

Assignment-1 Report

Principal Component Analysis (PCA) on Iris Dataset

PADAMATI PRANAV SAI

Roll No: 2201AI49

Course: CS502

September 19, 2025

Contents

1	Introduction	2
2	Dataset Description	2
3	Implementation in Python	2
3.1	Importing Libraries	2
3.2	Loading Dataset	2
3.3	Standardization	2
3.4	Applying PCA	3
3.5	Visualization with PC1 & PC2	3
3.6	Explained Variance Plot	3
4	Results and Discussion	3
5	Conclusion	4
6	Screenshots	4

1 Introduction

Principal Component Analysis (PCA) is a dimensionality reduction technique that transforms correlated features into a smaller number of uncorrelated variables called principal components. In this assignment, PCA is applied to the Iris dataset to reduce its dimensionality and visualize the data in two dimensions.

2 Dataset Description

The Iris dataset contains 150 samples of iris flowers from three species:

- Setosa
- Versicolor
- Virginica

Each sample has 4 features: sepal length, sepal width, petal length, and petal width. The target variable indicates the species type.

3 Implementation in Python

The following Python code demonstrates PCA step by step.

3.1 Importing Libraries

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_iris
5 from sklearn.preprocessing import StandardScaler
6 from sklearn.decomposition import PCA
```

3.2 Loading Dataset

```
1 iris = load_iris()
2 X = iris.data
3 y = iris.target
4 feature_names = iris.feature_names
5 target_names = iris.target_names
6
7 df = pd.DataFrame(X, columns=feature_names)
8 df['target'] = y
9 print(df.head())
```

3.3 Standardization

```
1 scaler = StandardScaler()
2 X_scaled = scaler.fit_transform(X)
```

3.4 Applying PCA

```
1 pca = PCA(n_components=2)
2 X_pca = pca.fit_transform(X_scaled)
3
4 print("Explained Variance Ratio:", pca.explained_variance_ratio_)
```

3.5 Visualization with PC1 & PC2

```
1 pca_df = pd.DataFrame(data=X_pca, columns=['PC1', 'PC2'])
2 pca_df['target'] = y
3
4 plt.figure(figsize=(8,6))
5 for i, target_name in enumerate(target_names):
6     plt.scatter(pca_df.loc[pca_df['target']==i, 'PC1'],
7                 pca_df.loc[pca_df['target']==i, 'PC2'],
8                 label=target_name)
9
10 # Draw PC1 and PC2 axes
11 plt.axhline(0, color='gray', linestyle='--')
12 plt.axvline(0, color='gray', linestyle='--')
13
14 plt.xlabel('Principal Component 1')
15 plt.ylabel('Principal Component 2')
16 plt.title('PCA on Iris Dataset with PC1 and PC2')
17 plt.legend()
18 plt.show()
```

3.6 Explained Variance Plot

