

```
In 1 1 # Importing required libraries
2 import numpy as np
3 import pandas as pd
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import LabelEncoder
6 import seaborn as sns
7 import matplotlib.pyplot as plt
8 from sklearn.metrics import accuracy_score
```

Executed at 2023.10.29 19:29:06 in 409ms

```
In 2 1 # Data Collection for experiment
2 df = pd.read_csv("../Iris.csv")
3 # Display the first few rows of the dataset to inspect its structure and content.
4 print("First 5 rows of the dataset:-\n", df.head())
```

Executed at 2023.10.29 19:29:06 in 40ms

First 5 rows of the dataset:-

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

```
In 3 1 # Check the dimensions of the dataset (number of rows and columns).
2 row, col = df.shape
3 print("No. of rows in the dataset: ", row)
4 print("No. of column in the dataset: ", col)
```

Executed at 2023.10.29 19:29:06 in 36ms

No. of rows in the dataset: 150

No. of column in the dataset: 5

```
In 4 1 # Identify the data types of each column (numeric, categorical, text, etc.).
2 print("Data types of each column:\n", df.dtypes)
```

Executed at 2023.10.29 19:29:06 in 54ms

Data types of each column:

SepalLengthCm	float64
SepalWidthCm	float64
PetalLengthCm	float64
PetalWidthCm	float64
Species	object
dtype:	object

```
In 5 1 # Display the number of missing values in each column
2 missingValues = df.isnull().sum()
3 print("Missing values per column:-")
4 print(missingValues)
```

Executed at 2023.10.29 19:29:06 in 50ms

```
Missing values per column:-
SepalLengthCm    0
SepalWidthCm     0
PetalLengthCm    0
PetalWidthCm     0
Species          0
dtype: int64
```

```
In 6 1 # Numerical Columns
2 numerical_columns = df.select_dtypes(include=['int64', 'float64']).columns
3 # Categorical Columns
4 categorical_columns = df.select_dtypes(include='object').columns
Executed at 2023.10.29 19:29:06 in 256ms
```

```
In 7 1 # Finding Unique categories of species column
2 print("Types of Species: ", df['Species'].unique())
Executed at 2023.10.29 19:29:06 in 284ms

Types of Species: ['Iris-setosa' 'Iris-versicolor' 'Iris-virginica']
```

