

Machine learning

Unit : 3

-Prof. Mayur Ghorpade

ML Unit: III

Supervised Learning: Regression

Q. Explain and Compute training, error and generalization error.

Training Error:

- In machine learning, training a predictive model means finding a function which maps a set of value x to a value y .
- Training error is calculated as:

$$E_{\text{train}} = \frac{1}{n} \sum_{i=1}^n \text{error}(f_D(x_i), y_i)$$

where,

$n \rightarrow$ No. of training example.

$f_D(x_i) \rightarrow$ Predictive value.

$y_i \rightarrow$ Actual or True value.

$(f_D(x_i), y_i) \rightarrow$ These two values are same or not, if not then find difference.

Generalization error:

- For Supervised learning applications in machine learning and statistical learning theory, generalization error is a measure of how accurately an algorithm is able to predict outcome values for previously unseen data.

$$E_{\text{gen}} = \int \text{error}(f_D(x_i), y_i) P(y, x) dx$$

where,

$f_D(x_i) \rightarrow$ Predictive value

$y_i \rightarrow$ Actual value

$P(x, y) \rightarrow$ How often we expect to see such x & y .

Q. Explain terms Bias, Variance, Generalization.

Bias:

- Bias is a phenomenon that skews the result of an algorithm in favour or against an idea.
- Bias is considered a systematic error that occurs in the ML model itself due to incorrect assumption in the ML process.
- Low Bias :- A low Bias model will make fewer assumptions about the form of the target function.
 - Closely match the training data set.
- High Bias :
 - ~~makes~~ makes high bias makes more assumptions.
 - It cannot perform well on new data.

Variance:

- Variance tells that how much a random variable is different from its expected value.
- Variance error are either low variance or high variance.
- Low Variance :
 - Small variation in the prediction of target funn
- High Variance :
 - Large variation in the prediction of target funn,
 - High Variance model leads to overfitting,
 - Increases model complexities,

Different combination of Bias-Variance:

- Low Bias, Low Variance : Ideal ML Model (Not possible)
- Low Bias, High Variance : Occurs when the model learns with a large no. of parameters. (Overfitting).
- High Bias, Low Variance : Occurs when few no. of parameter (Underfitting)
- High Bias, High Variance : Prediction are inconsistent & also inaccurate on average.

Generalization:

- Generalization is a term used to describe a model's ability to react to new data.
- After being trained on a training set, a model can digest new data and make accurate prediction.
- Model is accurately trained that time generalization not possible.
- It will make inaccurate prediction when given new data, it making model useless even if able to give accurate prediction for the training data is called overfitting.
- Underfitting, happens when a model has not been trained enough on the data.

Q. Explain overfitting with the techniques to reduce the overfitting?
→ Overfitting:

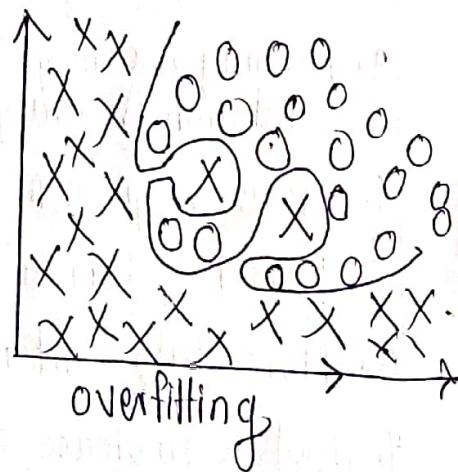
- When a model gets trained with so much data, it starts learning noise & inaccurate data entries in our dataset.
- So High Variance in test data.
- It causes overfitting.
- Simply model does not make accurate prediction on testing data called overfitting.

Reasons of overfitting:

- High Variance & Low Bias
- The model is too complex
- size of training data.

Techniques to reduce overfitting:

1. Increase Training Data.
2. Reduce model complexity.
3. Early stopping during the training phase.
4. Ridge Regularization & Lasso Regularization.
5. Use dropout for neural networks to tackle overfitting.
6. Ensembling



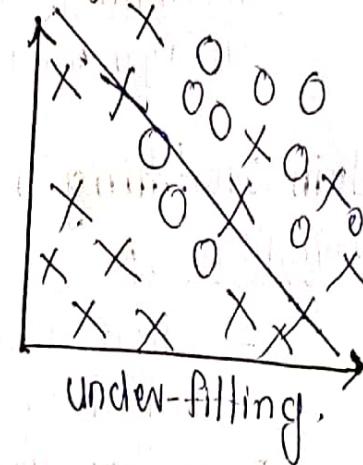
Q. Explain term Underfitting, with techniques to reduce underfitting.

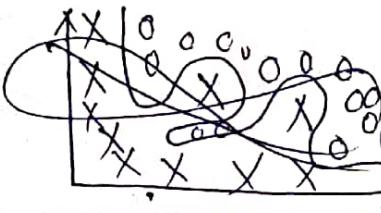
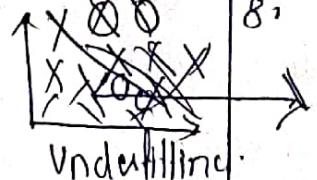
→ Underfitting:

- Machine learning algorithm is said to have underfitting when it cannot capture the underlying trend of the data.
- It only performs well on training data but performs poorly on testing data.
- Underfitting destroys accuracy of ml model.
- Reasons for Underfitting:
 - High Bias and low variance.
 - Size of the training data is not enough.
 - The model is too simple.

Techniques to reduce underfitting:

1. Increase model complexity.
2. Remove noise from the data.
3. Increase the no. of features.
4. Increasing training time of model.



Underfitting	Overfitting
<ol style="list-style-type: none">1. High Bias, Low Variance2. Performs poorly on training data, well on testing data.3. Training Accuracy → Poor Validation Accuracy → Poor4. More complex model5. Less regularization6. More data can't help.7. Need to increase complexity.8.   Underfitting.	<ol style="list-style-type: none">1. High Variance, low Bias2. Performs poorly on testing data, well on training data.3. Training Accuracy → Good Validation Accuracy → Poor4. More simple model.5. More regularization6. More data can help.7. Need to reduce complexity.8.  overfitting.

• What is Regression and List the type of Regression.

→ Regression :-

- Regression is supervised learning techniques.
- It is used to find
- Regression analysis is a statistical method to model the relationship between a dependent and independent variables with one or more variable.
- It is used for prediction, forecasting, time series modelling.
- example: linear regression.

Types of Regression:

- linear Regression
- Logistic Regression
- Polynomial Regression
- Ridge Regression
- Lasso Regression

Q. Explain Linear Regression.

→ Linear Regression:-

- linear regression is statistical regression method which is used for predictive analysis.
- It is one of the very easy and simple algorithm which works on regression and shows the relation between the continuous variables.
- It is used for solving the regression problem in ML.
- Linear regression shows the linear relationship between the independent variable (X-axes) and the dependent variable (Y-axes), hence, called linear regression.

$$Y = aX + b$$

where:

Y → Dependent variables.

X → Independent variable

a & b → linear coefficient

Application:

- Analyzing trends & sales estimates,
- Salary forecasting,
- Real estate prediction.
- Arriving at GTAs in traffic.

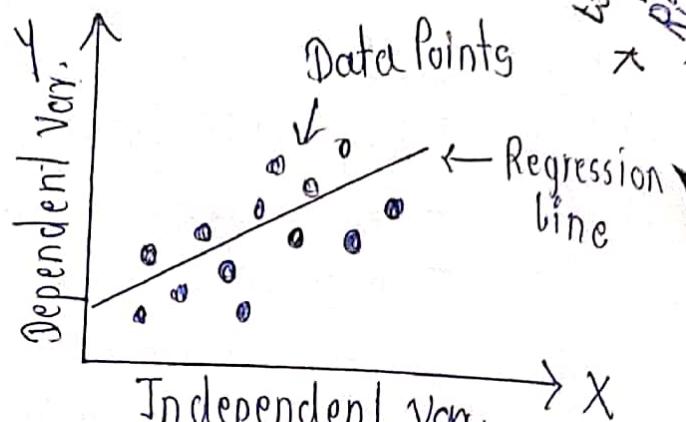


fig. Linear Regression.

Q. Explain Lasso Regression.

→ Lasso Regression:

- Lasso regression is regularization technique to reduce the complexity of the model.
- It is similar to Ridge regression except that penalty ~~that~~ term contains only the absolute weights instead of the square of weights.
- Since it takes absolute value, hence it can shrink the slope to 0, whereas Ridge Regression can only shrink it near to 0.
- It is also called L1 Regularization.
- Equation

$$L(x, y) = \min \left(\sum_{i=1}^n (y_i - w_i x_i)^2 + \lambda \sum_{i=1}^n |w_i| \right)$$

- Application:

Lasso has been applied in economics and finance.

Explain Ridge Regression.

