

VISVESVARAYA TECHNOLOGICAL UNIVERSITY

“JnanaSangama”, Belgaum -590014, Karnataka.



LAB REPORT on

Machine Learning (20CS6PCMAL)

Submitted by

KHUSHIL M SINDHWAD
(1BM19CS072)

in partial fulfillment for the award of the degree of
BACHELOR OF ENGINEERING
in
COMPUTER SCIENCE AND ENGINEERING



B.M.S. COLLEGE OF ENGINEERING

(Autonomous Institution under VTU)

BENGALURU-560019

May-2022 to July-2022

**B. M. S. College of Engineering,
Bull Temple Road, Bangalore 560019**
(Affiliated To Visvesvaraya Technological University, Belgaum)
Department of Computer Science and Engineering



CERTIFICATE

This is to certify that the Lab work entitled "**MACHINE LEARNING**" carried out by **KHUSHIL M SINDHWAD (1BM19CS072)**, who is bonafide student of **B. M. S. College of Engineering**. It is in partial fulfillment for the award of **Bachelor of Engineering in Computer Science and Engineering** of the Visvesvaraya Technological University, Belgaum during the year 2022. The Lab report has been approved as it satisfies the academic requirements in respect of a **Machine Learning - (20CS6PCMAL)** work prescribed for the said degree.

Name of the Lab-Incharge
Designation
Department of CSE
BMSCE, Bengaluru

Dr. Asha G R
Assistant Professor
Department of CSE
BMSCE, Bengaluru

Index Sheet

Sl. No.	Experiment Title	Page No.
1	Find S Algorithm	4
2	Candidate Elimination Algorithm	5
3	ID3 Algorithm	6
4	Naïve Bayes	9
5	Linear Regression	11
6	Bayesian Network	12
7	K Means	15
8	EM Algorithm	24
9	KNN Algorithm	26
10	Locally Weighted Regression	27

Course Outcome

CO1	Ability to apply the different learning algorithms.
CO2	Ability to analyse the learning techniques for given dataset.
CO3	Ability to design a model using machine learning to solve a problem.
CO4	Ability to conduct practical experiments to solve problems using appropriate machine learning techniques.

1. Find S Algorithm

Implement and demonstrate the FIND-S algorithm for finding the most specific hypothesis based on a given set of training data samples

```
In [1]: import numpy as np
import pandas as pd

In [2]: from google.colab import drive
drive.mount("/content/drive")

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

In [3]: data = pd.read_csv("/content/drive/MyDrive/finddata.csv")
print(data,"\n")

   Time Weather Temperature Company Humidity Wind Goes
0 Morning Sunny Warm Yes Mild Strong Yes
1 Evening Rainy Cold No Mild Normal No
2 Morning Sunny Moderate Yes Normal Normal Yes
3 Evening Sunny Cold Yes High Strong Yes

In [4]: d = np.array(data)[:, :-1]
print("\n The attributes are: ",d)
target = np.array(data)[:, -1]
print("\n The target is: ",target)

The attributes are: [['Morning' 'Sunny' 'Warm' 'Yes' 'Mild' 'Strong']
 ['Evening' 'Rainy' 'Cold' 'No' 'Mild' 'Normal']
 ['Morning' 'Sunny' 'Moderate' 'Yes' 'Normal' 'Normal']
 ['Evening' 'Sunny' 'Cold' 'Yes' 'High' 'Strong']]

The target is: ['Yes' 'No' 'Yes' 'Yes']

In [5]: def findS(c,t):
    for i, val in enumerate(t):
        if val == "Yes":
            specific_hypothesis = c[i].copy()
            break

    The attributes are: [['Morning' 'Sunny' 'Warm' 'Yes' 'Mild' 'Strong']
 ['Evening' 'Rainy' 'Cold' 'No' 'Mild' 'Normal']
 ['Morning' 'Sunny' 'Moderate' 'Yes' 'Normal' 'Normal']
 ['Evening' 'Sunny' 'Cold' 'Yes' 'High' 'Strong']]

    The target is: ['Yes' 'No' 'Yes' 'Yes']

    def findS(c,t):
        for i, val in enumerate(t):
            if val == "Yes":
                specific_hypothesis = c[i].copy()
                break

            for i, val in enumerate(c):
                if t[i] == "Yes":
                    for x in range(len(specific_hypothesis)):
                        if val[x] != specific_hypothesis[x]:
                            specific_hypothesis[x] = '?'
                    else:
                        pass

    return specific_hypothesis

In [6]: print("\n The final hypothesis is:",findS(d,target))

The final hypothesis is: ['?' 'Sunny' '?' 'Yes' '?' '?']
```

