

Table Of Contents

1. Concept of K-Means Clustering
2. Math Intuition Behind K-Means
3. Cluster Building Process
4. Edge Case Scenarios of K-Means
5. Challenges and Improvements in K-Means

Main Requirements for clustering algorithms

- ① It should be Scalable.
- ② It should be able to deal with attributes of different types.
- ③ It should be able to discover arbitrary shape clusters.
- ④ It should have an inbuilt ability to deal with noise & outliers.
- ⑤ The clusters should not vary with the order of input records.
- ⑥ It should be able to handle data of high dimension.
- ⑦ It should be easy to interpret & use.

10 K-means

old algorithm

Proposed in 1957 by Stuart Lloyd, also

In 1965 by Edward W. Fuzzy. 10

Sometimes called

Lloyd-Fuzzy Algo.

→ K-means is a clustering approach in which data is grouped into K distinct non overlapping clusters based on their distances from the K centres.

→ The value of 'K' needs to be specified first & then the algorithm assigns the points to exactly one cluster.

Theory of K-means

→ Let C_1, C_2, \dots, C_k be the k -clusters.

→ Then we can write

$$C_1 \cup C_2 \cup C_3 \cup \dots \cup C_k = \{1, 2, 3, \dots, n\} \rightarrow$$

$$\text{and also, } C_k \cap C_{k'} = \emptyset \rightarrow$$

i.e.,
Each data point
has been assigned
to a cluster

This means, clusters are non-overlapping.

→ The idea behind the K-means approach is that within-cluster variation amongst the points should be minimum.

The within-cluster variance is denoted by $\underline{w(C_k)}$

→ Hence, according to statement above,

we need to minimize the variance for all clusters mathematically,

$$\text{minimize}_{C_1, C_2, \dots, C_k} \left\{ \sum_{k=1}^K w(C_k) \right\}$$

WCSS →
within cluster
summation or variance

→ The next step is to define the criterion for measuring the within cluster variance.

Generally, the criterion is the Euclidean distance between two data points.

$$w(C_k) = \frac{1}{|C_k|} \sum_{i \in C_k} \sum_{j \in C_k} (x_{ij} - \bar{x}_{kj})^2$$

$$P = (1) w$$

→ the above formula says that we are
 ① calculating the distance b/w all the
 points in a cluster, ② then we are repeating
 it for all the k-clusters (so two summations)

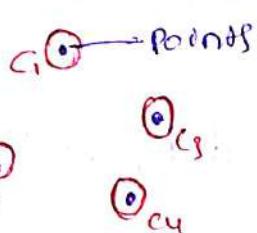
Eg:-

Let's take 4 points

case ① for

4 points

4 clusters



→ now, each point acts as a cluster and so if it is centroid or not.

∴ Distance b/w the point & centroid (same point) is zero (0).

→ And it is same for all the 4 clusters in this.

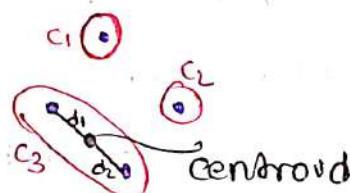
$$\text{e.g., } w(c_1) = 0 + 0 + 0 + 0 = 0 \\ c_1 \ c_2 \ c_3 \ c_4$$

(Summation of distances.)

case ② for

4 points

3 clusters



→ Now we need distance b/w Point & the centroid of cluster if present in.

$$w(c_3) = c_1 + c_2 + c_3$$

$$w(c_3) = 0 + 0 + (d_1 + d_2)$$

Individual

sum of distance in a cluster

$$\sum_{i=1}^3 (x - c)^2$$

Centroid

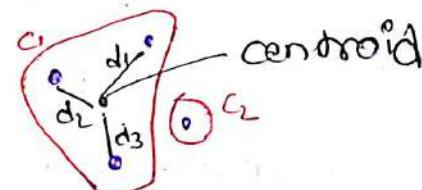
③ All the clusters sum

$$\sum_{i=1}^3 \sum_{j=1}^3 (x_j - c_i)^2$$

case ③ for

4 points

2 clusters



$$w(c_4) = c_1 + c_2$$

$$w(c_4) = (d_1 + d_2 + d_3) + 0$$

distance b/w point & centroid Summation of distance in a cluster.

Summation of the individual summations of distances in cluster.

case ④ for 4 points

1 cluster



$$w(c_1) = c_1$$

$$w(c_1) = (d_1 + d_2 + d_3 + d_4)$$

Goal is to minimize

$$w(c_1) + w(c_2) + w(c_3) + w(c_4)$$

↳ low

Variance
for all
clusters,

So, we got

$$w(c_4) = 0$$

$$w(c_3) = d_1 + d_2$$

$$w(c_2) = d_1 + d_2 + d_3$$

$$w(c_1) = d_1'' + d_2'' + d_3'' + d_4''$$

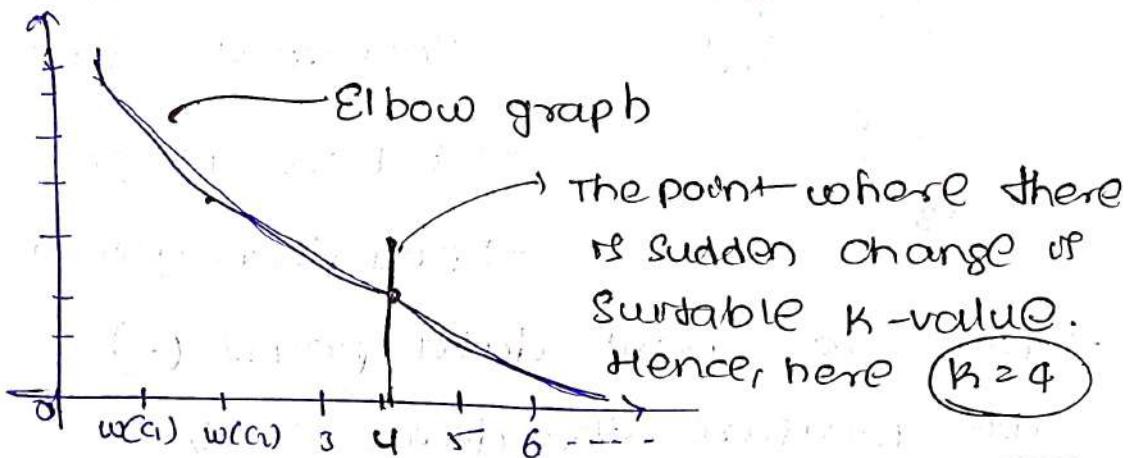
→ By observation, we
can say that

As no. of clusters
increases, $w(c_k)$
value decreases &
tends to '0'

$$\therefore w(c_k) \propto \frac{1}{\text{no. of cluster}(k)}$$

$$w(c_1) \succ w(c_2) \succ w(c_3) \succ \dots \succ w(c_n)$$

and the graph goes as -



within cluster summation or square (WCSS)

→ this is how we find the value of 'k'
in k-means

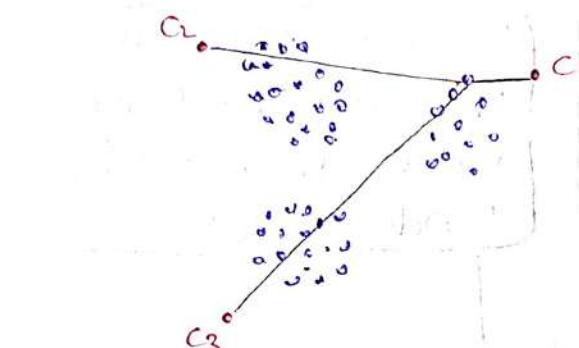
But, here is the problem in k-means,

* After selecting k value, it takes
k-points as centroids and starts
building the clusters normally.

