

# INDEX

NAME: Shikha Singh STD.: VI SEC.: D ROLL NO.: 202

## Week-1

### PANDAS (Reading CSV)

import pandas as pd

df2 = pd.read\_csv("/content/sample\_data/california\_housing-test.csv")

df2.head()

Output :

	longitude	latitude	housing_medianage	total_rooms	total_bath
0	-122.05	37.37	270	3645.0	661.0
	population	households	median_income	median_house_value	
	1537.0	606.0	6.6085	344700.0	

df2.to\_csv('w:csv')

colnames = ['sepal-length-incm', 'sepal-width-in.cm',  
 'petal-length-in-cm', 'petal-width-in(cm)', 'class']

wrl = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"

df = pd.read\_csv(wrl, names = colnames)  
 df.head()

OUTPUT:

sepal-length-in-cm	sepal-width-in-cm	petal-length-in-cm
5.1	3.5	1.4
0.2	0.2	0.2

class  
Iris-setosa

Week-2

1. Performance measure
  2. Get the data
  - 2.1 Download the data
- Step 1: Use the data
- ```
import os
import tarfile
import wget
```

Download Root = "https://seawater.usgs.gov/datasets/hudson-north/"

Housing\_PTH = os.path.join("data", "01")

Housing\_URL = Download\_Root + "datasets/housing/housing.tgz"

df. fetch\_housing\_data (housing\_url = HOUSING\_URL,  
housing\_path = HOUSING\_PATH):

os.makedirs(name=housing\_path, exist\_ok=True)

tgz\_path = os.path.join(housing\_path, "housing.tgz")

wget.download(url=housing\_url, out=housing\_path,  
filename=tgz\_path)

housing\_tgz = tarfile.open(name=tgz\_path)

housing\_tgz.extractall(path=housing\_path)

housing\_tgz.close()

fetch\_housing\_data()

import pandas as pd

df = load\_housing\_data ( housing\_path = Housing\_path):

data\_path = os.path.join(housing\_path, "housing.csv")

return pd.read\_csv(data\_path)

2) Take a quick look at the data structure

housing = load\_housing\_data()

housing.head()

housing.info()

housing[["ocean\_proximity"]].value\_counts()

housing.describe()

What's

3) PLOTTING DRAFT

- import matplotlib.pyplot as plt
- import seaborn as sns
- housing.hist(bins=50, figsize=(20, 15))
- plt.show()

4) Create a Test Set

import numpy as np

def split\_train\_test(data, test\_ratio=0.2):

shuffled\_indices = np.random.permutation(len(data))

test\_set\_size = int(len(data) \* test\_ratio)

test\_indices = shuffled\_indices[:test\_set\_size]

train\_indices = shuffled\_indices[test\_set\_size:]

return data.iloc[train\_indices], data.iloc[test\_indices]

divide & visualize

~~train\_set, test\_set = np.split(data, train\_test(data=housing))~~

~~len(train\_set), len(test\_set)~~

~~from zlib import~~

~~from zlib import crc32~~

def test\_set\_check(identifier, test\_ratio=0.2):

~~total\_size = 2\*\*32~~

~~hex\_repr = crc32(np.int64(identifier)) & 0xffffffff~~

`in-test = here_repair < (test_rate * total_size)`  
return in-test,  
(test\_set\_check('i') for i in range(10))]

Q.S) Coordinates to Index

```
def from_z_to_n(z):  
    if z >= 0:  
        n = 2*z  
    else:  
        n = -2*z - 1  
    return n
```

```
def cantor_pairing(n1, n2):  
    n = ((n1 + n2) * (n1 + n2 + 1)) // 2 + n2  
    return n
```

```
def lat_lon_to_index(lat, lon):  
    lat, lon = int(lat * 100), int(lon * 100)  
    lat, lon = from_z_to_n(lat), from_z_to_n(lon)  
  
    index = cantor_pairing(lat, lon)  
    return np.int64(index)
```

newing['id'] = newing.apply(lambda row: lat\_lon\_to\_index  
(row['latitude'], row['longitude']), axis=1)

newing['id'].value\_counts()

