

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

A Comparative Analysis of PCA for Dimensionality Reduction in Diverse Datasets

PRABIN BOHARA¹, PRABIN SHARMA POUDEL², AND SANTOSH PANDEY³

¹Institute of Engineering, Thapathali Campus, Kathmandu, Nepal (e-mail: prabinbohara10@gmail.com)

²Institute of Engineering, Thapathali Campus, Kathmandu, Nepal (e-mail: prabinsharmapoudel@gmail.com)

³Institute of Engineering, Thapathali Campus, Kathmandu, Nepal (e-mail: suntoss.pandey@gmail.com)

* All Authors contributed equally to this work.

ABSTRACT High-dimensional datasets pose significant challenges in data analysis and machine learning tasks. With the increasing volume and complexity of data in various domains, deriving insights from such data is becoming challenging, mentally and computationally. Principal Component Analysis (PCA) enables feature extraction, data compression, and easier visualization, facilitating faster training and improving model performance. This report presents a comparative analysis of PCA applied to diverse datasets, including Random Numbers, Iris, and Covertype, with the aim of addressing these challenges. In our study, we applied PCA by standardizing the data, computing the covariance matrix, performing eigenvalue-eigenvector decomposition, and reducing the dimensionality. Our work showcases the versatility of PCA in addressing real-world challenges associated with high-dimensional data analysis and machine learning. Insights gained from this study can guide future applications of PCA for enhancing data understanding, modeling, and decision-making processes.

INDEX TERMS Dimensionality reduction, Eigendecomposition, Principal Component Analysis

I. INTRODUCTION

THE digital revolution has created an era of unprecedented data accumulation, and this immense data footprint presents a formidable challenge for data analysis. Principal Component Analysis (PCA) is a technique that addresses the complexities of high-dimensional data by transforming it into a lower-dimensional representation. By projecting the data onto principal components, PCA aims to capture the essential trends and patterns inherent in the dataset while minimizing the distance between the original data and its transformed representation. This reduction in complexity not only enables more efficient data analysis but also preserves a substantial amount of meaningful information and the underlying structure of the dataset. As a result, PCA has emerged as a powerful tool for navigating and extracting knowledge from the vast realms of data that the digital age has bestowed upon us.

In this report, we explore the application of PCA to three distinct datasets: Random Numbers following a Gaussian distribution, Iris Flower, and the Covertype identification. The primary objective is to investigate the effectiveness of PCA in addressing the challenges associated with high-

dimensional data analysis and machine learning. By applying a standardized approach that includes data standardization, covariance matrix computation, eigenvalue-eigenvector decomposition, and dimensionality reduction, we aim to assess the impact of PCA on data compression, visualization, and model performance. In addition to our customized implementation of Principal Component Analysis (PCA), we also leverage the power of existing PCA libraries for comparative analysis. By employing PCA directly through the library, we aim to evaluate the consistency and performance of our results and explore the benefits of using established PCA implementations in terms of computational efficiency and ease of use.

II. METHODOLOGY

A. THEORY

The variance of data over multiple dimensions is reduced to few prominent dimensions with the application of PCA. This technique figure outs the high quality new components that captures the maximum notion of overall dataset. The number of newly generated components (k) is dependent upon the developer's configuration whose cumulative proportion of

variance is above the required threshold. If an instance of principal component is obtained, the new co-ordinate of dataset is obtained by projecting the original dataset to unit vector of a principal component.

Considering the original points of dataset is spreading over 2D plane (XY-plane) represented by 'x'. Let, the new dimension given by unit vector ' \mathbf{u} ' captures the maximum variance of the original data. So, the new co-ordinates of points 'x' in Figure 1 can be found out by projecting it to the assumed principal component ' \mathbf{u} '.

So, the variance in original dataset with respect to X-axis and Y-axis is given by the formulas:

$$\text{Variance1} = \text{var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (1)$$

$$\text{var}(y) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 \quad (2)$$

However, the relative variance between two or more dimensions gives the relation of variation amongst each other. The co-variance(cov) of dataset in X dimension with respect to Y dimension is given as:

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \quad (3)$$

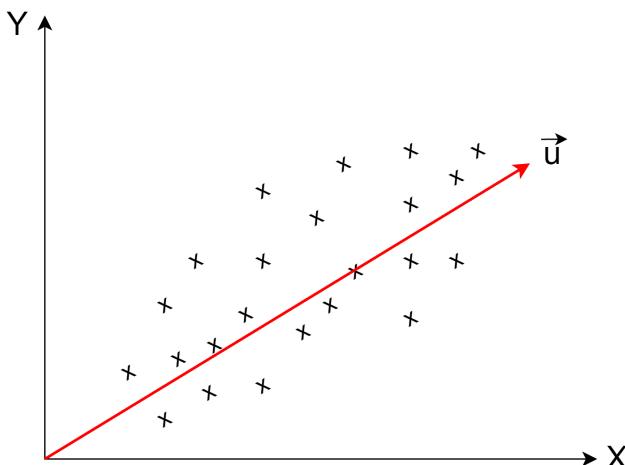


FIGURE 1: Distribution of original Dataset

Thus, the covariance matrix captures all the spread in own dimension and relation between others in remaining dimensions. So, from the relation above,

$$\text{var}(x) = \text{cov}(x, x) \quad (4)$$

$$\text{cov}(x, y) = \text{cov}(y, x) \quad (5)$$

$$\text{Covariance Matrix} = \begin{bmatrix} \text{var}(x) & \text{cov}(x, y) \\ \text{cov}(y, x) & \text{var}(y) \end{bmatrix}$$

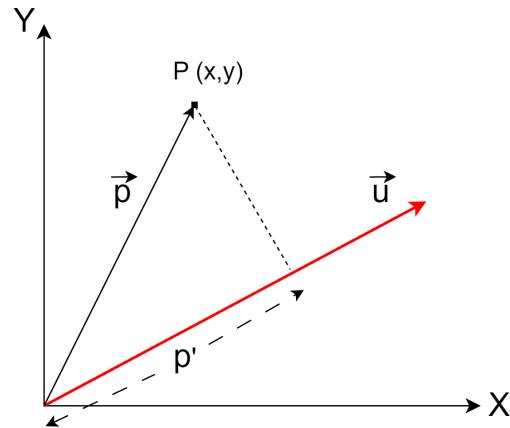


FIGURE 2: Projection of point(P) on assumed principal component

Taking $P(x,y)$ as an arbitrary point in Figure 2, the projection of p over the principal component (u) is calculated as:

$$p' = \frac{\mathbf{u} \cdot \mathbf{p}}{\|\mathbf{u}\|} = \mathbf{u} \cdot \mathbf{p} = \mathbf{u}^T \mathbf{p} \quad (6)$$

Now, the variance of data points with respect to new principal component (u) is given as:

$$\text{Variance2} = \frac{1}{n} \sum_i (u^T x_i - u^T x_m)^2 \quad (7)$$

Thus, the ultimate goal of PCA is to find out the unit vector (u) that maximizes the Variance2 which is an optimization problem. The principal unit vector (u) can be formulated by Rayleigh Quotient which happens to be an Eigenvector with the maximum Eigenvalue of the covariance matrix of the original dataset. This is because linear transformation by covariance matrix to the greatest Eigenvector results in the component that captures the maximum variance of original data. Thus, the Eigendecomposition by the application of transformation makes no linear changes rather than scaling by a factor called Eigenvalue (λ).

B. SYSTEM BLOCK DIAGRAM

The system architecture in Figure 3 shows various blocks in sequence. Diverse datasets were collected and data visualization techniques like scatter plots, histograms, and heatmaps were employed to gain insights into the data's distribution and relationships. Under preprocessing standardization or normalization was done to bring the data to a common scale, ensuring that different features contribute equally to the analysis. Afterward, covariance matrix computation was done to measure the linear relationships between pairs of features. Eigenvalues and eigenvectors of the covariance matrix were computed, so that the top k eigenvectors could be selected to find out the most significant directions of variation in the data. Both the selected eigenvectors and the training

data were arranged row-wise before employing the matrix multiplication between them to obtain the projected final value which could be visualized.

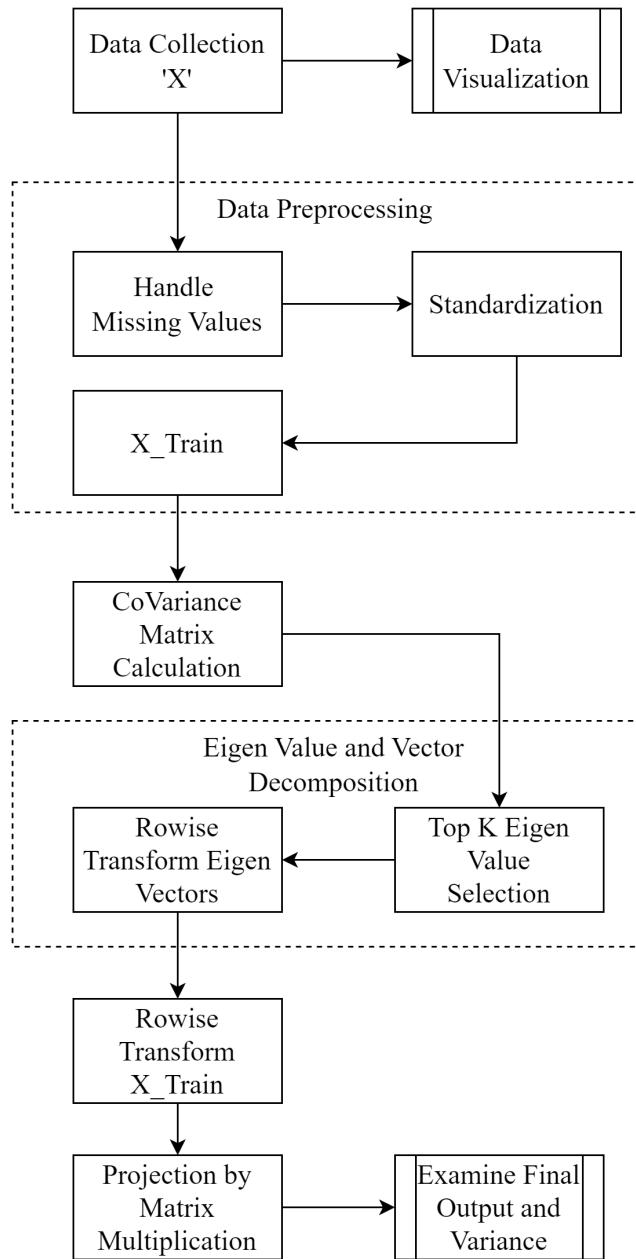


FIGURE 3: System Architecture of PCA

C. INSTRUMENTATION

The implementation of Principal Component Analysis (PCA) in our study involved the use of several libraries and tools to facilitate data visualization and mathematical operations.

