

## INDEX

Name	SARAAG	Sub.	ML lab
Std.:	6C	Div.	Roll No. 190
Telephone No.		E-mail ID.	
Blood Group.		Birth Day.	

- 1) Python program to import and export data using pandas library.

```
import pandas as pd
airb-data = pd.read_csv("data/listing-austin.csv")
airb-data.head()
```

O/P:-

	id	name	host-id	host-name	neighbourhood-group	neighbourhood
2065		NW Austin	2466	paddy	NULL	787022
	latitude	longitude	room-type	price		
	30.27752	-97.7137	private	179		

- 2) Reading data from the URL

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

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

# Read data from URL

```
iris-data = pd.read_csv(url, names=col-names)
```

iris-data.head()

	<u>0/p</u>	sepal-length-in-cm	sepal-width-in-cm	petal-length-in-cm	petal-width-in-cm
0		5.1	3.5	1.4	0.2
1		4.9	3	1.4	0.2
2		4.7	3.2	1.3	0.2
3		4.6	3.1	1.5	0.2
4		5	3.5	1.9	0.2

class

iris-setosa

iris-setosa

iris-setosa

iris-setosa

iris-setosa

3) Exporting data frame to a csv file

~~iris-data.to.csv("cleaned-iris-data.csv")~~

## END TO END ML

\* Get the data

- ① import os  
import tarfile  
import urllib
- ② Download-root = "https://raw.githubusercontent.com/ageron/m12/madel"  
housing-path = os.path.join("data", "01")  
housing-url = Download-root + "datasets/housing/housing.tgz"
- ③ def 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")  
 urllib.request.urlretrieve(url=housing-url, filename=tgz-path)  
 housing\_tgz = tarfile.open(name=tgz-path)  
 housing\_tgz.extractall(path=housing-path)  
 housing\_tgz.close()
- ④ fetch\_housing\_data()
- ⑤ import pandas as pd  
def load\_housing\_data(housing\_path=housing-path):  
 data\_path = os.path.join(housing-path, "housing.csv")  
 return pd.read\_csv(data-path)
- ⑥ housing = load\_housing\_data()
- ⑦ housing.head()

\* Create a Test Set

```
① import numpy as np
② 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]
    return data.iloc[train_indices], data.iloc[test_indices]
```

③ train-set, test-set = split-train(data=housing)  
len(train-set), len(test-set)

Olp- (16512, 4128)

④ from zlib import crc32  
def test-set-check(identifier, test\_ratio=.2);  
 total\_size = 2\*\*32  
 hex\_repr = crc32(np.int64(identifier)) & 0xffffffff  
 in\_test = hex\_repr < (test\_ratio \* total\_size)  
 return in\_test

⑤ def split\_train\_test\_by\_id(data, test\_ratio, id\_column):

```
ids = data[id-column]
```

```
in-test-set = idr.apply(lambda id: test_set.check(id, tot_rols))
```

return data.loc[nint-test-set], data.loc[in-test-set]

⑥ train-set, test-set = split-train-test-by-id (data = housing-with-id)  
test-ratio = 0.2, id-column = "index"  
train-set.shape, test-set.shape

This step divides data set into test and train sets with 20% of data in test dataset. The indices are selected randomly.

Week -3

\* Describe & Visualize the Data to Gain Insight

① housing = sbat.trainset.copy(); housing.shape  
(16512, 10)

② housing.plot(kind='scatter', x='longitude', y='latitude')  
plt.show()

③ housing[['population', 'median\_house\_value']].corr()

	population	median_house_value
population	1.000000	-0.026882
median_house_value	-0.026882	1.000000

④ housing.describe()

⑤ housing.hist(bins=50, figsize=(20,15))\nplt.show()

The descriptive statistics of dataset are presented including mean, standard deviation and percentiles of numerical attributes

\* Select & Train Model

① from sklearn.linearmodel import LinearRegression

② lin-reg = LinearRegression()

