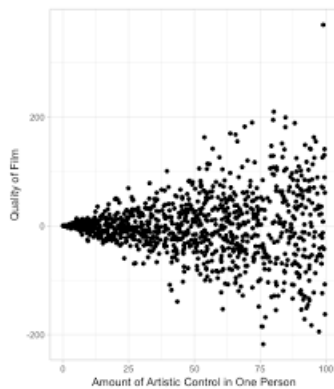


# **DATA SCIENCE INTERVIEW PREPARATION (30 Days of Interview Preparation) # DAY 03**

## Q1. How do you treat heteroscedasticity in regression?

Heteroscedasticity means unequal scattered distribution. In regression analysis, we generally talk about the heteroscedasticity in the context of the error term. Heteroscedasticity is the systematic change in the spread of the residuals or errors over the range of measured values. Heteroscedasticity is the problem because *Ordinary least squares (OLS)* regression assumes that all residuals are drawn from a random population that has a constant variance.



What causes Heteroscedasticity?

Heteroscedasticity occurs more often in datasets, where we have a large range between the largest and the smallest observed values. There are many reasons why heteroscedasticity can exist, and a generic explanation is that the error variance changes proportionally with a factor.

We can categorize Heteroscedasticity into two general types:-

**Pure heteroscedasticity:-** It refers to cases where we specify the correct model and let us observe the non-constant variance in residual plots.

**Impure heteroscedasticity:-** It refers to cases where you incorrectly specify the model, and that causes the non-constant variance. When you leave an important variable out of a model, the omitted effect is absorbed into the error term. If the effect of the omitted variable varies throughout the observed range of data, it can produce the telltale signs of heteroscedasticity in the residual plots.

### *How to Fix Heteroscedasticity*

Redefining the variables:

If your model is a cross-sectional model that includes large differences between the sizes of the observations, you can find different ways to specify the model that reduces the impact of the size

differential. To do this, change the model from using the raw measure to using rates and per capita values. Of course, this type of model answers a slightly different kind of question. You'll need to determine whether this approach is suitable for both your data and what you need to learn.

**Weighted regression:**

It is a method that assigns each data point to a weight based on the variance of its fitted value. The idea is to give small weights to observations associated with higher variances to shrink their squared residuals. Weighted regression minimizes the sum of the weighted squared residuals. When you use the correct weights, heteroscedasticity is replaced by homoscedasticity.

## Q2. What is multicollinearity, and how do you treat it?

**Multicollinearity** means independent variables are highly correlated to each other. In regression analysis, it's an important assumption that the regression model should not be faced with a problem of multicollinearity.

If two explanatory variables are highly correlated, it's hard to tell, which affects the dependent variable. Let's say Y is regressed against X1 and X2 and where X1 and X2 are highly correlated. Then the effect of X1 on Y is hard to distinguish from the effect of X2 on Y because any increase in X1 tends to be associated with an increase in X2.

Another way to look at the multicollinearity problem is: Individual t-test P values can be misleading. It means a P-value can be high, which means the variable is not important, even though the variable is important.

**Correcting Multicollinearity:**

- 1) Remove one of the highly correlated independent variables from the model. If you have two or more factors with a high VIF, remove one from the model.
- 2) Principle Component Analysis (PCA) - It cut the number of interdependent variables to a smaller set of uncorrelated components. Instead of using highly correlated variables, use components in the model that have eigenvalue greater than 1.
- 3) Run PROC VARCLUS and choose the variable that has a minimum (1-R<sup>2</sup>) ratio within a cluster.
- 4) Ridge Regression - It is a technique for analyzing multiple regression data that suffer from multicollinearity.
- 5) If you include an interaction term (the product of two independent variables), you can also reduce multicollinearity by "centering" the variables. By "centering," it means subtracting the mean from the values of the independent variable before creating the products.

### **When is multicollinearity not a problem?**

- 1) If your goal is to predict Y from a set of X variables, then multicollinearity is not a problem. The predictions will still be accurate, and the overall R<sup>2</sup> (or adjusted R<sup>2</sup>) quantifies how well the model predicts the Y values.
- 2) Multiple dummy (binary) variables that represent a categorical variable with three or more categories.

### **Q3. What is market basket analysis? How would you do it in Python?**

Market basket analysis is the study of items that are purchased or grouped in a single transaction or multiple, sequential transactions. Understanding the relationships and the strength of those relationships is valuable information that can be used to make recommendations, cross-sell, up-sell, offer coupons, etc.

Market Basket Analysis is one of the key techniques used by large retailers to uncover associations between items. It works by looking for combinations of items that occur together frequently in transactions. To put it another way, it allows retailers to identify relationships between the items that people buy.

### **Q4. What is Association Analysis? Where is it used?**

*Association analysis uses a set of transactions to discover rules that indicate the likely occurrence of an item based on the occurrences of other items in the transaction.*

The technique of association rules is widely used for retail basket analysis. It can also be used for classification by using rules with class labels on the right-hand side. It is even used for outlier detection with rules indicating infrequent/abnormal association.

Association analysis also helps us to identify cross-selling opportunities, for example, we can use the rules resulting from the analysis to place associated products together in a catalog, in the supermarket, or the Webshop, or apply them when targeting a marketing campaign for product B at customers who have already purchased product A.

Association rules are given in the form as below:

$A \Rightarrow B[\text{Support}, \text{Confidence}]$  The part before  $\Rightarrow$  is referred to as if (Antecedent) and the part after  $\Rightarrow$  is referred to as then (Consequent).

Where A and B are sets of items in the transaction data, a and B are disjoint sets.

Computer  $\Rightarrow$  Anti-virus Software [Support=20%, confidence=60%] Above rule says:

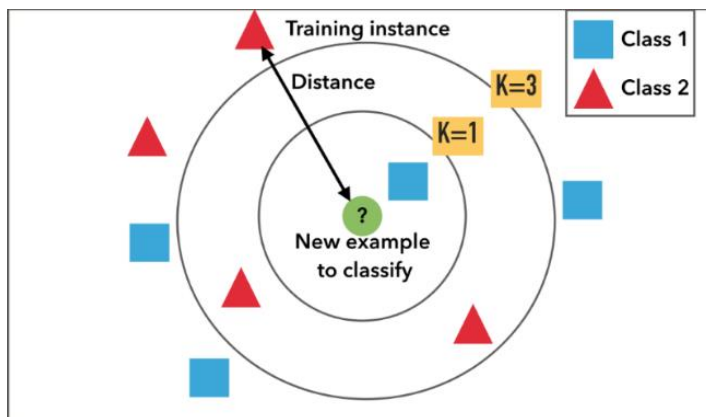
1. 20% transaction show Anti-virus software is bought with purchase of a Computer
2. 60% of customers who purchase Anti-virus software is bought with purchase of a Computer

An example of Association Rules \* Assume there are 100 customers

1. 10 of them bought milk, 8 bought butter and 6 bought both of them 2 .bought milk  $\Rightarrow$  bought butter
2. support =  $P(\text{Milk} \& \text{Butter}) = 6/100 = 0.06$
3. confidence =  $\text{support}/P(\text{Butter}) = 0.06/0.08 = 0.75$
4. lift =  $\text{confidence}/P(\text{Milk}) = 0.75/0.10 = 7.5$

## Q5. What is KNN Classifier ?

KNN means **K-Nearest Neighbour** Algorithm. It can be used for both classification and regression.



It is the simplest machine learning algorithm. Also known as **lazy learning** (why? Because it does not create a generalized model during the time of training, so the testing phase is very important where it does the actual job. Hence Testing is very costly - in terms of time & money). Also called an instance-based or memory-based learning

In k-NN classification, the output is a class membership. An object is classified by a plurality vote of its neighbors, with the object being assigned to the class most common among its k nearest neighbors (k is a positive integer, typically small). If k = 1, then the object is assigned to the class of that single nearest neighbor.

In **k-NN regression**, the output is the property value for the object. This value is the average of the values of k nearest neighbors.

#### Distance functions

Euclidean  $\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$

Manhattan  $\sum_{i=1}^k |x_i - y_i|$

Minkowski  $\left( \sum_{i=1}^k (|x_i - y_i|)^q \right)^{1/q}$

All three distance measures are only valid for continuous variables. In the instance of categorical variables, the Hamming distance must be used.

#### Hamming Distance

$$D_H = \sum_{i=1}^k |x_i - y_i|$$

$$x = y \Rightarrow D = 0$$

$$x \neq y \Rightarrow D = 1$$

X	Y	Distance
Male	Male	0
Male	Female	1

**How to choose the value of K:** K value is a hyperparameter which needs to choose during the time of model building

Also, a small number of neighbors are most flexible fit, which will have a low bias, but the high variance and a large number of neighbors will have a smoother decision boundary, which means lower variance but higher bias.

We should choose an odd number if the number of classes is even. It is said the most common values are to be 3 & 5.

## Q6. What is Pipeline in sklearn ?

A pipeline is what chains several steps together, once the initial exploration is done. For example, some codes are meant to transform features—normalize numerically, or turn text into vectors, or fill up missing data, and they are transformers; other codes are meant to predict variables by fitting an algorithm,

such as random forest or support vector machine, they are estimators. Pipeline chains all these together, which can then be applied to training data in block.

