

Supervised Learning

Shahadat Hoshen

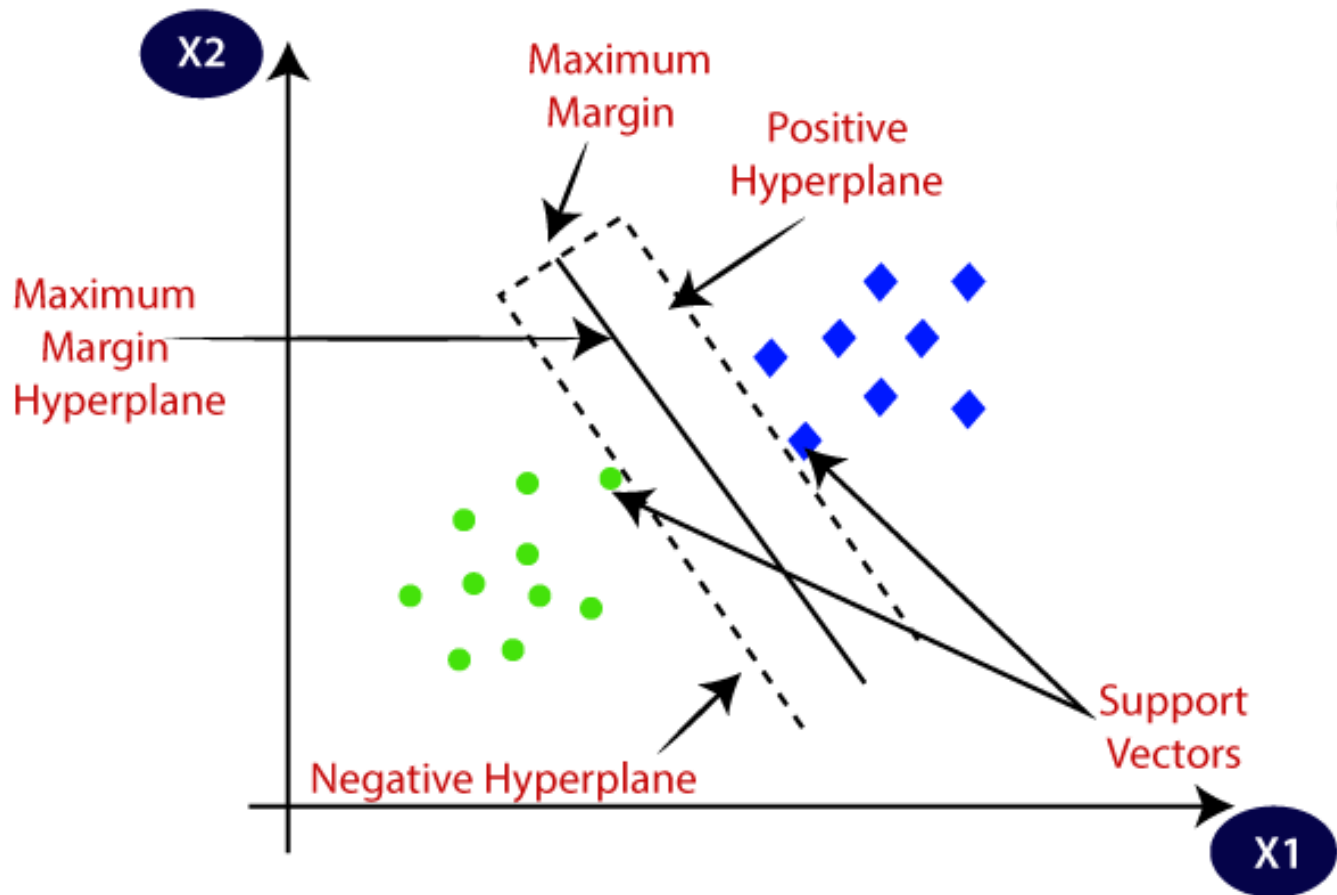
Lecturer,

Dept. of CSE, NUBTK

Support Vector Machines

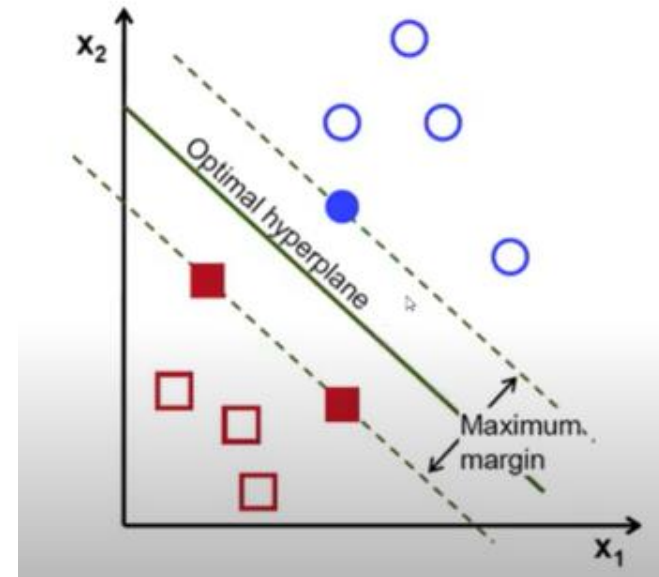
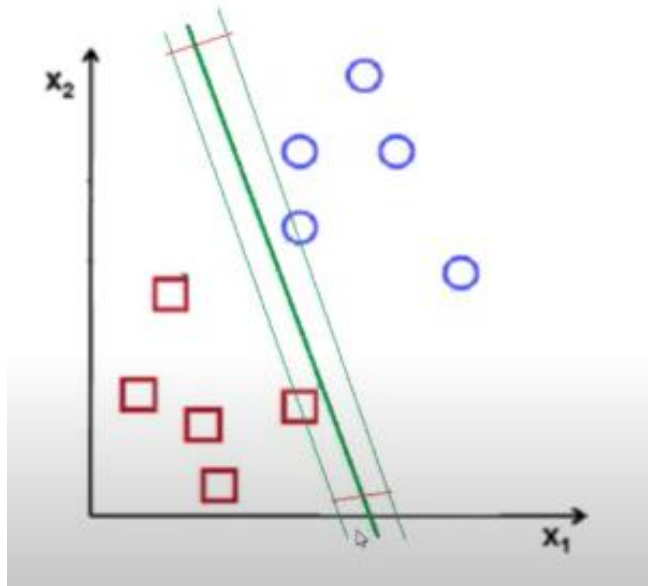
- ▶ Support Vector Machines, a supervised machine learning algorithm used for classification and regression tasks and effective in dealing with **complex, high-dimensional** datasets.
- ▶ SVM performs classification by first transforming the training records into higher dimensions. Then, within the new dimension, it searches for the best decision boundary that separates the training records of one class from others.
- ▶ The main objective of SVMs is to find the **best hyperplane** that separates the data points of different classes in a way that maximizes the margin, or the distance, between the classes. This hyperplane serves as the decision boundary for classifying new, unseen data points

