**Machine Learning**

**(IV CSE – I SEM.)**

**A.Y.: 2022 – 2023**

**UNIT – II**

**Supervised Learning**

Learning a Class from Examples, Linear, Non-linear, Multi-class and Multi-label classification, Decision Trees: ID3, Classification and Regression Trees (CART), Regression: Linear Regression, Multiple Linear Regression, Logistic Regression.
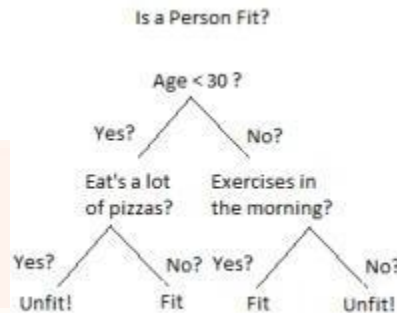
# UNIT – II
# Supervised Learning

**Decision Tree**

   **Introduction** Decision Trees are a type of Supervised Machine Learning (that is you explain what the input is and what the corresponding output is in the training data) where the data is continuously split according to a certain parameter. The tree can be explained by two entities, namely **decision nodes and leaves**. The leaves are the decisions or the final outcomes. And the decision nodes are where the data is split.



An example of a decision tree can be explained using above binary tree. Let's say you want to predict whether a person is fit given their information like age, eating habit, and physical activity, etc. The decision nodes here are questions like 'What's the age?', 'Does he exercise?', and 'Does he eat a lot of pizzas'? And the leaves, which are outcomes like either 'fit', or 'unfit'. In this case this was a binary classification problem (a yes no type problem). There are two main types of Decision Trees:

1. **Classification trees** (Yes/No types)

What we have seen above is an example of classification tree, where the outcome was a variable like 'fit' or 'unfit'. Here the decision variable is **Categorical**.

2. **Regression trees** (Continuous data types)

Here the decision or the outcome variable is **Continuous**, e.g. a number like 123. **Working** Now that we know what a Decision Tree is, we'll see how it works internally. There are many algorithms out there which construct Decision Trees, but one of the best is called as **ID3 Algorithm**. ID3 Stands for **Iterative Dichotomiser 3**. Before discussing the ID3 algorithm, we'll go through few definitions. **Entropy** Entropy, also called as Shannon Entropy is denoted by H(S) for a finite set S, is the measure of the amount of uncertainty or randomness in data.

$$H(S) = \sum_{x \in X} p(x) \log_2 \frac{1}{p(x)}$$

Intuitively, it tells us about the predictability of a certain event. Example, consider a coin toss whose probability of heads is 0.5 and probability of tails is 0.5. Here the entropy is the highest possible, since there's no way of determining what the outcome might be. Alternatively, consider a coin which has heads on both the sides, the entropy of such an event can be predicted

perfectly since we know beforehand that it'll always be heads. In other words, this event has **no randomness** hence it's entropy is zero. In particular, lower values imply less uncertainty while higher values imply high uncertainty. **Information Gain** Information gain is also called as Kullback-Leibler divergence denoted by IG(S,A) for a set S is the effective change in entropy after deciding on a particular attribute A. It measures the relative change in entropy with respect to the independent variables.

*Information Gain Formula*

where IG(S, A) is the information gain by applying feature A. H(S) is the Entropy of the entire set, while the second term calculates the Entropy after applying the feature A, where P(x) is the probability of event x. Let's understand this with the help of an example Consider a piece of data collected over the course of 14 days where the features are Outlook, Temperature, Humidity, Wind and the outcome variable is whether Golf was played on the day. Now, our job is to build a predictive model which takes in above 4 parameters and predicts whether Golf will be played on the day. We'll build a decision treeto do that using **ID3 algorithm.**

| Day | Outlook | Temperature | Humidity | Wind | Play Golf |
|-----|---------|-------------|----------|------|-----------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

**ID3**

ID3 Algorithm will perform following tasks recursively

1. Create root node for the tree
2. If all examples are positive, return leaf node „positive"
3. Else if all examples are negative, return leaf node „negative"
4. Calculate the entropy of current state H(S)
5. For each attribute, calculate the entropy with respect to the attribute „x" denoted by H(S, x)
6. Select the attribute which has maximum value of IG(S, x)
7. Remove the attribute that offers highest IG from the set of attributes

8. Repeat until we run out of all attributes, or the decision tree has all leaf nodes.

Now we'll go ahead and grow the decision tree. The initial step is to calculate H(S), the Entropy of the current state. In the above example, we can see in total there are 5 No's and 9 Yes's.

