

EXPERIMENT -1

Aim: Implement and demonstrate the FIND-S algorithm for finding the most specific hypothesis based on a given set of training data samples. Read the training data from a .CSV file.

Description:

- The concept learning approach in machine learning, can be formulated as “Problem of searching through a predefined space of potential hypotheses for the hypothesis that best fits the training examples”.
- Find-S algorithm for concept learning is one of the most basic algorithms of machine learning.

Algorithm:

1. Initialize h to the most specific hypothesis in H
2. For each positive training instance x
 - For each attribute constraint a_i in h :
 - If the constraint a_i in h is satisfied by x then do nothing
 - Else replace a_i in h by the next more general constraint that is satisfied by x
3. Output hypothesis h

Source code:

```

import io
import csv
import pandas as pd
num_attributes = 6
a = []
print("The given dataset is")
with open('E:\MLDS\First.csv', 'r') as csv_file:
    reader = csv.reader(csv_file)
    for row in reader:
        a.append(row)
    print(row)
print("The initial hypothesis is")
hypothesis = [0]*num_attributes

```

```

print (hypothesis)
# Take the first 'yes' row attributes into hypothesis
for j in range (0, num_attributes):
    hypothesis [j] = a [1] [j]
print (hypothesis)
print ("Finds: Finding maximal specific hypothesis")
for i in range (1, len (a)):
    if a[i] [num_attributes] == 'Yes' or a[i] [num_attributes] == 'yes':
        for j in range (0, num_attributes):
            if a[i] [j] != hypothesis [j]:
                hypothesis [j] = '?'
            else:
                hypothesis [j] = a[i] [j]
print ("For Training instance No: {} the hypothesis is {}".format(i), hypothesis).

```

Output: The given dataset is:

['sky', 'air temp', 'humidity', 'wind', 'water', 'Forecast', 'Enjoy sport']

['sunny', 'warm', 'normal', 'strong', 'warm', 'same', 'yes']

['sunny', 'warm', 'high', 'strong', 'warm', 'same', 'yes']

['Sunny', 'cold', 'High', 'strong', 'warm', 'change', 'No']

['sunny', 'warm', 'high', 'strong', 'cool', 'change', 'yes']

The initial hypothesis is

['o', 'o', 'o', 'o', 'o', 'o']

['sunny', 'warm', 'normal', 'strong', 'warm', 'same']

Find S: find maximal specific hypothesis.

For Training instance No: 1 the hypothesis is

['sunny', 'warm', 'normal', 'strong', 'warm', 'same']

for Training instance No: 2 the hypothesis is

['sunny', 'warm', '?', 'strong', 'warm', 'same']

for Training instance NO: 3 the hypothesis is

['sunny', 'warm', '?', 'strong', 'warm', 'same']

for Training instance No: 4 the hypothesis is

['sunny', 'warm', '?', 'strong', '?', '?']

Viva Voce Questions

1) What is machine learning?

Machine learning is an application of AI. Machine learning uses data to train and find accurate results. The capability of a machine to imitate intelligent human behaviour

2) Differentiate between Supervised and Unsupervised learning.

- * Supervised learning allows you to collect data (or) produce a data output from the previous experience
- * unsupervised learning helps you to find all kind of unknown patterns in data.

3) What is hypothesis in Find-s algorithm?

The find-s algorithm finds the most specific hypothesis that fits all the positive examples.

4) Which algorithm can be used for identifying suitable hypothesis?

The find-s algorithm can be used for identifying suitable hypothesis.

5) Why find-s algorithm does not consider negative example?

Find-s algorithm ignore the negative examples as it will give negative end result.

```

Source code: import os
import csv # upload
import pandas as pd
df = pd.read_csv("week2.csv")
print(df)

def entropy(probs):
    import math
    return sum(-prob * math.log(prob, 2) for prob in probs)

def entropy_of_list(a_list):
    from collections import Counter
    cnt = Counter(x for x in a_list)
    num_instances = len(a_list)
    probs = [x / num_instances for x in cnt.values()]
    return entropy(probs)

total_entropy = entropy_of_list(df['play Tennis'])

print(total_entropy)

def information_gain(df, split_attribute_name, target_attribute_name,
                     trace=0):
    df_split = df.groupby(split_attribute_name)
    for name, group in df_split:
        nobs = len(group.index) * 1.0
        df_agg_ent = df_split.agg({target_attribute_name:
            [entropy_of_list, lambda x: len(x) / nobs]})(target_attribute_name)
        avg_info = sum(df_agg_ent['entropy_of_list'] * df_agg_ent['lambda'])
        old_entropy = entropy_of_list(df[target_attribute_name])
        if trace:
            print("Old Entropy = ", old_entropy)
            print("Avg Info = ", avg_info)
            print("Info Gain = ", old_entropy - avg_info)
    return old_entropy - avg_info

def id3DT(df, target_attribute_name, attribute_names, default_class=None):
    from collections import Counter
    cnt = Counter(x for x in df[target_attribute_name])

