

# Lab-1

Date 21/3/24  
Page 1

⇒ Write a python program to import and export data using pandas library functions

Code:

```
import pandas as pd
iris_data = pd.read_csv("iris.csv")
iris_data.head()
iris.columns = col_names
```

	se	sw	pl	pw	class
0	5.1	3.5	1.4	0.2	Iris - setosa
1	4.9	3.0	1.3	0.2	Iris - setosa
2	4.7	3.2	1.4	0.2	Iris - setosa
3	4.6	3.1	1.5	0.2	Iris - setosa
4	5.0	3.6	1.4	0.4	Iris - setosa

Code:

```
url = "https://archive.ics.uci.edu/ml/machine-learning-
database/iris/iris.data"
```

```
col_names = ["sepal-length-in-cm",
             "sepal-width-in-cm",
             "petal-length-in-cm",
             "petal-width-in-cm",
             "class"]
```

```
iris_data = pd.read_csv(url, names=col_names)
iris_data.head()
```

~~Output:~~

	sepal-length-in-cm	sepal-width-in-cm	petal-length-in-cm	petal-width-in-cm	class
0	5.1	3.5	1.4	0.2	Iris-setosa

```
iris_data.to_csv("cleaned-iris-data.csv")
```

## Get the Data

→ Download the data

```
import os  
import tarfile  
import urllib
```

- Using these libraries, download and extract the housing data and load the data into "housing.csv".

```
housing.head()  
housing.info()  
housing.describe()
```

- Using the above commands, inspect the attributes of the loaded housing dataset.
- Using matplotlib and seaborn libraries, by plotting the histplot detect the outliers

## Create the Test Set

- Splitting the dataset on test ratio 0.2, i.e., training data is 80% of dataset and testing data is 20% of dataset
- Stratified Sampling is when random chosen data are representative of a whole target population. Each homogeneous subgroup is called strata.

Discover and visualize the Data to gain insights

Visualize the data using matplotlib and seaborn libraries.

Calculating the standard correlation coefficient of every pair of columns

Prepare the data for machine learning algorithms

Data cleaning, handling text and categorical data, custom transformers, feature scaling, transformation pipelines etc., are done here

Select and Train model

At first linear regression model is used to train but the model is overfitting the data.

To tackle this, DecisionTreeRegressor model is used as it is capable of finding non-linear relationships within the data.

But, the decision tree model is also overfitting so badly that it performs worse than the linear regression model.

At last Random forest regressor model is used. It is much better.

Fine tune your model

At last fine tuning of model is carried out, evaluating on the test set and then launch, monitor and maintaining the system.

→ python implementation of linear regression

import numpy as np

import matplotlib.pyplot as plt

def estimate\_coef(x, y):

n = np.size(x)

m\_x = np.mean(x)

m\_y = np.mean(y)

ss\_xy = np.sum(y\*x) - n\*m\_y\*m\_x

ss\_xx = np.sum(x\*x) - n\*m\_x\*m\_x

b\_1 = ss\_xy / ss\_xx

b\_0 = m\_y - b\_1 \* m\_x

return (b\_0, b\_1)

def plot\_regression\_line(x, y, b):

plt.scatter(x, y, color="m", marker="o", s=30)

y\_pred = b[0] + b[1]\*x

plt.plot(x, y\_pred, color="g")

plt.xlabel('x')

plt.ylabel('y')

def main():

x = np.array([0, 1, 2, ..., 9])

y = np.array([1, 3, 2, ..., 12])

b = estimate\_coef(x, y)

print(b)

plot\_regression\_line(x, y, b)

Output :

(b<sub>0</sub>, b<sub>1</sub>) = (1.2363..., 1.16964...)

## → multiple linear regression

```
from sklearn.model_selection import train_test_split  
import matplotlib.pyplot as plt  
import numpy as np  
from sklearn import datasets, linear_model,  
metrics
```

```
data_url = "http://lib.stat.cmu.edu/datasets/boston"  
raw_df = pd.read_csv(data_url, sep="\t+", skiprows=22,  
header=None)
```

```
x = np.hstack([raw_df.values[:, 2:], raw_df.values  
[1::2, :2]])
```

```
y = rawdf.values[1::2, 2]
```

```
x_train, x_test, y_train, y_test = train_test_split(x, y,  
test_size=0.4, random_state=1)
```

```
reg = linear_model.LinearRegression()  
reg.fit(x_train, y_train)
```

```
print("Coefficients:", reg.coef_)
```

```
print("Variance score: {}" .format(reg.score(x_test,  
y_test)))
```

```
plt.style.use('fivethirtyeight')
```

