

Practical No. 01 01

Aim: - Perform Data preprocessing On given Dataset

Apparatus

pandas | numpy | Scipy, matplotlib, python machine learning.

Theory :-

Data preprocessing is process of preparing the raw data and making it suitable for a machine learning model. It is first and crucial step while creating a machine learning model.

When creating a machine learning project it is not always a case that we come across the clean, fit and formatted data. And while doing any operation with data, it is mandatory to clean it and put in a formatted way. So, for this, we use data preprocessing task.

Behavior of data :-

A real-world data generally contains noises, missing values, and maybe in an unusable format which cannot be directly used for machine learning. Data processing helps to increases the accuracy & efficiency of a machine learning model.

Boxplot

It is a type of chart that depicts a group of numerical data through their quartiles. It is a simple way to visualize the shape of the data.

Boxplot gives a five-number summary of a set of data:

- minimum - It is the minimum value in the dataset excluding outliers.

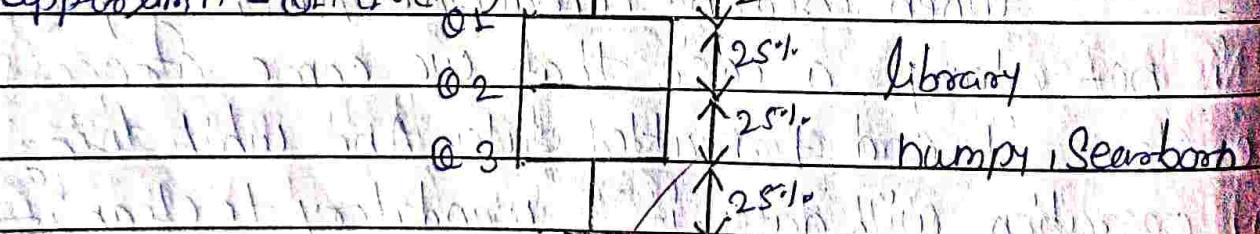
First Quartile (Q1) - 25% of data lies below it.

Median (Q2) - It is mid point of dataset.

Third Quartile (Q3) - 75% of the data lies below it.

- maximum - It is the maximum value in the dataset excluding outliers.