```

```

if len(cnt) == 1:
    return next(iter(cnt))
elif dt.empty or (not attribute_names):
    return default_class
else:
    default_class = max(cnt.keys())
    gainz = [information_gain(df, attr, target_attribute_name) for attr in
             attribute_names]
    index_of_max = gainz.index(max(gainz))
    best_attr = attribute_names[index_of_max]
    tree = {best_attr: {}}
    remaining_attributes_names = [i for i in attribute_names
                                   if i != best_attr]
    for attr_val, data_subset in df.groupby(best_attr):
        subtree = id3DT(data_subset, target_attribute_name,
                         remaining_attributes_names, default_class)
        tree[best_attr][attr_val] = subtree
    return tree

attribute_names = list(df.columns)
attribute_names.remove('playTennis')

from pprint import pprint
tree = id3DT(df, 'playTennis', attribute_names)
print('The Resultant Decision Tree is')
pprint(tree)

attribute = next(iter(tree))
print("Best Attribute: ", attribute)
print("Tree keys: ", tree[attribute].keys())

def classify(instance, tree, default=None):
    attribute = next(iter(tree))
    print("Key: ", tree.keys())

```

```

print("Attribute", attribute)
if instance[attribute] in tree[attribute].keys():
    result = tree[attribute][instance[attribute]]
    print("Instance Attribute", instance[attribute], "Treekeys:", tree[attribute].keys())
    if isinstance(result, dict):
        return classify(instance, result)
    else:
        return result.
else:
    return default.

```

tree 1 = { 'outlook': ['Rainy', 'Sunny'], 'Temperature': ['mild', 'Hot', 'Humidity': ['Normal', 'High'], 'Wind': ['weak', 'strong']] }

df2 = pd.DataFrame(tree 1)

df2['predicted'] = df2.apply(classify, axis=1, args=(tree, 'No'))

print(df2) ~~output~~

| Day | Outlook | Temp | Humidity | Wind | Play Tennis |
|-----|----------|------|----------|--------|-------------|
| 0 | Sunny | Hot | High | Weak | No |
| 1 | Sunny | Hot | High | Strong | No |
| 2 | Overcast | Hot | High | Weak | Yes |
| 3 | Rain | Mild | High | Weak | Yes |
| 4 | Rain | Cool | Normal | Weak | Yes |
| 5 | Rain | Cool | Normal | Strong | No |
| 6 | Overcast | Cool | Normal | Strong | Yes |
| 7 | Sunny | Mild | High | Weak | No |
| 8 | Sunny | Cool | Normal | Weak | Yes |
| 9 | Rain | Mild | Normal | Weak | Yes |
| 10 | Sunny | Mild | Normal | Weak | Yes |
| 11 | Overcast | Mild | High | Strong | Yes |
| 12 | Overcast | Hot | Normal | Strong | Yes |
| 13 | Rain | Mild | High | Weak | Yes |

Entropy:- Entropy(S) = $\sum_{i=1}^n p_i \log_2 p_i$

$$\rightarrow \text{Entropy}(S) = -(p+) \log_2 (p+) - (p-) \log_2 (p-)$$

\rightarrow for out 100% we have 9 "Yes" and 5 "No"

$$\text{Entropy}(S) = -(9/14) + \log_2 (9/14) - (5/14) - \log_2 (5/14) = 0.940$$

$$\text{Temp: } S = [9, 5] = 0.94 \Rightarrow -(9/14) \log_2(9/14) - (5/14) \log_2(5/14)$$

$$\text{Output: } S_{Hot} = [2, 2] = 1.0 \Rightarrow -(2/14) \log_2(2/14) - (2/14) \log_2(2/14)$$

$$S_{Mild} = [4, 2] = 0.98 \Rightarrow -(4/14) \log_2(4/14) - (2/14) \log_2(2/14)$$

$$S_{Cool} = [3, 1] = 0.81 \Rightarrow -(3/14) \log_2(3/14) - (1/14) \log_2(1/14)$$