| Yes | No | Total |
|---|---|---|
| 9 | 5 | 14 |

$$Entropy(S) = \sum_{x \in X} p(x) \log_2 \frac{1}{p(x)}$$

$$Entropy(S) = -\left(\frac{9}{14}\right) \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right)$$

$$= 0.940$$

where "x" are the possible values for an attribute. Here, attribute "Wind" takes two possible values in the sample data, hence x = {Weak, Strong} we'll have to calculate:

1. $H(S_{weak})$
2. $H(S_{strong})$
3. $P(S_{weak})$
4. $P(S_{strong})$
5. $H(S) = 0.94$ which we had already calculated in the previous example

Amongst all the 14 examples we have **8 places where the wind is weak and 6 where the wind is Strong**.

| Wind = Weak | Wind = Strong | Total |
|---|---|---|
| 8 | 6 | 14 |

$$P(S_{weak}) = \frac{Number\ of\ Weak}{Total}$$

$$= \frac{8}{14}$$

$$P(S_{strong}) = \frac{Number\ of\ Strong}{Total}$$

$$= \frac{6}{14}$$

Now out of the 8 Weak examples, 6 of them were "Yes" for Play Golf and 2 of them were "No" for "Play Golf". So, we have,

$$Entropy(S_{weak}) = -\left(\frac{6}{8}\right) \log_2\left(\frac{6}{8}\right) - \left(\frac{2}{8}\right) \log_2\left(\frac{2}{8}\right)$$

$$= 0.811$$

Similarly, out of 6 Strong examples, we have **3 examples where the outcome was "Yes" for Play Golf and 3 where we had "No" for Play Golf**.

$$Entropy(S_{strong}) = -\left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right) - \left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right)$$
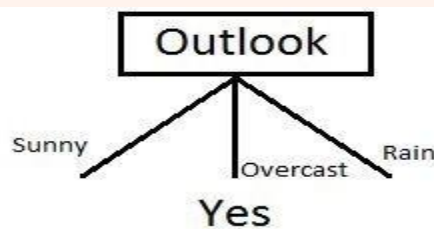
$$= 1.000$$

Remember, here half items belong to one class while other half belong to other. Hence we have perfect randomness. Now we have all the pieces required to calculate the Information Gain,

$$IG(S, Wind) = H(S) - \sum_{i=0}^{n} P(x) * H(x)$$

$$IG(S, Wind) = H(S) - P(S_{weak}) * H(S_{weak}) - P(S_{strong}) * H(S_{strong})$$

$$= 0.940 - \left(\frac{8}{14}\right)(0.811) - \left(\frac{6}{14}\right)(1.00)$$

$$= 0.048$$

Which tells us the Information Gain by considering „Wind" as the feature and give us information gain of **0.048**. Now we must similarly calculate the Information Gain for all the features.

$$IG(S, Outlook) = 0.246$$
$$IG(S, Temperature) = 0.029$$
$$IG(S, Humidity) = 0.151$$
$$IG(S, Wind) = 0.048 \text{ (Previous example)}$$

We can clearly see that IG(S, Outlook) has the highest information gain of 0.246, **hence we chose Outlookattribute as the root node**. At this point, the decision tree looks like.

Here we observe that whenever the outlook is Overcast, Play Golf is always 'Yes', it's no coincidence by any chance, the simple tree resulted because of **the highest information gain is given by the attribute Outlook**. Now how do we proceed from this point? We can simply apply **recursion**, you might want to look at the algorithm steps described earlier. Now that we've used Outlook, we've got three of them remaining Humidity, Temperature, and Wind. And, we had three possible values of Outlook: Sunny, Overcast, Rain. Where the Overcast node already ended up having leaf node 'Yes', so we're left with two subtrees to compute: Sunny and Rain.
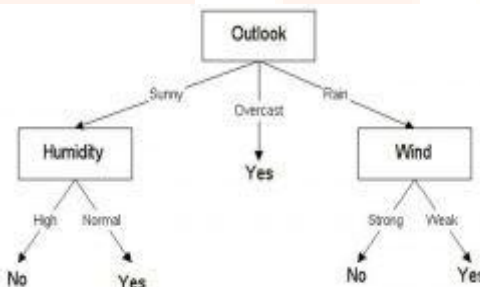
Next step would be computing $H(S_{sunny})$.

