# SNS COLLEGE OF TECHNOLOGY
### (AN AUTONOMOUS INSTITUTION)
COIMBATORE – 35
**DEPARTMENT OF COMPUTER SIENCE AND ENGINEERING**

## UNIT II      SUPERVISED LEARNING

Introduction - Linear Models for Regression – Linear Regression Models and Least Squares
Subset Selection
Shrinkage Methods – Derived Input Directions
Linear Models for Classification- Discriminant Analysis
Logistic Regression
Separating Hyperplanes

## Regression Analysis in Machine learning

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.

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:

- **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**.
- **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**.
- **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.
- **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.
- **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:
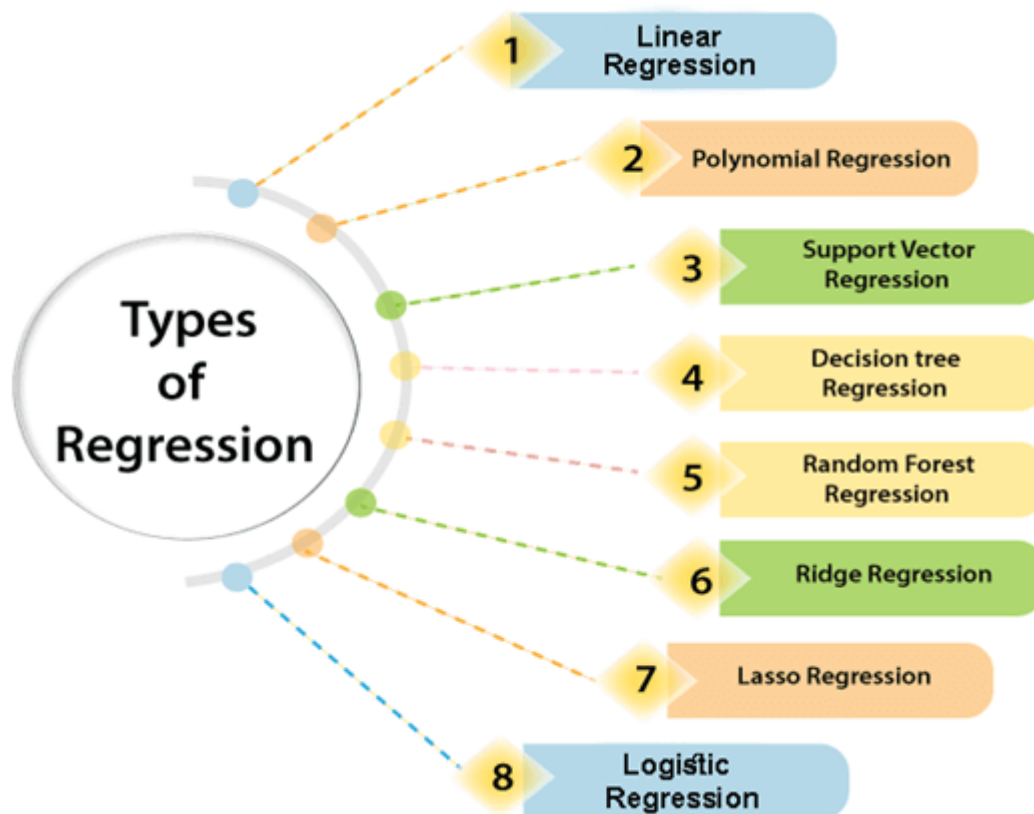
- Regression estimates the relationship between the target and the independent variable.
- It is used to find the trends in data.

