**KAMMAVARI SANGHAM® 1952**
**K.S. SCHOOL OF ENGINEERING AND MANAGEMENT**

**Department of Computer Science & Engineering**

# Artificial Intelligence and Machine Learning Laboratory 18CSL76



# Lab Manual

**Academic Year 2021 - 2022**

K.S. Group of Institutions
No. 15, Mallasandra, off. Kanakapura Road,
Bengaluru - 560 109
www.kssem.edu.in

| ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING LABORATORY (Effective from the academic year 2018 - 2019) | | | |
|---|---|---|---|
| **SEMESTER – VII** | | | |
| **Course Code** | **18CSL76** | **CIE Marks** | 40 |
| **Number of Contact Hours/Week** | 0:0:2 | **SEE Marks** | 60 |
| **Total Number of Lab Contact Hours** | 36 | **Exam Hours** | 03 |
| **Credits – 2** | | | |

**Course Learning Objectives:** This course (18CSL76) will enable students to:

- Implement and evaluate AI and ML algorithms in and Python programming language.

**Descriptions (if any):**

**Installation procedure of the required software must be demonstrated, carried out in groups and documented in the journal.**

**Programs List:**

1. Implement A* Search algorithm.
2. Implement AO* Search algorithm.
3. For a given set of training data examples stored in a .CSV file, implement and demonstrate the Candidate-Elimination algorithm to output a description of the set of all hypotheses consistent with the training examples.

4. Write a program to demonstrate the working of the decision tree based ID3 algorithm. Use an appropriate data set for building the decision tree and apply this knowledge to classify a new sample.

5. Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets.

6. Write a program to implement the naïve Bayesian classifier for a sample training data set stored as a .CSV file. Compute the accuracy of the classifier, considering few test data sets.

7. Apply EM algorithm to cluster a set of data stored in a .CSV file. Use the same data set for clustering using k-Means algorithm. Compare the results of these two algorithms and comment on the quality of clustering. You can add Java/Python ML library classes/API in the program.

8. Write a program to implement k-Nearest Neighbour algorithm to classify the iris data set. Print both correct and wrong predictions. Java/Python ML library classes can be used for this problem.

9. Implement the non-parametric Locally Weighted Regression algorithm in order to fit data points. Select appropriate data set for your experiment and draw graphs

**Laboratory Outcomes**: The student should be able to:
- Implement and demonstrate AI and ML algorithms.
- Evaluate different algorithms.

**Conduct of Practical Examination:**
- Experiment distribution
    - o For laboratories having only one part: Students are allowed to pick one experiment from the lot with equal opportunity.
    - o For laboratories having PART A and PART B: Students are allowed to pick one experiment from PART A and one experiment from PART B, with equal opportunity.
- Change of experiment is allowed only once and marks allotted for procedure to be made zero of the changed part only.
- Marks Distribution *(Courseed to change in accoradance with university regulations)*
    - q) For laboratories having only one part – Procedure + Execution + Viva-Voce: 15+70+15 =
      100  Marks
    - r) For laboratories having PART A and PART B
        - i. Part A – Procedure + Execution + Viva = 6 + 28 + 6 = 40 Marks
        - ii. Part B – Procedure + Execution + Viva = 9 + 42 + 9 = 60 Marks

**Machine learning**

Machine learning is a subset of artificial intelligence in the field of computer science that often uses statistical techniques to give computers the ability to "learn" (i.e., progressively improve performance on a specific task) with data, without being explicitly programmed. In the past decade, machine learning has given us self-driving cars, practical speech recognition, effective web search, and a vastly improved understanding of the human genome.

**Machine learning tasks**

Machine learning tasks are typically classified into two broad categories, depending on whether there is a learning "signal" or "feedback" available to a learning system:

1. **Supervised learning:** The computer is presented with example inputs and their desired outputs, given by a "teacher", and the goal is to learn a general rule that maps inputs to outputs. As special cases, the input signal can be only partially available, or restricted to special feedback:
2. **Semi-supervised learning:** the computer is given only an incomplete training signal: a training set with some (often many) of the target outputs missing.
3. **Active learning:** the computer can only obtain training labels for a limited set of instances (based on a budget), and also has to optimize its choice of objects to acquire labels for. When used interactively, these can be presented to the user for labeling.
4. **Reinforcement learning:** training data (in form of rewards and punishments) is given only as feedback to the program's actions in a dynamic environment, such as driving a vehicle or playing a game against an opponent.
5. **Unsupervised learning:** No labels are given to the learning algorithm, leaving it on its own to find structure in its input. Unsupervised learning can be a goal in itself (discovering hidden patterns in data) or a means towards an end (feature learning).

| Supervised learning | Un Supervised learning | Un Supervised learning |
|---|---|---|
| Find-s algorithm EM algorithm | Find-s algorithm EM algorithm | Locally weighted Regression algorithm |
| Candidate elimination algorithm | K means algorithm | |
| Decision tree algorithm | | |
| Back propagation Algorithm | | |
| Naïve Bayes Algorithm | | |
| K nearest neighbor algorithm(lazy learning algorithm) | | |

## Machine learning applications

In **classification,** inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more (multi-label classification) of these classes. This is typically tackled in a supervised manner. Spam filtering is an example of classification, where the inputs are email (or other) messages and the classes are "spam" and "not spam".

In **regression,** also a supervised problem, the outputs are continuous rather than discrete.

In **clustering,** a set of inputs is to be divided into groups. Unlike in classification, the groups are not known beforehand, making this typically an unsupervised task.

**Density estimation** finds the distribution of inputs in some space.

**Dimensionality reduction** simplifies inputs by mapping them into a lower dimensional space. Topic modeling is a related problem, where a program is given a list of human language documents and is tasked with finding out which documents cover similar topics.

## Machine learning Approaches

### 1. Decision tree learning

Decision tree learning uses a decision tree as a predictive model, which maps observations about an item to conclusions about the item's target value.

## 2. Association rule learning

Association rule learning is a method for discovering interesting relations between variables in large databases.

## 3. Artificial neural networks

An artificial neural network (ANN) learning algorithm, usually called "neural network" (NN), is a learning algorithm that is vaguely inspired by biological neural networks. Computations are structured in terms of an interconnected group of artificial neurons, processing information using a connectionist approach to computation. Modern neural networks are non-linear statistical data modeling tools. They are usually used to model complex relationships between inputs and outputs, to find patterns in data, or to capture the statistical structure in an unknown joint probability distribution between observed variables.

## 4. Deep learning

Falling hardware prices and the development of GPUs for personal use in the last few years have contributed to the development of the concept of deep learning which consists of multiple hidden layers in an artificial neural network. This approach tries to model the way the human brain processes light and sound into vision and hearing. Some successful applications of deep learning are computer vision and speech Recognition.

## 5. Inductive logic programming

Inductive logic programming (ILP) is an approach to rule learning using logic Programming as a uniform representation for input examples, background knowledge, and hypotheses. Given an encoding of the known background knowledge and a set of examples represented as a logical database of facts, an ILP system will derive a hypothesized logic program that entails all positive and no negative examples. Inductive programming is a related field that considers any kind of programming languages for representing hypotheses (and not only logic programming), such as functional programs.

## 6. Support vector machines

Support vector machines (SVMs) are a set of related supervised learning methods used for classification and regression. Given a set of training examples, each marked as belonging to one of two categories, an SVM training algorithm builds a model that predicts whether a new example falls into one category or the other.

## 7. Clustering

Cluster analysis is the assignment of a set of observations into subsets (called clusters) so that observations within the same cluster are similar according to some pre designated criterion or criteria, while observations drawn from different clusters are dissimilar. Different clustering techniques make different assumptions on the structure of the data, often defined by some similarity metric and evaluated for example by internal compactness (similarity between members of the same cluster) and separation between different clusters. Other methods are based on estimated density and graph connectivity. Clustering is a method of unsupervised learning, and a common technique for statistical data analysis.

## 8. Bayesian networks

A Bayesian network, belief network or directed acyclic graphical model is a probabilistic graphical model that represents a set of random variables and their conditional independencies via a directed acyclic graph (DAG). For example, a Bayesian network could represent the probabilistic relationships between diseases and symptoms. Given symptoms, the network can be used to compute the probabilities of the presence of various diseases. Efficient algorithms exist that perform inference and learning.

## 9. Reinforcement learning

Reinforcement learning is concerned with how an agent ought to take actions in an environment so as to maximize some notion of long- term reward. Reinforcement learning algorithms attempt to find a policy that maps states of the world to the actions the agent ought to take in those states. Reinforcement learning differs from the supervised learning problem in that correct input/output pairs are never presented, nor sub-optimal actions explicitly corrected.