Table where the value of Outlook is Sunny looks like:

| Temperature | Humidity | Wind | Play Golf |
|---|---|---|---|
| Hot | High | Weak | No |
| Hot | High | Strong | No |
| Mild | High | Weak | No |
| Cool | Normal | Weak | Yes |
| Mild | Normal | Strong | Yes |

$$H(S_{sunny}) = \left(\frac{3}{5}\right)\log_2\left(\frac{3}{5}\right) - \left(\frac{2}{5}\right)\log_2\left(\frac{2}{5}\right) = 0.96$$

As we can see the **highest Information Gain is given by Humidity**. Proceeding in the same way with $S_{rain}$

will give us Wind as the one with highest information gain. The final Decision Tree looks something like this. The final Decision Tree looks something like this.
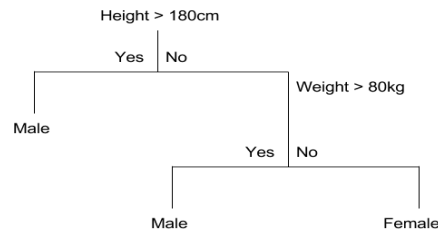


## Classification and Regression Trees

### Classification Trees

A classification tree is an algorithm where the target variable is fixed or categorical. The algorithm is then used to identify the "class" within which a target variable would most likely fall.

An example of a classification-type problem would be determining who will or will not subscribe to a digital platform; or who will or will not graduate from high school.
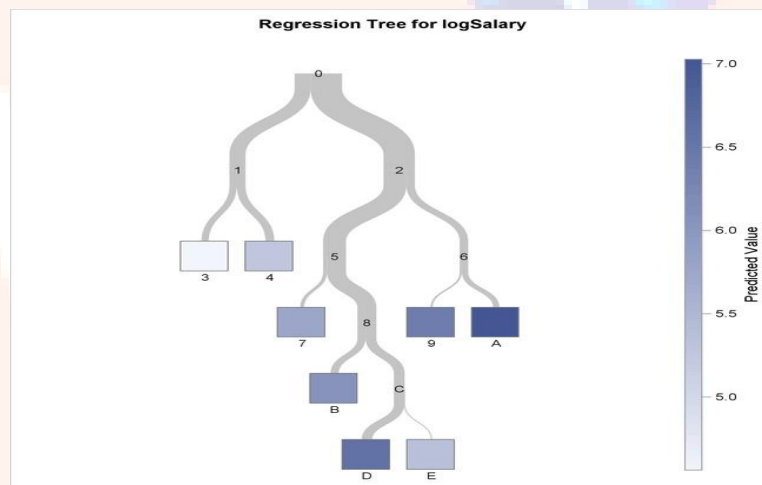
These are examples of simple binary classifications where the categorical dependent variable can assume only one of two, mutually exclusive values. In other cases, you might have to predict among a number of different variables. For instance, you may have to predict which type of smartphone a consumer may decide to purchase.

In such cases, there are multiple values for the categorical dependent variable. Here's what a classic classification tree looks like

Height > 180cm
Yes | No
Male
Weight > 80kg
Yes | No
Male    Female

## Regression Trees

A regression tree refers to an algorithm where the target variable is and the algorithm is used to predict it's value. As an example of a regression type problem, you may want to predict the selling prices of a residential house, which is a continuous dependent variable.

This will depend on both continuous factors like square footage as well as categorical factors like the style of home, area in which the property is located and so on.



## When to use Classification and Regression Trees

Classification trees are used when the dataset needs to be split into classes which belong to the response variable. In many cases, the classes Yes or No.

In other words, they are just two and mutually exclusive. In some cases, there may be more than two classes in which case a variant of the classification tree algorithm is used.

Regression trees, on the other hand, are used when the response variable is continuous. For instance, if the response variable is something like the price of a property or the temperature of the day, a regression tree is used.