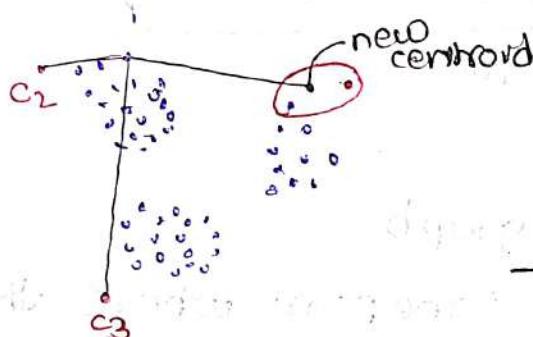
But the question is where do we
place those k-points?

* cluster building process

let's say we got: $K=3$



→ after, randomly placing the K points,



→ get taken a point from our data (\cdot)

and measured the distance to all the K -centroid point (\cdot),

and whichever distance is smaller, the point goes that cluster.

→ so, here it goes to c_1 .

→ and then gets a new centroid.

→ now for next data point (\cdot)

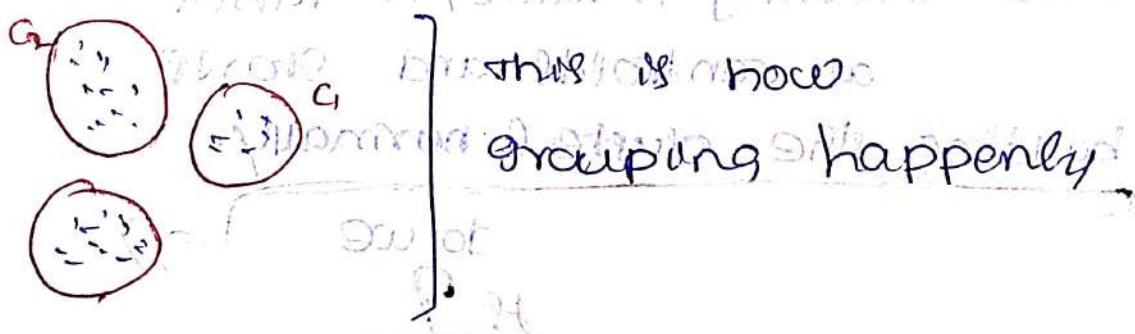
we measure the distance of

point to the new centroid

→ and in this way, cluster keeps building and centroid gets updated.

→ this process continues until

all the data points are finished.

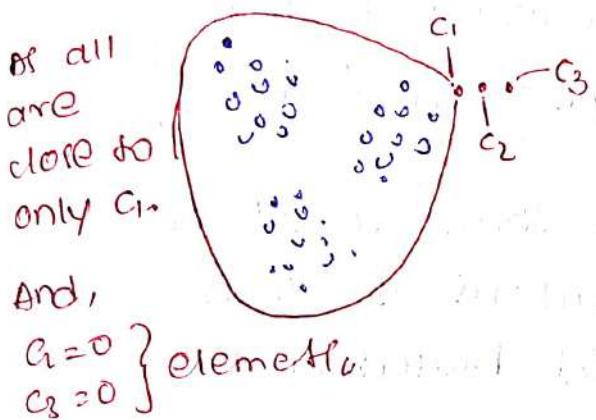


edge case scenario

problem with k-means

In the above example, we took $K=3$ and then it generated 3 random points and build cluster.

So, what if all the 3 points are close?



In this case,
all data points
would go to
cluster ① (c_1)

which is not what
we are expecting!!

Algorithm

- ① find wcss for k -groups possible ($\max K = n$)
- ② Build elbow graph, and understand k -value (point where there's distinction)
- ③ Randomly assign k -centres.
- ④ calculate distance of all points from all k -centres and allocate the points to cluster based on shortest distance to "cluster centroid".
- ⑤ The model's inertia is the mean squared distance b/w each instance & closest centroid.
- ⑥ Goal is to have model with low inertia.
- ⑦ Once all points assigned to cluster recompute centroid.

⑥ Repeat steps 4 & 5 until the location of the centroids stop changing and the cluster allocation of the points becomes constants

Note:

- K-means is only recommended for spherical data or symmetrical data
- For non-spherical data there won't be fair chance for k clusters to form clusters, and it will become highly biased clusters.

* Challenges and Improvements in K-Means

- ① We need to specify the no. of clusters beforehand.
- ② It is required to run the algorithm multiple times to avoid a suboptimal solution.
- ③ K-means does not behave very well when the clusters have varying sizes, different densities or non-spherical shapes.
- ④ An improvement to K-means is proposed in 2003 by Charles Elkan. It uses Euclidean distance calculations which is achieved by employing the TRIANGLE INEQUALITY.

On a triangle with sides a, b, c the straight line is always shortest

- ⑤ Another important improvement to k-means algorithm called "k-means++", was proposed in 2006 by David Arthur & Sergei.
- they introduced a smart initialization step that tends to select centroids that are distant from one another, and this makes the k-means algorithm much less likely to converge to a suboptimal solution.
- ⑥ yet another important variant of the k-means algorithm was proposed in a 2010 paper by David Sculley.
- instead of using the full dataset at each iteration, the algorithm is capable of using mini-batches, moving the centroids just slightly at each iteration. This speeds up the algorithm typically by a factor of 3 or 4 and makes it possible to cluster huge datasets that do not fit in memory. Scikit-learn implements this algorithm in the "MiniBatchKMeans" class.

II Hierarchical Clustering

Bottom-up / Agglomerative approach

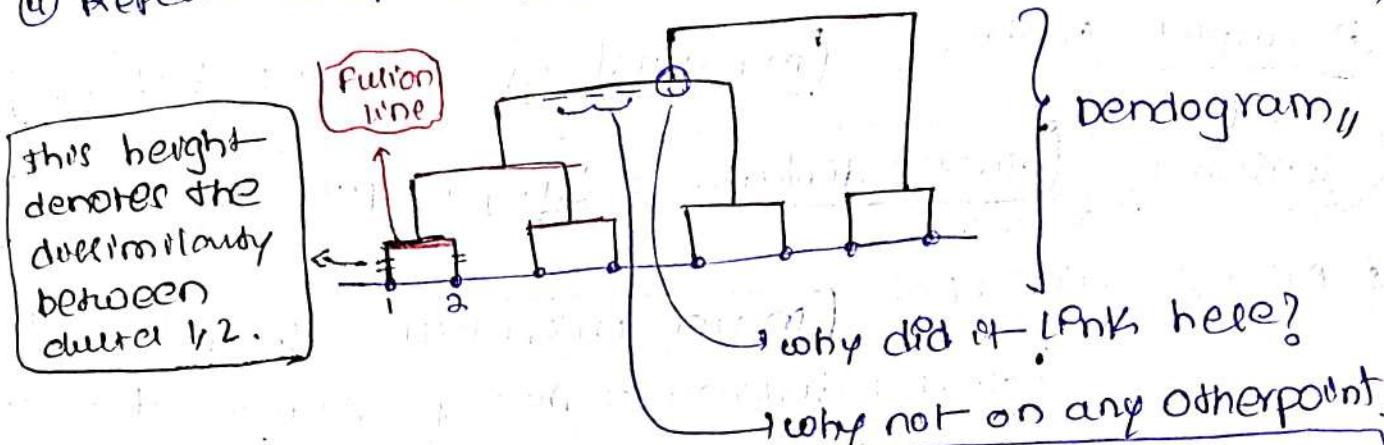
- Initially, all points are separate clusters.
- one main disadvantage of K-means is that we need us to pre-enter the number of clusters (k).
 - Hierarchical clustering is an alternative approach which does not need us to give the value of k beforehand and also, it creates a beautiful tree-based structure for visualization.
 - we start by defining any sort of similarity between the datapoints. Generally, we consider Euclidean distance. The points which are closer to each other are more similar than the points which are farther away.