10. **Similarity and metric learning**

In this problem, the learning machine is given pairs of examples that are considered similar and pairs of less similar objects. It then needs to learn a similarity function (or a distance metric function) that can predict if new objects are similar. It is sometimes used in Recommendation systems.

11. **Genetic algorithms**

A genetic algorithm (GA) is a search heuristic that mimics the process of natural selection, and uses methods such as mutation and crossover to generate new genotype in the hope of finding good solutions to a given problem. In machine learning, genetic algorithms found some uses in the 1980s and 1990s. Conversely, machine learning techniques have been used to improve the performance of genetic and evolutionary algorithms.

12. **Rule-based machine learning**

Rule-based machine learning is a general term for any machine learning method that identifies, learns, or evolves "rules" to store, manipulate or apply, knowledge. The defining characteristic of a rule-based machine learner is the identification and utilization of a set of relational rules that collectively represent the knowledge captured by the system. This is in contrast to other machine  learners that commonly identify a singular model that can be universally applied to any instance in order to make a prediction. Rule-based machine learning approaches include learning classifier systems, association rule learning, and artificial immune systems.

13. **Feature selection approach**

Feature selection is the process of selecting an optimal subset of relevant features for use in model construction. It is assumed the data contains some features that are either redundant or irrelevant, and can thus be removed to reduce calculation cost without incurring much loss of information. Common optimality criteria include accuracy, similarity and information measures.

### 1. Implement A* Search algorithm.

```
#from collections import deque

class Graph:
    def __init__(self, adjac_lis):
        self.adjac_lis = adjac_lis

    def get_neighbors(self, v):
        return self.adjac_lis[v]

    # This is heuristic function which is having equal values for all nodes
    def h(self, n):
        H = {'A': 1,'B': 1,'C': 1,'D': 1}

        return H[n]

    def a_star_algorithm(self, start, stop):
        open_lst = set([start])
        closed_lst = set([])

        poo = {}
        poo[start] = 0

        par = {}
        par[start] = start

        while len(open_lst) > 0:
            n = None

            # it will find a node with the lowest value of f() -
            for v in open_lst:
                if n == None or poo[v] + self.h(v) < poo[n] + self.h(n):
                    n = v;

            if n == None:
                print('Path does not exist!')
                return None

            # if the current node is the stop
            # then we start again from start
            if n == stop:
                reconst_path = []

                while par[n] != n:
                    reconst_path.append(n)
```

```
            n = par[n]

        reconst_path.append(start)

        reconst_path.reverse()

        print('Path found: {}'.format(reconst_path))
        return reconst_path

    # for all the neighbors of the current node do
    for (m, weight) in self.get_neighbors(n):

      # if the current node is not presentin both open_lst and closed_lst
        # add it to open_lst and note n as it's par
        if m not in open_lst and m not in closed_lst:
            open_lst.add(m)
            par[m] = n
            poo[m] = poo[n] + weight

        else:
            if poo[m] > poo[n] + weight:
                poo[m] = poo[n] + weight
                par[m] = n

                if m in closed_lst:
                    closed_lst.remove(m)
                    open_lst.add(m)

    open_lst.remove(n)
    closed_lst.add(n)

  print('Path does not exist!')
  return None

adjac_lis = {
   'A': [('B', 1), ('C', 3), ('D', 7)],
   'B': [('D', 5)],
   'C': [('D', 12)]
}
graph1 = Graph(adjac_lis)
graph1.a_star_algorithm('A', 'D')
```

**OUTPUT**

```
Path found: ['A', 'B', 'D']

['A', 'B', 'D']
```

## 2. Implement AO* Search algorithm.

```python
#from collections import deque

class Graph:
    def __init__(self, adjac_lis):
        print('adjac_lis values are ==== ' , adjac_lis)
        self.adjac_lis = adjac_lis

    def get_neighbors(self, v):
        print(" v values are --------", v)
        return self.adjac_lis[v]

    # This is heuristic function which is having equal values for all nodes
    def h(self, n):
        H = {'A': 1,'B': 1,'C': 1,'D': 1}
        #print(H)
        return H[n]

    def a_star_algorithm(self, start, stop):
        open_lst = set([start])
        closed_lst = set([])

        print("open list initially ))))))))", open_lst)
        print("closed list initially ))))))))))))))", closed_lst)

        poo = {}
        poo[start] = 0

        par = {}
        par[start] = start

        print("par value &&&&&&&&&&", par)

        while len(open_lst) > 0:
            n = None

            # it will find a node with the lowest value of f() -
            for v in open_lst:
                if n == None or poo[v] + self.h(v) < poo[n] + self.h(n):
                    print("n value ", n)
                    print("v value ", v)
                    n = v;

            if n == None:
                print('Path does not exist!')
```

```
            return None

    # if the current node is the stop
    # then we start again from start
    if n == stop:
        reconst_path = []
        print("reconst_path", reconst_path)

        while par[n] != n:
            reconst_path.append(n)
            n = par[n]
            print("reconst_path within while", reconst_path)
        reconst_path.append(start)

        reconst_path.reverse()

        print('Path found: {}'.format(reconst_path))
        print("reconst_path", reconst_path)
        return reconst_path

    # for all the neighbors of the current node do
    for (m, weight) in self.get_neighbors(n):
      # if the current node is not presentin both open_lst and closed_lst
        # add it to open_lst and note n as it's par
        if m not in open_lst and m not in closed_lst:
            print("open list  within while and if ",open_lst)
            print("closed  list  within while and if ",closed_lst)
            open_lst.add(m)
            par[m] = n
            poo[m] = poo[n] + weight

            print("open list  within while and if after ... ",open_lst)
            print("closed list  within while and if after ... ",closed_lst)
            print("poo values within if ", poo)

        else:
            if poo[m] > poo[n] + weight:
                poo[m] = poo[n] + weight
                par[m] = n

                print("poo values within else ", poo)

                if m in closed_lst:
                    closed_lst.remove(m)
                    open_lst.add(m)
                    print("closed list within while and if ",closed_lst)

    print("open===========----  list within while and if ",open_lst)
```

```
        print("closed ========= list within while and if ",closed_lst)

        open_lst.remove(n)
        closed_lst.add(n)

    print('Path does not exist!')
    return None

adjac_lis = {
  'A': [('B', 1), ('C', 3), ('D', 7)],
  'B': [('D', 5)],
  'C': [('D', 12)]
}
graph1 = Graph(adjac_lis)
# print("graph1 is ", graph1)
graph1.a_star_algorithm('A', 'D')

adjac_lis = {
  'A': [('B', 1), ('C', 3), ('D', 7)],
  'B': [('D', 5)],
  'C': [('D', 12)]
}
graph1 = Graph(adjac_lis)
# print("graph1 is ", graph1)
graph1.a_star_algorithm('A', 'D')
```

### OUTPUT

```
adjac_lis values are ====  {'A': [('B', 1), ('C', 3), ('D', 7)], 'B': [('D'
, 5)], 'C': [('D', 12)]}
open list initially ))))))) {'A'}
closed list initially ))))))))))) set()
par value &&&&&&&&& {'A': 'A'}
n value  None
v value  A
 v values are -------- A
open list  within while and if  {'A'}
closed  list  within while and if  set()
open list  within while and if after ...  {'B', 'A'}
closed list  within while and if after ...  set()
poo values within if  {'A': 0, 'B': 1}
open list  within while and if  {'B', 'A'}
closed  list  within while and if  set()
open list  within while and if after ...  {'B', 'A', 'C'}
closed list  within while and if after ...  set()
```

```
poo values within if  {'A': 0, 'B': 1, 'C': 3}
open list  within while and if  {'B', 'A', 'C'}
closed  list  within while and if  set()
open list  within while and if after ...  {'D', 'B', 'A', 'C'}
closed list  within while and if after ...  set()
poo values within if  {'A': 0, 'B': 1, 'C': 3, 'D': 7}
open=============---  list within while and if  {'D', 'B', 'A', 'C'}
closed ========= list within while and if  set()

n value  None
v value  D
n value  D
v value  B
 v values are -------- B
poo values within else  {'A': 0, 'B': 1, 'C': 3, 'D': 6}
open=============---  list within while and if  {'D', 'B', 'C'}
closed ========= list within while and if  {'A'}
n value  None
v value  D
n value  D
v value  C
 v values are -------- C
open=============---  list within while and if  {'D', 'C'}
closed ========= list within while and if  {'B', 'A'}
n value  None
v value  D
reconst_path []
reconst_path within while ['D']
reconst_path within while ['D', 'B']
Path found: ['A', 'B', 'D']
reconst_path ['A', 'B', 'D']
```