Support Vector Machines



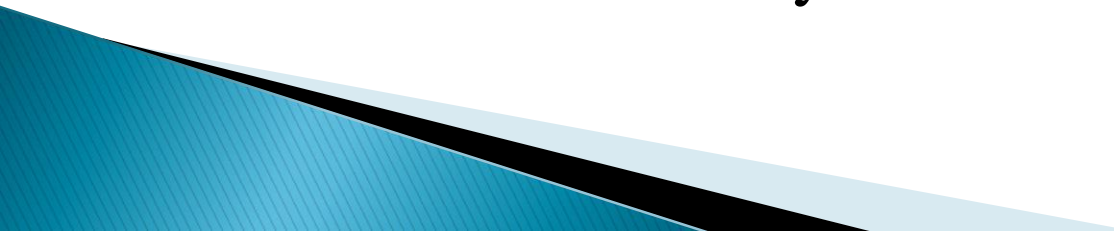
- ▶ Maximum marginal hyperplane is calculated from the records that fall in the hyperplane pair (called “Support Vectors”) to be used as the decision boundary.

Support Vector Machines



- ▶ From the Figure, we see that the records distributed between two class values are linearly separable; however, there can be an infinite number of lines separating them.
- ▶ The target of an SVM is to find the best line (best decision boundary) that will help minimize classification error for unseen records.

Non-linear SVM

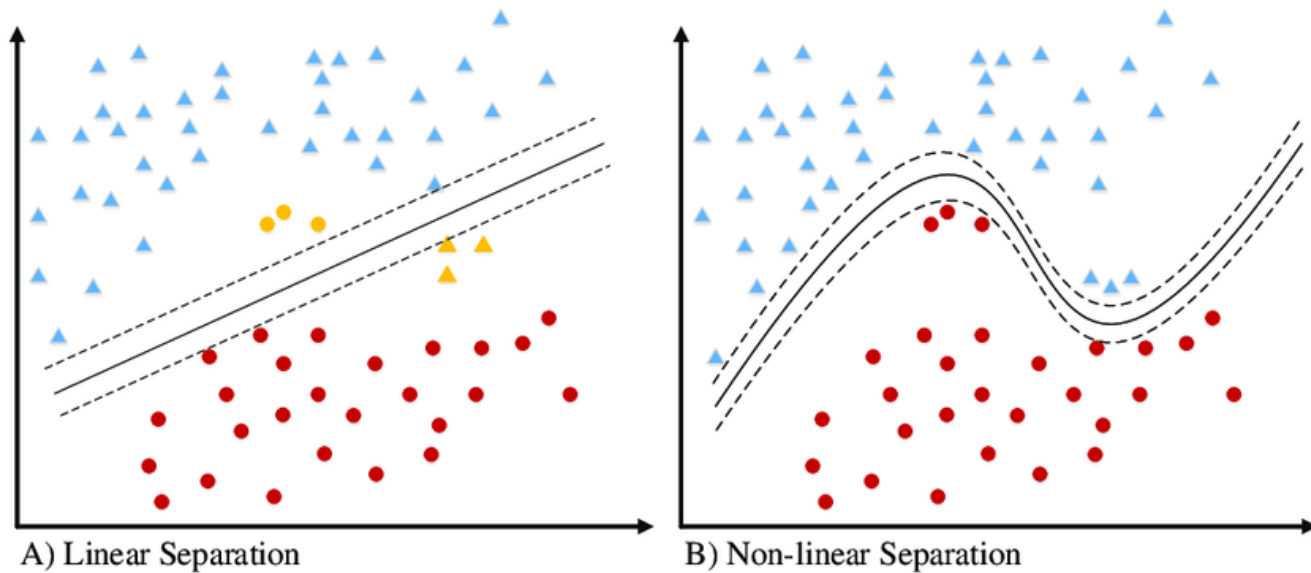
- ▶ Non-linear SVMs extend the concept of SVMs to handle non-linear decision boundaries by mapping the input features into a higher-dimensional space.
 - ▶ This is done through a process called the kernel trick.
 - ▶ The kernel trick allows the algorithm to compute the dot product of the data points in this higher-dimensional space without explicitly calculating the transformation, saving computational resources.
 - ▶ The choice of kernel function determines the shape of the decision boundary.
- 

Different types of kernel

1. Linear Kernel
2. Polynomial Kernel
3. Radial Basis Function (RBF) or Gaussian Kernel

Linear vs Non-linear SVM

When we can easily separate data with hyperplane by drawing a straight line is **Linear SVM**. When we cannot separate data with a straight line we use **Non-Linear SVM**.