Algorithm

- ① Begin with n observations and a measure (such as Euclidean distance) of all the $n(n-1)/2$ pairwise dissimilarities (distances).
Treat each observation as its own cluster.
So, initially we have n clusters.
- ② compare all the distances and put two closest points/clusters in the same cluster.
The dissimilarity between these two clusters indicate the height in the dendrogram at which "fusing line" should be placed.

③ compute the new pairwise inter-cluster dissimilarity among the remaining cluster.

④ Repeat steps 2 & 3 till we have only one cluster left,



Linkage methods

→ Based on pairwise distances, we can now compute a linkage matrix.

The linkage matrix is simply a table listing which pairs of points are merged at what step & what distance.

→ we can cut dendrogram to form flat clusters.

④ we know, arbitrary HC starts with clusters consisting of individual points.

→ later it compares the cluster with each other and merge the two "closest clusters".

→ since clusters are pair of points,

there are many different kinds of linkage methods.

↳ most common is Dissimilarity

Notelet

① Single Linkage

cluster distance = smallest pairwise distance

② complete Linkage

minimal intercluster distance

less sensitive
to outliers

cluster distance = largest pairwise distance

③ Average Linkage

mean intercluster dissimilarity

cluster distance = average pairwise distance

④ centroid Linkage

can result in undesirable envelopes

cluster distance = distance between the centroids of clusters

⑤ ward's Linkage

[Before & after merging]

cluster distance = minimize the variance in cluster.

simple linkage

minimal intercluster dissimilarity

→ single linkage can result in extended,

drawn-out clusters in which single observations are fused one-at-a-time.

→ cluster distance is the smallest distance b/w any point in cluster ① & any point in cluster ②

→ high sensitivity to outliers when forming flat clusters.

→ works well for low-noise data with unusual structures

⑩ DBSCAN (Density Based Spatial Clustering of Applications with Noise.)

→ It is an unsupervised machine learning algorithm. This algorithm defines clusters as continuous regions of high density.

Key words :-

① Epsilon :- (EPS)

This is the distance till which we look for the neighbouring points.

② min-points :-

The min no. of points specified by the user.

③ core point :-

If the no. of points inside the epsilon radius of a point is greater than or equal to the min-points then it is called a core point.

④ Border point :-

If the no. of points inside the epsilon radius of a point is less than the min points and it lies within the epsilon radius region of a core point, then it is called border point.

⑤ Noise :-

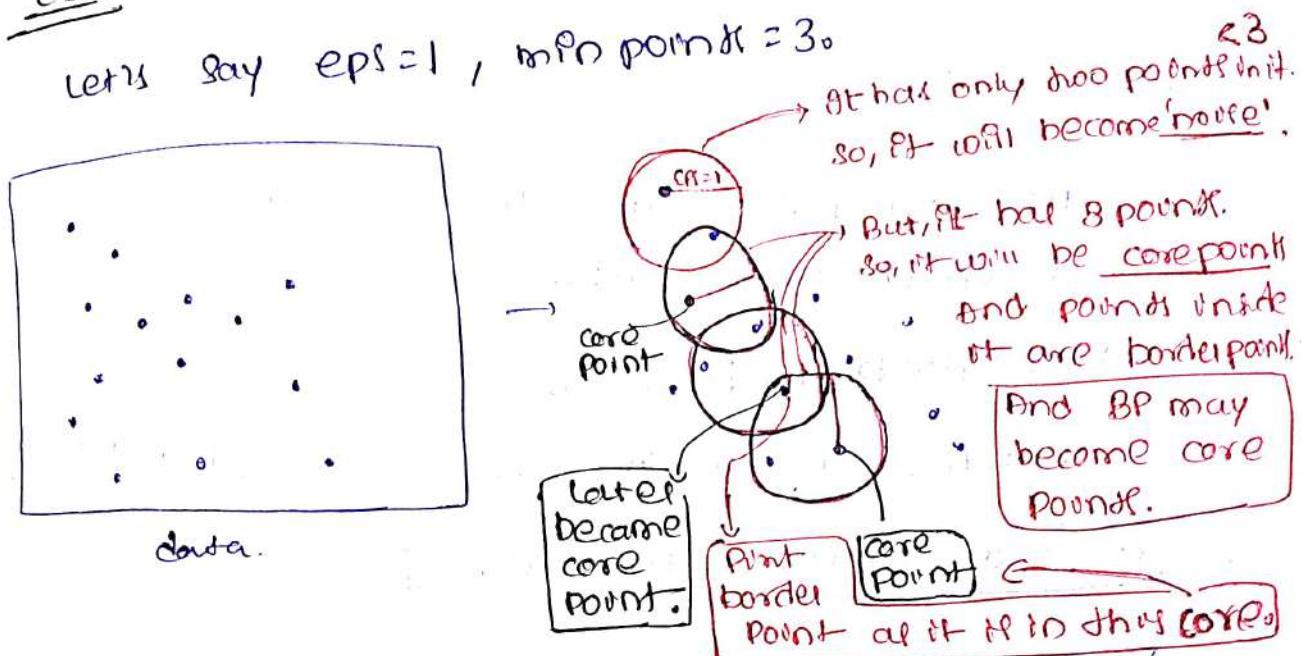
A point which is neither core nor a border point is a noise point.

Algorithm steps

- ① The algorithm starts with a random point in the dataset which has not been visited yet and its neighbouring points are identified based on the ϵ value.
- ② If the point contains \geq points than min points then the cluster formation starts and this point becomes a core point, else it's considered as Noise.
 - The point to note here is that a point initially classified as noise can later become a border point or it's in the ϵ radius of a core point.
- ③ If the point is not a core point, then all its neighbours become a part of cluster.
If the points in the neighbourhood turn out to be core points, then their neighbours are also part of cluster.
- ④ Repeat the steps above until all points are classified into different clusters or noise.

This algo works well if all clusters are dense enough, and they are well separated by low dense regions.

Ex 6



→ In short, DBSCAN is a very simple yet powerful algorithm, capable of identifying any no. of clusters of any shape; it is robust to outliers, and it has just two hyperparameters. (In practice, it needs only one parameter eps & min samples)

→ However, if the density varies significantly across one cluster, it can be impossible for it to capture all the clusters properly. Moreover its computational complexity is roughly $O(n \log n)$, making it pretty close to linear with regard to the no. of instances.

→ //

* Evaluating or clustering

Cluster validity

The validation of clusters created us a troublesome task.

→ the problem here is

clusters are "in the eyes of the beholder".

• good clusters will have

Problems in tool

- High inter class similarity

- low inter class similarity.

→ difficult to measure, subjective judgement

Approach of cluster validation

- External - compare your cluster to ground truth.

- Internal - Evaluating the cluster without

- reference to the external data

- Reliability - the clusters are not formed by chance (randomly) -

→ some statistical framework can be used

① External measures

→ no. of objects in the data P^o or P_1, P_2, \dots, P_m
the set of ground truth clusters :- C_1, C_2, \dots, C_n
the set of cluster formed by the algorithm.
the incidence matrix $N \times N$ matrix.

→ $P_{ij} = 1$, if the two points O_i^o & O_j^o belong to
the same cluster in the ground truth
else, $P_{ij} = 0$.