1) Pandas

A powerful Python library used for data manipulation, Pandas provides data structures and functions which help us

efficiently work with structured data, such as tabular data, and time series. The exact significance of Pandas is that it's commonly used for data preprocessing, cleaning, exploration, and analysis. It offers efficient data structures, like DataFrame that help us in easy indexing and manipulation of data. We used Pandas to load and read data from CSV files, such as the cancer dataset. It was used to inspect the data, extract features, and preprocess the dataset by handling missing values, encoding categorical variables, and scaling numerical data. For instance, we used label encoding for the target feature in the Cancer dataset, and the use of Pandas made the manipulation process convenient. Additionally, Pandas was used to select specific columns, filter rows based on conditions, and perform descriptive statistics on the data.

2) Matplotlib

Matplotlib is a data visualization tool, that helps us gain insights in visual form. It helps us gain insights, communicate findings, and present results. It offers flexible options for creating line plots, scatter plots, bar plots, histograms, heat maps, and more. We used Matplotlib to create scatter plots, box plots, histograms, and heat maps to visualize different aspects of the data and analyze patterns and relationships. Matplotlib's functions, such as scatter(), boxplot(), hist(), and imshow(), were utilized for these visualizations

3) NumPy

NumPy or Numerical Python provides a multidimensional array object and a collection of functions for performing mathematical operations on arrays efficiently. The general usefulness of NumPy can't be overstated as it offers high-performance array operations, such as mathematical, logical, shape manipulation, and linear algebra operations. We used NumPy to perform various array operations, such as element-wise operations, matrix multiplication, transposing matrices, and calculating eigenvalues and eigenvectors. NumPy's functions, like np.dot(), np.transpose(), and np.linalg.eig(), were used for these tasks. For the application of PCA in particular, np.linalg.eig() was a lot helpful since the eigen vectors primarily defined the direction of the components while eigenvalues gave the relative magnitude of the components. This mathematical insight is what PCA directly relies on, thus the use of NumPy made the implementation much easier for us.

By utilizing these libraries, we were able to seamlessly integrate data visualization and mathematical operations into our PCA implementation. This instrumentation allowed us to effectively analyse the datasets, visualize the results, and gain valuable insights from the dimensionality reduction process.

D. WORKING PRINCIPLE

The working principle of Principal Component Analysis (PCA) involves a pipeline of steps from the input data to the output representation. Here, we provide a detailed description of the PCA pipeline:

Initially, random data was sampled from the Gaussian distribution. Thus, the dataset obtained did not contain any categorical values. To visualize this set of values, a scatter plot is shown in Figure 4. Similarly, for the IRIS, three classes of Iris flowers were represented in the target variable and the dataset consisted of four features: sepal length, sepal width, petal length, and petal width. The initial data can be visualized using a pair plot, as shown in Figure 5, and using a scatterplot in Figure 6. Likewise, the Cancer dataset consists of various features related to cancer tumors, such as radius, texture, perimeter, area, smoothness, and more. It can be visualized using Figure 7.

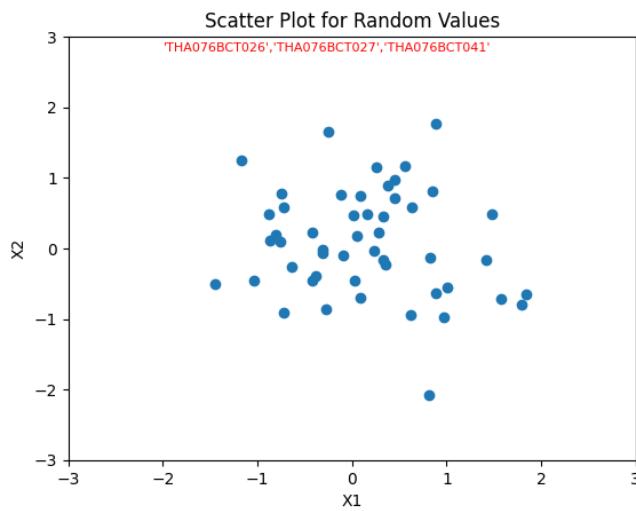


FIGURE 4: Scatter Plot for Random Values

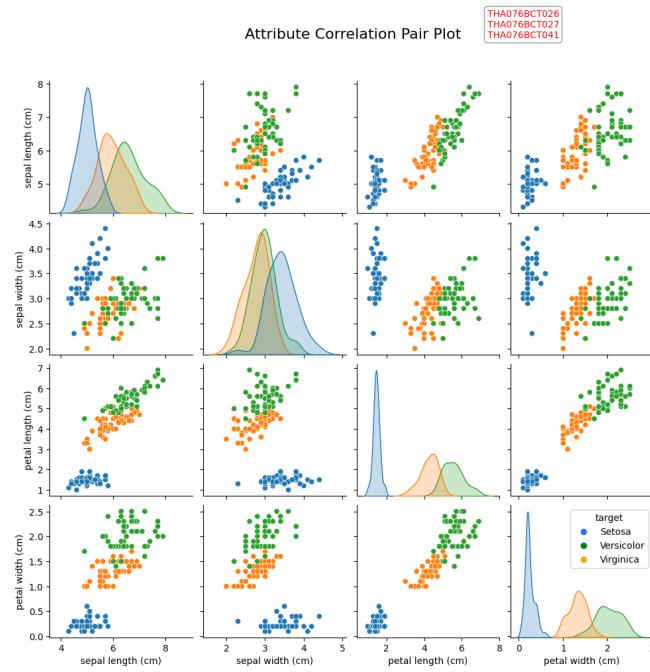


FIGURE 5: Pair Plot for Iris Dataset

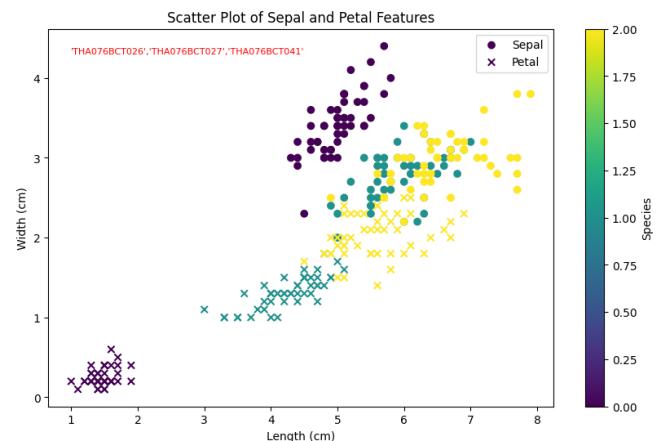


FIGURE 6: Sepal and Petal Features before Preprocessing

1) Data Pre-processing

The input dataset is pre-processed to ensure its suitability for PCA. This typically involves steps such as handling missing values, scaling the data, and removing any outliers. In our implementation, we applied data pre-processing techniques specific to each dataset, such as standardization or normalization. Encoding the categorical target variable was necessary to convert it into a numerical form. A technique like label encoding was used to represent the two classes (malignant and benign) in a numeric format. The processed dataset can be visualized using the given figures.

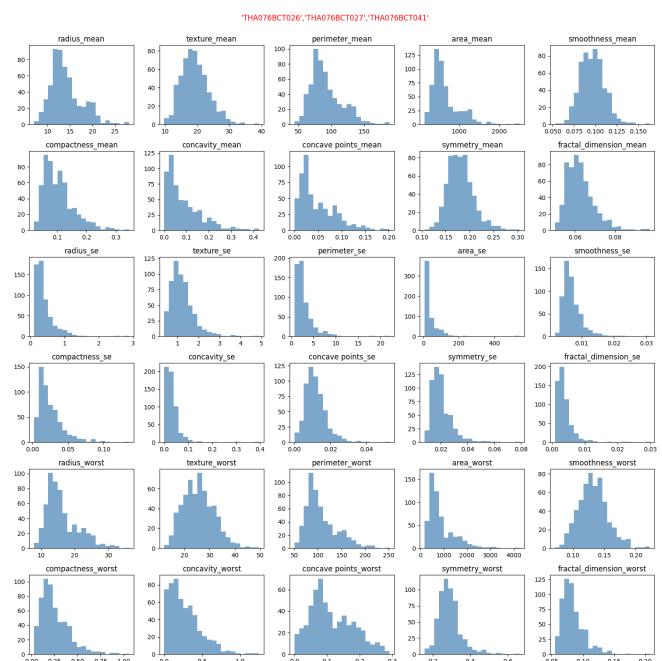


FIGURE 7: Histogram for each feature of Cancer Dataset

2) Covariance Matrix Computation

The covariance matrix is computed based on the pre-processed dataset. The covariance matrix captures the relationships between the variables and provides information about the variance and covariance of the data. This step involves calculating pairwise covariances between the features in the dataset. The covariance matrix can be used to identify highly correlated features.