$$\text{upper limit} = Q1 + 1.5(IQR)$$



$$\text{lower limit} = Q1 - 1.5(IQR)$$

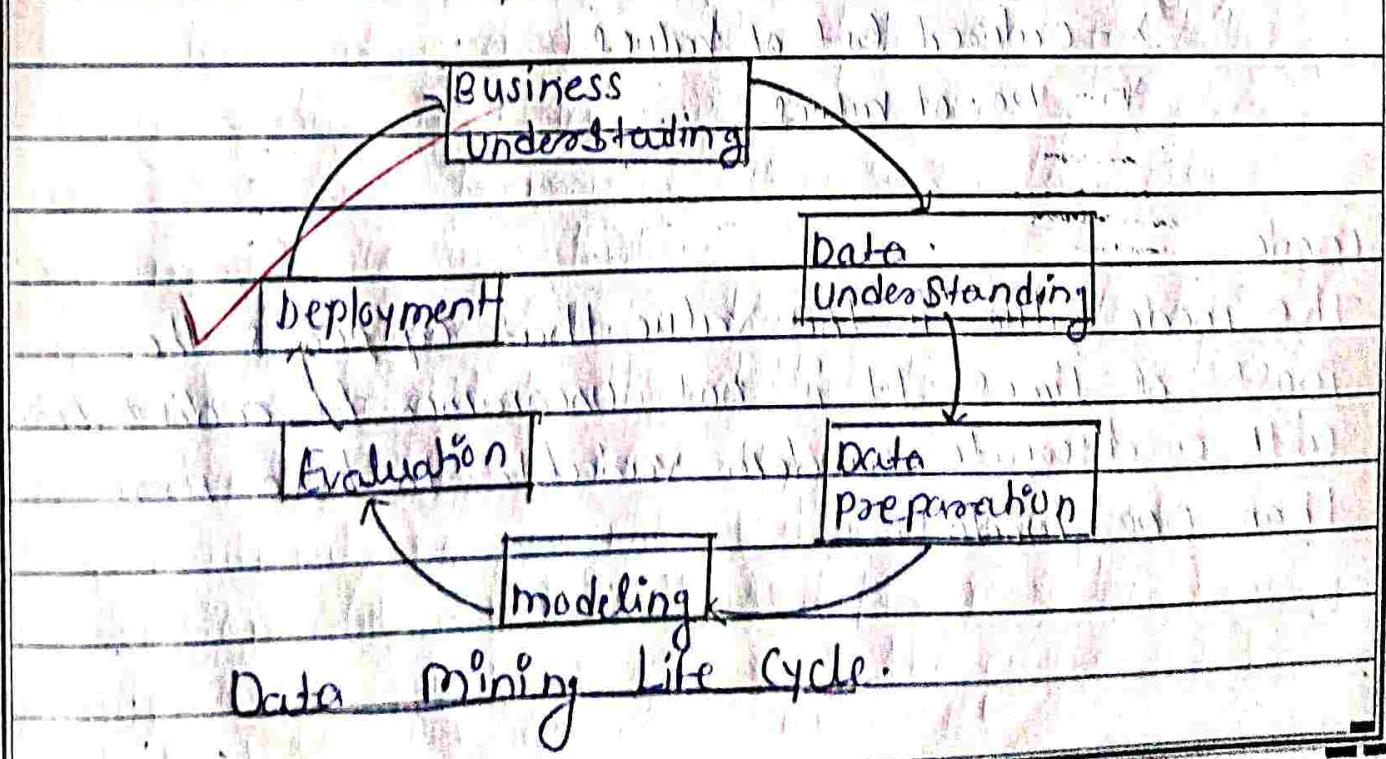
Histogram

Histogram groups the data in bins and is the fastest way to get idea about the distribution of each attribute in dataset.

- It provides visual count of the number of observations in each bin used for visualization.

- From the shape of the bin, we can easily observe the distribution line, whether it is Gaussian, Skewed or exponential.
- Histogram also helps us to see possible outliers.

CRISP-DM (Cross-Industry Standard process) for Standard process for Data mining, is an industry proven way to guide your data mining efforts. As methodology, It includes descriptions of the typical phases of a project, the tasks involved with each phase, and an explanation of relationship between these tasks. CRISP-DM provides an overview of the data mining life cycle. The CRISP-DM model is flexible and can be customized easily.



Mean

The mean value is the average value. To find the mean, find the sum of all values, and divide by the number of values.

formula: $\frac{\sum x_i}{N}$

median

The median of a distribution with a random variable depends on whether the number terms in distribution is even or odd. If the number of terms is odd, then the median is the value of the middle term even or equal number. It is difficult to find the median if

formula: $x_{\lceil \frac{n+1}{2} \rceil}$ if n is odd

$$\frac{x_{\lceil \frac{n}{2} \rceil} + x_{\lfloor \frac{n}{2} \rfloor}}{2} \text{ if } n \text{ is even.}$$

x - ordered list of values

n - no. of values

mode

The mode value is the value that appears the most of times. It is not uncommon for a distribution with a discrete random variable to have more than one mode.

Standard Deviation

Standard Deviation is a number that describes how spread out the values are. A low standard deviation means that the most of the numbers are close to the mean (average) value. A High Standard deviation means values are spread out over a wider range.

$$\text{Standard Deviation} = \sqrt{\frac{\sum (x_i - \mu)^2}{N}}$$

- Use the Numpy std() method to find the Standard Deviation.