~~week 3~~

Part A

Step 3

prepare the data for Mr Algo

Data Cleaning

```
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='median')
housing_num = housing.drop(['ocean_proximity'], axis=1)
imputer.fit(housing_num)
```

Handling text and categorical attributes

```
housing_cat = housing[['ocean_proximity']]
housing_cat = pd.get_dummies(housing_cat)
housing_cat['value_counts()']
```

Customer Transformer

Class combined Attributes added (BaseEstimator, TransformerMixin):

```
def __init__(self, add_bedrooms_per_room=True):
    self.add_bedrooms_per_room = add_bedrooms_per_room
def fit(X, y=None):
    return self
```

~~def transform(self, X, y=None):~~

~~rooms\_per\_household = X[:, room\_ix] / X[:, household\_ix]~~

~~population\_per\_household = X[:, population\_ix] / X[:, household\_ix]~~

~~if self.add\_bedrooms\_per\_room:~~

~~bedrooms\_per\_room = X[:, bedroom\_ix]] / X[:, room\_ix]~~

~~return np.c[X, room\_per\_household, population\_per\_household, per~~

~~household, bedrooms\_per\_room]~~

~~else:~~

return np.c\_[X, scores - per-household, 'Reputation per household']

### Transformation pipelines

num\_pipeline = Pipeline([('imputer', SimpleImputer(strategy='median')), ('attribs\_adder', CombinedAttributesAdder()), ('stdscaler', StandardScaler())])  
housing\_num\_to = num\_pipeline.fit\_transform(housing\_num)

housing\_num\_to.shape

### Select and Train a Model

clf display\_scores(scores):

print("scores:", scores)

print("Mean", scores.mean())

print("Standard deviation", scores.std())

### Fit Best Model

param\_grid = [{}{'n\_estimators': [3, 10, 30], 'max\_features': [2, 4, 6, 8]}, {'bootstrap': [False], 'n\_estimators': [3, 5], 'max\_features': [2, 3, 4]}]

~~Go back to your system on the test set.~~

final\_model = grid\_search.best\_estimator

X\_test = data\_test.re.drop(['label = median\_house\_value', 'incis = 1'])

y\_test = data\_test['median\_house\_value'].copy()

X\_test\_prepared = full\_pipeline.transform(X=X\_test)

final\_predictions = final\_model.predict(X=X\_test\_prepared)

final\_mse = mean\_squared\_error ( $y_{true} = y_{test}$ ,  $y_{predicted} = final\_prediction$ )

final\_rmse =  $\sqrt{final\ mse}$   
final\_rss.

import pandas as pd

import numpy as np

import matplotlib.pyplot as plt

import seaborn as sns

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

from pandas.core.common import linear regression.

df\_sal = pd.read\_csv ('/content/drive/MyDrive/Scalor.  
data/salaryData.csv')

df\_sal.head()

plt.title ('Salary Distribution Plot')

sns.distplot (df\_sal['salary'])

plt.show()

plt.scatter (df\_sal['yrsExperience'], df\_sal['Salary'])

plt.title ('Salary vs Experience')

plt.xlabel ('years of Experience')

plt.ylabel ('Salary')

plt.box (false)

plt.show()

## \* Splitting variables

x = df\_sal.iloc[:, :-1]

y = df\_sal.iloc[:, -1:]

x-train, x-test, y-train, y-test = train-test-split  
(x, y, test\_size=0.2, random\_state=0)

regressor = linear regression  
regressor = fit (x-train, y-train)

y-pred-test = regressor.predict(x-test)

y-pred-train = regressor.predict(x-train)

plt.scatter(x-train, y-train, color='lightblue')

plt.plot(x-train, y-pred-train, color='darkblue')

plt.title('Salary vs Experience (Training set)')

plt.xlabel('years of experience')

plt.ylabel('Salary')

plt.legend(['x-train', 'y-test', 'x-train/y-train'])

title = 'Sal vs Exp'; loc = 'best', facecolor='white')

plt.show()

plt.show()

Print (of coefficients: regressor.coef\_?)

Print (of intercept: regressor.intercept\_?)