3) Eigenvalue-Eigenvector Decomposition

The covariance matrix is then decomposed into its eigenvalues and eigenvectors. This step is performed using numerical methods such as the eigenvalue decomposition provided by libraries like NumPy. The eigenvalues represent the variance explained by each eigenvector, while the eigenvectors define the principal components.

4) Selection of Principal Components

The principal components are selected based on the eigenvalues. The eigenvectors with higher eigenvalues capture more variance in the data and are chosen as the principal components. The number of principal components to retain can be determined by analysing the cumulative proportion of variance explained

5) Dimensionality Reduction

The dimensionality of the data is reduced by projecting it onto the selected principal components. The original data is transformed into a lower-dimensional space by multiplying it with the matrix of selected eigenvectors.

6) Output Representation

The transformed data, represented in terms of the principal components, serves as the output of the PCA process. This lower-dimensional representation retains the most important information from the original data while reducing the complexity and dimensionality.

By following this pipeline, PCA enables the extraction of meaningful features, data compression, and easier visualization. The reduced-dimensional representation obtained through PCA can be used for various tasks such as data exploration, clustering, classification, or regression, facilitating more efficient and accurate data analysis and model training.

III. RESULTS

A. RANDOM NUMBERS DATASET

The dataset comprises 20 random samples drawn from a Gaussian distribution, which are visually scattered throughout Figure 4. Before extracting information from these samples, it is necessary to perform PCA analysis. Since the samples follow a Gaussian distribution, they are transformed into a normal distribution by multiplying them with a 2D matrix. This transformation aligns the random samples in a single direction, as depicted in Figure 8.

After computing the eigenvalues and eigenvectors, the principal component (PC1) is determined and visualized in

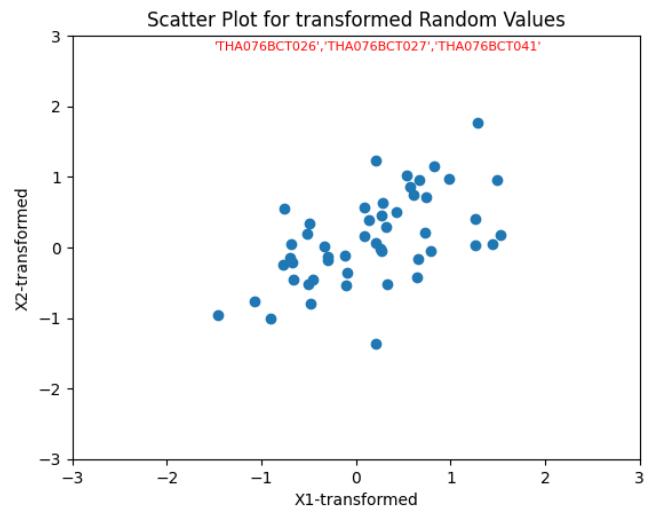


FIGURE 8: Scatter Plot for Transformed Random Values

Figure 9. By selecting only one principal component, all the data points become aligned in a singular direction.

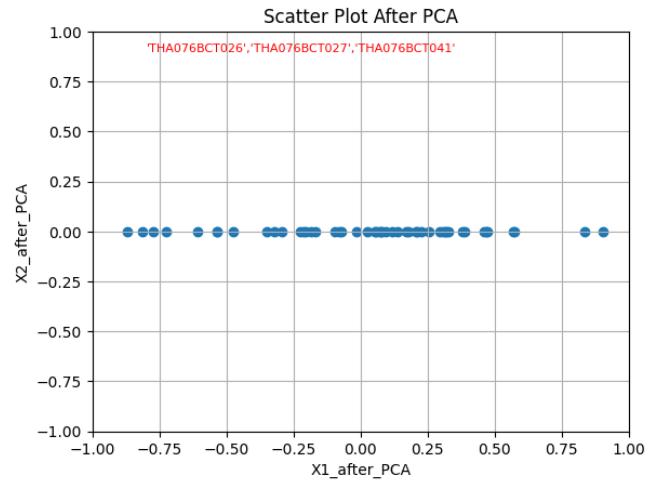


FIGURE 9: Scatter Plot After PCA on Random Values

B. IRIS DATASET

The Iris dataset imported from scikit-learn's library is widely used in machine learning and statistics. It includes measurements of four features from three species of Iris flowers: Setosa, Versicolor, and Virginica. The goal is to classify the flowers based on these features. The dataset has 150 samples, with 50 for each species. It is popular for beginners to learn classification algorithms due to its manageable size and well-structured nature.

Principal Component Analysis (PCA) was applied to extract the most informative components, and scatter plots visualized their directions and magnitudes. Covariance matrices were also examined for different PCA combinations. PCA was further compared using scikit-learn, revealing that both

explicitly specifying two principal components and retaining 95% of the information required the same number of components. Two models were trained for the PCA from the library, 1 model was trained for the PCA from scratch and 1 model was trained without applying PCA so that we could compare the performance of the model. It was observed that the respective model accuracies were 0.93, 0.93, 0.96, and 0.93.

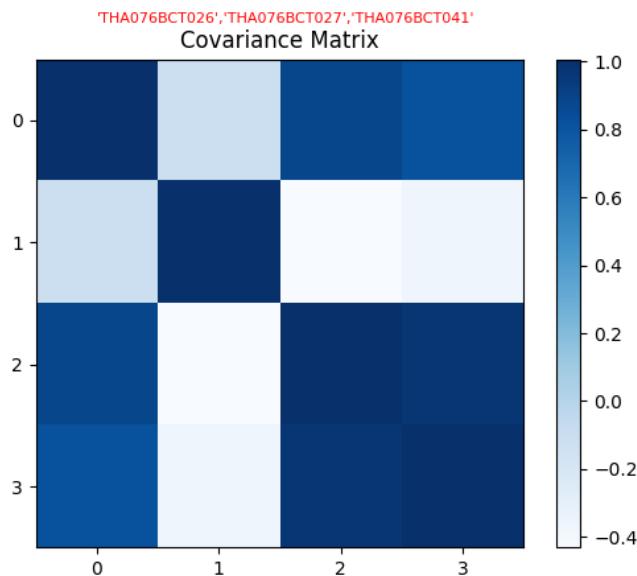


FIGURE 10: Covariance Matrix of Iris Dataset

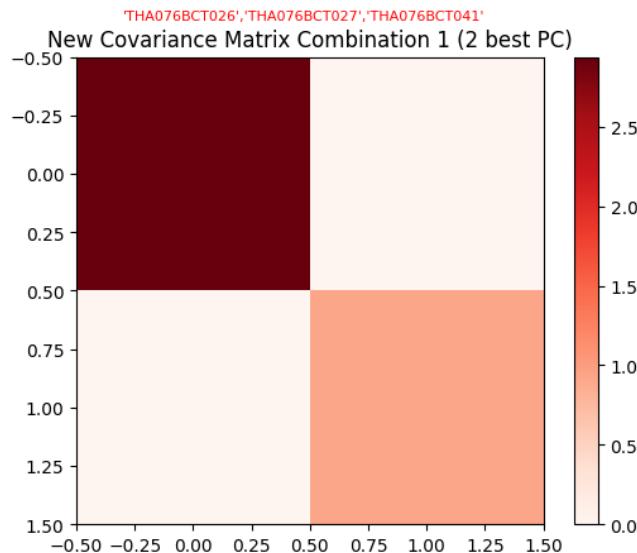


FIGURE 11: Covariance Matrix for two high variance components for Iris Dataset

C. CANCER DATA

The Cancer data was downloaded from Kaggle public dataset that contains data of 570 cancer cells with 30 features to

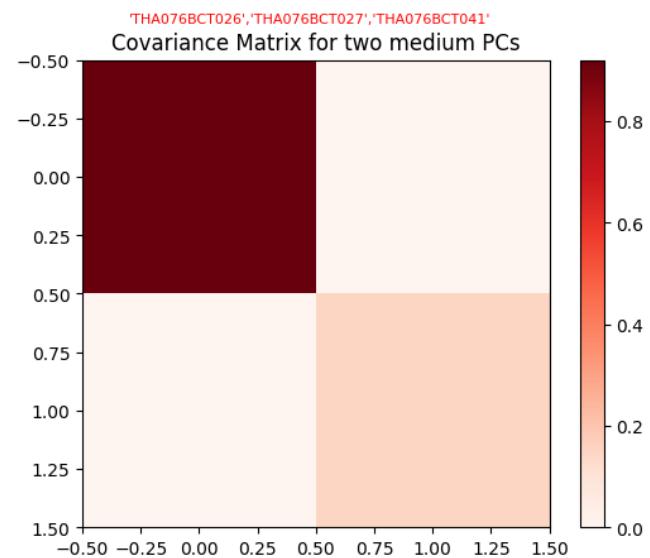


FIGURE 12: Covariance Matrix for two medium variance components for Iris Dataset

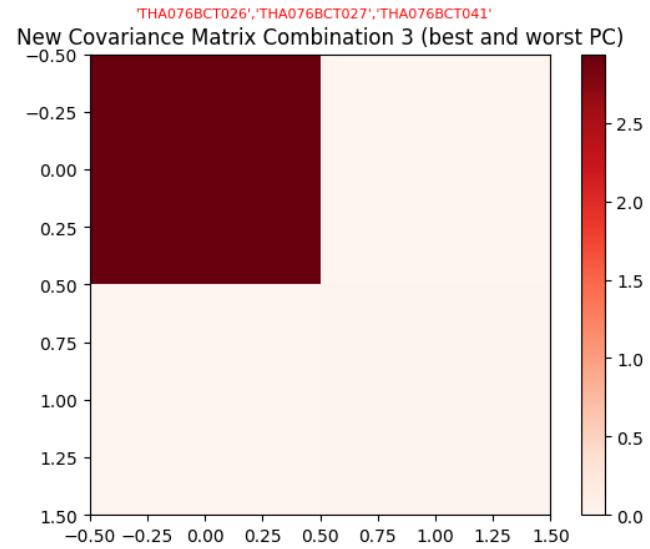


FIGURE 13: Covariance Matrix for one high and one low variance components for Iris Dataset

determine whether the cancer cells data are benign or malignant.