Linear vs Non-linear SVM

Linear SVM	Non-Linear SVM
It can be easily separated with a linear line.	It cannot be easily separated with a linear line.
Data is classified with the help of hyperplane.	We use Kernels to make non-separable data into separable data.
Data can be easily classified by drawing a straight line.	We map data into high dimensional space to classify.

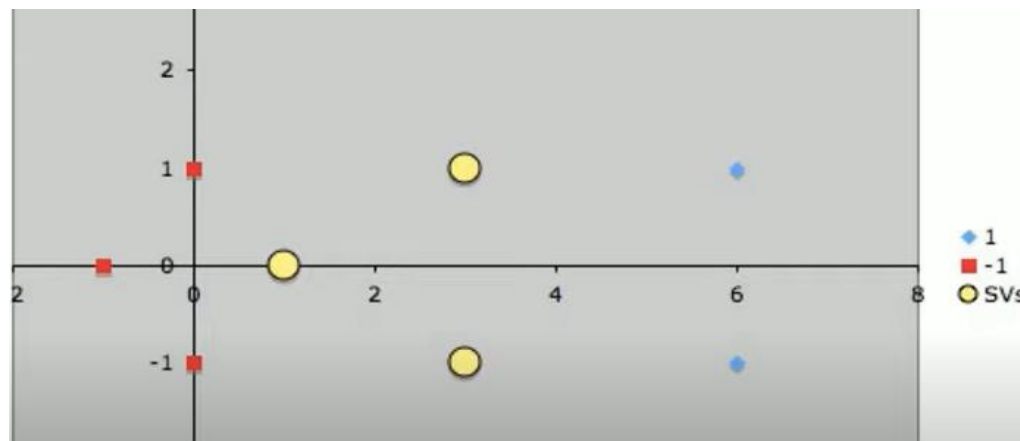
Problem:1 (Linear SVM)

Here are the data point:

$$\begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 1 \end{pmatrix}, \begin{pmatrix} 6 \\ -1 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix}$$

Find out the appropriate decision boundary among these data.

Solution:



Solution

By inspection, it should be obvious that there are **three** support vectors,

$$\left\{ s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, s_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, s_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix} \right\}$$

Each vector is augmented with a 1 as a bias input

$$\text{So, } s_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \text{ then } \tilde{s}_1 = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

Similarly,

$$s_2 = \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \tilde{s}_2 = \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \quad s_3 = \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \tilde{s}_3 = \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix}$$

Linear SVM Solve

$$\alpha_1 \tilde{s}_1 \cdot \tilde{s}_1 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_1 + \alpha_3 \tilde{s}_3 \cdot \tilde{s}_1 = -1$$

$$\alpha_1 \tilde{s}_1 \cdot \tilde{s}_2 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_2 + \alpha_3 \tilde{s}_3 \cdot \tilde{s}_2 = +1$$

$$\alpha_1 \tilde{s}_1 \cdot \tilde{s}_3 + \alpha_2 \tilde{s}_2 \cdot \tilde{s}_3 + \alpha_3 \tilde{s}_3 \cdot \tilde{s}_3 = +1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} = -1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} = 1$$

$$\alpha_1 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_2 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} + \alpha_3 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} = 1$$

Linear SVM Solve

$$\alpha_1(1 + 0 + 1) + \alpha_2(3 + 0 + 1) + \alpha_3(3 + 0 + 1) = -1$$

$$\alpha_1(3 + 0 + 1) + \alpha_2(9 + 1 + 1) + \alpha_3(9 - 1 + 1) = 1$$

$$\alpha_1(3 + 0 + 1) + \alpha_2(9 - 1 + 1) + \alpha_3(9 + 1 + 1) = 1$$

$$2\alpha_1 + 4\alpha_2 + 4\alpha_3 = -1$$

$$4\alpha_1 + 11\alpha_2 + 9\alpha_3 = 1$$

$$4\alpha_1 + 9\alpha_2 + 11\alpha_3 = 1$$

$$\alpha_1 = -3.5$$

$$\alpha_2 = 0.75$$

$$\alpha_3 = 0.75$$

Linear SVM Solve

$$\begin{aligned}\tilde{w} &= \sum_i \alpha_i \tilde{s}_i \\ &= -3.5 \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ 1 \\ 1 \end{pmatrix} + 0.75 \begin{pmatrix} 3 \\ -1 \\ 1 \end{pmatrix} \\ &= \begin{pmatrix} 1 \\ 0 \\ -2 \end{pmatrix}\end{aligned}$$

We can equate the last entry in \tilde{w} as the hyperplane offset b and write the separating