Ridge Regression:

- Ridge Regression is the one the most robust versions of linear regression in which a small amount of bias is introduced so that we can get better long term predictions.
- The amount of bias added to the model is known as Ridge Regression penalty.
- We can compute this penalty term by multiplying with the lambda to the squared weight of each individual features.
- Equation,

$$L(x, y) = \min \left(\sum_{i=1}^n (y_i - w_i x_i)^2 + \lambda \sum_{j=1}^n (w_j)^2 \right)$$

- Ridge regression is a regularization technique, which is used to reduce the complexity of the model.
- It is also called as L2 regularization.
- It helps to solve the problems if we have more parameters than samples.

Q. Explain Gradient Descent Algorithm.

Gradient Descent Algorithm:

- Gradient Descent is an optimization algorithm which is used to minimize the cost function for many ML algorithms.
- Gradient Descent algorithm is used for updating the parameters of the learning models.
- Types:
 - Batch Gradient Descent:
used for small training dataset.
 - Mini-Batch Gradient Descent:
for large training set.
 - Stochastic Gradient Descent: Process one training dataset set at a time

∴ Parameters will be updated after each iteration.

classmate

Q. Explain MAE, MSE, RMSE, R².

→ MAE (Mean Absolute Error):

- MAE represents the average of the absolute difference between the actual & predictive value in the dataset.

$$MAE = \frac{1}{N} \sum |y - \hat{y}|$$

Where,

N → No. of data Points

y → Actual output

\hat{y} → Predictive output

- Advantages:

- MAE you get in the same unit as the output variable.
- It is most Robust to outliers.

- Disadvantages:

- Graph of MAE is not differentiable so we apply different optimizers like GDA.

→ MSE (Mean Squared Error):

- MSE represent the average of the squared difference betn the original and predicted values in the dataset.
- It measures the variance of the residuals.
- We perform squared to avoid the cancellation of negative terms, & it is the benefit of MSE.

$$MSE = \frac{1}{n} \sum (y - \hat{y})^2$$

- Advantages:

- Graph of MSE is differentiable.

- Disadvantages:

- The value you get after calculating MSE is a squared unit of output.
- It is not robust to outliers.

RMSF (Root Mean Squared Error):

- RMSF is the square root of MSE.
- It measures standard deviation of residuals.

Advantages:

- Output value you get is in the same unit as the required output variable.

Disadvantages:

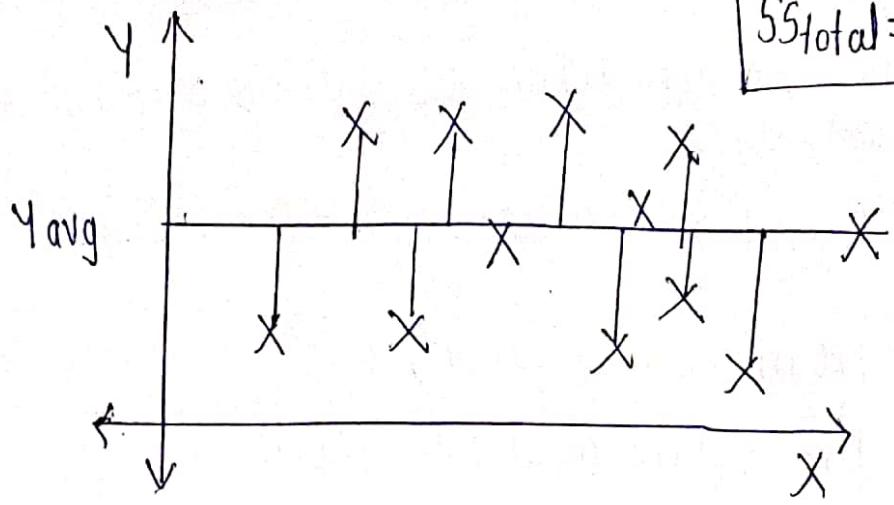
- Not robust to outliers.

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$$

$$\text{RMSE} = \sqrt{\text{MSE}}$$

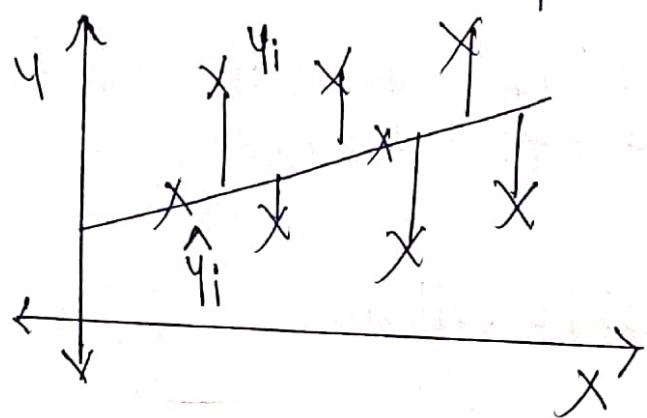
* R-Squared :-

- R-squared is a statistical measure that represents the goodness of fit of a regression model.
- The ideal value for R-squared is 1.
$$R^2 = \frac{1 - \text{MSE(model)}}{1 - \text{MSE(base line)}}$$
- R-squared is a comparison of the residual sum of squares (SS_{res}) with the total sum of squares (SS_{tot}).
- The total sum of squares is calculated by summation of squares of perpendicular distance b/w data points & the average line.



$$SS_{\text{Total}} = \sum (y_i - y_{\text{avg}})^2$$

The residual sum of squares is calculated by the summation of squares of perpendicular distn bet'n data points & the best fitted line.



$$SS_{\text{Res}} = \sum (y_i - \hat{y}_i)^2$$

$$R^2 = 1 - \frac{SS_{\text{Res}}}{SS_{\text{Total}}}$$

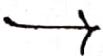
SS_{Res} → Residual Sum of squares
 SS_{Total} → total sum of squares.

Q. Why we require evaluation metrics?

- - ML model cannot have 100 per cent eff. otherwise the model is known as a biased model which further includes the concept of overfitting & underfitting.
- It is necessary to obtain the accuracy on training data, But it is also important to get a genuine & approximate result on unseen data.
- So to build & deploy a generalized model we require to evaluate the model on different metrics which helps us to better optimize the performance, fine-tune it, & obtain a better result.

Following table shows the midterm & final exam grades obtained for student in a database course. Use the method of least squares using regression to predict the final exam grade of a student who received 86 in the mid-term exam.

Midterm ex (X)	72	50	81	74	94	86	59	83	86	33	88	81
Final Exam (Y)	84	63	77	78	90	75	49	79	77	82	74	90



S1 No	X	Y	XY	X^2
1	72	84	6048	5184
2	50	63	2650	2500
3	81	77	6237	6561
4	74	78	5772	5476
5	94	90	8460	8836
6	86	75	6450	7396
7	69	49	2841	3481
8	83	79	6557	6889
9	86	77	6622	7396
10	33	82	1716	1089
11	88	74	6512	7744
12	81	90	7200	6561
Total	887	878	67209	69113

Eqⁿ of Reg. line is $y' = ax + b$

$$a = \frac{n \sum XY - \sum x \sum y}{n \sum X^2 - (\sum x)^2} = \text{[redacted]} 0.65$$

$$b = \frac{1}{n} (\sum y - n \bar{x} \bar{y}) = 25.12$$

$$y' = 0.65x + 25.12$$

The final ex. grade of a student who received 86 in the mid term exam.

$$y' = 0.65 \times 86 + 25.12$$

$$y' = 81.02$$

Q. Given the following data for the sales of car of an automobile company for six consecutive years. Predict the sales for next two consecutive years.

Years (X)	2013	2014	2015	2016	2017	2018
Sales (Y)	110	100	250	275	230	300

→ we will take $t = x - 2013$

Sr. No	t	Y	tY	t^2
0	0	110	0	0
1	1	100	100	1
2	2	250	500	4
3	3	275	825	9
4	4	230	920	16
5	5	300	1500	25
Total	15	1265	3845	55

Eqn. of regression line is

$$y' = at + b \quad a = 39, b = 113.33$$

eqn. of line become

$$y' = 39t + 113.33$$

The sales of company for next two years,

$$X = 2019, t = 6, \quad y' = 39 \times 6 + 113.33 = 347.33$$

$$X = 2020, t = 7, \quad y' = 39 \times 7 + 113.33 = 386.33$$

Unit : 4

MC

UNIT - 4] Supervised Learning - Classification.

Q. What is classification? List the types of classification?

→ Classification :-

- Classification algorithm is the supervised learning technique that is used to identify the category of new observations on the basis of training data.
- In classification, a ~~green~~ program learns from the given dataset of observations & then classifies the new observations into no. of classes or groups.

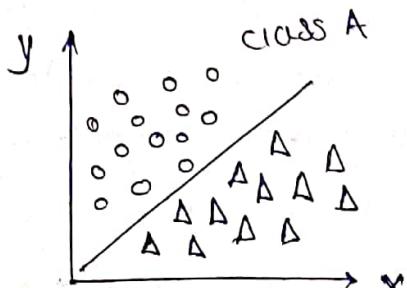


fig. classification.

