

Making Sense of Unstructured Data

Week 3

Learning Objectives of the Session

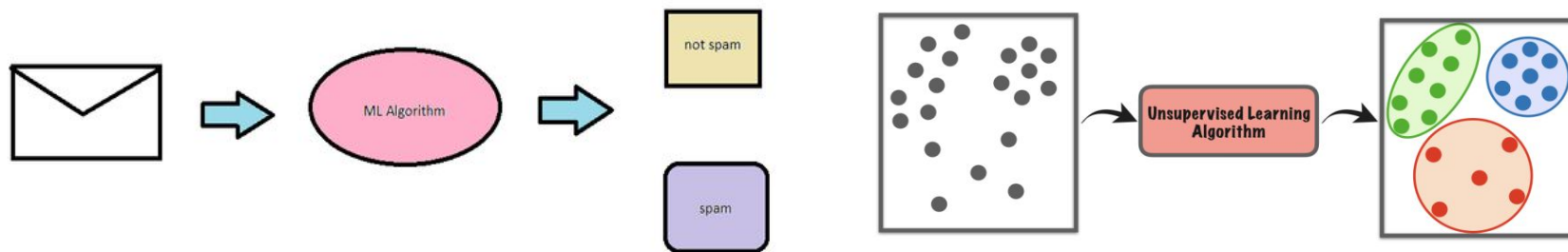
- Introduction to Supervised and Unsupervised learning.
- Clustering
- K-Means clustering algorithm
- K- Medoids
- GMM
- Spectral clustering
- Principal Component Analysis

Discussion Questions

1. What is the difference between Supervised learning and Unsupervised learning?
2. What is clustering and what are the most common clustering techniques?
3. Why and when to use clustering?
4. How does clustering help in finding hidden groupings and patterns?
5. What is K-Means clustering?
6. Why is K-Means clustering so popular? What are its assumptions?
7. Why converting features into continuous values and scaling the data is important while performing K-Means clustering?
8. How to determine the value of K in K-Means clustering
9. What is Spectral clustering?
10. How does PCA help in Dimensionality reduction?

Supervised vs. Unsupervised Learning

Supervised learning algorithms are trained using labeled data. Example - Predict whether a new email is SPAM or NOT SPAM, predict whether a customer will churn or not etc.



Unsupervised learning algorithms are trained using unlabeled data. It learns on itself and finds patterns in the data to form different groups within the dataset. For example: Customer segmentation, image segmentation etc.

What is clustering?

Cluster analysis, or clustering, is an unsupervised machine learning technique. It involves automatically discovering natural grouping in data (groups data according to the notion of similarity).

The most common types of clustering are K-means clustering, LDA Clustering, Spectral Clustering, Modularity Clustering, Hierarchical Clustering, DBSCAN clustering etc.

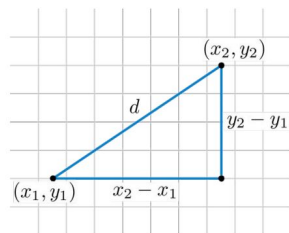
Unsupervised Learning: Clustering

Cluster analysis, or clustering, is an unsupervised machine learning technique. It involves automatically discovering natural grouping in data (groups data according to the notion of similarity).

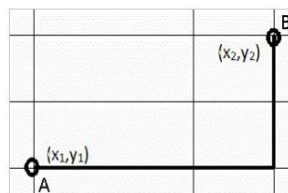
Notion of Similarity: In Clustering we can find similarity between two data points and group them together based on the distance between them.

There are different distance measures which can be used to find similarity, some of them are:

1. Euclidean distance
2. Manhattan distance
3. Chebyshev distance

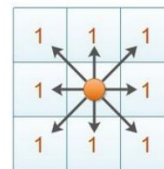


Euclidean distance



$$D_n = |X_2 - X_1| + |Y_2 - Y_1|$$

Manhattan distance



$$\max(|x_1 - x_2|, |y_1 - y_2|)$$

Chebyshev distance

Why and When to Use Clustering

Why

Retail, finance, and marketing are some of the key domains that use clustering methods to analyze their data. This can help them in gaining further insights into their customers.