③ lin-reg.fit(x=housing['prepared'], y=housing['label'])

④ from sklearn.metrics import mean\_squared\_error

⑤ housing\_predictions = lin-reg.predict(housing['prepared'])

\* Fine tune the model

- ① from sklearn.model\_selection import GridSearchCV
- ② forest\_reg = RandomForestRegressor()
- ③ grid\_search.fit(X=housing['rooms-prepared'], y=housing['labels'])

GridSearchCV

estimator: RandomForestRegressor

\* Launch, Monitor and maintain your system

The final model is prepared, monitoring code is written to check the system's performance. The model is retrained on a regular basis with fresh data to avoid slow degradation

S  
H/4/2024

## # Linear Regression

- ① from sklearn.model\_selection import train\_test\_split  
from sklearn.linear\_model import LinearRegression  
from sklearn.metrics import mean\_squared\_error
- ② # Load data  
data = pd.read\_csv("content/50\_startups.csv")  
data = pd.get\_dummies(data, columns=['State'], drop\_first=True)
- ③ # Split data into features & target variable  
 $X = \text{data}$   
# Split data into training & testing sets (80% train, 20% test)  
 $X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train\_test\_split}(X, y,$   
 $\text{test\_size} = 0.2, \text{random\_state} = 42)$
- ④ linear\_reg = LinearRegression()  
linear\_reg.fit(X\_train, y\_train)  
 $y_{\text{pred\_linear}} = \text{linear\_reg.predict}(X_{\text{test}})$   
 $\text{mse\_linear} = \text{mean_squared_error}(y_{\text{test}}, y_{\text{pred\_linear}})$
- ⑤ print("Mean Squared Error: ", mse\_linear)

O/P:

Mean Squared Error: 82010363.04436099

## \* Multiple Regression

- ① data-url = "http://lib.stat.cmu.edu/datasets/Boston"  
raw-df = pd.read\_csv(data-url, sep = '|st', skiprows = 22, header = None)
- ② data = pd.DataFrame(data)
- ③ X\_train, X-test, y-train, y-test = train-test-split(X, y, test-size = 0.2, random-state = 42)
- ④ # Creating & Training Regression Model  
model = MultipleRegression()  
model.fit(X-train, y-train)
- ⑤ plt.scatter(y-test, y-pred)  
plt.xlabel("True Values")  
plt.ylabel("Prediction")  
plt.title("True Values vs Prediction")  
plt.show()

## ID3 Decision Tree

- ① import pandas as pd  
from sklearn.datasets load\_iris  
from sklearn.model\_selection import train\_test\_split
- ② # Load the dataset  
iris = ~~load~~ load\_iris()  
 $X = \text{iris.data}$   
 $y = \text{iris.target}$
- ③  $X_{\text{train}}, X_{\text{test}}, y_{\text{train}}, y_{\text{test}} = \text{train\_test\_split}(X, y,$   
 $\text{test\_size} = 0.2, \text{random\_state} = 42)$   
model = DecisionTreeClassifier(criterion = 'entropy', random\_state = 42)
- ④ model.fit( $X_{\text{train}}, y_{\text{train}}$ )
- ⑤  $y_{\text{pred}} = \text{model.predict}(X_{\text{test}})$   
accuracy = accuracy\_score( $y_{\text{test}}, y_{\text{pred}}$ )  
print("Accuracy: ", accuracy)  
O/P: Accuracy: 1.0
- ⑥ new\_sample = [5.1, 3.5, 1.4, 0.2]  
predicted\_class = model.predict(new\_sample)  
predicted\_species = iris.target\_names[predicted\_class][0]  
print("Predicted species: " predicted\_species)  
O/P: Predicted species: setosa