Types :-

① Binary classification :-

- The problem has only two possible outcome.
- e.g. Yes or No, 0 or 1, etc.

② Multi-class classifier :-

- The problem has more than two possible outcome.
- e.g. classification of music.

Types of ML classification Algorithm :-

- KNN
- SVM
- Naive Bayes
- Decision tree
- random forest classification.

Q. Explain KNN Algorithm ?

→ K-nearest neighbour :-

- KNN is one of the simplest machine learning algorithm based on supervised learning technique.
- KNN algorithm assumes the similarity between the new data & available data & put the new data into category that most similar to the available categories.
- KNN algorithm stores all the available data and ~~then~~ classified a new data point based on the similarity.
- This means when new data appears then it can be easily classified into a well suited category by KNN algorithm.
- It is used for both classification as well as regression.
- It is also called as a lazy learned algorithm.

Steps :-

Step - 1] Select the number K of neighbors.

Step - 2] calculate the euclidean distance of K number of neighbors.

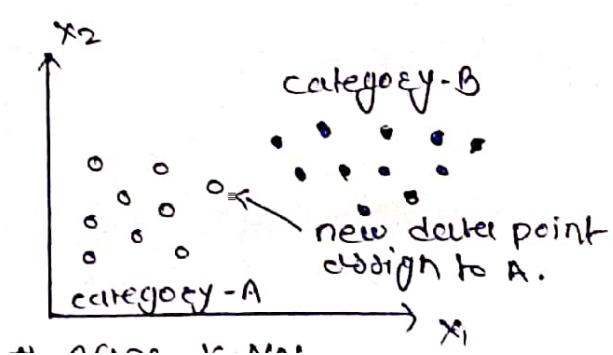
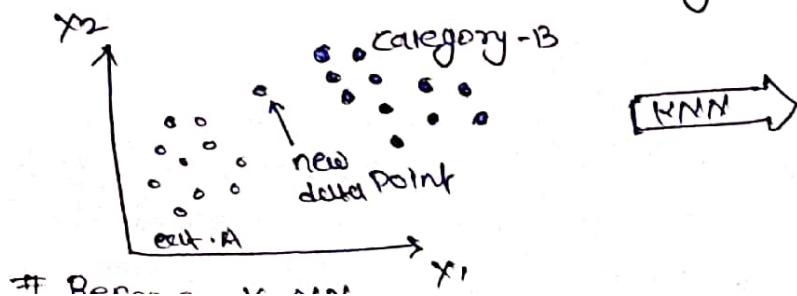
$$\text{Euclidean distance b/w A \& B} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

Step - 3] Take the K nearest neighbors as per the calculated Euclidean distance.

Step - 4] Among these K neighbors, count the number of the data points in each category.

Step - 5] Assign the new data points to that category for which the number of neighbors is maximum.

Step - 6] Our model is ready.



Advantages of K-NN algorithm :-

- It is simple to implement.
- It is robust to the noisy training data.
- More effective on large training data.

Disadvantages :-

- Always needs to determine the value of K, which may be complex some time.
- High computation cost.

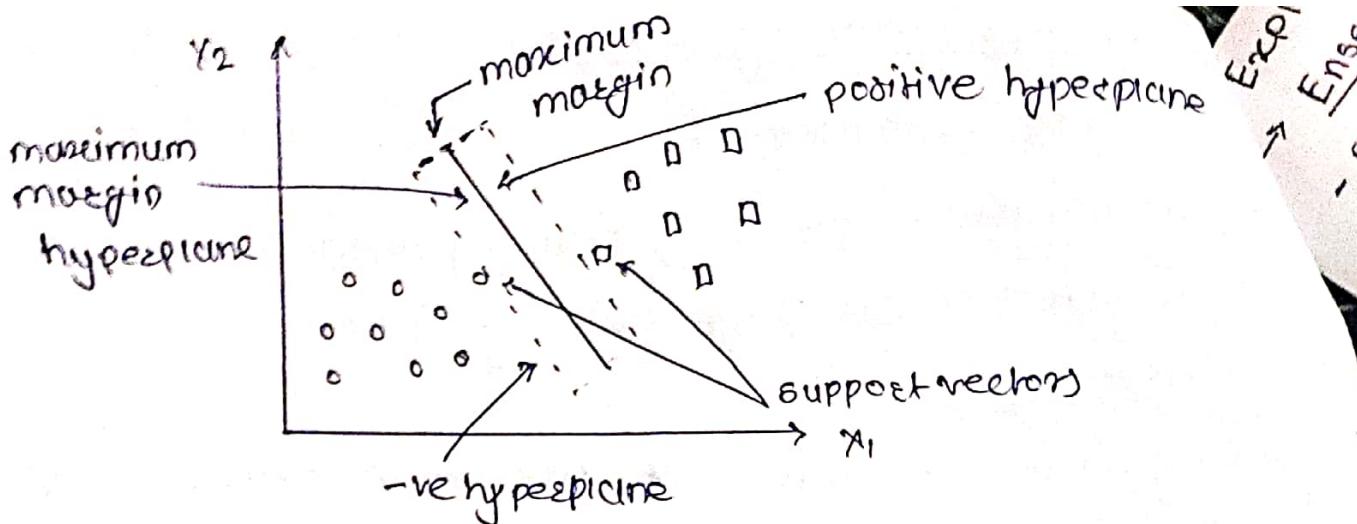
Applications :-

- finance - credit rating, loan management, stock market.
- medicine

Q. Define SVM - Explain how margin is computed & optimal hyper plane is decided ?

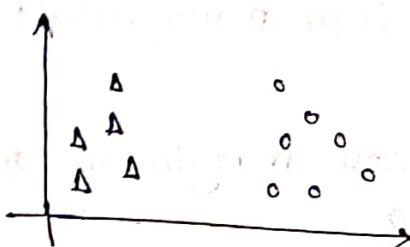
→ SVM -

- Support Vector machine is one of the most popular supervised learning algorithm.
- used in classification as well as regression problems. but primarily used for classification.
- the goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes.
- so that we can easily put new data point in the correct category in future.
- this best decision boundary is called a hyper plane.
- SVM chooses the extreme points that help in creating the hyperplane.
- these extreme are called as support vectors.
- the boundary which devide classes & helps to find the support vector in example is called Optimal Decision Boundary

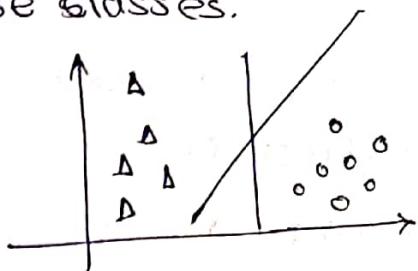


Example :-

Suppose we have a dataset has two tags (green & blue), & then classify this by using SVM.

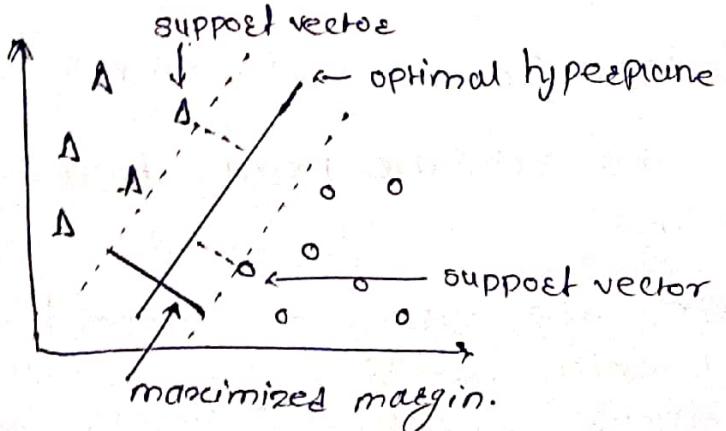


here it is 2-d space so using straight line, we can easily separate these two classes, but it has multiple lines that can separate, these classes.



Hence the SVM helps to find the best line or decision boundary, this best boundary is called hyperplane.