The factors analyzed through clustering can have a big impact on sales and customer satisfaction, making it an invaluable tool to improve overall business performance.

When

- When you are starting with a large and unstructured dataset.
- When you do not know how many or which classes your data is divided into.
- When annotating your data can be very expensive.



How does clustering help in finding hidden groupings and patterns?

Clustering methods simply try to group similar patterns into clusters whose members are more similar to each other (according to some **distance measure**) than to members of other clusters.

K-Means Clustering

K-Means Clustering is an iterative **algorithm** that divides the unlabeled dataset into **k** different **clusters** in such a way that each point in the dataset belongs to only one group that has similar properties.

The algorithm starts with initial estimates for the **K** centroids, which can either be randomly generated or randomly selected from the data set.

The algorithm then iterates between two steps:

1. Data assignment step:

Each centroid defines one of the clusters. In this step, each data point based on the squared Euclidean distance is assigned to its nearest centroid. If c_i is the collection of centroids in set C , then each data point x is assigned to a cluster based on where $\text{dist}(\cdot)$ is the standard (L2) Euclidean distance. **Min dist(C,x)**

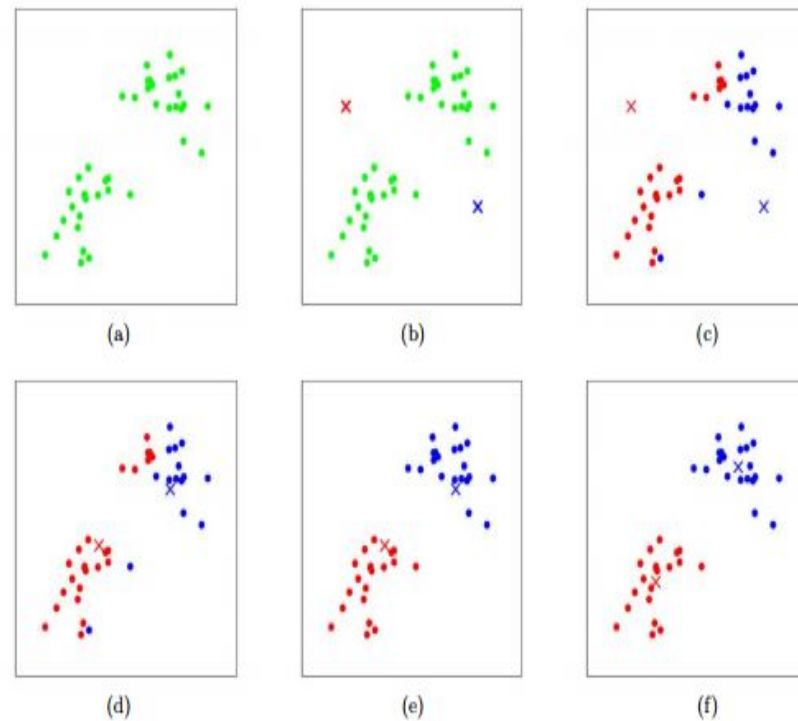
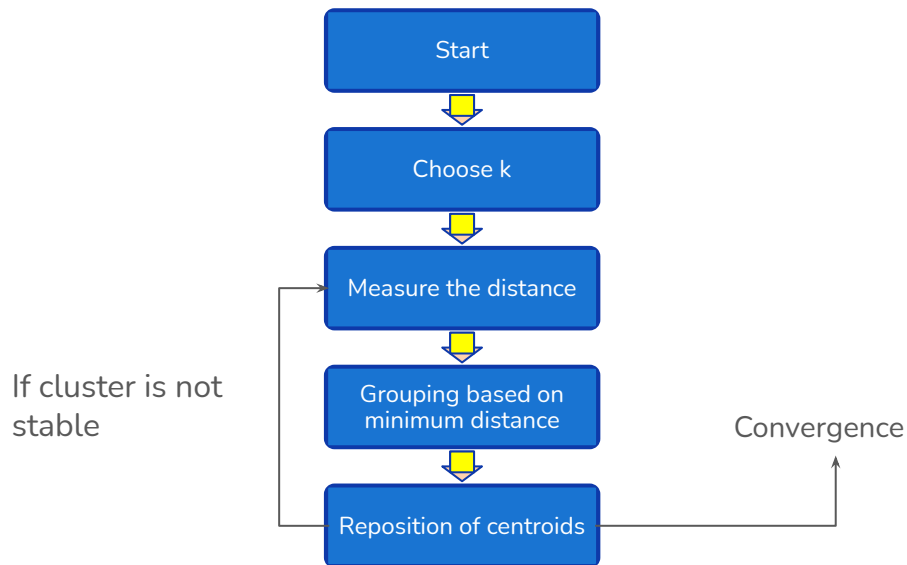
2. Centroid update step:

Centroids are recomputed by taking the mean of all data points assigned to that centroid cluster.

The algorithm iterates between step 1 and 2 until a stopping criteria is met.

This algorithm may converge on a local optimum. Assessing more than one run of the algorithm with randomized starting centroids may give a better outcome.

K-Means Clustering Steps



Advantages and Disadvantages of using K-Means Clustering

Advantages:

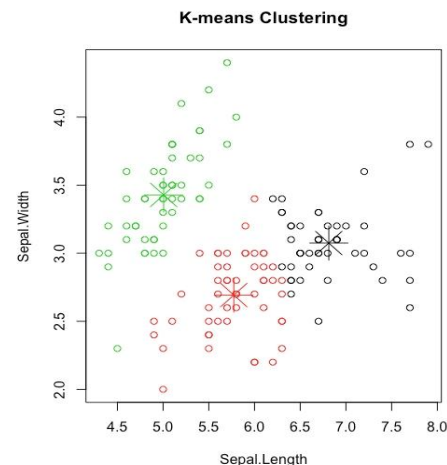
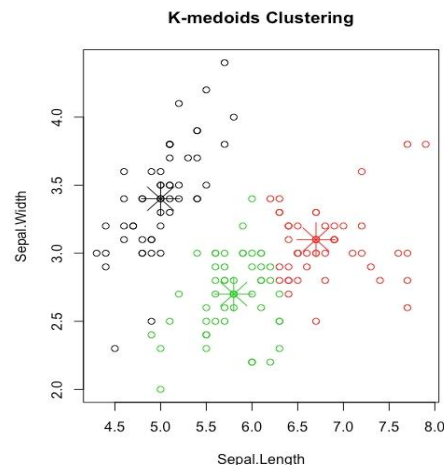
- K-means is relatively simple to implement
- It scales to large datasets
- It guarantees convergence
- It can easily adapt to new examples

Disadvantages:

- It is difficult to identify the value of K
- k-means has trouble clustering data where clusters are of varying sizes and density
- It can easily get affected by outliers
- It assumes the cluster shape to be spherical in nature and does not perform well on arbitrary data
- It depends on the initial values assigned to the centroids and gives different results for different initializations

Alternative to K Means - PAM (K Medoids) clustering

- The problem with K means is that the final centroids are not interpretable i.e. centroids are not actual points but the means of the points present in the cluster.
- The idea behind K Medoids clustering is to make the final centroids as actual data points so that they are interpretable.
- In K Medoids, we only change one step from K Means which is to update the centroids. In this process if there are m points in a cluster, swap the previous centroids with all other $(m-1)$ points from the cluster and finalize the point as new centroid which has minimum loss.
- Because of this, unlike K Means it is robust to outliers and converges fast.
- You can see in this image that the centroids in K Medoids are the actual data points represented as the cross, unlike K Means.



