



UNIT-II: REGRESSION, BAYESIAN LEARNING, . SUPPORT VECTOR  
MACHINE

Table of Content

- 1) *REGRESSION: Linear Regression*
- 2) *REGRESSION :Logistic Regression*
- 3) *BAYESIAN LEARNING - Bayes theorem*
- 4) *Concept learning*
- 5) *Bayes Optimal Classifier,*
- 6) *Naïve Bayes classifier*
- 7) *Bayesian belief networks*
- 8) *SUPPORT VECTOR MACHINE: Introduction*
- 9) *Types of support vector kernel – (Linear kernel, polynomial kernel, and Gaussian kernel),*
- 10) *Hyperplane – (Decision surface)*
- 11) *Properties of SVM*
- 12) *Issues in SVM.*



## Regression

Regression analysis is a statistical method to model the relationship between a dependent (target) and independent (predictor) variables with one or more independent variables. More specifically, Regression analysis helps us to understand how the value of the dependent variable is changing corresponding to an independent variable when other independent variables are held fixed. It predicts continuous/real values such as temperature, age, salary, price, etc..

Regression algorithms are used to solve regression problems in which there is a linear relationship between input and output variables. These are used to predict continuous output variables, such as market trends, weather prediction, etc.

It predicts the continuous output variables based on the independent input variable. like the prediction of house prices based on different parameters like house age, distance from the main road, location, area, etc.

We can understand the concept of regression analysis using the below example: 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 the advertisement of \$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.

Regression is a supervised learning technique which helps in finding the correlation between variables and enables us to predict the continuous output variable based on the one or more predictor variables. It is mainly used for **prediction, forecasting, time series modeling, and determining the causal-effect relationship between variables**.

In Regression, we plot a graph between the variables which best fits the given datapoints, using this plot, the machine learning model can make predictions about the data. In simple words, "**Regression shows a line or curve that passes through all the datapoints on target-predictor graph in such a way that the vertical distance between the datapoints and the regression line is minimum.**" The distance between datapoints and line tells whether a model has captured a strong relationship or not.

Some examples of regression can be as:

- Prediction of rain using temperature and other factors
- Determining Market trends
- Prediction of road accidents due to rash driving.



### **Terminologies Related to the Regression Analysis:**

- o **Dependent Variable:** The main factor in Regression analysis which we want to predict or understand is called the dependent variable. It is also called target variable.
- o **Independent Variable:** The factors which affect the dependent variables or which are used to predict the values of the dependent variables are called independent variable, also called as a predictor.
- o **Outliers:** Outlier is an observation which contains either very low value or very high value in comparison to other observed values. An outlier may hamper the result, so it should be avoided.
- o **Multicollinearity:** If the independent variables are highly correlated with each other than other variables, then such condition is called Multicollinearity. It should not be present in the dataset, because it creates problem while ranking the most affecting variable.
- o **Underfitting and Overfitting:** If our algorithm works well with the training dataset but not well with test dataset, then such problem is called Overfitting. And if our algorithm does not perform well even with training dataset, then such problem is called underfitting.

### **Why do we use Regression Analysis?**

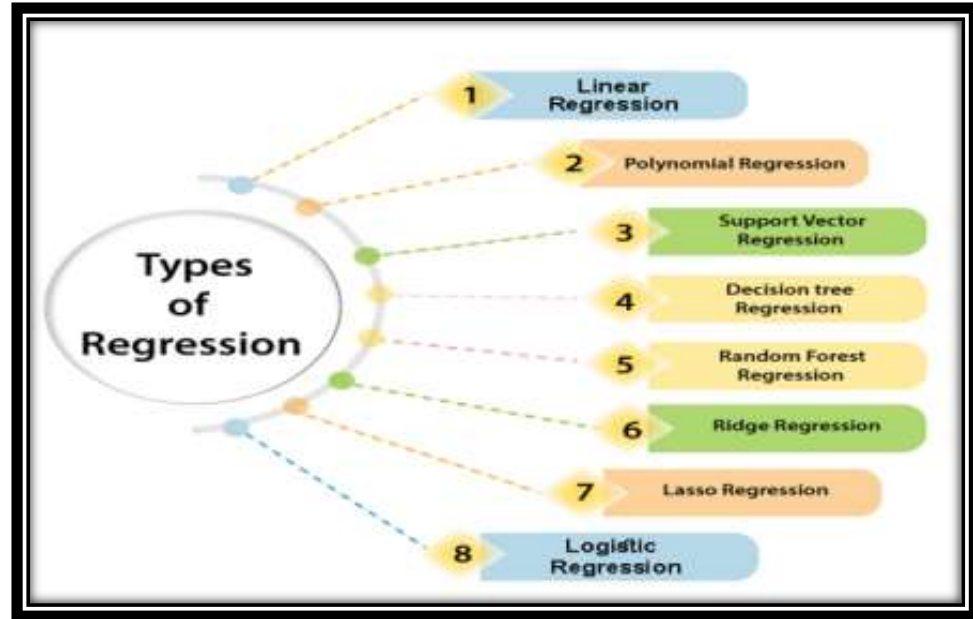
As mentioned above, Regression analysis helps in the prediction of a continuous variable. There are various scenarios in the real world where we need some future predictions such as weather condition, sales prediction, marketing trends, etc., for such case we need some technology which can make predictions more accurately. So for such case we need Regression analysis which is a statistical method and used in machine learning and data science. Below are some other reasons for using Regression analysis:

- o Regression estimates the relationship between the target and the independent variable.
- o It is used to find the trends in data.
- o It helps to predict real/continuous values.
- o By performing the regression, we can confidently determine the most important factor, the least important factor, and how each factor is affecting the other factors