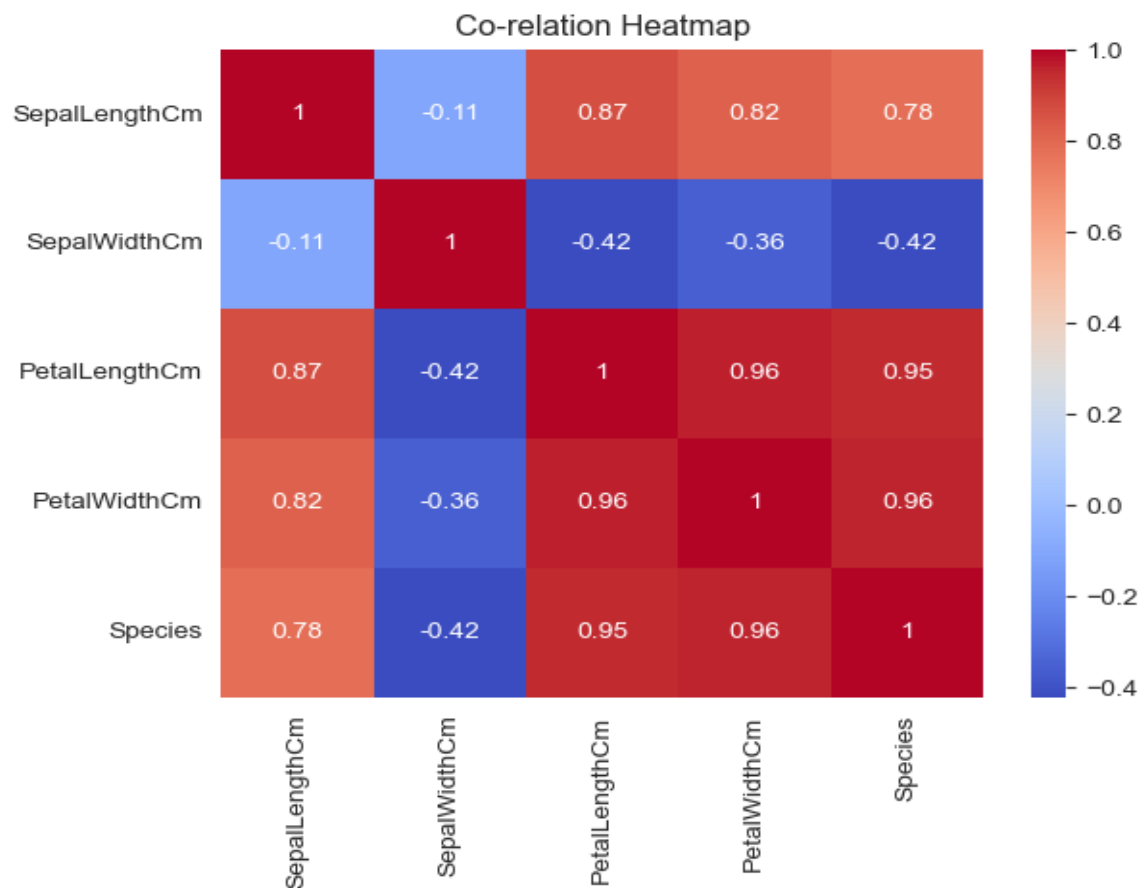
```
In 8 1 # Label Encoding
2 data_encoded_label = df.copy()
3 print("Categorical columns Before Label Encoding:-\n", data_encoded_label[categorical_columns].head())
Executed at 2023.10.29 19:29:06 in 279ms
```

```
Categorical columns Before Label Encoding:-
Species
0  Iris-setosa
1  Iris-setosa
2  Iris-setosa
3  Iris-setosa
4  Iris-setosa
```

```
In 9 1 # Fit the label encoder to the categorical columns
2 label_encoder = LabelEncoder()
3 for column in categorical_columns:
4     data_encoded_label[column] = label_encoder.fit_transform(data_encoded_label[column])
5
6 # Displays the first few rows of the encoded data
7 print("Categorical columns After Label Encoding:-\n", data_encoded_label[categorical_columns].tail())
Executed at 2023.10.29 19:29:06 in 284ms
```

```
Categorical columns After Label Encoding:-
Species
145      2
146      2
147      2
148      2
149      2
```

```
In 10 1 # Plotting the Co-relation between different features
2 sns.heatmap(data_encoded_label.corr(), annot=True, cmap='coolwarm')
3 plt.title("Co-relation Heatmap")
4 plt.show()
Executed at 2023.10.29 19:29:07 in 561ms
```



```
In 11 1 # Features and target variable
2 X = data_encoded_label.drop(columns=['Species']) # Features
3 Y = data_encoded_label['Species'] # Target Variable
Executed at 2023.10.29 19:29:07 in 16ms
```

```
In 12 1 # Split the dataset into training and testing sets(80% training, 20%testing)
2 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.2, random_state=42)
3 # Display the size of the training and testing sets
4 print(f'Training set size: {X_train.shape[0]} samples \nTest set size: {X_test.shape[0]} samples')
Executed at 2023.10.29 19:29:07 in 10ms

Training set size: 120 samples
Test set size: 30 samples
```

```
In 13 1 from sklearn.svm import SVC
2
3 # Create an SVM model
4 svm_model = SVC(kernel='rbf')
Executed at 2023.10.29 19:29:07 in 135ms
```

```
In 14 1 # Train the model
2 svm_model.fit(X_train, y_train)
3 # Making Predictions
4 y_predict_rbf = svm_model.predict(X_test)
Executed at 2023.10.29 19:29:07 in 61ms
```

```
In 15 1 # Accuracy scores
2 acc_rbf = accuracy_score(y_test, y_predict_rbf)
3 # Print or visualize the accuracy scores
4 print(f'Accuracy RBF: {acc_rbf}')
```