Humidity:

$$S = [9, 5] = 0.94 \Rightarrow -(9/14) \log_2(9/14) - (5/14) \log_2(5/14)$$

$$S_{High} = [3, 4] = \dots \Rightarrow -(3/14) \log_2(3/14) - (4/14) \log_2(4/14)$$

$$S_{Normal} [6, 1] = \dots \Rightarrow -(6/14) \log_2(6/14) - (1/14) \log_2(1/14)$$

Wind:

$$S = [9, 5] = 0.94 \Rightarrow -(9/14) \log_2(9/14) - (5/14) \log_2(5/14)$$

$$S_{Strong} [3, 3] = 1.0 \Rightarrow -(3/14) \log_2(3/14) - (3/14) \log_2(3/14)$$

$$S_{Weak} [6, 2] = 0.8 \Rightarrow -(6/14) \log_2(6/14) - (2/14) \log_2(2/14)$$

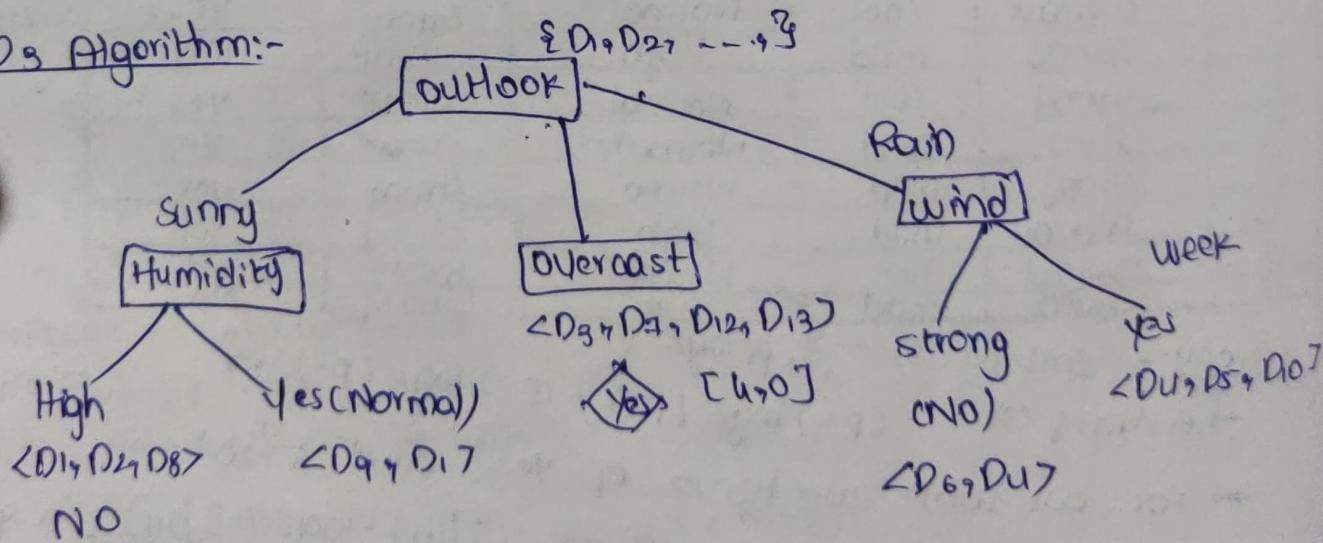
Information Gain:-

$$\rightarrow \text{Gain}(S, \text{outlook}) = \text{Entropy}(S) - \sum_{i \in \{\text{sunny, overcast, rain}\}} p_i \text{Entropy}(S|A_i)$$

csv | is) entropy.

$$\text{Gain}(S, \text{outlook}) = \frac{\text{entropy}(S) - 5/14 \text{entropy}(S, \text{sunny}) - 4/14 \text{entropy}(S, \text{overcast}) - 5/14 \text{entropy}(S, \text{rain})}{\text{entropy}(S)} =$$

$$\rightarrow \text{Gain}(S, \text{outlook}) = 0.94 - 4/14 (1.0) - 6/14 (0.9183) - 5/14 (0.9183) = 0.0289.$$

ID3 Algorithm:-

Viva Voce Questions

1) Define decision tree.