→ $C_{ij} = 1$, if the two points O_i^o & O_j^o belong to
the same cluster in the ground truth
else, $C_{ij} = 0$

now there can be following scenarios

1) $C_{ij} = P_{ij} = 1 \rightarrow$ Both the points belong to the
same cluster for both our algorithm &
ground truth (Agree) — SD

2) $C_{ij} = P_{ij} = 0 \rightarrow$ Both the points don't belong to
same cluster for both our algorithm &
ground truth (Agree) — DD

3) $C_{ij} = 1$ but $P_{ij} = 0 \rightarrow$ the points belong to the
same cluster for our algorithm but different
cluster in ground truth (Disagree) — SD

4) $C_{ij} = 0$ but $P_{ij} = 1 \rightarrow$ the points don't belong to
same cluster for our Algo but same
cluster in ground truth (Disagree) — DS

just like accuracy score,

$$\text{Rand Index} = \frac{\text{Total Agree}}{\text{Total Disagree}} = \frac{(SS + DD)}{(SS + DD + DS + SD)}$$

→ the disadvantage of this is it can be dominated by DD.

$$\text{Jaccard coefficient} = \frac{SS}{(SS + DS + SD)}$$

A higher value of Rand Index & Jaccard coefficient mean that the clusters generated by our algorithm mostly agree to the ground truth.

confusion matrix

n = no. of points

m_i = points in cluster i

C_j = points in class j

n_{ij} = points in cluster i coming from class j .

$$P_{ij} = \frac{n_{ij}}{m_i}$$

Probability of element from cluster i to be assigned to class j .

	class 1	class 2	class 3
cluster 1	n_{11}/P_{11}	n_{12}/P_{12}	n_{13}/P_{13}
cluster 2	n_{21}/P_{21}	n_{22}/P_{22}	n_{23}/P_{23}
cluster 3	n_{31}/P_{31}	n_{32}/P_{32}	n_{33}/P_{33}
	c_1	c_2	c_3
			n

→ Entropy of cluster i ,

$$e_i^o = - \sum P_{ij} (\log P_{ij})$$

→ For entering clustering algorithm, entropy can be generalized

$$C = \sum m_i e_i^o$$

→ purity or cluster i , $P_i^o = \max(P_{ij}^o)$

and for entire cluster it is $P(c) = \frac{1}{n} \cdot \sum P_i^o$

→ Purity, is the overall percentage of data points "clustered correctly."

→ A high value of purity score means that our clustering algorithm performs well against the ground truth.

Note:

→ In calculating External measure, most of time we don't have ground truths.

① Internal measures

There are the methods we use to measure the quality of clusters without external reference. There are two aspects of it.

Cohesion →

How closely the objects in the same cluster are related to each other. It is the within-cluster sum of squared distances. It is the same metric that we used to calculate for K-means algorithm.

$$WCSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_j - \mu_i)^2$$

Separation

How different objects in different clusters are and how different a well-separated cluster is from other clusters.

It is the between cluster sum of squared distance.

$$BSS = \sum c_i (m - m_i)^2$$

where, c is the size of individual cluster & m is centroid or all data points.

Note

$BSS + WSS$ is always a constant

→ the silhouette can be calculated as

$$s(x) = \frac{b(m) - a(m)}{\max\{a(m), b(m)\}}$$

where,

$a(m)$ is avg. distance of x from all other points in same cluster.

$b(m)$ is avg. distance of x from all other points in other cluster.

And silhouette coefficient is,

$$SC = \frac{1}{N} \sum s(m)$$

Higher SC means that the inter cluster similarity is less and the intra cluster dissimilarity is more. (Good clustering)

inter cluster similarity is less and the intra cluster dissimilarity is more.

* The curse of Dimensionality

→ In the real world datasets, we often encounter data with hundreds or even thousands of features. Problems with more features/dimensions.

① As the majority of the machine learning algorithms rely on the calculation of distance for model building, and as the number of dimensions increases, it becomes more & more computationally intensive to create a model out of it.

Eg:-

→ To calculate the distance b/w two points, in just one dimension, we can just subtract the coordinates of one point from another.

→ For 2D, $= \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$

→ And for ND, it becomes $\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + \dots + (z_1 - z_2)^2}$

→ This is the effort for just two points, just imagine the no. of calculations involved for all datapoints.

② Hard to visualize the relationship b/w features! If we have an n-dimensional dataset, the only solution left is to create either a 2D/3D graph out of it. Suppose we have 1000 features in the dataset, if we plan to create 2D graphs, we would need $(1000 \times 999)/2 = 499500$ combinations!! And we know it is not possible to spend time to analyze that many graphs to understand the relationship b/w the variables.

→ And like this, when we have huge no. of features, we need to ask the questions like -

- ① Are all the features really contributing to decision making.
- ② Is there a way to come to the same conclusion using lesser no. of features.
- ③ Is there a way to combine features to create a new feature and drop the old ones.
- ④ Is there a way to remodel features in a way to make them visually comprehensible.

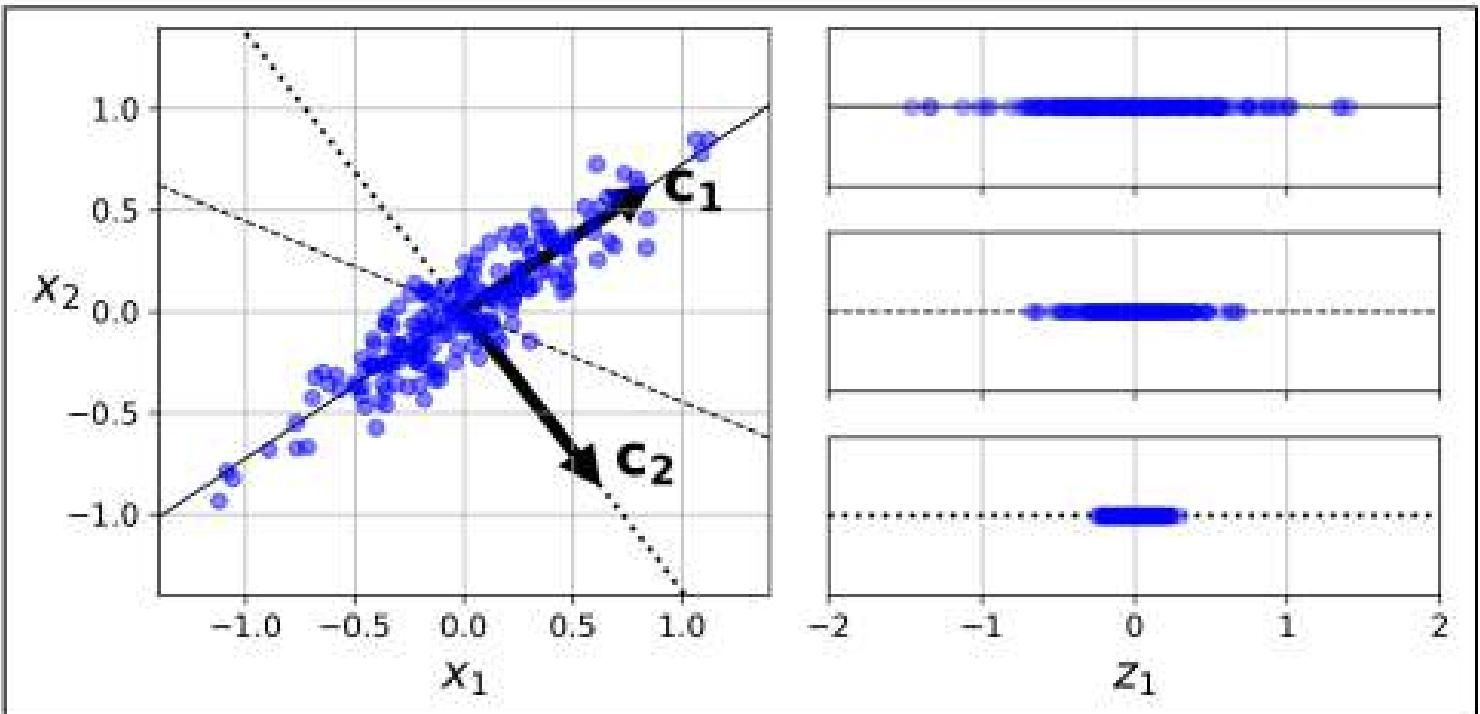
The answer to all the above questions is:

Dimensionality Reduction Technique

→ Dimensionality reduction is a feature selection technique using which we reduce the no. of features to be used for making a model without losing significant amount of information compared to original dataset.