Variance

Variance is another number that indicates how spread out values are, in fact, if you take square root of the variance, you get the Standard Deviation.

$$\text{Variance} = \frac{\sum (x_i - \mu)^2}{N}$$

- Use Numpy var() method to find Variance.

Scatter- plot-

In a dataset, for k set of variables/columns (x_1, x_2, \dots, x_k), the Scatter plot all the pairwise scatter between different variable in the form of a matrix.

It gives Answer of following questions like

- i) Are there any outlier in the dataset?
- ii) Is there any clustering by groups present in the dataset on the basis of particular variable.

practical No. 02

Aim :- Perform Data cleaning On given Dataset.

Apparatus :- System Library-numpy, pandas, seaborn, winsorize
Theory :-

Data cleaning is a crucial process in Data mining. It carries an important part in the building of a model. Data cleaning can be regarded as the process needed, but everyone often neglects it. Data quality is main issue in quality information management. Data quality problems occur anywhere in information systems. These problems are solved by Data cleaning. Data cleaning is fixing or removing incorrect, corrupted, incorrectly formatted, duplicate, or incomplete data within a dataset. If data is incorrect outcomes and algorithms are unreliable, even though they may look correct.

Generally, data cleaning reduces errors and improves data quality. Correcting errors in data and eliminating bad records can be a time-consuming and tedious process. But it cannot be ignored. Data mining has various techniques that are suitable for data cleaning.

Null Values

The values which is unknown or missing are called Null values. These values lead to inaccurate prediction. So, the machine learning developer deal with data cleaning, it removes or replace this null values.

`isnull()` → This will return boolean value for every column in data frame.

`isnull.sum()` → This code will give you total number of null values in each features in the data frame.

How to treat Null Values.

- i) you can drop the missing values with the method `df.filter()` ~~df.drop()~~. This will drop all the row which contain the missing value.
- ii) you can fill the missing values with method `df.fillna`

Winsorization

is the process of replacing the extreme values of statistical data in order to limit the effect of the outliers on the calculations or the results obtained by using that data. The mean value calculated after the such replacement of the extreme values is called winsorized mean.

Removing duplicate or irrelevant observations

Remove duplicate values from your dataset, including unwanted / irrelevant values. Duplicate values will be happen most often during data collection when you combine data sets from multiple places. Scrape data, receive data. These are opportunities to create duplicate that data.

dup

duplicated() → function indicate duplicate series values. The duplicated values are indicated as True values in the resulting series.

drop_duplicates() → helps in removing duplicates from the pandas DataFrame in python.

Outliers

In the series normal peaks and valleys occur outside of time frame when that seasonal sequence is normal or as a combination of time series that is in an outlier state as a group. Outliers detecting using visualization Once we are able to visualize the outlier then it becomes quite easy to decide what actions we can take.

We can remove outlier using the winsorization method.

Cause of outliers

- i) A lot of mispredictions.
- ii) Bad performance on train and test data.
- iii) Accuracy decrease.

One Hot Encoding

One approach to solve the Non-numerical part data problem can be Label Encoding where we will assign a numerical value to these labels for example male and female mapped to 0 and 1. But this can add bias in model performance and also since the data is string labels it will start high preference to the female parameter. Ideally both labels are equally important in the dataset. To deal with this issue we will use One Hot Encoding technique.

Ex. pass ATKT fail Column labels. In table

	Pass	ATKT	fail
1	1	0	0
2	1	0	0
3	0	1	0
4	0	0	1

Label Encoding

Label Encoding refers to converting the labels into a numeric form so as to convert them into the machine-readable form. machine learning algorithms can then decide in a better way how those labels must be operated. It is an important pre-processing step for the Structured dataset in Supervised Learning.

Example:-

Suppose we have Column Height in some dataset. After applying Label encoding, the Height column is converted: where 0 is the label for tall, 1 is label for medium and 2 is label for short height.

