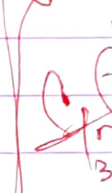
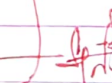
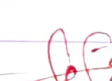


Sl. No.	Date	Title	Page No.	Teacher Sign / Remarks
1.	21/3/24	lab 1 — NP 21/3/2024	0-2	
2	28/3/24	ML-lab-2 exploration of Github Projection	2-4	
3	18/4/24	Linear regression Multiple regression	4-6 6-8	
4	25/4/24	Decision tree	8-10	
5	9/8/24	logistic regression KNN	10-12	
6	23/5/24	K-means svm PCA	12-14	
7.	30 may 24	random forest ensemble method Boosting ensemble method	14-16	
<p>Do upload everything — </p> <p>Completed everything — </p>				

lab:1

```
# import csv file from disc
```

```
import pandas as pd
```

```
airbnb_data = pd.read_csv("/content/cleaned-data.csv")
```

```
airbnb_data.head()
```

unnamed: 0	0	Sepal length in cm	Sepal width	Petal length	Petal width	class
0	0	5.1	3.5	1.4	0.2	iris-Setosa
1	1					
2	2					
3	3					
4	4					

```
# import from url
```

```
import pandas as pd
```

```
url = "https://archive.ics.uci.edu/ml/data"
```

```
# define columns
```

```
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()
```

	sepal-length	sepal-width	petal-length	petal-width	class
0	5.1	3.5	1.4	0.2	iris-Setosa
1					
2					
3					
4					

```
# export the dataframe to the csv file
```

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

ALD
21/2/2024

lab 1

```
# import csv file from disc
```

```
import pandas as pd
```

```
airbnb_data = pd.read_csv("/content/cleaned-data.csv")
```

```
airbnb_data.head()
```

Unnamed: 0	0	Sepal length in cm	Sepal width	Petal length	Petal width	class
0	0	5.1	3.5	1.4	0.2	iris-setosa
1	1					
2	2					
3	3					
4	4					

```
# import from url
```

```
import pandas as pd
```

```
url = "https://archive.ics.uci.edu/ml/data"
```

```
# define columns
```

```
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()
```

	sepal-length	sepal-width	petal-length	petal-width	class
0	5.1	3.5	1.4	0.2	iris-setosa
1					
2					
3					
4					

```
# export the dataframe to the csv file
```

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

21/5/2024

lab-2. HANDS ON LEARNING WITH SCIKIT-LEARN-KERN & TF

①. Framing the problem & looking at the bigger picture

about the data: California sales data to build model of housing prices in the state. features of the data include

↳ Population

MEAN

↳ Median income

↳ Median housing price for each blk group in California

- The problem statement involves creating or building a pipeline so that this model output can be fed into another model that attempt to increase the ROI of the company on investing on a given district
- As most of the data is numerical in nature, except one feature, the most obvious method of choice is linear regression & the performance measure chosen to visualize or interpret the data is root mean square

$x^{(i)}$ → is a vector containing all of the input feature values excluding labels

- also

$y^{(i)}$ is the desired output of input $x^{(i)}$

x matrix → contains all feature values excluding labels & targets

- describe ignores null values
- To plot histogram, you need to use `.hist()` with parameter, bin & fig size
- The `crc32` function converts a variable-length string into 8-character string that is a text representation of the hex value of 32-bit Binary seq

get the data

Download the data.

~~import os~~

import urllib

import urllib

use these diff lib to download & extract web-based dataset that is housing dataset

retrieve the data into housing dataset. Retrieve that data into 'housing-tgs' and read the data into housing.csv

housing.head()

housing.info()

housing.describe()

→ using these commands, inspect the attributes of the dataset.

import matplotlib.pyplot as plt

import seaborn as sns

housing.hist(bins=50, figsize=(20,15))

plt.show()

→ using the matplotlib & seaborn libraries detecting the outliers

lab-3

Create testset

- splitting the dataset on 1:1 ratio = 0.2 i.e., training data is 20% of the dataset
- Stratified Sampling is when random chosen data are representation of a whole target population. each homogeneous subgroup is called strata

Discover & validate the data to gain insight.

- Visualize the data using matplotlib & seaborn libraries
- Calculating the standard corr co-eff of every pair of columns

Prepare the data for machine learning algorithm.

- Data cleaning, handling text & categorical data, custom Transformers, feature scaling, transformation pipelines etc are done here

Select & train model.