A decision tree is a type of supervised machine learning used to categorize (or) make predictions based on how a previous set of questions were answered.

2) What is the ID3 algorithm and how do we use it in decision tree regression?

- * ID₃ used a top-down greedy approach to build a decision tree.
- * It can be used to construct decision tree for regression by replacing information gain with standard deviation reduction.

3) How decision tree is constructed using ID3 algorithm?

The ID₃ algorithm is a classic machine learning algorithm used to build a decision tree.

The id₃ algorithm uses entropy and information gain to construct tree.

4) What are the advantages of ID3 algorithm?

Simple and easy to understand.

- * Fast training time.
- * Handles categorical data well.

5) What are the different types of decision trees?

- * ID₃
- * CART (classification and regression tree)
- * Random forest
- * C4.5 ID₃.

EXPERIMENT -3

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

Description:

- Given a dataset X, y , we attempt to find a linear model $h(x)$ that minimizes residual sum of squared errors. The solution is given by Normal equations.
- Linear model can only fit a straight line, however, it can be empowered by polynomial features to get more powerful models. Still, we have to decide and fix the number and types of features ahead.
- Alternate approach is given by locally weighted regression.
- Given a dataset X, y , we attempt to find a model $h(x)$ that minimizes residual sum of weighted squared errors.
- The weights are given by a kernel function which can be chosen arbitrarily and in my case I chose a Gaussian kernel.
- The solution is very similar to Normal equations, we only need to insert diagonal weight matrix W .

Algorithm:

```
def local_regression(x0, X, Y, tau):
    # add bias term
    x0 = np.r_[1, x0]
    X = np.c_[np.ones(len(X)), X]
    # fit model: normal equations with kernel
    xw = X.T * radial_kernel(x0, X, tau)
    beta = np.linalg.pinv(xw @ X) @ xw @ Y
    # predict value
    return x0 @ beta
def radial_kernel(x0, X, tau):
    return np.exp(np.sum((X - x0) ** 2, axis=1) / (-2 * tau * tau))
```

Source code:

```
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))
    print("WEIGHTS", weights)
    for j in range(m):
        weights[j, :] = np.exp(-((xmat[j, :] - point) ** 2) / (k * k))
```

```

diff = point - X[j]
weights[j,j] = np.exp(diff * T / (-2.0 * K * k))
print("WEIGHTS", weights)
return weights

from numpy.core.memmap import ndarray
def local_weight(point, Xmat, ymat, K):
    print("XMAT", Xmat)
    print("YMAT", ymat)
    print("K", K)
    wei = kernel(point, Xmat, K)
    w = (X.T * (wei * X)).I * (X.T * (wei * ymat.T))
    print("W", w)
    return w

def localWeightRegression(Xmat, ymat, K):
    m, n = np.shape(Xmat)
    ypred = np.zeros(m)
    print("m", m)
    for i in range(m):
        ypred[i] = Xmat[i] * local_weight(Xmat[i], Xmat, ymat, K)
        print("ypred", ypred[i])
    return ypred

data = pd.read_csv("hotel-bill.csv")
bill = np.array(data.total_bill)
print("Bill", bill)
tip = np.array(data.tip)
print("Tip", tip)

```

```

mbill = np.mat(cbill)
print ("mbill", mbill)
mtip = np.mat(ctip)
print ("mtip", mtip)

```

```

m = np.shape(mbill)[1]
one = np.mat(np.ones(m))
X = np.hstack((one.T, mbill.T))
print ("X", X)