The Principal Component Analysis (PCA) was done on the dataset and Scree Plot was plotted and the plot highlights the valuable insights into the explained variance by each principal component. By examining the scree plot, we can identify the principal components that contribute the most to the overall variance. Therefore, by analyzing the scree plot, we determined the optimal number of principal components to retain for further analysis. This decision is crucial as it balances the trade-off between capturing sufficient variance and minimizing the dimensionality of the dataset.

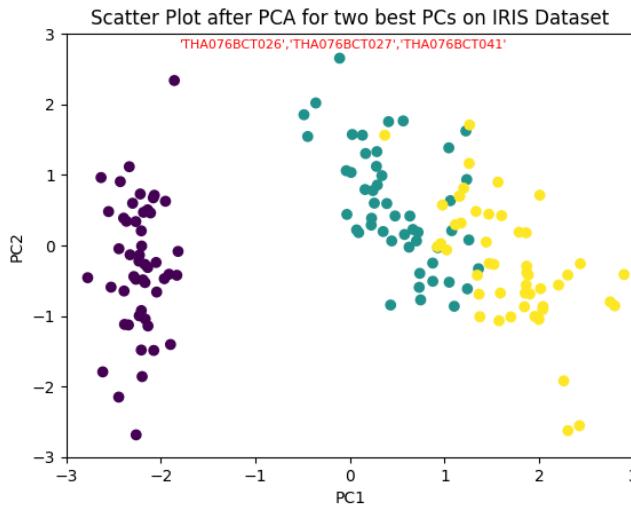


FIGURE 14: Scatter Plot after PCA for two high variance components for Iris Dataset

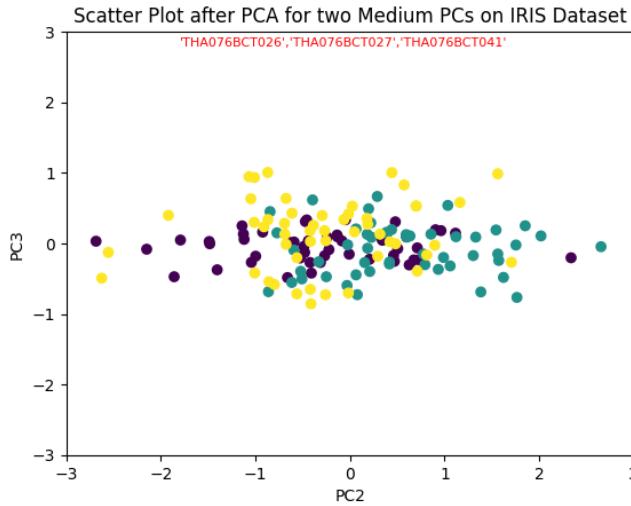


FIGURE 15: Scatter Plot after PCA for two medium variance components for Iris Dataset

The covariance in the processed data of cancer dataset is shown in the Figure 20 and the scree plot can be seen in the Figure 21.

Likewise, PCA was applied from the library to visualize the difference. Again, two instances were created for a similar motivation, that is one by explicitly specifying two components and another by specifying a 0.95 retention value. This time it was observed that to retain 95% of the information, 10 principal components were required.

IV. DISCUSSION AND ANALYSIS

A. RANDOM NUMBERS DATASET

The data was standardized by multiplying it with a 2×2 random uniform matrix, where elements were sampled from a uniform distribution. This transformation resulted in aligning

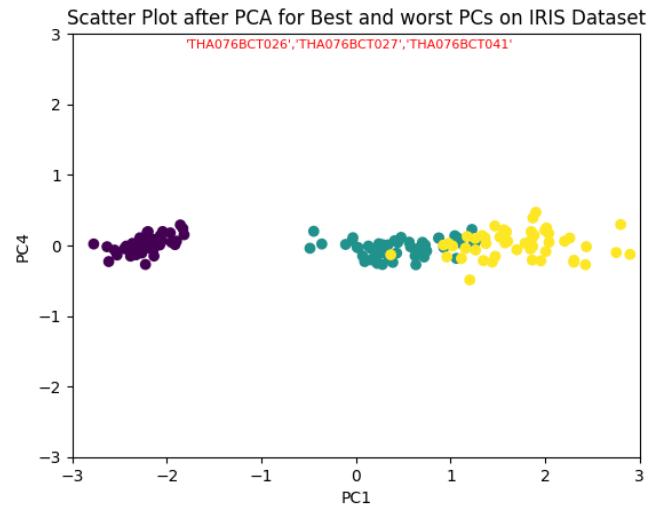


FIGURE 16: Scatter Plot after PCA for one high and one low variance components for Iris Dataset

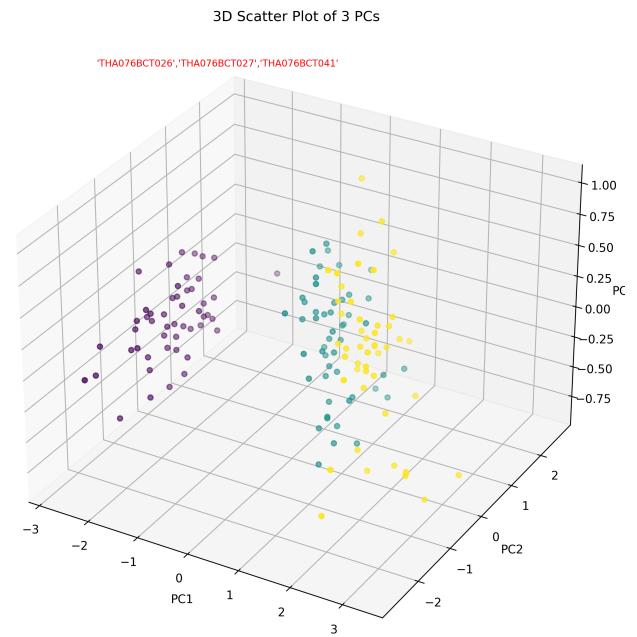


FIGURE 17: 3D Scatter Plot of 3 PCs for Iris Dataset

the random samples along the principal axes of a normal distribution. The purpose of this alignment was to simplify the subsequent PCA analysis by orienting the data in a specific direction. Applying PCA to the dataset of 20 random samples from a Gaussian distribution with 2 features has allowed us to transform the data and visualize it in a lower-dimensional space. The principal components obtained through PCA capture the major sources of variance in the data, providing a concise representation for further analysis.

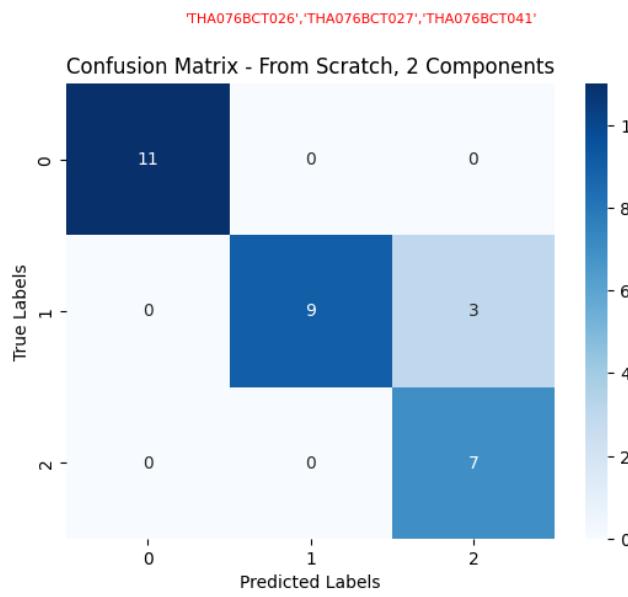


FIGURE 18: Confusion Matrix for the model after applying PCA from Scratch for Iris

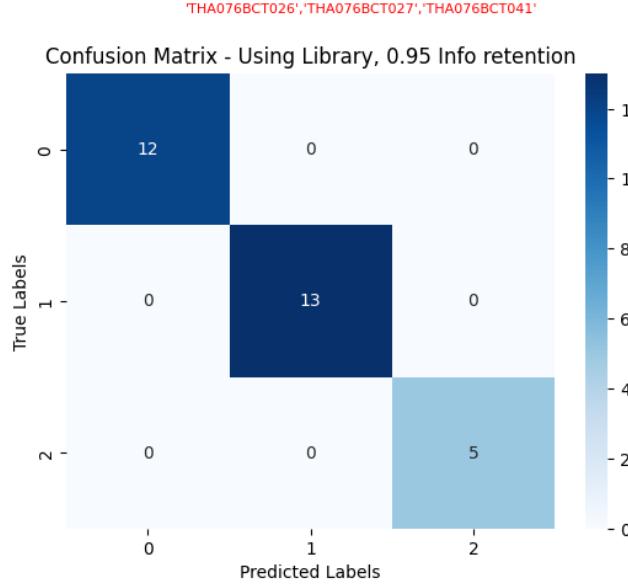


FIGURE 19: Confusion Matrix for the model after applying PCA using Library for Iris

B. IRIS DATASET

During the Principal Component Analysis on the IRIS dataset, it was observed that using principal components beyond the first and second ones in the PCA analysis of the Iris dataset resulted in a mixed-up or differently clustered pattern. This is due to limitations of the third and fourth components (described as medium PCs) in capturing significant variations. The first and second components were more effective in representing the dataset's structure and differentiating between Iris species. Lower-order components show