In other words, regression trees are used for prediction-type problems while classification trees are used for classification-type problems.
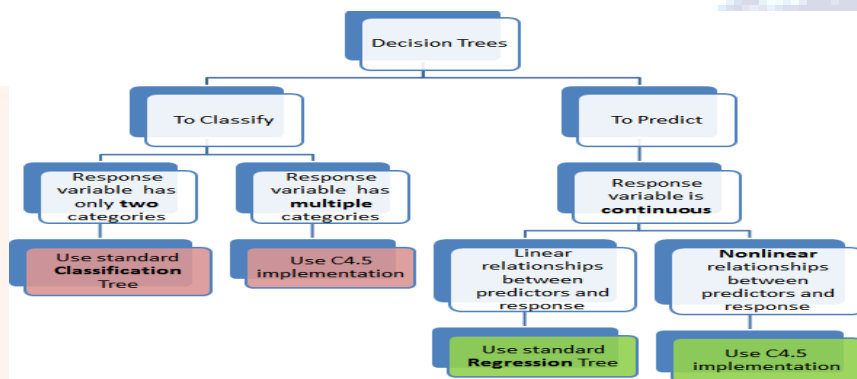
## How Classification and Regression Trees Work

A classification tree splits the dataset based on the homogeneity of data. Say, for instance, there are two variables; income and age; which determine whether or not a consumer will buy a particular kind of phone.

If the training data shows that 95% of people who are older than 30 bought the phone, the data

gets split there and age becomes a top node in the tree. This split makes the data "95% pure". Measures of impurity like entropy or Gini index are used to quantify the homogeneity of the data when it comes to classification trees.

In a regression tree, a regression model is fit to the target variable using each of the independent variables. After this, the data is split at several points for each independent variable.

At each such point, the error between the predicted values and actual values is squared to get "A Sum of Squared Errors" (SSE). The SSE is compared across the variables and the variable or point which has the lowest SSE is chosen as the split point. This process is continued recursively.



## Advantages of Classification and Regression Trees

The purpose of the analysis conducted by any classification or regression tree is to create a set of if-else conditions that allow for the accurate prediction or classification of a case.

### (i) The Results are Simplistic

The interpretation of results summarized in classification or regression trees is usually fairly simple. The simplicity of results helps in the following ways.

- It allows for the rapid classification of new observations. That's because it is much simpler to evaluate just one or two logical conditions than to compute scores using complex nonlinear equations for each group.
- It can often result in a simpler model which explains why the observations are either classified or predicted in a certain way. For instance, business problems are much easier to explain with if-then statements than with complex nonlinear equations.

### (ii) Classification and Regression Trees are Nonparametric & Nonlinear

The results from classification and regression trees can be summarized in simplistic if-then conditions. This negates the need for the following implicit assumptions.

- The predictor variables and the dependent variable are linear.
- The predictor variables and the dependent variable follow some specific nonlinear link function.
- The predictor variables and the dependent variable are monotonic.

Since there is no need for such implicit assumptions, classification and regression tree methods are well suited to data mining. This is because there is very little knowledge or assumptions that

can be made beforehand about how the different variables are related.

As a result, classification and regression trees can actually reveal relationships between these variables that would not have been possible using other techniques.

**(iii) Classification and Regression Trees Implicitly Perform Feature Selection**

Feature selection or variable screening is an important part of analytics. When we use decision trees, the top few nodes on which the tree is split are the most important variables within the set. As a result, feature selection gets performed automatically and we don't need to do it again.


**Limitations of Classification and Regression Trees**

Classification and regression tree tutorials, as well as classification and regression tree ppts, exist in abundance. This is a testament to the popularity of these decision trees and how frequently they are used. However, these decision trees are not without their disadvantages.

There are many classification and regression trees examples where the use of a decision tree has not led to the optimal result. Here are some of the limitations of classification and regression trees.

**(i) Overfitting**

Overfitting occurs when the tree takes into account a lot of noise that exists in the data and comes up with an inaccurate result.

**(ii) High variance**

In this case, a small variance in the data can lead to a very high variance in the prediction, thereby affecting the stability of the outcome.

**(iii) Low bias**

A decision tree that is very complex usually has a low bias. This makes it very difficult for the model to incorporate any new data.


**What is a CART in Machine Learning?**

A Classification and Regression Tree (CART) is a predictive algorithm used in machine learning. It explains how a target variable's values can be predicted based on other values.

It is a decision tree where each fork is a split in a predictor variable and each node at the end has a prediction for the target variable.

The CART algorithm is an important decision tree algorithm that lies at the foundation of machine learning. Moreover, it is also the basis for other powerful machine learning algorithms like bagged decision trees, random forest and boosted decision trees.

**Summing up**

The Classification and regression tree (CART) methodology is one of the oldest and most fundamental algorithms. It is used to predict outcomes based on certain predictor variables.

They are excellent for data mining tasks because they require very little data pre-processing. Decision tree models are easy to understand and implement which gives them a strong advantage when compared to other analytical models.