Example of a pipeline that imputes data with the most frequent value of each column, and then fit a decision tree classifier.

```
from sklearn.pipeline import Pipeline
steps = [('imputation', Imputer(missing_values='NaN', strategy = 'most_frequent', axis=0)),
        ('clf', DecisionTreeClassifier())]
pipeline = Pipeline(steps)
clf = pipeline.fit(X_train,y_train)``
```

Instead of fitting to one model, it can be looped over several models to find the best one.

```
classifiers = [ KNeighborsClassifier(5), RandomForestClassifier(), GradientBoostingClassifier()]
```

for clf in classifiers:

```
steps = [('imputation', Imputer(missing_values='NaN', strategy = 'most_frequent', axis=0)),
        ('clf', clf)]
```

```
pipeline = Pipeline(steps)
```

I also learned the pipeline itself can be used as an estimator and passed to cross-validation or grid search.

```
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
kfold = KFold(n_splits=10, random_state=seed)
results = cross_val_score(pipeline, X_train, y_train, cv=kfold)
print(results.mean())
```

## Q7. What is Principal Component Analysis(PCA), and why we do?

The main idea of principal component analysis (PCA) is to reduce the dimensionality of a data set consisting of many variables correlated with each other, either heavily or lightly, while retaining the variation present in the dataset, up to the maximum extent. The same is done by transforming the variables to a new set of variables, which are known as the principal components (or simply, the PCs) and are orthogonal, ordered such that the retention of variation present in the original variables decreases as we move down in the order. So, in this way, the 1st principal component retains maximum variation that was present in the original components. The principal components are the eigenvectors of a covariance matrix, and hence they are orthogonal.

Main important points to be considered:

1. Normalize the data
2. Calculate the covariance matrix
3. Calculate the eigenvalues and eigenvectors
4. Choosing components and forming a feature vector
5. Forming Principal Components

## Q8. What is t-SNE?

(t-SNE) t-Distributed Stochastic Neighbor Embedding is a non-linear dimensionality reduction algorithm used for exploring high-dimensional data. It maps multi-dimensional data to two or more dimensions suitable for human observation. With the help of the t-SNE algorithms, you may have to plot fewer exploratory data analysis plots next time you work with high dimensional data.

## Q9. VIF(Variation Inflation Factor),Weight of Evidence & Information Value. Why and when to use?

### Variation Inflation Factor

It provides an index that measures how much the variance (the square of the estimate's standard deviation) of an estimated regression coefficient is increased because of collinearity.

$VIF = 1 / (1 - R\text{-Square of } j\text{-th variable})$  where  $R^2$  of  $j$ th variable is the coefficient of determination of the model that includes all independent variables except the  $j$ th predictor.

Where  $R\text{-Square of } j\text{-th variable}$  is the multiple  $R^2$  for the regression of  $X_j$  on the other independent variables (a regression that does not involve the dependent variable  $Y$ ).

If  $VIF > 5$ , then there is a problem with multicollinearity.

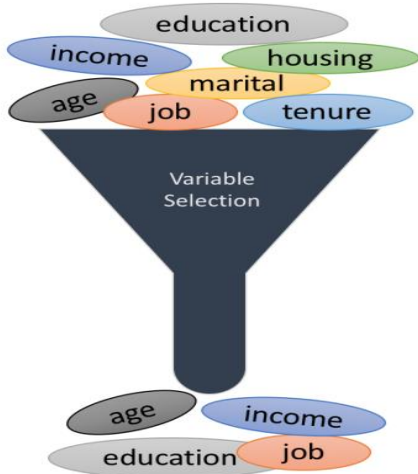
### Understanding VIF

If the variance inflation factor of a predictor variable is 5 this means that variance for the coefficient of that predictor variable is 5 times as large as it would be if that predictor variable were uncorrelated with the other predictor variables.

In other words, if the variance inflation factor of a predictor variable is 5 this means that the standard error for the coefficient of that predictor variable is 2.23 times ( $\sqrt{5} = 2.23$ ) as large as it would be if that predictor variable were uncorrelated with the other predictor variables.



Weight of evidence (WOE) and information value (IV) are simple, yet powerful techniques to perform variable transformation and selection.



The formula to create WOE and IV is

$$WOE = \ln \left( \frac{\text{Event\%}}{\text{Non Event\%}} \right)$$

$$IV = \sum (\text{Event\%} - \text{Non Event\%}) * \ln \left( \frac{\text{Event\%}}{\text{Non Event\%}} \right)$$

Here is a simple table that shows how to calculate these values.

Variable Name	Min. Value	Max. Value	Count	# Event	# Non Event	Event%	Non event%	WOE	Event% - Non event%	IV
Age	10	20	1200	150	1050	28.3%	19.0%	0.3992	9.3%	0.03718
Age	21	30	900	120	780	22.6%	14.1%	0.4733	8.5%	0.04040
Age	31	40	1090	110	980	20.8%	17.7%	0.1580	3.0%	0.00479
Age	41	50	1460	100	1360	18.9%	24.6%	-0.2650	-5.7%	0.01517
Age	50	inf	1410	50	1360	9.4%	24.6%	-0.9582	-15.2%	0.14525
<b>Total</b>			<b>6060</b>	<b>530</b>	<b>5530</b>					<b>0.24279</b>

The IV value can be used to select variables quickly.

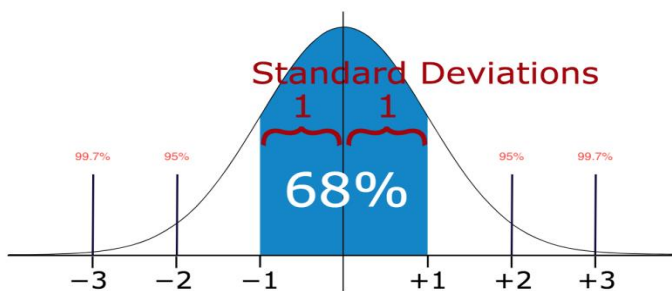
Information Value (IV)	Predictive Power
< 0.02	useless for prediction
0.02 to 0.1	weak predictor
0.1 to 0.3	medium predictor
0.3 to 0.5	strong predictor
> 0.5	suspicious or too good to be true

## Q10: How to evaluate that data does not have any outliers ?

In statistics, outliers are data points that don't belong to a certain population. It is an abnormal observation that lies far away from other values. An outlier is an observation that diverges from otherwise well-structured data.

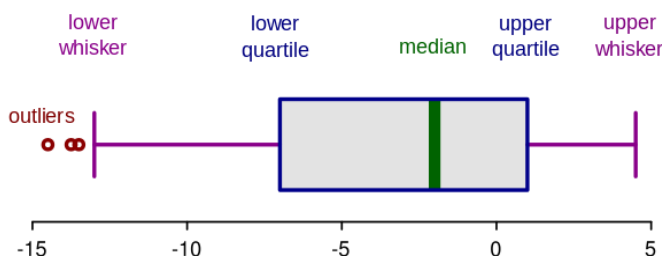
### Detection:

**Method 1 — Standard Deviation:** In statistics, If a data distribution is approximately normal, then about 68% of the data values lie within one standard deviation of the mean, and about 95% are within two standard deviations, and about 99.7% lie within three standard deviations.

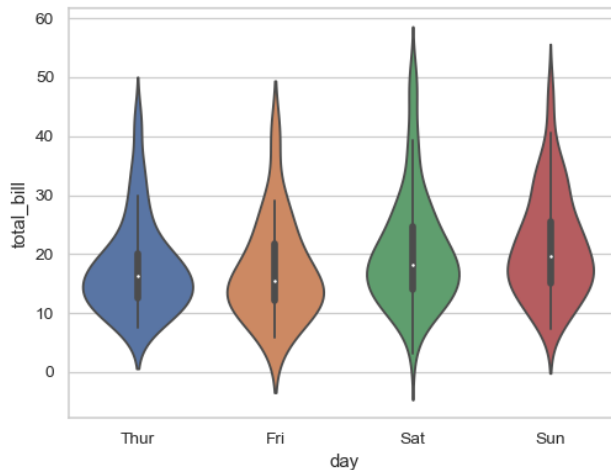


Therefore, if you have any data point that is more than 3 times the standard deviation, then those points are very likely to be anomalous or outliers.

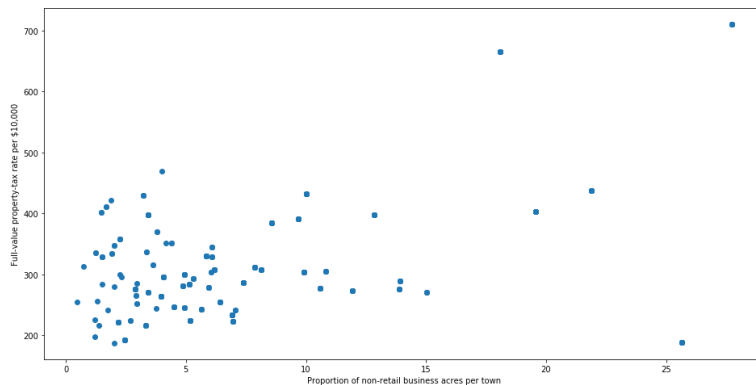
**Method 2 — Boxplots:** Box plots are a graphical depiction of numerical data through their quantiles. It is a very simple but effective way to visualize outliers. Think about the lower and upper whiskers as the boundaries of the data distribution. Any data points that show above or below the whiskers can be considered outliers or anomalous.



**Method 3 - Violin Plots:** Violin plots are similar to box plots, except that they also show the probability density of the data at different values, usually smoothed by a kernel density estimator. Typically a violin plot will include all the data that is in a box plot: a marker for the median of the data, a box or marker indicating the interquartile range, and possibly all sample points if the number of samples is not too high.