2. Candidate Elimination Algorithm

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.

```

In [14]: import numpy as np
import pandas as pd

#to read the data in the csv file
data = pd.DataFrame(data=pd.read_csv('enjoysport.csv'))
print(data,"\n")

#making an array of all the attributes
concepts = np.array(data.iloc[:,0:-1])
print("The attributes are: ",concepts)

#segregating the target that has positive and negative examples
target = np.array(data.iloc[:, -1])
print("\n The target is: ",target)

#training function to implement candidate_elimination algorithm
def learn(concepts, target):
    specific_h = concepts[0].copy()
    print("Initialization of specific hypothesis and general hypothesis:")
    print("S0: ",specific_h)
    general_h = [[ "?" for i in range(len(specific_h))] for i in range(len(specific_h))]
    print("G0: ",general_h)
    for i, h in enumerate(concepts):
        if target[i] == "yes":
            for x in range(len(specific_h)):
                if h[x] != specific_h[x]:
                    specific_h[x] = "?"
                    general_h[x][x] = "?"
        if target[i] == "no":
            for x in range(len(specific_h)):
                if h[x] != specific_h[x]:
                    general_h[x][x] = specific_h[x]
                else:
                    general_h[x][x] = "?"
    print("The steps of Candidate Elimination Algorithm ", i+1)
    print(f"S{i+1}: ", specific_h)
    actual_general_h= [i for i in general_h if i != ["?" for i in range(len(specific_h))]]
    print(f"G{i+1}: ",actual_general_h)
    print(general_h)
    print(actual_general_h)
    actual_general_h= [i for i in general_h if i != ["?" for i in range(len(specific_h))]]

1 sunny warm high strong warm same yes
2 rainy cold high strong warm change no
3 sunny warm high strong cool change yes

The attributes are: [['sunny' 'warm' 'normal' 'strong' 'warm' 'same']
['sunny' 'warm' 'high' 'strong' 'warm' 'same']
['rainy' 'cold' 'high' 'strong' 'warm' 'change']
['sunny' 'warm' 'high' 'strong' 'cool' 'change']]

The target is: ['yes' 'yes' 'no' 'yes']

In [15]: s_final, g_final = learn(concepts, target)
print("The final Specific hypothesis is: ", s_final, sep="\n")
print("The final General Hypothesis is: ", g_final, sep="\n")

Initialization of specific hypothesis and general hypothesis:
S0: ['sunny' 'warm' 'normal' 'strong' 'warm' 'same']
G0: [[ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?']]
The steps of Candidate Elimination Algorithm 1
S1: ['sunny' 'warm' 'normal' 'strong' 'warm' 'same']
G1: []
[[ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?']]
The steps of Candidate Elimination Algorithm 2
S2: ['sunny' 'warm' '?' 'strong' 'warm' 'same']
G2: []
[[ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?']]
The steps of Candidate Elimination Algorithm 3
S3: ['sunny' 'warm' '?' 'strong' 'warm' 'same']
G3: [[ 'sunny', '?', '?', '?', '?', '?'], [ '?', 'warm', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [[ 'sunny', '?', '?', '?', '?', '?'], [ '?', 'warm', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?']]]
The steps of Candidate Elimination Algorithm 4
S4: ['sunny' 'warm' '?' 'strong' '?' '?']
G4: [[ 'sunny', '?', '?', '?', '?', '?'], [ '?', 'warm', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [[ 'sunny', '?', '?', '?', '?', '?'], [ '?', 'warm', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'],
      [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?'], [ '?', '?', '?', '?', '?', '?']]]
The final Specific hypothesis is:
['sunny' 'warm' '?' 'strong' '?' '?']
The final General Hypothesis is:
[[ 'sunny', '?', '?', '?', '?', '?'], [ '?', 'warm', '?', '?', '?', '?']]
```

3. ID3 Algorithm

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.

```
In [32]: import pandas as pd  
import numpy as np  
  
from sklearn.datasets import load_iris  
  
data = load_iris()  
  
In [33]: df = pd.DataFrame(data.data, columns = data.feature_names)  
  
In [34]: df  
  
Out[34]:  
    sepal length (cm)  sepal width (cm)  petal length (cm)  petal width (cm)  
0             5.1          3.5            1.4           0.2  
1             4.9          3.0            1.4           0.2  
2             4.7          3.2            1.3           0.2  
3             4.6          3.1            1.5           0.2  
4             5.0          3.6            1.4           0.2  
...           ...          ...           ...           ...  
145            6.7          3.0            5.2           2.3  
146            6.3          2.5            5.0           1.9  
147            6.5          3.0            5.2           2.0  
148            6.2          3.4            5.4           2.3  
149            5.9          3.0            5.1           1.8  
  
150 rows × 4 columns
```

```
In [20]: df['Species'] = data.target  
#replace this with the actual names  
  
target = np.unique(data.target)  
print(target)  
  
target_names = np.unique(data.target_names)  
print(target_names)  
targets = dict(zip(target, target_names))  
print(targets)  
df['Species'] = df['Species'].replace(targets)  
  
[0 1 2]  
['setosa' 'versicolor' 'virginica']  
{0: 'setosa', 1: 'versicolor', 2: 'virginica'}  
  
In [21]: x = df.drop(columns="Species")  
y = df["species"]  
  
In [22]: feature_names = x.columns  
labels = y.unique()  
  
In [36]: from sklearn.model_selection import train_test_split  
  
x_train, test_x, y_train, test_lab = train_test_split(x,y,test_size = 0.4,random_state = 42)  
# print(x_train)  
# print(test_x)  
# print(y_train)  
# print(test_lab)
```

```
In [24]: from sklearn.tree import DecisionTreeClassifier  
clf = DecisionTreeClassifier(max_depth=3, random_state=42, criterion='entropy')
```

```
In [25]: clf.fit(x_train, y_train)
```

```
Out[25]: DecisionTreeClassifier(criterion='entropy', max_depth=3, random_state=42)
```

```
In [35]: test_pred = clf.predict(test_x)  
test_pred
```

```
Out[35]: array(['versicolor', 'setosa', 'virginica', 'versicolor', 'versicolor',  
       'setosa', 'versicolor', 'virginica', 'versicolor', 'versicolor',  
       'virginica', 'setosa', 'setosa', 'setosa', 'setosa', 'versicolor',  
       'virginica', 'versicolor', 'versicolor', 'virginica', 'setosa',  
       'virginica', 'setosa', 'virginica', 'virginica', 'virginica',  
       'virginica', 'versicolor', 'setosa', 'setosa', 'setosa', 'setosa',  
       'versicolor', 'setosa', 'virginica', 'versicolor',  
       'setosa', 'setosa', 'virginica', 'versicolor',  
       'versicolor', 'setosa', 'setosa', 'versicolor', 'versicolor',  
       'virginica', 'versicolor', 'virginica', 'versicolor', 'virginica',  
       'versicolor', 'setosa', 'virginica', 'versicolor', 'setosa',  
       'setosa', 'setosa', 'versicolor'], dtype=object)
```

```
In [27]: from sklearn import metrics
```

```
import seaborn as sns
```

```
import matplotlib.pyplot as plt
```

```
confusion_matrix = metrics.confusion_matrix(test_lab, test_pred)
```

```
In [28]: confusion_matrix
```

```
Out[28]: array([[23, 0, 0],
```

```
               [0, 19, 0],
```

```
               [0, 1, 17]], dtype=int64)
```

```
In [29]: matrix_df = pd.DataFrame(confusion_matrix)
```

```
ax = plt.axes()
```

```
sns.set(font_scale=1.3)
```

```
plt.figure(figsize=(10,7))
```

```
sns.heatmap(matrix_df, annot=True, fmt="g", ax=ax, cmap="magma")
```

```
ax.set_title("Confusion Matrix - Decision Tree")
```