### **Types of Regression**

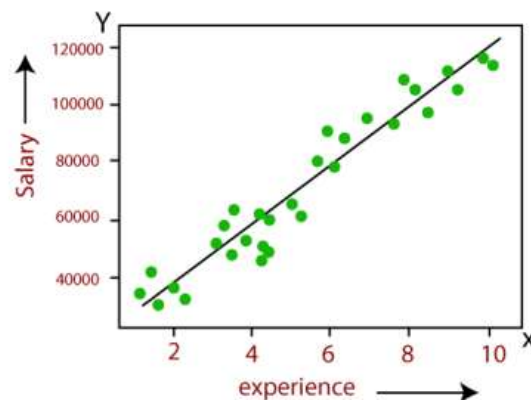
There are various types of regressions which are used in data science and machine learning. Each type has its own importance on different scenarios, but at the core, all the regression methods analyze the effect of the independent variable on dependent variables. Here we are discussing some important types of regression which are given below:

- o Linear Regression
- o Polynomial Regression
- o Decision Tree Regression
- o Ridge Regression
- o Logistic Regression
- o Support Vector Regression
- o Random Forest Regression
- o Lasso Regression



## Linear Regression

- o Linear regression is a statistical regression method which is used for predictive analysis.
  - o It is one of the very simple and easy algorithms which works on regression and shows the relationship between the continuous variables.
  - o It is used for solving the regression problem in machine learning.
  - o Linear regression shows the linear relationship between the independent variable (X-axis) and the dependent variable (Y-axis), hence called linear regression.<sup>34</sup>
  - o If there is only one input variable (x), then such linear regression is called simple linear regression. And if there is more than one input variable, then such linear regression is called multiple linear regression.
  - o The relationship between variables in the linear regression model can be explained using the below image.
- Here we are predicting the salary of an employee on the basis of the year of experience.





Below is the mathematical equation for Linear regression:

$$Y = aX + b$$

Here, Y = dependent variables (target variables),

X = Independent variables (predictor variables),

a and b are the linear coefficients

Some popular applications of linear regression are:

- o Analyzing trends and sales estimates
- o Salary forecasting
- o Real estate prediction
- o Arriving at ETAs in traffic.

Linear regression is a type of supervised machine learning algorithm that computes the linear relationship between the dependent variable and one or more independent features by fitting a linear equation to observed data.

Linear regression is a widely used statistical technique for modeling the relationship between a dependent variable and one or more independent variables. It assumes a linear relationship between the independent and dependent variables. This technique is extensively used in various fields such as economics, finance, social sciences, and machine learning.

Linear regression is a fundamental supervised machine learning algorithm that computes the linear relationship between a dependent variable (also known as the target variable) and one or more independent features (input variables). It achieves this by fitting a linear equation to observed data.

### **Why Linear Regression is Important?**

The interpretability of linear regression is a notable strength. The model's equation provides clear coefficients that elucidate the impact of each independent variable on the dependent variable, facilitating a deeper understanding of the underlying dynamics. Its simplicity is a virtue, as linear regression is transparent, easy to implement, and serves as a foundational concept for more complex algorithms.

Linear regression is not merely a predictive tool; it forms the basis for various advanced models. Techniques like regularization and support vector machines draw inspiration from linear regression, expanding its utility. Additionally, linear regression is a cornerstone in assumption testing, enabling researchers to validate key assumptions about the data.

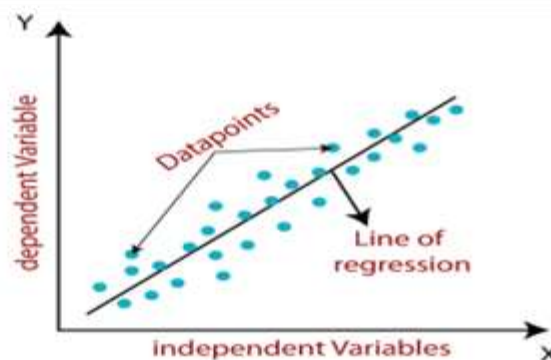


Linear regression has several strengths:

1. **Interpretability:** The model's equation provides clear coefficients that explain the impact of each independent variable on the dependent variable. This transparency facilitates a deeper understanding of the underlying dynamics.
2. **Simplicity:** Linear regression is easy to implement and serves as a foundational concept for more complex algorithms.
3. **Basis for Advanced Models:** Techniques like regularization and support vector machines draw inspiration from linear regression, expanding its utility.

Linear regression is one of the easiest and most popular Machine Learning algorithms. It is a statistical method that is used for predictive analysis. Linear regression makes predictions for continuous/real or numeric variables such as sales, salary, age, product price, etc. Linear regression algorithm shows a linear relationship between a dependent (y) and one or more independent (x) variables, hence called as linear regression. Since linear regression shows the linear relationship, which means it finds how the value of the dependent variable is changing according to the value of the independent variable. The linear regression model provides a sloped straight line representing the relationship between the variables.

Consider the below image



Mathematically, we can represent a linear regression as:

$$y = a_0 + a_1x + \epsilon$$

**Here,**

Y= Dependent Variable (Target Variable)

X= Independent Variable (predictor Variable)

$a_0$ = intercept of the line (Gives an additional degree of freedom)

$a_1$  = Linear regression coefficient (scale factor to each input value).

$\epsilon$  = random error

The values for x and y variables are training datasets for Linear Regression model representation:



### Types of Linear Regression

Linear regression can be further divided into two types of the algorithm:

o Simple Linear Regression:

If a single independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Simple Linear Regression.

o Multiple Linear regression:

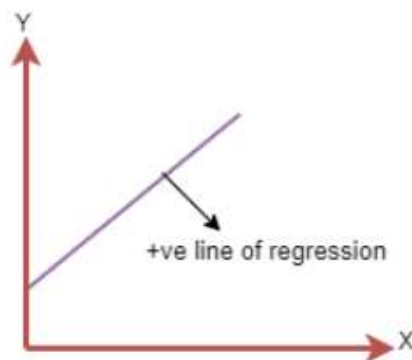
If more than one independent variable is used to predict the value of a numerical dependent variable, then such a Linear Regression algorithm is called Multiple Linear Regression.