Out[1]:  ['A', 'B', 'D']

**3. For a given set of training data examples stored in a .CSV file, implement and demonstrate the Candidate-Elimination algorithm to output a description of the set of all hypotheses consistent with the training examples.**

```python
import numpy as np
import pandas as pd

# Loading Data from a CSV File
data = pd.DataFrame(data = pd.read_csv("C:/Users/mbacc/Documents/dataset/Training.csv"))
data

# Separating concept features from Target
concepts = np.array(data.iloc[:,0:-1])
concepts
target = np.array(data.iloc[:,-1])
target
def learn(concepts, target):
    specific_h = concepts[0].copy()
    general_h = [["?" for i in range(len(specific_h))] for i in range(len(specific_h))]

    # The learning iterations
    for i, h in enumerate(concepts):

        # Checking if the hypothesis has a positive target
        if target[i] == "Yes":
            for x in range(len(specific_h)):

                # Change values in S & G only if values change
                if h[x] != specific_h[x]:
                    specific_h[x] = '?'
                    general_h[x][x] = '?'

        # Checking if the hypothesis has a positive target
        if target[i] == "No":
            for x in range(len(specific_h)):
                print(f"specific={specific_h[x]}")
                # For negative hyposthesis change values only  in G
                if h[x] != specific_h[x]:
                    general_h[x][x] = specific_h[x]
                    #print(f"general{x}={general_h[x][x]}")
                else:
                    general_h[x][x] = '?'
    # find indices where we have empty rows, meaning those that are unchanged
    indices = [i for i,val in enumerate(general_h) if val == ['?', '?', '?', '?', '?', '?']]
```

```
    for i in indices:
        # remove those rows from general_h
        general_h.remove(['?', '?', '?', '?', '?', '?'])

    # Return final values
    return specific_h, general_h
s_final, g_final = learn(concepts, target)
s_final
g_final
```

## Dataset

Training.csv

| Sky | Airtemp | Humidity | Wind | Water | Forecast | WaterSport |
|-----|---------|----------|------|-------|----------|------------|
| Sunny | Warm | Normal | Strong | Warm | Same | Yes |
| Sunny | Warm | High | Strong | Warm | Same | Yes |
| Cloudy | Cold | High | Strong | Warm | Change | No |
| Sunny | Warm | High | Strong | Cool | Change | Yes |

### OUTPUT

```
specific=Sunny
specific=Warm
specific=?
specific=Strong
specific=Warm
specific=Same

Out[7]:
[['Sunny', '?', '?', '?', '?', '?'], ['?', 'Warm', '?', '?', '?', '?']]
```

**4. Write a program to demonstrate the working of the decision tree based ID3 algorithm. Use an appropriate data set for building the decision tree and apply this knowledge to classify a new sample.**

# import libraries

```
import numpy as np
import pandas as pd
```

#load dataset

```
from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()

data.data
data.feature_names
data.target
data.target_names
```

# create dataframe

```
df = pd.DataFrame(np.c_[data.data, data.target], columns=[list(data.feature_names)+['target']])
df.head()

df.tail()

row1=df.iloc[3]
row1

df.shape
```

#Split Data

```
X = df.iloc[:, 0:-1]
y = df.iloc[:, -1]

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2020)


print('Shape of X_train = ', X_train.shape)

print('Shape of y_train = ', y_train.shape)

print('Shape of X_test = ', X_test.shape)
print('Shape of y_test = ', y_test.shape)
```

#Train Decision Tree Classification Model

```
from sklearn.tree import DecisionTreeClassifier


classifier = DecisionTreeClassifier(criterion='gini')
classifier.fit(X_train, y_train)

classifier.score(X_test, y_test)

classifier_entropy = DecisionTreeClassifier(criterion='entropy')
classifier_entropy.fit(X_train, y_train)

classifier_entropy.score(X_test, y_test)
```

#Feature Scaling
```
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()

sc.fit(X_train)


X_train_sc = sc.transform(X_train)
X_test_sc = sc.transform(X_test)


classifier_sc = DecisionTreeClassifier(criterion='gini')

classifier_sc.fit(X_train_sc, y_train)


classifier_sc.score(X_test_sc, y_test)
```

#Predict Cancer
```
patient1 = [17.99,
 0.38,
 122.8,
 1001.0,
 0.1184,
 0.2776,
 0.3001,
 0.1471,
 0.2419,
```

0.07871,

1.095,

0.9053,

8.589,

153.4,

0.006399,

0.04904,

0.05373,

0.01587,

0.03003,

0.006193,

25.38,

17.33,

184.6,

2019.0,

0.1622,

0.6656,

0.7119,

0.2654,

0.4601,
0.1189]

```python
patient1 = np.array([[patient1])
patient1
classifier.predict(patient1)
data.target_names
pred = classifier.predict(patient1)


if pred[0] == 0:
  print('Patient has Cancer (malignant tumor)')
else:
  print('Patient has no Cancer (malignant benign)')
```
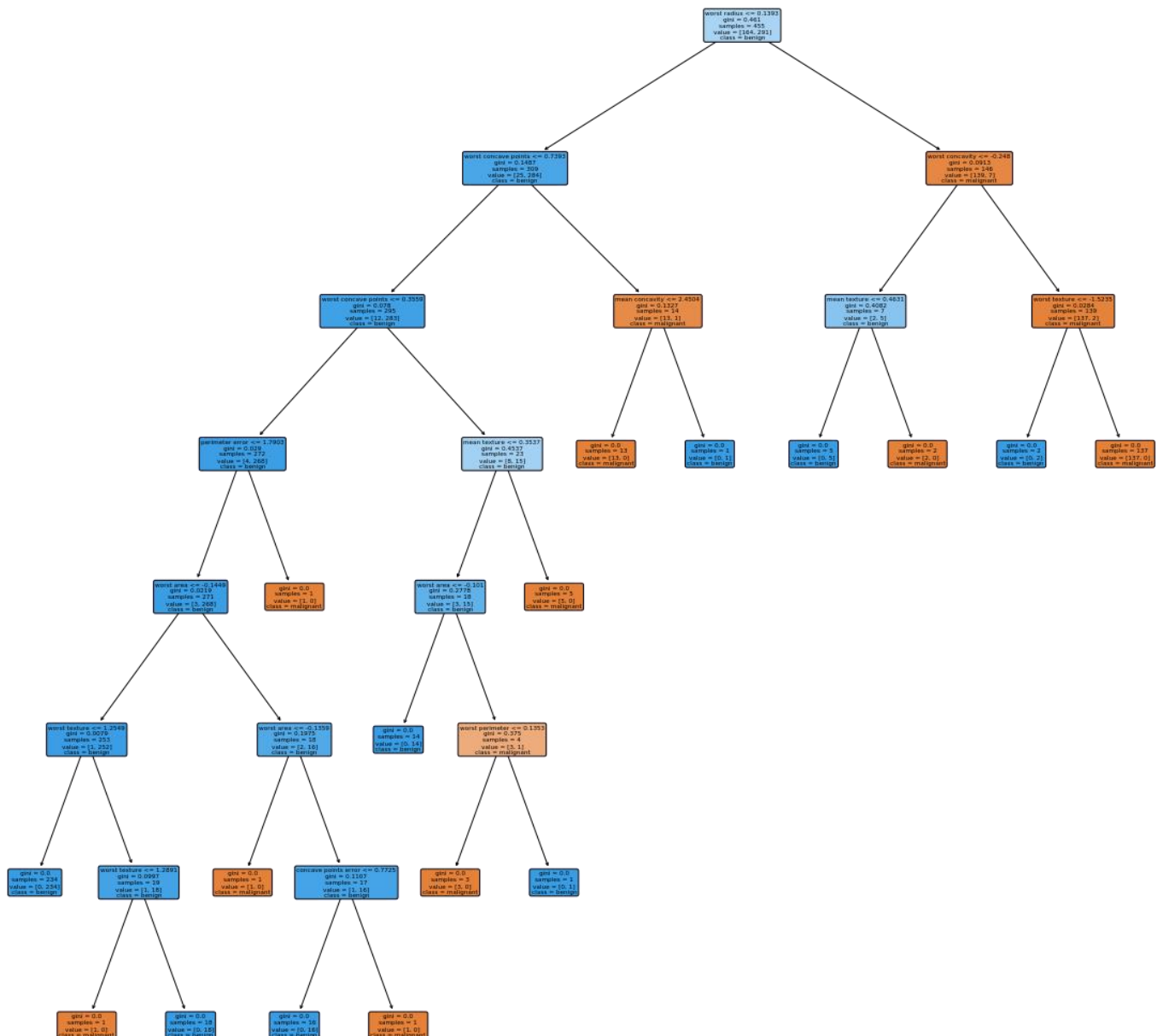
**OUTPUT**

```
Patient has Cancer (malignant tumor)
```

```
from sklearn.tree import plot_tree
import matplotlib.pyplot as plt
```