```

```

ypred = local_weight_regression (X, mtip, 0.5)
print ("ypred", ypred)
SortIndex = X[:,1].argsort(0)
print ("SortIndex", SortIndex)
xsort = X[SortIndex][:,0]
print ("xsort", xsort)

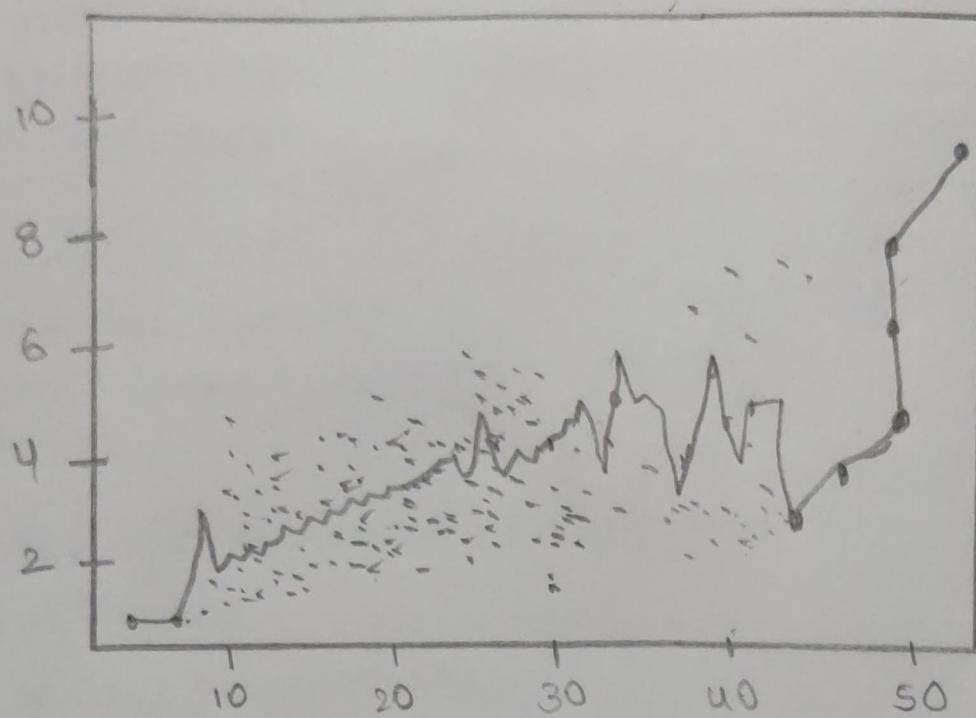
```

```

fig = plt.figure()
ax = fig.add_subplot(1,1,1)
ax.scatter(cbill, tip, color='green')
ax.plot(xsort[:,1], ypred[SortIndex], color='red', linewidth=5)
plt.xlabel('Total bill')
plt.ylabel('Tip')
plt.show();

```

Output:



Viva Voce Questions

1) What is locally weighted regression algorithm?

It is a non-parametric algorithm. The model will not learn a fixed set of parametric as is done in ordinary linear regression.

2) What are the advantages of locally weighted regression?

It allows to improve the over all performance of regression methods.

3) What type of learning is used in locally weighted regression?

It is a supervised learning algorithm is used in locally weighted regression. As it is a non-parametric algorithm.

4) What is logistic regression?

Logistic regression is used to calculate or predict the prob of binary ($Y|N$) event occurring.

5) What is weighted regression analysis?

It is a method that you can use when the least square assumption of constant variance are unresidual is violated.

Find the mean closest to the item

Assign item to mean

Update mean

Source code:

```

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

iris = datasets.load_iris()
X = pd.DataFrame(iris.data)
print(X.shape)
X.columns = ['sepal-length', 'sepal-width', 'petal-length',
              'petal-width']
y = pd.DataFrame(iris.target)
y.columns = ['target']
plt.figure(figsize=(10, 7))
colormap = np.array(['red', 'lime', 'black'])
plt.subplot(1, 2, 1)
plt.scatter(X['Sepal-Length'], X['Sepal-Width'], c=colormap[y['target']])
plt.title('Sepal')
plt.subplot(1, 2, 2)
plt.scatter(X['petal-length'], X['petal-width'], c=colormap[y['target']])
plt.title('petal')

```

```

model = KMeans(n_clusters=3)
model.fit(x)
print(model.labels_)
plt.subplot(1, 2, 1)
plt.scatter(x.sepal_Length, x.sepal_Width, c=colormap[y.target], s=100)
plt.title('Real classification')
plt.subplot(1, 2, 2)
plt.scatter(x.Sepal_Length, x.sepal_Width, c=colormap[model.labels_], s=100)
plt.title('KMEANS classification')

print(sm.accuracy_score(y, model.labels_))
sm.confusion_matrix(y, model.labels_)

from sklearn.naive_bayes import GaussianNB
clf = GaussianNB()
clf.fit(x,y)
y_cluster_gmm = clf.predict(x)
print(y_cluster_gmm)
plt.subplot(1, 2, 1)
plt.scatter(x.petal_Length, x.petal_Width, c=colormap[y.target], s=100)
plt.title('Real classification')
plt.subplot(1, 2, 2)
plt.scatter(x.petal_Length, x.petal_Width, c=colormap[y_cluster_gmm], s=100)
plt.title('KMEANS classification')