Linear Regression Line:

A linear line showing the relationship between the dependent and independent variables is called a regression line. A regression line can show two types of relationship:

o Positive Linear Relationship:

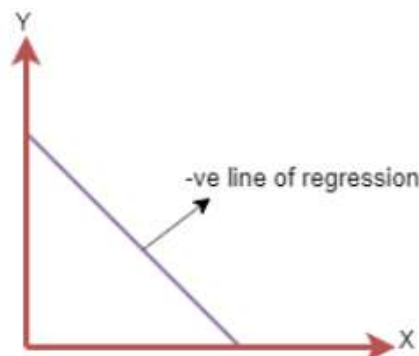
If the dependent variable increases on the Y-axis and independent variable increases on X-axis, then such a relationship is termed as a Positive linear relationship



The line equation will be:  $Y = a_0 + a_1X$

o **Negative Linear Relationship:**

If the dependent variable decreases on the Y-axis and independent variable increases on the X-axis, then such a relationship is called a negative linear relationship.



The line of equation will be:  $Y = -a_0 + a_1X$



### Finding the best fit line:

When working with linear regression, our main goal is to find the best fit line that means the error between predicted values and actual values should be minimized. The best fit line will have the least error.

The different values for weights or the coefficient of lines ( $a_0$ ,  $a_1$ ) gives a different line of regression, so we need to calculate the best values for  $a_0$  and  $a_1$  to find the best fit line, so to calculate this we use cost function.

Cost function

- o The different values for weights or coefficient of lines ( $a_0$ ,  $a_1$ ) gives the different line of regression, and the cost function is used to estimate the values of the coefficient for the best fit line.

- o Cost function optimizes the regression coefficients or weights. It measures how a linear regression model is performing.

- o We can use the cost function to find the accuracy of the mapping function, which maps the input variable to the output variable. This mapping function is also known as Hypothesis function.

### How Does Linear Regression Work?

1. **Hypothesis:** We assume that the relationship between the dependent variable and the independent features can be represented by a linear equation:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$$

- o Here, ( $y$ ) is the dependent variable, ( $x_1, x_2, \dots, x_n$ ) are the independent features, and ( $\beta_0, \beta_1, \dots, \beta_n$ ) are the coefficients.
2. **Cost Function:** We minimize the **cost function** (usually mean squared error) to find the best-fitting line.
  3. **Gradient Descent:** We use optimization techniques like gradient descent to adjust the coefficients iteratively until we converge to the optimal values.

### Assumptions of Linear Regression

1. **Linearity:** The relationship between the variables is linear.
2. **Independence:** The residuals (differences between predicted and actual values) are independent.
3. **Homoscedasticity:** The variance of residuals is constant across all levels of the independent variables.
4. **Normality:** The residuals follow a normal distribution.





## Evaluation Metrics

Common evaluation metrics for linear regression include:

- **Mean Absolute Error (MAE):** Average absolute differences between predicted and actual values.
- **Mean Squared Error (MSE):** Average squared differences.
- **Root Mean Squared Error (RMSE):** Square root of MSE.

## Simple Linear Regression in Machine Learning

Simple Linear Regression is a type of Regression algorithms that models the relationship between a dependent variable and a single independent variable. The relationship shown by a Simple Linear Regression model is linear or a sloped straight line, hence it is called Simple Linear Regression.

The key point in Simple Linear Regression is that the dependent variable must be a continuous/real value. However, the independent variable can be measured on continuous or categorical values.

Simple Linear regression algorithm has mainly two objectives:

- o Model the relationship between the two variables. Such as the relationship between Income and expenditure, experience and Salary, etc.
- o Forecasting new observations. Such as Weather forecasting according to temperature, Revenue of a company according to the investments in a year, etc.

Simple Linear Regression Model:

The Simple Linear Regression model can be represented using the below equation:

$$y = a_0 + a_1x + \epsilon$$

Where,

$a_0$  = It is the intercept of the Regression line (can be obtained putting  $x=0$ )

$a_1$  = It is the slope of the regression line, which tells whether the line is increasing or decreasing.

$\epsilon$  = The error term. (For a good model it will be negligible.)

## Multiple Linear Regressions

In the previous topic, we have learned about Simple Linear Regression, where a single Independent/Predictor(X) variable is used to model the response variable (Y). But there may be various cases in which the response variable is affected by more than one predictor variable; for such cases, the Multiple Linear Regression algorithm is used.



Moreover, Multiple Linear Regression is an extension of Simple Linear regression as it takes more than one predictor variable to predict the response variable.

We can define it as: “Multiple Linear Regression is one of the important regression algorithms which models the linear relationship between a single dependent continuous variable and more than one independent variable.”

Example:

Prediction of CO<sub>2</sub> emission based on engine size and number of cylinders in a car.

Some key points about MLR:

- o For MLR, the dependent or target variable(Y) must be the continuous/real, but the predictor or independent variable may be of continuous or categorical form.
- o Each feature variable must model the linear relationship with the dependent variable.
- o MLR tries to fit a regression line through a multidimensional space of data-points.

MLR equation: In Multiple Linear Regression, the target variable(Y) is a linear combination of multiple predictor variables  $x_1, x_2, x_3, \dots, x_n$ . Since it is an enhancement of Simple Linear Regression, so the same is applied for the multiple linear regression equation, the equation becomes:

$$Y = b_0 + b_1x_1 + b_2x_2 + b_3x_3 + \dots + b_nx_n \dots\dots\dots (a)$$

Where,

Y= Output/Response variable

$b_0, b_1, b_2, b_3, \dots, b_n$ ...= Coefficients of the model.

$x_1, x_2, x_3, x_4, \dots$ = Various Independent/feature variable

Assumptions for Multiple Linear Regression:

- o A linear relationship should exist between the Target and predictor variables.
- o The regression residuals must be normally distributed.
- o MLR assumes little or no multicollinearity (correlation between the independent variable) in data.