- It helps to predict real/continuous values.
- 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:
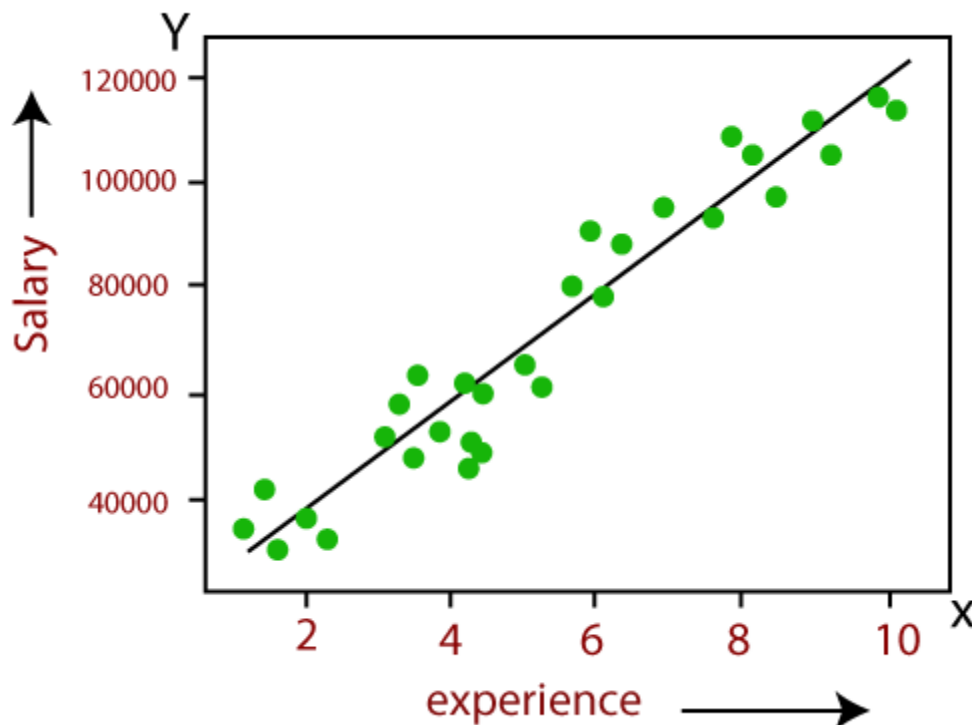
- **Linear Regression**
- **Logistic Regression**
- **Polynomial Regression**
- **Support Vector Regression**
- **Decision Tree Regression**
- **Random Forest Regression**
- **Ridge Regression**
- **Lasso Regression:**



### Linear Regression:

- Linear regression is a statistical regression method which is used for predictive analysis.
- It is one of the very simple and easy algorithms which works on regression and shows the relationship between the continuous variables.

- It is used for solving the regression problem in machine learning.
- Linear regression shows the linear relationship between the independent variable (X-axis) and the dependent variable (Y-axis), hence called linear regression.
- 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**.
- 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:

1. 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:

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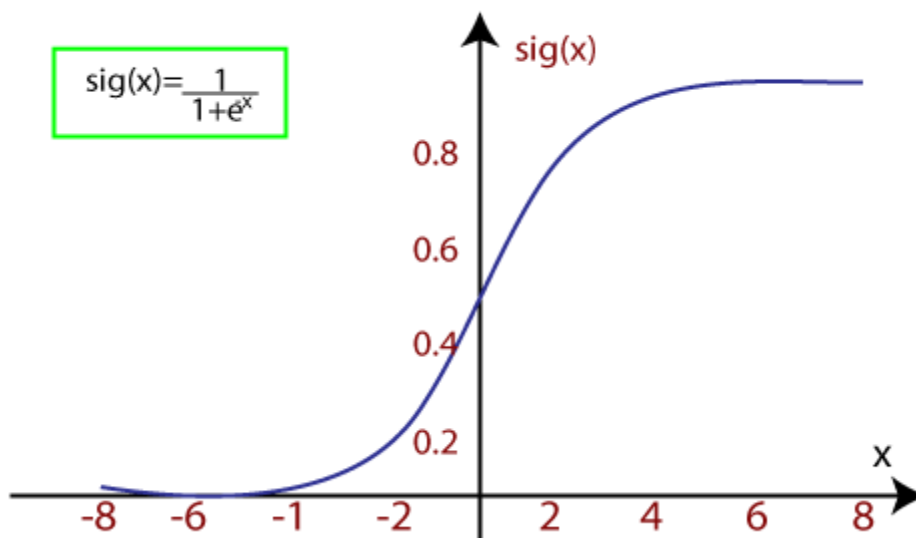- **Arriving at ETAs in traffic.**

