

```
import pandas as pd
import numpy as np
iris = pd.read_csv('Iris-flower-dataset.csv')
```

species	sepal.length.cm	sepal.width	petal.length	petal.width
0	5.1	3.5	1.4	0.2
1	5.2	3.7	1.5	0.2
2	5.4	3.9	1.6	0.2
3	5.5	4.0	1.7	0.2
4	5.6	4.1	1.8	0.2

```
# import from url
import pandas as pd
```

```
url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.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
5.1	3.5	1.4	0.2	iris-setosa
5.2	3.7	1.5	0.2	iris-setosa
5.4	3.9	1.6	0.2	iris-setosa
5.5	4.0	1.7	0.2	iris-setosa
5.6	4.1	1.8	0.2	iris-setosa

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

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

2/13/2024

## lab 2 HANDS ON LEARNING WITH SCIKIT LEARN VERSION 2.11

### 1. Framing the problem & looking at the bigger picture

about the data: California house 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 attempts 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

$y^{(i)}$  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 local web-based dataset that is housing dataset

retrieving 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=10, figsize=(10,10))

plt.show()

→ using the matplotlib & seaborn libraries detecting the outliers

### lab-3.

#### Create testset

- splitting the dataset on test 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 homogenous subgroup is called strata

#### Discover & validate the data to gain insights.

- 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.

#### fine tune your model.

- At last fine 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 method.

Launch Monitor to maintain your System

we can automate the process by

↳ collecting fresh data regularly & labelling it

↳ writing script to train model to the find tune 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-eff : \nb_0 = {} \nb_1 = {}".format(b))
```

Output :

(b\_0, b\_1) = (1.2363 ... 1.16969 ...)

## 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, multivariate
```

```
data_wt1 = 'wt1'
```

```
row_of = pd.read_csv('data_wt1', sep=';', skiprows=20, header=
```

```
= None)
```

```
x = np.stack([row_of.values[:, 2, 1], row_of.values[:, 2, 2],
```

```
y = row_of.values[:, 1, 2, 2])
```

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

```
random_state=1)
```

```
reg = LinearModel.LinearRegression()
```

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

```
print('coeff =', reg.coef)
```

```
print('variance score = %.3' % r_squared(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='blue', s=10, label='test data')
```

```
plt.legend(y=0, x_min=0, x_max=0, x_min=50, y_max=50)
```

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

```
plt.title('regression core')
```

```
plt.show()
```

week 4

Decision tree ID3

import numpy as np

import pandas as pd

df = pd.read\_csv('data.csv')

df.head()

df.info()

df.describe()

def find\_entropy(df):

target = df.key[0][1]

entropy = 0

value = df[target].unique()

for value in value:

fraction = df[target].value\_counts()[value] / len(df[target])

entropy += fraction \* np.log2(fraction)

return entropy

def buildTree(df, tree = None):

target = df.keys[0][1]

node = find\_miner(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(subtable[target], return\_counts=True)

if len(count) == 1:

tree[node][value] = df\_value[0]

else:

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

return tree

tree = buildTree(df)

import print



Decision tree (sukman)

```
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[6:-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)
```

*Spl*  
*23*  
*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('prophet.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'],
```

```
knn = KNeighborsClassifier(n_neighbors=1, test_size=0.50)
```

```
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("Logic Regression model accuracy (in %):")
```

*Sc*  
*Ans*  
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, 23, 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

    kmeans = kmeans(n\_clusters=i)

    kmeans.fit(data)

    inertia.append(kmeans.inertia\_)