**Regression**

**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:
  o Prediction of rain using temperature and other factors
  o Determining Market trends
  o 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**.

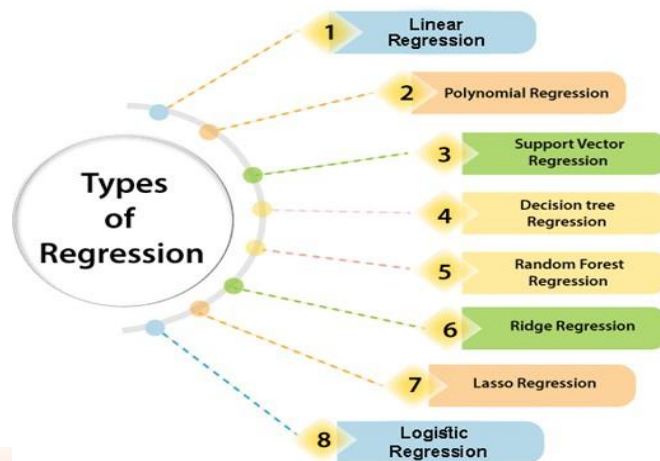**Why do we use Regression Analysis?**

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, theleast important factor, and how each factor is affecting the other factors**.
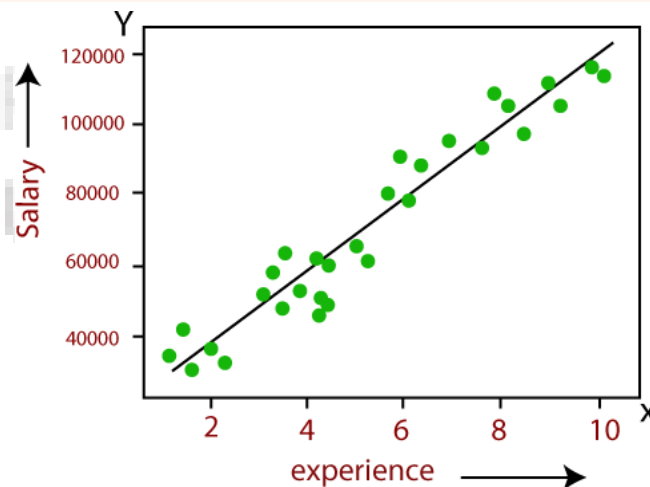
**Types of Regression**

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 **Logistic Regression**
- o **Polynomial Regression**
- o **Support Vector Regression**
- o **Decision Tree Regression**
- o **Random Forest Regression**
- o **Ridge 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 relationshipbetween 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 thedependent variable (Y-axis), hence called linear regression.

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

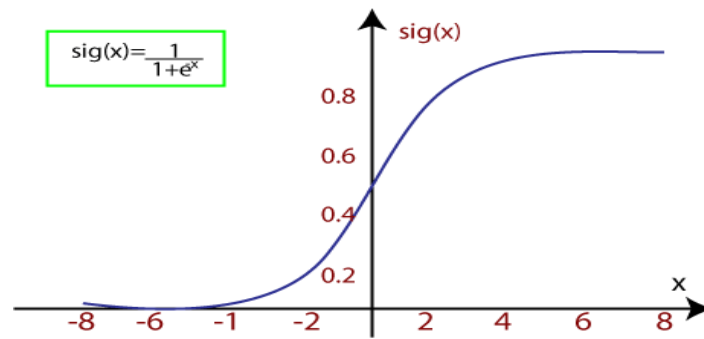- o Arriving at ETAs in traffic.

**Logistic Regression:**
- o 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 0or 1.

- o Logistic regression algorithm works with the categorical variable such as 0 or 1, Yes or No, True or False, Spam or not spam, etc.

- o It is a predictive analysis algorithm which works on the concept of probability.

- o Logistic regression is a type of regression, but it is different from the  linear regression algorithm in the term how they are used.

- o 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

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

function can be represented as:
- o f(x)= Output between the 0 and 1 value.

- o x= input to the function

- o e= base of natural logarithm.

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

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

- o 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.

