

I N D E X

Name Saikhumar Polle Patil Std 12 Sec 6C

Roll No. 182 Subject ML - Lab School/College

School/College Tel. No. Parents Tel. No.

Sl. No.	Date	Title	Page No.	Teacher Sign. Remarks
1	21-3-24	import & export using pandas	1-2	S.T.
2	4-4-24	End-End ML Project	3-5	
3	11-4-24	visualizing graphical data	6-7	
4	18-4-24	Simple linear regression	8-9	
5	25-4-24	Decision Tree algorithm	10-11	
6	2-5-24	KNN classifier	12-13	
7	9-5-24	logistic regression Model	14-15	
8	23-5-24	(A) SVM	15-16	
	23-5-24	(B) Dimensionality reduction	17	
	23-5-24	(C) K-means clustering	18-19	
a.	30-5-24	(A) Random forest ensemble method	20-21	
	30-5-24	(B) Boosting algorithm.	22-23	

Lab -1

⇒ Write a Python Program to Import and export data using Pandas, IPbrrary Functions.

```
Import Pandas as pd
```

```
airbnb_data = pd.read_csv("data/listings-austin.csv")  
airbnb_data.head()
```

Output :

id	name	host-id	host-name
2265	ken-retreat	2466	Paddy
5245	Ecological	2466	Paddy
3265	neighborhood george	7892	Private room
1234	latitude	3234	Executive apt
5482	Longitude	7892	Private room

Reading data From URL

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

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

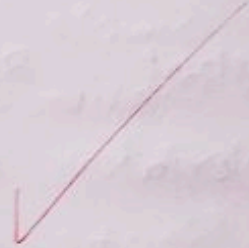
```
iris_data = url.read_csv(url, names=url-names)
```

```
iris_data.read()
```


Output :

	Sepal-length-in-cm	Sepal-width-in-cm
0	5.1	3.5
1	4.9	3

Sepal-length-in-cm	Sepal-width-in-cm	class
4.4	0.2	iris-setosa
1.4	0.2	iris-setosa



Lab-2

End-End - ML - Project

1. Look at the big picture.

Our first task is to use the California census data to build a model of the housing prices in the state. The data includes features such as.

- Population
- Median Income
- Median housing price for each ~~housing~~ block group in California.

A block is the smallest geographical unit for which census data is published.

Frame the problem

The problem is a supervised learning problem, and because we're predicting the median housing price for a district, then this is a Regression task.

Performance measure: A typical measure for regression problem is Root Mean Squared Error.

$$RMSE(X, h) = \sqrt{\frac{1}{n} \sum_{i=1}^n (h(x^{(i)}) - y^{(i)})^2}$$

Assumption Checking

Lastly, it is good practice to list & verify the list of assumptions made by us or others, this can help us catch serious mistakes early on.

①. Get the data.

→ Download the data: It is preferable to create some util functions to automate the process of downloading / extracting web-based data sets. So, we need to import some of the python libraries such as os, tarfile, urllib.

Then we download the data using Fetch-housing-data. Pandas library is used to import the data from url's abstracted.

→ Take a quick look at the data structure.

command: housing.head()

Each row represents one district & has the 10 attributes.

→ The `info()` method is used to take a quick look at the data. It gives the info such as:

The number of rows, Name of column, data type of each column in the dataset.

The matplotlib.pyplot and seaborn libraries are imported in order to plot the graphs to gain the most insightful view into the data.

→ The Histograms shows the distribution of a particular column.

⇒ Create a Test set: creating a test set is easy. randomly select a subset of the data, typically 20% (or less if the dataset is large).

- housing['id'].value_counts()

- housing.groupby('by = ['longitude', 'latitude']).
count()

from sklearn.model_selection import train_test_split
~~train-set, test-set = train_test_split(housing,~~
~~test_size=0.2, random_state=42).~~

Sp.9
4/4/21

③ Discover & visualizing the data to gain insight
visualizing the graphical Data.