Hyperplane equation $y = wx + b$

$$w = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, b = -2$$

Linear SVM Solve



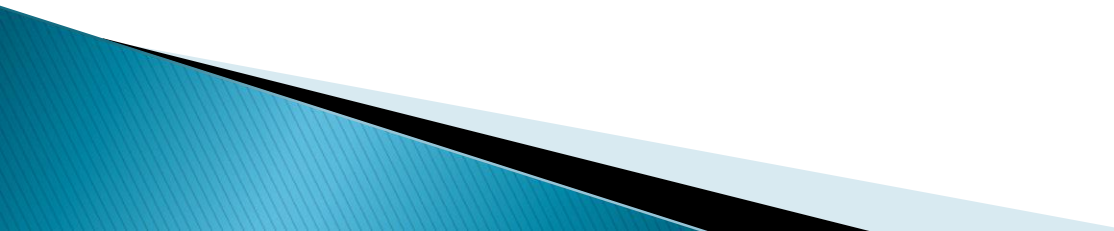
Problem: 2

Here are some data point:

$$\begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 2 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} 2 \\ -1 \end{pmatrix}, \begin{pmatrix} 4 \\ 0 \end{pmatrix}, \begin{pmatrix} 5 \\ 1 \end{pmatrix}, \begin{pmatrix} 5 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 0 \end{pmatrix}$$

Find out the appropriate decision boundary among these data.

Ensemble learning

- ▶ Ensemble learning refers to the technique of combining the predictions of multiple models to improve overall performance and generalization.
 - ▶ The idea is that by aggregating the predictions of diverse models, the ensemble can often achieve better results than any individual model.
 - ▶ Ensemble methods are commonly used in various machine learning tasks, including classification, regression, and anomaly detection.
 - ▶ There are two main types of ensemble learning: **bagging** and **boosting**.
- 

Ensemble learning

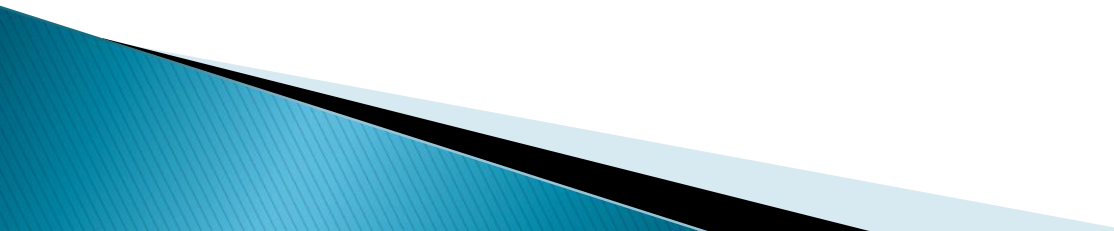
- ▶ **Bagging (bootstrap aggregating)**

Creating a different training subset from sample training data with replacement is called Bagging. The final output is based on majority voting.

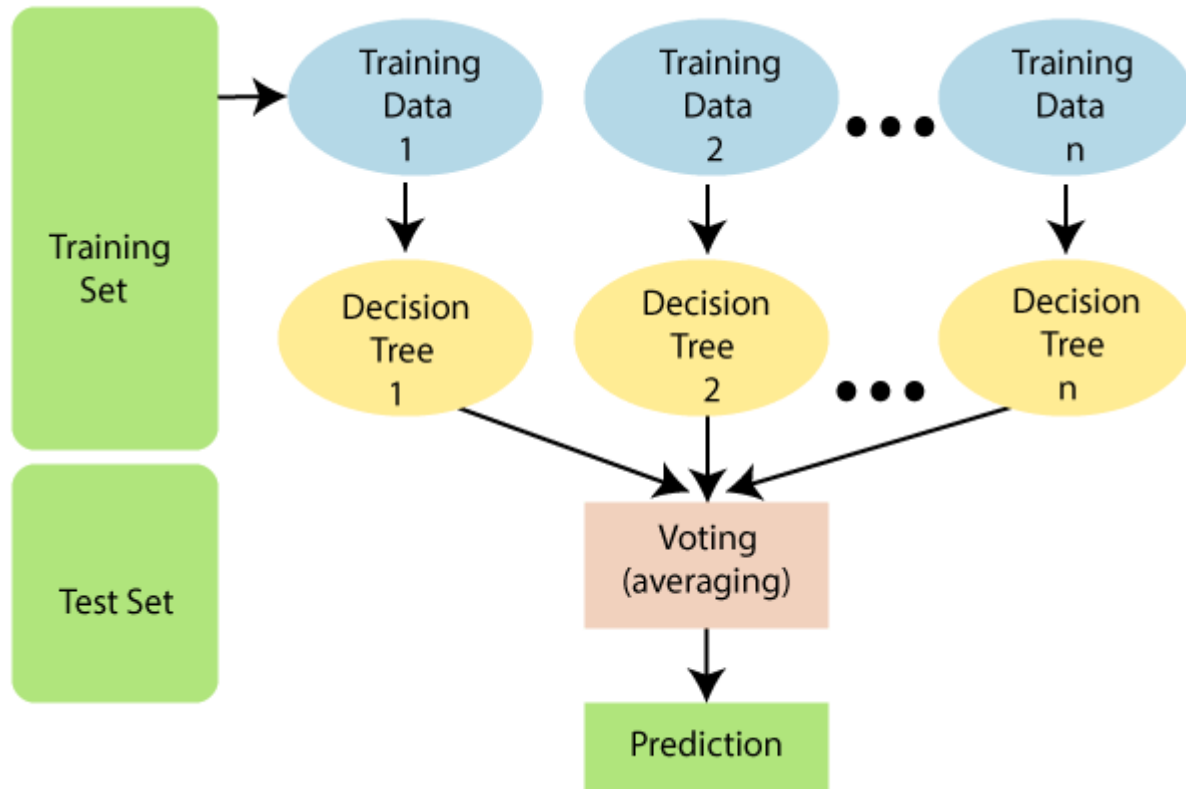
- ▶ **Boosting**

Combining weak learners into strong learners by creating sequential models such that the final model has the highest accuracy is called Boosting. Example: ADA BOOST, XG BOOST.