## play-tennis.csv dataset

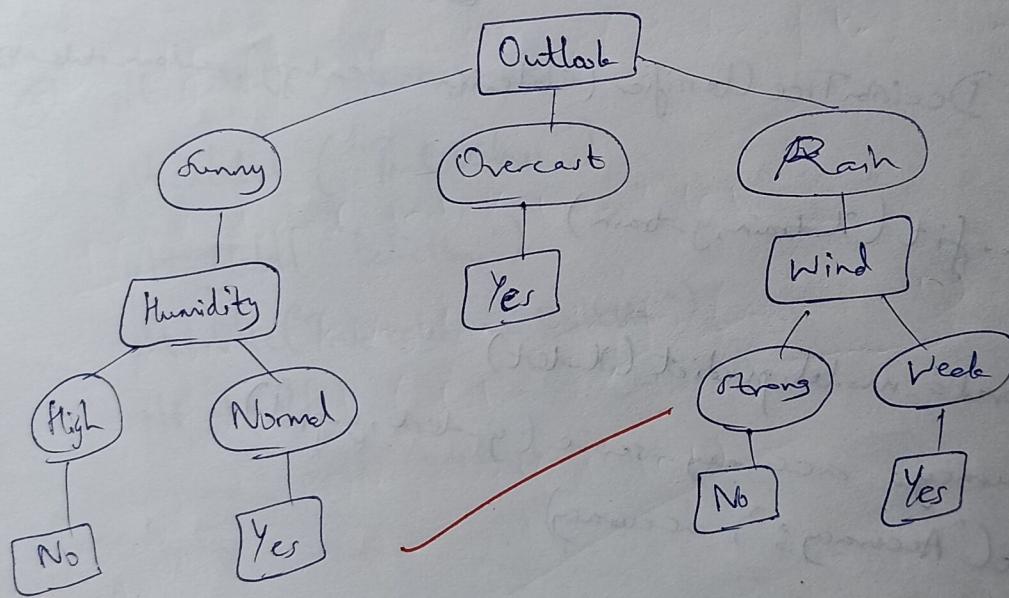
⑦ tree = buildTree(df)

⑧ import pprint

pprint.pprint(tree)

O/P:

```
{'outlook': { 'Overcast': 'Yes',
    'Rain': { 'wind': { 'Strong': 'No', 'Weak': 'Yes' } },
    'Sunny': { 'humidity': { 'High': 'No', 'Normal': 'Yes' } } }}
```



25/4/2024

Logistic Regression

- ① import numpy as np  
import pandas as pd  
from sklearn.preprocessing import StandardScaler  
from sklearn.linear\_model import LogisticRegression
- ② mushroom\_df = pd.read\_csv("Datasets/mushroom.csv")  
mushroom\_df.head()

Index	Donuts	Shop	Gold element	5'th colour
0	0	2	2	10
1	1	2	2	10
2	0	2	2	10
3	1	6	2	10

- ③  $X = \text{mushroom\_df}([["\text{Donuts}"], "\text{shop}"]])$
- ④  $Y = \text{mushroom\_df}["\text{donuts}"]$

- ⑤  $X_{\text{train}} = \text{scaler}.fit\text{-transform}(X_{\text{train}})$
- ⑥  $Y_{\text{train}} = \text{scaler}.fit\text{-transform}(Y_{\text{test}})$

~~accuracy > 0.63~~

$$([x_1, x_2], [y_1, y_2])$$

# KNN classification

Data:

- ① import pandas as pd
- import matplotlib.pyplot as plt
- import numpy as np
- ② df = pd.read\_csv("probl.csv")  
df.head()
- ③ scalar = StandardScaler()  
scalar.fit(df.drop("Target"))  
scalar.transform(df.drop("Target"))
- ④ df\_fit = pd.DataFrame(scalar.transform(df), columns=df.columns)
- ⑤ X-train, Y-train, X-test
- ⑥ KNN = KNeighborsClassifier(n\_neighbors=1)  
KNN.fit(X-train, Y-train)  
pred = KNN.predict(X-test)