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

```
housing.plot(kind = 'scatter', x = 'longitude', y = 'latitude',
alpha = 0.4, s = housing['population']/100, label =
population, figsize = (10, 4), c = 'median-house-value',
 cmap = plt.get_cmap(name = 'set'), colorbar = True)
plt.legend()
```

Looking For Correlations.

```
corr-matrix = housing.corr()
```

Prepare the data For Machine learning algorithm

```
housing = start-train-set.drop("median-house-value")
housing-labels = start-train-set("median-house-value").
copy()
```

* Data cleaning

```
housing.dropna(subset = 'total-bedrooms')
```

```
computer = simpleInputs(strategy = 'median')
```

```
housing-num = housing.drop("clean-proximity")
inputs.fit(housing-num)
```

• Handling Test and categorical Attributes

```
housing - cat = housing[['clean - proximity']]  
housing - cat.head(10)
```

* select & train a model

```
from sklearn.linear_model import LinearRegression
```

```
lin-reg = LinearRegression()
```

```
tree-reg = DecisionTreeRegressor
```

```
Forest-reg = RandomForestRegressor
```

* Fine-tune Your Model

```
Forest-reg = RandomForestRegressor()
```

```
grid-search = GridSearchCV(estimator=Forest-reg)
```

```
Param-grid = Param-grid, scoring = 'neg-mean-square-error',  
cv=5, return-train-score=True, n-jobs=1)
```

* Launch Monitor & Maintain your System

We can hot-load the model with in web application ~~or~~ alternatively, wrap the model around ~~PH~~ own API end-point & design web component separately.

Simple Linear Regression

```
import pandas as pd
```

```
import numpy as np
```

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

```
import seaborn as sns
```

Import data

```
df_sal = pd.read_csv('content/salary-data.csv')
```

```
df.head()
```

Analyze data

```
df_sal.describe()
```

 "mean, count, std, min, maxDistribution :

```
plt.title("Salary Distribution plot")
```

```
sns.distplot(df_sal['salary'])
```

```
plt.show()
```

Scatter Plot [Salary vs Experience]

```
plt.scatter(df_sal['Years Experience'], df_sal['salary'],  
            color = 'lightcoral')
```

```
plt.title('Salary vs experience')
```

```
plt.xlabel('Years of experience')
```

```
plt.ylabel('Salary')
```

```
plt.box(False)
```

split Data

x = df.sal.iloc[:, :1] // independent
y = df.sal.iloc[:, 1:] // dependent

x_train, x_test, y_train, y_test = train_test_split(
x, y, test_size=0.2, random_state=0)

regressor = LinearRegression()

regressor.fit(x_train, y_train)

y_pred_test = regressor.predict(x_test)

y_pred_train = regressor.predict(x_train)

Prediction on training set

plt.scatter(x_train, y_train)

plt.plot(x_train, y_pred_train)

plt.title('Salary vs experience')

plt.xlabel('Years of experience')

plt.ylabel('Salary')

plt.show()

Prediction on test set

plt.scatter(x_test, y_test)

~~plt.plot(x_train, y_pred_train)~~

~~plt.title('Salary vs experience')~~

~~plt.xlabel('Years of experience')~~

~~plt.ylabel('Salary')~~

Lab-5

Decision TD1	Algorithm	Implementation.
--------------	-----------	-----------------

⇒ Reading the dataset

import numpy as np

import pandas as pd

from numpy import log2 as log

df = pd.read_csv("dataset.csv")

calculating entropy

$$\text{Entropy} = -\frac{P}{P+n} \log_2\left(\frac{P}{P+n}\right) - \frac{n}{(P+n)} \log_2\left(\frac{n}{P+n}\right)$$

def find_entropy(df):

target = df.keys()[0]

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

Average Information gain

def average_information(df, attributes):

target = df.keys()[0]

target_variables = df[target].unique()

variables = df[attributes].unique()

entropy2 = 0

for variable in variables

entropy2 = 0

for target_variable in target_variables:

```

num = len(df[attribute][df[attribute] == variable])
df[target] = target_variable
den = len(df[attribute][df[attribute] == variable])
fraction = num / den
entropy += fraction * log(fraction)
fraction2 = den / len(df)
entropy2 += -fraction2 * entropy
return abs(entropy2)