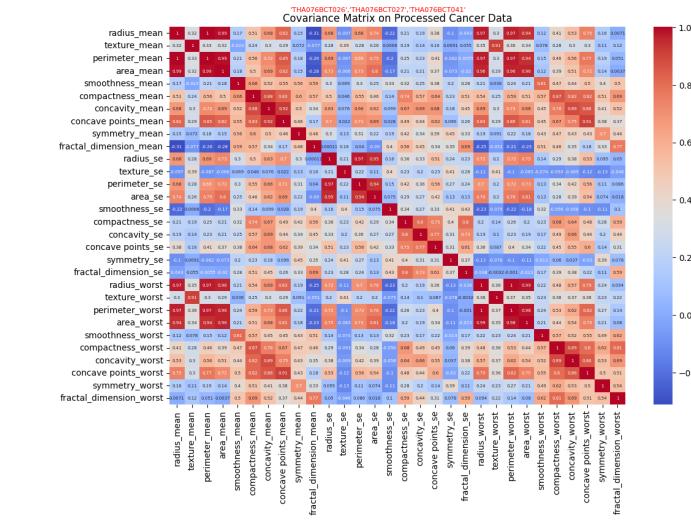


FIGURE 20: Covariance Matrix on Preprocessed Cancer Data

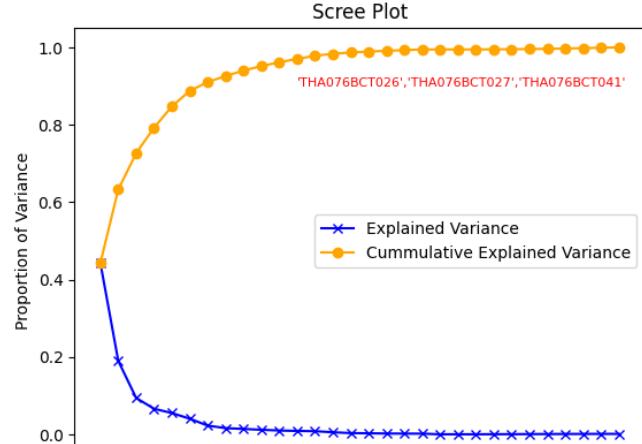


FIGURE 21: Scree Plot of Variance for Cancer Data

clear separations, while higher-order ones may be influenced by noise or irrelevant features. Careful selection of principal components is crucial for extracting meaningful information from the dataset.

C. CANCER DATA

In the PCA analysis of the Cancer dataset, the Scree plot was utilized to examine the variance explained by each principal component. Higher eigenvalues in the plot indicate a larger amount of captured variance. By analysing the Scree plot, we can identify the principal components that contribute the most to the overall variance. Usually, there is a significant drop in eigenvalues after a certain number of components, indicating that the additional components provide diminishing returns in terms of explained variance. Hence it indicated that the first few principal components tend to

be the most informative. Analyzing the Scree plot assists in determining the optimal number of components to retain, striking a balance between capturing sufficient variance and reducing dimensionality. This process is vital for accurately distinguishing between the two types of cancer cells based on their features.

V. CONCLUSION

In conclusion, this report presented a comparative analysis of Principal Component Analysis (PCA) applied to diverse datasets, including Random Numbers, Iris, and Cancer datasets. The main objectives were to explore the capabilities of PCA in addressing real-world challenges associated with high-dimensional data analysis and machine learning. Through our analysis, we successfully applied PCA to each dataset, following a standardized pipeline that involved data pre-processing, covariance matrix computation, eigenvalue-eigenvector decomposition, selection of principal components, and dimensionality reduction. By leveraging the power of PCA, we were able to reduce the dimensionality of the datasets while preserving a significant amount of meaningful information and the underlying structure. The results demonstrated the effectiveness of PCA in achieving the objectives of the lab. Additionally, PCA provided a valuable means of data compression, allowing us to represent the data in a lower-dimensional space without significant loss of information.

Furthermore, PCA greatly enhanced the visualization of the datasets by projecting them onto the principal components. We observed that PCA enabled us to extract the most important features from the datasets, leading to improved data understanding and facilitating more efficient model training, and supporting better decision-making processes. Overall, the objectives of the lab were successfully met through the application of PCA to the datasets. The comparative analysis highlighted the versatility and usefulness of PCA in handling high-dimensional data, deriving meaningful insights, and enhancing data analysis and machine learning tasks. The findings from this study can guide future applications of PCA in various domains and contribute to the advancement of data analysis techniques.

REFERENCES

- [1] Abdi, H., & Williams, L. J. (2010). Principal component analysis. Wiley interdisciplinary reviews: computational statistics, 2(4), 433-459. doi:10.1002/wics.101.
- [2] I. T. Jolliffe, "Principal component analysis and factor analysis," *Principal Component Analysis*, pp. 115-128, 1986.
- [3] K. Pearson, "On lines and planes of closest fit to systems of points in space," *Philosophical Magazine*, vol. 2, no. 11, pp. 559-572, 1901.
- [4] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, no. 4, pp. 433-459, 2010.
- [5] M. Ringnér, "What is principal component analysis?," *Nature Biotechnology*, vol. 26, no. 3, pp. 303-304, 2008.
- [6] scikit-learn developers, "sklearn.decomposition.PCA," scikit-learn documentation, [Online]. Available: <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>.
- [7] J. Brownlee, "How to Use Principal Component Analysis for Feature Selection," Machine Learning Mastery, [Online]. Available:

<https://machinelearningmastery.com/principal-components-analysis-for-feature-selection-in-python/>.

- [8] D. Pattnaik, "Principal Component Analysis (PCA) - Step by Step Tutorial," Towards Data Science, [Online]. Available: <https://towardsdatascience.com/principal-component-analysis-pca-step-by-step-tutorial-81b6387c55c1>.
- [9] A. Dey, "PCA — Machine Learning Mystery," Machine Learning Mystery, [Online]. Available: <https://machinelearningmastery.com/pca-machine-learning-mystery/>.
- [10] A. Géron, "Principal Component Analysis (PCA) Explained," Medium, [Online]. Available: <https://towardsdatascience.com/principal-component-analysis-pca-explained-and-implemented-eaab7cb73b72>.



PRABIN BOHARA is an enthusiastic individual currently pursuing a Bachelor's degree in Computer Engineering, set to graduate in 2024. He has a genuine passion for coding, web technologies, and machine learning. Prabin's primary focus revolves around AI research and exploring the various applications of data science in general. He actively contributes to open-source projects and continuously expands his knowledge through academic endeavors. Prabin actively participates

in webinars, both nationally and internationally, as well as competing in diverse technology competitions. He also shares his expertise by conducting workshops on various technology topics. With a strong drive for excellence, Prabin aims to make valuable contributions in the field of AI while keeping up with the latest advancements.



PRABIN SHARMA POUDEL, an undergraduate senior at IOE's Thapathali Campus, has forever been in pursuit of knowledge and wisdom. His involvement extends beyond technical confines, with involvement in non-technical organizations like AYON. With a vision to become a scholarly figure, he aims to bridge multiple disciplines of knowledge. His current field of interest encompasses AI in conflict resolution, Computer vision, and AI in creative industries Beyond academics, Prabin finds joy in playing chess, composing music, and journaling his thoughts through poetry.



SANTOSH PANDEY is a fourth-year Bachelor's student at IOE, Thapathali Campus, with a passion for coding and technology. Currently working as an intern at NAAMII organization, he is driven by the goal of applying his skills to earn a living through coding. With a keen interest in full-stack development, graphic designing, mobile app development, and computer vision, Santosh is constantly exploring new technologies and honing his skills in these areas. In his free time, he actively participates in coding competitions and enjoys participating in marathons. Santosh's dedication, enthusiasm, and continuous learning make him a promising professional in the field of technology and coding.

CODE

A. RANDOM NUMBERS DATASET

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from numpy import random
4
5 # Generates and visualizes a scatter plot
6 # of two sets of 50 random data points
7 # each. [randn=from normal
8 # distribution.]
9 # x1 and x2 represent the 2 axes
10 x1 = np.random.randn(50)
11 x2 = np.random.randn(50)
12 print(x1)
13 print(x2)
14 plt.scatter(x1,x2)
15 plt.xlabel("X1")
16 plt.ylabel("X2")
17 plt.text(-2, 2.8, "'THA076BCT026',",
18 'THA076BCT027', 'THA076BCT041"',
19 fontsize=8,color='red')
20 plt.title('Scatter Plot for Random Values
21 ')
22 plt.ylim(-3,3)
23 plt.xlim(-3,3)
24 plt.savefig('Scatter Plot for Random
25 Values', bbox_inches='tight')
26
27 # Mean and standard deviation
28 mean1 = np.mean(x1)
29 mean2 = np.mean(x2)
30 sd1 = np.std(x1)
31 sd2 = np.std(x2)
32 print(mean1,sd1)
33 print(mean2,sd2)
34
35 # 2x2 random matrix with elements sampled
36 # from a uniform distribution (range
37 # [0, 1])..
38 random_matrix = np.random.rand(2,2)
39 print(random_matrix)
40
41 # Converts the lists x1 and x2 into NumPy
42 # arrays, then creates a data matrix
43 # by column stacking the arrays.
44 x1_array = np.array(x1)
45 x2_array = np.array(x2)
46 data_matrix = np.column_stack((x1_array,
47 x2_array))
48 print(data_matrix.shape)
49 print(data_matrix)
50
51 # **STEP 1: Reshape the data by
52 # multiplication with a random matrix
53 # ** To visualize the linear
54 # transformation, click on the link :
55 # https://www.geogebra.org/m/YCza8TAH?
56 # fbclid=IwAR3LK2VOQoobx-qzH-
57 # LLhlyzVEhi8Qsb8Of_VlzeiqWb3FcjYibVPbL0NGo
58
59 transformed_data_by_random_matrix =
60 data_matrix @ random_matrix
61 print(transformed_data_by_random_matrix.
62 shape)

```