Random Forest

- ▶ Random Forest is an ensemble learning algorithm that combines the predictions of multiple decision trees to improve the overall accuracy and can handle complex problems.
 - ▶ It contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset.
 - ▶ Rather than depending on one tree it takes the prediction from each tree and based on the majority votes of predictions, predicts the final output. For regression, it's the average of the predictions.
 - ▶ It uses bagging methods.
- 

Random Forest



Regression

- ▶ Regression determines the statistical relationship between a dependent variable and one or more independent variables which are used to predict real or continuous values.
- ▶ The ultimate goal of the regression algorithm is to plot a best-fit line or a curve between the data.
- ▶ There are several types of regression. Linear Regression, Multiple Linear Regression, and Polynomial Regression.
- ▶ In simple words, "Regression shows a line or curve that passes through all the data points on target-predictor graph in such a way that the vertical distance between the data points and the regression line is minimum."
- ▶ The distance between data points and line tells whether a model has captured a strong relationship or not.

Self-study: Classification vs Regression

Regression

- ▶ **Example:** Suppose there is a marketing company A, who does various advertisement every year and get sales on that. The below list shows the advertisement made by the company in the last 5 years and the corresponding sales:

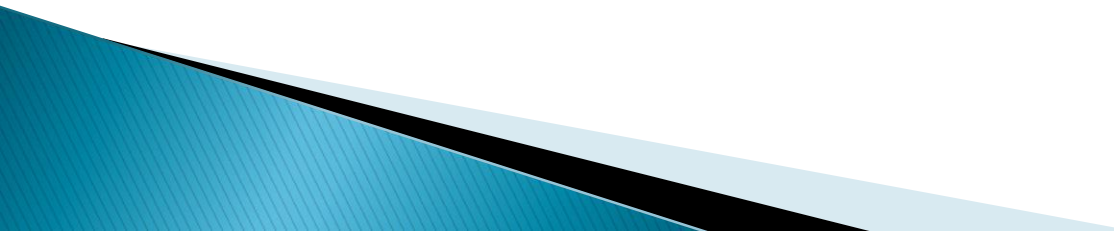
Advertisement	Sales
\$90	\$1000
\$120	\$1300
\$150	\$1800
\$100	\$1200
\$130	\$1380
\$200	??

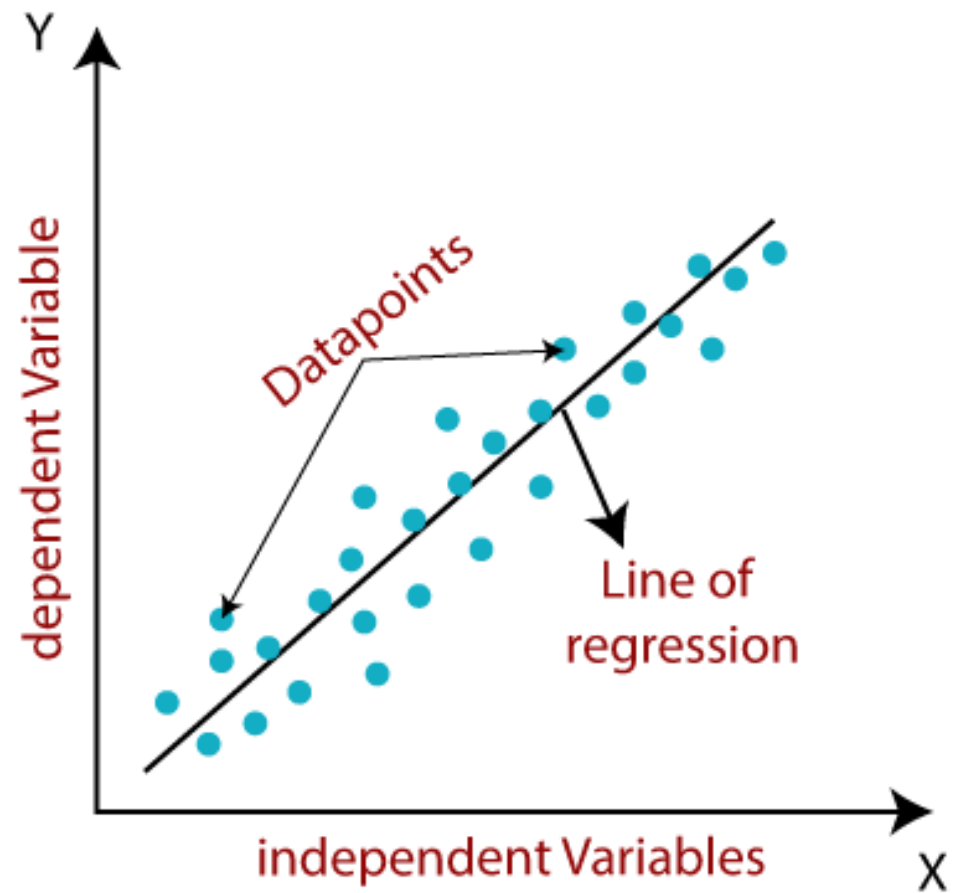
- ▶ Now, the company wants to do an advertisement for \$200 in the year 2019 and wants to know the prediction about the sales for this year. So to solve such type of prediction problems in machine learning, we need regression analysis.