### 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.

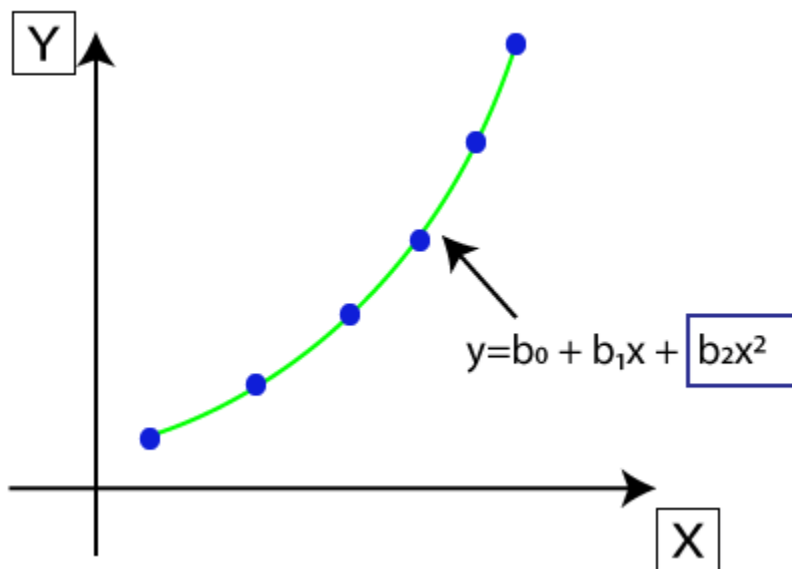When we provide the input values (data) to the function, it gives the S-curve as follows:



- 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.

There are three types of logistic regression:

- **Binary(0/1, pass/fail)**
- **Multi(cats, dogs, lions)**
- **Ordinal(low, medium, high)**

### Polynomial Regression:

- Polynomial Regression is a type of regression which models the **non-linear dataset** using a linear model.
- It is similar to multiple linear regression, but it fits a non-linear curve between the value of x and corresponding conditional values of y.
- Suppose there is a dataset which consists of datapoints which are present in a non-linear fashion, so for such case, linear regression will not best fit to those datapoints. To cover such datapoints, we need Polynomial regression.
- I**n Polynomial regression, the original features are transformed into polynomial features of given degree and then modeled using a linear model.** Which means the datapoints are best fitted using a polynomial line.



$$y = b_0 + b_1 x + \boxed{b_2 x^2}$$

- The equation for polynomial regression also derived from linear regression equation that means Linear regression equation $Y = b_0 + b_1 x$, is transformed into Polynomial regression equation $Y = b_0 + b_1 x + b_2 x^2 + b_3 x^3 + ..... + b_n x^n$.
- Here Y is the **predicted/target output, $b_0$, $b_1$,... $b_n$ are the regression coefficients**. x is our **independent/input variable**.
- The model is still linear as the coefficients are still linear with quadratic

### Note: This is different from Multiple Linear regression in such a way that in Polynomial regression, a single element has different degrees instead of multiple variables with the same degree.
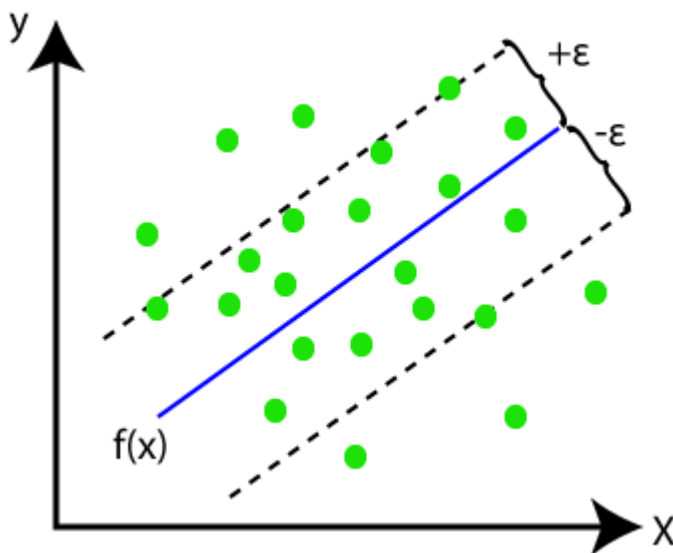
### Support Vector Regression:

Support Vector Machine is a supervised learning algorithm which can be used for regression as well as classification problems. So if we use it for regression problems, then it is termed as Support Vector Regression.