→ In other words, a dimensionality reduction technique projects a data, of higher dimension to a lower dimension subspace.

→ DR technique shall be used before feeding the data to a machine learning algorithm, it reduces the space in which the distances are calculated, thereby improving ML algorithm performance.



Understanding Principal Component Analysis

Table Of Contents

1. Idea Behind PCA
2. What are Principal Components
3. Eigen Decomposition Approach
4. Singular Value Decomposition Approach
5. Why do we maximize Variance
6. What is Explained Variance Ratio
7. How to select optimal no.of Prinicpal Components
8. Understanding Scree plot
9. Issues with PCA
10. Understanding Kernel PCA

* Principal Component Analysis

- The most popular dimensionality reduction technique is the principal component analysis, it is an unsupervised algorithm used for feature selection.

Idea behind PCA

- PCA transforms and fits the data from a higher dimensional space to a lower dimensional subspace, this results into an entirely new coordinate system or the points where the first axis explains the most variance in the data (first principle component), this is the idea behind PCA.

what are principal components

- Principal components are the axis, that explain the maximum variance in the data, the first principal component explains the most variance, the second a bit less and so on..
- Each new axis / principal component found using PCA is a linear combination of the original features
- All the principal components are uncorrelated and orthogonal to each other.



Math Behind PCA

→ Principal component analysis can be implemented in two approaches, each with its unique focus and method. They are as follows:

① Eigen decomposition method

② singular value decomposition method,

Let's derive both them.

① Eigen Decomposition Method

→ Eigen decomposition is a process that decomposes a square matrix A into a set of eigenvectors and eigenvalues.

→ For a given square matrix A , the eigenvectors are the non-zero vectors that satisfy the equation,

$$A = \lambda V$$

where, V is the eigen vector and λ is the corresponding eigen value.

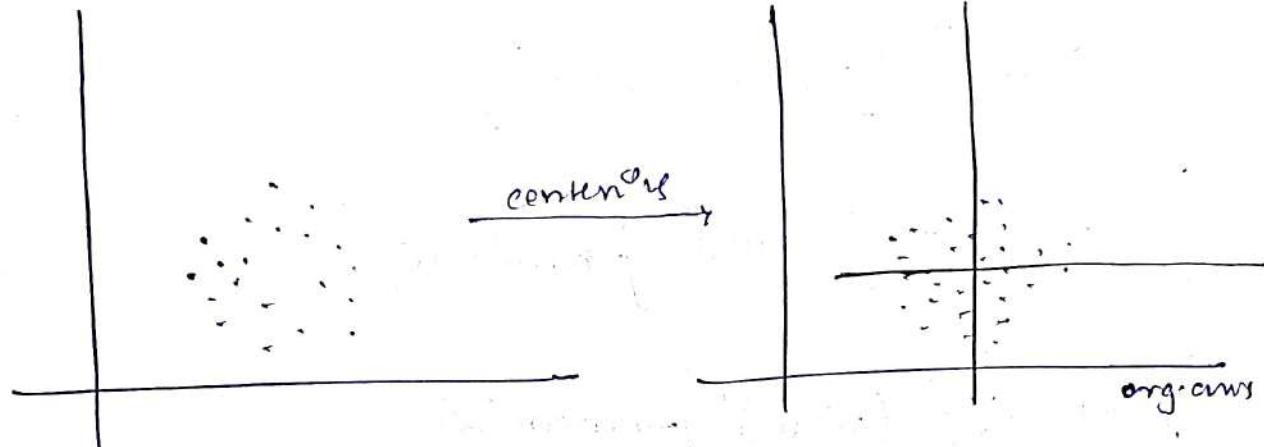
Steps for Eigen-decomposition approach

- ① Standardize columns of A , so that each feature has a mean of zero, this ensures that each feature contributes equally to the analysis.

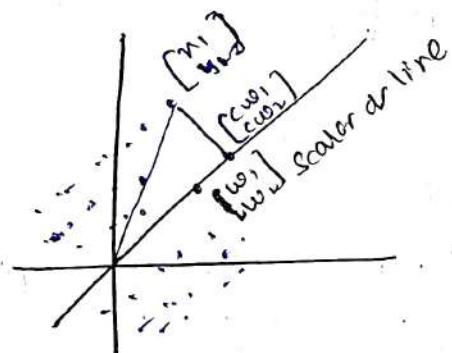
- ② compute co-variance matrix from the normalized data, $\Sigma = A^T A / (m-1)$
- ③ Perform eigen decomposition of Σ (covariance matrix), to obtain eigenvalues & eigenvectors.
- ④ choose the top K eigen vectors comput associated with the largest eigenvalues to construct the principal components as $A X_K$, K th eigen vector.

Proof

Step 1



we have, $D = \{n_1, n_2, \dots, n_m\} \quad n_i \in \mathbb{R}^d$



the projection point will be, $(n^T w) \cdot w$, considering w is a scalar or the best fit line/ the line that has least reconstruction error.

$$\text{Error}(\text{line}, \text{dataset}) = \sum_{i=1}^m \text{error}(\text{line}, n_i)$$

$$= \sum_{i=1}^m \|n_i - (n^T w) w\|^2$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n \| n - (n^T \omega) \cdot \omega \|^2$$

This is the mean of the residuals,
we can think of this as the avg. error.

$$= \frac{1}{n} \sum_{i=1}^n \left[(n - (n^T \omega) \cdot \omega)^T \cdot (n - (n^T \omega) \cdot \omega) \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[n^T n_i - (n^T \omega)^2 - (n^T \omega)^2 + (n^T \omega)^2 \right]$$

$$= \frac{1}{n} \sum_{i=1}^n \left[n^T n_i - (n^T \omega)^2 \right]$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n [n^T n_i - (n^T \omega)^2] \rightarrow \min \omega$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n (n^T \omega)^2 \rightarrow \max \omega$$

$$F(\omega) = \frac{1}{n} \sum_{i=1}^n (n^T \omega)(n^T n_i) = \frac{1}{n} \sum_{i=1}^n \omega^T (n_i \cdot n_i^T) \omega$$

$$F(\omega) = \omega^T \left(\frac{1}{n} \sum_{i=1}^n n_i \cdot n_i^T \right) \omega$$

we know
this is covariance matrix

Each column corresponds
to one data point.

$$\therefore \max(\omega) = \omega^T C \cdot \omega$$

* where, C - covariance matrix
 ω - eigen vector
 corresponding to
 max eigen value
 of C,

→ so, basically we can say that, the line which represents the data points or minimum reconstruction error is same as the line which minimizes the $w^T C w$, as C - covariance matrix,
 and the line is given by eigen vector corresponding to min. eigen value of C ,
 and this line is the principal component!!

★ First principal component (w_1) = $\arg \min_w (w^T C w)$

→ now, any line which minimizes the sum of errors/residuals, must also be orthogonal to w_1 .

Hence, second principal component will be orthogonal to the first $w_2^T w_1 = 0$ ★

→ By, continuing this procedure,
 we get $\{w_1, w_2 \dots w_d\}$

such that

$$\begin{aligned} \|w_k\| &= 1 \quad \forall k \text{ and} \\ w_i^T w_j &= 0 \quad \forall i \neq j \end{aligned}$$

Residuals after d rounds

$$\forall i \rightarrow n_i - [(n_i^T w_1) w_1 + (n_i^T w_2) w_2 + \dots + (n_i^T w_d) w_d] = r_i$$

$$\forall i \rightarrow n_i = (n_i^T w_1) w_1 + (n_i^T w_2) w_2 + \dots + (n_i^T w_d) w_d$$

→ if original dataset is of $(d \times n)$, then based on above
 so, after K rounds of PCA, we get $(d \times K) + (K \times n)$

Eg

+ suppose we have an initial dataset with 50 features and 100 samples. $d=50$, $n=100$.