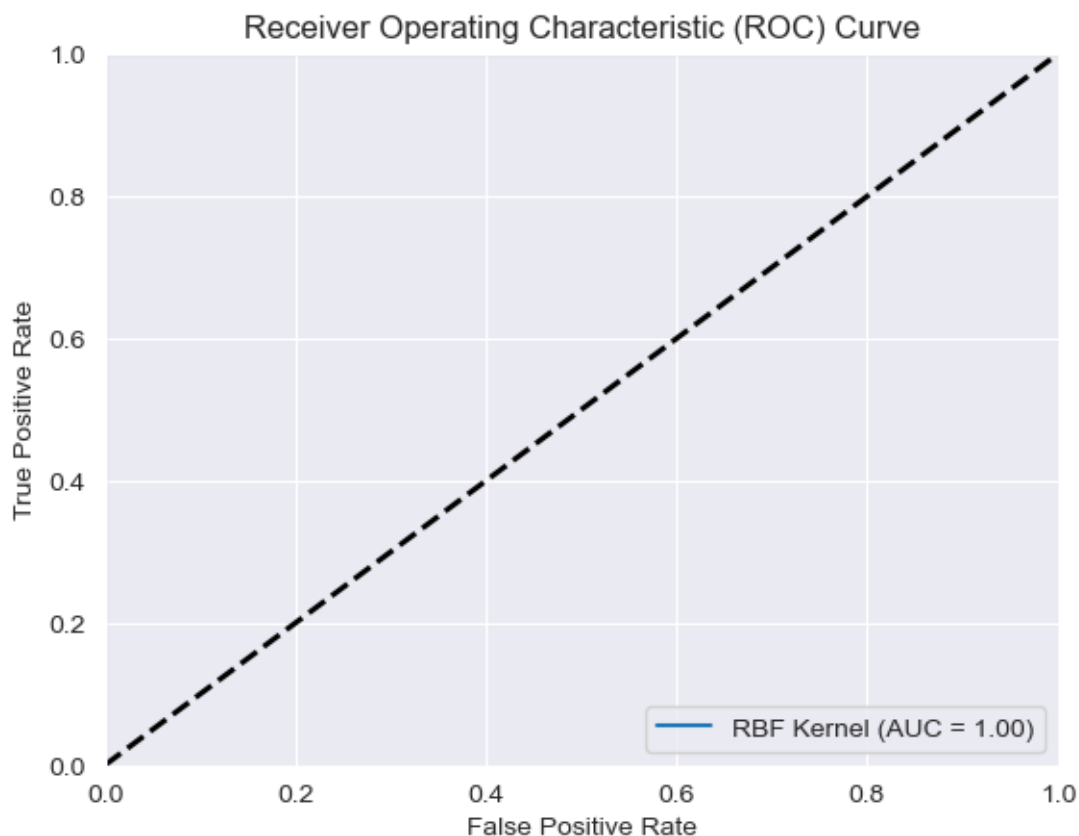
Accuracy RBF: 1.0

```

In 16 1  from sklearn.metrics import roc_curve, auc
      2  from sklearn.preprocessing import label_binarize
      3
      4  # Binarize the output
      5  y_test_bin = label_binarize(y_test, classes=[0, 1, 2])
      6
      7  # Get decision function values for each model
      8  y_score_rbf = svm_model.decision_function(X_test)
      9  # Compute ROC curve and ROC area for each class
     10  fpr_rbf, tpr_rbf, _ = roc_curve(y_test_bin.ravel(), y_score_rbf.ravel())
     11  # Compute AUC for each class
     12  auc_rbf = auc(fpr_rbf, tpr_rbf)
     13
     14  # Plot ROC curves
     15  plt.plot(fpr_rbf, tpr_rbf, label=f'RBF Kernel (AUC = {auc_rbf:.2f})')
     16  plt.plot([0, 1], [0, 1], 'k--', lw=2)
     17  plt.xlim([0.0, 1.0])
     18  plt.ylim([0.0, 1.0])
     19  plt.xlabel('False Positive Rate')
     20  plt.ylabel('True Positive Rate')
     21  plt.title('Receiver Operating Characteristic (ROC) Curve')
     22  plt.legend(loc="lower right")
     23  plt.show()