~~```
plt.scatter(reg.predict(x_train), reg.predict(x-  
train) - y_train, color="green",  
s=10, label="Train data")
```~~

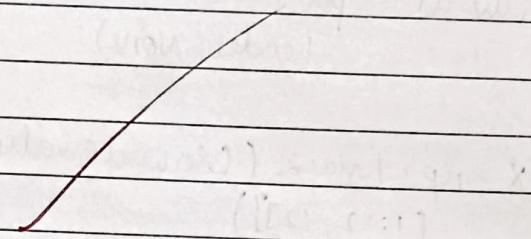
```
plt.scatter( reg.predict(x-test),  
            reg.predict(x-test) - y-test,  
            color = "green"blue, s=10,  
            label = "test data")
```

```
plt.hlines(y=0, xmin=0, xmax=50, linewidth=2)
```

```
plt.legend(loc='upper right')
```

```
plt.title("residual errors")
```

```
plt.show()
```



Decision tree - ID3

```

import numpy as np
import pandas as pd
eps = np.finfo(float).eps
from numpy import log2 as log
from google.colab import drive
drive.mount('/content/drive')
path = 'drive/My Drive/datasets/play-tennis.csv'
df = pd.read_csv(path)
df.head()
df = df.drop('day', axis=1)
  
```

```

def find_entropy(df):
    target = df.keys()[-1]
    entropy = 0
    values = df[target].unique()
    for value in values:
        fraction = df[target].value_counts()[value] / len(df[target])
        entropy += -fraction * np.log2(fraction)
    return entropy
  
```

```

def avg_information(df, attribute):
    target = df.keys()[-1]
    target_variables = df[target].unique()
    variables = df[attribute].unique()
    entropy2 = 0
    for variable in variables:
        entropy = 0
        for target_variable in target_variables:
            df_sub = df[(df[attribute] == variable) & (df[target] == target_variable)]
            entropy += len(df_sub) / len(df) * find_entropy(df_sub)
        entropy2 += len(df_sub) / len(df) * entropy
    return entropy2
  
```

$\text{num} = \text{len}(\text{df}[\text{attribute}][\text{df}[\text{attribute}] == \text{variable}])$   
 $\text{den} = \text{len}(\text{df}[\text{attribute}][\text{df}[\text{attribute}] == \text{target\_variable}])$   
 $\text{fraction} = \text{num}/\text{den}$   
 $\text{entropy} += -\text{fraction} * \log(\text{fraction} + \epsilon)$   
 $\text{fraction2} = \text{den}/\text{len}(\text{df})$   
 $\text{entropy2} += -\text{fraction2} * \text{entropy}$   
 return  $\text{abs}(\text{entropy2})$

```

def find_winner(df):
    IG = []
    for key in df.keys()[:-1]:
        IG.append(find_entropy(df) - avg_information)
    print(IG)
    info = df.keys()[:-1][np.argmax(IG)]
    print(info)
    return df.keys()[:-1][np.argmax(IG)]
    
```

```

def get_subtable(df, node, value):
    res = df[df[node] == value].reset_index(drop=True)
    return res
    
```

```

def build_tree(df, tree=None):
    target = df.keys()[-1]
    node = find_winner(df)
    attr_value = np.unique(df[node])
    if tree is None:
        tree = {}
        tree[node] = {}
    for value in attr_value:
        subtable = get_subtable(df, node, value)
        tree[node][value] = subtable
    return tree
    
```

```

for value in attr_value:
    subtable = get_subtable(df, node, value)
    class_label, counts = np.unique(subtable[target])
    if counts[0] == 1:
        tree[node][value] = class_label
    else:
        tree[node][value] = build_tree(subtable, tree)
    
```

```
if len(counts) == 1:  
    tree[node][value] = dvalue[0]
```

else:

```
    tree[node][value] = buildTree(subtable)
```

```
return tree
```

```
tree = buildTree(df)
```

```
import pprint
```

```
pprint.pprint(tree)
```

Output:

```
key))  
{'outlook': {'overcast': 'Yes',  
             'Rain': {'wind': {'strong': 'No', 'weak':  
                               'Yes'}}},  
  'sunny': {'humidity': {'high': 'No', 'Normal':  
                               'Yes'}}}
```

✓ 09-05-2023

## Logistic regression

```
import pandas as pd
```

#

```
data = pd.read_csv("cet-data.csv")
```

```
from sklearn.model_selection import train_test_split  
from matplotlib import pyplot as plt  
%matplotlib inline
```

```
plt.scatter(data['CET-Score'], data['Admitted'],  
marker='.', color='purple')
```