Output

$$[(10, 3) [2, 6]]$$

Support vector machine

- ① import pandas as pd  
from sklearn import datasets  
from sklearn.model\_selection import train\_test\_split  
from sklearn.svm import SVC  
from sklearn.metrics import accuracy\_score
- ② iris = datasets.load\_iris()  
X = iris.data  
y = iris.target
- ③ X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y, test\_size=0.2, random\_state=42)
- ④ svm\_model = SVC(kernel='linear')  
svm\_model.fit(X\_train, y\_train)  
y\_pred = svm\_model.predict(X\_test)
- ⑤ accuracy = accuracy\_score(y\_test, y\_pred)  
print("Accuracy of SVM model: ", accuracy)

## K-Means Algorithm

- ① import matplotlib.pyplot as plt  
from sklearn.cluster import KMeans
- ② data = pd.read\_csv("data.csv")  
 $x = \text{data}.\text{values}$
- ③ kmeans = KMeans(n\_clusters=3, random\_state=42)  
kmeans.fit(x)
- ④ data['cluster'] = labels  
data.to\_csv('clustered-data.csv', index=False)
- ⑤ plt.scatter(x[:, 0], x[:, 1], c=labels, cmap='viridis')  
plt.show()

## PCA

- ① import pandas as pd  
from sklearn import datasets  
from sklearn.decomposition import PCA  
import matplotlib.pyplot as plt
- ② iris = datasets.load\_iris()  
 $x = \text{iris}.\text{data}$   
 ~~$y = \text{iris}.\text{target}$~~

③  $pca = PCA(n\_components=2)$

$X_pca = pca.fit\_transform(X)$

④  $pca\_df = pd.DataFrame([date=X_pca, columns=['Principal Component 1', 'Principal Component 2']])$

$pca\_df['target'] = y$

⑤ plt.figure(figsize=(8,6))

plt.scatter(pca\_df['pc1'], pca\_df['pc2'], c=pca\_df['Target'],  
cmap='viridis')

plt.colorbar()

plt.show()

## Random Forest Ensemble Method

- ① import pandas as pd  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.metrics import accuracy\_score
  - ② iris = load\_iris()  
 $X = \text{iris}.\text{data}$   
 $y = \text{iris}.\text{target}$
  - ③  $X\text{-train}, X\text{-test}, y\text{-train}, y\text{-test} = \text{train-test-split}(X, y,$   
 $\text{test-size} = 0.2, \text{random-state} = 42)$   
# Build Random Forest Model  
rf\_model = RandomForestClassifier(n\_estimators = 100, random\_state = 42)
  - ④ # Train the model  
rf\_model.fit(X\_train, y\_train)
  - ⑤ # Predict using trained model  
accuracy = accuracy\_score(y\_test, y\_pred)  
print("Accuracy of Random Forest Model:", accuracy)

01P:

~~Accuracy of Random Forest model: 1.0~~

## Boosting Ensemble Algorithm

① import pandas as pd

```
from sklearn.ensemble import AdaBoostClassifier
```

```
from sklearn.tree import DecisionTreeClassifier
```

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

② iris = load\_iris()

X = iris.data

y = iris.target

③ X\_train, X\_test, y\_train, y\_test = train\_test\_split(X, y,  
test\_size=0.2, random\_state=42)

# Build AdaBoost Model using DecisionTreeClassifier as base

```
base_estimator = DecisionTreeClassifier(max_depth=1, random_state=42)
```

```
ada_model = AdaBoostClassifier(base_estimator=base_estimator,  
n_estimators=50, random_state=42)
```

④ ada\_model.fit(X\_train, y\_train)

⑤ y\_pred = ada\_model.predict(X\_test)

⑥ accuracy = accuracy\_score(y\_test, y\_pred)

print("Accuracy of AdaBoost model:", accuracy)

O/P:-

Accuracy of AdaBoost model : 1.0

ES  
30/5/2024