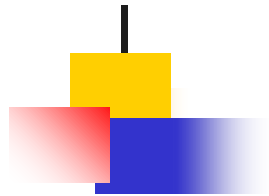
M Step: maximizes log likelihood

$$\boldsymbol{\mu}_k^{(t+1)} = \frac{\sum_n r_{nk}^{(t)} \mathbf{y}_n}{r_k^{(t)}}$$

$$\boldsymbol{\Sigma}_k^{(t+1)} = \frac{\sum_n r_{nk}^{(t)} (\mathbf{y}_n - \boldsymbol{\mu}_k^{(t+1)}) (\mathbf{y}_n - \boldsymbol{\mu}_k^{(t+1)})^T}{r_k^{(t)}}$$

$$\pi_k^{(t+1)} = \frac{1}{N} \sum_n r_{nk}^{(t)} = \frac{r_k^{(t)}}{N}$$

$r_k^{(t)} \triangleq \sum_n r_{nk}^{(t)}$ is the weighted number of points assigned to cluster k

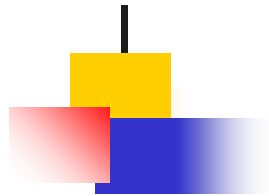


Relation to K-Means

we fix $\Sigma_k = \mathbf{I}$ and $\pi_k = 1/K$ for all the clusters (so we just have to estimate the means μ_k)

we approximate the E step, by replacing the soft responsibilities with hard cluster assignments

$$z_n^* = \operatorname{argmax}_k r_{nk}$$



EM clustering