```

Information gain

```
def find_min(df):
```

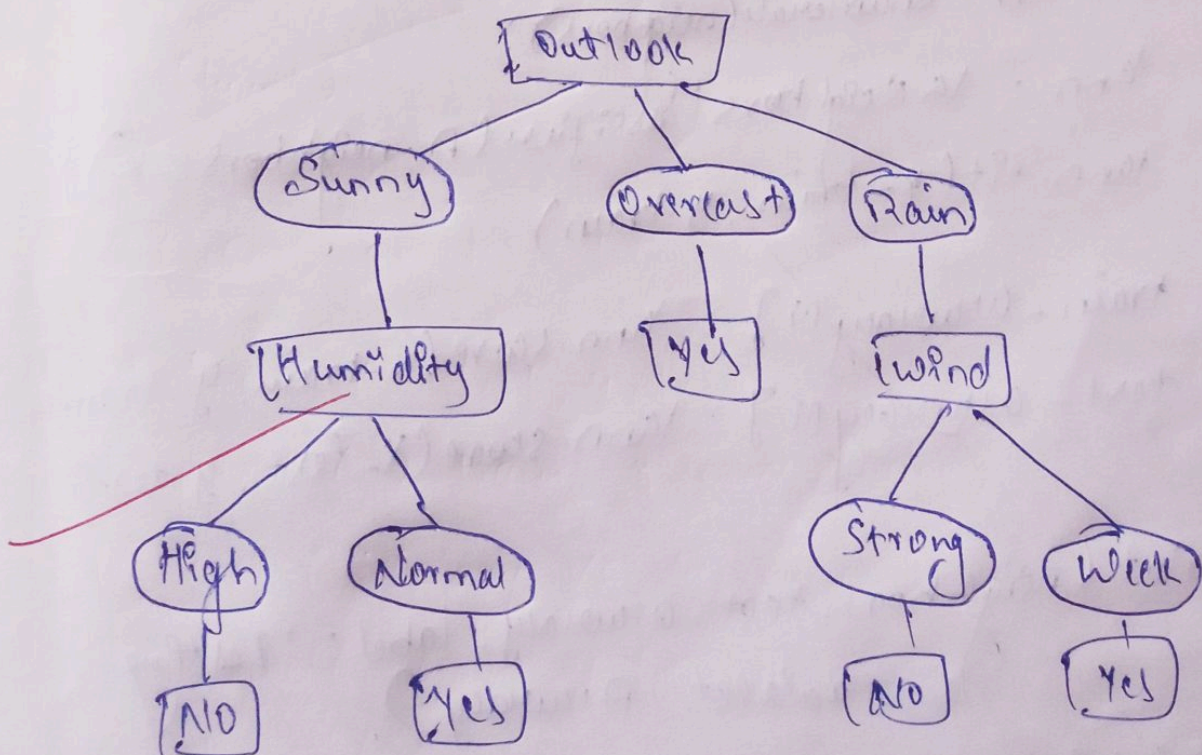
```
    Rg = []
```

```
    for key in df.keys()[1:-1]:
```

```
        Rg.append(find_entropy(df) - average_information(df, key))
```

```
    return df.keys()[1:-1][np.argmax(Rg)]
```

Decision Tree



⇒ KNN classifier

```

from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import train_test_split
from sklearn.datasets import load_iris
import numpy as np
import matplotlib.pyplot as plt

```

```
irisData = load_iris()
```

```
x = irisData.data
```

```
y = irisData.target
```

```

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

```

```
neighbors = np.arange(1, 9)
```

```
train_accuracy = np.empty(len(neighbors))
```

```
test_accuracy = np.empty(len(neighbors))
```

```
for i, k in enumerate(neighbors):
```

```
    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)
```

```
plt.plot(neighbors, test_accuracy, label='Testing  
dataset Accuracy')
```

```
plt.plot(neighbors, train_accuracy, label='Training
```

plt.legend()

plt.xlabel('n-neighbors')

plt.ylabel('Accuracy')

plt.show()