Support Vector Regression is a regression algorithm which works for continuous variables. Below are some keywords which are used in **Support Vector Regression**:

- **Kernel:** It is a function used to map a lower-dimensional data into higher dimensional data.
- **Hyperplane:** In general SVM, it is a separation line between two classes, but in SVR, it is a line which helps to predict the continuous variables and cover most of the datapoints.
- **Boundary line:** Boundary lines are the two lines apart from hyperplane, which creates a margin for datapoints.
- **Support vectors:** Support vectors are the datapoints which are nearest to the hyperplane and opposite class.
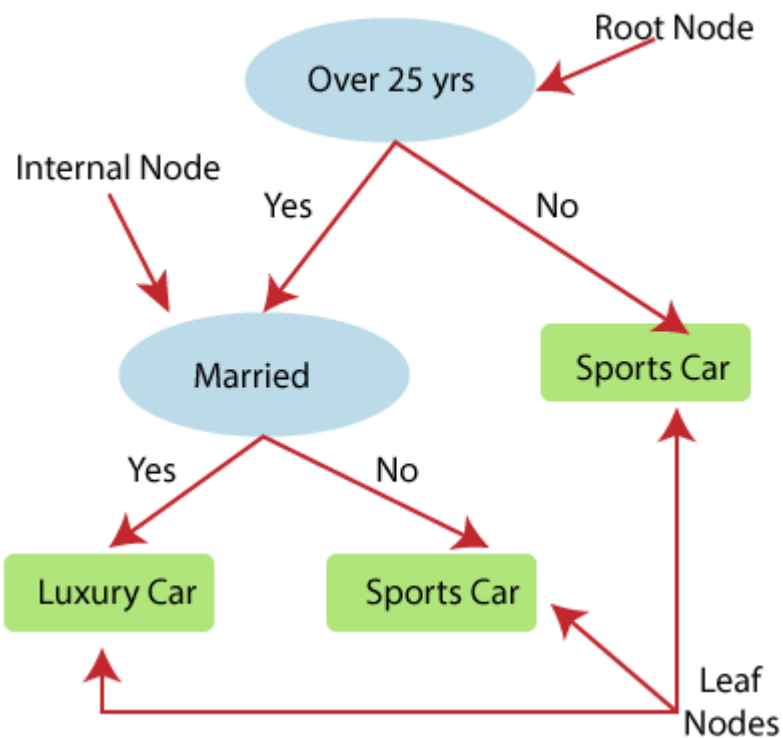
In SVR, we always try to determine a hyperplane with a maximum margin, so that maximum number of datapoints are covered in that margin. *The main goal of SVR is to consider the maximum datapoints within the boundary lines and the hyperplane (best-fit line) must contain a maximum number of datapoints*. Consider the below image:



Here, the blue line is called hyperplane, and the other two lines are known as boundary lines.

### Decision Tree Regression:

- Decision Tree is a supervised learning algorithm which can be used for solving both classification and regression problems.
- It can solve problems for both categorical and numerical data
- Decision Tree regression builds a tree-like structure in which each internal node represents the "test" for an attribute, each branch represent the result of the test, and each leaf node represents the final decision or result.
- A decision tree is constructed starting from the root node/parent node (dataset), which splits into left and right child nodes (subsets of dataset). These child nodes are further divided into their children node, and themselves become the parent node of those nodes. Consider the below image:
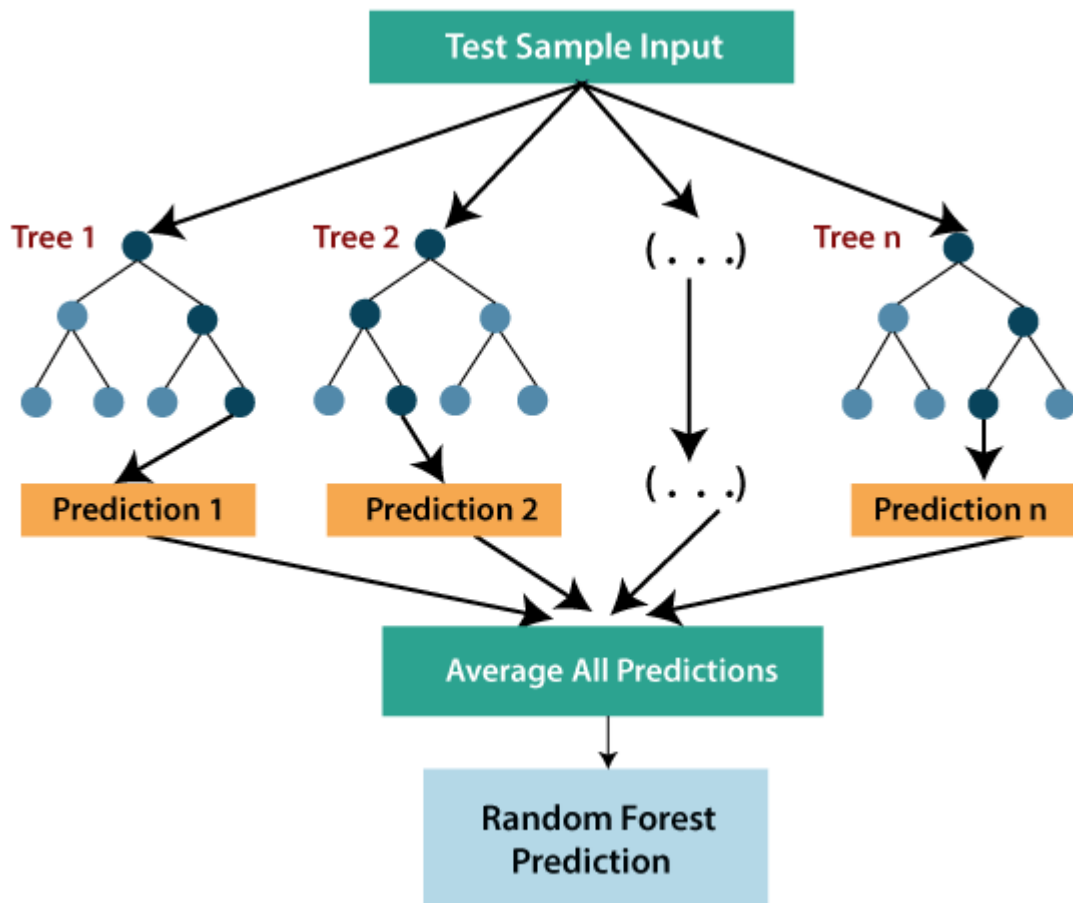
Above image showing the example of Decision Tee regression, here, the model is trying to predict the choice of a person between Sports cars or Luxury car.