Normalization

Normalization is no mandate for all datasets available in machine learning. It is used whenever the attributes of the dataset have different ranges. It helps to enhance the performance and flexibility of machine learning model. It simply contains values in a range of 0 to 1:

formula :-

$$\text{Normalization } (x_n) = \frac{(x - x_{\min})}{(x_{\max} - x_{\min})}$$

x_{\max} = maximum value of a feature
 x_{\min} = minimum value of a feature

Standardization

Standardization or z-score normalization is the transformation of features by subtracting from mean and dividing by standard deviation. This is often called z-score.

$$x_{\text{new}} = \frac{(x - \text{mean})}{\text{std. deviation}}$$

Standardization can be helpful in cases where the data follows Gaussian distribution.

Practical NO. 03

Aim :- Apply hierarchical clustering on given datasets

Apparatus :- System, python libraries - numpy, pandas, sklearn, scipy.

Hierarchical clustering:

In other clustering algorithm like KMeans Clustering we face some challenges (which are a predetermined number of clusters and it always tries to create of the same size). To solve these two challenges, we can use the hierarchical clustering algorithm. It is unsupervised machine learning algorithm.

The hierarchical clustering technique has two approaches.

1) Agglomerative

is a bottom-up approach, in which the algorithm starts with taking all data points as single clusters and merging them until one cluster is left.

2) Divisive

Divisive algorithm is the reverse of the agglomerative algorithm as it is top-down approach.

Agglomerative hierarchical clustering

is a popular example of HCA. To group the data into clusters, it follows the bottom-up approach. Its bottom-up Approach, in which the algorithm starts with taking all the data points as single clusters and merging them until one cluster is left.

The working of AHC will be as follows:

Step-1 → Create each data as single cluster.

Step-2 → Take two closest data points as clusters and merge them to form one cluster.

Step-3 → Again, take two clusters and merge them together to form one cluster.

Step-4 → Repeat 3 until only one cluster left.

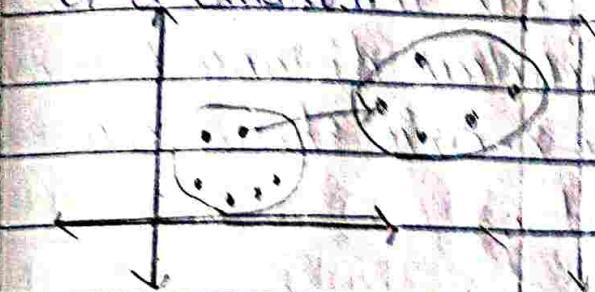
Step-5 → Once all clusters combined into one big cluster, develop the dendrogram to divide the cluster as per the problem.

The way of to decide the rule for clustering we use linkage methods.

Linkage methods are given below:

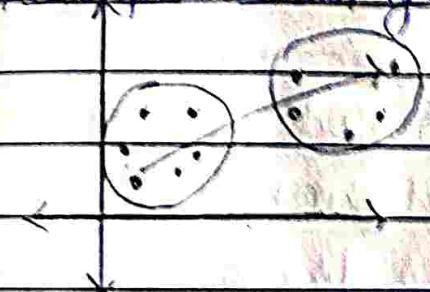
1) Single linkage

It is the shortest distance between the closest points of a cluster.



2) Complete linkage

It is the farthest distance between the two points of two different clusters.



3) Average linkage

It is the linkage method in which the distance between each pair of datasets is added up and then divided by the total number of the dataset to calculate average distance of cluster.

4) Centroid linkage

~~It is the linkage method in which the distance between the centroid of the cluster is calculated.~~

Dendrogram

The dendrogram is a tree-like structure that is mainly used to store each step of memory that the HC algorithm performs. In the dendrogram plot, the y-axis shows Euclidean distance between the data points, and the x-axis shows all the points of given dataset.

Practical No. 04

Aim :- Apply K-means clustering on given Dataset.

Apparatus:- Syder, numpy, pandas, matplotlib.

Theory :- Sklearn.