**Method 4 - Scatter Plots:** A scatter plot is a type of plot or mathematical diagram using Cartesian coordinates to display values for typically two variables for a set of data. The data are displayed as a collection of points, each having the value of one variable determining the position on the horizontal axis and the value of the other variable determining the position on the vertical axis.



The points which are very far away from the general spread of data and have a very few neighbors are considered to be outliers

## Q11: What you do if there are outliers?

Following are the approaches to handle the outliers:

1. Drop the outlier records
2. Assign a new value: If an outlier seems to be due to a mistake in your data, you try imputing a value.
3. If percentage-wise the number of outliers is less, but when we see numbers, there are several, then, in that case, dropping them might cause a loss in insight. We should group them in that case and run our analysis separately on them.

## Q12: What are the encoding techniques you have applied with Examples ?

In many practical data science activities, the data set will contain categorical variables. These variables are typically stored as text values". Since machine learning is based on mathematical equations, it would cause a problem when we keep categorical variables as is.

Let's consider the following dataset of fruit names and their weights.

Some of the common encoding techniques are:

**Label encoding:** In label encoding, we map each category to a number or a label. The labels chosen for the categories have no relationship. So categories that have some ties or are close to each other lose such information after encoding.

**One - hot encoding:** In this method, we map each category to a vector that contains 1 and 0 denoting the presence of the feature or not. The number of vectors depends on the categories which we want to keep. For high cardinality features, this method produces a lot of columns that slows down the learning significantly.

## Q13: Tradeoff between bias and variances, the relationship between them.

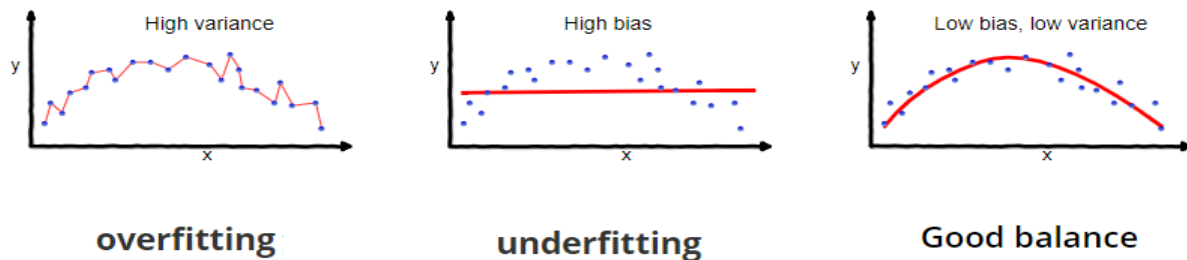
Whenever we discuss model prediction, it's important to understand prediction errors (bias and variance). The prediction error for any machine learning algorithm can be broken down into three parts:

- Bias Error
- Variance Error
- Irreducible Error

The irreducible error cannot be reduced regardless of what algorithm is used. It is the error introduced from the chosen framing of the problem and may be caused by factors like unknown variables that influence the mapping of the input variables to the output variable.

**Bias:** Bias means that the model favors one result more than the others. Bias is the simplifying assumptions made by a model to make the target function easier to learn. The model with high bias pays very little attention to the training data and oversimplifies the model. It always leads to a high error in training and test data.

**Variance:** Variance is the amount that the estimate of the target function will change if different training data was used. The model with high variance pays a lot of attention to training data and does not generalize on the data which it hasn't seen before. As a result, such models perform very well on training data but have high error rates on test data.



So, the end goal is to come up with a model that balances both Bias and Variance. This is called *Bias Variance Trade-off*. To build a good model, we need to find a good balance between bias and variance such that it minimizes the total error.

## Q14: What is the difference between Type 1 and Type 2 error and severity of the error?

### Type I Error

A Type I error is often referred to as a "false positive" and is the incorrect rejection of the true null hypothesis in favor of the alternative.

In the example above, the null hypothesis refers to the natural state of things or the absence of the tested effect or phenomenon, i.e., stating that the patient is HIV negative. The alternative hypothesis states that the patient is HIV positive. Many medical tests will have the disease they are testing for as the alternative hypothesis and the lack of that disease as the null hypothesis.

A Type I error would thus occur when the patient doesn't have the virus, but the test shows that they do. In other words, the test incorrectly rejects the true null hypothesis that the patient is HIV negative.

### **Type II Error**

A Type II error is the inverse of a Type I error and is the false acceptance of a null hypothesis that is not true, i.e., a false negative. A Type II error would entail the test telling the patient they are free of HIV when they are not.

Considering this HIV example, which error type do you think is more acceptable? In other words, would you rather have a test that was more prone to Type I or Types II error? With HIV, the momentary stress of a false positive is likely better than feeling relieved at a false negative and then failing to take steps to treat the disease. Pregnancy tests, blood tests, and any diagnostic tool that has serious consequences for the health of a patient are usually overly sensitive for this reason – they should err on the side of a false positive.

But in most fields of science, Type II errors are seen as less serious than Type I errors. With the Type II error, a chance to reject the null hypothesis was lost, and no conclusion is inferred from a non-rejected null. But the Type I error is more serious because you have wrongly rejected the null hypothesis and ultimately made a claim that is not true. In science, finding a phenomenon where there is none is more egregious than failing to find a phenomenon where there is.

## **Q15: What is binomial distribution and polynomial distribution?**

**Binomial Distribution:** A binomial distribution can be thought of as simply the probability of a SUCCESS or FAILURE outcome in an experiment or survey that is repeated multiple times. The binomial is a type of distribution that has two possible outcomes (the prefix “bi” means two, or twice). For example, a coin toss has only two possible outcomes: heads or tails, and taking a test could have two possible outcomes: pass or fail.

**Multinomial/Polynomial Distribution:** Multi or Poly means many. In probability theory, the multinomial distribution is a generalization of the binomial distribution. For example, it models the probability of counts of each side for rolling a k-sided die n times. For n independent trials each of which leads to success for exactly one of k categories, with each category having a given fixed success probability, the multinomial distribution gives the probability of any particular combination of numbers of successes for the various categories

## Q16: What is the Mean Median Mode standard deviation for the sample and population?

**Mean** It is an important technique in statistics. Arithmetic Mean can also be called an average. It is the number of the quantity obtained by summing two or more numbers/variables and then dividing the sum by the number of numbers/variables.

**Mode** The mode is also one of the types for finding the average. A mode is a number that occurs most frequently in a group of numbers. Some series might not have any mode; some might have two modes, which is called a bimodal series.

In the study of statistics, the three most common ‘averages’ in statistics are mean, median, and mode.

**Median** is also a way of finding the average of a group of data points. It’s the middle number of a set of numbers. There are two possibilities, the data points can be an odd number group, or it can be an even number group.

If the group is odd, arrange the numbers in the group from smallest to largest. The median will be the one which is exactly sitting in the middle, with an equal number on either side of it. If the group is even, arrange the numbers in order and pick the two middle numbers and add them then divide by 2. It will be the median number of that set.

**Standard Deviation (Sigma)** Standard Deviation is a measure of how much your data is spread out in statistics.

## Q17: What is Mean Absolute Error ?

**What is Absolute Error?** Absolute Error is the amount of error in your measurements. It is the difference between the measured value and the “true” value. For example, if a scale states 90 pounds, but you know your true weight is 89 pounds, then the scale has an absolute error of  $90 \text{ lbs} - 89 \text{ lbs} = 1 \text{ lbs}$ .

This can be caused by your scale, not measuring the exact amount you are trying to measure. For example, your scale may be accurate to the nearest pound. If you weigh 89.6 lbs, the scale may “round up” and give you 90 lbs. In this case the absolute error is  $90 \text{ lbs} - 89.6 \text{ lbs} = .4 \text{ lbs}$ .

**Mean Absolute Error** The Mean Absolute Error(MAE) is the average of all absolute errors. The formula is: mean absolute error

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |x_i - x|$$

Where,

$n$  = the number of errors,  $\Sigma$  = summation symbol (which means “add them all up”),  $|x_i - x|$  = the absolute errors. The formula may look a little daunting, but the steps are easy:

Find all of your absolute errors,  $x_i - x$ . Add them all up. Divide by the number of errors. For example, if you had 10 measurements, divide by 10.

## Q18: What is the difference between long data and wide data?

There are many different ways that you can present the same dataset to the world. Let's take a look at one of the most important and fundamental distinctions, whether a dataset is wide or long.

The difference between wide and long datasets boils down to whether we prefer to have more columns in our dataset or more rows.

**Wide Data** A dataset that emphasizes putting additional data about a single subject in columns is called a wide dataset because, as we add more columns, the dataset becomes wider.

**Long Data** Similarly, a dataset that emphasizes including additional data about a subject in rows is called a long dataset because, as we add more rows, the dataset becomes longer. It's important to point out that there's nothing inherently good or bad about wide or long data.

In the world of data wrangling, we sometimes need to make a long dataset wider, and we sometimes need to make a wide dataset longer. However, it is true that, as a general rule, data scientists who embrace the concept of tidy data usually prefer longer datasets over wider ones.

## Q19: What are the data normalization method you have applied, and why?

**Normalization** is a technique often applied as part of data preparation for machine learning. The goal of normalization is to change the values of numeric columns in the dataset to a common scale, without distorting differences in the ranges of values. For machine learning, every dataset does not require normalization. It is required only when features have different ranges.