Output

	sepal-length	sepal-width	petal-len	petal-wid	Species
0	5.1	3.5	1.4	0.2	Irps-setosa
1	4.9	3.0	1.4	0.2	Irps-setosa
2	4.7	3.2	1.3	0.2	Irps-setosa
3	4.6	3.1	1.5	0.2	Irps-setosa
4	4.6	3.6	1.4	0.2	Irps-setosa

②. Build Logistic Regression Model

```

from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
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=0.20, 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)
  
```

Output:

Logistic Regression model accuracy : 95.61

*Sel
2/6/24*

Q. Build Support vector machine model for a given dataset

=> From sklearn.datasets import load_iris

From sklearn.model-selection import train-test-split

From sklearn.preprocessing import StandardScaler

From sklearn.~~preprocess~~svm import SVC

From sklearn.metrics import accuracy_score,
classification_report

iris = load_iris()

x, y = iris.data, iris.target

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

scaler = StandardScaler()

x_train_scaled = scaler.fit_transform(x_train)

x_test_scaled = scaler.transform(x_test)

Svm-model = SVC(kernel = 'rbf', C = 1.0, gamma = 'scale',
random_state = 42)

Svm-model.fit(x_train_scaled, y_train)

=> ~~$\frac{SVC}{SVC(random_state = 42)}$~~

~~y_pred = Svm-model.predict(x_test_scaled)~~

~~accuracy = accuracy_score(y_test, y_pred)~~

~~print("Accuracy:", accuracy)~~

print("\n Classification Report:")

print(classification_report(y_test, y_pred, target_names = iris.target_names))

Classification Report:

	Precision	recall	F1-score	support
Setosa	1.00	1.00	1.00	10
Versicolor	1.00	1.00	1.00	9
Virginica	1.00	1.00	1.00	11

A Accuracy

macro avg	1.00	1.00	1.00	30
weighted avg	1.00	1.00	1.00	30
				30

3. Implement dimensionality reduction using principal component Analysis (PCA) method

- ⇒ Import the libraries.
- ⇒ Load the iris dataset
- ⇒ split the dataset into training & testing sets.
- ⇒ Preprocess the data.

⇒ Apply PCA For dimensionality reduction.

PCA = PCA(n-components) # Reduce to 2 principal components

X_train_pca = PCA.fit_transform(X_train_scaled)

X_test_pca = PCA.transform(X_test_scaled)

Svm-model = SVC(kernel = 'rbf', C=1.0, gamma = 'scale', random_state = 42)

y_pred = Svm-model.predict(X_test_pca)

accuracy = accuracy_score(y_test, y_pred)

print("Accuracy with PCA:", accuracy)

print("Classification Report with PCA:")

print(classification_report(y_test, y_pred, target_names = iris.target_names))

⇒ Accuracy with PCA = 0.9

Classification Report with PCA:

	Precision	Recall	F1-score	Support
Accuracy			0.90	30
macro avg	0.90	0.90	0.90	30
weighted avg	0.90	0.90	0.90	30

③. Build k-means algorithm to cluster a set of data stored in a .csv file.

⇒ import numpy as np

class KMeans:

def __init__(self, n_clusters, max_iter=300):

self.n_clusters = n_clusters

self.max_iter = max_iter

def fit(self, x):

centroids_indices = np.random.choice(len(x),

size=self.n_clusters, replace=False)

self.centroids = x[centroids_indices]

for _ in range(self.max_iter):

labels = self.assign_clusters(x)

new_centroids = np.array([x[labels == k]

mean(axis=0) for k in range(self.n_clusters)])

if np.all(self.centroids == new_centroids):

break

self.centroids = new_centroids

np.random.seed(42)

x = np.random.rand(100, 2)

kmeans = KMeans(n_clusters=3)

kmeans.fit(x)

labels = kmeans. - assign. clusters (x).