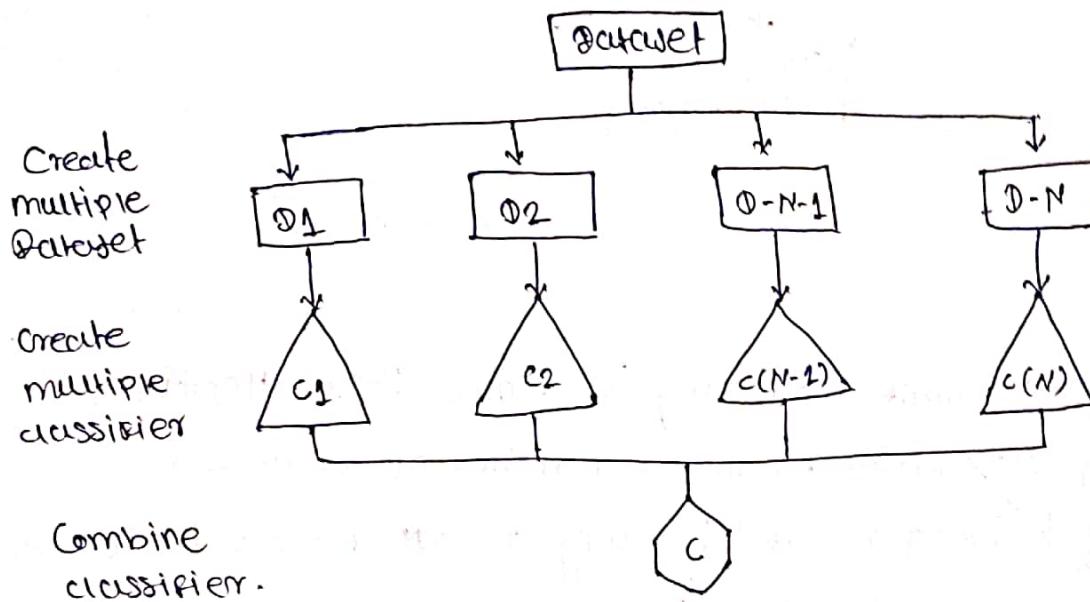
The data points which are closest to the line is called support vector.



Explain : Ensemble learning & explain the term bagging & boosting?

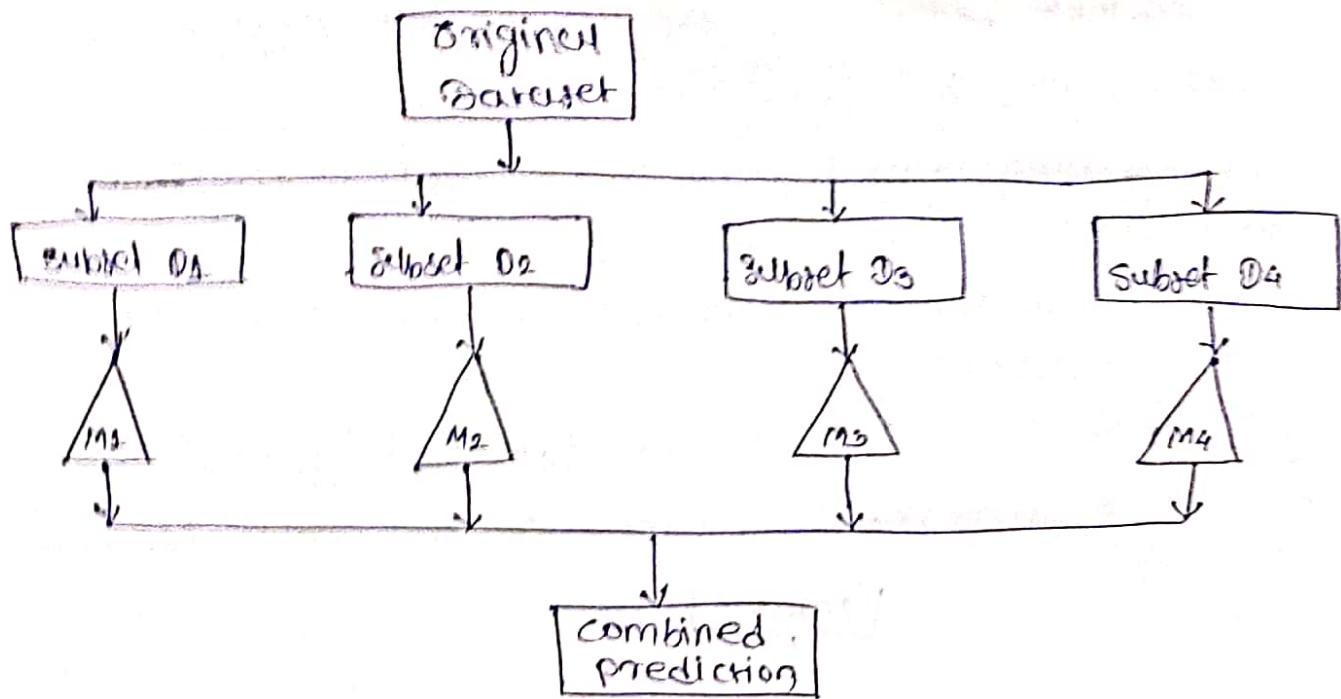
→ Ensemble Learning :-

- Ensemble learning helps improve ML results by combining several models.
- The model that contribute to the ensemble are called as ensemble members.
- They may or same type or of different types.
- They may or may not be trained on the same training data.



Bagging :-

- Bagging is the combine approach of Bootstrap & aggregation.
- It is used to reduce the variance of a decision tree.
- Implementation steps of Bagging :
 1. multiple subset are created from the original data set with equal tuples, selecting observations with replacement.
 2. A base model is created on each of these subsets.
 3. Each model is learned in parallel from each training set and independent of each other.
 4. The final predictions are determined by combining the predictions from all the models.



Boosting

- Boosting is an ensemble modeling technique that attempts to build a strong classifier from the number of weak classifiers.
- It is done by building a model by using a weak models in series.
- Firstly, model is created by using training data.
- Then second model built which tries to correct the errors present in the first model.
- This procedure is continued & models are added until either the complete training data set is predicted correctly or the maximum no. of models are added.
- AdaBoost was the first really successful boosting algorithm developed for the purpose of binary classification.

④

Explain random forest algorithm in detail?

Randon forest :-

- Random forest is a supervised ML algorithm that is used for classification & regression.
- It is based on ensemble learning.
- Random Forest is a classifier that contains a number of decision tree on various subsets of the given datasets & takes the average to improve the predictive accuracy of that dataset.
- Steps :-

1. select random k data points from the training set.
2. Build the decision trees associated with the selected data points.
3. choose the number N for decision trees that you want to build.
4. Repeat step 1 & 2.
5. For new data points , find the predictions of each decision tree , & assign the new data points to the category that wins the majority vote.

Advantages :-

- Random Forest is capable of performing both classification & regression.
- It is capable of handling large datasets with high dimensionality.
- prevent the overfitting problem.

Disadvantages :-

- It is not more suitable for regression tasks.

Applications :-

- 1) Banking.
- 2) medicine.
- 3) Land use.
- 4) marketing.

Q. Difference bet? Binary classification vs multiclass classif?

Parameters	Binary classification	multi-class classification
No. of classes	It is a classification of 2 groups: i.e. classifies objects in at most two classes.	there can be any no. of classes in it, i.e. classifies the object into more than two classes.
Algorithms used	<ul style="list-style-type: none"> • Logistic regression • K-nearest neighbor • decision tree • SVM • naive bay 	<ul style="list-style-type: none"> • KNN • decision tree • naive bay • Random forest • gradient boosting.
Example	<ul style="list-style-type: none"> • email spam detection • churn prediction 	<ul style="list-style-type: none"> • face classification • optical character recognition.

Q. explain variant of multiclass classification : one vs one
one vs all?

→ One vs one :-

- In one vs one reduction, one trains $\frac{k(k-1)}{2}$ binary classifiers for a k-way multiclass problem.
- each receives the samples of a pair of classes from the original training set and learns to distinguish these two classes.
- At predictive time, a voting scheme is applied:

$\frac{K(K-1)}{2}$ classifiers are applied to an unseen sample & the class that gets the highest number of +1 predictions, gets the predicted by combined classifier.

- Like one vs all one suffers from ambiguities in that some regions of its input space may receive the same number of votes.

one vs all :- (One vs rest)

- one vs all strategy involves training single classifier per class with the samples of that class as positive samples & all other samples are negative

#- Inputs :

1] L , a learned.

2] Samples X .

3] Labels y where $y_i \in \{1, \dots, K\}$ is the label for the sample x_i .

Output :

1] a list of classifiers f_k for $k \in \{1, \dots, K\}$

Procedure :

i] construct a new label vector Z where $z_i = y_i$, if $y_k = K$ & $z_i = 0$ otherwise.

ii] Apply L to X , Z to obtain f_k .

Q. Explain terms : Accuracy, Precision, Recall, Fscore?

→ confusion matrix:

Actual value

		positive	negative
positive	TP		FP
negative	FN		TN

Accuracy :-

- Accuracy is the fraction of predictions our model got right.

$$\text{Accuracy} = \frac{\text{No. of correct prediction}}{\text{Total number of prediction.}}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

Precision :-

- It is the ratio of correctly positive samples (True positive) to a total number of classified positive samples (either correct or incorrect).

$$\text{Precision} = \frac{TP}{TP + FP}$$

Recall :- (sensitivity)

- The recall is calculated as the ratio b/w the number of positive samples correctly classified as positive to the total number of positive samples.

$$\text{Recall} = \frac{TP}{TP + FN}$$

Fscore :-

- Precision and Recall are combined into a single measure to form F measure (F-score).

$$\text{Fscore} = \frac{2 * \text{precision} * \text{recall}}{\text{precision} + \text{recall}}$$

what is cross validation in ML?