- Random forest is one of the most powerful supervised learning algorithms which is capable of performing regression as well as classification tasks.
- The Random Forest regression is an ensemble learning method which combines multiple decision trees and predicts the final output based on the average of each tree output. The combined decision trees are called as base models, and it can be represented more formally as:

```
g(x) = f₀(x) + f₁(x) + f₂(x) + ....
```

$$g(x) = f_0(x) + f_1(x) + f_2(x) + ....$$

- Random forest uses **Bagging or Bootstrap Aggregation** technique of ensemble learning in which aggregated decision tree runs in parallel and do not interact with each other.
- With the help of Random Forest regression, we can prevent Overfitting in the model by creating random subsets of the dataset.

### Ridge Regression:

- Ridge regression is one of 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.
- The equation for ridge regression will be:

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

- A general linear or polynomial regression will fail if there is high collinearity between the independent variables, so to solve such problems, Ridge regression can be used.
- 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.

### Lasso Regression:

- Lasso regression is another regularization technique to reduce the complexity of the model.
- It is similar to the Ridge Regression except that penalty term contains only the absolute weights instead of a square of weights.

- Since it takes absolute values, hence, it can shrink the slope to 0, whereas Ridge Regression can only shrink it near to 0.
- It is also called as **L1 regularization**. The equation for Lasso regression will be:
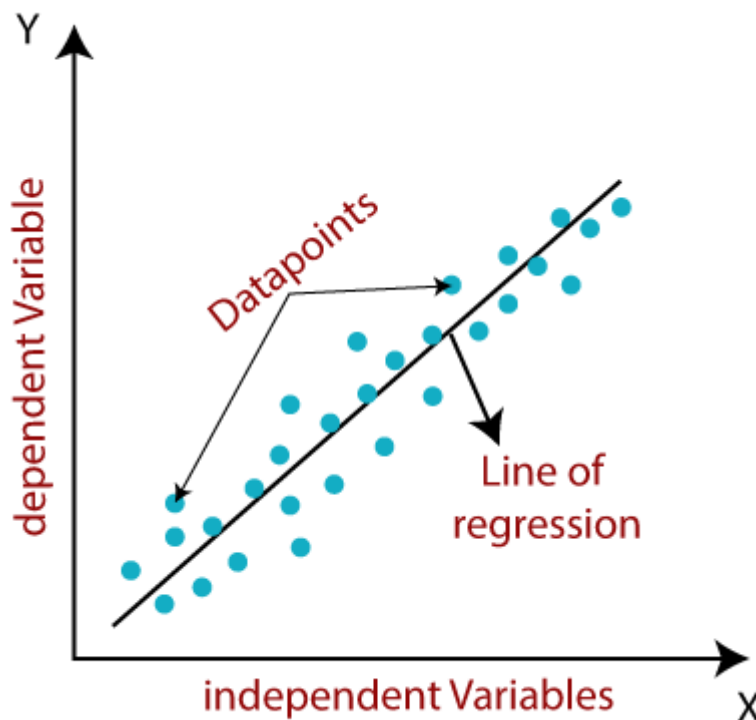
$$L(x, y) = Min\left( \sum_{i=1}^{n} (y_i - w_i x_i)^2 + \lambda \sum_{i=1}^{n} |w_i| \right)$$

