

Experiment-05.

Dimensionality reduction and Impact Analysis.

problem statement:- Evaluate the effect of dimensionality reduction on classification accuracy.

dataset: winequality (Red + white)

[import pandas as pd

import]

Loading the Dataset:-

from google.colab import files

uploaded = files.upload().

choosing file:- wine quality.csv.

Reading the file :-

import pandas as pd.

df = pd.read_csv('wine quality.csv')

print(df).

Output:-

	fixed acidity	volatile acidity	- - - - -	quality	Id
0	7.4	0.700		5	0
1	7.8	0.880		5	1
2	7.8	0.760		5	2
3	11.2	0.280	- - - - -	6	3
;	;	;		;	;
;	;	;		;	;
;	;	;		;	;
;	;	;		;	;
;	;	;		;	;
1142	5.9	0.645		5	1597

LABORATORY CONTINUOUS EVALUATION			
S.NO.	PROGRAM	Marks	Marks Scored
1.	Preparedness	5	
2.	Coding	5	
3.	Testing	5	
4.	Viva	5	
	TOTAL	20	

Faculty Signature with Date

Checking null values :-
print(df.isnull().sum())

output:-	fixed acidity	0
	volatile acidity	0
	citric acid	0
	residual sugar	0
	chlorides	0
	free sulphur dioxide	0
	total sulphur dioxide	0
	density	0
	pH	0
	sulphates	0
	alcohol	0
	quality	0
	Id.	0

Load and preprocess dataset

```
import pandas as pd.
```

```
from sklearn.model_selection
```

```
import train_test_split
```

```
from sklearn.preprocessing import
```

```
StandardScalar.
```

```
df = pd.read_csv("winequality.csv")
```

```
df = df.drop(columns=["Id"])
```

```
x = df.drop(["quality"], axis=1)
```

```
y = df["quality"]
```

```
print("Datashape:", df.shape)
```

```
print("class distribution:\n", y.value_counts())
```

```
x_train, x_test, y_train, y_test = train_test_split(
```

```
x_train, x_test, y_train, y_test, test_size=0.2, random_state=42,
```

```
stratify=y).
```

```
scalar = StandardScalar()
```

```
x_train_scaled = scalar.fit_transform(x_train)
```

```
x_test_scaled = scalar.transform(x_test).
```

Output:- Dataset shape : (143, 12)

class distribution :

quality.