In simple words, when multiple attributes are there, but attributes have values on different scales, this may lead to poor data models while performing data mining operations. So they are normalized to bring all the attributes on the same scale, usually something between (0,1).

It is not always a good idea to normalize the data since we might lose information about maximum and minimum values. Sometimes it is a good idea to do so.

*For example*, ML algorithms such as Linear Regression or Support Vector Machines typically converge faster on normalized data. But on algorithms like K-means or K Nearest Neighbours, normalization could



be a good choice or a bad depending on the use case since the distance between the points plays a key role here.

person_name	Salary	Year_of_experience	Expected Position Level
Aman	100000	10	2
Abhinav	78000	7	4
Ashutosh	32000	5	8
Dishi	55000	6	7
Abhishek	92000	8	3
Avantika	120000	15	1
Ayushi	65750	7	5

The attributes salary and year\_of\_experience are on different scale and hence attribute salary can take high priority over attribute year\_of\_experience in the model.

## Types of Normalisation :

**1 Min-Max Normalization:** In most cases, standardization is used feature-wise

$$\hat{X}[:, i] = \frac{X[:, i] - \min(X[:, i])}{\max(X[:, i]) - \min(X[:, i])}$$

**2 Z-score normalization** In this technique, values are normalized based on a mean and standard deviation of the data

$$v' = \frac{v - \bar{A}}{\sigma_A}$$

$v'$ ,  $v$  is new and old of each entry in data respectively.  $\sigma_A$ ,  $\bar{A}$  is the standard deviation and mean of  $A$  respectively.

standardization (or Z-score normalization) is that the features will be rescaled so that they'll have the properties of a standard normal distribution with

$\mu=0$  and  $\sigma=1$  where  $\mu$  is the mean (average) and  $\sigma$  is the standard deviation from the mean; standard scores (also called z scores) of the samples are calculated as follows:

$$z = (x - \mu) / \sigma$$

## Q20: What is the difference between normalization and Standardization with example?

In ML, every practitioner knows that feature scaling is an important issue. The two most discussed scaling methods are **Normalization** and **Standardization**. Normalization typically means it rescales the values into a range of [0,1].

It is an alternative approach to Z-score normalization (or standardization) is the so-called Min-Max scaling (often also called “normalization” - a common cause for ambiguities). In this approach, the data is scaled to a fixed range - **usually 0 to 1**. Scikit-Learn provides a transformer called **MinMaxScaler** for this. A Min-Max scaling is typically done via the following equation:

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

**Example with sample data: Before Normalization:** Attribute Price in Dollars Storage Space Camera

- Attribute Price in Dollars Storage Space Camera
- Mobile 1 250 16 12
- Mobile 2 200 16 8
- Mobile 3 300 32 16
- Mobile 4 275 32 8
- Mobile 5 225 16 16

**After Normalization: (Values ranges from 0-1 which is working as expected)**

- Attribute Price in Dollars Storage Space Camera
- Mobile 1 0.5 0 0.5
- Mobile 2 0 0 0
- Mobile 3 1 1 1
- Mobile 4 0.75 1 0
- Mobile 5 0.25 0 1

Standardization (or Z-score normalization) typically means rescales data to have a mean of 0 and a standard deviation of 1 (unit variance) Formula: **Z or X<sub>new</sub> = (x - μ) / σ** where μ is the mean (average), and σ is the standard deviation from the mean; standard scores (also called z scores) Scikit-Learn provides a transformer called StandardScaler for standardization **Example:** Let's take an approximately normally distributed set of numbers: 1, 2, 2, 3, 3, 3, 4, 4, and 5. Its mean is 3, and its standard deviation: 1.22. Now, let's subtract the mean from all data points. we get a new data set of: -2, -1, -1, 0, 0, 0, 1, 1, and 2. Now, let's divide each data point by 1.22. As you can see in the picture below, we get: -1.6, -0.82, -0.82, 0, 0, 0, 0.82, 0.82, and 1.63

# **DATA SCIENCE INTERVIEW PREPARATION (30 Days of Interview Preparation)**

## **# DAY 04**

## Q1. What is upsampling and downsampling with examples?

The classification data set with skewed class proportions is called an imbalanced data set. Classes which make up a large proportion of the data sets are called majority classes. Those make up smaller proportions are minority classes.

Degree of imbalance Proportion of Minority Class

1>> Mild 20-40% of the data set

2>> Moderate 1-20% of the data set

3>> Extreme <1% of the data set

If we have an imbalanced data set, first try training on the true distribution. If the model works well and generalises, you are done! If not, try the following up sampling and down sampling technique.

### 1. Up-sampling

Upsampling is the process of randomly duplicating observations from the minority class to reinforce its signal.

First, we will import the resampling module from Scikit-Learn:

Module for resampling Python

1- From sklearn.utils import resample

Next, we will create a new Data Frame with an up-sampled minority class. Here are the steps:

1- First, we will separate observations from each class into different Data Frames.

2- Next, we will resample the minority class with replacement, setting the number of samples to match that of the majority class.

3- Finally, we'll combine the up-sampled minority class Data Frame with the original majority class Data Frame.

### 2-Down-sampling

Downsampling involves randomly removing observations from the majority class to prevent its signal from dominating the learning algorithm.

The process is similar to that of sampling. Here are the steps:

1-First, we will separate observations from each class into different Data Frames.

2-Next, we will resample the majority class without replacement, setting the number of samples to match that of the minority class.

3-Finally, we will combine the down-sampled majority class Data Frame with the original minority class Data Frame.

**Q2. What is the statistical test for data validation with an example,**

**Chi-square, ANOVA test, Z statics, T statics, F statics,**

**Hypothesis Testing?**

Before discussing the different statistical test, we need to get a clear understanding of what a null hypothesis is. A null hypothesis proposes that has no significant difference exists in the set of a given observation.

**Null:** Two samples mean are equal. Alternate: Two samples mean are not equal.

For rejecting the null hypothesis, a test is calculated. Then the test statistic is compared with a critical value, and if found to be greater than the critical value, the hypothesis will be rejected.

**Critical Value:-**

Critical values are the point beyond which we reject the null hypothesis. Critical value tells us, what is the probability of N number of samples, belonging to the same distribution. Higher, the critical value which means lower the probability of N number of samples belonging to the same distribution.

Critical values can be used to do hypothesis testing in the following way.

1. Calculate test statistic
2. Calculate critical values based on the significance level alpha
3. Compare test statistics with critical values.

IMP-If the test statistic is lower than the critical value, accept the hypothesis or else reject the hypothesis.

**Chi-Square Test:-**

A chi-square test is used if there is a relationship between two categorical variables.

Chi-Square test is used to determine whether there is a significant difference between the expected frequency and the observed frequency in one or more categories. Chi-square is also called the non-parametric test as it will not use any parameter

## **2-Anova test:-**

ANOVA, also called an analysis of variance, is used to compare multiples (three or more) samples with a single test.

Useful when there are more than three populations. Anova compares the variance within and between the groups of the population. If the variation is much larger than the within variation, the means of different samples will not be equal. If the between and within variations are approximately the same size, then there will be no significant difference between sample means. Assumptions of ANOVA: 1-All populations involved follow a normal distribution. 2-All populations have the same variance (or standard deviation). 3-The samples are randomly selected and independent of one another.

ANOVA uses the mean of the samples or the population to reject or support the null hypothesis. Hence it is called parametric testing.

## **3-Z Statics:-**

In a z-test, the samples are assumed to be normal distributed. A z score is calculated with population parameters as “population mean” and “population standard deviation” and it is used to validate a hypothesis that the sample drawn belongs to the same population.

The statistics used for this hypothesis testing is called z-statistic, the score for which is calculated as  $z = (x - \mu) / (\sigma / \sqrt{n})$ , where  $x$  = sample mean  $\mu$  = population mean  $\sigma / \sqrt{n}$  = population standard deviation If the test statistic is lower than the critical value, accept the hypothesis or else reject the hypothesis

## **4- T Statics:-**

A t-test used to compare the mean of the given samples. Like z-test, t-test also assumed a normal distribution of the samples. A t-test is used when the population parameters (mean and standard deviation) are unknown.

There are three versions of t-test

1. Independent samples t-test which compare means for two groups
2. Paired sample t-test which compares mean from the same group at different times
3. Sample t-test, which tests the mean of the single group against the known mean. The statistic for hypothesis testing is called t-statistic, the score for which is calculated as  $t = (x_1 - x_2) / (\sigma / \sqrt{n_1} + \sigma / \sqrt{n_2})$ , where

$x_1$  = It is mean of sample A,  $x_2$  = mean of sample B,

$n_1$  = size of sample 1  $n_2$  = size of sample 2

## 5- F Statics:-

The F-test is designed to test if the two population variances are equal. It compares the ratio of the two variances. Therefore, if the variances are equal, then the ratio of the variances will be 1.

The F-distribution is the ratio of two independent chi-square variables divided by their respective degrees of freedom.

$F = s_1^2 / s_2^2$  and where  $s_1^2 > s_2^2$ .

If the null hypothesis is true, then the F test-statistic given above can be simplified. This ratio of sample variances will be tested statistic used. If the null hypothesis is false, then we will reject the null hypothesis that the ratio was equal to 1 and our assumption that they were equal.

## Q3. What is the Central limit theorem?

### Central Limit Theorem

Definition: The theorem states that as the size of the sample increases, the distribution of the mean across multiple samples will approximate a Gaussian distribution (Normal). Generally, sample sizes equal to or greater than 30 are considered sufficient for the CLT to hold. It means that the distribution of the sample means is normally distributed. The average of the

sample means will be equal to the population mean. This is the key aspect of the theorem.