# Linear Regression in Machine Learning

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 (y) 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:

$y = a_0 + a_1 x + \varepsilon$

**Here,**

Y= Dependent Variable (Target Variable)
X= Independent Variable (predictor Variable)
a0= intercept of the line (Gives an additional degree of freedom)
a1 = Linear regression coefficient (scale factor to each input value).
$\varepsilon$ = 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:

- **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.
- **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.

## Types of Linear Regression:

1. Simple Linear Regression:

   - Used when there is only one independent variable.
   - The model can be written as:

$$y = \beta_0 + \beta_1 x + \epsilon$$

   where:

   - $y$ is the dependent variable (output),
   - $x$ is the independent variable (input),
   - $\beta_0$ is the y-intercept (constant term),
   - $\beta_1$ is the coefficient (slope),
   - $\epsilon$ is the error term (residuals).

2. Multiple Linear Regression:

   - Used when there are multiple independent variables.
   - The model can be written as:

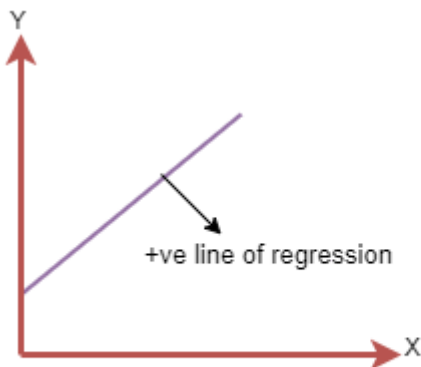$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_n x_n + \epsilon$$

   where:

   - $y$ is the dependent variable,
   - $x_1, x_2, \ldots, x_n$ are the independent variables,
   - $\beta_0$ is the intercept,
   - $\beta_1, \beta_2, \ldots, \beta_n$ are the coefficients,
   - $\epsilon$ is the error term.

# 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:

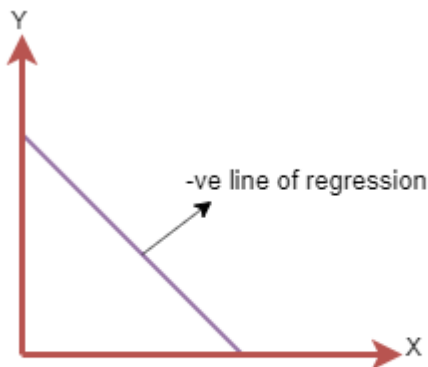- **Positive Linear Relationship:**
  If the dependent variable increases on the Y-axis and independent variable increases on X-axis, then such a relationship is termed as a Positive linear relationship.



The line equation will be: $Y = a_0 + a_1 x$

- **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: $Y = -a_0 + a_1 x$

# Finding the best fit line: (Cost function, Residuals, Gradient Descent)

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-

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

For Linear Regression, we use the **Mean Squared Error (MSE)** cost function, which is the average of squared error occurred between the predicted values and actual values. It can be written as:

For the above linear equation, MSE can be calculated as:

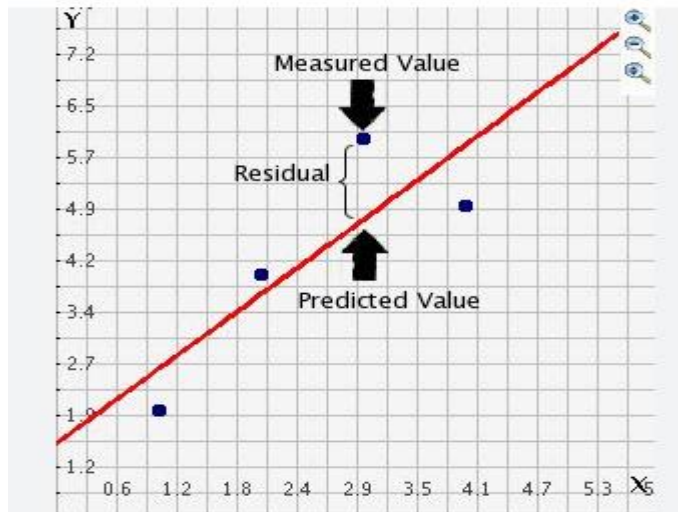$$MSE = 1\frac{1}{N}\sum_{i=1}^{n}(y_i - (a_1x_i + a_0))^2$$

**Where,**

N=Total number of observation
Yi = Actual value
($a1x_i+a_0$)= Predicted value.