Data dimension = $50 \times 100 = 1500$ data points to represent entire dataset.

After 5 PCA rounds, $k=5$

Data dimension = $(50 \times 5) + (5 \times 100) = 750$

So, now we only need half of datapoints to represent entire dataset

+ $\{w_1, w_2, w_3, w_4, w_5\} \rightarrow$ Representation,

Common for
dataset

+ $n^o \rightarrow [n_i^T w_1 \ n_i^T w_2 \ n_i^T w_3 \ n_i^T w_4 \ n_i^T w_5] \rightarrow$

Data
Point
Specific

Hence, to reconstruct n^o we don't need 1500 numbers like naive way

[we only need 5 principal components]

+ And through Eigen decomposition approach, we calculate the covariance matrix, and find its eigenvalues & eigenvectors. And define the principal components from eigenvectors. we can reconstruct by

A^TXK^T

② singular value decomposition method

→ SVD method factorizes the dataset in such a way that it becomes a product of the multiplication of 3 individual matrices.

$$X(\text{original data}) = U * \Sigma * V^T$$

where, U is the matrix that contains the principal components.

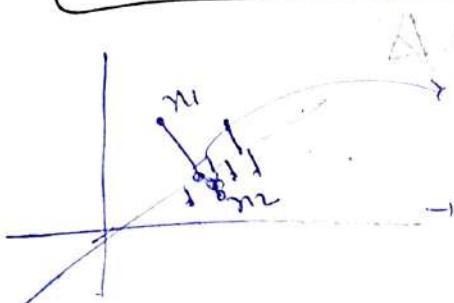
steps in SVD approach

- ① Standardize columns of A , to ensure each feature contributes equally to analysis.
- ② compute SVD of the centered data matrix to obtain matrices U , Σ , and V^T , such that,
$$X = U \Sigma V^T$$
 where U contains left singular vectors, Σ holds the singular values & V^T includes the right singular vectors.
- ③ select principal components; the principal components are obtained by selecting the columns for U corresponding to the largest singular values in Σ , these principal components represent the directions of maximum variance in the data, etc.

Note:-

→ what does high variance mean?

when making a proxy,
what if two points get same proxy.



Now, how do we distinguish
 $m_1 \& m_2$.

→ And how do we reconstruct
 $m_1 \& m_2$.

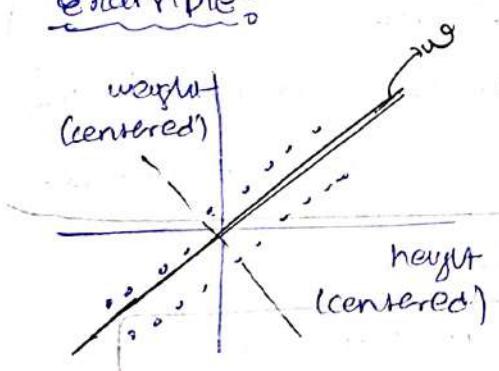
→ So, if all the datapoints are crowded,
then we can't reconstruct and the
information is lost.

This happens due to low variance / small variance.

Hence, we want directions where projections not crowded.

~~✓~~ Variance has to be high,
~~✓~~ so the datapoints will be spread out
and reconstruction of points can be done
properly.

Example:



- height gives some
information weight

→ we have two directions now
 w_1, w_2 ..

→ And these new directions are
de-correlated, every direction is
giving a new information
from the other.

→ Hence, we can say these
directions are **de-correlated**.

~~✓~~ This's what the algorithm is
essentially doing. & it's called
Principal Component Analysis.

* Once we get principal components, project the original data onto the subspace spanned by the selected principal components, effectively reducing the dimensionality of the data.

(*) What is the optimal no. of PCs needed?

* Before we ask this qn, we need to ask how much of the information in a given dataset is lost by projecting the observations onto first few principal components, i.e.,

the proportion of variance explained by each principal component

Explained variance ratio = $\frac{\text{variance by } m\text{th PC}}{\text{total variance}}$

$$\text{EVR} = \sum_{j=1}^n \left(\sum_{i=1}^p w_{ij} z_{ij} \right)^2$$

$$\sum_{j=1}^n \sum_{i=1}^p w_{ij}^2$$

$$\rightarrow \text{EVR of PC1} = \frac{\text{distance of PC1 point}}{\text{distance of PC1 point + PC2 + ... + PCn}}$$

If EVR of PC1 = 0.91, then that means PC1 explains 91% of the variance of data.

→ Now that we have a metrics to get the variance explained, we can choose the no. of dimensions that add up to a sufficiently large proportion of the variance (eg- 95%).

Eg:-

$$\text{EV}R \text{ or } PC_1 = 0.88$$

$$\text{EV}R \text{ or } PC_2 = 0.10$$

$$\text{EV}R \text{ or } PC_3 = 0.05$$

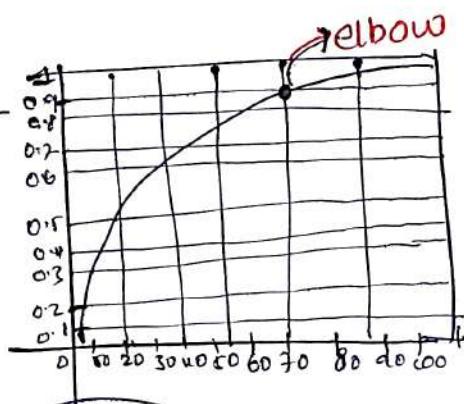
→ Now, if we want to preserve (99%) variance, then we can just take $PC_1 \& PC_2$, as they sum up to 0.94, 94% variance.

→ Instead if we want to preserve (99.5%) of variance, then we should take $PC_1, PC_2 \& PC_3$. as the sum up to 0.99, 99%. 99.5%.

Scree plots

→ Scree plots are the graphs that convey how much variance is explained by corresponding principal components.

→ So, by eyeballing the screeplot and looking for a point at which the proportion of variance explained by each subsequent PC drops off. This is often referred as



elbow in screeplots.

Issue/concern with PCA

1. Time complexity

→ in particular,
covariance matrix.

→ usually, it takes

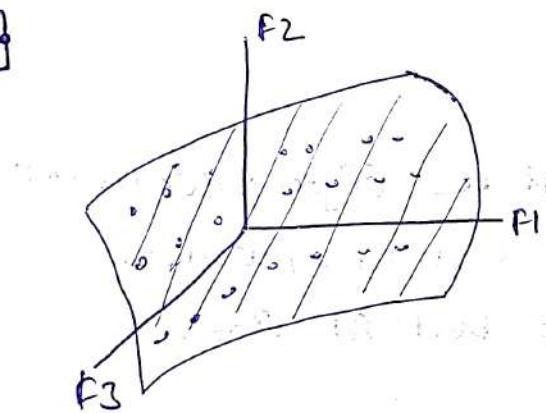
$$O(d^3)$$

→ if no. of features were large, [like 32000 for image]
then d^3 is going to be
very very large.

Finding the eigen vectors
and eigen values of
covariance matrix is
the time consuming step.

Eg: Face recognition (Eigen faces)

2.

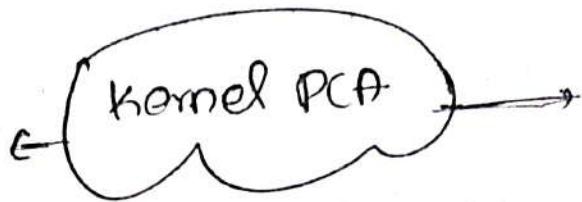


Data may not necessarily lie in a low-dimensional linear subspace

→ features will not be linearly related.

Notes

- we will have same solutions for both the above issues, solving one will solve the other.