```
X_train, X_test, y_train, y_test = train_test_split(data  
[['CET-Score']], data['Admitted'], train_size=0.8)
```

```
from sklearn.linear_model import LogisticRegression  
model = LogisticRegression()  
model.fit(X_train, y_train)  
y_predicted = model.predict(X_test)  
model.score(X_test, y_test)  
print(y_predicted)  
print(X_test)
```

```
from sklearn.linear_model import LinearRegression  
model = LinearRegression()  
model.fit(X_train, y_train)  
print("coefficient (m):", model.coef_)  
print("Intercept (b):", model.intercept_)
```

import math

def Sigmoid(x):

return 1 / (1 + math.exp(-x))

x = 0

Sigmoid\_value = sigmoid(x)

print ("Sigmoid value at x = ", x, ":", sigmoid\_value)

import math

def prediction\_function (CET\_Score):

m = 0.042

b = -1.53

z = m \* CET\_Score + b

y = sigmoid(z)

return y

predicted\_probability = prediction\_function(500)

print ("predicted probability when CET-score is 500: "

predicted\_probability)

Output:

predicted probability when CET-score is 500: 0.9999

## KNN implementation:

```

import numpy as np
import pandas as pd
from google.colab import drive
drive.mount('/content/drive')
dataset = pd.read_csv('/content/drive/MyDrive/iris.csv')
dataset.head()

```

| id | sepal length | sepal width | petal length | petal width | Species     |
|----|--------------|-------------|--------------|-------------|-------------|
| 1  | 5.1          | 3.5         | 1.4          | 0.2         | Iris-setosa |
| 2  | 4.9          | 3.0         | 1.4          | 0.2         | Iris-setosa |
| 3  | 4.7          | 3.2         | 1.3          | 0.2         | Iris-setosa |
| 4  | 4.6          | 3.1         | 1.5          | 0.2         | Iris-setosa |
| 5  | 5.0          | 3.6         | 1.4          | 0.2         | Iris-setosa |

dataset.groupby('Species').size()

| Species         | Count |
|-----------------|-------|
| Iris-setosa     | 50    |
| Iris-versicolor | 50    |
| Iris-virginica  | 50    |

feature\_columns = ['sepal length', 'sepal width', 'petal length',  
                   'petal width']

X = dataset[feature\_columns].values

y = dataset['Species'].values

from sklearn.preprocessing import LabelEncoder  
 le = LabelEncoder()

y = le.fit\_transform(y)

from sklearn.model\_selection import train\_test\_split

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y,

test\_size=0.2, random\_state=0)

```
import matplotlib.pyplot as plt  
import seaborn as sns  
%matplotlib inline
```

```
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.metrics import confusion_matrix, accuracy_ score  
from sklearn.model_selection import cross_val_score
```

```
classifier = KNeighborsClassifier(n_neighbors=3)
```

```
classifier.fit(X_train, y_train)
```

```
y_pred = classifier.predict(X_test)
```

```
accuracy = accuracy_score(y_test, y_pred) * 100
```

```
print('Accuracy of our model is equal' +  
      str(round(accuracy, 2)) + '%!')
```

output: Accuracy of our model is equal 96.67%

# Lab-7

123/5/2024

## k-mean implementation

```
import pandas as pd  
data = pd.read_csv("iris.csv")  
data.head()
```

```
import numpy as np  
import seaborn as sns  
import matplotlib.pyplot as plt  
from sklearn.datasets import load_iris  
from sklearn.cluster import KMeans
```

```
x, y = load_iris(return_X_y=True)
```

```
sse = []
```

```
for k in range(1, 11):
```

```
KM = KMeans
```

```
Kmeans = KMeans(n_clusters=3, random_state=2)
```

```
Kmeans.fit(x)
```

```
pred = Kmeans.predict(x)
```

pred

## SVM :

```
from sklearn.datasets import load_breast_cancer
import matplotlib.pyplot as plt
from sklearn.inspection import DecisionBoundaryDisplay
from sklearn.svm import SVC
```

cancer = load\_breast\_cancer()

X = cancer.data[:, :2]

y = cancer.target

SVM = SVC(kernel="rbf", gamma=0.5, C=1.0)

SVM.fit(X, y)

DecisionBoundaryDisplay.from\_estimator(

SVM,

X,

response\_method="predict",

cmap=plt.cm.Spectral,

alpha=0.8,

xlabel=cancer.feature\_names[0],

ylabel=cancer.feature\_names[1],

)

~~plt.scatter(X[:, 0], X[:, 1], c=y, s=20, edgecolor='k')~~

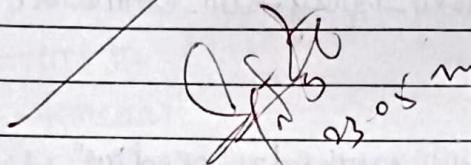
~~plt.show()~~

## PCA