#Visualising the graph

plt.figure(figsize = (20,20))

dec_tree = plot_tree(decision_tree=classifier_sc, feature_names = data.feature_names,
        class_names =['malignant', 'benign'] , filled = True , precision = 4, rounded = True)

**OUTPUT**

## 5. Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets.

```python
import numpy as np

X = np.array(([2, 9], [1, 5], [3, 6]), dtype=float)    # X = (hours sleeping, hours studying)
y = np.array(([92], [86], [89]), dtype=float)          # y = score on test

# scale units
X = X/np.amax(X, axis=0)       # maximum of X array
y = y/100                      # max test score is 100

print(X)
print(y)

class Neural_Network(object):
  def __init__(self):
                    # Parameters
    self.inputSize = 2
    self.outputSize = 1
    self.hiddenSize = 3
                  # Weights
    self.W1 = np.random.randn(self.inputSize, self.hiddenSize)
                                    # (3x2) weight matrix from input to hidden layer
    self.W2 = np.random.randn(self.hiddenSize, self.outputSize)
                                    # (3x1) weight matrix from hidden to output layer
  def forward(self, X):
                          # forward propagation through our network
    self.z = np.dot(X, self.W1)       # dot product of X (input) and first set of 3x2 weights
    self.z2 = self.sigmoid(self.z)    # activation function
    self.z3 = np.dot(self.z2, self.W2)  # dot product of hidden layer (z2) and second set of 3x1 weights
    o = self.sigmoid(self.z3)         # final activation function
    return o

  def sigmoid(self, s):
    return 1/(1+np.exp(-s))           # activation function

  def sigmoidPrime(self, s):
    return s * (1 - s)                # derivative of sigmoid

  def backward(self, X, y, o):
                          # backward propgate through the network
    self.o_error = y - o              # error in output
```

```
      self.o_delta = self.o_error*self.sigmoidPrime(o)          # applying derivative of sigmoid to

      self.z2_error = self.o_delta.dot(self.W2.T)     # z2 error: how much our hidden layer weights
contributed to output error
      self.z2_delta = self.z2_error*self.sigmoidPrime(self.z2)       # applying derivative of sigmoid to z2
error

      self.W1 += X.T.dot(self.z2_delta)        # adjusting first set (input --> hidden) weights
      self.W2 += self.z2.T.dot(self.o_delta)  # adjusting second set (hidden --> output) weights

   def train (self, X, y):
      o = self.forward(X)
      self.backward(X, y, o)

NN = Neural_Network()
for i in range(1000):                           # trains the NN 1,000 times
   print ("\nInput: \n" + str(X))
   print ("\nActual Output: \n" + str(y))
   print ("\nPredicted Output: \n" + str(NN.forward(X)))
   print ("\nLoss: \n" + str(np.mean(np.square(y - NN.forward(X)))))       # mean sum squared loss)
   NN.train(X, y)
```

## OUTPUT

```
Input:
[[0.66666667 1.          ]
 [0.33333333 0.55555556]
 [1.         0.66666667]]

Actual Output:
[[0.92]
 [0.86]
 [0.89]]

Predicted Output:
[[0.52902915]
 [0.52826157]
 [0.53631347]]

Loss:
0.12933425267706222

Input:
[[0.66666667 1.          ]
 [0.33333333 0.55555556]
 [1.         0.66666667]]
```

```
Actual Output:
[[0.92]
 [0.86]
 [0.89]]

Predicted Output:
[[0.5855257 ]
 [0.58103597]
 [0.60227742]]

Loss:
0.09082609014210542

Input:
[[0.66666667 1.        ]
 [0.33333333 0.55555556]
 [1.         0.66666667]]

Actual Output:
[[0.92]
 [0.86]
 [0.89]]

Predicted Output:
[[0.89499911]
 [0.86577542]
 [0.91033262]]

Loss:
0.0003572717610598643

Input:
[[0.66666667 1.        ]
 [0.33333333 0.55555556]
 [1.         0.66666667]]

Actual Output:
[[0.92]
 [0.86]
 [0.89]]

Predicted Output:
[[0.89500363]
 [0.86577342]
 [0.91032839]]

Loss:
0.0003571313973137464
```

## 6. Write a program to implement the naïve Bayesian classifier for a sample training data set stored as a .CSV file. Compute the accuracy of the classifier, considering few test data sets.

```python
# import necessary libarities
import pandas as pd
from sklearn import tree
from sklearn.preprocessing import LabelEncoder
from sklearn.naive_bayes import GaussianNB

# load data from CSV
data = pd.read_csv('C:/Users/mbacc/Documents/dataset/tennisdata.csv')
print("THe first 5 values of data is :\n",data.head())

# obtain Train data and Train output
X = data.iloc[:,:-1]
print("\nThe First 5 values of train data is\n",X.head())
y = data.iloc[:,-1]
print("\nThe first 5 values of Train output is\n",y.head())

# Convert then in numbers
le_outlook = LabelEncoder()
X.Outlook = le_outlook.fit_transform(X.Outlook)

le_Temperature = LabelEncoder()
X.Temperature = le_Temperature.fit_transform(X.Temperature)

le_Humidity = LabelEncoder()
X.Humidity = le_Humidity.fit_transform(X.Humidity)

le_Windy = LabelEncoder()
X.Windy = le_Windy.fit_transform(X.Windy)

print("\nNow the Train data is :\n",X.head())

le_PlayTennis = LabelEncoder()
y = le_PlayTennis.fit_transform(y)
print("\nNow the Train output is\n",y)

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.20)

classifier = GaussianNB()
classifier.fit(X_train,y_train)

from sklearn.metrics import accuracy_score
print("Accuracy is:",accuracy_score(classifier.predict(X_test),y_test))
```

### Dataset

tennisdata.csv

| Outlook | Temperature | Humidity | Windy | PlayTennis |
|---|---|---|---|---|
| Sunny | Hot | High | FALSE | No |
| Sunny | Hot | High | TRUE | No |
| Overcast | Hot | High | FALSE | Yes |
| Rainy | Mild | High | FALSE | Yes |
| Rainy | Cool | Normal | FALSE | Yes |
| Rainy | Cool | Normal | TRUE | No |
| Overcast | Cool | Normal | TRUE | Yes |
| Sunny | Mild | High | FALSE | No |
| Sunny | Cool | Normal | FALSE | Yes |
| Rainy | Mild | Normal | FALSE | Yes |
| Sunny | Mild | Normal | TRUE | Yes |
| Overcast | Mild | High | TRUE | Yes |
| Overcast | Hot | Normal | FALSE | Yes |
| Rainy | Mild | High | TRUE | No |

### OUTPUT

```
The first 5 values of data is :
    Outlook Temperature Humidity  Windy PlayTennis
0    Sunny         Hot     High  False         No
1    Sunny         Hot     High   True         No
2 Overcast         Hot     High  False        Yes
3    Rainy        Mild     High  False        Yes
4    Rainy        Cool   Normal  False        Yes

The First 5 values of train data is
    Outlook Temperature Humidity  Windy
0    Sunny         Hot     High  False
1    Sunny         Hot     High   True
2 Overcast         Hot     High  False
3    Rainy        Mild     High  False
4    Rainy        Cool   Normal  False

The first 5 values of Train output is
 0     No
1     No
2    Yes
3    Yes
4    Yes
Name: PlayTennis, dtype: object
```

```
Now the Train data is :
    Outlook   Temperature   Humidity   Windy
0        2             1          0       0
1        2             1          0       1
2        0             1          0       0
3        1             2          0       0
4        1             0          1       0

Now the Train output is
  [0 0 1 1 1 0 1 0 1 1 1 1 1 0]
Accuracy is: 1.0
```