Cross Validation :-

Cross validation is the technique for validating the model efficiency.

- It is a technique how a statistical model generalized to an independent dataset.
- Methods used for cross validation :-

- ① validation set approach
- ② K-fold cross validation
- ③ leave-P out cross validation
- ④ leave one out cross validation

Q. Explain terms Micro Average - precision, Recall & Fscore.

→ Micro-Average precision :-

- micro precision measure the precision of the aggregated contributions of all classes.
- It is micro average precision.
- micro averaging is used when a problem has more than 2 classes that can be true

$$\text{micro precision} = \frac{(TP) \text{ sum}}{(TP) \text{ sum} + (FP) \text{ sum}}$$

Micro average Recall :-

- Ratio betw. the no. of tve samples correctly classified as positive to the total number of positive samples.

$$\text{Recall} = \frac{TP}{TP + FN}$$

Micro F1 Score :-

- micro F1 score is used to assess the quality of multilabel / binary problems
- It measures the F1-score of the aggregated contributions of all classes.

- Micro F1 score is defined as the harmonic mean of the precision & recall.

$$\text{Micro F1-Score} = \frac{2}{\left[\frac{(\text{micro precision}) \cdot (\text{micro Recall})}{(\text{micro precision}) + (\text{micro Recall})} \right]}$$

Q. Explain terms macro average precision, Recall & F-score?

→ Macro Average precision:

- we calculate the precision for each class separately in an one vs all way.
- Then we take average of all precision value.

E.g. 3 classes a, b, c.

we calculate P_a , P_b , P_c

macro average will be, $\frac{P_a + P_b + P_c}{3}$

Macro Average Recall:

- macro recall measures the average recall per class.
- Recall is a metric used in binary classification problem

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

Macro - Average F1 score :-

- The macro F1 score is computed using the arithmetic mean of all the per class F_1 Scores.
- It is used for overall performance.
- Macro average f score is the harmonic mean.

$$F_1 \text{ Score} = \frac{2}{\left[\frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} \right]}$$

Unit: 5

Unit V: Unsupervised Learning

Q. Define Clustering. Explain k-means clustering algorithm.

Clustering:-

- Clustering is a technique in which the data points are arranged in similar group dynamically without any pre-assignment of groups.
- Cluster means group of the similar things.

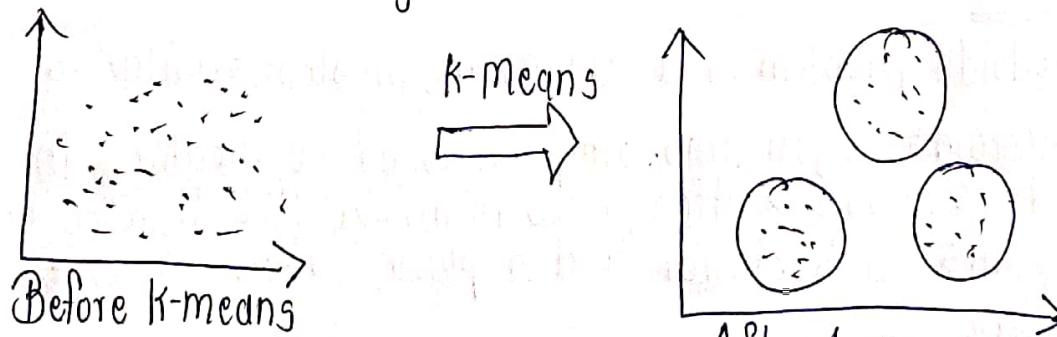
* k-means Clustering Algorithm :-

- It is unsupervised learning algorithm.
- k-means is one of the popular clustering algorithm.
- It helps to form clusters from the given dataset.
- Here $k \rightarrow$ defines the number predefined cluster.

Working :

1. Determine the number of k to decide the no. of cluster.
2. Select random k points or centroids.
3. Assign each data points to nearest centroid.
4. Calculate the variance & place new centroid of each cluster
5. Repeat Step 3.
6. If any reassignment occurs, then go to step 4. else, go to END.
7. The model is ready.

FINISH



Advantages :- Relatively efficient

- Adopt easily to new example.

Disadvantages :-

- Difficult to predict k -value.
- With Global Cluster didn't work well.

Problem: Apply k-means algorithm on given data for $k=3$

$C_1(2)$, $C_2(15)$ & $C_3(38)$ as initial cluster centres.

Use Explain
Hindi

Data: 2, 4, 6, 3, 31, 12, 15, 16, 38, 35, 19, 21, 23, 25, 30.

$$C_1 = 2, C_2 = 15, C_3 = 38$$

The numbers which are close to mean are grouped into resp. clusters.

$$k_1 = \{2, 4, 6, 3\}, k_2 = \{12, 15, 16, 14, 21, 23, 25\}, k_3 = \{31, 35, 30\}$$

Again calculate new mean for new clusters group.

$$C_1 = 3.75, C_2 = 18, C_3 = 32$$

New clusters

$$k_1 = \{2, 4, 6, 3\}, k_2 = \{12, 15, 16, 14, 21, 23, 25\}, k_3 = \{31, 35, 30\}$$

$$C_1 = 3.75, C_2 = 18, C_3 = 32$$

clusters remains unchanged.

final clusters.

$$k_1 = \{2, 4, 6, 3\}, k_2 = \{12, 15, 16, 14, 21, 23, 25\}, k_3 = \{31, 35, 30\}$$

Q. Explain the concept of k-medoids.

→ k-medoids :-

- k-medoids problem is a clustering problem similar to k-means.
- k-medoids algorithms are partitioned i.e breaking the dataset up into group & attempt to minimize the distance between the points in a cluster and a point which is center of that cluster.
- k-medoids chooses actual data points as center.
- It is also called Partitioning Around Medoid.

$$E = |P_i - C_j|$$

medoids

medoids

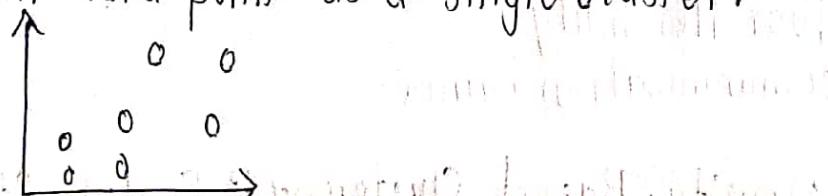
Explain Hierarchical Clustering Algorithm.

+ Hierarchical Clustering :-

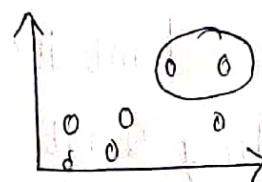
- It is unsupervised machine learning algorithm.
- Hierarchical clustering is also called as Hierarchical Cluster Analysis (HCA).
- Group of similar objects into groups is called cluster.
- Two types of HC:
 - ① Divisive: (Top-down) Also called as DIANA,
 - ② Agglomerative (Bottom-up): Also called as AGNES.
- We represents the hierarchy of cluster in the form of tree called as dendrogram.

Steps:

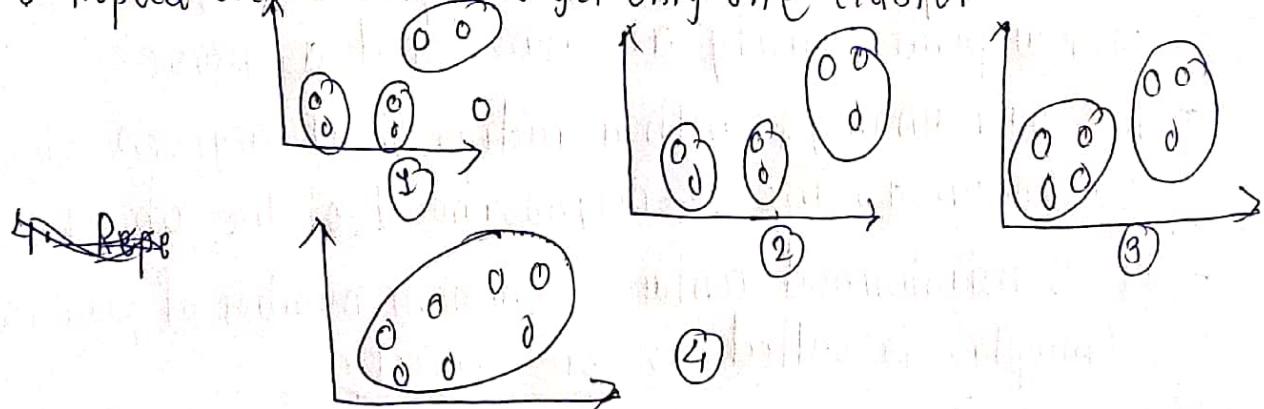
1. Create each data point as a single cluster.



2. Merge two nearest clusters.



3. Repeat step 2 until we get only one cluster.



- 4 Once all cluster are combined into one big cluster, then develop the dendrogram to divide the cluster as per problem.