Assumptions:

1. The data must follow the randomization condition. It must be sampled randomly
2. Samples should be independent of each other. One sample should not influence the other samples
3. Sample size should be no more than 10% of the population when sampling is done without replacement
4. The sample size should be sufficiently large. The mean of the sample means is denoted as:

$$\mu_{\bar{X}} = \mu$$

Where,

$\mu_{\bar{X}}$  = Mean of the sample means  $\mu$  = Population mean and, the standard deviation of the sample mean is denoted as:

$$\sigma_{\bar{X}} = \sigma / \sqrt{n}$$

Where,

$\sigma_{\bar{X}}$  = Standard deviation of the sample mean  $\sigma$  = Population standard deviation  $n$  = sample size

A sufficiently large sample size can predict the characteristics of a population accurately. For Example, we shall take a uniformly distributed data:

Randomly distributed data: Even for a randomly (Exponential) distributed data the plot of the means is normally distributed.

The advantage of CLT is that we need not worry about the actual data since the means of it will always be normally distributed. With this, we can create component intervals, perform T-tests and ANOVA tests from the given samples.

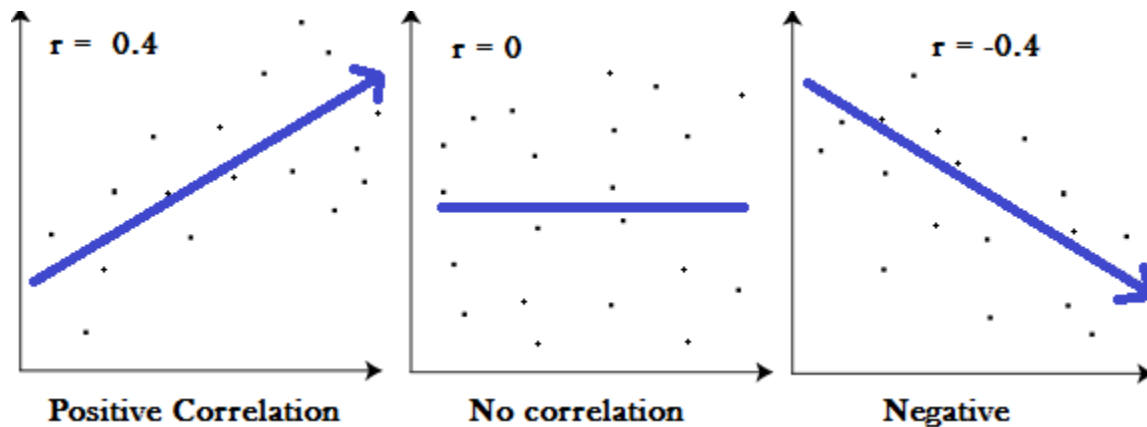
#### **Q4. What is the correlation and coefficient?**

##### **What is the Correlation Coefficient?**

The correlation coefficient is a statistical measure that calculates the strength of the relationship between the relative movements of two



variables. We use it to measure both the strength and direction of a linear relationship between two variables the values range between -1.0 and 1.0. A calculated number greater than 1.0 or less than -1.0 means that there was an error in the correlation measurement. A correlation of -1.0 shows a perfect negative correlation, while a correlation of 1.0 shows a perfect positive correlation.



Correlation coefficient formulas are used to find how strong a relationship is between data. The formulas return a value between -1 and 1, where:

1 indicates a strong positive relationship. -1 indicates a strong negative relationship. A result of zero indicates no relationship at all.

### Meaning

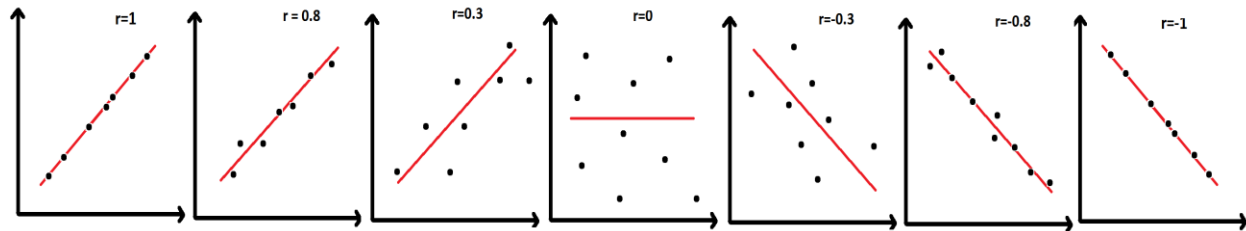
1. A correlation coefficient of 1 means that for every positive increase in one variable, there is a positive increase in a fixed proportion in the other. For example, shoe sizes go up in (almost) perfect correlation with foot length.
2. A correlation coefficient of -1 means that for every positive increase in one variable, there is a negative decrease of a fixed proportion in the other. For example, the amount of gas in a tank decreases in (almost) perfect correlation with speed.
3. Zero means that for every increase, there isn't a positive or negative increase. The two just aren't related.

## What is a Negative Correlation?

Negative correlation is a relationship between two variables in which one variable increases as the other decreases, and vice versa. In statistics, a perfect negative correlation is represented by the value  $-1$ . Negative correlation or inverse correlation is a relationship between two variables whereby they move in opposite directions. If variables  $X$  and  $Y$  have a negative correlation (or are negatively correlated), as  $X$  increases in value,  $Y$  will decrease; similarly, if  $X$  decreases in value,  $Y$  will increase.

## What Is Positive Correlation?

Positive correlation is a relationship between two variables in which both variables move in tandem—that is, in the same direction. A positive correlation exists when one variable decreases as the other variable decreases or one variable increases while the other increases.



We use the correlation coefficient to measure the strength and direction of the linear relationship between two numerical variables  $X$  and  $Y$ . The correlation coefficient for a sample of data is denoted by  $r$ .

## Pearson Correlation Coefficient

Pearson is the most widely used correlation coefficient. Pearson correlation measures the linear association between continuous variables. In other words, this coefficient quantifies the degree to which a relationship between two variables can be described by a line. Formula developed by Karl Pearson over 120 years ago is still the most widely used today. The formula for the correlation ( $r$ ) is

#### Correlation Coefficient Formula

$$r = \frac{n(\sum xy) - (\sum x)(\sum y)}{\sqrt{[n\sum x^2 - (\sum x)^2][n\sum y^2 - (\sum y)^2]}}$$

Where n is the number of pairs of data;

$\bar{x}$  and  $\bar{y}$  are the sample means of all the x-values and all the y-values, respectively; and  $s_x$  and  $s_y$  are the sample standard deviations of all the x- and y-values, respectively.

1. Find the mean of all the x-values and mean of all y-values.
2. Find the standard deviation of all the x-values (call it  $s_x$ ) and the standard deviation of all the y-values (call it  $s_y$ ). For example, to find  $s_x$ , you would use the following equation:
3. For each of the n pairs (x, y) in the data set, take
4. Add up the n results from Step 3.
5. Divide the sum by  $s_x * s_y$ .
6. Divide the result by  $n - 1$ , where n is the number of (x, y) pairs. (It's the same as multiplying by 1 over  $n - 1$ .) This gives you the correlation, r.

### Q5: What is the difference between machine learning and deep learning?

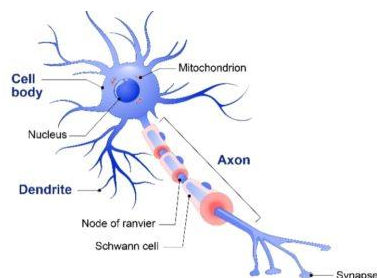
#### Machine Learning | deep learning

Machine Learning is a technique to learn from that data and then apply what has been learnt to make an informed decision | The main difference between deep and machine learning is, machine learning models become better progressively but the model still needs some guidance. If a machine-learning model returns an inaccurate prediction then the programmer needs to fix that problem explicitly but in the case of deep learning, the model does it by himself.

- >Machine Learning can perform well with small size data also | Deep Learning does not perform as good with smaller datasets.
- >Machine learning can work on some low-end machines also | Deep Learning involves many matrix multiplication operations which are better suited for GPUs
- >Features need to be identified and extracted as per the domain before pushing them to the algorithm | Deep learning algorithms try to learn high-level features from data.
- >It is generally recommended to break the problem into smaller chunks, solve them and then combine the results | It generally focusses on solving the problem end to end
- >Training time is comparatively less | Training time is comparatively more
- >Results are more interpretable | Results Maybe more accurate but less interpretable
- > No use of Neural networks | uses neural networks
- > Solves comparatively less complex problems | Solves more complex problems.

#### **Q6: What is perceptron and how it is related to human neurons?**

If we focus on the structure of a biological neuron, it has dendrites, which are used to receive inputs. These inputs are summed in the cell body and using the Axon it is passed on to the next biological neuron as shown below.