• Input — $\{n_1, \dots, n_n\}$ $n_i \in \mathbb{R}^d$,

kernel $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$

Step ① compute $K \in \mathbb{R}^{n \times n}$,

where $K_{ij} = K(n_i, n_j)$ $\forall i, j$.

And then we need to center the kernel

Step ② compute β_1, \dots, β_d eigen vectors &

$n_{1,1}, \dots, n_{d,d}$ eigen values.

Step and normalize to get,

$$\alpha_k = \frac{\beta_k}{\sqrt{n_{kk}}}$$

Step ③

$$w_k = \phi(n) \cdot \alpha_k$$

Eigen computation of kernel matrix.
Hence, thus defeats the purpose,
as it needs $\phi(n)$, which we
tried to avoid all long!

we cannot
reconstruct the
eigen vectors or
the covariance
matrix

→ But we can still compute the
compressed representation.

$$\phi(n_i)^T \cdot w_k = \phi(n_i)^T \cdot \left(\sum_{j=1}^n \phi(n_j) \alpha_{kj} \right) = \sum_{j=1}^n \alpha_{kj} \phi(n_i)^T \phi(n_j)$$

$$\phi(m_i)^T w_k = \sum_{j=1}^n \alpha_{kj} \cdot k_{pj}$$

→ for downstream task,
this is usually
good enough.

→ typically, the reason why we do kernel PCA
is to get these projections and throw away
original points and only retain these
projections along each, or these data points.

→ this is where dimensionality reduction happens.

(eg)

original datapoint was 100 dimensions,
but then you only care about top 5.

→ so we map the 100 dimensional data into
a 5 dimensional projection onto these
5 most important directions.

→ that becomes a 5 dimensional representation
of your 100 dimensional data

→ once you learn these representations
in a low dimensional space, as put the
projections onto these eigen directions.
you can use these projection values
at your data point and work this for a
downstream task, (supervised learning task)

$$(\phi(m_i)^T w_k) = (\phi(m_i) \cdot \frac{w}{\|w\|}) \cdot \frac{w}{\|w\|} = \text{scalar}$$

modified

Step ③

compute $\sum_{j=1}^n \alpha_{kj} \cdot k_{ij} + k$

$$\begin{matrix} n \\ \in \mathbb{R}^d \end{matrix} \rightarrow \left[\sum_{j=1}^n \alpha_{1j} \cdot k_{1j}, \sum_{j=1}^n \alpha_{2j} \cdot k_{2j}, \dots, \sum_{j=1}^n \alpha_{Lj} \cdot k_{Lj} \right]$$

→ It might increase the dimension,
nevertheless it will allow you to capture
more complicated relationships.

Note

→ After Step 2, we have to center the kernel.

→ Given, $k \in \mathbb{R}^{n \times n}$

$$k_{ij} = k(n_i^o, n_j^o) + i_j^o$$

create a new kernel k^c — centered.

$$k_{ij}^c = k_{ij} - \theta_i l_j^T - l_i \theta_j^T + p l_i \cdot l_j^T$$

where,

$$\theta_i^o = \frac{1}{n} \sum_{p=1}^n k_{ip} \cdot k_p$$

$$p = \frac{1}{n} \sum_{p,j} k_{pj}$$

→ Now, we can say we are able to
solve both the issues by using kernel PCA.

– Supervised Algorithms –

Regression Models

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
Linear Regression	Linear Regression models a linear relationship between input variables and a continuous numerical output variable. The default loss function is the mean square error (MSE).	<ul style="list-style-type: none">1. Fast training because there are few parameters.2. Interpretable/Explainable results by its output coefficients.	<ul style="list-style-type: none">1. Assumes a linear relationship between input and output variables.2. Sensitive to outliers.3. Typically generalizes worse than ridge or lasso regression.
Polynomial Regression	Polynomial Regression models nonlinear relationships between the dependent, and independent variable as the n-th degree polynomial.	<ul style="list-style-type: none">1. Provides a good approximation of the relationship between the dependent and independent variables.2. Capable of fitting a wide range of curvature.	<ul style="list-style-type: none">1. Poor interpretability of the coefficients since the underlying variables can be highly correlated.2. The model fit is nonlinear but the regression function is linear.3. Prone to overfitting.
Support Vector Regression	Support Vector Regression (SVR) uses the same principle as SVMs but optimizes the cost function to fit the most straight line (or plane) through the data points. With the kernel trick it can efficiently perform a non-linear regression by implicitly mapping their inputs into high-dimensional feature spaces.	<ul style="list-style-type: none">1. Robust against outliers.2. Effective learning and strong generalization performance.3. Different Kernel functions can be specified for the decision function.	<ul style="list-style-type: none">1. Does not perform well with large datasets.2. Tends to underfit in cases where the number of variables is much smaller than the number of observations.
Gaussian Process Regression	Gaussian Process Regression (GPR) uses a Bayesian approach that infers a probability distribution over the possible functions that fit the data. The Gaussian process is a prior that is specified as a multivariate Gaussian distribution.	<ul style="list-style-type: none">1. Provides uncertainty measures on the predictions.2. It is a flexible and usable non-linear model which fits many datasets well.3. Performs well on small datasets as the GP kernel allows to specify a prior on the function space.	<ul style="list-style-type: none">1. Poor choice of kernel can make convergence slow.2. Specifying specific kernels requires deep mathematical understanding

Classification Models

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
SVM	In its simplest form, support vector machine is a linear classifier. But with the kernel trick, it can efficiently perform a non-linear classification by implicitly mapping their inputs into high-dimensional feature spaces. This makes SVM one of the best prediction methods.	<ul style="list-style-type: none"> 1. Effective in cases with a high number of variables. 2. Number of variables can be larger than the number of samples. 3. Different Kernel functions can be specified for the decision function. 	<ul style="list-style-type: none"> 1. Sensitive to overfitting, regularization is crucial. 2. Choosing a "good" kernel function can be difficult. 3. Computationally expensive for big data due to high training complexity. 4. Performs poorly if the data is noisy (target classes overlap).
Negrest Neighbors	Nearest Neighbors predicts the label based on a predefined number of samples closest in distance to the new point.	<ul style="list-style-type: none"> 1. Successful in situations where the decision boundary is irregular. 2. Non-parametric approach as it does not make any assumption on the underlying data. 	<ul style="list-style-type: none"> 1. Sensitive to noisy and missing data. 2. Computationally expensive because the entire set of n. points for every execution is required.
Logistic Regression (and its extensions)	The logistic regression models a linear relationship between input variables and the response variable. It models the output as binary values (0 or rather than numeric values.	<ul style="list-style-type: none"> 1. Explainable & Interpretable. 2. Less prone to overfitting using regularization. 3. Applicable for multi-class predictions. 	<ul style="list-style-type: none"> 1. Makes a strong assumption about the relationship between input and response variables. 2. Multicollinearity can cause the model to easily overfit without regularization.
Linear Discriminant Analysis	The linear decision boundary maximizes the separability between the classes by finding a linear combination of features.	<ul style="list-style-type: none"> 1. Explainable & Interpretable. 2. Applicable for multi-class predictions. 	<ul style="list-style-type: none"> 1. Multicollinearity can cause the model to overfit. 2. Assuming that all classes share the same covariance matrix. 3. Sensitive to outliers. 4. Doesn't work well with small class sizes.