```
import pandas as pd  
import numpy as np  
import matplotlib.pyplot as plt  
%matplotlib inline  
from sklearn.decomposition import PCA  
from sklearn.preprocessing import StandardScaler  
  
from sklearn.datasets import load_breast_cancer  
data = load_breast_cancer()  
data.keys()
```

~~print (data ['feature\_names'])~~

~~cb~~



```
df1 = pd.DataFrame (data ['data'], columns = data ['  
feature_names'])
```

scaling = StandardScaler()

scaling.fit (df1)

scaled\_data = scaling.transform(df1)

principal = PCA (n\_components=3)

principal.fit (scaled\_data)

x = principal.transform (scaled\_data)

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

plt.scatter (x[:, 0], x[:, 1], c = data ['target'], cmap='plasma')

plt.xlabel ('pc1')

plt.ylabel ('pc2')

## Random forest (classification)

Dataset: Titanic dataset

```
import pandas as pd
import warnings
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score, classification_report
warnings.filterwarnings('ignore')
```

url = "./content/titanic.csv"

titanic\_data = pd.read\_csv(url)

titanic\_data

X = titanic\_data[['pclass', 'sex', 'Age', 'sibsp', 'parch', 'Fare']]

y = titanic\_data['survived']

X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)

rf\_classifier = RandomForestClassifier(n\_estimators=100, random\_state=42)

rf\_classifier.fit(X\_train, y\_train)

y\_pred = rf\_classifier.predict(X\_test)

accuracy = accuracy\_score(y\_test, y\_pred)

print(f"Accuracy: {accuracy:.2f}")

Output: Accuracy: 0.80

## AdaBoost

```
from sklearn.datasets import load_iris  
from sklearn.ensemble import AdaBoostClassifier.  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score
```

```
iris = load_iris()
```

```
X = iris.data
```

```
y = iris.target
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,  
test_size=0.4, random_state=42).
```

```
adaBoost_clf = AdaBoostClassifier(n_estimators=30,  
learning_rate=1.0, random_state=42)  
adaBoost_clf.fit(X_train, y_train)
```

```
y_pred = adaBoost_clf.predict(X_test)
```

```
accuracy = accuracy_score(y_test, y_pred)  
print("Accuracy:", accuracy)
```

Output:

Accuracy: 0.9666667

Accuracy: 0.9666667

# Lab-9

Date 05/06  
Page 1

ANN:

```
import numpy as np
X = np.array(([2, 9], [1, 5], [3, 6]), dtype=float)
y = np.array(([92], [86], [89]), dtype=float)
X = X/np.amax(X, axis=0)
y = y/100
```

epoch = 5000

lr = 0.1

inputlayer\_neurons = 2

hiddenlayer\_neurons = 3

output\_neurons = 1

```
wh = np.random.uniform(size=(inputlayer_neurons, hiddenlayer_neurons))
```

```
bh = np.random.uniform(size=(1, hiddenlayer_neurons))
```

```
wout = np.random.uniform(size=(hiddenlayer_neurons, output_neurons))
```

```
bout = np.random.uniform(size=(1, output_neurons))
```

def sigmoid(x):

return 1/(1+np.exp(-x))

def derivatives\_sigmoid(x):

return x\*(1-x)

for i in range(epoch):

hinpl = np.dot(X, wh)

hinp = hinpl + bh

hlayer\_act = sigmoid(hinp)

outinp1 = np.dot(hlayer\_act, wout)

outinp = outinp1 + bout

output = sigmoid(outinp)

$$E0 = y - \text{output}$$

outgrad = derivatives\_sigmoid(output)

$$d\_output = E0 * outgrad$$

$$EH = doutput \cdot \text{dot}(wout.T)$$

hiddengrad = derivatives\_sigmoid(hlayer\_act)

$$d\_hiddenlayer = EH * hiddengrad.$$

$$wout += hlayer\_act \cdot T \cdot \text{dot}(d\_output) * lr$$

$$wh += X \cdot T \cdot \text{dot}(d\_hiddenlayer) * lr$$

print ("Input: \n" + str(x))

print ("Actual output: \n" + str(y))

print ("predicted output: \n", output)

Input:

[0.666667 1. ...]

[0.333333 0.555556]

[1. 0.6666667]

actual output:

[0.92]

[0.86]

[0.89]

predicted output

[0.86729245]

[0.845565]

[0.8640413]