K-means Clustering Algorithm

is an ~~semi~~ unsupervised learning algorithms that is used to solve the clustering problems in machine learning or data science. It groups the unlabeled dataset into different clusters. Here k defines the number of pre-defined clusters that need to be created in the process, as if $k=2$, there will be two clusters, and for $k=3$, there will be three clusters and so on.

It is an iterative algorithm that divides the unlabeled dataset into k different. It is a centroid based algorithm, where each cluster is associated with a centroid. The main aim of this algorithm is to minimize the sum of distance between the dataset point and their corresponding clusters.

Working of k-means algorithm.

Step-1 → Select the number of k to decide no. of clusters.

Step-2 → Select random k point as centroids.

Step-3 → Assign each data point to their closest centroid which will form the predefined k .

Step-4 → calculate the variance and place a new centroid of each cluster.

Step-5 → Repeat the third steps, which means data points to the new closest centroid of each other.

Step-6 → if any reassignment occurs, then go to 4 else to FINISH.

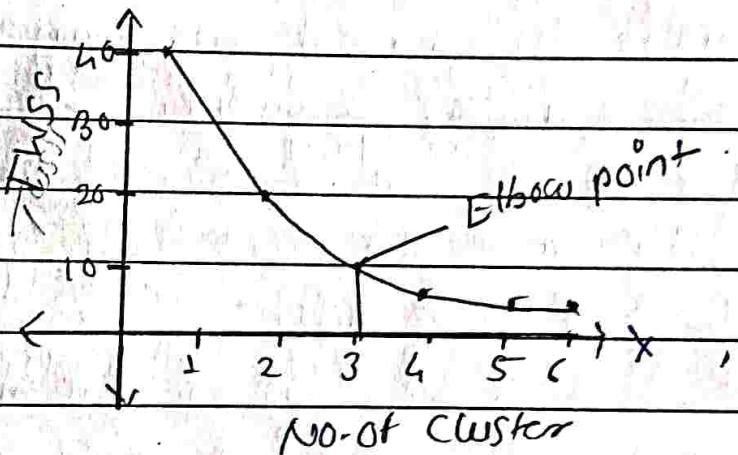
Step-7 → Model is Ready.

Elbow Curve.

The elbow method is a graphical representation of finding the optimal k in a k-means clustering by finding the sum of the square distance (Total within sum of squares) the sum of the square distance

between points in a cluster and the cluster centroid. The elbow graph shows TWSS values (on the y-axis) corresponding to the different values of k (on the x-axis). When we see an elbow shape in the graph, we pick the k -value where the elbow gets created. We can call this point the Elbow point.

Elbow Curve.



TWSS (Total Within Sum of Square)

TWSS values were calculated for all the clustering methods used. We compute the distance between each data point in the cluster and the centroid of the cluster for every cluster. This distance is referred to as ~~Within-cluster sum of squared errors (WSS)~~. We then sum these distances to obtain what we call the total within-cluster sum of squared error (TWSS).

$$T_{WSS} = (n - (k))^2$$

Euclidean Distance

It is a distance measure between a pair of samples in an n-dimensional feature space.

$$\text{Euclidean distance} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$

Practical No. 05

Aim :- Apply Simple Linear Regression Algorithm of given data set.

Apparatus :- Syder, Python, liblinear, scikit-learn, Scipy, Matlab, Octave, R, SPSS, Minitab, SPSS Modeler.

Simple Linear Regression

is a type of regression algorithms that model relationship between a dependent variable and a single independent variable. The relationship shown by a simple linear regression model is linear sloped straight line, hence it is called Simple Linear Regression. The key point in simple linear regression is that the dependent variable can be measured on continuous or categorical values.

Simple Linear Regression algorithm has mainly two objectives:

- 1) Model the relationship between the two variable such as the relationship between income and expenditure; experience and salary etc.