Output:

coefficient: [(9312.57512673,)]

Intercept: [26780.09915063])

DY/M

Date: \_\_\_\_\_  
Page: 9

## Week - 4 (Decision Tree)

import numpy as np

import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt

%matplotlib inline

import sklearn.datasets as datasets

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

from sklearn.tree import DecisionTreeClassifier

from sklearn.metrics import accuracy\_score,

roc\_auc\_score, roc\_curve

from sklearn.tree import plot\_tree

from sklearn.tree import DecisionTreeClassifier

from sklearn.model\_selection import GridSearchCV

RandomizedSearchCV

from sklearn.metrics import confusion\_matrix

from sklearn.metrics import classification\_report

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

df = pd.read\_csv(url, header=None, names=['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)', 'Species'])

~~df.head()~~

df.info()

df.isnull().sum()

df.describe().T

cols = df.columns[0:-1]

for i in cols:

sns.boxplot(y=df[i])

plt.show()

$q_1 = 0.4$  [ 'sepal width (cm)' ]. quantile (0.25)

$q_3 = 0.7$  [ 'sepal width (cm)' ]. quantile (0.75)

$$iqr = q_3 - q_1$$

$df = df [ (df [ 'sepal width (cm)' ] \geq q_1 + 1.5 * iqrs) \& (df [ 'sepal width (cm)' ] <= q_3 + 1.5 * iqrs) ]$

df.shape

sns. boxplot ( y = df [ 'sepal width (cm)' ] )

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, random\_state = 1)

dt = DecisionTreeClassifier ( max\_depth = 3, min

samples\_leaf = 10, random\_state = 1 )

dt.fit ( X, y )

from sklearn

from IPython. display import Image

from sklearn.tree import export\_graphviz

import pydotplus

features = X.columns

dot\_data = export\_graphviz ( dt, out\_file = None, feature\_names = features )

graph = pydotplus. graph\_from\_dot\_data ( dot\_data )

Image ( graph. create\_png () )

gini = 0.0  
samples = 47  
value = [47, 0, 0]

Petal width (cm) <= 0.8  
gini = 0.666  
samples = 146  
value = [47, 49, 50]

False  
Petal width (cm) <= 1.75  
gini = 0.5  
samples = 99  
value = [50, 49, 47]

(0.25)  
0.75)

= 91-1.5 →  
93+1.5 →  
1.75)

st split  
→)

3 min

wiz)

e = None,

use data)

se width  
(n) <= 1.75

= 0.5

b41=99  
(0,1,45)

dt = DecisionTreeClassifier(random\_state=1)

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

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

y\_prob = dt.predict\_proba(x\_test)

print('Accuracy of Decision Tree Train:', accuracy\_score(y\_pred\_train, y\_train))

print('Accuracy of Decision Tree Test:', accuracy\_score(y\_pred, y\_test))

print(classification\_report(y\_test, y\_pred))

dt = DecisionTreeClassifier(random\_state=1)

params = {'max\_depth': [2, 3, 4, 5]}

'min\_samples\_split': [2, 3, 4, 5]

'min\_samples\_leaf': [1, 2, 3, 4, 5]

gsearch = GridSearchCV(dt, param\_grid=params, cv=3)

gsearch.fit(X, y)

gsearch.best\_params\_

dt = DecisionTreeClassifier(\*+ gsearch.best\_params\_, \*  
random\_state=1)

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

~~y\_pred\_train = dt.predict(x\_train)~~

~~y\_prob\_train = dt.predict\_proba(x\_train)[:, 1]~~

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

y\_prob = dt.predict\_proba(x\_test)[:, 1]

print('confusion matrix - Train:', '\n',

confusion\_matrix(y\_train, y\_pred\_train))

print('In'; 'confusion matrix - Test:', '\n',

12

confusion-matrix ( g-fest, y-pred )

print(classification\_report(y-test, y-pred))

part ('Accuracy of Decision' (Gee-Town: ),

accuracy score ( $\hat{y}$ -predtrain,  $y$ -train))

point (\* Accuracy of Decision Tree-Test: )

accuracy score(y-pred, y-test))

Accuracy of Decision Rule from 20.9901960.