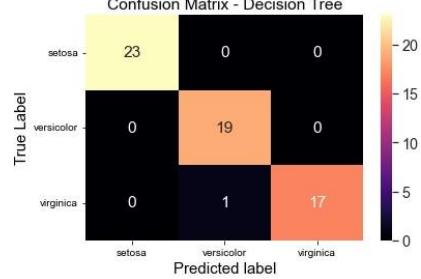
```
ax.set_xlabel("Predicted label", fontsize=15)
```

```
ax.set_xticklabels(['']+labels)
```

```
ax.set_ylabel("True Label", fontsize=15)
```

```
ax.set_yticklabels(list(labels), rotation=0)
```

```
plt.show()
```



```
<Figure size 720x504 with 0 Axes>
```

```
In [30]: clf.score(test_x, test_lab)
```

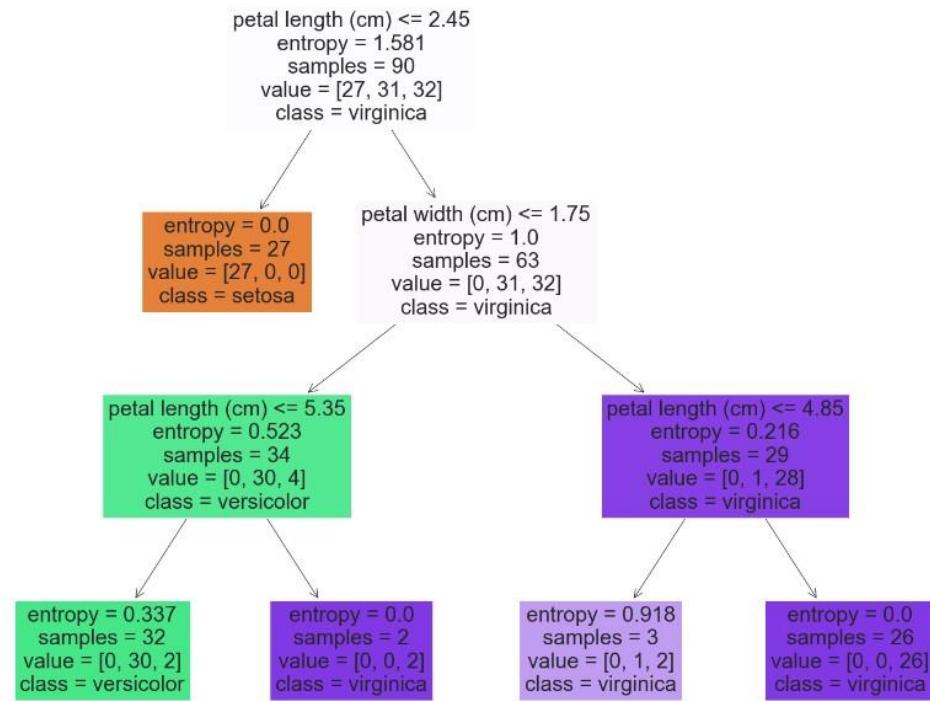
```
Out[30]: 0.9833333333333333
```

```
In [31]: from sklearn import tree
```

```
fig = plt.figure(figsize=(25,20))
```

```
_ = tree.plot_tree(clf,  
                   feature_names=data.feature_names,  
                   class_names=data.target_names,  
                   filled=True)
```

```
In [31]: from sklearn import tree
fig = plt.figure(figsize=(25,20))
_ = tree.plot_tree(clf,
                   feature_names=data.feature_names,
                   class_names=data.target_names,
                   filled=True)
```



In []:

4. Naïve Bayes

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

```
In [10]: import pandas as pd
data = pd.read_csv('PlayTennis.csv')
data

Out[10]:
   PlayTennis Outlook Temperature Humidity Wind
0           No     Sunny        Hot    High  Weak
1           No     Sunny        Hot    High  Strong
2          Yes  Overcast       Hot    High  Weak
3          Yes      Rain       Mild    High  Weak
4          Yes      Rain      Cool    Normal  Weak
5           No      Rain      Cool    Normal  Strong
6          Yes  Overcast      Cool    Normal  Strong
7           No     Sunny       Mild    High  Weak
8          Yes     Sunny      Cool    Normal  Weak
9          Yes      Rain       Mild    Normal  Weak
10         Yes     Sunny       Mild    Normal  Strong
11         Yes  Overcast      Mild    High  Strong
12         Yes  Overcast       Hot    Normal  Weak
13           No      Rain       Mild    High  Strong
```

```
In [11]: y = list(data['PlayTennis'].values)
X = data.iloc[:,1:6].values

print(f'Target Values: {y}')
print(f'Features: \n{X}')

Target Values: ['No', 'No', 'Yes', 'Yes', 'Yes', 'No', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']
Features:
[['Sunny', 'Hot', 'High', 'Weak'],
 ['Sunny', 'Hot', 'High', 'Strong'],
 ['Overcast', 'Hot', 'High', 'Weak'],
 ['Rain', 'Mild', 'High', 'Weak'],
 ['Rain', 'Cool', 'Normal', 'Weak'],
 ['Rain', 'Cool', 'Normal', 'Strong'],
 ['Overcast', 'Cool', 'Normal', 'Strong'],
 ['Sunny', 'Mild', 'High', 'Weak'],
 ['Sunny', 'Cool', 'High', 'Weak'],
 ['Rain', 'Mild', 'Normal', 'Weak'],
 ['Sunny', 'Mild', 'Normal', 'Strong'],
 ['Overcast', 'Mild', 'High', 'Strong'],
 ['Overcast', 'Hot', 'Normal', 'Weak'],
 ['Rain', 'Mild', 'High', 'Weak']]
```

```

['Rain' 'Mild' 'Normal' 'Weak']
['Sunny' 'Mild' 'Normal' 'Strong']
['Overcast' 'Mild' 'High' 'Strong']
['Overcast' 'Hot' 'Normal' 'Weak']
['Rain' 'Mild' 'High' 'Strong']

In [12]: y_train = y[:8]
y_val = y[8:]

X_train = X[:8]
X_val = X[8:]

print(f"Number of instances in training set: {len(X_train)}")
print(f"Number of instances in testing set: {len(X_val)}")

Number of instances in training set: 8
Number of instances in testing set: 6

In [17]: class NaiveBayesClassifier:
    def __init__(self, X, y):
        self.X, self.y = X, y
        self.N = len(self.X)
        self.dim = len(self.X[0])
        self.attrs = [[] for _ in range(self.dim)]
        self.output_dom = {}
        self.data = []
        for i in range(len(self.X)):
            for j in range(self.dim):
                if not self.X[i][j] in self.attrs[j]:
                    self.attrs[j].append(self.X[i][j])
            if not self.y[i] in self.output_dom.keys():
                self.output_dom[self.y[i]] = 1
            else:
                self.output_dom[self.y[i]] += 1
            self.data.append((self.X[i], self.y[i]))
        print(self.attrs)
        print(self.output_dom)
    #     print(self.data)
    def classify(self, entry):
        solve = None
        max_arg = -1
        for y in self.output_dom.keys():
            prob = self.output_dom[y]/self.N
            for i in range(self.dim):
                cases = [x for x in self.data if x[0][i] == entry[i] and x[1] == y]
                n = len(cases)
                prob *= n/self.N
            if prob > max_arg:
                max_arg = prob
                solve = y
        return solve

In [18]: nbc = NaiveBayesClassifier(X_train, y_train)

total_cases = len(y_val)

good = 0
bad = 0
predictions = []

for i in range(total_cases):
    predict = nbc.classify(X_val[i])
    predictions.append(predict)

    if y_val[i] == predict:
        good += 1
    else:
        bad += 1

print('Predicted values:', predictions)
print('Actual values:', y_val)
print()
print('Total number of testing instances in the dataset:', total_cases)
print('Number of correct predictions:', good)
print('Number of wrong predictions:', bad)
print()
print('Accuracy of Bayes Classifier:', good/total_cases)

[['Sunny'], ['Hot'], ['High'], ['Weak']]
{'No': 1}
[['Sunny'], ['Hot'], ['High'], ['Weak', 'Strong']]
{'No': 2}
[['Sunny', 'Overcast'], ['Hot'], ['High'], ['Weak', 'Strong']]
{'No': 2, 'Yes': 1}
[['Sunny', 'Overcast', 'Rain'], ['Hot', 'Mild'], ['High'], ['Weak', 'Strong']]
{'No': 2, 'Yes': 2}
[['Sunny', 'Overcast', 'Rain'], ['Hot', 'Mild', 'Cool'], ['High', 'Normal'], ['Weak', 'Strong']]
{'No': 2, 'Yes': 3}
[['Sunny', 'Overcast', 'Rain'], ['Hot', 'Mild', 'Cool'], ['High', 'Normal'], ['Weak', 'Strong']]
{'No': 3, 'Yes': 3}
[['Sunny', 'Overcast', 'Rain'], ['Hot', 'Mild', 'Cool'], ['High', 'Normal'], ['Weak', 'Strong']]
{'No': 3, 'Yes': 4}
[['Sunny', 'Overcast', 'Rain'], ['Hot', 'Mild', 'Cool'], ['High', 'Normal'], ['Weak', 'Strong']]
{'No': 4, 'Yes': 4}

Predicted values: ['No', 'Yes', 'No', 'Yes', 'Yes', 'No']
Actual values: ['Yes', 'Yes', 'Yes', 'Yes', 'Yes', 'No']

Total number of testing instances in the dataset: 6
Number of correct predictions: 4
Number of wrong predictions: 2

Accuracy of Bayes Classifier: 0.6666666666666666

```

5. Linear Regression

Implement the Linear Regression algorithm in order to fit data points. Select appropriate data set for your experiment and draw graphs.

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

In [2]: dataset = pd.read_csv('salary_data.csv')
X = dataset.iloc[:, :-1].values
y = dataset.iloc[:, 1].values

In [3]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=1/3, random_state=0)

In [4]: # Fitting Simple Linear Regression to the Training set
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(X_train, y_train)

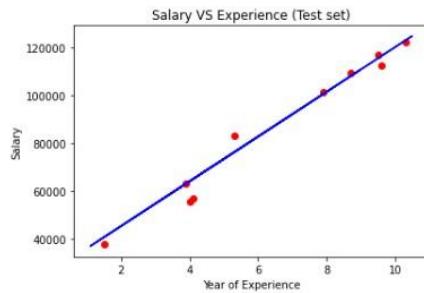
Out[4]: LinearRegression()

In [5]: # Predicting the Test set results
y_pred = regressor.predict(X_test)

In [6]: # Visualizing the Training set results
viz_train = plt
viz_train.scatter(X_train, y_train, color='red')
viz_train.plot(X_train, regressor.predict(X_train), color='blue')
viz_train.title('Salary VS Experience (Training set)')
viz_train.xlabel('Year of Experience')
viz_train.ylabel('Salary')
viz_train.show()
```



```
In [7]: # Visualizing the Test set results
viz_test = plt
viz_test.scatter(X_test, y_test, color='red')
viz_test.plot(X_train, regressor.predict(X_train), color='blue')
viz_test.title('Salary VS Experience (Test set)')
viz_test.xlabel('Year of Experience')
viz_test.ylabel('Salary')
viz_test.show()
```



```
In [ ]:
```

6. Write a program to construct a Bayesian network considering training data. Use this model to make predictions.

7/10/22, 11:18 PM

BayesianNetwork.ipynb - Colaboratory

```
import numpy as np
import pandas as pd
import csv
from pgmpy.estimators import MaximumLikelihoodEstimator
from pgmpy.models import BayesianModel
from pgmpy.inference import VariableElimination

/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
import pandas.util.testing as tm

[< >]

trainingData = pd.read_csv('/content/bayesian-dataset.csv')
trainingData = trainingData.replace('?',np.nan)
print('The sample instances from the dataset are:')
print(trainingData.head())
print('\n Attributes and datatypes: ')
print(trainingData.dtypes)