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

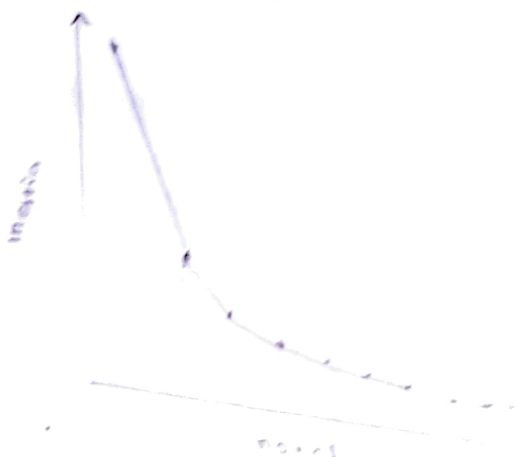
plt.title('Elbow method')

plt.xlabel('Number of clusters')

plt.ylabel('Inertia')

plt.show()

## Elbow method



```

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()

```

4 compute eigen values & Eigen vector

```

eigenvalues = np.linalg.eig(c)

```

```

print ('values : ', eigenvalues)

```

```

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

```

```

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

```

5 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

```

6 determine no. of principal components

```

n_comp = np.argmax(explained_var > 0.95) + 1

```

```

n_comp

```

*Done*

A PCA component or unit matrix

$u = \text{eigenvector } U^T, \text{ '0-component'}$

$PCA_{comp} = P \cdot D \cdot K \cdot Trans(u)$

$invar = covar(U^T \cdot Trans(X_{train}))$

$covar = [U^T P(1), U^T P(2)]$

)

# Plot weights

$PU$  figure (figure = (5,7))

$SM$  : heatmap (PCA - component)

$fit\_hml(U^T PCA_{comp} = u')$

$plt.show()$

# find projection in PCA

$Z \cdot PCA = Z \otimes PCA$  - component

$Z \cdot PCA_{comp} = [U^T P(1), U^T P(2), U^T P(2), \dots, U^T P(2)]$ ,  $axis=1$ ,  $hplow=1$

Gen



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',  
                      index = 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('2nd PCA')

plt.show()

pca.components\_

*Sen*



## SVM Support Vector Machine

from sklearn.svm import SVC, LinearSVC  
import matplotlib.pyplot as plt

from sklearn.svm import SVC

from sklearn.svm import SVC

svm = SVC(kernel='rbf', gamma=0.5)

X = concv : data[:, :2]

y = ~~conc~~ : data[:, 2]

svm = SVC(kernel='rbf', gamma=0.5)

svm.fit(X, y)

Decision boundary display from sklearn

svm

X,

response\_method = 'predict',

color = plt.cm.Spectral,

alpha = 0.8,

alpha = concv, feature\_names = ['x1', 'x2'],

label = concv, feature\_names = ['x1', 'x2'],

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

edgecolor='r',

plt.show()

gpr

## lab-8

- 3) implement random forest ensemble method on given dataset  
4) implement Boosting algorithm on given dataset

a) 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
```

f)

data.head(), data.tail()

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

y = data['species']

x\_train, x\_test, y\_train, y\_test = train\_test\_split(x, y, test\_size=0.3, random\_state=42)

clf = RandomForestClassifier(n\_estimators=100)

clf.fit(x\_train, y\_train)

y\_pred = clf.predict(x\_test)

from sklearn.metrics

print("Accuracy: ", metrics.accuracy\_score(y\_test, y\_pred))

Confusion matrix = metrics.confusion\_matrix(y\_test, y\_pred)

print(confusion\_matrix)

Output: ['setosa', 'versicolour', 'virginica']

['sepal length', 'sepal width', 'petal length', 'petal width']

Accuracy ~ 0.988

conf-mat =

[[0 0 0]

[0 15 1]

[0 0 15]]



# Implement Boosting

import pandas as pd  
from sklearn.model\_selection import train\_test\_split  
from sklearn.metrics import accuracy\_score  
from sklearn.metrics import roc\_auc\_score

cricket\_d1 = pd.read\_csv('cricket/cricket\_data.csv')

n = cricket\_d1.shape[0] # column = ['Player Name', 'Player Type']

y = cricket\_d1['Player Type']

n\_train, n\_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.predict(x\_test)

boosting.score = accuracy\_score(y\_test, boosting.predict(x\_test))

print("Boosting Accuracy: ")

boosting.score

Boosting Accuracy: 0.85

20.05.21