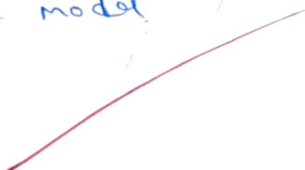
- at first linear regression model is used to train but the model is overfitting the data.
- To tackle this, Decision tree Regressor model is used as it is capable of finding non-linear relationships within the data. But the decision tree model is overfitting so badly that it performs worse than the linear regression model.
- At last Randomforest regressor model is used it is much better.

Finalize your model

- At last final tuning of model is carried out evaluating on the test set & then launch monitor and maintaining the system
- evaluate your system on the test set by using mean-square error

Launch Monitor & maintain your System

we can automate the process by

- ↳ collecting fresh data regularly & labelling it
 - ↳ writing script to train model to the finetune the hyper parameters
 - ↳ writing script to evaluate the model
- 

Week - 4

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, 3, 4, 5, 6, 7, 8, 9])
```

```
    y = np.array([1, 3, 2, 5, 7, 8, 9, 10, 12])
```

```
    b = estimate_coef(x, y)
```

```
    print("Estimate co-ef : \nb_0 = {} \nb_1 = {}".format(b))
```

output :

```
(b_0, b_1) = (1.4363, 1.6969)
```


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, multi_linear

data_w1 = 'w1'
row_of = pd.read_csv(data_w1, sep=';', skiprows=20, header=
    = None)
x = np.stack([row_of.values[:, 2, 1], row_of.values[:, 2, 2]])
y = row_of.values[:, 2, 2]
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=
    random_state=1)
reg = linear_model.LinearRegression()
reg.fit(x_train, y_train)
print("coeff =", reg.coef_)
print("variance score > 0.3", format(reg.score(x_test, y_test)))
plt.style.use('f5m thirty eight')
plt.scatter(reg.predict(x_train), reg.predict(x_train), y_train,
    wla="green", s=10, label='train data')
plt.scatter(reg.predict(x_test), reg.predict(x_test), y_test,
    color='blue', s=10, label='test data')
plt.xlim(y=0, x_min=0, x_max=0, x_max=50, figure=1)
plt.legend(loc='upper right')
plt.title('residual score')
plt.show()
```

week - 5

Decision tree ID3

```
import numpy as np
import pandas as pd

df = pd.read_csv("wt1")

df.head()
df.info()
df.describe()

def find_entropy(df):
    target = df.key()[-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 buildTree(df, tree = None):
    target = df.keys()[-1]
    node = find_minim(df)
    att = np.unique(df[node])
    if tree is None:
        tree = {}
        tree[node] = {}
    for value in att:
        sub = get_subtable(df, node, value)
        df_value, count = np.unique(sub[target], return_counts=True)
        if len(count) == 1:
            tree[node][value] = 0
        else:
            tree[node][value] = buildTree(subtable)
    return tree

tree = buildTree(df)
import pprint
pprint.print(tree)
```

Decision tree (sklearn):

```
import pandas as pd
```

```
import numpy as np
```

```
import sklearn.model_selection import DecisionTree  
from sklearn.tree import plot_tree
```

```
df = pd.read_csv('path')
```

```
df.head()
```

```
df.info()
```

```
df.isnull().sum()
```

```
cols = df.columns[0:-1]
```

```
for i in cols:
```

```
    sns.boxplot(y=df[i])
```

```
plt.show()
```

```
x = df.drop('species', axis=1)
```

```
y = df['species']
```

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

```
dt = DecisionTreeClassifier(max_depth=3)
```

```
dt.fit(x, y)
```

```
y_pred_train = dt.predict(x_train)
```

```
y_pred = dt.predict(x_test)
```

```
accuracy_score(y_pred, y_test)
```

Sp23
09-05-21

KNN.

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
```

```
df = pd.read_csv('prokto.csv')
```

```
df.head()
```

```
from sklearn.preprocessing import StandardScaler
```

```
Scaler = StandardScaler()
```

```
Scaler.fit(df.drop('Target', axis=1))
```

```
Scaled_features = scaler.transform(df.drop('Target', axis=1))
```

```
df_feat = pd.DataFrame(Scaled_features, columns=df.columns[:-1])
```

```
df_feat.head()
```

```
from sklearn import *
```

```
x_train, x_test, y_train, y_test = train_test_split(Scaled_features,
                                                    df['Target'],
```

