**Day 04**

# 30 Days of Machine Learning

## Types of Regression Trees

# Prepared By



## Ahmed Ali

*Quaid-e-Awam University of Engineering Science & Technology Nawabshah*
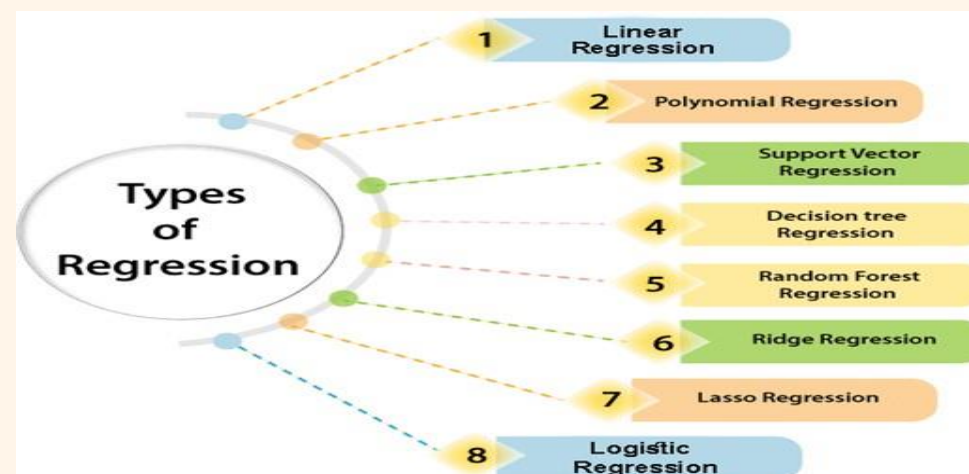
# Types of Regression Tree

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.
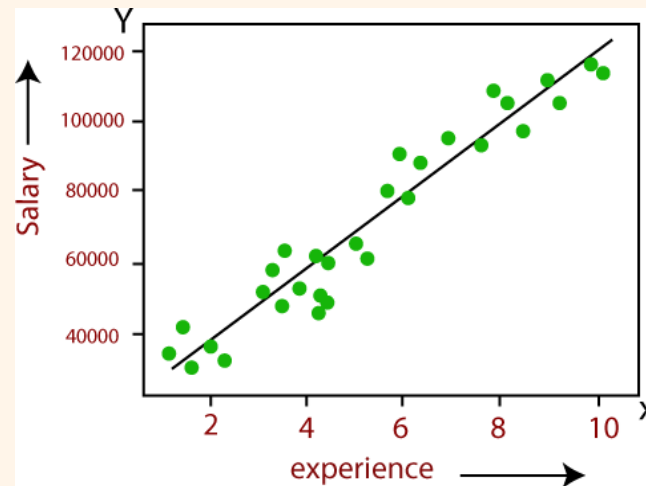
It is one of the very simple and easy algorithms which works on regression and shows the relationship between the continuous variables.

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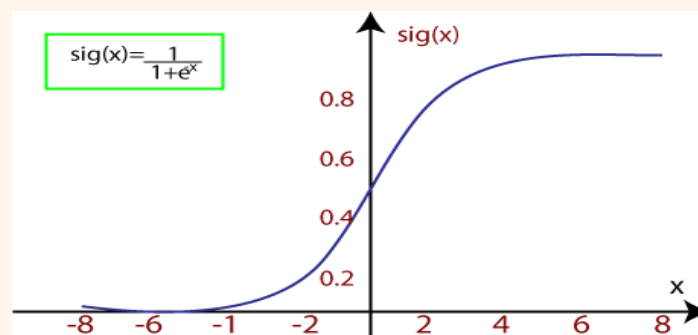
$$f(x) = \frac{1}{1 + e^{-x}}$$

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.

Ahmed Ali

**There are three types of logistic regression:**

- Binary (0/1, pass/fail)
- Multi (cats, dogs, lions)
- 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.

Mathematically, we can represent a linear regression as:

$$y = a0 + a1x + \epsilon$$

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). $\epsilon$ = random error
The values for x and y variables are training datasets for Linear Regression model representation.

Ahmed Ali

## *Types of Linear Regression*

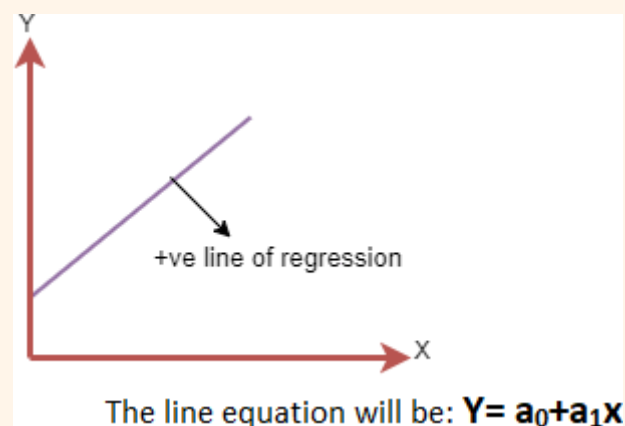Linear regression can be further divided into two types of the algorithm:

- **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.

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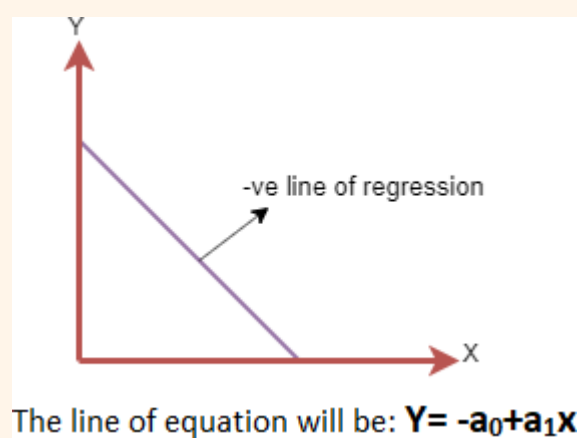
- *Positive Linear Relationship:*
  If the dependent variable increases on the Y-axis and independent variable increases on X-axis, then such a relationship is termed as a Positive linear relationship.



+ve line of regression

The line equation will be: $Y = a_0 + a_1x$

- *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.



-ve line of regression

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-

- 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:

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

## Where:
N=Total number of observations
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.

## *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.

Ahmed Ali

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.

## 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 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)

*UP NEXT:*
*MULTIPLE REGRESSION.*

Ahmed Ali