

# UNIT 1:

## 1. What is hypothesis testing ?

Hypothesis testing is a statistical method that is used in making statistical decisions using experimental data. Hypothesis Testing is basically an assumption that we make about the population parameter.

Ex : you say an average student in class is 40 or a boy is taller than girls.

All those examples we assume need some statistical way to prove those. we need some mathematical conclusion whatever we are assuming is true.

## 2. Why do we use it ?

Hypothesis testing is an essential procedure in statistics. A hypothesis test evaluates two mutually exclusive statements about a population to determine which statement is best supported by the sample data. When we say that a finding is statistically significant, it's thanks to a hypothesis test.

**Type I error:** When we reject the null hypothesis, although that hypothesis was true. Type I error is denoted by  $\alpha$ . In hypothesis testing, the normal curve that shows the critical region is called the alpha region

**Type II errors:** When we accept the null hypothesis but it is false. Type II errors are denoted by  $\beta$ . In Hypothesis testing, the normal curve that shows the acceptance region is called the beta region.

## 3. What is a decision tree in machine learning?

Decision trees are a way of modeling decisions and outcomes, mapping decisions in a branching structure. Decision trees are used to calculate the potential success of different series of decisions made to achieve a specific goal.

Decision tree are made up of different nodes. The root node is the start of the decision tree, which is usually the whole dataset within machine learning. Leaf nodes are the endpoint of a branch, or the final output of a series of decisions.

#### **4. What is Overfitting?**

- Overfitting & underfitting are the two main errors/problems in the machine learning model, which cause poor performance in Machine Learning.
- Overfitting occurs when the model fits more data than required, and it tries to capture each and every datapoint fed to it. Hence it starts capturing noise and inaccurate data from the dataset, which degrades the performance of the model.
- An overfitted model doesn't perform accurately with the test/unseen dataset and can't generalize well.
- An overfitted model is said to have low bias and high variance.

## **UNIT 2:**

### **1. What is Gibbs sampling in machine learning?**

Gibbs sampling is a Markov Chain Monte Carlo (MCMC) algorithm where each random variable is iteratively resampled from its conditional distribution given the remaining variables. It's a simple and often highly effective approach for performing posterior inference in probabilistic models.

## 2. What is the EM Algorithm in Machine Learning?

The EM algorithm can be used to determine the local maximum likelihood (MLE) parameters or maximum a posteriori (MAP) parameters for latent variables (unobservable variables that need to be inferred from observable variables) in a statistical model. It is used to predict these values or determine data that is missing or incomplete, provided that you know the general form of probability distribution associated with these latent variables.

## 3. Logistic Regression in Machine Learning?

Logistic regression predicts the output of a categorical dependent variable. Therefore the outcome must be a categorical or discrete value. It can be either Yes or No, 0 or 1, true or False, etc. but instead of giving the exact value as 0 and 1, **it gives the probabilistic values which lie between 0 and 1.**

## 4. Support Vector Machine Algorithm

Support Vector Machine or SVM is one of the most popular Supervised Learning algorithms, which is used for Classification as well as Regression problems. However, primarily, it is used for Classification problems in Machine Learning.

The goal of the SVM algorithm is to create the best line or decision boundary that can segregate n-dimensional space into classes so that we can easily put the new data point in the correct category in the future. This best decision boundary is called a hyperplane.

Types of SVM

**SVM can be of two types:**

- **Linear SVM:** Linear SVM is used for linearly separable data, which means if a dataset can be classified into two classes by using a single straight line, then

such data is termed as linearly separable data, and classifier is used called as Linear SVM classifier.

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

## UNIT 3:

### 1. Clustering in Machine Learning

A way of grouping the data points into different clusters, consisting of similar data points. The objects with the possible similarities remain in a group that has less or no similarities with another group."

It does it by finding some similar patterns in the unlabelled dataset such as shape, size, color, behavior, etc., and divides them as per the presence and absence of those similar patterns.

It is an unsupervised learning method, hence no supervision is provided to the algorithm, and it deals with the unlabeled dataset.

#### Connectivity-Based Clustering (Hierarchical Clustering)

Hierarchical Clustering is a method of unsupervised machine learning clustering where it begins with a pre-defined top to bottom hierarchy of clusters. It then proceeds to perform a decomposition of the data objects based on this hierarchy, hence obtaining the clusters.

#### k-Means Clustering

K-means clustering uses "centroids", K different randomly-initiated points in the data, and assigns every data point to the nearest centroid. After every point has been assigned, the centroid is moved to the average of all of the points assigned to it.

## **Supervised Machine Learning**

Supervised learning is a process of providing input data as well as correct output data to the machine learning model. The aim of a supervised learning algorithm is to find a mapping function to map the input variable( $x$ ) with the output variable( $y$ ).

# **UNIT 4:**

## **1.What is an Artificial Neural Network?**

An Artificial neural network is usually a computational network based on biological neural networks that construct the structure of the human brain. Similar to how a human brain has neurons interconnected to each other, artificial neural networks also have neurons that are linked to each other in various layers of the networks. These neurons are known as nodes.

## **2.What is Backpropagation?**

Backpropagation is the essence of neural network training. It is the method of fine-tuning the weights of a neural network based on the error rate obtained in the previous epoch (i.e., iteration). Proper tuning of the weights allows you to reduce error rates and make the model reliable by increasing its generalization.

## **3.DEEP NEURAL NETWORK**

The main purpose of a neural network is to receive a set of inputs, perform progressively complex calculations on them, and give output to solve real world problems like classification. We restrict ourselves to feed forward neural networks.

#### **4.CNN**

A convolutional neural network is a specific kind of neural network with multiple layers. It processes data that has a grid-like arrangement then extracts important features. One huge advantage of using CNNs is that you don't need to do a lot of pre-processing on images.

#### **5. What Is a Recurrent Neural Network (RNN)?**

RNN works on the principle of saving the output of a particular layer and feeding this back to the input in order to predict the output of the layer.

## **UNIT 5:**

### **1. RANDOM FOREST**

Random Forest is a classifier that contains a number of decision trees on various subsets of the given dataset and takes the average to improve the predictive accuracy of that dataset." Instead of relying on one decision tree, the random forest takes the prediction from each tree and based on the majority votes of predictions, and it predicts the final output.

The greater number of trees in the forest leads to higher accuracy and prevents the problem of overfitting.

## **2 What Is Bagging in Machine Learning?**

Bagging, also known as Bootstrap aggregating, is an ensemble learning technique that helps to improve the performance and accuracy of machine learning algorithms. It is used to deal with bias-variance trade-offs and reduces the variance of a prediction model.

## **3.BOOTSTING**

Boosting is an ensemble modeling technique that attempts to build a strong classifier from the number of weak classifiers. It is done by building a model by using weak models in series. Firstly, a model is built from the training data. Then the second model is built which tries to correct the errors present in the first model. This procedure is continued and models are added until either the complete training data set is predicted correctly or the maximum number of models are added