```
                                                    test_size=0.50)
```

```
knn = KNeighborsClassifier(n_neighbors=1)
```

```
knn.fit(x_train, y_train)
```

```
pred = knn.predict(x_test)
```

```
print(confusion_matrix(y_test, pred))
```

```
print(classification_report(y_test, pred))
```

```
error_rate = []
```

```
for i in range(1, 40):
```

```
    knn = KNeighborsClassifier(n_neighbors=i)
```

```
    knn.fit(x_train, y_train)
```

```
    pred_i = knn.predict(x_test)
```

```
    error_rate.append(np.mean(pred_i != y_test))
```

```
plt.figure(figsize=(10, 6))
```

```
plt.plot(range(1, 40), error_rate, color='blue'
```

```
        linestyle='dashed', marker='o')
```


Logistic Regression

```
from sklearn.datasets import load_breast_cancer
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
```

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

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
                                                    random_state=23)
```

```
clf = LogisticRegression(random_state=0)
```

```
clf.fit(x_train, y_train)
```

```
y_pred = clf.predict(x_test)
```

```
acc = accuracy_score(y_test, y_pred)
```

```
print("Logistic Regression model accuracy (in %):", acc*100)
```

```
digits = datasets.load_digits()
```

```
x = digits.data
```

```
y = digits.target
```

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

```
reg = LinearModel.LogisticRegression()
```

```
reg.fit(x_train, y_train)
```

```
y_pred = reg.predict(x_test)
```

```
print("Logistic Regression model accuracy (in %):",
```

Signature
23.05.21

k-means clustering

import matplotlib.pyplot as plt

x = [4, 5, 10, 4, 3, 11, 14, 6, 10, 12]

y = [21, 19, 24, 17, 16, 25, 24, 22, 21, 21]

plt.scatter(x, y)

plt.show()



from sklearn.cluster import KMeans

data = list(zip(x, y))

inertia = []

for i in range(1, 11):

 kmeans = KMeans(n_clusters=i)

 kmeans.fit(data)

 inertia.append(kmeans.inertia_)

plt.plot(range(1, 11), inertia, marker='o')

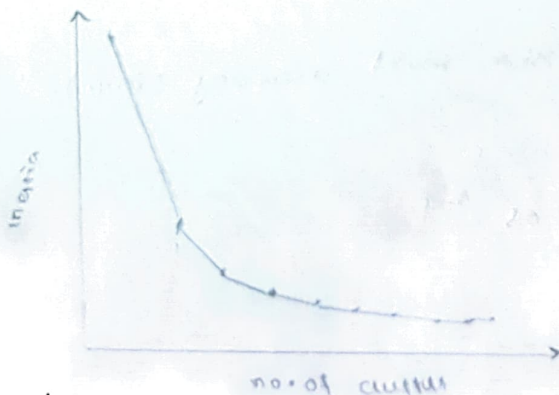
plt.title('Elbow method')

plt.xlabel('Number of clusters')

plt.ylabel('Inertia')

plt.show()

Elbow method

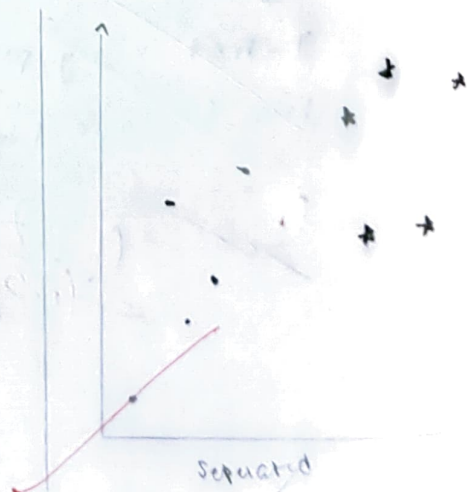


kmeans = KMeans(n_clusters=2)

kmeans.fit(data)

plt.scatter(x, y, c=kmeans.labels_)

plt.show()



Jan

How Principal Component Analysis (PCA)

```
import pandas as pd
```

```
import numpy as np
```

```
from sklearn.datasets import load_breast_cancer
```

```
cancer = load_breast_cancer (as_frame = True)
```

```
df = cancer.frame
```

```
print ('original dataframe shape : ', df.shape)
```

```
x = df[cancer.feature_names]
```

```
print ('input dataframe shape : ', x.shape)
```

```
x_mean = x.mean()
```

```
x_std = x.std()
```

```
z = (x - x_mean) / x_std
```

```
c = z.cov()
```

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

```
import seaborn as sns
```

```
sns.heatmap(c)
```

```
plt.show()
```

" compute eigen values & Eigen vectors

```
eigenvalues = np.linalg.eig(c)
```

```
print ('values : ', eigenvalues)
```

```
print ('values shape : ', eigenvalues.shape)
```

```
print ('Eigen vector shape : ', eigenvalues.shape)
```