**7. Apply EM algorithm to cluster a set of data stored in a .CSV file. Use the same data set for clustering using *k*-Means algorithm. Compare the results of these two algorithms and comment on the quality of clustering. You can add Java/Python ML library classes/API in the program.**

```python
from sklearn.cluster import KMeans
from sklearn import preprocessing
from sklearn.mixture import GaussianMixture
from sklearn.datasets import load_iris
import sklearn.metrics as sm
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt

dataset=load_iris()
# print(dataset)

X=pd.DataFrame(dataset.data)
X.columns=['Sepal_Length','Sepal_Width','Petal_Length','Petal_Width']

y=pd.DataFrame(dataset.target)
y.columns=['Targets']
# print(X)

plt.figure(figsize=(14,7))
colormap=np.array(['red','lime','black'])

# REAL PLOT
plt.subplot(1,3,1)
plt.scatter(X.Petal_Length,X.Petal_Width,c=colormap[y.Targets],s=40)
plt.title('Real')

# K-PLOT
plt.subplot(1,3,2)
model=KMeans(n_clusters=3)
model.fit(X)
predY=np.choose(model.labels_,[0,1,2]).astype(np.int64)
plt.scatter(X.Petal_Length,X.Petal_Width,c=colormap[predY],s=40)
plt.title('KMeans')

# GMM PLOT
scaler=preprocessing.StandardScaler()
scaler.fit(X)
xsa=scaler.transform(X)
xs=pd.DataFrame(xsa,columns=X.columns)
```

```
gmm=GaussianMixture(n_components=3)
gmm.fit(xs)

y_cluster_gmm=gmm.predict(xs)
plt.subplot(1,3,3)
plt.scatter(X.Petal_Length,X.Petal_Width,c=colormap[y_cluster_gmm],s=40)
plt.title('GMM Classification')
```

**OUTPUT**

```
Text(0.5, 1.0, 'GMM Classification')
```

**8. Write a program to implement k-Nearest Neighbour algorithm to classify the iris data set. Print both correct and wrong predictions. Java/Python ML library classes can be used for this problem.**

```python
from sklearn.datasets import load_iris
from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
import numpy as np

dataset=load_iris()
#print(dataset)
X_train,X_test,y_train,y_test=train_test_split(dataset["data"],dataset["target"],random_state=0)

kn=KNeighborsClassifier(n_neighbors=1)
kn.fit(X_train,y_train)

for i in range(len(X_test)):

  x=X_test[i]
  x_new=np.array([x])
  prediction=kn.predict(x_new)

print("TARGET=",y_test[i],dataset["target_names"][y_test[i]],"PREDICTED=",prediction,
dataset["target_names"][prediction])
print(kn.score(X_test,y_test))
```

**OUTPUT**

```
TARGET= 1 versicolor PREDICTED= [2] ['virginica']
0.9736842105263158
```

## 9. Implement the non-parametric Locally Weighted Regression algorithm in order to fit data points. Select appropriate data set for your experiment and draw graphs.

```python
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np

def kernel(point,xmat, k):
    m,n =np.shape(xmat)
    weights = np.mat(np.eye((m)))
    for j in range(m):
        diff = point - X[j]
    weights[j,j] = np.exp(diff*diff.T/(-2.0*k**2))
    return weights

def localWeight(point,xmat,ymat,k):
    wei = kernel(point,xmat,k)
    W=(X.T*(wei*X)).I*(X.T*(wei*ymat.T))
    return W

def localWeightRegression(xmat,ymat,k):
    m,n = np.shape(xmat)
    ypred = np.zeros(m)
    for i in range(m):
        ypred[i] = xmat[i]*localWeight(xmat[i],xmat,ymat,k)
    return ypred
# load data points
data = pd.read_csv('C:/Users/mbacc/Documents/dataset/tips.csv')
bill = np.array(data.total_bill)
tip = np.array(data.tip)
#preparing and add 1 in bill
mbill =np.mat(bill)
mtip = np.mat(tip)
m= np.shape(mbill)[1]
one = np.mat(np.ones(m))
X= np.hstack((one.T,mbill.T))
print(X.shape)
#set k here
ypred = localWeightRegression(X,mtip,0.5)
SortIndex = X[:,1].argsort(0)
xsort = X[SortIndex][:,0]
fig=plt.figure()
ax=fig.add_subplot(1,1,1)
ax.scatter(bill,tip,color='green')
```
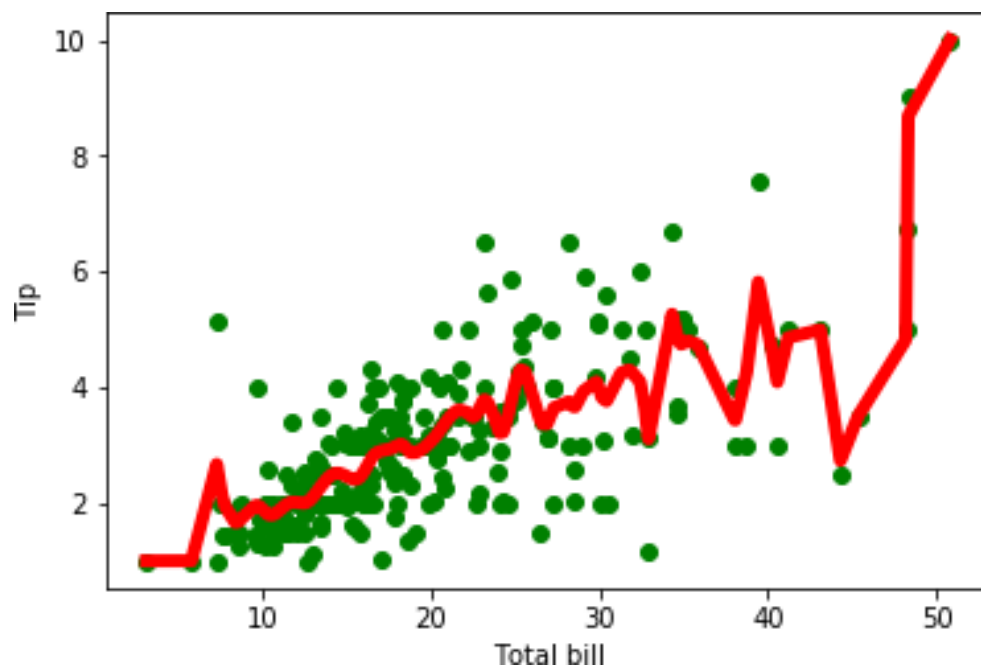
ax.plot(xsort[:,1],ypred[SortIndex],color='red',linewidth=5)
plt.xlabel('Total bill')
plt.ylabel('Tip')
plt.show();

## Dataset

*Add Tips.csv (256 rows)*

## OUTPUT

(244, 2)

Additional Programs :