The sample instances from the dataset are:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope \
0    63     1     1      145    233     1      2       150      0      2.3      3
1    67     1     4      160    286     0      2       108      1      1.5      2
2    67     1     4      120    229     0      2       129      1      2.6      2
3    37     1     3      130    250     0      0       187      0      3.5      3
4    41     0     2      130    204     0      2       172      0      1.4      1

   ca  thal  heartdisease
0   0     6             0
1   3     3             2
2   2     7             1
3   0     3             0
4   0     3             0

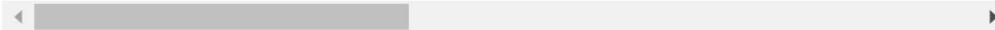
Attributes and datatypes:
age           int64
sex           int64
cp            int64
trestbps     int64
chol          int64
fbs           int64
restecg      int64
thalach      int64
exang          int64
oldpeak      float64
slope          int64
ca            object
thal          object
heartdisease int64
dtype: object

model = BayesianModel([('age','heartdisease'),('sex','heartdisease'),('exang','heartdisease')])
print('\n Learning CPD using Maximum likelihood estimators')
model.fit(trainingData,estimator=MaximumLikelihoodEstimator)
print('\n Inferencing with Bayesian Network:')
HeartDiseasetest_infer = VariableElimination(model)
```

```
/usr/local/lib/python3.7/dist-packages/pgmpy/models/BayesianModel.py:10: FutureWarning:  
FutureWarning,
```

```
Learning CPD using Maximum likelihood estimators
```

```
Inferencing with Bayesian Network:
```



```
print('\n 1. Probability of HeartDisease given evidence = restecg (Rest ECG): 1')  
q1 = HeartDiseasetest_infer.query(variables = ['heartdisease'], evidence={'restecg':1})  
print(q1)
```

```
1. Probability of HeartDisease given evidence = restecg (Rest ECG): 1
```

```
Finding Elimination Order: : 4/4 [00:00<00:00,
```

```
100% 8.00it/s]
```

```
Eliminating: age: 100% 4/4 [00:00<00:00, 4.88it/s]
```

heartdisease	phi(heartdisease)
heartdisease(0)	0.1012
heartdisease(1)	0.0000
heartdisease(2)	0.2392
heartdisease(3)	0.2015
heartdisease(4)	0.4581

```
print('\n 2. Probability of HeartDisease given evidence = chol (Cholestorol): 100 ')  
q2 = HeartDiseasetest_infer.query(variables = ['heartdisease'], evidence={'chol':100})  
print(q2)
```

```
2. Probability of HeartDisease given evidence = chol (Cholestorol): 100
/usr/local/lib/python3.7/dist-packages/pgmpy/factors/discrete/DiscreteFactor.py:537:
UserWarning,
Finding Elimination Order: : 4/4 [00:00<00:00,
100% 7.39it/s]

Eliminating: age: 100% 4/4 [00:00<00:00, 3.70it/s]

+-----+
| heartdisease | phi(heartdisease) |
+=====+
| heartdisease(0) | 1.0000 |
+-----+
| heartdisease(1) | 0.0000 |
+-----+
| heartdisease(2) | 0.0000 |
+-----+
| heartdisease(3) | 0.0000 |
+-----+
```

7. Apply k-Means algorithm to cluster a set of data stored in a .CSV file.

(using built in)

7/10/22, 11:03 PM K-MeansClustering(using libraries).ipynb - Colaboratory

▼ K-Means Clustering

```
import pandas as pd
from sklearn.cluster import KMeans
from sklearn.preprocessing import MinMaxScaler
from matplotlib import pyplot as plt
%matplotlib inline

from google.colab import drive
drive.mount('/content/drive')

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.remount()

df = pd.read_csv('/content/drive/MyDrive/income.csv')
df.head(10)
```

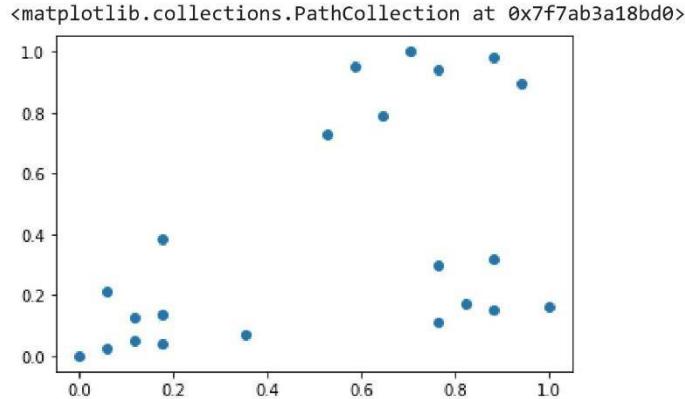
	Name	Age	Income(\$)
0	Rob	27	70000
1	Michael	29	90000
2	Mohan	29	61000
3	Ismail	28	60000
4	Kory	42	150000
5	Gautam	39	155000
6	David	41	160000
7	Andrea	38	162000
8	Brad	36	156000
9	Angelina	35	130000

```
scaler = MinMaxScaler()
scaler.fit(df[['Age']])
df[['Age']] = scaler.transform(df[['Age']])

scaler.fit(df[['Income($)']])
df[['Income($)']] = scaler.transform(df[['Income($)']])
df.head(10)
```

	Name	Age	Income(\$)
0	Rob	0.058824	0.213675
1	Michael	0.176471	0.384615
2	Mohan	0.176471	0.136752
3	Ismail	0.117647	0.128205
4	Kory	0.941176	0.897436
5	Gautam	0.764706	0.940171
6	David	0.882353	0.982906
7	Andrea	0.705992	1.000000