### Applications:

- **Predictive Modeling:** Predicting sales based on advertising expenditure, predicting house prices based on various features.
- **Forecasting:** Forecasting future demand for products or services.
- **Risk Assessment:** Assessing the risk factors associated with loan defaults, insurance claims, etc.
- **Causal Inference:** Determining the effect of one variable on another in controlled experiments or observational studies.



## Logistic Regression:

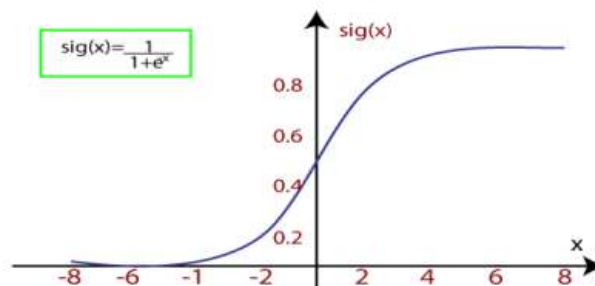
- Logistic regression is another supervised learning algorithm which is used to solve the classification problems. In classification problems, we have dependent variables in a binary or discrete format such as 0 or 1.
- Logistic regression algorithm works with the categorical variable such as 0 or 1, Yes or No, True or False, Spam or not spam, etc.
- It is a predictive analysis algorithm which works on the concept of probability.
- Logistic regression is a type of regression, but it is different from the linear regression algorithm in the term how they are used.
- Logistic regression uses sigmoid function or logistic function which is a complex cost function.

This sigmoid function is used to model the data in logistic regression. The function can be represented as:

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

- $f(x)$  = Output between the 0 and 1 value.
- $x$  = input to the function
- $e$  = base of natural logarithm.

When we provide the input values (data) to the function, it gives the S-curve as follows:



- It uses the concept of threshold levels, values above the threshold level are rounded up to 1, and values below the threshold level are rounded up to 0.



Logistic regression is a **supervised machine learning algorithm** primarily used for **binary classification tasks**. Its goal is to predict the **probability that an instance belongs to a given class** (e.g., spam or not spam, disease or no disease). Unlike linear regression, which predicts continuous values, logistic regression predicts the output of a **categorical dependent variable**. Here are some key points:

1. **Binary Classification:** Logistic regression is commonly used when the outcome is binary (e.g., yes/no, 0/1, true/false).
2. **Sigmoid Function (Logistic Function):** To map predicted values to probabilities, logistic regression uses the sigmoid function. This function ensures that the predicted values lie between 0 and 1. The sigmoid curve resembles an “S” shape.
3. **Threshold Value:** Logistic regression assigns instances to classes based on a threshold value (usually 0.5). If the predicted probability is greater than the threshold, it belongs to Class 1; otherwise, it belongs to Class 0.

#### Types of Logistic Regression

1. **Binomial Logistic Regression:**
  - Only two possible types of the dependent variable (e.g., pass/fail, spam/not spam).
  - Predicts probabilities for a binary outcome.
2. **Multinomial Logistic Regression:**
  - Three or more possible **unordered** types of the dependent variable (e.g., “cat,” “dog,” “sheep”).
  - Used for multi-class classification.
3. **Ordinal Logistic Regression:**
  - Three or more possible **ordered** types of dependent variables (e.g., “low,” “medium,” “high”).
  - Applicable when the outcome has an inherent order.

#### Assumptions:

**Linearity of Log Odds:** The relationship between the independent variables and the logit of the outcome is linear.

**Independence of Errors:** The observations are independent of each other.

**Absence of Multicollinearity:** The independent variables are not highly correlated with each other.

**Large Sample Size:** The sample size should be large enough to ensure stable estimates.

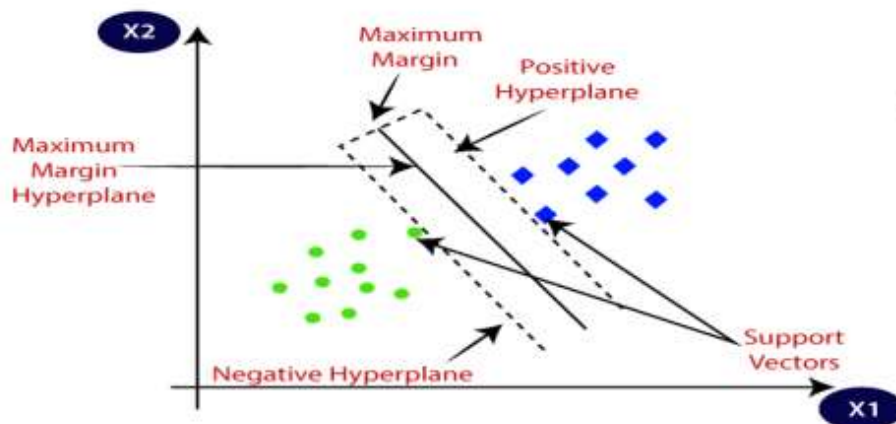


### **Applications:**

1. **Medical Research:** Predicting the likelihood of a disease based on patient characteristics.
2. **Marketing:** Predicting customer churn, likelihood of purchase, etc.
3. **Finance:** Predicting credit risk, likelihood of default, etc.
4. **Social Sciences:** Predicting voter preferences, survey responses, etc.