5	483
6	462
7	143
4	33
8	16
3	6

Name : count

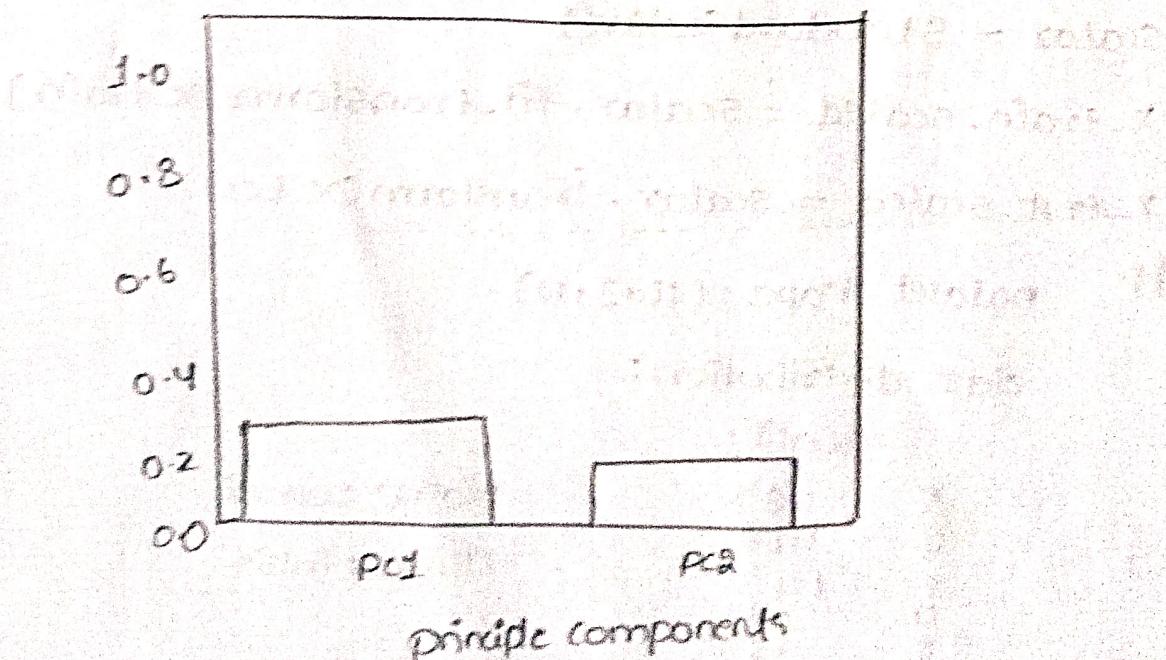
Dtype : int64

(ii) Apply PCA and LDA for reducing dimensionality.

code & (PCA)

```
import matplotlib.pyplot as plt.  
from sklearn.decomposition import PCA.  
pca = PCA(n_components=2)  
X_train_pca = pca.fit_transform(X_train_scaled)  
X_test_pca = pca.fit_transform(X_test_scaled).  
plt.figure(figsize=(6,4))  
plt.bar([1,2], pca.explained_variance_ratio_, tick_label=  
["PC1", "PC2"], color=[ "#1f77b4", "#ff7f0e"])  
plt.title("PCA - variance explained")  
plt.ylabel("Explained variance Ratio")  
plt.xlabel("principal components")  
plt.ylim(0,1)  
plt.show()
```

PCA variance explained.



Register No. :

Experiment No. :

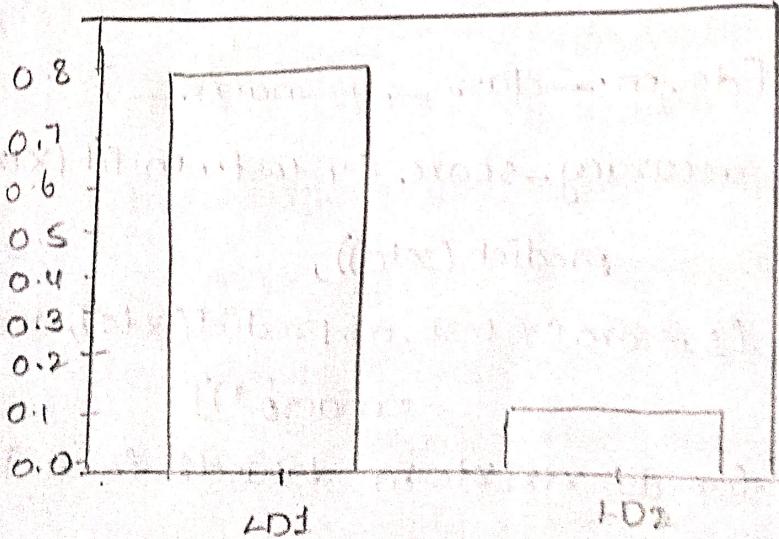
Date :

code(LDA).

 $\text{lda} = \text{LDA}(\text{n_components}=2)$ $\text{lda} = \text{fit}(\text{x_train_scaled}, \text{y_train})$ $\text{explained_var_ratio} = \text{lda.explained_variance_ratio}$

LABORATORY CONTINUOUS EVALUATION			
S.No.	PROGRAM	Max.Marks	Marks Scored
1.	Preparedness	5	
2.	Coding	5	
3.	Testing	5	
4.	Viva	5	
	TOTAL	20	

Faculty Signature with Date

`plt.figure(figsize=(6,4))``plt.bar(["LD1", "LD2"], explained_var_ratio, color=["#FFFFE0", "#FFD966"])``plt.title("LDA - explained variance ratio")``plt.ylabel("proportion of variance")``plt.show()``print("explained variance ratio (LD1, LD2):", explained_var_ratio)``LDA - explained variance Ratio.``explained variance ratio (LD1, LD2): [0.832755, 0.08835]`

2) (Testing) Train logistic regression and Decision-tree on:

- original dataset
- PCA-reduced dataset
- LDA-reduced dataset

Code % from sklearn.linear_model import LogisticRegression.

from sklearn.tree import DecisionTreeClassifier.

from sklearn.metrics import accuracy_score, f1_score.

models = [LogisticRegression(max_iter=2000,

random_state=42),

DecisionTreeClassifier(random_state=42)]

datasets = [("original", x_train_scaled, x_test_scaled)

("PCA", x_train_pca, x_test_pca),

("LDA", x_train_lda, x_test_lda)].

results = [[ds, m.class_name,

accuracy_score(y_test, m.fit(xtr, y_train)

predict(xte)),

f1_score(y_test, m.predict(xte), average="macro")]]

for ds, xtr, xte in datasets for m in models]

for r in results:

print(r).

Output %

[('original', LogisticRegression, 0.6244, 0.2841)]

[('original', DecisionTreeClassifier, 0.633, 0.3501)]



[PCA, Logistic Regression, 0.550, 0.2013]

[PCA, Decision Tree classifier, 0.55, 0.26]

[LDA, Logistic Regression, 0.628, 0.2818]

[LDA, Decision Tree classifier, 0.585, 0.301]

LABORATORY CONTINUOUS EVALUATION			
S.No.	PROGRAM	Max.Marks	Marks Scored
1.	Preparedness	5	
2.	Coding	5	
3.	Testing	5	
4.	Viva	5	
	TOTAL	20	

Faculty Signature with Date

(i) original dataset :

models = {

"Log Reg": LogisticRegression(max_iter=2000, random_state=42),

"Decision Tree": DecisionTreeClassifier(random_state=42)

}

results = []

for name, model in models.items():

y_pred = model.fit(x_train_scaled, y_train).predict(x_test_scaled);

results.append([name, accuracy_score(y_test, y_pred)])

f1_score(y_score, y_pred, average="macro")])

sns. heatmap(confusion_matrix(y_test, y_pred),

annot=True, fmt="d", cmap="Blues")

plt.title(f"matrix = {name}"); plt.show()

df = pd.DataFrame(results, columns=["model", "Accuracy", "F1 score"])

print(df)

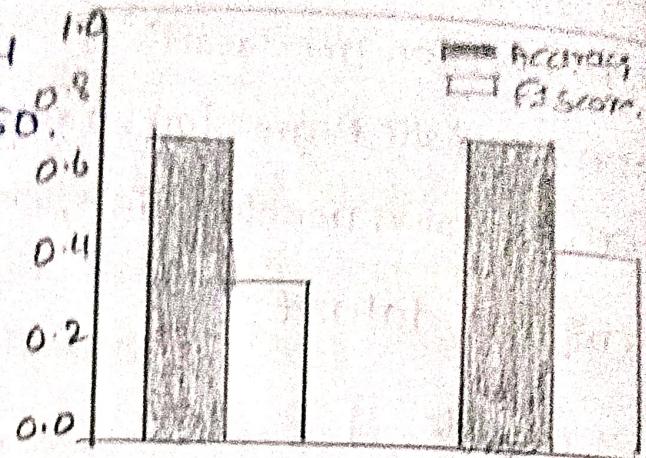
df.set_index("Model")[["Accuracy", "F1 score"]].plot(kind="bar",

Ylim=(0,1), figsize=(8,11), title="originalDataset") plt.show()

Output 8

model	accuracy	f1score
LogReg	0.694	0.684
DecisionTree	0.633	0.350

Original Dataset



pca Reduced Dataset

```

code: results = []
for name, model in models.items():
    model.fit(x_train_pca, y_train)
    y_pred = model.predict(x_test_pca)
    results.append({name: accuracy_score(y_test, y_pred),
                    f1-score(y_test, y_pred, average="macro")})
df = pd.DataFrame(results, columns = ["model", "Accuracy", "F1 Score"])
print(df)
df.set_index("model")[["Accuracy", "F1 Score"]].plot(kind="bar", ylim=[0,1], figsize=(6,4), title="PCA Dataset")
plt.show()