Application of regression

- Forecasting continuous outcomes like house prices, stock prices, or sales.
- Predicting the success of future retail sales or marketing campaigns to ensure resources are used effectively.
- Predicting customer or user trends, such as on streaming services or e-commerce websites.
- Analyzing datasets to establish the relationships between variables and output.
- Predicting interest rates or stock prices from a variety of factors.
- Creating time series visualizations.

Simple Linear Regression

- ▶ Linear regression finds the linear relationship between the dependent variable and one independent variable using a best-fit straight line.
 - ▶ Generally, a linear model predicts by simply computing a weighted sum of the input features, plus a constant called the bias term (also called the intercept term).
 - ▶ In this technique, the dependent variable is continuous, the independent variable(s) can be continuous or discrete, and the nature of the regression line is linear.
- 



Example-1

Consider the following table consisting of the five weeks' sales data.

Apply linear regression technique to predict sales of the 7th and 12th week.

x_i (Week)	y_j (Sales in Thousands)
1	1.2
2	1.8
3	2.6
4	3.2
5	3.8

Linear Equation is: $y = a_0 + a_1 * x$

$$a_1 = \frac{(\overline{xy}) - (\bar{x})(\bar{y})}{\overline{x^2} - \bar{x}^2}$$

$$a_0 = \bar{y} - a_1 * \bar{x}$$

	x_i (Week)	y_j (Sales in Thousands)	x_i^2	$x_i * y_j$
	1	1.2	1	1.2
	2	1.8	4	3.6
	3	2.6	9	7.8
	4	3.2	16	12.8
	5	3.8	25	19
Sum	15	12.6	55	44.4
Average	$\bar{x} = 3$	$\bar{y} = 2.52$	$\overline{x^2} = 11$	$\overline{xy} = 8.88$

$$a_1 = \frac{(\overline{xy}) - (\bar{x})(\bar{y})}{\overline{x^2} - \bar{x}^2} = \frac{8.88 - 3 * 2.52}{11 - 3^2} = 0.66$$

$$a_0 = \bar{y} - a_1 * \bar{x} = 2.52 - 0.66 * 3 = 0.54$$

Regression equation is

$$y = a_0 + a_1 * x$$

$$y = 0.54 + 0.66 * x$$

The predicted 7th week sale (when $x = 7$) is,

$$y = 0.54 + 0.66 * 7 = 5.16$$

the predicted 12th week sale (when $x = 12$) is,

$$y = 0.54 + 0.66 * 12 = 8.46$$

Example-2

SUBJECT	AGE X	GLUCOSE LEVEL Y
1	43	99
2	21	65
3	25	79
4	42	75
5	57	87
6	59	81
7	55	?

By using linear regression, predict the glucose level for age 55.

Multiple Linear regression

- Multiple linear regression refers to a technique that uses two or more independent variables to predict the outcome of a dependent variable.
- It achieves a better fit in the comparison to simple linear regression when multiple independent variables are involved.
- The equation for multiple linear regression is similar to the equation for a simple linear equation, i.e., $y(x) = p_0 + p_1x_1$ plus the additional weights and inputs for the different features which are represented by b_nx_n .
- The formula for multiple linear regression would look like,

$$y(x) = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Multiple Linear regression

- The formula for multiple linear regression would look like,

$$y(x) = b_0 + b_1x_1 + b_2x_2 + \dots + b_nx_n$$

Where,

$$\hat{b}_1 = \frac{(\sum x_2^2)(\sum x_1 y) - (\sum x_1 x_2)(\sum x_2 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$\hat{b}_2 = \frac{(\sum x_1^2)(\sum x_2 y) - (\sum x_1 x_2)(\sum x_1 y)}{(\sum x_1^2)(\sum x_2^2) - (\sum x_1 x_2)^2}$$

$$b_0 = \bar{y} - b_1\bar{x}_1 - b_2\bar{x}_2$$

$$\bullet \sum x_1^2 = \sum X_1^2 - (\sum X_1)^2 / n$$

$$\bullet \sum x_2^2 = \sum X_2^2 - (\sum X_2)^2 / n$$

$$\bullet \sum x_1 y = \sum X_1 y - (\sum X_1 \sum y) / n$$

$$\bullet \sum x_2 y = \sum X_2 y - (\sum X_2 \sum y) / n$$

$$\bullet \sum x_1 x_2 = \sum X_1 X_2 - (\sum X_1 \sum X_2) / n$$

Multiple Linear regression: Example

Suppose we have the following dataset containing the height, width and price of carpets, Predict the price of the carpet whose height 65 and width 23 using multiple linear regression.

Length (x_1)	Width (x_2)	Price (y)
60	22	140
62	25	155
67	24	159
70	20	179
71	15	192
72	14	200
75	14	212
78	11	215

Multiple Linear regression: Example

y	X_1	X_2	X_1^2	X_2^2	X_1y	X_2y	X_1X_2
140	60	22	3600	484	8400	3080	1320
155	62	25	3844	625	9610	3875	1550
159	67	24	4489	576	10653	3816	1608
179	70	20	4900	400	12530	3580	1400
192	71	15	5041	225	13632	2880	1065
200	72	14	5184	196	14400	2800	1008
212	75	14	5625	196	15900	2968	1050
215	78	11	6084	121	16770	2365	858
181.5	69.375	18.125	38767	2823	101895	25364	9859

Sum

1452	555	145	38767	2823	101895	25364	9859
------	-----	-----	-------	------	--------	-------	------