Expectation maximization in GMM Clustering

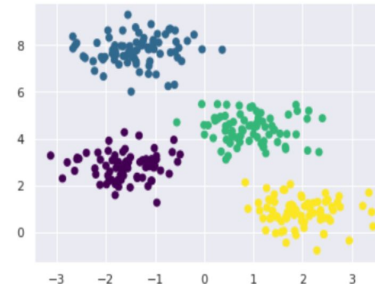
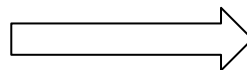
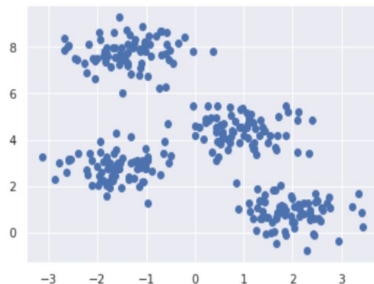
In GMM, we need the parameters of each Gaussian (variance, mean etc.) in order to cluster our data but we need to know which sample belongs to what Gaussian in order to estimate those very same parameters.

That is where we need EM algorithm. There are two steps involved in this algorithm:

1. **The E-step:** It estimates the probability that a given observation to be in a cluster/distribution. This value will be high when the point is assigned to the right cluster and lower otherwise.
2. **The M-step:** In this step we want to maximize the likelihood that each observation came from the distribution

After that we reiterate these two steps and updates the probabilities of an observation to be in a cluster.

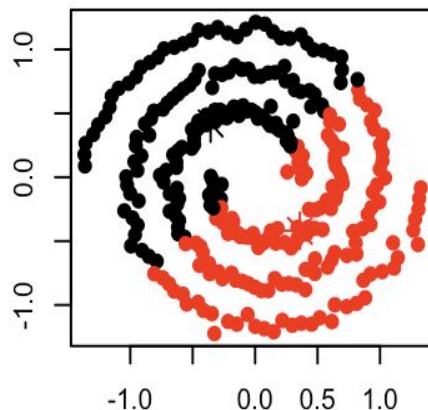
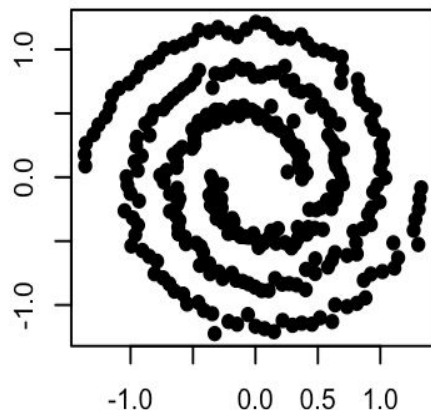
Example of
GMM clustering



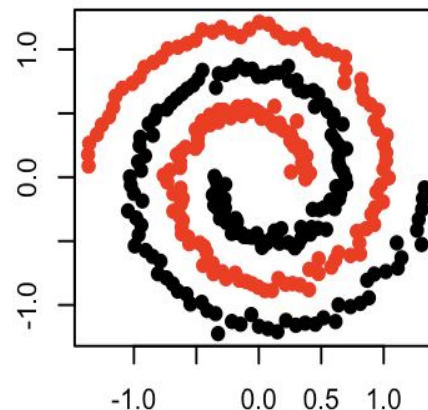
Spectral Clustering

In spectral clustering, **data points are treated as nodes of a graph**. Thus, spectral clustering is a graph partitioning problem. The nodes are then mapped to a low-dimensional space that can be easily segregated to form clusters. No assumption is made about the shape/form of the clusters. **The goal** of spectral clustering is to cluster data that is connected but not necessarily compact or clustered within convex boundaries.