Both Regression and Classification Models

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
Decision Trees	Decision Tree models learn on the data by making decision rules on the variables to separate the classes in a flowchart like a tree data structure. They can be used for both regression and classification.	1. Explainable and interpretable. 2. Can handle missing values.	1. Prone to overfitting. 2. Can be unstable with minor data drift. 3. Sensitive to outliers.
Random Forest	Random Forest classification models learn using an ensemble of decision trees. The output of the random forest is based on a majority vote of the different decision trees.	1. Effective learning and better generalization performance. 2. Can handle moderately large datasets. 3. Less prone to overfit than decision trees.	1. Large number of trees can slow down performance. 2. Predictions are sensitive to outliers. 3. Hyperparameter tuning can be complex.
Gradient Boosting	An ensemble learning method where weak predictive learners are combined to improve accuracy. Popular techniques include XGBoost, LightGBM and more.	1. Handling of multicollinearity. 2. Handling of non-linear relationships. 3. Effective learning and strong generalization performance. 4. XGBoost is fast and is often used as a benchmark algorithm.	1. Sensitive to outliers and can therefore cause overfitting. 2. High complexity due to hyperparameter tuning. 3. Computationally expensive.
Ridge Regression	Ridge Regression penalizes variables with low predictive outcomes by shrinking their coefficients towards zero. It can be used for classification and regression.	1. Less prone to overfitting. 2. Best suited when data suffers from multicollinearity. 3. Explainable & Interpretable.	1. All the predictors are kept in the final model. 2. Doesn't perform feature selection.
Lasso Regression	Lasso Regression penalizes features that have low predictive outcomes by shrinking their coefficients to zero. It can be used for classification and regression.	1. Good generalization performance. 2. Good at handling datasets where the number of variables is much larger than the number of observations. 3. No need for feature selection.	1. Poor interpretability/explainability as it can keep a single variable from a set of highly correlated variables.
AdaBoost	Adaptive Boosting uses an ensemble of weak learners that is combined into a weighted sum that represents the final output of the boosted classifier.	1. Explainable & Interpretable. 2. Less need for tweaking parameters. 3. Usually outperforms Random Forest.	1. Less prone to overfitting as the input variables are not jointly optimized. 2. Sensitive to noisy data and outliers.

– Unsupervised Algorithms –

Clustering Algorithms

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
K-Means	Most common clustering approach which assumes that the closer data points are to each other, the more similar they are. It determines K clusters based on Euclidean distances.	<ul style="list-style-type: none"> 1. Scales to large datasets 2. Interpretable & explainable results 3. Can generate tight clusters 	<ul style="list-style-type: none"> 1. Requires defining the expected number of clusters in advance. 2. Not suitable to identify clusters with non-convex shapes.
DBSCAN	Density-Based Spatial Clustering of Applications with Noise can handle non-linear cluster structures, purely based on density. It can differentiate and separate regions with varying degrees of density, thereby creating clusters.	<ul style="list-style-type: none"> 1. No assumption on the expected number of clusters. 2. Can handle noisy data and outliers 3. No assumptions on the shapes and sizes of the clusters 4. Can identify clusters with different densities 	<ul style="list-style-type: none"> 1. Requires optimization of two parameters 2. Can struggle in case of very high dimensional data
HDBSCAN	Family of the density-based algorithms and has roughly two steps: finding the core distance of each point, and expands clusters from them. It extends DBSCAN by converting it into a hierarchical clustering algorithm.	<ul style="list-style-type: none"> 1. No assumption on the expected number of clusters 2. Can handle noisy data and outliers. 3. No assumptions on the shapes and sizes of the clusters. 4. Can identify clusters with different densities 	<ul style="list-style-type: none"> 1. Mapping of unseen objects in HDBSCAN is not straightforward. 2. Can be computationally expensive
Agglomerative Hierarchical Clustering	Uses hierarchical clustering to determine the distance between samples based on the metric, and pairs are merged into clusters using the linkage type.	<ul style="list-style-type: none"> 1. There is no need to specify the number of clusters. 2. With the right linkage, it can be used for the detection of outliers. 3. Interpretable results using dendograms. 	<ul style="list-style-type: none"> 1. Specifying metric and linkages types requires good understanding of the statistical properties of the data 2. Not straightforward to optimize 3. Can be computationally expensive for large datasets
OPTICS	Family of the density-based algorithms where it finds core sample of high density and expands clusters from them. It operates with a core distance (e) and reachability distance.	<ul style="list-style-type: none"> No assumption on the expected number of clusters. 2. Can handle noisy data and outliers. 3. No assumptions on the shapes and sizes of the clusters. 4. Can identify clusters with different densities. 5. Not required to define fixed radius as in DBSCAN. 	<ul style="list-style-type: none"> 1. It only produces a cluster ordering. 2. Does not work well in case of very high dimensional data. 3. Slower than DBSCAN.

Dimensionality Reduction Techniques

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
PCA	Principal Component Analysis (PCA) is a feature extraction approach that uses a linear function to reduce dimensionality in datasets by minimizing information loss.	<ul style="list-style-type: none"> 1. Explainable Interpretable results. 2. New unseen datapoints can be mapped into the existing PCA space 3. Can be used as dimensionality reduction technique as preliminary step to other machine learning tasks 4. Helps reduce overfitting 5. Helps remove correlated features 	<ul style="list-style-type: none"> 1. Sensitive to outliers 2. Requires data standardization
t-SNE	t-distributed Stochastic Neighbor Embedding is a non-linear dimensionality reduction method that converts similarities between data points to joint probabilities using the Student t-distribution in the low-dimensional space	<ul style="list-style-type: none"> 1. Helps preserve the relationships seen in high dimensionality 2. Easy to visualise the structure of high dimensional data in 2 or 3 dimensions 3. Very effective for visualizing clusters or groups of data points and their relative proximities 	<ul style="list-style-type: none"> 1. The cost function is not convex: different initializations can get different results. 2. Computationally intensive for large datasets. 3. Default parameters do not always achieve the best results
UMAP	Uniform Manifold Approximation and Projection (UMAP) constructs a high-dimensional graph representation of the data then optimizes a low-dimensional graph to be as structurally similar as possible.	<ul style="list-style-type: none"> 1. It can be used as general-purpose dimension reduction technique as a preliminary step to other machine learning tasks. 2. Can be very effective for visualizing clusters or groups of data points and their relative proximities. 3. Able to handle high dimensional sparse datasets 	<ul style="list-style-type: none"> 1. Default parameters do not always achieve the best results
ICA	Independent Component Analysis (ICA) is a linear dimensionality reduction method that aims to separate a multivariate signal into additive subcomponents under the assumption that independent components are non-gaussian. Where PCA "compresses" the data, ICA "separates" the information.	<ul style="list-style-type: none"> 1. Can separate multivariate signals into its subcomponents 2. Clear aim of the method: only applicable if there are multiple independent generators of information to uncover. 3. Can extract hidden factors in the data by transforming a set of variables to new set that maximally independent. 	<ul style="list-style-type: none"> 1. Without any prior knowledge, determination of the number of independent components or sources can be difficult. 2. PCA is often required as a pre-processing step.

Association Rules

ALGORITHM	DESCRIPTION & APPLICATION	ADVANTAGES	DISADVANTAGES
Apriori algorithm	The Apriori algorithm uses the join and prune step iteratively to identify the most frequent itemset in the given dataset. Prior knowledge (apriori) of frequent itemset properties is used in the process.	1. Explainable & interpretable results. 2. Exhaustive approach based on the confidence and support.	1. Requires defining the expected number of clusters or mixture components in advance 2. The covariance type needs to be defined for the mixture of component
FP-growth algorithm	Frequent Pattern growth (FP-growth) is an improvement on the Apriori algorithm for finding frequent itemsets. It generates a conditional FP-Tree for every item in the data.	1. Explainable & interpretable results. 2. Smaller memory footprint than the Apriori algorithm	1. More complex algorithm to build than Apriori 2. Can result in many (incremental) overlapping/trivial itemsets
FP-Max Algorithm	A variant of Frequent pattern growth that is focused on finding maximal itemsets.	1. Explainable & Interpretable results. 2. Smaller memory footprint than the Apriori and FP-growth algorithms	1. More complex algorithm to build than Apriori