

I N D E X

NAME: Shravani STD.: _____ SEC.: _____ ROLL NO.: _____ SUB.: ML LAB

S. No.	Date	Title	Page No. Marks	Teacher's Sign / Remarks
1	21/03/2024	LAB 1	10	Avva 21/3/24
2	28/03/2024	LAB 2	10	Avva 28/3/24
3	04/04/2024	LAB 3	10	Avva 4/4/24
4	18/04/2024	LAB 4	10	Avva 18/4/24
5	25/04/2024	LABS	10	Avva 25/4/24
6	09/05/2024	LAB 6	10	Avva 9/5/24
7	09/05/2024	LAB 6	10	Avva 9/5/24
8	23/05/2024	LAB 7 - 8	10	Avva 23/5/24
9	23/05/2024	LAB 7 - 9a 9b	10	Avva 23/5/24
10	30/05/2024	LAB 8 - 10	10	Avva 30/5/24
11	30/05/2024	LAB 8 - 11	10	Avva 30/5/24

Importing and Exporting a dataset

```
import pandas as pd
```

```
url = "http://..."
```

```
df = pd.read_csv(url, columns = ['sepal_length',  
'sepal_width', 'petal_length',  
'petal_width', 'class'])
```

```
df.head()
```

	sepal_length	sepal_width	petal_length	petal_width	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
:					

```
iris_data.to_csv('cleaned_iris_data.csv')
```

EWI
20/03/24

LAB - 2
End To end Project

- 1 Look at The big picture : Performance
- 2 Get The Data Measure

fetch_housing_data()

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)

housing = load_housing_data()

housing.head()

housing.info()

housing.describe()

import matplotlib.pyplot as plt

import seaborn as sns

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

plt.show()

pd.cut

housing['income_cat'] = pd.cut(x =

housing['median_income'], bins =

[0, 1.5, 3, 4.5, 6, np.inf],
labels = [1, 2, 3, 4, 5])

housing['income_cat'].hist()

LAB-8

09/05/2024

Program 10:

K-Means Algorithm:

```

→ ins = pd.read_csv("iris.csv")
x = ins.iloc[:, [0, 1, 2, 3]].values
ins_outcome = pd.crosstab(index=ins["species"],
                           columns="count")

```

ins_outcome

```

sns.FacetGrid(ins, hue="species").map(sns.
    distplot, "petal_length").add_legend()
sns.FacetGrid(ins, hue="species").map(sns.
    distplot, "petal_width").add_legend()
sns.FacetGrid(ins, hue="species").map(
    sns.distplot, "sepal_length").add_legend()
plt.show()

sns.boxplot(x="species", y="petal_length",
            data=ins)
plt.show()

sns.set_style("whitegrid")
sns.pairplot(ins, hue="species", size=3)
plt.show()

from sklearn import KMeans
kmeans = KMeans(n_clusters=3,
                 init="k-means++", max_iter=300, n_init=10,
                 random_state=0)
y_kmeans = kmeans.fit_predict(x)

```

Program 11:

PCA

```

df = pd.read_csv("iris.csv")
df_features = df.drop([
    from sklearn.preprocessing import StandardScaler
standardized = StandardScaler().fit(df_features)
scaled_data = standardized.transform(df_features)

from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(scaled_data)
x_pca = pca.transform(scaled_data)
scaled_data.shape
x_pca.shape

def diag(x):
    if x == "M":
        return
    else:
        return df_diag = df.T["diagonal"]
        x_pca[1]
        df_pc = pd.DataFrame(
            columns=["PC1", "PC2"])
        plt.figure(figsize=(10, 7))
        sns.heatmap(df_pc)
        plt.title("Principal Component Analysis")

```

import statistics

housing = num.median().values

- from sklearn.preprocessing import OrdinalEncoder
ordinal_encoder = OrdinalEncoder()
ordinal_encoder.categories_
- attr_address = combined_attributes_Address[add_bedrooms_per_room = False]
housing_extra_attributes = attr_address.transform(housing.values)
- from sklearn.compose import ColumnTransformer
num_attributes = housing_num.columns.tolist()

5 select and Train a Model:

- from sklearn.linear_model import LinearRegression
lin_reg = LinearRegression()
lin_reg.fit(x = housing_prepared, y = housing_labels)
- housing_predictions = lin_reg.predict(housing_prepared)
lin_mse = mean_squared_error(housing_label, housing_predictions)
lin_rmse = np.sqrt(lin_mse)
- tree_reg = DecisionTreeRegressor()
tree_reg.fit(x = housing_prepared, y = housing_labels)
- forest_reg = RandomForestRegressor()
forest_reg.fit(x = housing_prepared, y = housing_labels)
forest_reg.fi

6 Fine - Tune Your Model:

- `grid_search.best_params_`
- `grid_search.best_estimator_`
- `cat_encoder = full_pipeline.named_transformers_['cat']`
- `final_model = grid_search.best_estimator_`
- `X_test_prepared = full_pipeline.transform(X=X)`
$$\text{final_rmse} = \text{np.sqrt}(\text{final_mse})$$

$$\text{final_rmse}$$
- $\text{squared_errors} = (\text{y_test} - \text{final_prediction})^2$