```

Decision Boundaries.

```

def plot_db(x, y, model, title):
    xx, yy = np.meshgrid(np.linspace(x[:, 0].min() - 1,
                                      x[:, 0].max() + 1, 200), np.linspace(x[:, 1].min() - 1,
                                      x[:, 1].max() + 1, 200))

```

2. `z = model.predict(np.c_[xx.ravel(), yy.ravel()]).reshape(xx.shape).`

`plt.scatter(x[:,0], x[:,1], c=y, cmap="coolwarm", edgecolor="k", s=20).`

`plt.title(title); plt.xlabel("PC1"); plt.ylabel("PC2"); plt.show()`

for name, model in models.items():

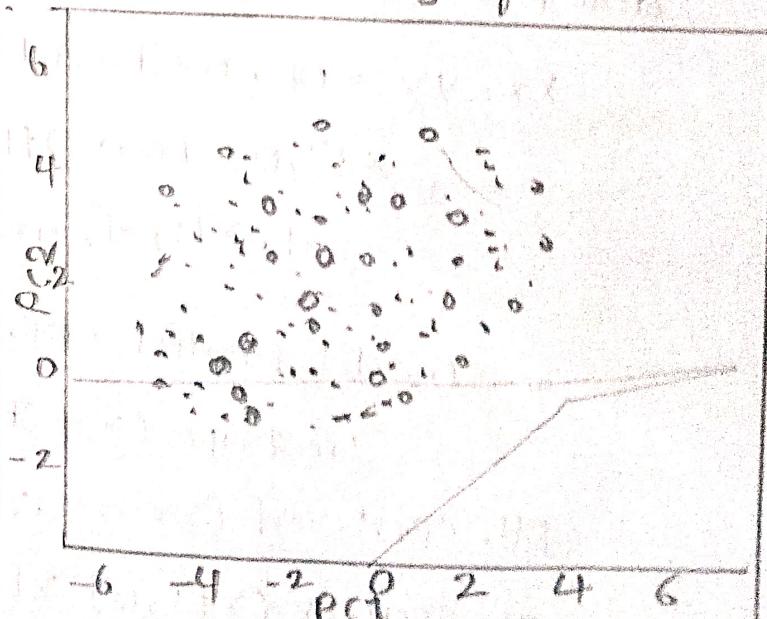
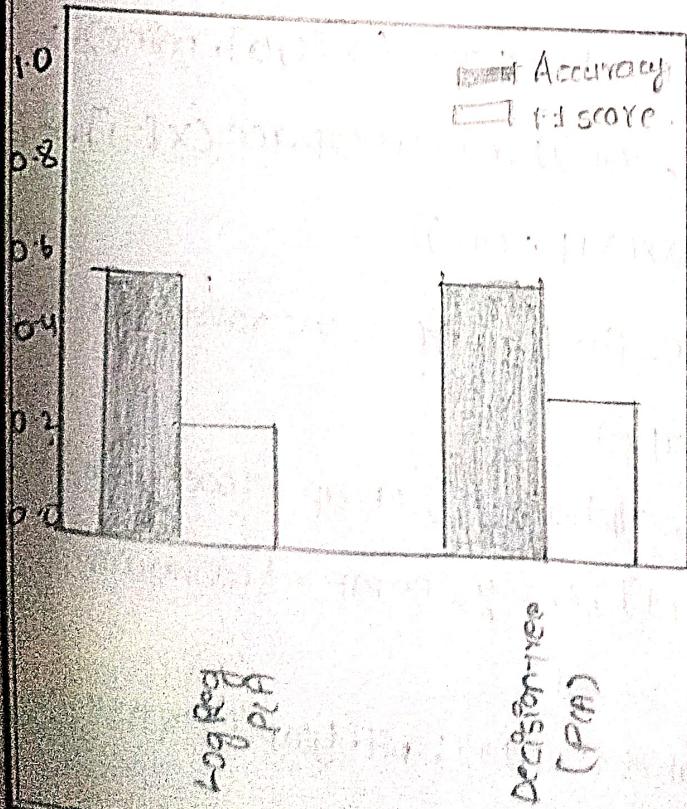
`plot_db(x-train-pca, y-train, model, f"PCA-{name}").`

outputs:

	model	Accuracy	F1 score
0	Log Reg(PCA)	0.550	0.201
1	DecisionTree(PCA)	0.550	0.261.