1. **Write a program to demonstrate the working of the decision tree based ID3 algorithm. Use an appropriate data set for building the decision tree and apply this knowledge to classify a new sample.**

```
import pandas as pd
import numpy as np
#Import the dataset and define the feature as well as the target datasets / columns#
dataset = pd.read_csv('playtennis.csv',
            names=['outlook','temperature','humidity','wind','class',])
#Import all columns omitting the fist which consists the names of the animals
#We drop the animal names since this is not a good feature to split the data on

attributes =('Outlook','Temperature','Humidity','Wind','PlayTennis')
def entropy(target_col):
    """
    Calculate the entropy of a dataset.
    The only parameter of this function is the target_col parameter which specifies the target
column
    """
    elements,counts = np.unique(target_col,return_counts = True)

    entropy = np.sum([(-counts[i]/np.sum(counts))*np.log2(counts[i]/np.sum(counts)) for i in
range(len(elements))])
    print('Entropy =', entropy)
    return entropy

def InfoGain(data,split_attribute_name,target_name="class"):
    #Calculate the entropy of the total dataset
    total_entropy = entropy(data[target_name])

    ##Calculate the entropy of the dataset

    #Calculate the values and the corresponding counts for the split attribute
    vals,counts= np.unique(data[split_attribute_name],return_counts=True)

    #Calculate the weighted entropy
    Weighted_Entropy =
np.sum([(counts[i]/np.sum(counts))*entropy(data.where(data[split_attribute_name]==vals[i]).dr
opna()[target_name]) for i in range(len(vals))])

    #Calculate the information gain
    Information_Gain = total_entropy - Weighted_Entropy
    return Information_Gain
```

```
def ID3(data,originaldata,features,target_attribute_name="class",parent_node_class = None):
#Define the stopping criteria --> If one of this is satisfied, we want to return a leaf node#

    #If all target_values have the same value, return this value

    if len(np.unique(data[target_attribute_name])) <= 1:
        return np.unique(data[target_attribute_name])[0]

    #If the dataset is empty, return the mode target feature value in the original dataset
    elif len(data)==0:
        return
np.unique(originaldata[target_attribute_name])[np.argmax(np.unique(originaldata[target_attribut
e_name],return_counts=True)[1])]

    elif len(features) ==0:
        return parent_node_class

    #If none of the above holds true, grow the tree!

    else:
        #Set the default value for this node --> The mode target feature value of the current node
        parent_node_class =
np.unique(data[target_attribute_name])[np.argmax(np.unique(data[target_attribute_name],return
_counts=True)[1])]

        #Select the feature which best splits the dataset
        item_values = [InfoGain(data,feature,target_attribute_name) for feature in features] #Return
the information gain values for the features in the dataset
        best_feature_index = np.argmax(item_values)
        best_feature = features[best_feature_index]

        #Create the tree structure. The root gets the name of the feature (best_feature) with the
maximum information
        #gain in the first run
        tree = {best_feature:{}}

        #Remove the feature with the best inforamtion gain from the feature space
        features = [i for i in features if i != best_feature]

        #Grow a branch under the root node for each possible value of the root node feature

        for value in np.unique(data[best_feature]):
            value = value
            #Split the dataset along the value of the feature with the largest information gain and
therwith create sub_datasets
```

```
            sub_data = data.where(data[best_feature] == value).dropna()

            #Call the ID3 algorithm for each of those sub_datasets with the new parameters --> Here
    the recursion comes in!
            subtree = ID3(sub_data,dataset,features,target_attribute_name,parent_node_class)

            #Add the sub tree, grown from the sub_dataset to the tree under the root node
            tree[best_feature][value] = subtree

        return(tree)

    def predict(query,tree,default = 1):

        #1.
        for key in list(query.keys()):
            if key in list(tree.keys()):
                #2.
                try:
                    result = tree[key][query[key]]
                except:
                    return default

                #3.
                result = tree[key][query[key]]
                #4.
                if isinstance(result,dict):
                    return predict(query,result)
                else:
                    return result

    def train_test_split(dataset):
        training_data = dataset.iloc[1:15].reset_index(drop=True)
        #We drop the index respectively relabel the index
        #starting form 0, because we do not want to run into errors regarding the row labels / indexes
        #testing_data = dataset.iloc[10:].reset_index(drop=True)
        return training_data  #,testing_data

    def test(data,tree):
        #Create new query instances by simply removing the target feature column from the original
    dataset and
        #convert it to a dictionary
        queries = data.iloc[:,:-1].to_dict(orient = "records")

        #Create a empty DataFrame in whose columns the prediction of the tree are stored
        predicted = pd.DataFrame(columns=["predicted"])
```

#Calculate the prediction accuracy
for i in range(len(data)):
    predicted.loc[i,"predicted"] = predict(queries[i],tree,1.0)

    print('The prediction accuracy is: ',(np.sum(predicted["predicted"] ==
data["class"])/len(data))*100,'%')
"""
Train the tree, Print the tree and predict the accuracy
"""
XX = train_test_split(dataset)
training_data=XX
#testing_data=XX[1]
tree = ID3(training_data,training_data,training_data.columns[:-1])
print(' Display Tree',tree)
print('len=',len(training_data))
test(training_data,tree)

## Dataset

tennis.csv

| Outlook | Temperature | Humidity | Windy | PlayTennis |
|---------|-------------|----------|-------|------------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rainy | Mild | High | Weak | Yes |
| Rainy | Cool | Normal | Weak | Yes |
| Rainy | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rainy | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rainy | Mild | High | Strong | No |

**OUTPUT**

```
Entropy = 0.9402859586706311
Entropy = 0.0
Entropy = 0.9709505944546686
Entropy = 0.9709505944546686
Entropy = 0.9402859586706311
Entropy = 0.8112781244591328
Entropy = 1.0
Entropy = 0.9182958340544896
Entropy = 0.9402859586706311
Entropy = 0.9852281360342515
Entropy = 0.5916727785823275
Entropy = 0.9402859586706311
Entropy = 1.0
Entropy = 0.8112781244591328
Entropy = 0.9709505944546686
Entropy = 1.0
Entropy = 0.9182958340544896
Entropy = 0.9709505944546686
Entropy = 1.0
Entropy = 0.9182958340544896
Entropy = 0.9709505944546686
Entropy = 0.0
Entropy = 0.0
Entropy = 0.9709505944546686
Entropy = 0.0
Entropy = 0.0
Entropy = 1.0
Entropy = 0.9709505944546686
Entropy = 0.0
Entropy = 0.0
Entropy = 0.9709505944546686
Entropy = 1.0
Entropy = 0.9182958340544896
 Display Tree {'Outlook': {'Overcast': 'Yes', 'Rainy': {'Wind': {'Strong': 'No', 'Weak':
'Yes'}}, 'Sunny': {'Humidity': {'High': 'No', 'Normal': 'Yes'}}}}
len= 14
The prediction accuracy is:  100.0 %
```

## 2. Build an Artificial Neural Network by implementing the Backpropagation algorithm and test the same using appropriate data sets.

```python
from math import exp
from random import seed
from random import random

# Initialize a network
def initialize_network(n_inputs, n_hidden, n_outputs):
        network = list()
        hidden_layer = [{'weights':[random() for i in range(n_inputs + 1)]} for i in
range(n_hidden)]
        network.append(hidden_layer)
        output_layer = [{'weights':[random() for i in range(n_hidden + 1)]} for i in
range(n_outputs)]
        network.append(output_layer)
        return network
# Calculate neuron activation for an input
def activate(weights, inputs):
        activation = weights[-1]
        for i in range(len(weights)-1):
                activation += weights[i] * inputs[i]
        return activation
# Transfer neuron activation
def transfer(activation):
        return 1.0 / (1.0 + exp(-activation))

# Forward propagate input to a network output
def forward_propagate(network, row):
        inputs = row
        for layer in network:
                new_inputs =  []
                for neuron in layer:
                        activation = activate(neuron['weights'], inputs)
                        neuron['output'] = transfer(activation)
                        new_inputs.append(neuron['output'])
                inputs = new_inputs
        return inputs
# Calculate the derivative of an neuron output
def transfer_derivative(output):
        return output * (1.0 - output)

# Backpropagate error and store in neurons
def backward_propagate_error(network, expected):
```

```
for i in reversed(range(len(network))):
            layer = network[i]
            errors = list()
            if i != len(network)-1:
                    for j in range(len(layer)):
                            error = 0.0
                            for neuron in network[i + 1]:
                                    error += (neuron['weights'][j] * neuron['delta'])
                            errors.append(error)
            else:
                    for j in range(len(layer)):
                            neuron = layer[j]
                            errors.append(expected[j] - neuron['output'])
            for j in range(len(layer)):
                    neuron = layer[j]
                    neuron['delta'] = errors[j] * transfer_derivative(neuron['output'])
# Update network weights with error
def update_weights(network, row, l_rate):
        for i in range(len(network)):
                inputs = row[:-1]
                if i != 0:
                        inputs = [neuron['output'] for neuron in network[i - 1]]
                for neuron in network[i]:
                        for j in range(len(inputs)):
                                neuron['weights'][j] += l_rate * neuron['delta'] * inputs[j]
                        neuron['weights'][-1] += l_rate * neuron['delta']
# Train a network for a fixed number of epochs
def train_network(network, train, l_rate, n_epoch, n_outputs):
        for epoch in range(n_epoch):
                sum_error = 0
                for row in train:
                        outputs = forward_propagate(network, row)
                        expected = [0 for i in range(n_outputs)]
                        expected[row[-1]] = 1
                        sum_error += sum([(expected[i]-outputs[i])**2 for i in
range(len(expected))])
                        backward_propagate_error(network, expected)
                        update_weights(network, row, l_rate)
                print('>epoch=%d, lrate=%.3f, error=%.3f' % (epoch, l_rate, sum_error))

# Test training backprop algorithm
seed(1)

dataset = [[2.7810836,2.550537003,0],
        [1.465489372,2.362125076,0],
```

```
        [3.396561688,4.400293529,0],
        [1.38807019,1.850220317,0],
        [3.06407232,3.005305973,0],
        [7.627531214,2.759262235,1],
        [5.332441248,2.088626775,1],
        [6.922596716,1.77106367,1],
        [8.675418651,-0.242068655,1],
        [7.673756466,3.508563011,1]]
n_inputs = len(dataset[0]) - 1
n_outputs = len(set([row[-1] for row in dataset]))
network = initialize_network(n_inputs, 2, n_outputs)
train_network(network, dataset, 0.5, 20, n_outputs)
for layer in network:
        print(layer)
```