```

44 # first column [:,0] as x-axis value.
45 # similarly 2nd column as y-axis
46 plt.scatter(
47     transformed_data_by_random_matrix
48     [:,0],
49     transformed_data_by_random_matrix
50     [:,1])
51 plt.xlabel("X1-transformed")
52 plt.ylabel("X2-transformed")
53 plt.ylim(-3,3)
54 plt.xlim(-3,3)
55 plt.text(-1.5, 2.8, "'THA076BCT026',",
56 'THA076BCT027', 'THA076BCT041"',
57 fontsize=8,color='red')
58 plt.title('Scatter Plot for transformed
59 Random Values')
60 plt.savefig('transformed Random Values',
61 bbox_inches='tight')
62
63
64 # **STEP 2: Calculate covariance matrix**
65 # The covariance matrix provides
66 # insights into the relationships
67 # between different dimensions of the
68 # transformed_data_by_random_matrix.
69 # See if elements in main diagonal are
70 # max and other diagonal are min. If
71 # not, we have to proceed further .
72 covariance_matrix = np.cov(
73     transformed_data_by_random_matrix.T)
74 covariance_matrix
75 #same as above
76 covariance_matrix = np.cov(np.transpose(
77     transformed_data_by_random_matrix))
78 print(covariance_matrix)
79
80
81 # **STEP 3: Eigen values and vector
82 # calculation**
83 # The eigenvalues represent the variance
84 # explained by each principal component
85 # (magnitude of new feature space)
86 # The eigenvectors define the directions
87 # of the principal components. ((
88 # direction of new feature space)
89 eigen_values, eigen_vectors = np.linalg.
90 eig(covariance_matrix)
91 print(eigen_values)
92 print(eigen_vectors)
93
94 # Just To check if the eigenvalues really
95 # represent the variance, lets find
96 # out the proportion of variance for
97 # each components
98 proportion1 = eigen_values[0] / np.sum(
99 eigen_values)
100 print(proportion1*100)
101 proportion2 = eigen_values[1] / np.sum(
102 eigen_values)
103 print(proportion2*100)
104
105
106 # **STEP 4: Transform the data with the
107 # selected eigen vector.**
108 #first write the eigenvectors row wise (
109 # do the transpose, from column wise to
110 # row wise)

```

```

80 arranged_eigen_vector = np.transpose(
81     eigen_vectors)
82 print(arranged_eigen_vector.shape)
83 print(arranged_eigen_vector)
84
85 #also transpose the (50,2) data matrix to
86 # (2,50)
87 transposed_data_matrix =
88     transformed_data_by_random_matrix.T
89 print(transposed_data_matrix.shape)
90
91 #Project data into first principal
92 # component (by matrix multiplication)
93 new_data = np.dot(arranged_eigen_vector,
94     transposed_data_matrix)
95 print(new_data.shape)
96 #print(new_data)
97
98 #transpose back to (50,2)
99 trans_new_data=new_data.T
100 print(trans_new_data.shape)
101
102 # again see if the elements in main
103 # diagonal are max
104 covariance_of_new_data = np.cov(new_data)
105 print(covariance_of_new_data)
106
107 #lets select the best eigen vectors as
108 # specified in the step 4
109 best_eigen_vector = np.transpose(
110     eigen_vectors[:,0])
111 print(best_eigen_vector.shape)
112 print(best_eigen_vector)
113 new_new_data = np.dot(best_eigen_vector,
114     transposed_data_matrix)
115 print(new_new_data.shape)
116 print(new_new_data)
117
118 # **STEP 6: plot the new data**
119 plt.scatter(new_data[0,:], new_data[1,:])
120 plt.ylim(-3,3)
121 zeros = np.zeros(50)
122 plt.scatter(new_data[1,:],zeros)
123 plt.xlabel("X1_after_PCA")
124 plt.ylabel("X2_after_PCA")
125 plt.title('Scatter Plot After PCA')
126 plt.text(-0.8, 0.9, "'THA076BCT026','
127     THA076BCT027','THA076BCT041'",
128     fontsize=8,color='red')
129 plt.ylim(-1,1)
130 plt.xlim(-1,1)
131 plt.grid()
132 plt.savefig('scatter plot after PCA',
133     bbox_inches='tight')
134
135 print(new_data.shape)
136 print(new_data)

```

B. IRIS DATASET

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 from numpy import random
5 import sklearn

```

```

6
7 # **STEP 1: Load the iris flower dataset
8
9 from sklearn.datasets import load_iris
10 iris = load_iris()
11
12 X = iris.data      # X is 2d array, where
13 # row=samples, column=features
14 y = iris.target    # y is 1d array,
15 # representing class label for all
16 # samples
17 print(X.shape)
18 print(y.shape)
19 print(iris.feature_names)
20 print("Target names:", iris.target_names)
21 print(X)
22
23 # **STEP 2: Create a dataframe and see
24 # the data statistics**
25
26 # assign meaningful columns(feature_names)
# to feature values for 150 samples(
# iris.data) in the dataframe.
27 df = pd.DataFrame(data=iris.data,columns=
28     iris.feature_names)
29 # add a new target column.
30 df['target']=iris.target
31 df.head(2)      # view 1st 2 data, both are
32 # setosa(numbering starts from 0)
33
34 #view 49th to 52nd dataset to see 2
# classes(setosa and vesicolor)
35 subset = df.iloc[48:52]
36 print(subset)
37
38 # statistics computed for each column
39 summary_stats = df.describe()
40 print(summary_stats)
41
42 # data visualization, box plot for each
# feature
43 df.plot(kind='box', figsize=(10, 6))
44 plt.title('Box Plot of Iris Dataset')
45 plt.text(1.5, 8, "'THA076BCT026','
46     THA076BCT027','THA076BCT041'",
47     fontsize=8,color='red')
48
49 plt.savefig('Box Plot of Iris Dataset',
50     bbox_inches='tight')
51 plt.show()
52
53 import seaborn as sns
54 # Set the size of the figure
55 plt.figure(figsize=(12, 7))
56 df['target'] = df['target'].replace({0: 'Setosa', 1: 'Versicolor', 2: 'Virginica'})
57
58 # Create a pair plot with hue based on
# the 'target' column in the dataframe
59 g=sns.pairplot(df, hue='target',plot_kws
# ={'legend': False})
60 g._legend.set_title('')
61 # Set the title of the plot
62 plt.suptitle('Attribute Correlation Pair
# Plot', fontsize=16, y=0.98)

```

```

54 # Set the style and color of the text box
55 textbox_props = dict(boxstyle='round',
56     facecolor='white', edgecolor='gray',
57     alpha=0.8)
58
59 # Add a text box with the highlighted
60 # feature names
61 plt.text(0.1, 8, "THA076BCT026\
62     nTHA076BCT027\nnTHA076BCT041",
63     fontsize=10, color='red', ha='center'
64     , bbox=textbox_props)
65 legend_elements = [plt.Line2D([0], [0],
66     marker='o', color='w', label='Setosa'
67     , markerfacecolor='royalblue',
68     markersize=8),
69     plt.Line2D([0], [0],
70     marker='o', color=
71     'w', label='
72     Versicolor',
73     markerfacecolor='g
74     ', markersize=8),
75     plt.Line2D([0], [0],
76     marker='o', color=
77     'w', label='
78     Virginica',
79     markerfacecolor='orange',
80     markersize=8)]
81
82 # Add the legend to the plot
83 plt.legend(handles=legend_elements, title
84     ='target')
85
86 # Adjust the spacing between subplots
87 plt.subplots_adjust(top=0.9)
88 # Set custom tick labels for the legend
89 # Save the plot to a file
90 plt.savefig('./plots/pair_plot.png',
91     bbox_inches='tight')
92
93 # Display the plot
94 plt.show()
95
96 # Visualization of data points based on
97 # the target class (purple=setosa,green
98 # =Versicolour,yellow=Virginica )
99 df['target']=iris.target
100 plt.figure(figsize=(10, 6))
101 plt.scatter(df['sepal length (cm)'], df['
102     sepal width (cm)'], c=df['target'],
103     label='Sepal')
104 plt.scatter(df['petal length (cm)'], df['
105     petal width (cm)'], c=df['target'],
106     marker='x', label='Petal')
107 plt.xlabel('Length (cm)')
108 plt.ylabel('Width (cm)')
109 plt.title('Scatter Plot of Sepal and
110     Petal Features')
111 plt.colorbar(label='Species')
112 plt.text(1, 4.3, "'THA076BCT026',"
113     "'THA076BCT027', 'THA076BCT041'",
114     fontsize=8,color='red')
115 plt.legend()
116 plt.savefig('./plots/Sepal and Petal
117     Features', bbox_inches='tight')
118 plt.show()