## Support Vector Machines

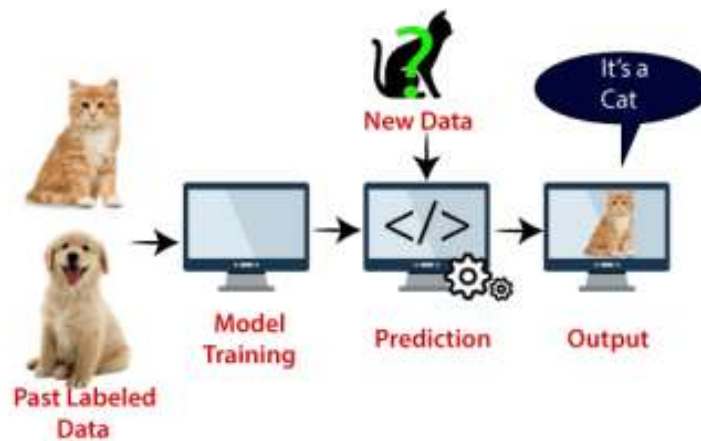
Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning. The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane. SVM chooses the extreme points/vectors that help in creating the hyperplane. These extreme cases are called as support vectors, and hence algorithm is termed as Support Vector Machine. Consider the below diagram in which there are two different categories that are classified using a decision boundary or hyperplane:



Example: SVM can be understood with the example that we have used in the KNN classifier. Suppose we see a strange cat that also has some features of dogs, so if we want a model that can accurately identify whether it is a cat or dog, so such a model can be created by using the SVM algorithm. We will first train our model with lots of images of cats and dogs so that it can learn about different features of cats and dogs, and then we test it with this strange creature. So as support vector creates a decision boundary between these two data (cat and dog) and choose extreme cases (support vectors), it will see the



extreme case of cat and dog. On the basis of the support vectors, it will classify it as a cat.  
Consider the below diagram:



**Support Vector Machine (SVM)** is a powerful **supervised machine learning algorithm** that can be used for both **classification** and **regression** tasks. The primary objective of SVM is to find the optimal hyperplane that best separates the data points belonging to different classes. This hyperplane is chosen in such a way that the margin, i.e., the distance between the hyperplane and the nearest data points (called support vectors), is maximized

Here are the key points about SVM:

1. **Objective:**

- SVM aims to find an **optimal hyperplane** in an N-dimensional space that can **separate data points** belonging to different classes.
- The hyperplane's goal is to maximize the **margin** between the closest points of different classes.

2. **Linear Separation:**

- When we have **two input features** (e.g.,  $x_1$  and  $x_2$ ), the hyperplane is just a **line**.
- For **three input features**, the hyperplane becomes a **2-D plane**.
- Beyond three features, it becomes challenging to visualize.

3. **Choosing the Best Hyperplane:**

- The best hyperplane is the one that represents the **largest separation** or margin between the two classes.
- SVM selects the hyperplane whose distance from the nearest data point on each side is **maximized**.
- If such a hyperplane exists, it's called the **maximum-margin hyperplane** (or hard margin).

4. **Handling Outliers:**



- SVM is **robust to outliers**.
- Even if there's an outlier, SVM finds the best hyperplane that maximizes the margin.

5. **Nonlinear Separation:**

- SVM can handle **nonlinear relationships** by using **kernel functions**.
- These functions map the data into a higher-dimensional space where linear separation becomes possible.

**SVM can be of two types:**

○ **Linear SVM:** Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then such data is termed as linearly separable data, and classifier is used called as Linear SVM classifier.

○ **Non-linear SVM:** Non-Linear SVM is used for non-linearly separated data, which means if a dataset cannot be classified by using a straight line, then such data is termed as non-linear data and classifier used is called as Non-linear SVM classifier.

**Key Concepts:**

1. **Hyperplane:**

- In SVM, a hyperplane is a decision boundary that separates the data into different classes.
- For binary classification, the hyperplane is a line. In higher dimensions, it becomes a plane (3D) or a hyperplane (more than 3 dimensions).

2. **Support Vectors:**

- Support vectors are the data points closest to the hyperplane and influence the position and orientation of the hyperplane.
- They are critical for defining the margin and determining the decision boundary.

3. **Margin:**

- The margin is the distance between the hyperplane and the nearest data points (support vectors).
- SVM aims to maximize the margin, as it helps improve the generalization ability of the model and reduces the risk of overfitting.

4. **Kernel Trick:**

- SVM can efficiently perform classification in nonlinearly separable datasets by using the kernel trick.
- Kernels transform the input data into higher-dimensional feature spaces, where it's easier to find a hyperplane that separates the classes.
- Popular kernels include linear, polynomial, radial basis function (RBF), and sigmoid kernels.

Advantages of SVM:

1. **Effective in High-Dimensional Spaces:** SVM performs well even when the number of features is much greater than the number of samples.



2. **Versatile:** SVM can be applied to various types of data, including text, images, and numerical data.
3. **Memory Efficient:** SVM uses a subset of training points (support vectors) in the decision function, making it memory efficient.
4. **Robustness:** SVM is less affected by noisy data due to its focus on maximizing the margin.

### Applications:

- SVMs are used for various tasks, including:
  - **Text classification**
  - **Image classification**
  - **Spam detection**
  - **Handwriting identification**
  - **Gene expression analysis**
  - **Face detection**
  - **Anomaly detection**

### Hyperplane and Support Vectors in the SVM algorithm:

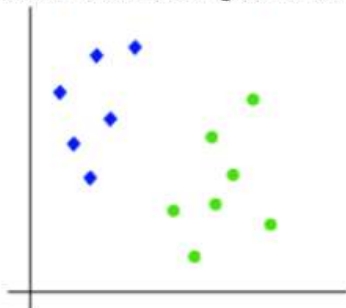
**Hyperplane:** There can be multiple lines/decision boundaries to segregate the classes in  $n$  dimensional space, but we need to find out the best decision boundary that helps to classify the data points. This best boundary is known as the hyperplane of SVM. The dimensions of the hyperplane depend on the features present in the dataset, which means if there are 2 features (as shown in image), then hyperplane will be a straight line. And if there are 3 features, then hyperplane will be a 2-dimension plane. We always create a hyperplane that has a maximum margin, which means the maximum distance between the data points.

**Support Vectors:** The data points or vectors that are the closest to the hyperplane and which affect the position of the hyperplane are termed as Support Vector. Since these vectors support the hyperplane, hence called a Support vector. How does SVM works?