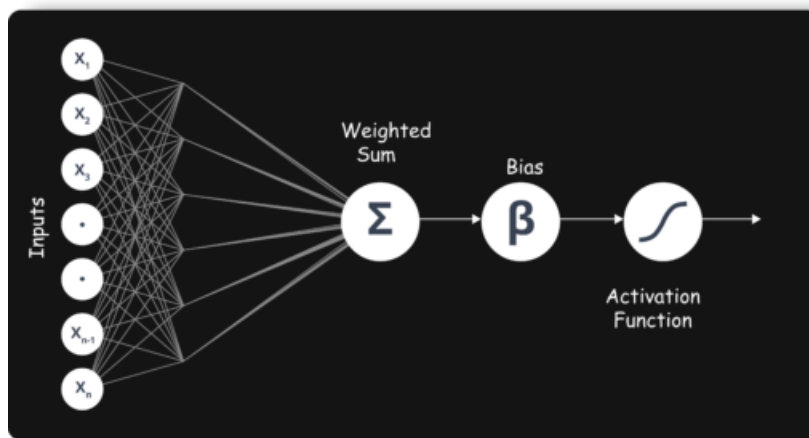


**Dendrite:** Receives signals from other neurons

**Cell Body:** Sums all the inputs

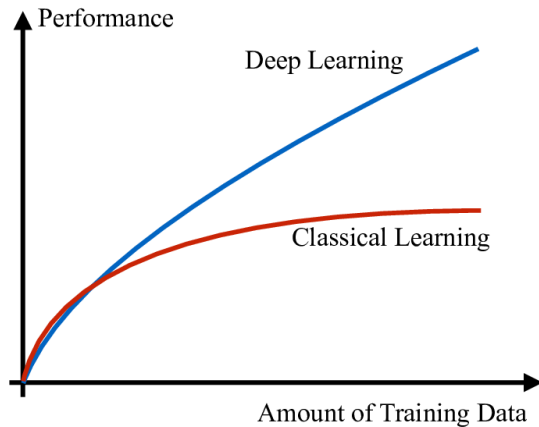
**Axon:** It is used to transmit signals to the other cells

Similarly, a perceptron receives multiple inputs, applies various transformations and functions and provides an output. A Perceptron is a linear model used for binary classification. It models a neuron, which has a set of inputs, each of which is given a specific weight. The neuron computes some function on these weighted inputs and gives the output.



### Q7: Why deep learning is better than machine learning?

Though traditional ML algorithms solve a lot of our cases, they are not useful while working with high dimensional data that is where we have a large number of inputs and outputs. For example, in the case of handwriting recognition, we have a large amount of input where we will have different types of inputs associated with different types of handwriting.



The second major challenge is to tell the computer what are the features it should look for that will play an important role in predicting the outcome as well as to achieve better accuracy while doing so.

#### **Q8: What kind of problem can be solved by using deep learning?**

Deep Learning is a branch of Machine Learning, which is used to solve problems in a way that mimics the human way of solving problems.

Examples:

- Image recognition
- Object Detection
- Natural Language processing- Translation, Sentence formations, text to speech, speech to text
- understand the semantics of actions

#### **Q9: List down all the activation function using mathematical**

**Expression and example. What is the activation function?**

**Activation functions** are very important for an Artificial Neural Network to learn and make sense of something complicated and the Non-linear complex functional mappings between the inputs and response variable. They introduce non-linear properties to our Network. Their main purposes are to convert an input signal of a node in an A-NN to an output signal.

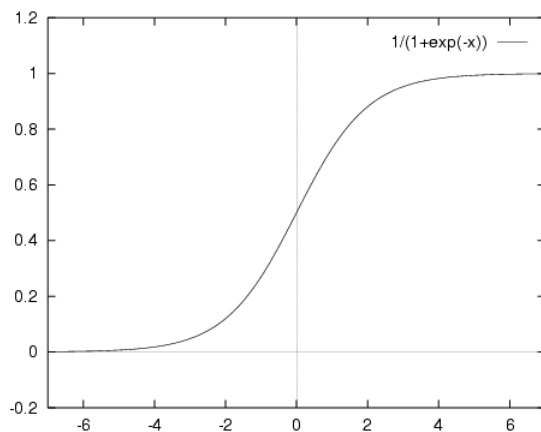
*So why do we need Non-Linearities?*

Non-linear functions are those, which have a degree more than one, and they have a curvature when we plot a Non-Linear function. Now we need a Neural Network Model to learn and represent almost anything and any arbitrary complex function, which maps inputs to outputs. Neural-Networks are considered Universal Function Approximations. It means that they can compute and learn any function at all.

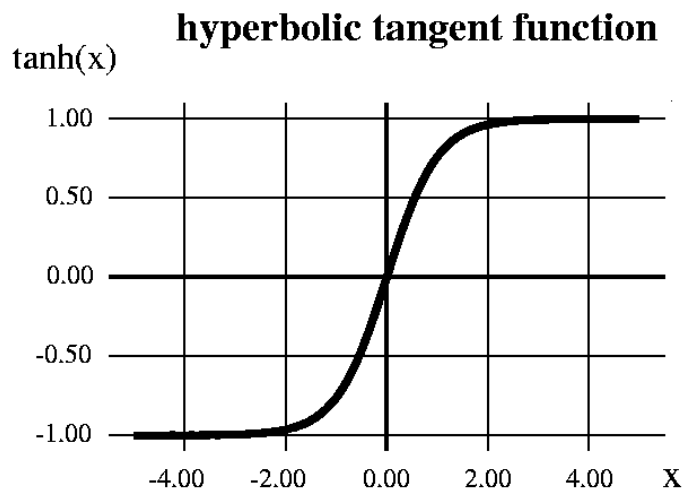
*Most popular types of Activation functions -*

- Sigmoid or Logistic
- Tanh — Hyperbolic tangent
- ReLu -Rectified linear units

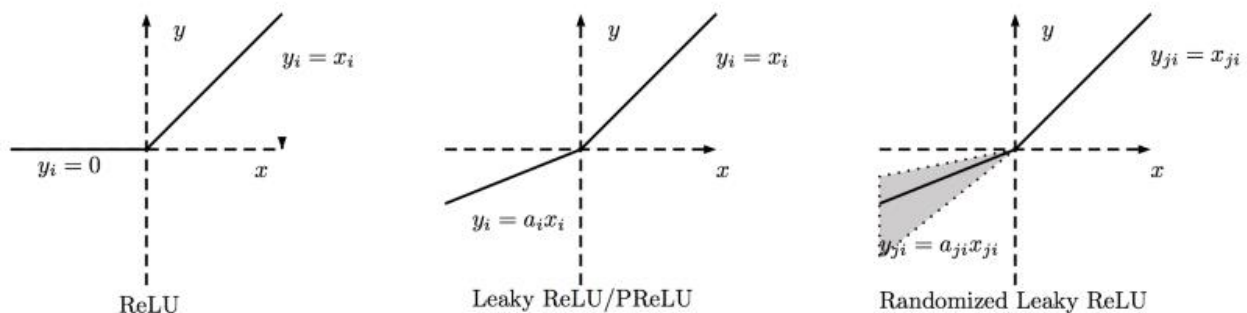
**Sigmoid Activation function:** It is a activation function of form  $f(x) = 1 / 1 + \exp(-x)$  . Its Range is between 0 and 1. It is an S-shaped curve. It is easy to understand.



**Hyperbolic Tangent function- Tanh :** It's mathematical formula is  $f(x) = 1 - \exp(-2x) / 1 + \exp(-2x)$ . Now it's the output is zero centred because its range in between -1 to 1 i.e.  $-1 < \text{output} < 1$  . Hence optimisation is easier in this method; Hence in practice, it is always preferred over Sigmoid function.



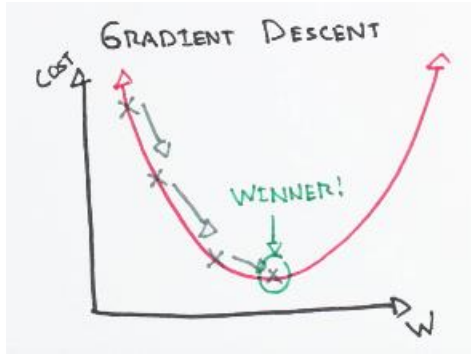
**ReLU- Rectified Linear units:** It has become more popular in the past couple of years. It was recently proved that it has six times improvement in convergence from Tanh function. It's  $R(x) = \max(0, x)$  i.e. if  $x < 0$ ,  $R(x) = 0$  and if  $x \geq 0$ ,  $R(x) = x$ . Hence as seen that mathematical form of this function, we can see that it is very simple and efficient. Many times in Machine learning and computer science we notice that most simple and consistent techniques and methods are only preferred and are the best. Hence, it avoids and rectifies the vanishing gradient problem. Almost all the deep learning Models use ReLU nowadays.



**Q10: Detail explanation about gradient decent using example and Mathematical expression?**



Gradient descent is an optimisation algorithm used to minimize some function by iteratively moving in the direction of steepest descent as defined by negative of the gradient. In machine learning, we used gradient descent to update the parameters of our model. Parameters refer to coefficients in the Linear Regression and weights in neural networks.



The size of these steps called the learning rate. With the high learning rate, we can cover more ground each step, but we risk overshooting the lower point since the slope of the hill is constantly changing. With a very lower learning rate, we can confidently move in the direction of the negative gradient because we are recalculating it so frequently. The Lower learning rate is more precise, but calculating the gradient is time-consuming, so it will take a very large time to get to the bottom.

## Math

Now let's run gradient descent using new cost function. There are two parameters in cost function we can control:  $m$  (weight) and  $b$  (bias). Since we need to consider that the impact each one has on the final prediction, we need to use partial derivatives. We calculate the partial derivative of the cost function concerning each parameter and store the results in a gradient.