**Residuals:** The distance between the actual value and predicted values is called residual.

 If the observed points are far from the regression line, then the residual will be high, and so cost function will high. If the scatter points are close to the regression line, then the residual will be small and hence the cost function.



## Residual Calculation:

For a data point $i$, the residual $e_i$ is calculated as:

$$e_i = y_i - \hat{y}_i$$

where:

- $y_i$ is the actual value of the dependent variable for the $i$-th observation,

- $\hat{y}_i$ is the predicted value from the regression model for the $i$-th observation.

## Interpretation of Residuals:

- **Positive Residual:** The actual value is greater than the predicted value, meaning the model under-predicted for this data point.

- **Negative Residual:** The actual value is less than the predicted value, meaning the model over-predicted for this data point.

- **Zero Residual:** The actual and predicted values are equal, indicating a perfect prediction for this data point.

### Gradient Descent:

- Gradient descent is used to minimize the MSE by calculating the gradient of the cost function.

- A regression model uses gradient descent to update the coefficients of the line by reducing the cost function.
- It is done by a random selection of values of coefficient and then iteratively update the values to reach the minimum cost function.

# Model Performance:

The Goodness of fit determines how the line of regression fits the set of observations. The process of finding the best model out of various models is called **optimization**. It can be achieved by below method:

**1. R-squared method:**

- R-squared is a statistical method that determines the goodness of fit.
- It measures the strength of the relationship between the dependent and independent variables on a scale of 0-100%.
- The high value of R-square determines the less difference between the predicted values and actual values and hence represents a good model.
- It is also called a **coefficient of determination,** or **coefficient of multiple determination** for multiple regression.
- It can be calculated from the below formula:

$$R\text{-squared} = \frac{\text{Explained variation}}{\text{Total Variation}}$$

# Assumptions of Linear Regression

Below are some important assumptions of Linear Regression. These are some formal checks while building a Linear Regression model, which ensures to get the best possible result from the given dataset.

- **Linear relationship between the features and target:**
  Linear regression assumes the linear relationship between the dependent and independent variables.
- **Small or no multicollinearity between the features:**
  Multicollinearity means high-correlation between the independent variables. Due to multicollinearity, it may difficult to find the true relationship between the predictors and target variables. Or we can say, it is difficult to determine which predictor variable is affecting the target variable and which is not. So, the model assumes either little or no multicollinearity between the features or independent variables.
- **Homoscedasticity Assumption:**
  Homoscedasticity is a situation when the error term is the same for all the values of independent variables. With homoscedasticity, there should be no clear pattern distribution of data in the scatter plot.
- **Normal distribution of error terms:**
  Linear regression assumes that the error term should follow the normal distribution pattern. If error terms are not normally distributed, then confidence intervals will become either too wide or too narrow, which may cause difficulties in finding coefficients.

It can be checked using the **q-q plot**. If the plot shows a straight line without any deviation, which means the error is normally distributed.
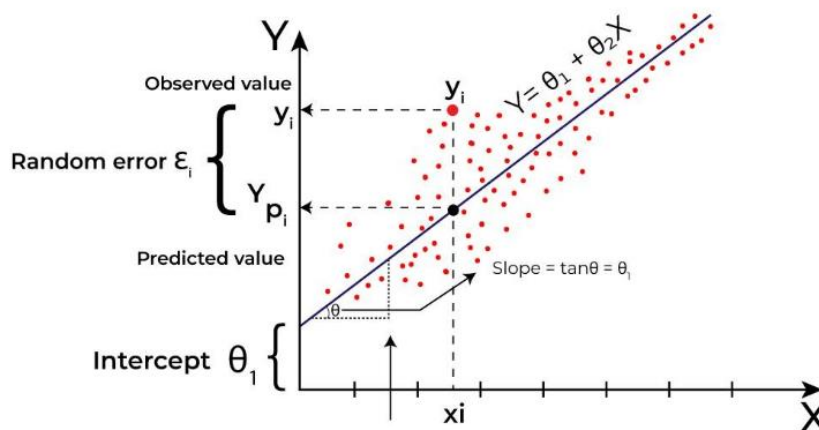
- **No autocorrelations:**
  The linear regression model assumes no autocorrelation in error terms. If there will be any correlation in the error term, then it will drastically reduce the accuracy of the model. Autocorrelation usually occurs if there is a dependency between residual errors.