```
plt.scatter(df['Age'], df['Income($)'])
```



▼ Finding Elbow Point

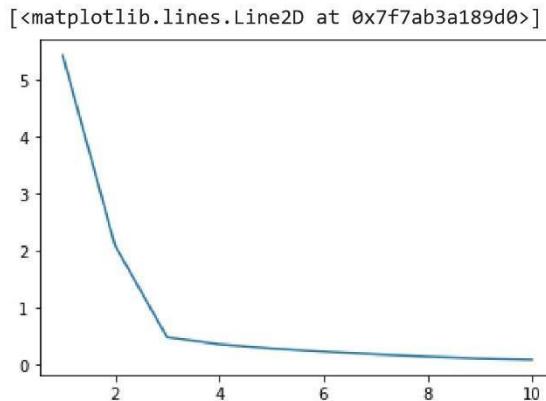
```
k_range = range(1, 11)
sse = []
for k in k_range:
    kmc = KMeans(n_clusters=k)
    kmc.fit(df[['Age', 'Income($)']])
    sse.append(kmc.inertia_)

sse
[5.434011511988178,
 2.091136388699078,
 0.4750783498553096,
 0.3491047094419566,
 0.27768187154369994,
 0.22020960864009398,
 0.1735559655531264,
 0.1327661931978319,
 0.10188787724979426,
 0.08026197041664467]
```

7/10/22, 11:03 PM

K-MeansClustering(using libraries).ipynb - Colaboratory

```
plt.xlabel = 'Number of Clusters'  
plt.ylabel = 'Sum of Squared Errors'  
plt.plot(k_range, sse)
```



- ▼ Therefore, the elbow point is 3

```
km = KMeans(n_clusters=3)  
km  
  
KMeans(n_clusters=3)  
  
y_predict = km.fit_predict(df[['Age', 'Income($)']])  
y_predict  
  
array([1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2],  
      dtype=int32)
```

```
df['cluster'] = y_predict  
df.head()
```

	Name	Age	Income(\$)	cluster
0	Rob	0.058824	0.213675	1
1	Michael	0.176471	0.384615	1
2	Mohan	0.176471	0.136752	1
3	Ismail	0.117647	0.128205	1
4	Kory	0.941176	0.897436	0

```
df0 = df[df.cluster == 0]  
df0
```

7/10/22, 11:03 PM

K-MeansClustering(using libraries).ipynb - Colaboratory

	Name	Age	Income(\$)	cluster
4	Kory	0.941176	0.897436	0
5	Gautam	0.764706	0.940171	0
6	David	0.882353	0.982906	0
7	Andrea	0.705882	1.000000	0
8	Brad	0.588235	0.948718	0
9	Angelina	0.529412	0.726496	0

```
df1 = df[df.cluster == 1]
df1
```



	Name	Age	Income(\$)	cluster
0	Rob	0.058824	0.213675	1
1	Michael	0.176471	0.384615	1
2	Mohan	0.176471	0.136752	1
3	Ismail	0.117647	0.128205	1
11	Tom	0.000000	0.000000	1
12	Arnold	0.058824	0.025641	1
13	Jared	0.117647	0.051282	1
14	Stark	0.176471	0.038462	1
15	Ranbir	0.352941	0.068376	1

```
df2 = df[df.cluster == 2]
df2
```

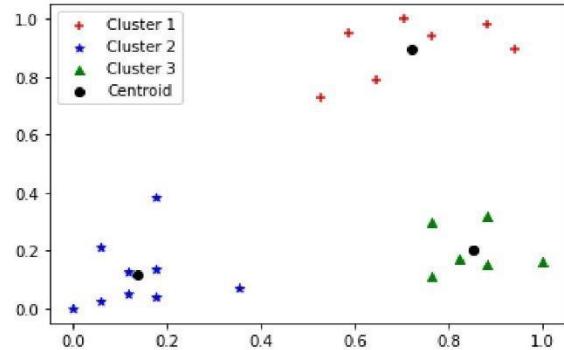
	Name	Age	Income(\$)	cluster
16	Dipika	0.823529	0.170940	2
17	Priyanka	0.882353	0.153846	2
18	Nick	1.000000	0.162393	2
19	Alia	0.764706	0.299145	2
20	Sid	0.882353	0.316239	2
21	Abdul	0.764706	0.111111	2

km.cluster_centers_

```
array([[0.72268908, 0.8974359 ],
       [0.1372549 , 0.11633428],
       [0.85294118, 0.2022792 ]])
```

```
p1 = plt.scatter(df0['Age'], df0['Income($)'), marker='+', color='red')
p2 = plt.scatter(df1['Age'], df1['Income($)'), marker='*', color='blue')
p3 = plt.scatter(df2['Age'], df2['Income($)'), marker='^', color='green')
c = plt.scatter(km.cluster_centers_[:,0], km.cluster_centers_[:,1], color='black')
plt.legend((p1, p2, p3, c),
           ('Cluster 1', 'Cluster 2', 'Cluster 3', 'Centroid'))
```

```
<matplotlib.legend.Legend at 0x7f7aad7cf610>
```



```
import math;
import sys;
import pandas as pd
import numpy as np
from random import choice
from matplotlib import pyplot
from random import shuffle, uniform;

def ReadData(fileName):
    f = open(fileName,'r')
    lines = f.read().splitlines()
    f.close()

    items = []

    for i in range(1,len(lines)):
        line = lines[i].split(',')
        itemFeatures = []

        for j in range(len(line)-1):
            v = float(line[j])
            itemFeatures.append(v)
        items.append(itemFeatures)

    shuffle(items)

    return items

def FindColMinMax(items):
    n = len(items[0])
    minima = [float('inf') for i in range(n)]
    maxima = [float('-inf') -1 for i in range(n)]

    for item in items:
        for f in range(len(item)):
            if(item[f] < minima[f]):
                minima[f] = item[f]

            if(item[f] > maxima[f]):
                maxima[f] = item[f]

    return minima,maxima

def EuclideanDistance(x,y):
    S = 0
    for i in range(len(x)):
        S += math.pow(x[i]-y[i],2)

    return math.sqrt(S)

def InitializeMeans(items, k, cMin, cMax):
```

7/10/22, 11:04 PM K-MeansClustering_hardcoded.ipynb - Colaboratory