Multiple Linear regression: Solution

- $\Sigma x_1^2 = \Sigma X_1^2 - (\Sigma X_1)^2 / n = 38,767 - (555)^2 / 8 = \mathbf{263.875}$
- $\Sigma x_2^2 = \Sigma X_2^2 - (\Sigma X_2)^2 / n = 2,823 - (145)^2 / 8 = \mathbf{194.875}$
- $\Sigma x_1 y = \Sigma X_1 y - (\Sigma X_1 \Sigma y) / n = 101,895 - (555 * 1,452) / 8 = \mathbf{1,162.5}$
- $\Sigma x_2 y = \Sigma X_2 y - (\Sigma X_2 \Sigma y) / n = 25,364 - (145 * 1,452) / 8 = \mathbf{-953.5}$
- $\Sigma x_1 x_2 = \Sigma X_1 X_2 - (\Sigma X_1 \Sigma X_2) / n = 9,859 - (555 * 145) / 8 = \mathbf{-200.375}$

$$\hat{b}_1 = \frac{(\Sigma x_2^2)(\Sigma x_1 y) - (\Sigma x_1 x_2)(\Sigma x_2 y)}{(\Sigma x_1^2)(\Sigma x_2^2) - (\Sigma x_1 x_2)^2}$$
$$= [(194.875)(1162.5) - (-200.375)(-953.5)] / [(263.875)(194.875) - (-200.375)^2] = \mathbf{3.148}$$

$$\hat{b}_2 = \frac{(\Sigma x_1^2)(\Sigma x_2 y) - (\Sigma x_1 x_2)(\Sigma x_1 y)}{(\Sigma x_1^2)(\Sigma x_2^2) - (\Sigma x_1 x_2)^2}$$
$$= [(263.875)(-953.5) - (-200.375)(1162.5)] / [(263.875)(194.875) - (-200.375)^2] = \mathbf{-1.656}$$

Multiple Linear regression: Solution

$$b_0 = \bar{y} - b_1\bar{X}_1 - b_2\bar{X}_2$$

$$b_0 = 181.5 - 3.148(69.375) - (-1.656)(18.125) = \mathbf{-6.867}$$

The estimated linear regression equation is: $\hat{y} = b_0 + b_1x_1 + b_2x_2$

In our example, it is $\hat{y} = \mathbf{-6.867 + 3.148x_1 - 1.656x_2}$

So, for height 65 and width 23 the price will be,

$$\hat{y} = -6.867 + (3.14 \times 65) - (1.65 \times 23) = \mathbf{159.665}$$

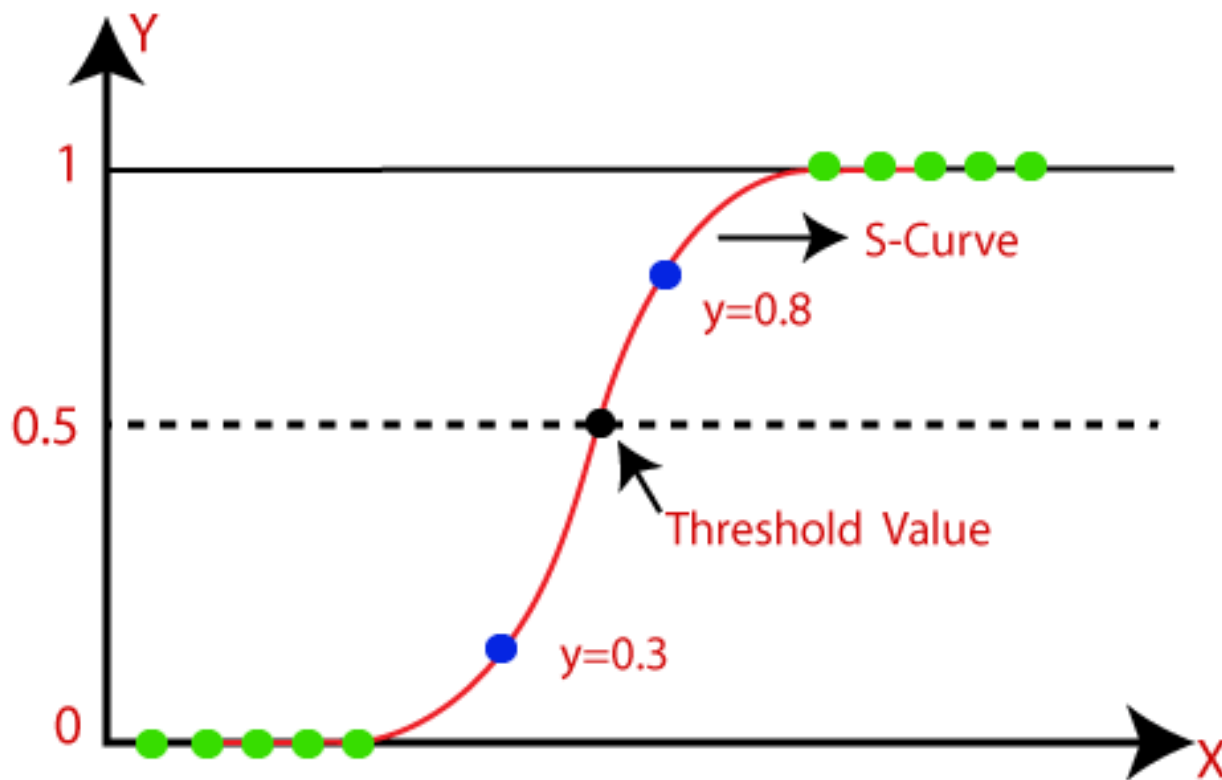
Logistic Regression

- ▶ Logistic regression is a supervised machine learning algorithm mainly used for classification tasks where the goal is to predict the probability that an instance belongs to a given class or not.
- ▶ It is a statistical method that is used for building machine learning models where the dependent variable is dichotomous: i.e. binary which works with categorical variables such as 0 or 1, Yes or No, True or False, Spam or not spam, etc.
- ▶ It has the ability to provide probabilities and classify new data using continuous and discrete datasets.