## Least Squares

**Least Square method** is a fundamental mathematical technique widely used in **data analysis, statistics, and regression modeling** to identify the **best-fitting curve or line** for a given set of data points. This method ensures that the overall error is reduced, providing a highly accurate model for predicting future data trends.

**How to find the best fit line in linear regression**

To find the best fit line in linear regression, we use methods like Ordinary Least Squares (OLS) to minimize the errors between the actual data points and the predictions made by the line. Here's a step-by-step guide on how to find the best fit line:



**Steps to Find the Best Fit Line:**

## 1. Formulate the Linear Model:

In linear regression, we assume the relationship between the dependent variable $y$ and the independent variable $x$ is linear. The equation of the line can be expressed as:

$$y = \beta_0 + \beta_1 x$$

where:

- $y$ is the dependent variable,

- $x$ is the independent variable,

- $\beta_0$ is the intercept,

- $\beta_1$ is the slope (the change in $y$ for each unit increase in $x$).

## 2. Determine the Cost Function:

The goal is to find the values of $\beta_0$ and $\beta_1$ that minimize the **sum of squared residuals** (or errors). The cost function (also called the Mean Squared Error, MSE) is:

$$J(\beta_0, \beta_1) = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

where:

- $y_i$ is the actual value of the dependent variable,

- $\hat{y}_i$ is the predicted value from the line,

- $n$ is the number of observations.

The smaller this cost, the better the fit of the line.

### 3. Find the Parameters Using Ordinary Least Squares (OLS):

The ordinary least squares (OLS) method calculates the values of the slope ($\beta_1$) and intercept ($\beta_0$) that minimize the cost function. The formulas for these parameters are:

- **Slope ($\beta_1$):**

$$\beta_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n}(x_i - \bar{x})^2}$$

Where:

- $\bar{x}$ is the mean of the independent variable $x$,

- $\bar{y}$ is the mean of the dependent variable $y$.

- **Intercept ($\beta_0$):**

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

These formulas give the values of the slope and intercept that minimize the sum of squared residuals, leading to the best fit line.

## 4. Calculate Predictions:

Once you have $\beta_0$ and $\beta_1$, you can predict $y$ values for any $x$ using the formula:

$$\hat{y} = \beta_0 + \beta_1 x$$

## Example Calculation:

Let's calculate the best fit line for a small dataset:

**Example Data:**

$$x = [1, 2, 3, 4, 5]$$

$$y = [2, 4, 5, 4, 5]$$

**Step 1: Calculate Means:**

$$\bar{x} = \frac{1 + 2 + 3 + 4 + 5}{5} = 3$$

$$\bar{y} = \frac{2 + 4 + 5 + 4 + 5}{5} = 4$$

**Step 2: Calculate Slope ($\beta_1$):**

$$\beta_1 = \frac{(1-3)(2-4) + (2-3)(4-4) + (3-3)(5-4) + (4-3)(4-4)}{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)}$$

$$\beta_1 = \frac{(-2)(-2) + (-1)(0) + (0)(1) + (1)(0) + (2)(1)}{(-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2}$$

$$\beta_1 = \frac{4 + 0 + 0 + 0 + 2}{4 + 1 + 0 + 1 + 4} = \frac{6}{10} = 0.6$$

Step 3: Calculate Intercept ($\beta_0$):

$$\beta_0 = \bar{y} - \beta_1\bar{x} = 4 - (0.6 \times 3) = 4 - 1.8 = 2.2$$

Step 4: Equation of the Best Fit Line:

Now we have:

$$\hat{y} = 2.2 + 0.6x$$

This is the equation of the best fit line for the given data.

**Extra reference link:**

**https://www.geeksforgeeks.org/least-square-method/**

**https://medium.com/@rndayala/linear-regression-a00514bc45b0**

**https://www.javatpoint.com/regression-analysis-in-machine-learning**

**https://www.javatpoint.com/linear-regression-in-machine-learning**