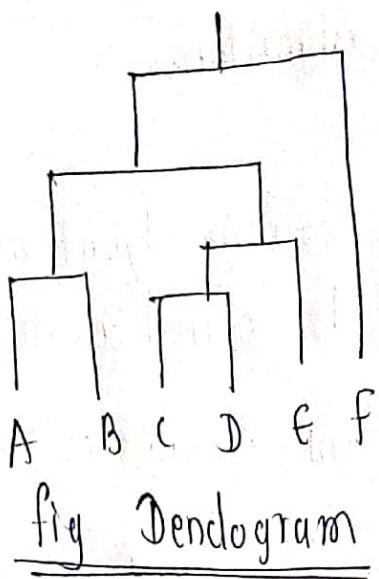


fig Dendrogram

Advantages :

- Easy to implement.
- clear path for advancement.

Disadvantages :

- Poor flexibility
- communication barrier.

Q. What is Density-Based Clustering ? Explain its working ?
 → Density-Based Clustering :

- It is most popular unsupervised ml methodologies.
- It is used for model building & ML algorithms.
- The data points in the region separated by two clusters of low point density are considered as noise.
- The surroundings with a radius ϵ of a given object are known as the ϵ -neighbourhood of the object.
- If ϵ -neighbourhood contains minimum number of points objects ($minpts$) is called as core object.

Working :-

1. Let set of object denoted by \mathcal{D} , object i is directly reachable from the object j only if it is located within the ϵ -neighborhood of j and j is a core object.

2. An object i is density reachable from the object J with respect to ϵ & minpts in a given set of objects.

$i \dots j$

$i = j$

$p_n = i$

such that i_{i+1} is directly reachable from i &

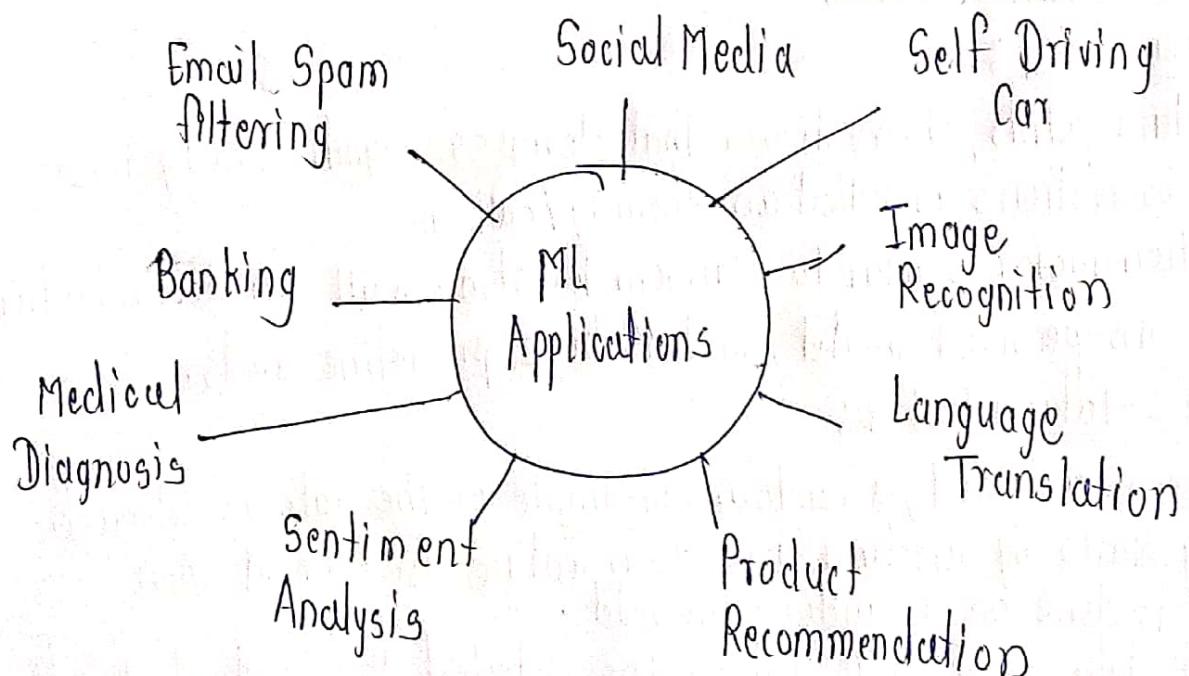
minpts.

3. O is belonging to D so i & j are density reachable from O with respect to ϵ and minpts.

Features of Density Based Clustering:

- It is scan method.
- It required density parameters as a termination condn,
- It is used to manage noise in data clusters.
- DBC is used to identify clusters of arbitrary size.

Q. Write short note on ML applications.



Q. What is outlier analysis.

- Outlier analysis is a fundamental issue in data mining.
- It is used to detect & remove anomalous object from data mining.
- The approach to detect outlier includes three methods.
 - i) Clustering:
 - k-means used
 - Partitions the data set into given number of clusters.
 - ii) Pruning:
 - It is based of on distance measure,
 - iii) Computing Outlier score:
 - for unpruned points this method used.
 - It is based on LDOF (Local distance based outlier factor)
 - LDOF tells how much a point is deviating from its neighbors.

Q. What are isolation factor.

→ Isolation factor :-

- Any data points/observations that deviates significantly from other observations is called an Anomaly/outlier.
- Isolation factor similar to Random forest are built on decision tree.
- It is unsupervised model, so don't have pre-defined labels.

Working of Isolation factor:

1. When given a dataset, a random subsample of the data is selected.
2. Branching starts by selecting a random feature first. And then branching is done on a random threshold.
3. Value of data point is less than the selected threshold, it goes to the left branch otherwise to the right. And thus a node is split into left & right branches.
4. This process from Step 2 is continued recursively till each data point is completely isolated.

Explain local Outlier factor. Mention its advantages & disadvantages

→ Local Outlier factor :-

- Local outlier factor (LOF) is an algorithm used for unsupervised outlier detection.
- It produce an anomaly score that represents data points which are outlier in the dataset.
- It does this by measuring the local density deviation.

Advantages :-

- Sometime it is difficult to find outliers. because of lack of training data. But for LOF is not required prior example.
- LOF can be applied in many other fields to solve problems like video streams, geographic data, etc.
- LOF work well for anomaly detection algorithm for multimodel distribution.

Disadvantages :-

- LOF score not always same it might vary for different datasets.
- In higher dimensions, accuracy of LOF algorithm get effected.

Q. Explain different evaluation metrics & score.

→ 1 Confusion matrix :-

		AV	
		P	N
PV	P	TP	FP
	N	FN	TN

$$\bullet \text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\bullet \text{Precision} = (\text{Positive Predictive Value}) \frac{TP}{TP + \cancel{TN} FP}$$

* Recall or sensitivity = $\frac{TP}{TP+FN}$

2. F1 Score :

F1 score is harmonic mean of precision & recall values.

$$F_1 = 2 \left[\frac{(\text{Precision}) \cdot (\text{recall})}{(\text{Precision}) + (\text{Recall})} \right]$$

3. RMSE

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2}$$

$$\text{RMSF} = \sqrt{\text{MSE}}$$

4. R-squared:

$$R^2 = 1 - \frac{\text{MSE}(\text{model})}{\text{MSE}(\text{base line})}$$

Q. Explain Elbow Method.

→ Elbow method :-

- In cluster analysis, the elbow method is a heuristic used in determining the number of clusters in a dataset.

Elbow Method Calculation:

1. Compute clustering algorithm for different values of k,
2. For each k calculate the total WSS.
3. Plot the curve of WSS according to the no. of clusters k.

- Elbow method does not always work well, if data is not very much clustered. Dataset ~~need not~~ did not have a clear elbow.
- We can have a fairly smooth curve, it is unclear what is the best value of k.

Difference between Intrinsic & extrinsic methods

Intrinsic Motivation	Extrinsic Motivation
1. Known as internal motivation.	1. Known as external motivation.
2. Keep you productive.	2. Just play with your mind.
3. Useful for whole life.	3. Can only be used for temporary purpose.
4. Makes your decision making strong.	4. Makes your instant decision may be bad.
5. Hard to find.	5. Easy to get.
6. Example: Your family Happiness.	6. Example: YouTube channel likes.

Unit : 6

Unit VI: Introduction to Neural Networks

Q. Explain Artificial Neural Networks?

→ Artificial Neural Networks :-

- ANN is a massively distributed processor made up of simple processing units which store experiential knowledge & make it available for use.
- Important features of ANN:-
 - ① Knowledge is required from its environment through a learning process.
 - ② Interneuron connections store the acquired knowledge.
 - ③ Artificial Neural Networks learns by examples.
- Advantages:-
 - ① Parallel processing capability,
 - ② Storing data on entire network,
 - ③ Capability to work with incomplete knowledge,
 - ④ Having memory distribution,
 - ⑤ Having a fault tolerance.