2) Forecasting new observations, such as weather forecasting according to temperature, Revenue of a company according to investment in a year etc.

Equation for representation of SLR:

$$y = B_0 + B_1 x + \epsilon$$

where,

B_0 = It is intercept of Regression Line.

B_1 = It is the slope of the Regression Line.

which tells whether the line is increasing or decreasing.

ϵ = The error term.

Confusion matrix

The confusion matrix is a matrix used to determine the performance of the classification models for a given set of test data. It can only be determined if the true values for the data are known. The matrix itself can be easily understood but the related terminologies may be confusing. Since it shows the errors in the model performance in the form of matrix, hence also known as an error matrix.

		Actual value		
		1	0	
Predicted value	1	TP	FP	
	0	FN	TN	

TP → True Negative-positive

FP → False positive

$$\text{Accuracy} = \frac{TP+TN}{TP+TN+FP+FN}$$

TN → True negative

$$\frac{TP+TN}{TP+TN+FP+FN}$$

FN → false negative

Regression line

It is used for relationship between dependent & independent values.

- 1) forecasting
- 2) Trend Analysis
- 3) To make prediction

Type's of Regression line

- 1) linear
- 2) logistic
- 3) Polynomial

Practical NO.06

Aim:- Apply Multiple Linear Regression algorithm on given dataset.

Apparatus:- Spyder, sklearn ~~for linear multiple regression~~

Theory:-

Multiple Linear Regression

is an extension of Simple Linear regression as it models the linear relationship between a single dependent variable and more than one predictor Variable to predict the response Variable. We can define it as: Multiple 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.

MLR equation

In multiple linear regression, the target Variable (y) is a linear combination of multiple predictor Variables $x_1, x_2, x_3, \dots, x_n$. Since it is an enhancement of simple linear regression. So, the same is applied for the multiple linear regression equation, the equation becomes

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

$$Y = b_0 + b_1 X_1 +$$

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

where

Y = Output

$b_0, b_1, b_2, \dots, b_n$ = coefficients of the model

X_1, X_2, X_3 = various independent

Variable.

Pairplot

pairplot visualization comes handy when you want to go for Exploratory data analysis (EDA) pairplot visualizes given data to find the relationship between them where the Variable can be continuous or categorical. pairwise relationship in data-set.

pairplot is a module of seaborn library which provides a high-level interface for drawing attractive and informative Statistical graphics.

pairplot is used to understand the best set of features to explain a relationship between two Variable or to form the most separated clusters.

Residual · Leverage plot (Diagnostic plots)

In linear or multiple regression, it is not enough to just fit the model into the dataset. But, it may not give the desired result. To apply the linear or multiple regression efficiently to the dataset.

Forms in Regression Diagnostic plots

Outlier

Outliers are the points that are distinct and deviant from the bulk of the dataset.

Leverage points

A leverage point is defined as an observation that has a value x that is far away from mean.

Influential Points

is defined as an observation that has a large influence on the fit of one model.

Residual VS Fitted plot

The residual can be calculated As

$$\text{res} = y_{\text{observed}} - y_{\text{predicted}}$$

Practical NO. 07

Aim :- Apply K-Nearest Neighbor (KNN) Algorithm on Given Dataset.

Apparatus :- Syder, numpy, pandas, matplotlib, Sklearn, k-neighbor classifier.

Theory :-

K-Nearest Neighbor (KNN) Algorithm

K-Nearest Neighbor is one of the simplest machine learning algorithms on supervised learning technique. It assumes the similarity between the new case / data and available cases and put the new case into the category that is most similar to the available categories.

K-NN algorithm stores all the available data and classifies a new data point based on the similarity. This means when new data appears then it can be easily classified into a well-defined category by using K-NN algorithm.

K-NN algorithm stores all can be used for Regression like recall for classification but mostly it is used for the classification problems. K-NN is a non-parametric algorithm which means it does not make any assumptions on underlying data.