```

```

89 # **STEP 3: Data preprocessing**
90
91 # Handling missing values:
92 # df.isna() returns a boolean, where True
93 # means missing values(NaN) and .sum()
94 # sums the values along each column
95 missing_values = df.isna().sum()
96 print("Columns with missing values:")
97 print(missing_values[missing_values > 0])
98
99 # Standardizing to eliminate the chance
100 # of PCA being influenced by magnitude
101 # of some features' values.
102 from sklearn.preprocessing import
103     StandardScaler
104 scaler = StandardScaler()
105 X = df.iloc[:,1:].values
106 df_scaled = pd.DataFrame(scaler.
107     fit_transform(df), columns=df.columns
108     )
109
110 # IMPORTANT STEP!
111 # Create new input and output after
112 # preprocessing X and y.
113 X_train = df_scaled.drop('target', axis
114     =1).values
115 y_train = df['target']
116 print("shape of new input X",X_train.
117     shape)
118 print("shape of new output Y",y_train.
119     shape)
120 print(X_train)
121 print(y_train)
122
123 df_plot=df_scaled.copy()
124 df_plot['target'] = df['target'].replace
125     ({0: 'Setosa', 1: 'Versicolor', 2: ' '
126     'Virginica'})
127 plt.figure(figsize=(12, 8))
128
129 # Create a pair plot with hue based on
130 # the 'target' column in the dataframe
131 g=sns.pairplot(df_plot, hue='target',
132     plot_kws={'legend': False})
133 g._legend.set_title('')
134
135 plt.suptitle('Attribute Correlation Pair
136     Plot of Transformed Data', fontsize
137     =16, y=0.98)
138
139 # Set the style and color of the text box
140 textbox_props = dict(boxstyle='round',
141     facecolor='white', edgecolor='gray',
142     alpha=0.8)
143
144 # Add a text box with the highlighted
145 # feature names
146 plt.text(1.5, 6, "THA076BCT026\
147     nTHA076BCT027\nnTHA076BCT041",
148     fontsize=10, color='red', ha='center'
149     , bbox=textbox_props)
150
151 # Adjust the spacing between subplots
152 plt.subplots_adjust(top=0.9)
153
154 legend_elements = [plt.Line2D([0], [0],
155     marker='o', color='w', label='Setosa'
156     , markerfacecolor='royalblue',
157     markersize=8),
158     plt.Line2D([0], [0],
159     marker='o', color=
160     'w', label='
161     Versicolor',
162     markerfacecolor='g
163     ', markersize=8),
164     plt.Line2D([0], [0],
165     marker='o', color=
166     'w', label='
167     Virginica',
168     markerfacecolor='orange',
169     markersize=8)]
170
171 plt.legend(handles=legend_elements, title
172     ='target')
173
174 plt.show()

```

```

131     markersize=8),
132         plt.Line2D([0], [0],
133                   marker='o', color=
134                   'w', label='
135                   Versicolor',
136                   markerfacecolor='g
137                   ', markersize=8),
138         plt.Line2D([0], [0],
139                   marker='o', color=
140                   'w', label='
141                   Virginica',
142                   markerfacecolor='
143                   orange',
144                   markersize=8])

145 # Add the legend to the plot
146 plt.legend(handles=legend_elements, title
147             ='target')
148 # Save the plot to a file
149 plt.savefig('./plots/
150             pair_plot_transformed.png',
151             bbox_inches='tight')

152 # Display the plot
153 plt.show()

154 # lets see the new plot after
155     standardization
156 plt.figure(figsize=(10, 6))
157 plt.scatter(df_scaled['sepal length (cm)'
158                 ], df_scaled['sepal width (cm)'], c=
159                 df_scaled['target'], label='Sepal')
160 plt.scatter(df_scaled['petal length (cm)'
161                 ], df_scaled['petal width (cm)'], c=
162                 df_scaled['target'], marker='x',
163                 label='Petal')
164 plt.xlabel('Length (cm)')
165 plt.ylabel('Width (cm)')
166 plt.title('Scatter Plot of Sepal and
167             Petal Features')
168 plt.colorbar(label='Species')
169 plt.ylim(-3,4)
170 plt.xlim(-3,3)
171 plt.text(-1, 3.5, "'THA076BCT026','
172             THA076BCT027','THA076BCT041'",
173             fontsize=8,color='red')
174 plt.legend()
175 plt.savefig('./plots/processed Sepal and
176             Petal Feature', bbox_inches='tight')
177 plt.show()

178 # **STEP 4: Covariance matrix calculation
179 .**
180
181 # X_train is already an array so no need
182     to do data_matrix = np.array(X)
183 data_matrix = X_train
184 data_matrix.shape

185 # covariance matrix will be 4*4 since x
186     =(150,4)
187 covariance_matrix = np.cov(X_train,
188                           rowvar=False)           # row of
189                           X_train != variable(so column of
190                           X_train means variable, and rows
191                           means observation)
192 print("Covariance matrix shape:",
```

```

193 covariance_matrix.shape)
194 print(covariance_matrix)

195 # visualize covariance matrix
196 plt.imshow(covariance_matrix, cmap='Blues
197             ', interpolation='nearest')
198 plt.colorbar()
199 plt.title('Covariance Matrix')
200 plt.text(0, -0.8, "'THA076BCT026','
201             THA076BCT027','THA076BCT041'",
202             fontsize=8,color='red')
203 plt.xticks(np.arange(4), np.arange(4))
204             #tick labels on axes changed from 0
205             to 4
206 plt.yticks(np.arange(4), np.arange(4))
207 plt.savefig('Iris Covariance Matrix',
208             bbox_inches='tight')
209 plt.show()

210 # **STEP 5: Eigen values and eigen vector
211     calculation**
212
213 eigenvalues, eigenvectors = np.linalg.eig
214     (covariance_matrix)
215 #descending sort
216 sorted_indices = np.argsort(eigenvalues)
217             [::-1]
218 eigenvalues = eigenvalues[sorted_indices]
219 eigenvectors = eigenvectors[:,sorted_
220             _indices]
221 print("eigen values are",eigenvalues)
222 print("eigen vectors are",eigenvectors)
223 print(eigenvalues.shape)
224 print(eigenvectors.shape)

225 # **STEP 6: See the variance explained by
226     each eigen values.**
227
228 # Proportion of variance for each eigen
229     values
230 pov_list=[]
231 for i in range (0,4):
232     pov=eigenvalues[i]/sum(eigenvalues)
233     pov_list.append(pov)
234     print("Proportion of variance for",
235             eigenvalues[i],"eigen value is",
236             pov)

237 # Plot the scree plot
238 plt.plot(range(1, len(pov_list) + 1),
239             pov_list, marker='o',color='blue')
240 varlegend= "'THA076BCT026','THA076BCT027
241             ','THA076BCT041'"
242 plt.text(2.5,0.7,varlegend,color='red')
243             #plt.text(x_position, y_position,
244             varlegend)
245 plt.xlabel('Principal Component')
246 plt.ylabel('Proportion of Variance')
247 plt.title('Scree Plot')
248 plt.show()

249 # **STEP 7: Select the eigen vectors as
250     needed to compute the final output**
251
252 print("Shape of x train is",X_train.shape
253             )
```

```

212 transposed_X_train = np.transpose(X_train)
213 transposed_X_train.shape
214 eigenvectors_transposed=np.transpose(
215     eigenvectors)
216 final_data = {}      # Dictionary to
217     # store the final projected data
218 # Define the combinations of selected
219     # components
220 selected_components = [[0], [0, 1],
221     [0,2], [0, 3], [1], [1,2], [1,3], [0,
222     1, 2, 3]]
223 count=len(selected_components)
224
225 # Iterate over the selected components
226 for i in range (0,count):
227     selected_eigenvectors =
228         eigenvectors_transposed[
229             selected_components[i]]
230             #select the
231             # combination of components from
232             # above list
233     final_data[i] = np.transpose(
234         selected_eigenvectors @
235         transposed_X_train) #projection
236     print(final_data[i].shape)
237
238 eigenvectors_transposed[
239     selected_components[2]]
240
241 # **STEP 8: Computing and visualizing
242     # covariance for the new data.**
243
244 new_covariance_matrices_dictionary = {}
245 # Compute covariance matrix for each
246     # projected data
247 for i in range(0, count):
248     print(" * Covariance matrix for
249         # combination no", i)
250     new_covariance_matrix = np.cov(
251         final_data[i], rowvar=False)
252     new_covariance_matrices_dictionary[i]
253         = new_covariance_matrix
254     # Print the shape of the covariance
255         # matrix
256     print(" --->Covariance matrix shape:
257         ", new_covariance_matrix.shape)
258
259 #Lets access the first combination(
260     # combination 0), which is (eigen[0])
261             # (best component)
262 new_covariance_matrices_dictionary[0]
263
264 #Lets access combination 1, which is
265     # eigen[0] and eigen[1]           (2
266         # best components)
267 new_covariance_matrices_dictionary[1]
268
269 #Lets access combination 3, which is
270     # eigen[0] and eigen[3]           (best
271         # and worst component)
272 new_covariance_matrices_dictionary[3]
273
274 #Lets access combination 5, which is
275     # eigen[1] and eigen[2]           (2

```

```

276 medium medium best components)
277 new_covariance_matrices_dictionary[5]
278
279 #Lets access combination 7, which is
280     # eigen[0],eigen[1],eigen[2] and eigen
281         # [3]           (all components combined)
282 new_covariance_matrices_dictionary[7]
283
284 # Heatmap plot, darker shade indicate
285     # higher values. (combination starts
286         # from 0 to 7)
287 plt.imshow(
288     new_covariance_matrices_dictionary
289     [1], cmap='Reds', interpolation='
290         nearest') #for smoother
291 interpolation of color on pixel
292 values, instead of 'nearest', we can
293 use 'bilinear','bicubic','spline16',
294 'spline36', 'hanning', 'hamming', '
295 hermite',etc..
296 plt.colorbar()
297 plt.title("New Covariance Matrix
298     Combination 1 (2 best PC)")
299 plt.text(-0.3, -0.66, "'THA076BCT026'",
300         'THA076BCT027','THA076BCT041"',
301         fontsize=8,color='red')
302 plt.savefig('Covariance Matrix for 2 best
303     PC', bbox_inches='tight')
304 plt.show()
305
306 plt.imshow(
307     new_covariance_matrices_dictionary
308     [3], cmap='Reds', interpolation='
309         nearest')
310 plt.colorbar()
311 plt.title("New Covariance Matrix
312     Combination 3 (best and worst PC)")
313 plt.text(-0.3, -0.66, "'THA076BCT026'",
314         'THA076BCT027','THA076BCT041"',
315         fontsize=8,color='red')
316 plt.savefig('Covariance Matrix for best
317     and worst PC', bbox_inches='tight')
318 plt.show()
319
320 plt.imshow(
321     new_covariance_matrices_dictionary
322     [5], cmap='Reds', interpolation='
323         nearest')
324 plt.colorbar()
325 plt.title("Covariance Matrix for two
326         medium PCs")
327 plt.text(-0.3, -0.66, "'THA076BCT026'",
328         'THA076BCT027','THA076BCT041"',
329         fontsize=8,color='red')
330 plt.savefig('Covariance Matrix for two
331         medium PCs', bbox_inches='tight')
332 plt.show()
333
334 plt.imshow(
335     new_covariance_matrices_dictionary
336     [7], cmap='Reds', interpolation='
337         nearest')
338 plt.colorbar()
339 plt.title("New Covariance Matrix
340     Combination 7 (All Components)")
341 plt.show()
342
343