Logistic Regression

- ▶ Logistic regression uses the **sigmoid function** which is a mathematical function used to map the predicted values to probabilities.
- ▶ It maps any real value into another value within a range of 0 and 1 and forms a curve like the "S" form that curve is called the Sigmoid function or the logistic function.
- ▶ Probability is either 0 or 1, depending on whether the event happens or not.
- ▶ For binary predictions, you can divide the population into two groups with a cut-off of 0.5. If the hypothesis is above 0.5 is considered to belong to group A, and below is considered to belong to Group B.

sigmoid function: $p(x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$

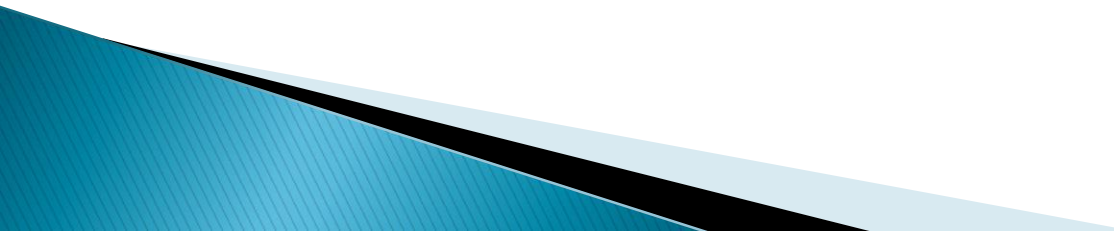


Type of Logistic Regression

On the basis of the categories, Logistic Regression can be classified into three types:

1. **Binomial:** In binomial Logistic regression, there can be only two possible types of dependent variables, such as 0 or 1, Pass or Fail, etc.
2. **Multinomial:** In multinomial Logistic regression, there can be 3 or more possible unordered types of the dependent variable, such as “cat”, “dogs”, or “sheep”. In this case, the softmax function is used in place of the sigmoid function.
3. **Ordinal:** In ordinal Logistic regression, there can be 3 or more possible ordered types of dependent variables, such as “low”, “Medium”, or “High”.

Application of LR

- In health care, logistic regression can be used to predict if a tumor is likely to be benign or malignant.
 - In the financial industry, logistic regression can be used to predict if a transaction is fraudulent or not.
 - In marketing, logistic regression can be used to predict if a targeted audience will respond or not.
- 

Example-1

- The dataset of pass or fail in an exam of 5 students is given in the table.
 - Use logistic regression as classifier to answer the following questions.
1. Calculate the probability of pass for the student who studied 33 hours.
 2. At least how many hours student should study that makes he will pass the course with the probability of more than 95%.

Hours Study	Pass (1) / Fail (0)
29	0
15	0
33	1
28	1
39	1

Assume the model suggested by the optimizer for odds of passing the course is,

$$\log(odds) = -64 + 2 * hours$$

1. Calculate the probability of pass for the student who studied 33 hours.

$$p = \frac{1}{1+e^{-z}} \quad s(x) = \frac{1}{1+e^{-x}}$$

$$z = -64 + 2 * 33 = -64 + 66 = 2$$

$$p = \frac{1}{1+e^{-2}} = 0.88$$

- That is, if student studies 33 hours, then there is **88% chance** that the student will pass the exam

Hours Study	Pass (1) / Fail (0)
29	0
15	0
33	1
28	1
39	1

$$\log(odds) = z = -64 + 2 * hours$$

2. At least how many hours student should study that makes he will pass the course with the probability of more than 95%.

- $p = \frac{1}{1+e^{-z}} = 0.95$
- $0.95 * (1 + e^{-z}) = 1$
- $0.95 * e^{-z} = 1 - 0.95$
- $e^{-z} = \frac{0.05}{0.95} = 0.0526$
- $\ln(e^{-z}) = \ln(0.0526)$

$$\ln(e^x) = x$$

$$-z = \ln(0.0526) = -2.94$$

$$z = 2.94$$

Hours Study	Pass (1) / Fail (0)
29	0
15	0
33	1
28	1
39	1

- $z = 2.94$
- $\log(\text{odds}) = z = -64 + 2 * \text{hours}$
- $2.94 = -64 + 2 * \text{hours}$
- $2 * \text{hours} = 2.94 + 64$
- $2 * \text{hours} = 66.94$
- $\text{hours} = \frac{66.94}{2}$
- **$\text{hours} = 33.47 \text{ Hours}$**

Hours Study	Pass (1) / Fail (0)
29	0
15	0
33	1
28	1
39	1

- The student should study **at least 33.47 hours**, so that he will pass the exam with more than 95% probability

Example-2

- ▶ Here 5 students' data is given, you know the number of hours they studied (x) and whether they passed (1) or failed (0) the exam (y).

Hours Studied (x)	Pass/Fail (y)
2	0
3	0
4	0
5	1
6	1

- ▶ Use logistic Regression to predict whether a student who studies 10 hours will pass or fail.

THANK YOU