```

```
print (sm.accuracy_score(y,y_cluster_gmm))
```

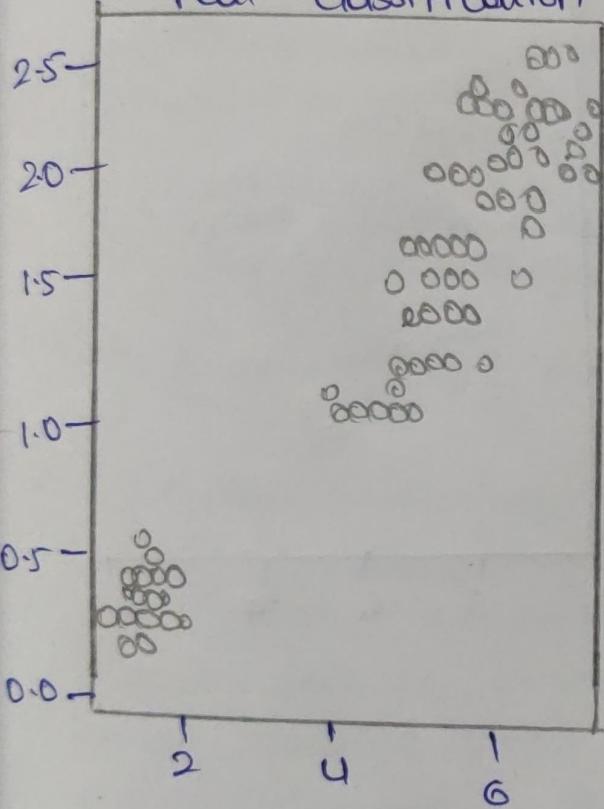
```
sm.confusion_matrix(y,y_cluster_gmm)
```

Output:

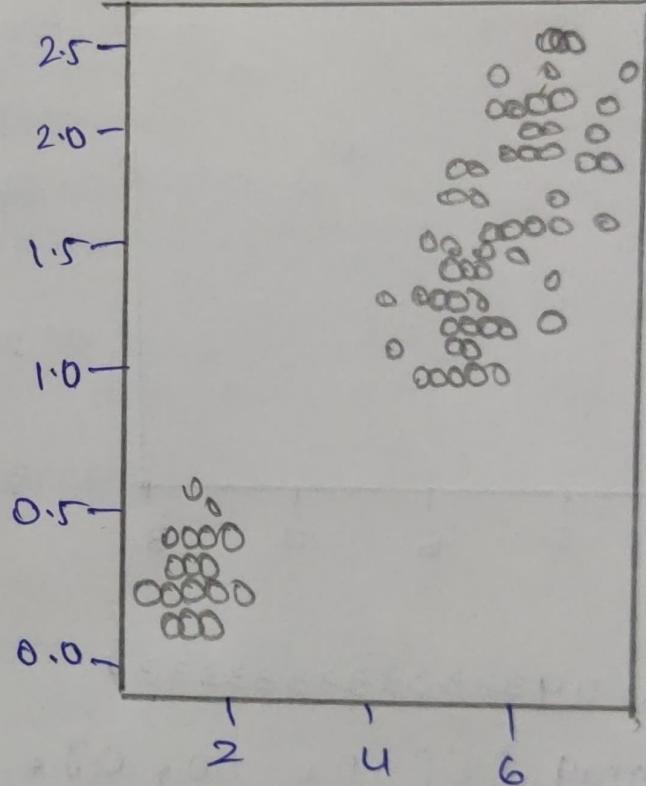
Gaussian matrix (n-components = 3)

Text (0.5, 1.0, 'kmeans classification')