### Linear SVM:

The working of the SVM algorithm can be understood by using an example. Suppose we have a dataset that has two tags (green and blue), and the dataset has two features  $x_1$  and  $x_2$ . We want a classifier that can classify the pair( $x_1$ ,  $x_2$ ) of coordinates in either green or blue. Consider the below image:

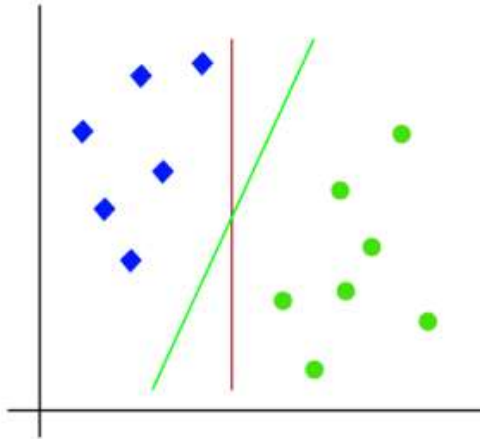
Figure 1: A 2D scatter plot showing data points in either green or blue. Consider the







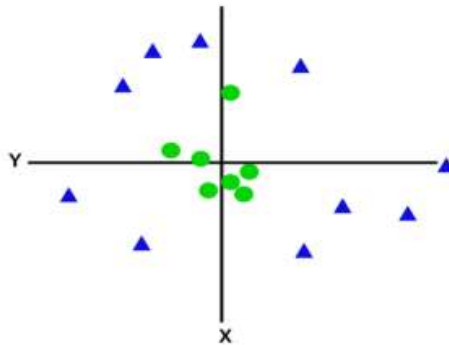
So as it is 2-d space so by just using a straight line, we can easily separate these two classes. But there can be multiple lines that can separate these classes. Consider the below image



Hence, the SVM algorithm helps to find the best line or decision boundary; this best boundary or region is called as a **hyperplane**. SVM algorithm finds the closest point of the lines from both the classes. These points are called support vectors. The distance between the vectors and the hyperplane is called as **margin**. And the goal of SVM is to maximize this margin. The **hyperplane** with maximum margin is called the **optimal hyperplane**.

### Non-Linear SVM:

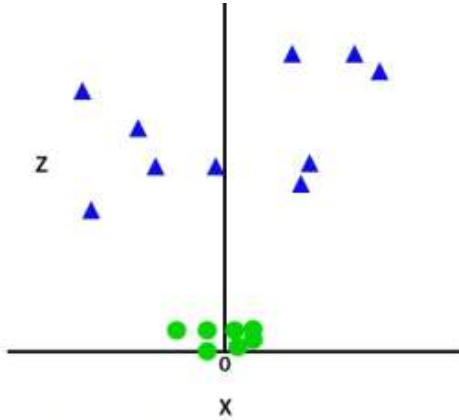
If data is linearly arranged, then we can separate it by using a straight line, but for non-linear data, we cannot draw a single straight line. Consider the below image:



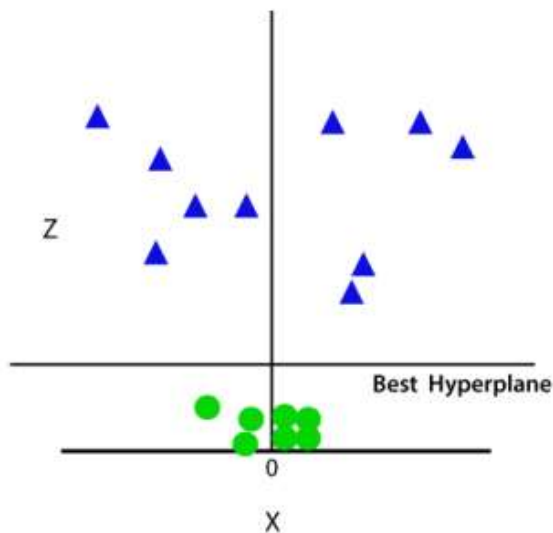
So to separate these data points, we need to add one more dimension. For linear data, we have used two dimensions x and y, so for non-linear data, we will add a third dimension z. It can be calculated as:

$$z = x^2 + y^2$$

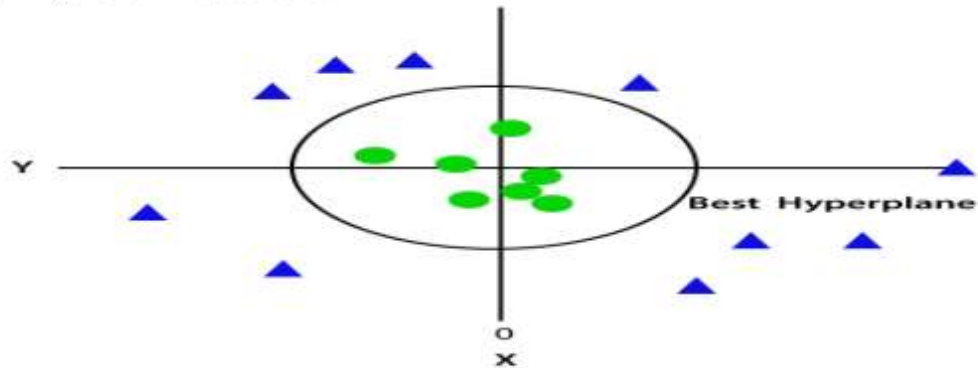
By adding the third dimension, the sample space will become as below image:



So now, SVM will divide the datasets into classes in the following way. Consider the below image:



Since we are in 3-d Space, hence it is looking like a plane parallel to the x-axis. If we convert it in 2d space with  $z=1$ , then it will become as:



Hence we get a circumference of radius 1 in case of non-linear data.



## SVM Kernels

In practice, SVM algorithm is implemented with kernel that transforms an input data space into the required form. SVM uses a technique called the kernel trick in which kernel takes a low dimensional input space and transforms it into a higher dimensional space. In simple words, kernel converts non separable problems into separable problems by adding more dimensions to it. It makes SVM more powerful, flexible and accurate.