7 Launch, Monitor and Maintain your system

6
4/4/24

Simple Linear Regression and Multiple Linear Regression

Simple Linear Regression

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

```
df_sal = pd.read_csv(r'C:\salary\Salary.csv')
```

```
df_sal.head()
```

```
df_sal.describe()
```

Salary Distribution Plot

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

```
plt.show()
```

```
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=42)~~

```
regressor = LinearRegression()
```

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

```
y_pred_test = regressor.predict(x_test)
```

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

print(f'coefficient : {regressor.coef_}')
print(f'intercept : {regressor.intercept_}')

```
df_start = pd.read_csv(r'C:\SO_Startups.csv')
df_start.head()
```

```
df_start.describe()
sns.distplot(df_start['Profit'])
plt.show()
```

```
regressor = LinearRegression()
regressor.fit(x_train, y_train)
y_pred = regressor.predict(x_test)
np.set_printoptions(precision=2)
result = np.concatenate((y_pred, reshape(
    len(y_pred), 1), y_test, reshape(
    len(y_test), 1)), axis=1)
```

result

Use an appropriate data set for building the decision tree (ID3) and apply this knowledge to classify a new sample.

```

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.metrics import classification_report

```

url = "....."

df = pd.read_csv(url, names=['sepal_length(cm)',
.....])

df.head()

~~df = 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 IPython.display import Image  
from sklearn.tree import export_graphviz  
!pip install pydotplus  
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())
```

```
dt = DecisionTreeClassifier(random_state=1)
```

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

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

$\text{dt2} = \text{DecisionTreeClassifier}(\text{random_state}=1)$
params = {
 'max_depth': [2, 3, 4, 5],
 'min_samples_split': [2, 3, 4, 5],
 'min_samples_leaf': [1, 2, 3, 4, 5]}
}

$\text{gsearch} = \text{GridSearchCV}(\text{dt2}, \text{param_grid}=\text{params}, \text{cv}=3)$

$\text{gsearch.fit}(X, y)$

$\text{gsearch.best_params_}$

$\text{dt2} = \text{DecisionTreeClassifier}(**\text{gsearch.best_params_}, \text{random_state}=1)$

$\text{dt2.fit}(X_{\text{Train}}, y_{\text{Train}})$

$y_{\text{pred_Train}} = \text{dt2.predict}(X_{\text{Train}})$

$y_{\text{prob_Train}} = \text{dt2.predict_proba}(X_{\text{Train}})$

$y_{\text{pred}} = \text{dt2.predict}(X_{\text{Test}})$

$y_{\text{prob}} = \text{dt2.predict_proba}(X_{\text{Test}})$

Output:

petal width(cm) ≤ 0.8
gini = 0.667
samples = 150
value = [50, 50, 50]

True

gini = 0.0
samples = 50
value = [50, 0, 0]

False

petal width(cm) ≤ 1.7
gini = 0.5
samples = 100
value = [0, 50, 50]

(random_state=1)

[2, 3, 4, 5]

[2, 3, 4, 5]

[1, 2, 3, 4, 5]

, param_grid=
cv = 3

* gsearch, best_

= 1
best (xTrain)
best_prob(xTrain)

Test
best (xTest)[:, 1]

<= 0.8
.667
150
0, 50, 50

False
petal width(cm) <= 1.75
gini = 0.5
samples = 100
vowc = [0, 50, 50]

Accuracy of Decision Tree-Train : 1.0
Accuracy of Decision Tree - Test : 0.955556

18-4-2024



Scanned with OKEN Scanner

```
import pandas as pd  
from matplotlib import pyplot as plt  
from sklearn.model_selection import train_test_split  
from sklearn.linear_model import LogisticRegression
```

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

```
X_Train, X_Test, y_Train, y_Test = train_test_split(df[['age']], df['bought_insurance'], test_size=0.2)  
print(X_Test[:])  
print(X_Test)
```

```
model = LogisticRegression()  
model.fit(X_Train, y_Train)  
y_pred = model.predict(X_Test)  
print(y_pred)
```

```
print(model.predict_proba(X_Test))  
print(model.score(X_Test, y_Test))
```

Output:

Accuracy = 0.5

```
import Math  
def sigmoid(z):  
    return 1 / (1 + Math.exp(-z))  
  
def predict(age):  
    z = 0.042 * age - 1.53  
    y = sigmoid(z)  
    return y  
  
print(predict(35))  
print(predict(43))
```

Output

Prediction: array([1, 0, 1, 0, 0, 0, 0, 10])

Score: 0.8888

Linear Reg score: 0.584321

Predictions: 0.485

0.5685



KNN Classifier Model

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

```
df = pd.read_csv('Iont')
df.head()
```

```
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)
```

```
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)
```

```
plt.scatter(neighbours, train_accuracy, color='blue')
plt.show()
```

KNN = K Neighbours Classifies (n_neighbors = 7)

knn.fit(x_train, y_train) report MA
knn.score(x_test, y_test) TTestReport.gn = db