PCA Dataset

PCA - LogReg(PCA)



LDA Reduced Dataset

```
results = []
for name, model in models.items():
    model.fit(X_train_lda, y_train)
    y_pred = model.predict(X_test_lda)
    results.append([name, accuracy_score(y_test, y_pred),
                    f1_score(y_test, y_pred, average="macro")])
df = pd.DataFrame(results, columns=[["model"], "Accuracy",
                                     ["f1 score"]])
print(df)
df.set_index("model")["Accuracy", "f1 score"].plot(kind="bar", ylim=(0,1), figsize=(6,4), title="LDA dataset")
plt.show()

def decision_boundary():
    def plot_db(X, Y, model, title):
        xx, yy = np.meshgrid(np.linspace(X[:, 0].min() - 1,
                                         X[:, 0].max() + 1, 200),
                             np.linspace(X[:, 1].min() - 1,
                                         X[:, 1].max() + 1, 200))
        z = model.predict(np.c_[xx.ravel(), yy.ravel()])
        z = z.reshape(xx.shape)
        plt.contourf(xx, yy, z, alpha=0.3, cmap="coolwarm")
        plt.scatter(X[:, 0], X[:, 1], c=y, cmap="coolwarm",
                    edgecolor="k", s=20)
        plt.title(title); plt.xlabel("LD1"); plt.ylabel("LD2");
decision_boundary()
```

pit show()

for name in model.items():

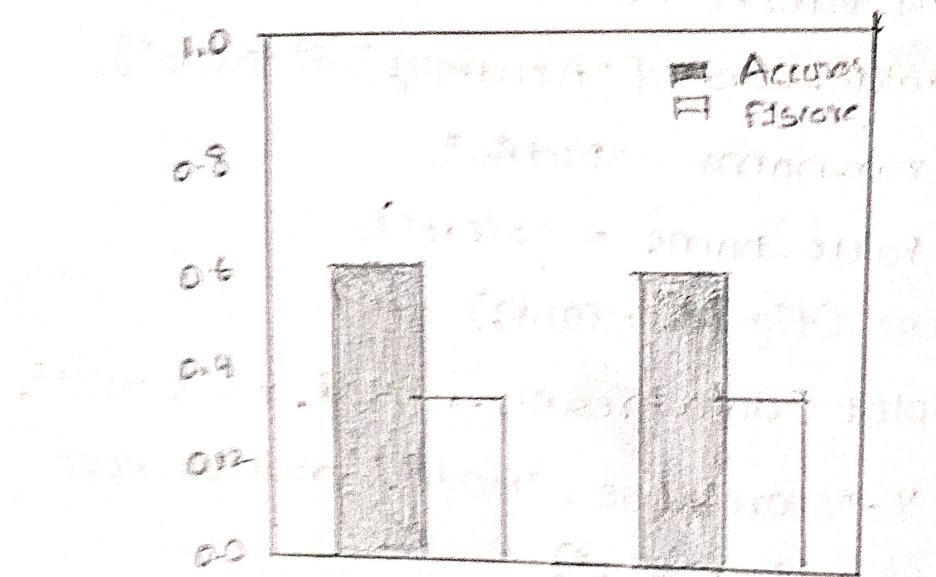
plot_db(x_train_lda, y_train,
model, f" LDA - {name}").

LABORATORY CONTINUOUS EVALUATION			
S.NO.	PROGRAM	Max Marks	Marks Secured
1.	Preparations	5	5
2.	Coding	5	5
3.	Testing	5	5
4.	Viva	5	5
	TOTAL		20