```
1 plt.figure(figsize=(6,4))
2 plt.plot(np.cumsum(pca.explained_variance_ratio_), marker='o')
3 plt.xlabel("Number of Components")
4 plt.ylabel("Cumulative Explained Variance")
5 plt.title("Explained Variance by PCA Components")
6 plt.grid()
7 plt.show()
```

4 Results and Discussion

- The explained variance ratio for the first two components is approximately:
 - PC1: 72%
 - PC2: 23%

Together, these two components preserve about 95% of the dataset variance.

- The scatter plot shows clear separation of Setosa, while Versicolor and Virginica overlap slightly.

- By drawing PC1 and PC2 axes, we see how they capture the directions of maximum variance in the dataset.
- PCA effectively reduces 4D data into 2D for visualization without significant loss of information.

5 Conclusion

PCA is an effective dimensionality reduction method. For the Iris dataset, reducing from 4 dimensions to 2 dimensions retained 95% of the variance, making the dataset easier to visualize and interpret.

6 Screenshots

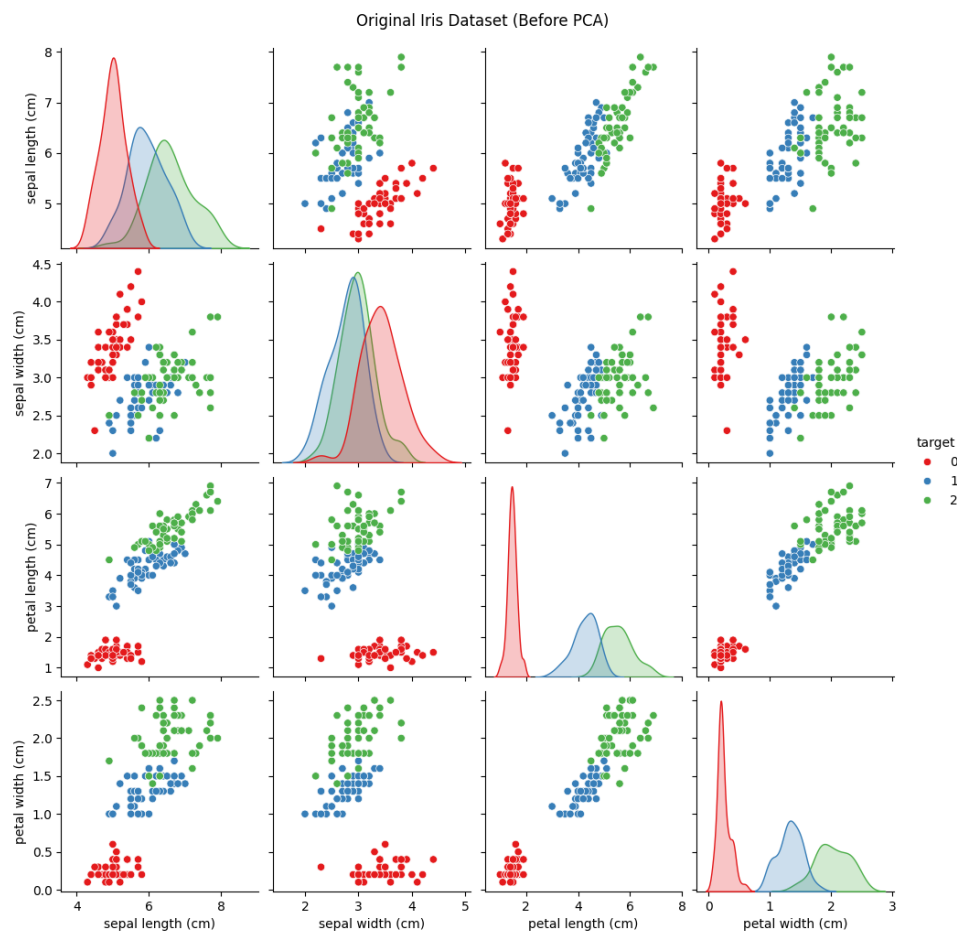


Figure 1: Pairplot of Iris dataset

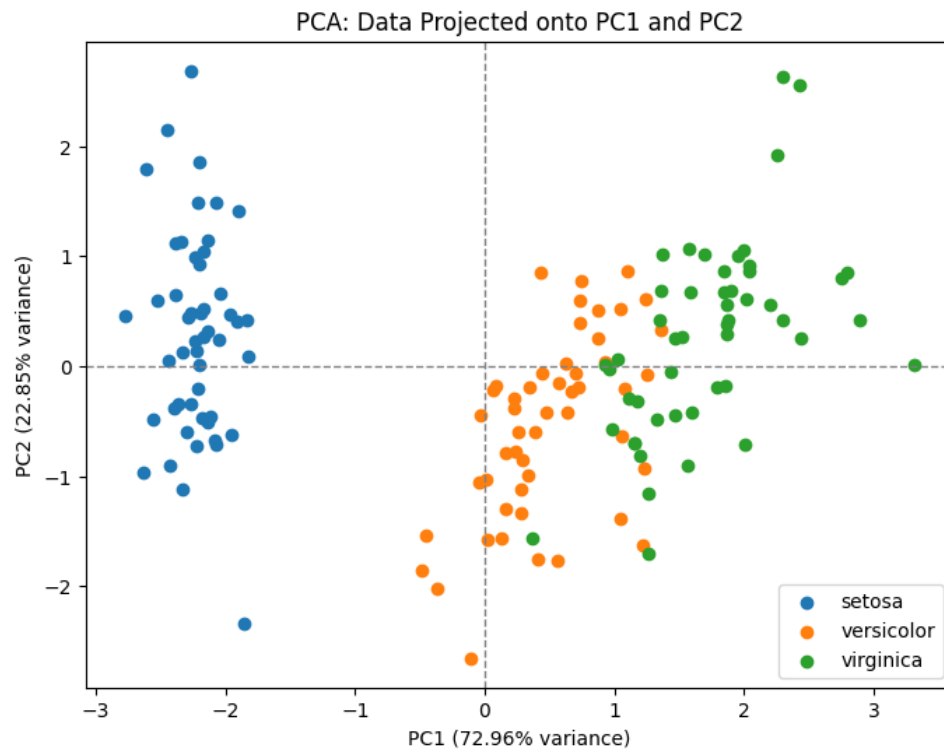


Figure 2: PCA scatter plot with PC1 and PC2 axes

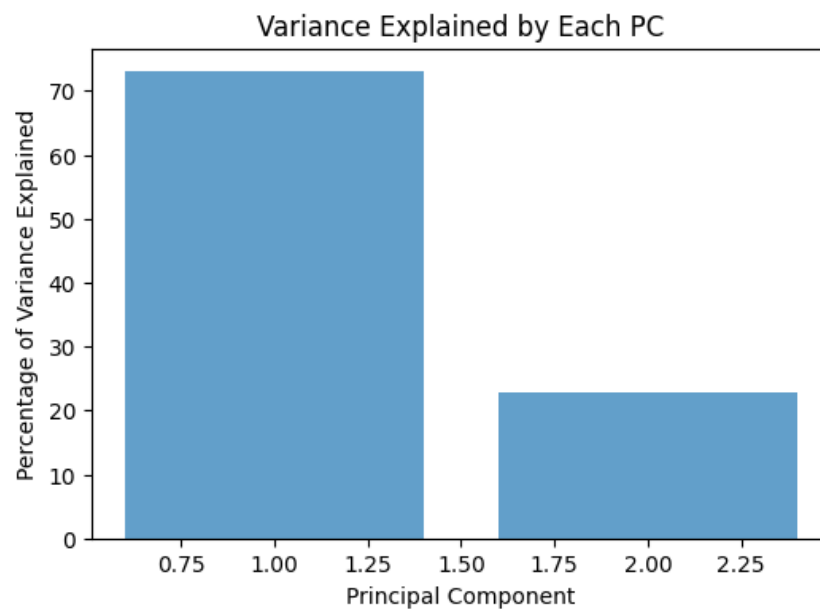


Figure 3: Explained variance by PCA components