Q. Explain Single layer neural networks.

→ Single layer Neural Networks;

- Single layer neural network contains input & output layer.
- The input layer receives the input signals and the output layer generates the output signals accordingly.

Let us consider a Single Layer Network

Now, the separating line is given by,

$$b + x_1 w_1 + x_2 w_2 = 0$$

if $w_2 \neq 0$ then:

$$w_2 = -\frac{w_1}{w_2}x_1 - \frac{b}{w_2}$$

Thus the requirement for the positive response is:

$$b + w_1x_1 + w_2x_2 > 0.$$

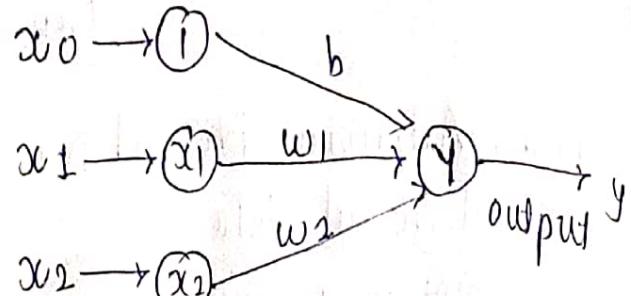
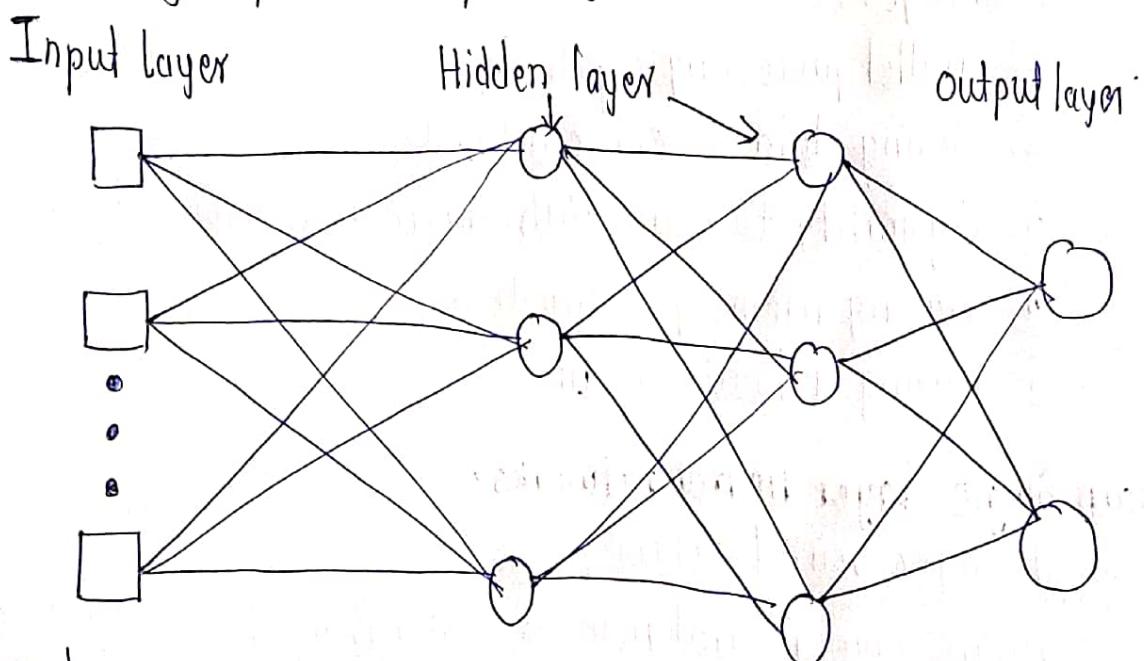


fig. SLNN

Q. Explain multilayer perceptron.

→ Multilayer Perceptron:

- This network is made up of multiple layers.
- The architecture of this class consists of one or more intermediate layers called as hidden layers, besides having input & output layer.



Advantages:

1. Can be applied to non-linear complex problem.
2. Works well with large input data.
3. Provides quick predictions after training.

Disadvantages :

1. Computations are difficult & time consuming.
2. functioning of model depends on training data.

Q. Explain Back Propagation Network (BPN) / learning.

→ BPN :

- Backpropagation is a widely used algorithm for training feedforward neural network.
- The network connected to back-propagation learning algorithm are also called as back-propagation network (BPN).

Working of BPN :

- ① It compares generated output to the desired output & generate error report if the result does not match the generated output vector.
- ② Then it adjust the weights according to bug report then get desired output.

Algorithm :

Steps:

1. Inputs X , arrive through the perconnected path.
2. The input is modeled using true weights, weights chosen randomly.
3. Calculate the output of each neuron.
4. Calculate error in output.

$$\text{Backpropagation Error} = \text{Actual Output} - \text{Desired Output}$$

5. From the output layer, go back to the hidden layer to adjust the weights to reduce the error.
6. Repeat the process until desired output is achieved.

Types of Backpropagation:

1. Static Propagation: Designed to map static input for static output.
2. Recurrent Backpropagation: It is also called as recursive backpropagation, used for fixed-point learning.

Advantages:-

- It is simple, easy, fast to program.
- Flexible and efficient.

Disadvantage:-

- Performance depends on input data.
- Spend too much time to training.

Q. Explain function link Artificial Neural Network (FLANN)?

→ FLANN :-

- Functional link artificial neural network (FLANN) is a single layer ANN.
- It has low computational complexity, which is used for different field of applications, such as system identification, pattern recognition, etc.
- It is used to model the reln b/w input & output variables.

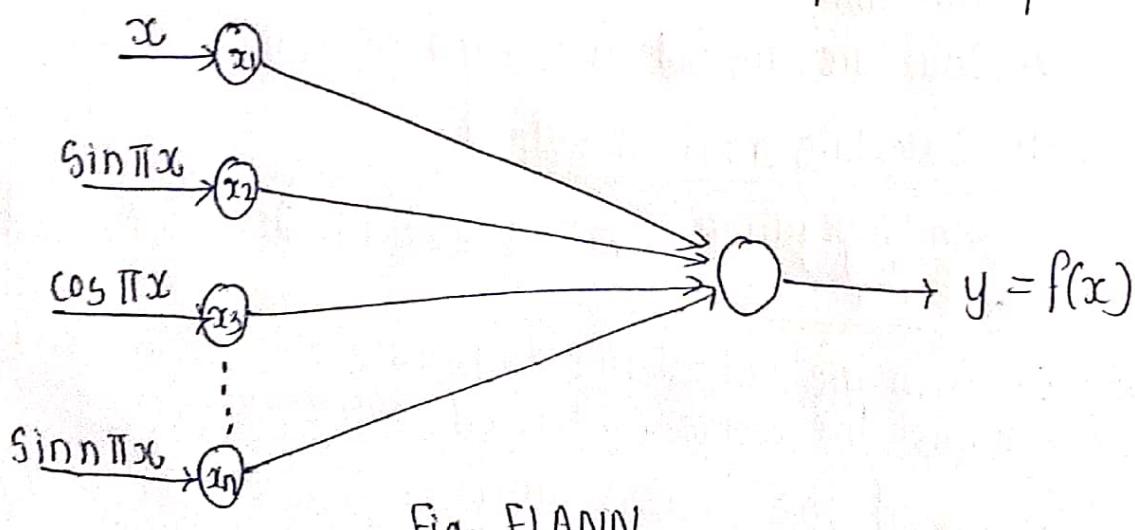
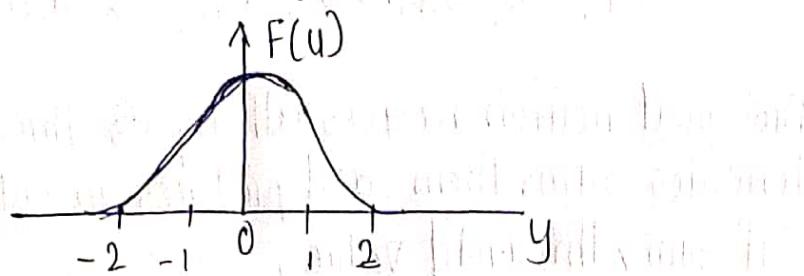


Fig. FLANN.