Real classification



kmeans classification



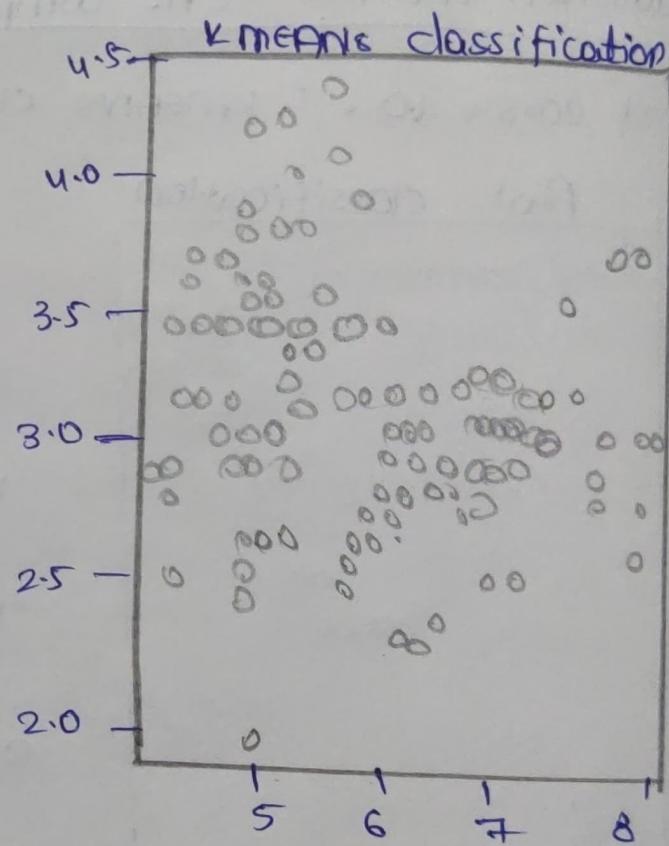
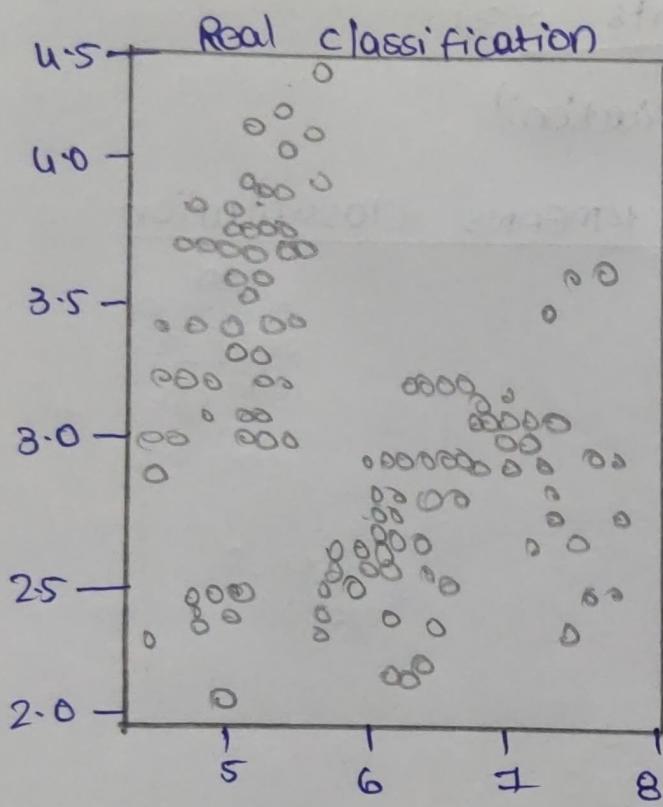
0.96

```
array([[50, 0, 0],  
       [0, 47, 3],  
       [0, 3, 4]])
```

Output:

K-means cn-clusters = 3)

Text (0.5, 1.0, 'kmeans classification').



0.0933333333333334

array [[0, 50, 0],
 [2, 0, 48],
 [36, 0, 14]]

Viva Voce Questions

1) What is EM algorithm in machine learning?

Expected maximization algorithm is defined as combo of various unsupervised machine learning also which used to define local maximum likelihood estimation for unobservable in statistical models.

2) What is K-means clustering algorithm and List out the applications?

K-means clustering is a type of unsupervised learning, which is used when you have unlabeled data. The goal of this algorithm is to find groups in the data, with the no. of groups represented by the variable k.

Applications:- public data, capturing Images.

3) What are the advantages and disadvantages of K-means clustering?

Advantage:-
 → Relatively simple to implement
 → scales to large data sets.

Disadvantages:- choosing k manually, being depends on initial value

4) What is the difference between K-means and EM?

K-means and EM in the sense that they allow model refining of an iterative process to find the best congestion. The K-means algorithm differs in the method used for calculating the Euclidean distance while calculating the distance b/w of two data items; and EM uses statistical methods.

5) What are the Applications of EM algorithm?

- * chosen in unsupervised data clustering and Psychiatric analysis
- * it has numerous uses in NLP, complex visual quantitative analysis of genetics.