## Math

Given the cost function:

$$f(m, b) = \frac{1}{N} \sum_{i=1}^n (y_i - (mx_i + b))^2$$

The gradient can be calculated as:

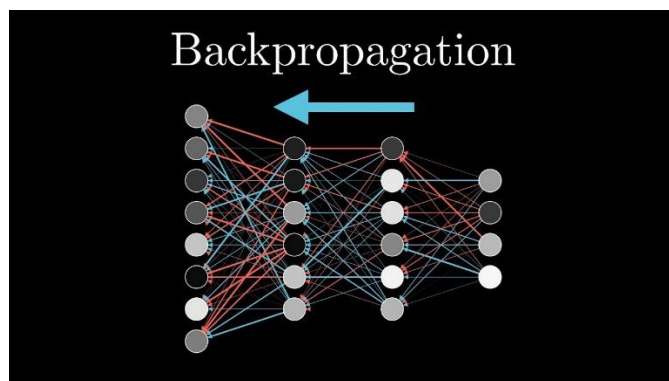
$$f'(m, b) = \begin{bmatrix} \frac{df}{dm} \\ \frac{df}{db} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum -2x_i(y_i - (mx_i + b)) \\ \frac{1}{N} \sum -2(y_i - (mx_i + b)) \end{bmatrix}$$

To solve for the gradient, we iterate by our data points using our new  $m$  and  $b$  values and compute the partial derivatives. This new gradient tells us about the slope of the cost function at our current position (current parameter values) and the directions we should move to update our parameters. The learning rate controls the size of our update.

### Q11: What is backward propagation?

**Back-propagation** is the essence of the neural net training and this method of fine-tuning the weights of a neural net based on the errors rate obtained in the previous epoch. Proper tuning of the weights allows us to reduce error rates and to make the model reliable by increasing its generalisation.

Backpropagation is a short form of "backward propagation of errors." This is the standard method of training artificial neural networks. This helps to calculate the gradient of a loss function with respects to all the weights in the network.



Most prominent advantages of Backpropagation are:

- Backpropagation is the fast, simple and easy to program.
- It has no parameters to tune apart from the numbers of input.
- It is the flexible method as it does not require prior knowledge about the network
- It is the standard method that generally works well.
- It does not need any special mentions of the features of the function to be learned.

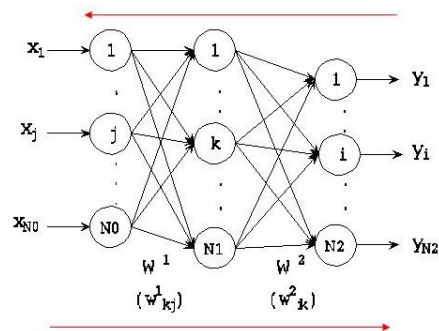
## The BackPropagation Algorithm

Main idea:

For each example in the training set:

- compute the output signal
- compute the error corresponding to the output level
- propagate the error back into the network and store the corresponding delta values for each layer
- adjust each weight by using the error signal and input signal for each layer

Computation of the error signal (BACKWARD)



Computation of the output signal (FORWARD)

### Q12: How we assign weights in deep learning?

We already know that in a neural network, weights are usually initialised randomly and that kind of initialisation takes a fair/significant amount of repetitions to converge to the least loss and reach the ideal weight matrix. The problem is, that kind of initialisation is prone to vanishing or exploding gradient problems.

*General ways to make it initialise better weights:*

ReLU activation function in the deep nets.

1. Generate a random sample of weights from a Gaussian distribution having mean 0 and a standard deviation of 1.
2. Multiply the sample with the square root of  $(2/n_i)$ . Where  $n_i$  is the number of input units for that layer.

b) Likewise, if you're using Tanh activation function :

1. Generate a random sample of weights from a Gaussian distribution having mean 0 and a standard deviation of 1.
2. Multiply the sample with the square root of  $(1/n_i)$  where  $n_i$  is several input units for that layer.

### **Q13: What is optimiser is deep learning, and which one is the best?**

Deep learning is an iterative process. With so many hyperparameters to tune or methods to try, it is important to be able to train models fast, to quickly complete the iterative cycle. This is the key to increase the speed and efficiency of a machine learning team.

Hence the importance of optimisation algorithms such as stochastic gradient descent, min-batch gradient descent, gradient descent with momentum and the Adam optimiser.

Adam optimiser is the best one.

Given an algorithm  $f(x)$ , it helps in either minimisation or maximisation of the value of  $f(x)$ . In this context of deep learning, we use optimisation algorithms to train the neural network by optimising the cost function  $J$ .

The cost function is defined as:

$$J(W, b) = \sum_{i=1}^m L(y'^i, y^i)$$

The value of the cost function  $J$  is the mean of the loss  $L$  between the predicted value  $y'$  and actual value  $y$ . The value  $y'$  is obtained during the forward propagation step and makes use of the Weights  $W$  and biases  $b$  of the network. With the help of optimisation algorithms, we minimise the value of Cost Function  $J$  by updating the values of trainable parameters  $W$  and  $b$ .

#### **Q14: What is gradient descent, mini-batch gradient descent, batch gradient decent, stochastic gradient decent and adam?**

##### **Gradient Descent**

it is an iterative machine learning optimisation algorithm to reduce the cost function, and help models to make accurate predictions.

Gradient indicates the direction of increase. As we want to find the minimum points in the valley, we need to go in the opposite direction of the gradient. We update the parameters in the negative gradient direction to minimise the loss.

$$\theta = \theta - \eta \nabla J(\theta; x, y)$$

Where  $\theta$  is the weight parameter,  $\eta$  is the learning rate, and  $\nabla J(\theta; x, y)$  is the gradient of weight parameter  $\theta$

##### **Types of Gradient Descent**

Different types of Gradient descents are

- Batch Gradient Descent or Vanilla Gradient Descent
- Stochastic Gradient Descent

- Mini batch Gradient Descent

### Batch Gradient Descent

In the batch gradient, we use the entire dataset to compute the gradient of the cost function for each iteration for gradient descent and then update the weights.

### Stochastic Gradient descent

Stochastic gradient descent, we use a single data point or example to calculate the gradient and update the weights with every iteration.

We first need to shuffle the datasets so that we get a completely randomised dataset. As the datasets are random and weights, are updated for every single example, an update of the weights and the cost functions will be noisy jumping all over the place

### Mini Batch Gradient descent

Mini-batch gradients is a variation of stochastic gradient descent where instead of a single training example, a mini-batch of samples are used.

Mini -batch gradient descent is widely used and converges faster and is more stable.

The batch size can vary depending upon the dataset.

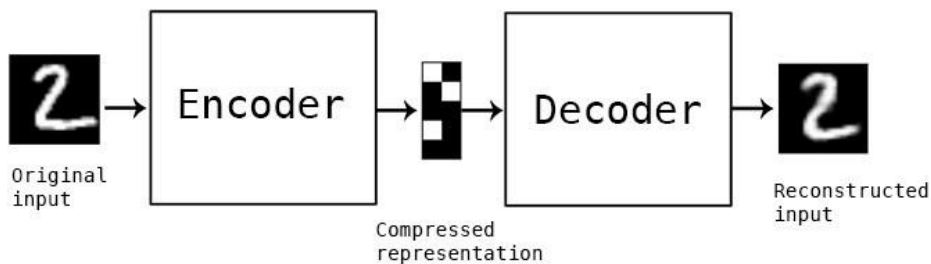
As we take batches with different samples, it reduces the noise which is a variance of the weights updates, and that helps to have a more stable converge faster.

### Q15: What are autoencoders?

An **autoencoder**, neural networks that have three layers:

An input layer, a hidden layer which is also known as encoding layer, and a decoding layer. This network is trained to reconstruct its inputs, which forces the hidden layer to try to learn good representations of the inputs.

An autoencoder neural network is an unsupervised Machine-learning algorithm that applies backpropagation, setting the target values to be equal to the inputs. An autoencoder is trained to attempts to copy its input to its output. Internally, it has a hidden layer which describes a code used to represent the input.



### **Autoencoder Components:**

Autoencoders consists of 4 main parts:

- 1- Encoder: In this, the model learns how to reduce the input dimensions and compress the input data into an encoded representation.
- 2- Bottleneck: In this, the layer that contains the compressed representation of the input data. This is the lowest possible dimension of the input data.
- 3- Decoder: In this, the model learns how to reconstruct the data from the encoded representation to be as close to the original inputs as possible.
- 4- Reconstruction Loss: In this method that measures how well the decoder is performing and how close the output is related to the original input.

### **Types of Autoencoders :**

1. Denoising auto encoder
2. Sparse auto encoder
3. Variational auto encoder (VAE)
4. Contractive auto encoder (CAE)

## Q16: What is CNN?

This is the simple application of a filter to an input that results in inactivation. Repeated application of the same filter to input results in a map of activations called a feature map, indicating the locations and strength of a detected feature in input, such as an image.

Convolutional layers are the major building blocks which are used in convolutional neural networks.

A convnets is the sequence of layers, and every layer transforms one volume to another through differentiable functions.

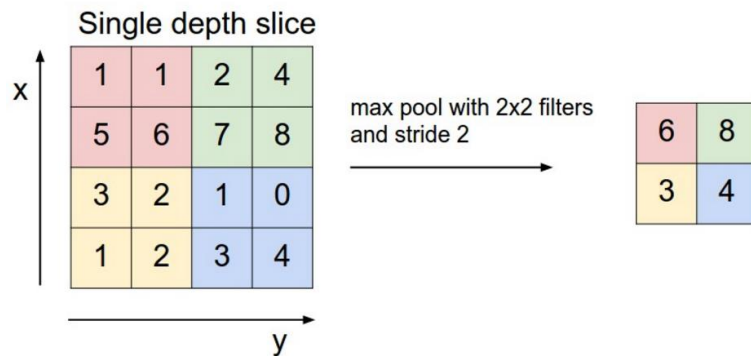
Different types of layers in CNN:

Let's take an example by running a convnets on of image of dimensions **32 x 32 x 3**.