K-means



Spectral clustering

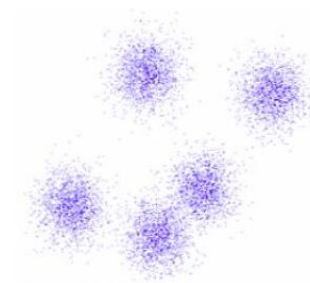


[Image Source](#)

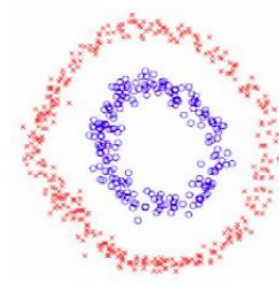
K-Means Vs Spectral Clustering

- **Compactness:** Points that lie close to each other fall in the same cluster and are compact around the cluster center. The closeness can be measured by the distance between the observations. E.g.: K-Means Clustering
- **Connectivity:** Points that are connected or immediately next to each other are put in the same cluster. Even if the distance between 2 points is less, if they are not connected, they are not clustered together. Spectral Clustering is a technique that follows this approach. K-means will fail to effectively cluster these, even when the true number of clusters K is known to the algorithm.

K-means, as a *data-clustering* algorithm, is ideal for discovering globular clusters where all the members of each cluster are in close proximity to each other (in the Euclidean sense)



Compactness



Connectivity

[Image Source](#)

Case Study Clustering

Dimensionality reduction

What is Dimensionality Reduction?

- Dimensionality reduction is the process to reduce the number of dimensions in the feature space.

Need for Dimensionality Reduction

- In the machine learning, we tend to add many features to get more accurate results. However, after a certain point the performance and robustness of the model starts decreasing and computational complexity starts increasing as we increase the number of features. This is called curse of dimensionality where the sample density decrease exponentially with the increase of dimensionality.
 - We use dimensionality reduction to transform the data into low dimensions while keeping most of the information intact.
 - It also helps us to visualize the high dimensional data to 2D & 3D.

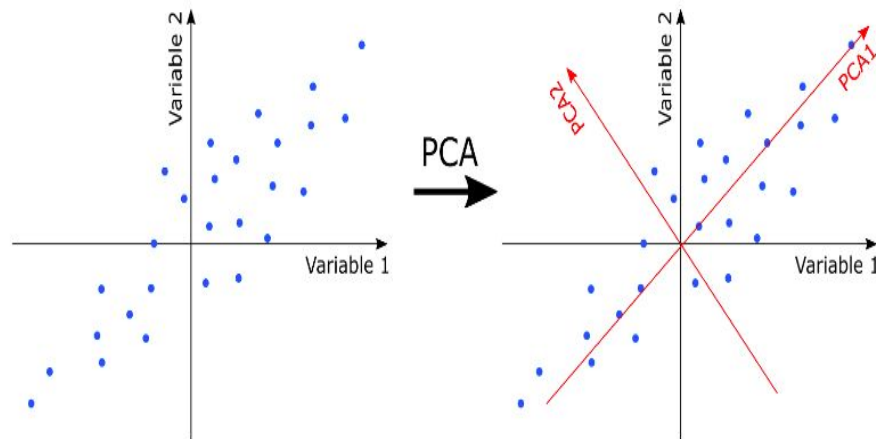
Principal Component Analysis

Principal Component Analysis, or PCA, is a method for reducing the dimensionality of data. This is used to transform your data into a new dimensional space by projecting the data on new axes. These axes are called Principal components

The selection of the principal components is such that they retain the maximum variation present in the original variables on the first principal component and the variation decreases as we move down the order. All Principal components are orthogonal axes to each other.

Steps for PCA:

- Begin by standardizing the data.
- Generate the covariance matrix
- Perform eigenvalue decomposition
- Sort the eigen pairs in descending order
- Order and select the largest one.



Case Study - PCA



Happy Learning !