```

Executed at 2023.10.29 19:29:07 in 340ms



AIM:-

Implement Principal Component Analysis Algorithm and use it to reduce dimension of iris dataset.

Consider the following instructions:

- Plot the magnitude of eigen values in sorted order.
- Plot the reconstructed data points along with the class labels using 1 and 2 PCs for reconstruction.
- Classify the dimension reduced dataset using Bayes Classifier.

```
In 17 1 # Step 1: Calculate Covariance Matrix
```

```
2 cov_matrix = np.cov(X, rowvar=False)
```

Executed at 2023.10.29 19:29:07 in 16ms

```
In 18 1 # Step 2: Compute Eigenvalues and Eigenvectors
```

```
2 eigenvalues, eigenvectors = np.linalg.eigh(cov_matrix)
```

Executed at 2023.10.29 19:29:07 in 10ms

```
In 19 1 # Step 3: Sort Eigenvalues
```

```
2 sorted_indices = np.argsort(eigenvalues)[::-1]
```

```
3 eigenvalues = eigenvalues[sorted_indices]
```

```
4 eigenvectors = eigenvectors[:, sorted_indices]
```

Executed at 2023.10.29 19:29:07 in 83ms

```
In 20 1 # Step 6: Plot Eigenvalues
```

```
2 plt.plot(eigenvalues, marker='o')
```

```
3 plt.title('Eigenvalues in Sorted Order')
```

```
4 plt.xlabel('Principal Component Index')
```

```
5 plt.ylabel('Eigenvalue Magnitude')
```

```
6 plt.show()
```

Executed at 2023.10.29 19:29:08 in 293ms



```
In 21 1 # Step 4: Choose Principal Components (1 and 2 in this case)
2 selected_components = eigenvectors[:, :2]
Executed at 2023.10.29 19:29:08 in 15ms
```

```
In 22 1 from sklearn.decomposition import PCA
2
3 # Step 5: Project Data using PCA with 2 components
4 pca = PCA(n_components=2)
5 X_pca = pca.fit_transform(X)
Executed at 2023.10.29 19:29:08 in 52ms
```

```
In 23 1 # Reconstruct the data using 1 and 2 PCs
2 X_reconstructed_1pc = np.dot(X_pca[:, :1], pca.components_[1, :])
3 X_reconstructed_2pc = np.dot(X_pca, pca.components_)
Executed at 2023.10.29 19:29:08 in 18ms
```

```
In 24 1 from sklearn.naive_bayes import GaussianNB
2
3 # Step 8: Bayes Classifier
4 X_train, X_test, y_train, y_test = train_test_split(X_pca, Y, test_size=0.3, random_state=42)
5 classifier = GaussianNB()
6 classifier.fit(X_train, y_train)
Executed at 2023.10.29 19:29:08 in 102ms
```