```
def InitializeMeans(items,k,cMin,cMax):
    f = len(items[0])
    means = [[0 for i in range(f)] for j in range(k)]

    for mean in means:
        for i in range(len(mean)):
            mean[i] = uniform(cMin[i]+1,cMax[i]-1)

    return means

def UpdateMean(n,mean,item):
    for i in range(len(mean)):
        m = mean[i]
        m = (m*(n-1)+item[i])/float(n)
        mean[i] = round(m,3)

    return mean

def FindClusters(means,items):
    clusters = [[] for i in range(len(means))]

    for item in items:
        index = Classify(means,item)
        clusters[index].append(item)

    return clusters

def Classify(means,item):

    minimum = float('inf');
    index = -1

    for i in range(len(means)):
        dis = EuclideanDistance(item,means[i])

        if(dis < minimum):
            minimum = dis
            index = i

    return index

def CalculateMeans(k,items,maxIterations=1000000):
    cMin, cMax = FindColMinMax(items)

    means = InitializeMeans(items,k,cMin,cMax)

    clusterSizes = [0 for i in range(len(means))]

    belongsTo = [0 for i in range(len(items))]

    for e in range(maxIterations):
        noChange = True;
        for i in range(len(items)):
            item = items[i];
            index = Classify(means.item)
```

7/10/22, 11:04 PM K-MeansClustering_hardcoded.ipynb - Colaboratory

```

        clusterSizes[index] += 1
        cSize = clusterSizes[index]
        means[index] = UpdateMean(cSize,means[index],item)

        if(index != belongsTo[i]):
            noChange = False
            belongsTo[i] = index

    if (noChange):
        break

return means

def CutToTwoFeatures(items,indexA,indexB):
    n = len(items)
    X = []
    for i in range(n):
        item = items[i]
        newItem = [item[indexA],item[indexB]]
        X.append(newItem)

    return X

def PlotClusters(clusters):
    n = len(clusters)
    X = [[] for i in range(n)]

    for i in range(n):
        cluster = clusters[i]
        for item in cluster:
            X[i].append(item)

    colors = ['r','b','g','c','m','y']

    for x in X:
        c = choice(colors)
        colors.remove(c)

    Xa = []
    Xb = []

    for item in x:
        Xa.append(item[0])
        Xb.append(item[1])

    pyplot.plot(Xa,Xb,'o',color=c)

    pyplot.show()

def main():
    items = ReadData('data.txt')
    k = 3
    items = CutToTwoFeatures(items,2,3)
    print(items)

```

```

7/10/22, 11:04 PM K-MeansClustering_hardcoded.ipynb - Colaboratory

means = CalculateMeans(k,items)
print("\nMeans = ", means)

clusters = FindClusters(means,items)

PlotClusters(clusters)
newItem = [1.5,0.2]
print(Classify(means,newItem))

if __name__ == "__main__":
    main()

❶ [[6.3, 1.8], [5.1, 1.9], [1.9, 0.4], [1.6, 0.4], [4.9, 1.8], [1.7, 0.4], [4.2, 1.2],
Means = [[5.583, 2.032], [1.462, 0.257], [4.258, 1.341]]


```

8. Apply EM algorithm to cluster a set of data stored in a .CSV file.

Compare the results of k-Means algorithm and EM algorithm.

7/10/22, 11:05 PM

EM_Kmeans_Comparison.ipynb - Colaboratory

```
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)
X.columns = ['Sepal_Length', 'Sepal_Width', 'Petal_Length', 'Petal_Width']

y = pd.DataFrame(iris.target)
y.columns = ['Targets']

model = KMeans(n_clusters=3)
model.fit(X)

plt.figure(figsize=(14,7))

colormap = np.array(['red', 'lime', 'black'])

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

# Plot the Models Classifications
plt.subplot(1, 2, 2)
plt.scatter(X.Petal_Length, X.Petal_Width, c=colormap[model.labels_], s=40)
plt.title('K Mean Classification')
plt.xlabel('Petal Length')
plt.ylabel('Petal Width')
print('The accuracy score of K-Mean: ',sm.accuracy_score(y, model.labels_))
print('The Confusion matrixof K-Mean: ',sm.confusion_matrix(y, model.labels_))

from sklearn import preprocessing
scaler = preprocessing.StandardScaler()
scaler.fit(X)
xsa = scaler.transform(X)
xs = pd.DataFrame(xsa, columns = X.columns)
#xs.sample(5)

from sklearn.mixture import GaussianMixture
gmm = GaussianMixture(n_components=3)
gmm.fit(xs)

y_gmm = gmm.predict(xs)
#v cluster gmm
```

https://colab.research.google.com/github/niths1271/MACHINE-LEARNING-LAB/blob/main/LAB%207/EM_Kmeans_Comparison.ipynb#printMode=true 1/3

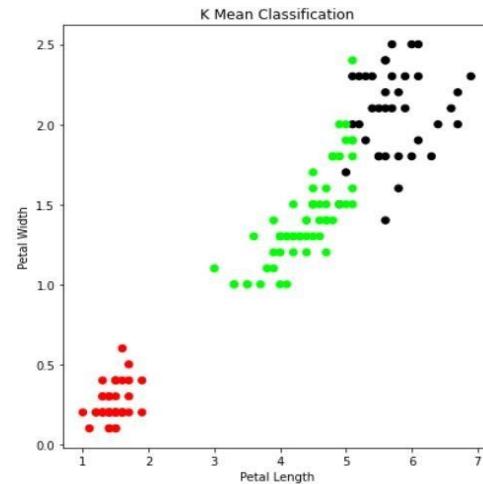
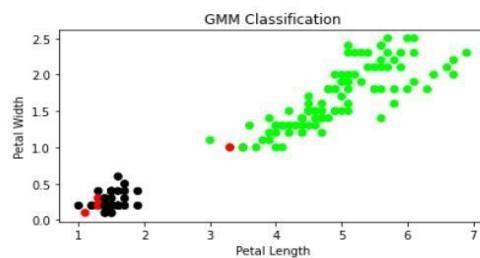
7/10/22, 11:05 PM
"j-----_B....."

EM_Kmeans_Comparison.ipynb - Colaboratory