Accuracy of Decision Tree-Test: 0.9848%

18

Week-5

import pandas as pd  
 from matplotlib import pyplot as plt  
 %matplotlib inline

```
df = pd.read_csv("insurance_data.csv")
df.head()
```

```
plt.scatter(df.age, df.bought_insurance, marker='+', color='red')
```

from sklearn.model\_selection

import train\_test\_split

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

(df[['age']], df.bought\_insurance, train\_size=0.2)

from sklearn.linear\_model

import LogisticRegression

model = LogisticRegression()

model.fit(X\_train, y\_train)

y\_predicted = model.predict(X\_test)

model.predict\_proba(X\_test)

model.score(X\_test, y\_test)

→ accuracy → 1.0

~~model.coef\_~~  
~~model.intercept\_~~ ] h-R

import math

def sigmoid(x):

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

def prediction function (age):

$$z = 0.042 \times \text{age} - 1.5396 \approx 0.04130133 \approx 0.042$$

and  $-1.52726963 \approx -1.53$

$y = \text{Aigmoid}(z)$

return y

age 35

prediction function (age)  $\rightarrow 0.485$

age = 43

prediction function (age)  $\rightarrow 0.569$

$w_0$   
 $x^{(0)}$   
 $y^{(0)}$

## Week - 6

Build KNN classification model for a given dataset.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
plt.style.use('ggplot')
```

```
df = pd.read_csv('~/content/diabetes.csv')
df.head()
```

df.shape

x = df.drop(['Outcome'], axis=1).values

y = df['Outcome'].values

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

x\_train, x\_test, y\_train, y\_test = train\_test\_split

(x, y, test\_size=0.4, random\_state=42, stratify=y)

```
from sklearn.neighbors import KNeighborsClassifier
```

neighbours = np.arange(1, 9)

train\_accuracy = np.empty(len(neighbours))

test\_accuracy = np.empty(len(neighbours))

for i, k in enumerate(neighbours):

Knn = KNeighborsClassifier(n\_neighbors=k)

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

train\_accuracy[i] = Knn.score(x\_train, y\_train)

test accuracy [i] = knn.score(x-test, y-test)

all-else ('KNN varying number of neighbours')

plt.plot(neighbours, test\_accuracy)

for i, k in enumerate(neighbours):

Knn = KNeighborsClassifier(n\_neighbors=k)

Knn.fit(x-train, y-train)

->

plt.scatter(neighbours, train\_accuracy, color='k')

plt.show()

Knn = KNeighborsClassifier(n\_neighbors=7)

Knn.fit(x-train, y-train)

Knn.score(x-test, y-test)

Output:

0.730519

week

week - 7

SVM Model

```
from sklearn import datasets
cancer = dataset.load_breast_cancer()
print(cancer.target)
from sklearn.model_selection
import train_test_split
```

```
X_train, X_test, y_train, y_test = train_test_split
(cancer.data, cancer.target, test_size=
0.3, random_state = 109)
```

```
from sklearn import SVM
clf = SVM.SVC(kernel='rbf')
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))
```

Output :

Accuracy : 0.9649

Q2

Q3

## week-8 Backpropagation Neural Network

```
db = np.loadtxt('.. / input/diabetes ..')  
db = np.random.shuffle(db)
```

```
y = db[:, 0]  
n = np.delete(db, [0], axis=1)
```

```
x_train, n-test, y-train, y-test = train_test_split  
(x, y, test_size=0.1)
```

```
print(np.shape(x_train), np.shape(x-test))
```

```
hidden_layer = np.zeros(72)
```

```
weights = np.random.random((len(n[0]), 72))
```

```
output_layer = np.zeros(2)
```

```
hidden_weights = np.random.random((72, 2))
```

```
def sum_function(weights, index_col_n):  
    result = 0
```

```
for i in range(0, len(n)):
```

```
    result += n[i] * weights[i][index_weight_col]
```

```
return result
```

```
def activate_layer(layer, weight, n):
```

```
for i in range(0, len(layer)):
```

```
layer[i] = 1.7789 * np.tanh((2.0 *
```

```
sum_function(weight, 1, n) / 3.0)
```

```
def backpropagation(hidden_layer, output_layer, one-hot  
encodding, learning_rate, n):
```