**OUTPUT**
>epoch=0, lrate=0.500, error=6.350
>epoch=1, lrate=0.500, error=5.531
>epoch=2, lrate=0.500, error=5.221
>epoch=3, lrate=0.500, error=4.951
>epoch=4, lrate=0.500, error=4.519
>epoch=5, lrate=0.500, error=4.173
>epoch=6, lrate=0.500, error=3.835
>epoch=7, lrate=0.500, error=3.506
>epoch=8, lrate=0.500, error=3.192
>epoch=9, lrate=0.500, error=2.898
>epoch=10, lrate=0.500, error=2.626
>epoch=11, lrate=0.500, error=2.377
>epoch=12, lrate=0.500, error=2.153
>epoch=13, lrate=0.500, error=1.953
>epoch=14, lrate=0.500, error=1.774
>epoch=15, lrate=0.500, error=1.614
>epoch=16, lrate=0.500, error=1.472
>epoch=17, lrate=0.500, error=1.346
>epoch=18, lrate=0.500, error=1.233
>epoch=19, lrate=0.500, error=1.132
[{'weights': [-1.4688375095432327, 1.850887325439514, 1.0858178629550297], 'output': 0.029980305604426185, 'delta': -0.005946604162323625}, {'weights': [0.37711098142462157, -0.0625909894552989, 0.2765123702642716], 'output': 0.9456229000211323, 'delta': 0.0026279652850863837}]
[{'weights': [2.515394649397849, -0.3391927502445985, -0.9671565426390275], 'output': 0.23648794202357587, 'delta': -0.04270059278364587}, {'weights': [-2.5584149848484263, 1.0036422106209202, 0.42383086467582715], 'output': 0.7790535202438367, 'delta': 0.03803132596437354}]

## 3. Write a program to implement the naïve Bayesian classifier for a sample training data set stored as a .CSV file. Compute the accuracy of the classifier, considering few test data sets.

```python
print("\nNaive Bayes Classifier for concept learning problem")
import csv
#import random
import math
#import operator
def safe_div(x,y):
    if y == 0:
        return 0
    return x / y

def loadCsv(filename):
        lines = csv.reader(open(filename))
        dataset = list(lines)
        for i in range(len(dataset)):
                dataset[i] = [float(x) for x in dataset[i]]
        return dataset

def splitDataset(dataset, splitRatio):
        trainSize = int(len(dataset) * splitRatio)
        trainSet = []
        copy = list(dataset)
        i=0
        while len(trainSet) < trainSize:
                #index = random.randrange(len(copy))

                trainSet.append(copy.pop(i))
        return [trainSet, copy]

def separateByClass(dataset):
        separated = {}
        for i in range(len(dataset)):
                vector = dataset[i]
                if (vector[-1] not in separated):
                        separated[vector[-1]] = []
                separated[vector[-1]].append(vector)
        return separated

def mean(numbers):
        return safe_div(sum(numbers),float(len(numbers)))

def stdev(numbers):
```

```
        avg = mean(numbers)
        variance = safe_div(sum([pow(x-avg,2) for x in numbers]),float(len(numbers)-1))
        return math.sqrt(variance)


def summarize(dataset):
        summaries = [(mean(attribute), stdev(attribute)) for attribute in zip(*dataset)]
        del summaries[-1]
        return summaries


def summarizeByClass(dataset):
        separated = separateByClass(dataset)
        summaries = {}
        for classValue, instances in separated.items():
                summaries[classValue] = summarize(instances)
        return summaries


def calculateProbability(x, mean, stdev):
        exponent = math.exp(-safe_div(math.pow(x-mean,2),(2*math.pow(stdev,2))))
        final = safe_div(1 , (math.sqrt(2*math.pi) * stdev)) * exponent
        return final


def calculateClassProbabilities(summaries, inputVector):
        probabilities = {}
        for classValue, classSummaries in summaries.items():
                probabilities[classValue] = 1
                for i in range(len(classSummaries)):
                        mean, stdev = classSummaries[i]
                        x = inputVector[i]
                        probabilities[classValue] *= calculateProbability(x, mean, stdev)
        return probabilities


def predict(summaries, inputVector):
        probabilities = calculateClassProbabilities(summaries, inputVector)
        bestLabel, bestProb = None, -1
        for classValue, probability in probabilities.items():
                if bestLabel is None or probability > bestProb:
                        bestProb = probability
                        bestLabel = classValue
        return bestLabel


def getPredictions(summaries, testSet):
        predictions = []
        for i in range(len(testSet)):
                result = predict(summaries, testSet[i])
                predictions.append(result)
        return predictions
```

```python
def getAccuracy(testSet, predictions):
        correct = 0
        for i in range(len(testSet)):
                if testSet[i][-1] == predictions[i]:
                        correct += 1
        accuracy = safe_div(correct,float(len(testSet))) * 100.0
        return accuracy

def main():
        filename = 'ConceptLearning.csv'
        splitRatio = 0.9
        dataset = loadCsv(filename)
        trainingSet, testSet = splitDataset(dataset, splitRatio)
        print('Split {0} rows into'.format(len(dataset)))
        print('Number of Training data: ' + (repr(len(trainingSet))))
        print('Number of Test Data: ' + (repr(len(testSet))))
        print("\nThe values assumed for the concept learning attributes are\n")
        print("OUTLOOK=> Sunny=1 Overcast=2 Rain=3\nTEMPERATURE=> Hot=1 Mild=2
Cool=3\nHUMIDITY=> High=1 Normal=2\nWIND=> Weak=1 Strong=2")
        print("TARGET CONCEPT:PLAY TENNIS=> Yes=10 No=5")
        print("\nThe Training set are:")
        for x in trainingSet:
                print(x)
        print("\nThe Test data set are:")
        for x in testSet:
                print(x)
        print("\n")
        # prepare model
        summaries = summarizeByClass(trainingSet)
        # test model
        predictions = getPredictions(summaries, testSet)
        actual = []
        for i in range(len(testSet)):
                vector = testSet[i]
                actual.append(vector[-1])
        # Since there are five attribute values, each attribute constitutes to 20% accuracy. So if all
attributes match with predictions then 100% accuracy
        print('Actual values: {0}%'.format(actual))
        print('Predictions: {0}%'.format(predictions))
        accuracy = getAccuracy(testSet, predictions)
        print('Accuracy: {0}%'.format(accuracy))

main()
```

### Dataset

ConceptLearning.csv

| 1 | 1 | 1 | 1 | 5 |
|---|---|---|---|---|
| 1 | 1 | 1 | 2 | 5 |
| 2 | 1 | 1 | 2 | 10 |
| 3 | 2 | 1 | 1 | 10 |
| 3 | 3 | 2 | 1 | 10 |
| 3 | 3 | 2 | 2 | 5 |
| 2 | 3 | 2 | 2 | 10 |
| 1 | 2 | 1 | 1 | 5 |
| 1 | 3 | 2 | 1 | 10 |
| 3 | 2 | 2 | 2 | 10 |
| 1 | 2 | 2 | 2 | 10 |
| 2 | 2 | 1 | 2 | 10 |
| 2 | 1 | 2 | 1 | 10 |
| 3 | 2 | 1 | 2 | 5 |
| 1 | 2 | 1 | 2 | 5 |
| 1 | 2 | 1 | 2 | 5 |

### OUTPUT

Naive Bayes Classifier for concept learning problem
Split 16 rows into
Number of Training data: 14
Number of Test Data: 2

The values assumed for the concept learning attributes are

OUTLOOK=> Sunny=1 Overcast=2 Rain=3
TEMPERATURE=> Hot=1 Mild=2 Cool=3
HUMIDITY=> High=1 Normal=2
WIND=> Weak=1 Strong=2
TARGET CONCEPT:PLAY TENNIS=> Yes=10 No=5

The Training set are:
[1.0, 1.0, 1.0, 1.0, 5.0]
[1.0, 1.0, 1.0, 2.0, 5.0]
[2.0, 1.0, 1.0, 2.0, 10.0]
[3.0, 2.0, 1.0, 1.0, 10.0]
[3.0, 3.0, 2.0, 1.0, 10.0]
[3.0, 3.0, 2.0, 2.0, 5.0]