Q. Explain Radial Basis Function (RBF) Network?

- - The Radical Basis function is a functional approximation neural network.
- It was developed by Powell.
- The network uses sigmoidal & Gaussian kernel Function.
- The response of Gaussian function is positive for all values of y & response decreases as $|y| \rightarrow 0$.
- Gaussian function $F(y) = e^{-y^2}$



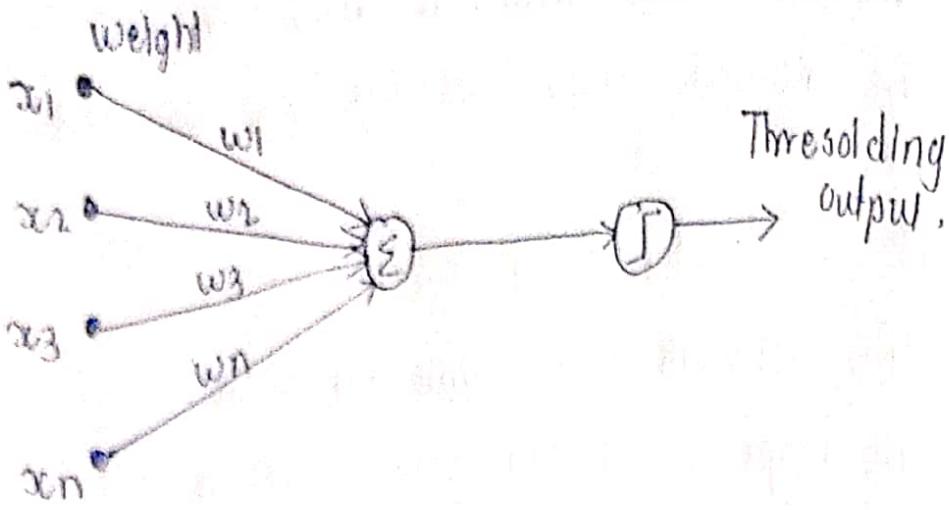
Advantages:

- Easy design.
- Good Generalization.
- Faster Training.
- Only one hidden layer.

Q. Explain Common activation functions used in neural network.

→ Activation function :-

- The activation function decides whether a neuron should be activated or not calculating the weighted sum & further adding bias to it.
- The purpose of the activation function is to introduce non-linearity into the output of a neuron.



- Here, $x_1, x_2, x_3, \dots, x_n \rightarrow n\text{-inputs}$.
 $w_1, w_2, w_3, \dots, w_n \rightarrow \text{weights}$.
- Biological neuron receives all inputs through the dendrites, sums them and produces an output.
 if sum > threshold value,
 then, Input signal pass onto cell body
 through synapse.

- $I = w_1x_1 + w_2x_2 + \dots + w_nx_n$

$$I = \sum_{i=1}^n w_i x_i$$

- Output:

$$y = f(I)$$

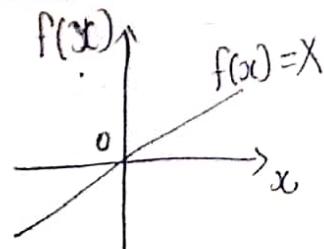
called Activation Function, or Transfer Function, Squared Function.

Different Activation function used in neural network;

① Linear function:

$$f(x) = x \quad \text{for all } x$$

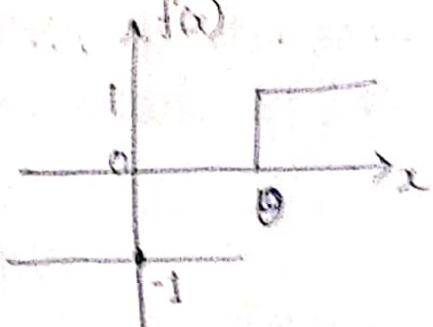
~~Here Output = Input~~ Input = Output.



② Bipolar Step function

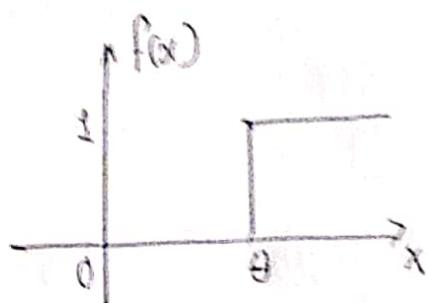
$$f(x) = \begin{cases} 1 & \text{if } x > 0 \\ -1 & \text{if } x < 0 \end{cases}$$

$\theta \rightarrow$ threshold value



③ Binary step function

$$f(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x < 0 \end{cases}$$



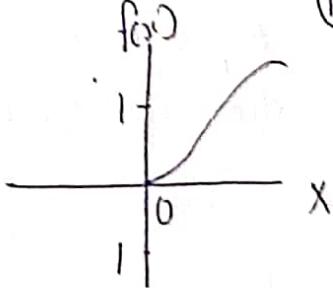
④ Sigmoidal function

$$S(x) = \frac{1}{1 + e^{-x}}$$

- Used for back-propagation nets.

- Two types:-

i) Binary Sigmoidal function (unipolar, logistic)



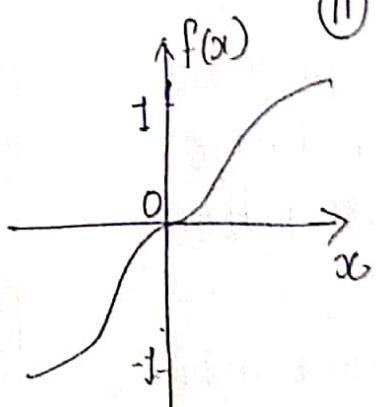
$$f(x) = \frac{1}{1 + e^{-\lambda x}} \quad \lambda \rightarrow \text{Slope parameter}$$

Range \rightarrow 0 to 1.

ii) Bipolar Sigmoid function

$$f(x) = \frac{1 - e^{-\lambda x}}{1 + e^{-\lambda x}}$$

Range \rightarrow -1 to +1



Q. Difference between ANN, RNN (Recurrent Neural Network) & CNN (convolution neural Network),

	ANN	RNN	CNN
1. Type of Data	Tabular data, Text data.	Sequence data	Image data.
2. Parameter sharing	No	Yes	Yes
3. Fixed length input	Yes	No	Yes
4. Recurrent connections	No	Yes	No
5. Vanished & exploding gradient	Yes	Yes	Yes.
6. Spatial reln	No	No	Yes
7. Performance	less powerful than CNN, RNN	less feature than CNN	more powerful than ANN, RNN
8. Application	Facial recognition	Text to Speech	Facial recognition.
9. Advantage	Having fault tolerance, ability to work with incomplete knowledge.	remembers time series prediction	High accuracy in image recognition
10 Disadvantage	Hardware dependency.	Gradient vanished, exploding gradient	Large training data needed.

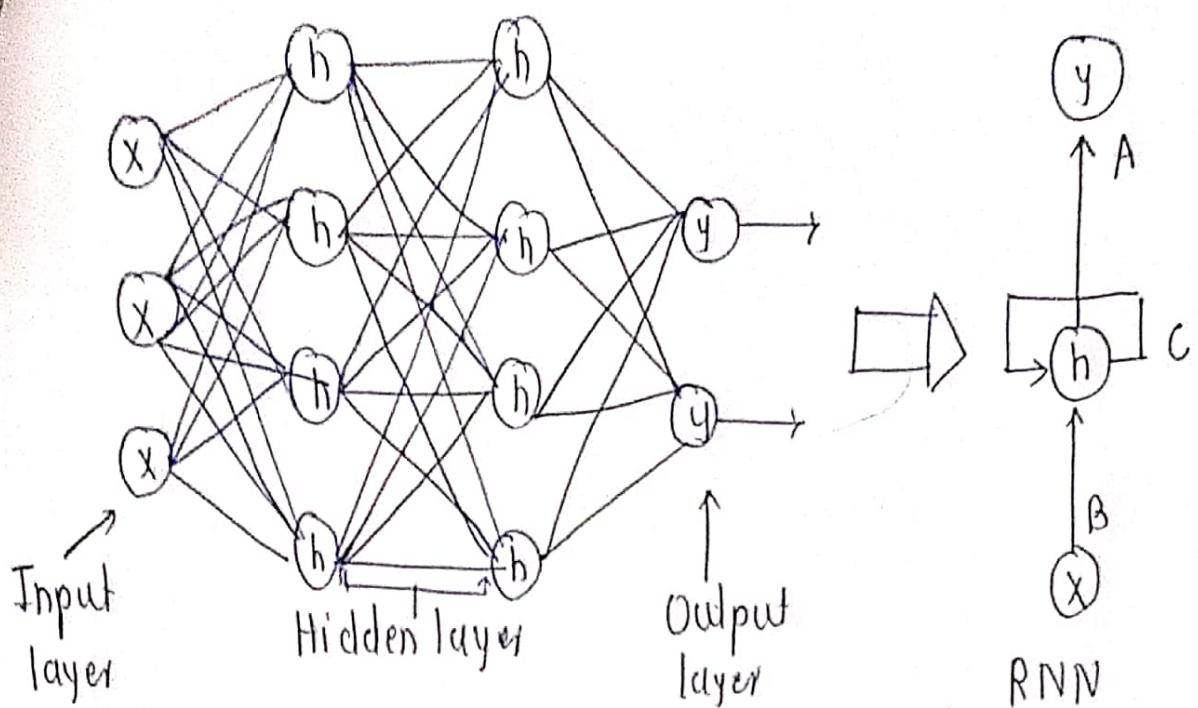


fig. RNN

RNN
(Recurrent Neural Network)

Thank you

~Prof. Mayur Ghorpade
~Prof. Suraj Sonawane