It is also called a lazy learner algorithm because it does not learn from the training set immediately instead it stores the dataset and at time of classification, it performs an action on the dataset. KNN algorithm at the training phase just stores the dataset and when it gets new data, then it classifies that data into a category that is much similar to new data.

Working of k-NN algorithm

Step-1 → Select the number k of the neighbors

Step-2 → Calculate the Euclidean distance of k numbers of neighbors.

Step-3 → Take the k nearest neighbors as per the calculated Euclidean distance.

Step-4 → Among these k neighbors, as per the calculated Euclidean count the number of the data points in each category.

Step-5 → Assign the new data points to that category for which the number of the neighbors is maximum.

Advantages of k-NN Algorithm:

- i) It is simple to implement.
- ii) It is robust to the noisy training data.
- iii) It can be more effective if the training data is large.

Disadvantages of kNN algorithm:

- i) Always needs to determine the value of k , which may be complex some time.
- ii) The computation cost is high because of calculating the distance between the data points for all the training samples.

Practical No. 08.

Aim:- Apply Decision Tree Algorithm on given datasets.

Apparatus:- Syder, numpy, pandas, matplotlib, sklearn, Decision tree classifier.

Theory:-

Decision Tree Algorithm

Decision Tree is supervised learning technique that can be used for both classification and regression problems, but mostly it is preferred for solving classification problems. It is a tree-structured classifier, where internal nodes represent the features of a dataset, branches represent the decision rules and each leaf node represents the outcome.

In a Decision tree there are two nodes which are the Decision Node and leaf node.

Decision nodes are used to make any decision and have multiple branches, whereas leaf nodes are the output of these decisions and do not contain any further branch.

These are various algorithm in machine learning so

choosing the best algorithm for given dataset and problem is the main point to remember while

Creating a machine learning model.

Decision Trees usually mimic human thinking ability while making a decision, so it is easy to understand; The logic behind the decision tree can be easily understand because it shows a tree-like structure.

Decision Tree Terminologies

- 1) Root node → is from where decision tree starts. It represents the entire dataset, which further gets divided into two or more homogeneous sets.
- 2) Leaf node → is the final output node, and tree cannot be segregated further after getting a leaf node.
- 3) Splitting → is the process of dividing the tree - decision node/root into sub-nodes.
- 4) Branch/SubTree → A tree formed by splitting the tree.
- 5) Pruning → is the process of removing the unwanted branches from the tree.
- 6) Parent/Child node → The root node of tree is called the parent node, and other nodes are called the child nodes.

Working of Decision Tree Algorithm:

- Step-1 Begin the tree with the root node. Say S, which is complete dataset.
- Step-2 Find the best attribute in the dataset using Attribute Selection (AMI).
- Step-3 Divide the S into subsets that contain possible values for the best attribute.
- Step-4 Create the decision tree node, which contains best attribute.
- Step-5 Recursively make new decision trees using the subsets of the dataset created.
- Step-6 Continue the decision tree nodes process until a stage is reached where you cannot further classify the nodes and the final node as a flat node.

Advantages of Decision Tree

- i) It is simple to understand.
- ii) It can be very useful for solving decision-related problems.
- iii) It helps to think about all the possible outcomes.

Disadvantages of Decision Tree

- i) If the decision tree contains lots of layers, which makes it complex.
- ii) It may be prone to overfitting issue.

Entropy

Information gain is defined as pattern observed in the dataset and reduction in the entropy.

Entropy is defined as the randomness or measuring the disorder of information being in machine learning. Entropy is metric that measures the unpredictability or impurity in the system.

gain = $\text{entropy before split} - \text{entropy after split}$

$$(G) = E_B - E_A$$

$$\text{Info}_B(D_j) = \sum_{i=1}^n p_i \log_2(p_i)$$

$$\text{Info}_D(D) = \sum_{j=1}^n |D_j| * \text{info}(D_j)$$

10