[2.0, 3.0, 2.0, 2.0, 10.0]
[1.0, 2.0, 1.0, 1.0, 5.0]
[1.0, 3.0, 2.0, 1.0, 10.0]
[3.0, 2.0, 2.0, 2.0, 10.0]
[1.0, 2.0, 2.0, 2.0, 10.0]
[2.0, 2.0, 1.0, 2.0, 10.0]
[2.0, 1.0, 2.0, 1.0, 10.0]
[3.0, 2.0, 1.0, 2.0, 5.0]

The Test data set are:
[1.0, 2.0, 1.0, 2.0, 5.0]
[1.0, 2.0, 1.0, 2.0, 5.0]


Actual values: [5.0, 5.0]%
Predictions: [5.0, 5.0]%
Accuracy: 100.0%

**4. Apply EM algorithm to cluster a set of data stored in a .CSV file. Use the same data set for clustering using *k*-Means algorithm. Compare the results of these two algorithms and comment on the quality of clustering. You can add Java/Python ML library classes/API in the program.**

```
import matplotlib.pyplot as plt
from sklearn import datasets
from sklearn.cluster import KMeans
import pandas as pd
import numpy as np


# import some data to play with
iris = datasets.load_iris()
X = pd.DataFrame(iris.data)
X.columns=['Sepal_Length','Sepal_Width','Petal_Length','Petal_Width']
y = pd.DataFrame(iris.target)
y.columns = ['Targets']


 # Build the K Means Model
model = KMeans(n_clusters=3)
model.fit(X)
# model.labels_ : Gives cluster no for which samples belongs to
# # Visualise the clustering results
plt.figure(figsize=(14,14))
colormap = np.array(['red', 'lime', 'black'])


 # Plot the Original Classifications using Petal features
plt.subplot(2, 2, 1)
plt.scatter(X.Petal_Length, X.Petal_Width, c=colormap[y.Targets], s=40)
plt.title('Real Clusters')
plt.xlabel('Petal Length')
plt.ylabel('Petal Width')
```

# Plot the Models Classifications

```
plt.subplot(2, 2, 2)

plt.scatter(X.Petal_Length, X.Petal_Width, c=colormap[model.labels_], s=40)

plt.title('K-Means Clustering')

plt.xlabel('Petal Length')

plt.ylabel('Petal Width')
```

# General EM for GMM

```
from sklearn import preprocessing
```

# transform your data such that its distribution will have a
# mean value 0 and standard deviation of 1.

```
scaler = preprocessing.StandardScaler()

scaler.fit(X)

xsa = scaler.transform(X)

xs = pd.DataFrame(xsa, columns = X.columns)

from sklearn.mixture import GaussianMixture

gmm = GaussianMixture(n_components=3)

gmm.fit(xs)

gmm_y = gmm.predict(xs)

plt.subplot(2, 2, 3)

plt.scatter(X.Petal_Length, X.Petal_Width, c=colormap[gmm_y], s=40)

plt.title('GMM Clustering')

plt.xlabel('Petal Length')

plt.ylabel('Petal Width')

print('Observation: The GMM using EM algorithm based clustering matched the true
labels more closely than the Kmeans.')
```

**OUTPUT :**

Observation: The GMM using EM algorithm based clustering matched the true labels more closely than the Kmeans.

## 5. Write a program to implement *k*-Nearest Neighbour algorithm to classify the iris data set. Print both correct and wrong predictions. Java/Python ML library classes can be used for this problem.

```
import csv import
random import
math import
operator
def loadDataset(filename, split, trainingSet=[], testSet=[]):
      with open(filename) as csvfile:
                 lines =
                 csv.reader(csvfile)
                 dataset = list(lines)
                 for x in
                    range(len(dataset)-1):
                    for y in range(4):
                       dataset[x][y] =
                    float(dataset[x][y])if
                    random.random() < split:
                       trainingSet.append(dataset
                 [x])else:
                    testSet.append(dataset[x])
def euclideanDistance(instance1, instance2, length):
    distance = 0
    for x in range(length):
          distance += pow((instance1[x] - instance2[x]), 2)
          return math.sqrt(distance)

def getNeighbors(trainingSet, testInstance, k):
      distances = []
      length = len(testInstance)-1
      for x in range(len(trainingSet)):
            dist = euclideanDistance(testInstance, trainingSet[x], length)
            distances.append((trainingSet[x], dist))
      distances.sort(key=operator.itemgetter(1))
      neighbors = []
      for x in range(k):
            neighbors.append(distances[x][0]
      return neighbors
def getResponse(neighbors):
      classVotes = {}
      for x in range(len(neighbors)):
            response = neighbors[x][-1]
            if response in classVotes:
                  classVotes[response] += 1
            else:
```

```
                        classVotes[response] = 1
          sortedVotes = sorted(classVotes.items(), key=operator.itemgetter(1), reverse=True)


                return sortedVotes[0][0]
def getAccuracy(testSet, predictions):
        correct = 0
        for x in range(len(testSet)):
                if testSet[x][-1] == predictions[x]:
                        correct += 1

        return (correct/float(len(testSet))) * 100.0
def main():
    # prepare data
    trainingSet=[]
    testSet=[]
    split = 0.7
    loadDataset('iris_data.csv', split, trainingSet, testSet)
    print ('\n Number of Training data: ' + (repr(len(trainingSet))))print
    (' Number of Test Data: ' + (repr(len(testSet))))
    # generate predictions
    predictions=[]
    k = 3
    print('\n The predictions are: ')
    for x in range(len(testSet)):
            neighbors = getNeighbors(trainingSet, testSet[x], k)
            result = getResponse(neighbors)
            predictions.append(result)
            print(' predicted=' + repr(result) + ', actual=' + repr(testSet[x][-1]))
    accuracy = getAccuracy(testSet, predictions)
    print('\n The Accuracy is: ' + repr(accuracy) + '%')
main()
```

**OUTPUT :**

Number of Training data: 101 Number of Test Data: 48


The predictions are:
predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'
 predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'

predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-setosa', actual='Iris-setosa'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-virginica', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-versicolor', actual='Iris-versicolor'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'
predicted='Iris-virginica', actual='Iris-virginica'

The Accuracy is: 97.91666666666666%

## 6.  Implement the non-parametric Locally Weighted Regression algorithm in order to fit data points. Select appropriate data set for your experiment and draw graphs.

```python
from math import ceil
import numpy as np from
scipy import linalg

def lowess(x, y, f, iterations):
    n = len(x)
    r = int(ceil(f * n))
    h = [np.sort(np.abs(x - x[i]))[r] for i in range(n)]
    w = np.clip(np.abs((x[:, None] - x[None, :]) / h), 0.0, 1.0)
    w = (1 - w ** 3) ** 3
    yest = np.zeros(n)
    delta =np.ones(n)
    for iteration in range(iterations):
        for i in range(n):
            weights = delta * w[:, i]
             b = np.array([np.sum(weights * y), np.sum(weights * y * x)])
             A = np.array([[np.sum(weights), np.sum(weights * x)],[np.sum(weights * x),
np.sum(weights * x * x)]])
            beta = linalg.solve(A, b)
            yest[i] = beta[0] + beta[1] * x[i]

             residuals = y - yest
             s = np.median(np.abs(residuals))
             delta = np.clip(residuals / (6.0 * s), -1, 1)
             delta = (1 - delta ** 2) ** 2

        return yest

import math
n = 100

= np.linspace(0, 2 * math.pi, n)
y = np.sin(x) + 0.3 * np.random.randn(n)
f =0.25
iterations=3
yest = lowess(x, y, f, iterations)

import matplotlib.pyplot as plt
plt.plot(x,y,"r.")
        plt.plot(x,yest,"b-")
```

## VIVA Questions

1. What is machine learning?
2. Define supervised learning
3. Define unsupervised learning
4. Define semi supervised learning
5. Define reinforcement learning
6. What do you mean by hypotheses
7. What is classification
8. What is clustering
9. Define precision, accuracy and recall
10. Define entropy
11. Define regression
12. How Knn is different from k-means clustering
13. What is concept learning
14. Define specific boundary and general boundary
15. Define target function
16. Define decision tree
17. What is ANN
18. Explain gradient descent approximation
19. State Bayes theorem
20. Define Bayesian belief networks
21. Differentiate hard and soft clustering
22. Define variance
23. What is inductive machine learning?
24. Why K nearest neighbor algorithm is lazy learning algorithm
25. Why naïve Bayes is naïve
26. Mention classification algorithms
27. Define pruning
28. Differentiate Clustering and classification
29. Mention clustering algorithms
30. Define Bias