```
output_derivation = np.zeros(2)
```

```
output_gradient = np.zeros(2)
```

for i in range(0, len(output\_layer)):

$$\text{output\_derivation}[i] = (1.0 - \text{output\_layer}[i]) * \text{output\_layer}[i]$$

for i in range(0, len(output\_layer)):

output\_gradient[i] : output\_derivation

$$[i] * (\text{one-hot\_encoding}[i] - \text{output\_layer}[i])$$

$$\text{hidden\_derivative} = \text{np.zeros}(72)$$

$$\text{hidden\_gradient} = \text{np.zeros}(72)$$

for i in range(0, len(hidden\_layer)):

$$\text{hidden\_derivative}[i] = (1.0 - \text{hidden_layer}[i]) * (1.0 + \text{hidden_layer}[i])$$

for i in range(0, len(hidden\_layer)):

$$\text{sum}^- = 0$$

for j in range(0, len(output\_gradient)):

$$\text{sum}^+ = \text{output\_gradient}[j] * \text{hidden\_weight}[i,j]$$

$$\text{hidden\_gradient}[i] = \text{sum}^+ * \text{hidden\_derivative}[i]$$

Vocabulary (learning\_rate, weight\_hidden\_gradient, n).

Random Boosting and Ada Boosting

dataset = pd.read\_csv(" ")

corr\_matrix = input\_df.corr().abs()

upper = corr\_matrix.where(np.triu)

np.ones(corr\_matrix.shape) \* 1).astype(np.float)

to\_drop = [column for column in upper columns

if any upper[column] > 0.95)]

print(to\_drop)

\* train, \* test, y\_train, y\_test =

Train Test Split ( $x, y$ , testsize = 0.1)  
n\_classes = 7

n\_estimators = 100

Random\_state = 13

names = ("RandomForestClassifier", "AdaBoostClassifier")

models = [RandomForestClassifier(n\_estimators = n\_estimators),  
AdaBoostClassifier(DecisionTreeClassifier(max\_depth = None), n\_estimators = n\_estimators)]

for counter Model, model in enumerate  
(models): model.fit(x\_train, y\_train)

$y_{pred} = \text{model.predict}(x_{test})$

Accuracy RandomForestClassifier: 0.768439

Accuracy AdaBoostClassifier: 0.7841337

✓  
plm

Week 9 - Build k-Means to cluster a set of  
stocal data in .csv file

```
iris = pd.read_csv(".../IRIS.csv")
x = iris.iloc[:, [0, 1, 2, 3]].values
iris_outcome = pd.DataFrame(index=iris['Species'],
                             columns=['Count'])
```

```
iris_outcome
sns.FacetGrid(iris, hue="Species").maps(sns.
    distplot, "petal.length").add_legend()
```

```
sns.FacetGrid(iris, hue="Species").map(sns.
    distplot, "petal.width").add_legend()
sns.FacetGrid(iris, hue="Species").map(plt.distplot,
    "sepal.length").add_length()
plt.show()
```

```
sns.boxplot(x="Species", y="Petal.Length",
             data=iris)
```

```
plt.show()
```

```
sns.set_style("whitegrid")
```

```
sns.pairplot(iris, hue="Species", size=3),
plt.show()
```

```
KMeans = KMeans(n_clusters=3,
                  init='k-means++', max_iter=300,
                  n_init=10, random_state=0)
```

```
y_means = KMeans.fit_predict(x)
```

## Program: 11 PCA

```
df = pd.read_csv("../data.csv")
df_feature = df.drop(['diagnosis'], axis=1)
```

from sklearn.preprocessing.

import StandardScaler

Standardized = StandardScaler()

Standardized.fit(df\_features)

Scaled\_data = Standardized.transform(df\_features)

from sklearn.decomposition

import PCA

pca = PCA(n\_components=3)

pca.fit(Scaled\_data)

xpca = pca.transform(Scaled\_data)

Scaled\_data.shape

xpca.shape

def diag(x):

if len(x) == 'M':

return 1

else:

return 0

df\_diag = df['diagnosis'].apply(diag)

x\_pca[:, 1]

df\_pc = pd.DataFrame(pca.components\_)

columns = df\_features.columns

plt.figure(figsize=(15, 8))

sns.heatmap(df\_pc, cmap='viridis')

plt.title('Principal components correlation  
with the features')