```
plt.subplot(2, 2, 3)
plt.scatter(X.Petal_Length, X.Petal_Width, c=colormap[y_gmm], s=40)
plt.title('GMM Classification')
plt.xlabel('Petal Length')
plt.ylabel('Petal Width')

print('The accuracy score of EM: ',sm.accuracy_score(y, y_gmm))
print('The Confusion matrix of EM: ',sm.confusion_matrix(y, y_gmm))
```

👤 The accuracy score of K-Mean: 0.8933333333333333
The Confusion matrix of K-Mean: [[50 0 0]
[0 48 2]
[0 14 36]]
The accuracy score of EM: 0.3533333333333333
The Confusion matrix of EM: [[5 0 45]
[2 48 0]
[0 50 0]]



9. Write a program to implement k-Nearest Neighbour algorithm to classify the iris dataset. Print both correct and wrong predictions.

7/10/22, 11:05 PM

KNN(using libraries).ipynb - Colaboratory

```
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import classification_report, confusion_matrix
from sklearn import datasets

iris=datasets.load_iris()

x = iris.data
y = iris.target

# print ('sepal-length', 'sepal-width', 'petal-length', 'petal-width')
# # print(x)
# print('class: 0-Iris-Setosa, 1- Iris-Versicolour, 2- Iris-Virginica')
# # print(y)

x_train, x_test, y_train, y_test = train_test_split(x,y,test_size=0.3)

#To Training the model and Nearest neighbors K=5
classifier = KNeighborsClassifier(n_neighbors=5)
classifier.fit(x_train, y_train)

#To make predictions on our test data
y_pred=classifier.predict(x_test)

print('Confusion Matrix')
print(confusion_matrix(y_test,y_pred))
print('Accuracy Metrics')
print(classification_report(y_test,y_pred))
```

Confusion Matrix
[[13 0 0]
 [0 13 2]
 [0 1 16]]
Accuracy Metrics

	precision	recall	f1-score	support
0	1.00	1.00	1.00	13
1	0.93	0.87	0.90	15
2	0.89	0.94	0.91	17

	accuracy		0.93	45
macro avg	0.94	0.94	0.94	45
weighted avg	0.93	0.93	0.93	45

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

7/10/22, 11:23 PM

Copy of LocallyWeightedRegression.ipynb - Colaboratory

```
import numpy as np
from bokeh.plotting import figure, show, output_notebook
from bokeh.layouts import gridplot
from bokeh.io import push_notebook

def local_regression(x0, X, Y, tau):# add bias term
    x0 = np.r_[1, x0] # Add one to avoid the loss in information
    X = np.c_[np.ones(len(X)), X]

    # fit model: normal equations with kernel
    xw = X.T * radial_kernel(x0, X, tau) # XTranspose * W

    beta = np.linalg.pinv(xw @ X) @ xw @ Y #@ Matrix Multiplication or Dot Product

    # predict value
    return x0 @ beta # @ Matrix Multiplication or Dot Product for prediction

def radial_kernel(x0, X, tau):
    return np.exp(np.sum((X - x0) ** 2, axis=1) / (-2 * tau * tau))
# Weight or Radial Kernel Bias Function

n = 1000
# generate dataset
X = np.linspace(-3, 3, num=n)
print("The Data Set ( 10 Samples) X :\n",X[1:10])
Y = np.log(np.abs(X ** 2 - 1) + .5)
print("The Fitting Curve Data Set (10 Samples) Y :\n",Y[1:10])
# jitter X
X += np.random.normal(scale=.1, size=n)
print("Normalised (10 Samples) X :\n",X[1:10])

domain = np.linspace(-3, 3, num=300)
print(" Xo Domain Space(10 Samples) :\n",domain[1:10])

def plot_lwr(tau):
    # prediction through regression
    prediction = [local_regression(x0, X, Y, tau) for x0 in domain]
    plot = figure(plot_width=400, plot_height=400)
    plot.title.text='tau=%g' % tau
    plot.scatter(X, Y, alpha=.3)
    plot.line(domain, prediction, line_width=2, color='red')
    return plot

show(gridplot([
    [plot_lwr(10.), plot_lwr(1.)],
    [plot_lwr(0.1), plot_lwr(0.01)]]))
```

The Data Set (10 Samples) X :
[-2.99399399 -2.98798799 -2.98198198 -2.97597598 -2.96996997 -2.96396396
-2.95795796 -2.95195195 -2.94594595]
The Fitting Curve Data Set (10 Samples) Y :
[2.13582188 2.13156806 2.12730467 2.12303166 2.11874898 2.11445659
2.11015444 2.10584249 2.10152068]
Normalised (10 Samples) X :

7/10/22, 11:23 PM

Copy of LocallyWeightedRegression.ipynb - Colaboratory

```
[-3.12282223 -2.9216174 -3.14051918 -3.09805236 -3.08215798 -2.88090494  
-3.05412007 -3.12734019 -2.98129254]  
Xo Domain Space(10 Samples) :  
[-2.97993311 -2.95986622 -2.93979933 -2.91973244 -2.89966555 -2.87959866  
-2.85953177 -2.83946488 -2.81939799]
```



```
from numpy import *
from os import listdir
import matplotlib
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np1
import numpy.linalg as np
from scipy.stats.stats import pearsonr

def kernel(point,xmat, k):
    m,n = np1.shape(xmat)
    weights = np1.mat(np1.eye((m)))
    for j in range(m):
        diff = point - X[j]
        weights[j,j] = np1.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 = np1.shape(xmat)
    ypred = np1.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('/content/tips.csv')
bill = np1.array(data.total_bill)
tip = np1.array(data.tip)

#preparing and add 1 in bill
mbill = np1.mat(bill)
mtip = np1.mat(tip) # mat is used to convert to n dimesiona to 2 dimensional array form
m= np1.shape(mbill)[1]
# print(m) 244 data is stored in m
one = np1.mat(np1.ones(m))
X= np1.hstack((one.T,mbill.T)) # create a stack of bill from ONE
#print(X)
#set k here
ypred = localWeightRegression(X,mtip,2)
SortIndex = X[:,1].argsort(0)
xsort = X[SortIndex][:,0]

fig = plt.figure()
ax = fig.add_subplot(1,1,1)
```

7/10/22, 11:25 PM

WeightedRegression.ipynb - Colaboratory

```
ax.scatter(bill,tip, color='blue')
ax.plot(xsort[:,1],ypred[SortIndex], color = 'red', linewidth=5)
plt.xlabel('Total bill')
plt.ylabel('Tip')
plt.show()
```