Out 24 ▾

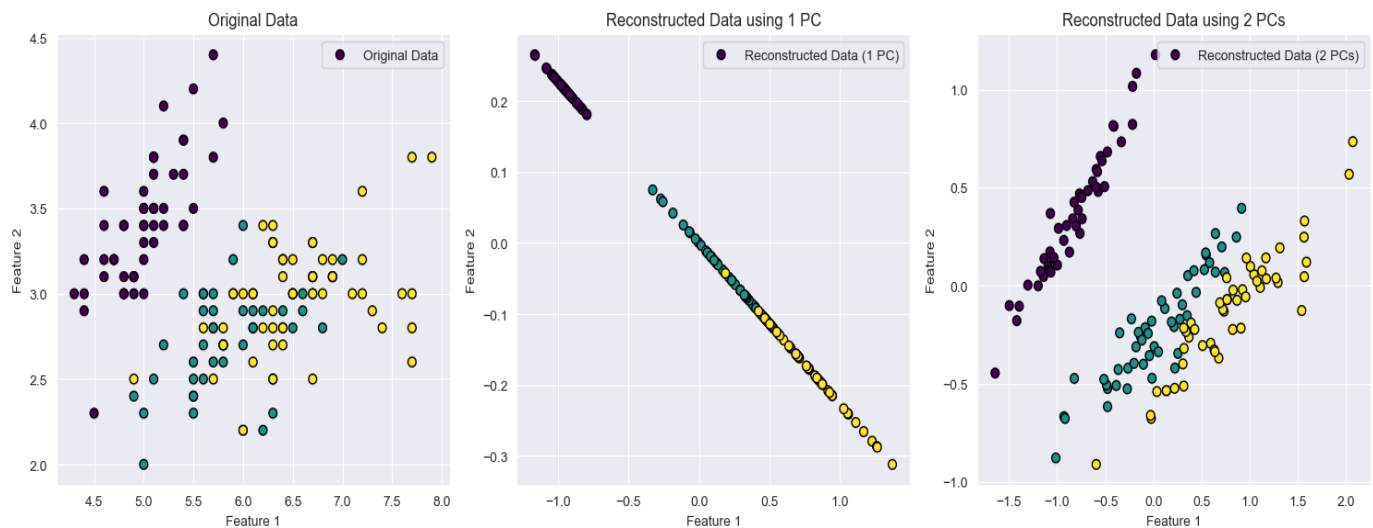
▾ GaussianNB
 GaussianNB()

```
In 25 1 # Step 9: Classify Data
2 y_predict = classifier.predict(X_test)
3 accuracy = accuracy_score(y_test, y_predict)
4 print(f'Accuracy of Bayes Classifier: {accuracy}')
```

Executed at 2023.10.29 19:29:08 in 83ms

Accuracy of Bayes Classifier: 0.8888888888888888

```
In 26 1 # Comparison Plot
2 plt.figure(figsize=(15, 5))
3
4 # Original Data
5 plt.subplot(1, 3, 1)
6 plt.scatter(X.values[:, 0], X.values[:, 1], c=Y, cmap='viridis', edgecolor='k', s=40, label='Original Data')
7 plt.title('Original Data')
8 plt.xlabel('Feature 1')
9 plt.ylabel('Feature 2')
10 plt.legend()
11
12 # Reconstructed Data using 1 PC
13 plt.subplot(1, 3, 2)
14 plt.scatter(X_reconstructed_1pc[:, 0], X_reconstructed_1pc[:, 1], c=Y, cmap='viridis', edgecolor='k',
15             s=40, label='Reconstructed Data (1 PC)')
16 plt.title('Reconstructed Data using 1 PC')
17 plt.xlabel('Feature 1')
18 plt.ylabel('Feature 2')
19 plt.legend()
20
21 # Reconstructed Data using 2 PCs
22 plt.subplot(1, 3, 3)
23 plt.scatter(X_reconstructed_2pc[:, 0], X_reconstructed_2pc[:, 1], c=Y, cmap='viridis', edgecolor='k',
24             s=40, label='Reconstructed Data (2 PCs)')
25 plt.title('Reconstructed Data using 2 PCs')
26 plt.xlabel('Feature 1')
27 plt.ylabel('Feature 2')
28 plt.legend()
29
30 plt.tight_layout()
31 plt.show()
Executed at 2023.10.29 19:29:09 in 1s 418ms
```



Submitted By,

Name:- Shankar Singh Mahanty

Regd. No:- 2101020758

Roll No:- CSE21238

Group:- 3

Sem:- 5th

Branch:- CSE