Print ("Cluster labels:", labels)

Output:

Cluster labels: [1 0 2 1 0 1 0 1 0 2 1 2 2 1 1 1 2
0 2 2 0 2 0 0 1 0 0 2 2 0 0 1
2 0 1 0 0 1 1 0 1 2 1 2 2 2]

S.P.
23/5/24

Lab-8
Q.6. Implement Random Forest ensemble method.

→ import numpy as np

from sklearn.datasets import load-iris

from sklearn.model-selection import train-test-split

from sklearn.ensemble import RandomForestClassifier

from sklearn.metrics import accuracy_score, classification_report

import

Load the Iris dataset

iris = load-iris()

x = iris.data

y = iris.target

split the dataset into training & testing sets.
x_train, x_test, y_train, y_test = train-test-split
(x, y, test-size = 0.3, random-state = 42)

Initialize Random Forest Classifier

rf_classifier = RandomForestClassifier(n_estimators = 100,
random_state = 42).

Train the Classifier.

rf_classifier.fit(x_train, y_train).

y_pred = rf_classifier.predict(x_test).

accuracy = accuracy_score(y_test, y_pred)

report = classification_report(y_test, y_pred)

print('Accuracy: {accuracy}').

print('Classification Report:')

print(report).

Output:

Accuracy : 1.0

Classification Report:

	Precision	recall	F1-score	Support
0	1.00	1.00	1.00	14
1	1.00	1.00	1.00	13
2	1.00	1.00	1.00	13

accuracy			1.00	45
macro avg	1.00	1.00	1.00	45
weighted avg.	1.00	1.00	1.00	45.

10. Implement Boosting ensemble method.

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

```
x, y = load_breast_cancer(return_x_y=True).
```

```
from sklearn.model_selection import train_test_split  
x_train, x_test, y_train, y_test = train_test_split(x, y,  
                                                    test_size = 0.3, random_state = 23)
```

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

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

```
dtree = DecisionTreeClassifier(max_depth=3, random_state=23)
```

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

```
dt_pred = dtree.predict(x_test)
```

```
dt_all = round(accuracy_score(y_test, dt_pred), 3)
```

```
print(f"Decision Tree classifier Accuracy score:", dt_all).
```

⇒ Decision Tree classifier Accuracy = 0.953

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

```
ada = AdaBoostClassifier(n_estimators=50, learning_rate=0.6)
```

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

```
ada_pred = ada.predict(x_test).
```

```
ada_all = round(accuracy_score(y_test, ada_pred), 3)
```

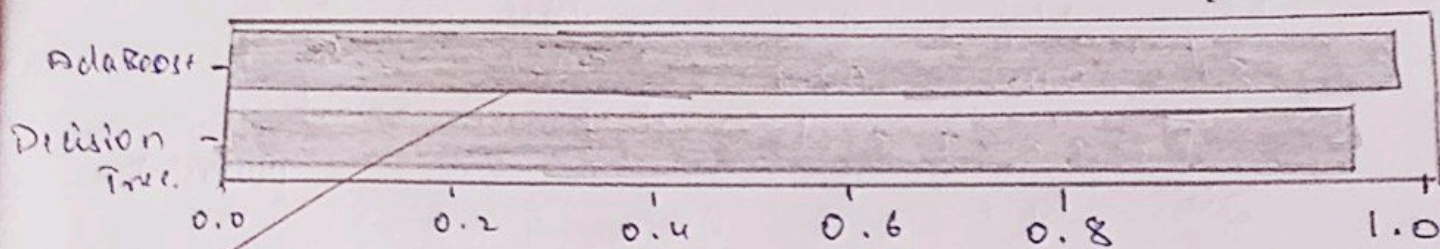
```
print(f"Decision tree AdaBoost Model Accuracy  
score:", ada_all)
```

⇒ Decision Tree AdaBoost Model Accuracy score:
0.982

```

import numpy as np
import matplotlib.pyplot as plt
plt.figure(figsize=(10,2))
plt.barh(np.arange(2), [dt-acc, ada-acc],
         tick_label = ['Decision Tree', 'AdaBoost'])

```



Sp. 1
20/5/24