Output

0.730519

SVM Model

from sklearn import datasets
cancer = datasets.load_breast_cancer()

print(cancer.target) 20895.gn = 0 = normal
1 = malignant

from sklearn.model_selection import train_test_

x_train, x_test, y_train, y_test = train_test_

split(cancer.data, cancer.target, test_size=

0.3, random_state = 42)

from sklearn import svm

clf = svm.SVC(kernel='linear')

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

ANN Program

```
db = np.load("input/duke.npy")
```

```
np.random.shuffle(db)
```

```
y = db[:, 0]
```

```
x = np.delete(db, [0], axis=1)
```

$x_{\text{Train}}, x_{\text{Test}}, y_{\text{Train}}, y_{\text{Test}} =$

```
train_test_split(x, y, test_size=0.1)
```

```
print(np.shape(x_train), np.shape(x_test))
```

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

```
weights_h = np.random.random([len(x[0]), 72])
```

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

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

```
def sum_function(weights, index_col, x):
```

```
result = 0
```

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

```
result += x[i] * weights[i][index_col]
```

locked_col

```
return result
```

```
def activate_layer(layer, weights, x):
```

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

```
layer[i] = 1.7159 * np.tanh(2.0 *
```

```
sum_function(weights, i, x) / 3.0)
```

def backpropagation(hidden_layer, output_layer, one_hot_encoding, learning_rate, n):

$$\text{output_derivative} = \text{np.zeros}(2)$$

$$\text{output_gradient} = \text{np.zeros}(2)$$

for i in range(0, len(output_layer)):

$$\text{output_derivative}[i] = (1.0 - \text{output_layer}[i]) * \text{output_layer}[i]$$

for i in range(0, len(output_layer)):

$$\text{output_gradient}[i] = \text{output_derivative}[i] * (\text{one_hot_encoding}[i] - \text{output_layer}[i])$$

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

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

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_layer)):

$$\text{sum_} += \text{output_gradient}[j] * \text{hidden_weights}[i][j]$$

$$\text{hidden_gradient}[i] = \text{sum_} * \text{hidden_derivative}[i]$$

recalculate(learning_rate, weights, hidden_gradient)

Random Boosting and Ada Boosting

dataset = pd.read_csv('...')

corr_matrix = input_x.corr().abs()

upper = corr_matrix.where(np.triu(np.ones(corr_matrix.shape), k=1).astype(np.bool))

to_drop = [column for column in upper.columns if any(upper[column] > 0.95)]

upper[upper[column] > 0.95]

list(set(to_drop) - set(['...']))

X_train, X_test, y_train, y_test =

train_test_split(X, y, test_size=0.1)

n_classes = 9 - 1 = 8

n_estimators = 100

random_state = 13

names = ['RandomForestClassifier',

'AdaBoostClassifier'

models = ['RandomForestClassifier',

'AdaBoostClassifier']

n_estimators = n_estimators,

max_depth = max_depth,

min_samples_leaf = min_samples_leaf,

n_estimators = n_estimators]

for counter, model in enumerate(models):

model.fit(X_train, y_train)

y_pred = model.predict(X_test)



Accuracy Random Forest Classifier : 0.768439
Accuracy AdaBoost Classifier : 0.7341337

~~accuracy = 0.768439~~
~~accuracy = 0.7341337~~
~~accuracy = 0.768439~~
~~accuracy = 0.7341337~~

~~accuracy = 0.768439~~
~~accuracy = 0.7341337~~
~~accuracy = 0.768439~~
~~accuracy = 0.7341337~~

~~accuracy = 0.768439~~
~~accuracy = 0.7341337~~

~~accuracy = 0.768439~~



Random Boosting and AdaBoosting

dataset = pd.read_csv('...')

corr_matrix = input_xi.to_numpy().astype('float')

upper = corr_matrix.where(np.triu(1) == 1).copy().triu(1).upper() # upper matrix shape

triu_idx = np.triu(1).astype(np.bool_)

to_drop = [column for column in

upper.columns if any(upper[upper[column] > 0.95].notnull().any())]

print(to_drop) # drop these two

X_train, X_test, y_train, y_test =

train_test_split(X, y, test_size=0.1)

n_classes = 5 # 0 = Squirrel - rabbit

n_estimators = 100 # n is 100

random_state = 13 #

names = ['RandomForestClassifier',

'AdaBoostClassifier'

models = [RandomForestClassifier(n_estimators=n_estimators),

n_estimators = n_estimators,

AdaBoostClassifier(max_depth=

None, n_estimators=n_estimators)]

for counter, model in enumerate(models):

model.fit(X_train, y_train)

y_pred = model.predict(X_test)

Program 11:

PCA

```
df = pd.read_csv("...\\data.csv")
```

```
df_features = 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)
```

```
x_pca = pca.transform(scaled_data)
```

```
scaled_data.shape
```

```
x_pca.shape
```

```
def diag(x):
```

```
    if x == 'M':  
        return 1
```

```
else:  
    return 0
```

```
df_diag = df[['diagnosis']].apply(diag)
```

```
x_pca[:]
```

```
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')
```