Faculty Signature with Date

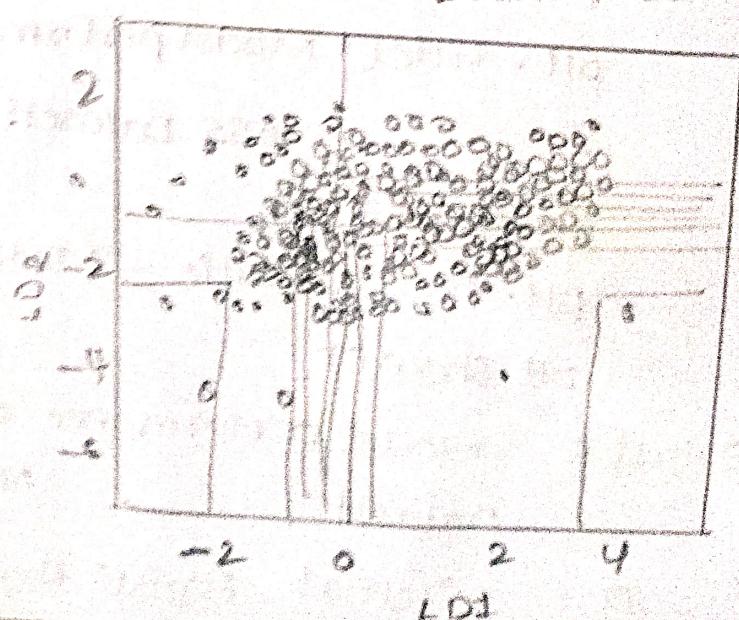
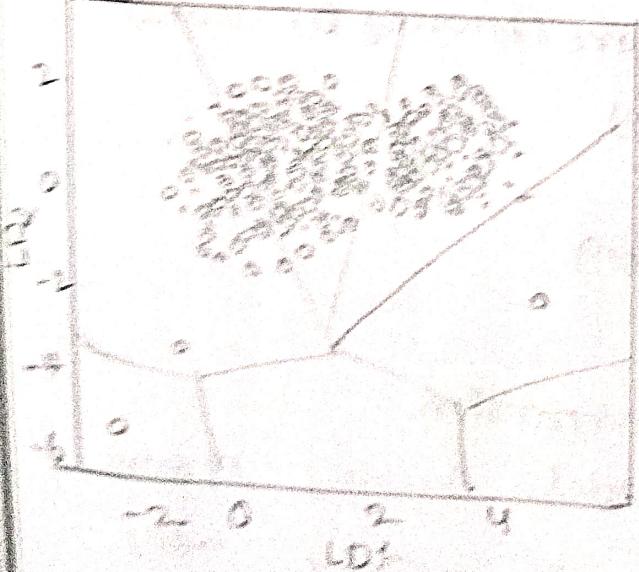
output :-

	model	Accuracy	f1 score
0	Log Reg (LDA)	0.628	0.287
1	DecisionTree(LDA)	0.585	0.3018



LDA - Verbose

model = LDA - DecisionTree(LDA)



Model performance comparison

```
import seaborn as sns  
results_df = pd.DataFrame(results, columns=[  
    "Dataset", "Model", "Accuracy", "F1 score"])  
print("In Model performance comparison.")  
print(results_df)  
results_melted = results_df.melt(  
    id_vars=["Dataset", "Model"],  
    value_vars=["Accuracy", "F1 score"],  
    var_name="Metric",  
    value_name="Score").  
plt.figure(figsize=(10, 6)).  
sns.barplot(data=results_melted, x="Dataset",  
            y="Score", hue="Model", errorbar=None,  
            palette="Set2").  
plt.ylabel("Score")  
plt.title("Model performance (Accuracy & F1 score)  
Across Datasets").
```

```
plt.ylim(0, 1)  
plt.legend(title="Model")  
plt.show().
```

Output :- Model performance comparison.

Dataset	Model	Accuracy	F1 score
0	Original Logistic Regression	0.624	0.284

1. original DecisionTreecla 0.633 0.350
2. PCA LogisticRegression 0.55 0.20
3. PCA DecisionTreecla 0.55 0.26
4. LDA LogisticRegression 0.62 0.28
5. LDA DecisionTreeclass 0.58 0.30.

LABORATORY CONTINUOUS EVALUATION			
S.No.	PROGRAM	Max.Marks	Marks Scored
1.	Preparedness	5	
2.	Coding	5	
3.	Testing	5	
4.	Viva	5	
	TOTAL	20	

Faculty Signature with Date

6 Model performance (Accuracy & F1 score) Across Datasets