## Types of Support Vector Kernel

Support Vector Machines (SVMs) use various kernel functions to transform input data into higher-dimensional spaces where it's easier to find a hyperplane that separates classes. The three main types of kernel functions used in SVM are:

### 1. Linear Kernel:

- The linear kernel is the simplest form of kernel and is represented by the dot product of the feature vectors.
- It works well when the data is linearly separable or when the number of features is very large compared to the number of samples.
- The decision boundary is a straight line in two dimensions and a hyperplane in higher dimensions.

### 2. Polynomial Kernel:

- The polynomial kernel calculates the dot product of two vectors raised to a power, which introduces nonlinearity into the decision boundary.
- It's useful when the data is not linearly separable in the original feature space but can be separated with a polynomial boundary in a higher-dimensional space.
- The degree of the polynomial determines the complexity of the decision boundary.

### 3. Gaussian Kernel (RBF - Radial Basis Function):

- The Gaussian kernel transforms data into an infinite-dimensional space using a Gaussian (radial basis) function.
- It's effective for capturing complex nonlinear relationships between features and can handle data that is not linearly separable in the original space.
- The kernel width parameter (sigma) controls the smoothness of the decision boundary, with smaller values leading to more complex boundaries.

## Hyperplane (Decision Surface)

In SVM, the hyperplane is the decision boundary that separates data points belonging to different classes. Key points regarding hyperplanes include:



- For binary classification, the hyperplane is a line (2D), a plane (3D), or a hyperplane (more than 3 dimensions).
- The goal of SVM is to find the hyperplane with the maximum margin, which is the distance between the hyperplane and the nearest data points (support vectors).
- Support vectors are the data points closest to the hyperplane and determine its position and orientation.

### Properties of SVM

Support Vector Machines possess several key properties that contribute to their effectiveness in classification tasks:

1. **Effective in High-Dimensional Spaces:** SVM performs well even when the number of features is much greater than the number of samples, making it suitable for high-dimensional data.
2. **Memory Efficient:** SVM uses a subset of training points (support vectors) in the decision function, making it memory efficient, especially for large datasets.
3. **Robustness to Noise:** SVM is less sensitive to noisy data because it focuses on maximizing the margin, which reduces the risk of overfitting.
4. **Versatility:** SVM can be applied to various types of data, including text, images, and numerical data, making it a versatile algorithm in machine learning.

### Issues in SVM

While SVM is a powerful algorithm, it has some limitations and potential issues:

1. **Choice of Kernel:** The performance of SVM can be sensitive to the choice of kernel function and its parameters. Selecting the appropriate kernel and tuning hyperparameters can be challenging.
2. **Computational Complexity:** Training an SVM model can be computationally expensive, especially for large datasets, as it involves solving a quadratic optimization problem.
3. **Scalability:** SVM may not scale well to very large datasets due to its computational complexity and memory requirements.
4. **Interpretability:** SVM models are often considered less interpretable compared to some other machine learning algorithms like decision trees or linear regression.
5. **Imbalanced Data:** SVM may not perform well with highly imbalanced datasets, where one class has significantly fewer samples than the other.



Addressing these issues often involves careful selection of hyperparameters, preprocessing of data, and consideration of alternative algorithms based on the specific characteristics of the dataset and the problem at hand.

## BAYESIAN LEARNING

The third family of machine learning algorithms is the probabilistic models. We have seen before that the k-nearest neighbour algorithm uses the idea of distance (e.g., Euclidian distance) to classify entities, and logical models use a logical expression to partition the instance space. In this section, we see how the probabilistic models use the idea of probability to classify new entities. Probabilistic models see features and target variables as random variables. The process of modelling represents and manipulates the level of uncertainty with respect to these variables. There are two types of probabilistic models: Predictive and Generative. Predictive probability models use the idea of a conditional probability distribution  $P(Y|X)$  from which  $Y$  can be predicted from  $X$ . Generative models estimate the joint distribution  $P(Y, X)$ . Once we know the joint distribution for the generative models, we can derive any conditional or marginal distribution involving the same variables. Thus, the generative model is capable of creating new data points and their labels, knowing the joint probability distribution. The joint distribution looks for a relationship between two variables. Once this relationship is inferred, it is possible to infer new data points. **Naïve Bayes** is an example of a probabilistic classifier.

**Bayes' theorem** is a fundamental concept in probability theory that plays a crucial role in various machine learning algorithms, especially in the fields of Bayesian statistics and probabilistic modeling. It provides a way to update probabilities based on new evidence or information.

Mathematically, Bayes' theorem can be expressed as:

$$P(A|B)=P(B)P(B|A) \cdot P(A)$$

Here are the key terms related to Bayes' theorem:

### 1. Likelihood ( $P(B|A)$ ):

- Represents the probability of observing the given evidence (features) given that the class is true.
- In the **Naive Bayes algorithm**, a key assumption is that features are conditionally independent given the class label. This assumption works well with discrete features.



- For example, in text classification, the likelihood represents how likely certain words appear in a specific class (e.g., spam or not spam).
- 2. **Prior Probability (P(A)):**
  - Represents the probability of a particular class before considering any features.
  - It is estimated from the training data.
- 3. **Evidence Probability (P(B)):**
  - The probability of observing the given evidence (features).
  - Serves as a normalization factor and is often calculated as the sum of the joint probabilities over all possible classes.
- 4. **Posterior Probability (P(A|B)):**
  - The updated probability of the class given the observed features.
  - It is what we are trying to predict or infer in a classification task.

### **Concept Learning:**

Concept learning is a process in machine learning where an algorithm learns to categorize objects or instances into different classes or concepts based on their attributes or features. Bayesian learning plays a significant role in concept learning by updating the probabilities of different concepts given observed data.

### **Bayesian Learning Process:**

**Prior Knowledge:** Start with prior beliefs or assumptions about the model parameters.

**Observation:** Collect data or evidence from the environment.

**Likelihood:** Calculate the likelihood of observing the data given the model parameters.

**Posterior Update:** Use Bayes' theorem to update the prior beliefs based on the observed data, resulting in the posterior distribution of the model parameters.