" sort eigen values in descending order

```
idx = eigenvalues.argsort()[::-1]
```

```
eigenvalues = eigenvalues[idx]
```

```
eigenvalues = eigenvalues[:, idx]
```

```
explained_var = np.cumsum(eigenvalues) / np.sum(eigenvalues)
```

```
explained_var
```

" determine no. of principal components

```
n_components = np.argmax(explained_var > 0.95) + 1
```

n_components

A PCA component or unit matrix

$u = \text{eigenvector}[i, : n\text{-component}]$

`pca_comp = pd.DataFrame(u,`

`index = cases ['feature_name'],`

`columns = ['PC1', 'PC2']`

`)`

Plot heatmap

`plt.figure(figsize=(5,7))`

`sm = heatmap(pca - component)`

`plt.title('PCA component')`

`plt.show()`

find projection in PCA

`z_pca = z @ pca_component`

`z_pca.rename(columns={'PC1': 'PCA1', 'PC2': 'PCA2'}, axis=1, inplace=True)`

`print(z_pca)`

See

PCA using sklearn

from sklearn.decomposition import PCA

pca = PCA(n_components = 2)

pca.fit(z)

x_pca = pca.transform(z)

dt_pca = pd.DataFrame(x_pca,
 columns = ['PC1', 'PC2'],

 format = '%f')

for i in range(n_components):

 print(dt_pca)

plt.figure(figsize = (8, 6))

plt.scatter(x_pca[:, 0], x_pca[:, 1],

 c = cancer['target'],

 cmap = 'plasma')

plt.xlabel('first PCA')

plt.ylabel('second PCA')

plt.show()

~~pca.components -~~

Gen

SVM Support vector machine

```
from sklearn.datasets import load_data set breast_cancer.  
import matplotlib.pyplot as plt
```

```
from sklearn.inspection import DecisionBoundaryDisplay
```

```
from sklearn.svm import SVC
```

```
cancer = load_data set_cancer()
```

```
X = cancer.data[:, :2]
```

```
y = cancer[:, 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.cool,
```

```
    alpha = 0.8,
```

```
    xlabel = cancer.feature_names[0],
```

```
    ylabel = cancer.feature_names[1],
```

```
    plt.scatter(X[:, 0], X[:, 1], c = y, s = 20,
```

```
    edgecolor = 'x')
```

```
plt.show()
```

Gen

Lab-8

a) implement random forest ensemble method on a given dataset

b) Implement Boosting ensemble method on given dataset

c) import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score

from sklearn import datasets

iris = datasets.load_iris()

import pandas as pd

data = pd.DataFrame({

 'sepal length': iris.data[:, 0],

 'sepal width': iris.data[:, 1],

 'petal length': iris.data[:, 2],

 'petal width': iris.data[:, 3],

 'species': iris.target

{})

data.head() model = selection import train_test_split

x = data[['sepal length', 'sepal width', 'petal length', 'petal width',
 'th']]

y = data['species']

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

from sklearn.ensemble import RandomForestClassifier

clf = RandomForestClassifier(n_estimators=100)

clf.fit(x_train, y_train)

y_pred = clf.predict(x_test)

from sklearn import metrics

print("Accuracy: ", metrics.accuracy_score(y_test, y_pred))

confusion_matrix = metrics.confusion_matrix(y_test, y_pred)

print(confusion_matrix)

output ['setosa' 'versicolor' 'virginica']

['sepal length' 'sepal width' 'petal length'

Accuracy = 0.985

conf-mat = [[10 0 0]

[0 15 1]

[0 4 15]

Implement boosting ensemble method on a dataset

```
import pandas as pd
```

```
from sklearn.model_selection import train_test_split
```

```
from sklearn.ensemble import GradientBoostingClassifier
```

```
from sklearn.metrics import accuracy_score
```

```
cricket_df = pd.read_csv('content/cricket dataset.csv')
```

```
x = cricket_df.drop(columns=['player Name', 'player type'])  
y = cricket_df['player type']
```

```
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size  
= 0.2, random_state = 42)
```

```
boosting = GradientBoostingClassifier(n_estimators=100, learning_rate  
= 0.1, random_state=42)
```

```
boosting.fit(x_train, y_train)
```

```
boosting_predictions = boosting.predict(x_test)
```

```
boosting_accuracy = accuracy_score(y_test, boosting_predictions)
```

```
print("Boosting Accuracy:", boosting_accuracy)
```

Output

Boosting Accuracy: 0.35

Gen
CSB
20.05.24