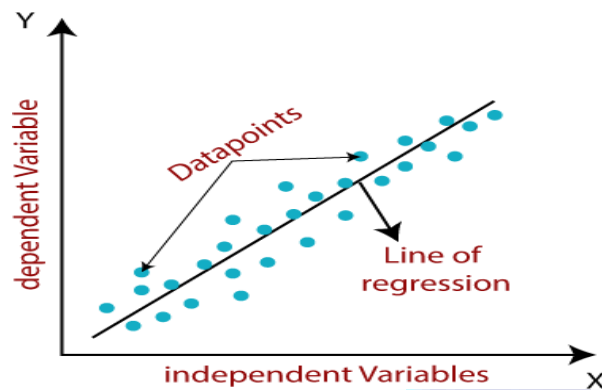There are three types of logistic regression:
- o **Binary(0/1, pass/fail)**

- o **Multi(cats, dogs, lions)**

- o **Ordinal(low, medium, high)**

**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.
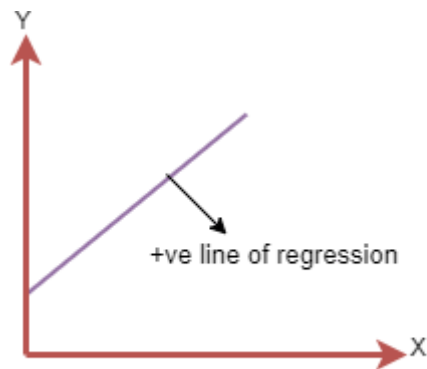
**Types of Linear Regression**

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

- o **Simple Linear Regression:**
  If a single independent variable is used to predict the value of a numerical dependent
  variable, then such aLinear 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, thensuch 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:
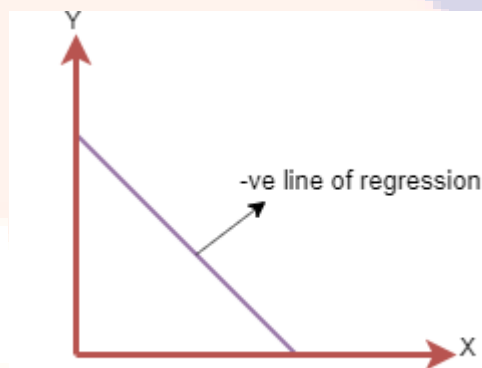
- o **Positive Linear Relationship:**
  If the dependent variable increases on the Y-axis and independent variable increases on
  X-axis, then such arelationship is termed as a Positive linear relationship.

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

- 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: $Y = -a_0 + a_1 x$

*Finding the best fitline:*

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 variableto 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 ofsquared error occurred between the predicted values and actual values. It can be written as:

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

$$\text{MSE} = 1\frac{1}{N} \sum_{i=1}^{n} (y_i - (a_1 x_i + a_0))^2$$

**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 arefar 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.

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

- o A regression model uses gradient descent to update the coefficients of the line by reducing the costfunction.

- o It is done by a random selection of values of coefficient and then iteratively update the values to reach theminimum cost function.

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

**1. R-squared method:**
- o R-squared is a statistical method that determines the goodness of fit.

- o It measures the strength of the relationship between the dependent and independent variables on a scale of0-100%.

- o The high value of R-square determines the less difference between the predicted values and actual valuesand hence represents a good model.

- o It is also called a **coefficient of determination,** or **coefficient of multiple determination** for multipleregression.

- o 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 aLinear Regression model, which ensures to get the best possible result from the given dataset.

- o **Linear relationship between the features and target:**
  Linear regression assumes the linear relationship between the dependent and independent variables.

- o **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.

- o **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.

- o **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.

- o **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.

## 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: $y = a_0 + a_1 x + \varepsilon$

Where,
$a0 =$ It is the intercept of the Regression line (can be obtained putting x=0)
$a1 =$ It is the slope of the regression line, which tells whether the line is increasing or decreasing. $\varepsilon =$ 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 $CO_2$ emission based on engine size and number of cylinders in a car.

**Some key points about MLR:**
- 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.

- Each feature variable must model the linear relationship with the dependent variable.

- 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 $x_1$, $x_2$, $x_3$, ...,$x_n$. Since it is an enhancement of Simple Linear Regression, so the same is applied for the multiple linear regression equation, the equation becomes:
$$Y = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + ...... b_n x_n \qquad (a)$$
**Where,**
Y= Output/Response variable
$b_0$, $b_1$, $b_2$, $b_3$ , $b_n$     = Coefficients of the model.
$x_1$, $x_2$, $x_3$, $x_4$, = Various Independent/feature variable
**Assumptions for Multiple Linear Regression:**
- A linear relationship should exist between the Target and predictor variables.
- The regression residuals must be normally distributed.
- MLR assumes little or no multicollinearity (correlation between the independent variable) in data.