**Decision Making:** Make decisions or predictions based on the updated posterior distribution.

### ***Bayesian Learning in Machine Learning:***

1. **Bayesian Regression:** In Bayesian regression, we model the posterior distribution of regression coefficients given the data, allowing us to quantify uncertainty in predictions.
2. **Bayesian Classification:** In Bayesian classification, we model the posterior probability of class labels given the input features, enabling us to make probabilistic predictions.
3. **Bayesian Optimization:** Bayesian optimization is a technique for global optimization of black-box functions that uses Bayesian inference to model the objective function and guide the search process.



4. **Bayesian Networks:** Bayesian networks are probabilistic graphical models that represent dependencies between random variables using a directed acyclic graph, allowing for efficient inference and reasoning under uncertainty.

### Advantages of Bayesian Learning:

**Incorporation of Prior Knowledge:** Bayesian learning provides a principled framework for incorporating prior knowledge or beliefs into the learning process.

**Uncertainty Estimation:** Bayesian methods enable the quantification of uncertainty in predictions, which is crucial for decision-making in uncertain environments.

**Flexibility:** Bayesian learning can accommodate various types of models, including regression, classification, and unsupervised learning, making it a versatile approach in machine learning.

**Robustness to Overfitting:** By incorporating prior beliefs and updating them based on observed data, Bayesian learning tends to be more robust to overfitting, especially in cases of limited data.

### Challenges in Bayesian Learning:

**Computational Complexity:** Bayesian inference can be computationally expensive, especially for complex models or large datasets, requiring approximation techniques like Markov chain Monte Carlo (MCMC) or variational inference.

**Model Specification:** Choosing appropriate prior distributions and model structures can be challenging and may require domain expertise.

**Interpretability:** Bayesian models can be less interpretable than simpler models like linear regression or decision trees, particularly when using complex probabilistic models.

### 1. Bayes Optimal Classifier:

The Bayes Optimal Classifier is a theoretical framework for classification tasks that aims to minimize the probability of misclassification. It classifies instances based on the class that has the highest posterior probability given the input features. Key points include:

- It assumes knowledge of the true underlying probability distributions of classes and features.
- It assigns each instance to the class with the highest posterior probability, computed using Bayes' theorem.
- In practice, the Bayes Optimal Classifier serves as a benchmark for evaluating the performance of other classifiers.
- The **Bayes Optimal Classifier** is a theoretical concept in machine learning that represents the best possible classifier for a given problem.



- It is based on **Bayes' theorem**, which describes how to update probabilities based on new evidence.
- The Bayes Optimal Classifier makes the most probable prediction for a new example, given the training dataset.
- However, in practice, it is often computationally expensive or even intractable to calculate directly.
- Simplifications such as the **Gibbs algorithm** or **Naive Bayes** can be used to approximate the outcome.
- **This model is also referred to as the Bayes optimal learner, Bayes classifier, Bayes optimal decision boundary, or Bayes optimal discriminant function**

## 2. Naïve Bayes Classifier:

The Naïve Bayes Classifier is a simple probabilistic classifier based on Bayes' theorem with an assumption of independence between features. Key points include:

- **It assumes that all features are conditionally independent given the class label.**
- **Despite its simplicity, Naïve Bayes often performs surprisingly well in practice, especially for text classification tasks.**
- **It's computationally efficient and requires a small amount of training data.**
- **Popular variants include Gaussian Naïve Bayes (for continuous features), Multinomial Naïve Bayes (for discrete features), and Bernoulli Naïve Bayes (for binary features).**
- The **Naïve Bayes classifier** is a simple yet effective classification algorithm based on **Bayes' Theorem**.
- It assumes that all features used to describe an observation are **conditionally independent**, given the class label.
- Despite this “naive” assumption, Naïve Bayes classifiers are widely used for their simplicity and efficiency.
- Applications include **text classification, spam filtering, sentiment detection**, and more.
- It works well with high-dimensional data, such as text data where each word represents a feature

## 3. Bayesian Belief Networks (BBNs):

Bayesian Belief Networks, also known as Bayesian Networks, are probabilistic graphical models that represent dependencies between random variables using a directed acyclic graph (DAG). Key points include:

- Nodes in the graph represent random variables, and edges represent probabilistic dependencies between them.





- Each node is associated with a conditional probability distribution that quantifies the probability of the node given its parents in the graph.
- Bayesian Networks provide a compact and interpretable representation of complex probability distributions.
- They are used for probabilistic inference, prediction, and decision-making under uncertainty.
- Learning the structure and parameters of Bayesian Networks from data is a challenging task, often requiring efficient algorithms for structure learning and parameter estimation.
- A **Bayesian Belief Network** (BBN) is a graphical representation of probabilistic relationships among random variables.
- BBNs are used to represent **uncertain knowledge** and make decisions based on that knowledge.
- They are a type of **Bayesian network**, which models probabilistic relationships between variables.
- In a BBN, nodes represent variables, and directed edges indicate their dependencies.
- BBNs are useful for reasoning under uncertainty and handling complex dependencies

### Advantages and Limitations:

- **Advantages:** Bayesian classifiers are robust, computationally efficient, and can handle high-dimensional data. They provide probabilistic predictions and can quantify uncertainty.
- **Limitations:** Naïve Bayes assumes independence between features, which may not hold in practice. Bayesian Belief Networks can become computationally expensive for large and complex networks, and learning their structure from data may be challenging.

### Applications:

- **Bayes Optimal Classifier:** Used as a theoretical benchmark for evaluating the performance of other classifiers.
- **Naïve Bayes Classifier:** Widely used in text classification, spam filtering, and other tasks where the independence assumption holds reasonably well.
- **Bayesian Belief Networks:** Applied in various domains such as healthcare, finance, and risk assessment for modeling complex probabilistic relationships and making informed decisions.