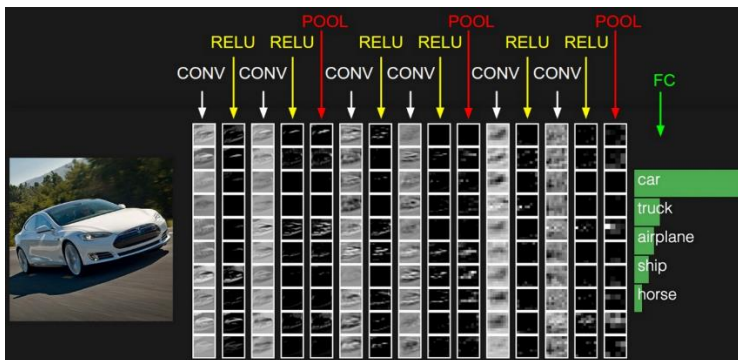
1. **Input Layer:** It holds the raw input of image with width 32, height 32 and depth 3.
2. **Convolution Layer:** It computes the output volume by computing dot products between all filters and image patches. Suppose we use a total of 12 filters for this layer we'll get output volume of dimension 32 x 32 x 12.
3. **Activation Function Layer:** This layer will apply the element-wise activation function to the output of the convolution layer. Some activation functions are RELU:  $\max(0, x)$ , Sigmoid:  $1/(1+e^{-x})$ , Tanh, Leaky RELU, etc. So the volume remains unchanged. Hence output volume will have dimensions 32 x 32 x 12.
4. **Pool Layer:** This layer is periodically inserted within the convnets, and its main function is to reduce the size of volume which makes the



computation fast reduces memory and also prevents overfitting. Two common types of pooling layers are max pooling and average pooling. If we use a max pool with  $2 \times 2$  filters and stride 2, the resultant volume will be of dimension  $16 \times 16 \times 12$ .



5. Fully-Connected Layer: This layer is a regular neural network layer that takes input from the previous layer and computes the class scores and outputs the 1-D array of size equal to the number of classes.

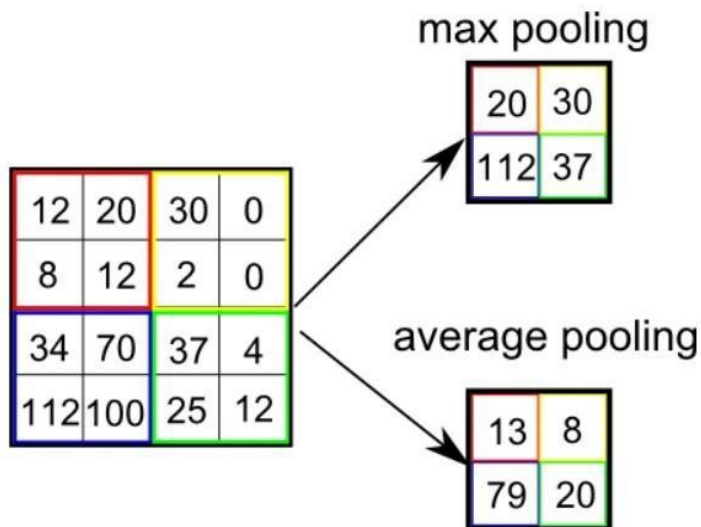


## Q17: What is pooling, padding, filtering operations on CNN?

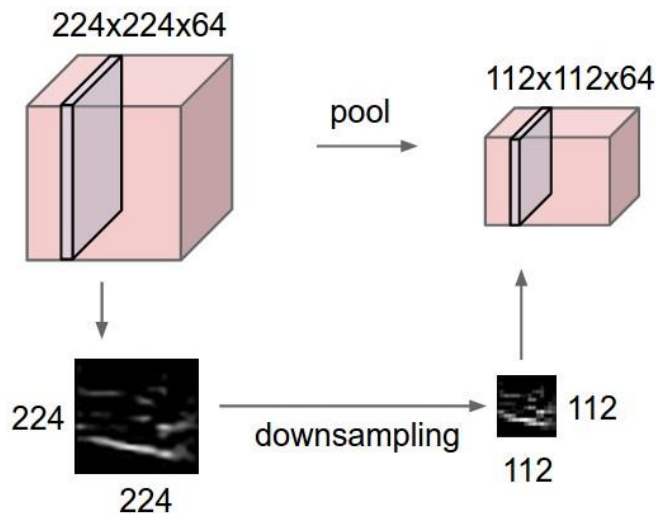
### Pooling Layer

It is commonly used to periodically insert a Pooling layer in-between successive Conv layers in a ConvNet architecture. Its function is to progressively reduce the spatial size of the representation to reduce the number of parameters and computation in the network, and hence to also

control overfitting. The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, using the MAX operation.

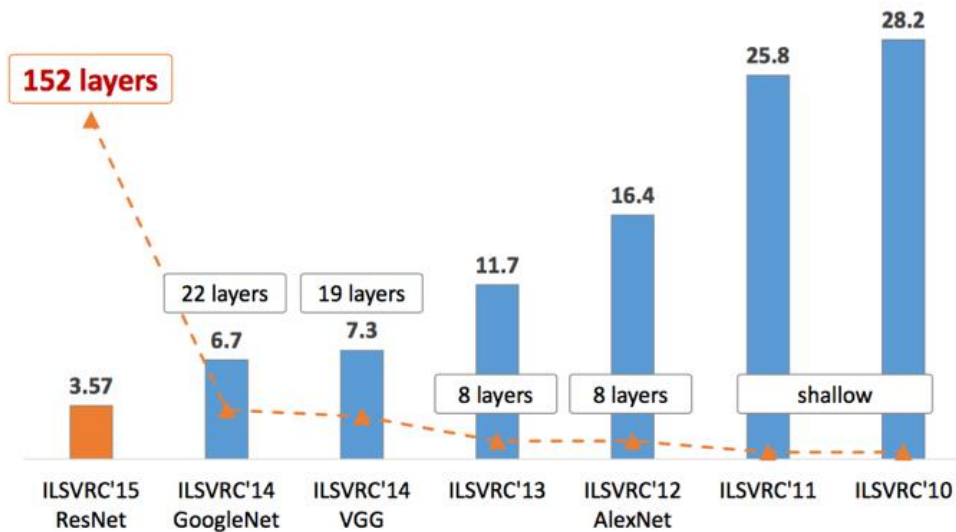


The most common form is a pooling layer with filters of size 2x2 applied with a stride of 2 downsamples every depth slice in the input by two along both width and height, discarding 75% of the activations. Every MAX operation would, in this case, be taking a max over four numbers (little 2x2 region in some depth slice). The depth dimension remains unchanged.



**Q18: What is the Evolution technique of CNN?**

It all started with LeNet in 1998 and eventually, after nearly 15 years, led to groundbreaking models winning the ImageNet Large Scale Visual Recognition Challenge which includes AlexNet in 2012 to GoogleNet in 2014 to ResNet in 2015 to an ensemble of previous models in 2016. In the last two years, no significant progress has been made, and the new models are an ensemble of previous groundbreaking models.



## LeNet in 1998

LeNet is a 7-level convolutional network by LeCun in 1998 that classifies digits and used by several banks to recognise the hand-written numbers on cheques digitised in 32x32 pixel greyscale input images.

## AlexNet in 2012

AlexNet: It is considered to be the first paper/ model, which rose the interest in CNNs when it won the ImageNet challenge in the year 2012. It is a deep CNN trained on ImageNet and outperformed all the entries that year.

## VGG in 2014

VGG was submitted in the year 2013, and it became a runner up in the ImageNet contest in 2014. It is widely used as a simple architecture compared to AlexNet.

## GoogleNet in 2014

In 2014, several great models were developed like VGG, but the winner of the ImageNet contest was GoogleNet.

GoogLeNet proposed a module called the inception modules that includes skipping connections in the network, forming a mini-module, and this module is repeated throughout the network.

## ResNet in 2015

There are 152 layers in the Microsoft ResNet. The authors showed empirically that if you keep on adding layers, the error rate should keep on decreasing in contrast to “plain nets” we're adding a few layers resulted in higher training and test errors.

### **Q19: How to initialise biases in deep learning?**

It is possible and common to initialise the biases to be zero since the random numbers in the weights provide the asymmetry braking. For ReLU non-linearities, some people like to use small constant value such as 0.01 for all biases because this ensures that all ReLU units fire in the beginning, therefore obtain, and propagate some gradient. However, it is unclear if this provides a consistent improvement (in fact some results seem to indicate that this performs worst) and it is more commonly used to use 0 bias initialisation.

### **Q20: What is learning Rate?**

#### **Learning Rate**

The learning rate controls how much we should adjust the weights concerning the loss gradient. Learning rates are randomly initialised.

Lower the values of the learning rate slower will be the convergence to global minima.

Higher values for the learning rate will not allow the gradient descent to converge

Since our goal is to minimise the function cost to find the optimised value for weights, we run multiples iteration with different weights and calculate the cost to arrive at a minimum cost

-----  
-----