```

```

283 # **STEP 9: Obtain final data and
284     visualize**
285
286 # Two best PCs
287 final_data[1].shape
#final_data[1]
288
289 plt.ylim(-3,3)
290 plt.xlim(-3,3)
291 # plt.title("2 best PCs")
292 plt.xlabel('PC1')
293 plt.ylabel('PC2')
294 plt.title("Scatter Plot after PCA for two
    best PCs on IRIS Dataset")
295 plt.text(-1.8, 2.8, "'THA076BCT026','
    THA076BCT027','THA076BCT041'",
    fontsize=8,color='red')
296 plt.scatter(final_data[1][:,0],final_data
    [1][:,1],c=y)
297 plt.savefig('Scatter Plot after PCA for
    two best PCs', bbox_inches='tight')
298
299 #Best and worst PCs
300 final_data[3].shape
#final_data[3]
301
302 plt.ylim(-3,3)
303 plt.xlim(-3,3)
304 plt.xlabel('PC1')
305 plt.ylabel('PC4')
306 plt.title("Scatter Plot after PCA for
    Best and worst PCs on IRIS Dataset")
307 plt.text(-1.8, 2.8, "'THA076BCT026','
    THA076BCT027','THA076BCT041'",
    fontsize=8,color='red')
308 plt.scatter(final_data[3][:,0],final_data
    [3][:,1],c=y)
309 plt.savefig('Scatter Plot after PCA for
    Best and worst PCs', bbox_inches='
    tight')
310
311 #Best and mediumly worse PCs
312 final_data[5].shape
#final_data[5]
313
314 plt.ylim(-3,3)
315 plt.xlim(-3,3)
316 plt.title("Scatter Plot after PCA for two
    Medium PCs on IRIS Dataset")
317 plt.text(-1.8, 2.8, "'THA076BCT026','
    THA076BCT027','THA076BCT041'",
    fontsize=8,color='red')
318 plt.xlabel('PC2')
319 plt.ylabel('PC3')
320 plt.scatter(final_data[5][:,0],final_data
    [5][:,1],c=y)
321 plt.savefig('Scatter Plot after PCA for
    two Medium PCs', bbox_inches='tight')
322
323 #All components
324 final_data[7].shape
#final_data[7]
325
326 plt.ylim(-3,3)
327 plt.xlim(-3,3)
328 plt.title("All PCs")
329 plt.xlabel('PC1')
330
331
332
333 plt.ylabel('PC2')
334 plt.scatter(final_data[7][:,0],final_data
    [7][:,1],c=y)
335
336 # Create a 3D scatter plot
337 from mpl_toolkits.mplot3d import Axes3D
338
339 fig = plt.figure()
340 fig = plt.figure(figsize=(8, 8))
341 ax = fig.add_subplot(111, projection='3d')
342 ax.scatter(final_data[7][:,0],final_data
    [7][:,1],final_data[7][:,2],c=y,
    marker='o')
343 ax.set_xlabel('PC1')
344 ax.set_ylabel('PC2')
345 ax.set_zlabel('PC3')
346 ax.set_title('3D Scatter Plot of 3 PCs')
347 ax.text(-6, 2.8, 1, "'THA076BCT026','
    THA076BCT027','THA076BCT041'",
    fontsize=8,color='red')
348 plt.tight_layout()
349 plt.savefig('3D Scatter Plot of 3 PCs',
    bbox_inches='tight',dpi=300)
350 plt.show()
351
352 # **STEP 10: PCA with library.**
353
354 from sklearn.decomposition import PCA
355
356 # Create instances of PCA with the
    desired number of components. (pca1 =
        2components. pca2=so as to retain
        95% info)
357 pca1 = PCA(n_components=2)
358 pca2 = PCA(0.95)
359
360 # Fit the PCA model to the data and
    transform the data
361 reduced_data1 = pca1.fit_transform(X)
362 reduced_data2 = pca2.fit_transform(X)
363
364 # Print the shape of the reduced data
365 print("Shape of reduced data:",
    reduced_data1.shape)
366 print("reduced through first",
    reduced_data1)
367
368 print("Shape of reduced data:",
    reduced_data2.shape)
369 print("reduced through second",
    reduced_data2)
370
371
372 # **STEP 11: Lets train 4 models and
    visualize the results.** 2 models are
    implemented after applying PCA using
    library. The third model is
    implemented after applying PCA using
    Math(using all the earlier above
    steps). The fourth model is trained
    without PCA.
373
374 # Create 4 sets of test-train for 4
    models.
375 from sklearn.model_selection import
    train_test_split

```

```

376 x_train1,x_test1,y_train1,y_test1 =
377     train_test_split(reduced_data1,y,
378         test_size=0.2) # using library, 2
379         components
380
381 x_train2,x_test2,y_train2,y_test2 =
382     train_test_split(reduced_data2,y,
383         test_size=0.2) # using library,
384         0.95 info
385
386 x_train3,x_test3,y_train3,y_test3 =
387     train_test_split(final_data[1],y,
388         test_size=0.2) # from scratch, 2
389         components
390
391 x_train4,x_test4,y_train4,y_test4 =
392     train_test_split(X_train,y,test_size
393         =0.2) # no PCA
394
395 #training of the model
396 # model 1: using library, 2 components
397 from sklearn.svm import SVC
398 svm1 = SVC()
399 svm1.fit(x_train1, y_train1)
400
401 #library, 0.95 info
402 svm2=SVC()
403 svm2.fit(x_train2, y_train2)
404
405 #from scratch, 2 components
406 svm3=SVC()
407 svm3.fit(x_train3, y_train3)
408
409 #No PCA
410 svm4=SVC()
411 svm4.fit(x_train4, y_train4)
412
413 # Accuracy check for 4 models
414 from sklearn.metrics import
415     accuracy_score
416
417 y_pred1 = svm1.predict(x_test1)
418 accuracy1 = accuracy_score(y_test1,
419     y_pred1)
420 print("Accuracy1:", accuracy1)
421
422 y_pred2 = svm2.predict(x_test2)
423 accuracy2 = accuracy_score(y_test2,
424     y_pred2)
425 print("Accuracy2:", accuracy2)
426
427 y_pred3 = svm3.predict(x_test3)
428 accuracy3 = accuracy_score(y_test3,
429     y_pred3)
430 print("Accuracy3:", accuracy3)
431
432 y_pred4 = svm4.predict(x_test4)
433 accuracy4 = accuracy_score(y_test4,
434     y_pred4)
435 print("Accuracy4:", accuracy4)
436
437 #predicton from 4 models
438 predictions1 = svm1.predict(x_test1)
439 predictions2 = svm2.predict(x_test2)
440 predictions3 = svm3.predict(x_test3)
441 predictions4 = svm4.predict(x_test4)
442 predictions1.shape
443
444 # Lets define a function to plot
445     confusion matrix

```

```

426 from sklearn.metrics import
427     confusion_matrix
428 import seaborn as sns
429 def plot_confusion_matrix(y_true, y_pred,
430     model_name, text):
431     cm = confusion_matrix(y_true, y_pred)
432     sns.heatmap(cm, annot=True, cmap=""
433         Blues", fmt="d")
434     plt.title(f"Confusion Matrix - {"
435         model_name}")
436     plt.xlabel("Predicted Labels")
437     plt.ylabel("True Labels")
438     plt.text(0.2, -0.5, "THA076BCT026",
439         'THA076BCT027','THA076BCT041",
440         fontsize=8,color='red')
441     plt.savefig(f'./plots/{model_name}.
442         png', bbox_inches='tight',dpi
443         =300)
444     plt.show()
445
446 # Plot confusion matrix for model 1
447 plot_confusion_matrix(y_test1,
448     predictions1, "UsingLibrary, 2
449     Components")
450
451 # for model 2
452 plot_confusion_matrix(y_test2,
453     predictions2, "Using Library, 0.95
454     Info retention")
455
456 #for model 3
457 plot_confusion_matrix(y_test3,
458     predictions3, "From Scratch, 2
459     Components")
460
461 #for model4
462 plot_confusion_matrix(y_test4,
463     predictions4, "No PCA")

```

C. CANCER DATASET

```

1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import sklearn
5 import seaborn as sns
6
7
8 # **STEP 1 : Load the dataset into a
8     dataframe and analyze the data**
9 #specify path of csv file relative to "/
9     content/drive" directory.
10 file_path = './dataset/Cancer_Data.csv'
11 df = pd.read_csv(file_path)
12 df.head(10)
13 df.shape
14
15 # Print the column names
16 print(df.columns)
17
18 # statistics computed for each column in
18     the df
19 summary_stats = df.describe()
20 print(summary_stats)
21
22 # see the unique values in diagnosis
22     column

```

```

23 unique_diagnosis = df['diagnosis'].unique()
24 print(unique_diagnosis)
25
26 # Looks like diagnosis is our target
27 # column. Lets assign X and y
28 # accordingly
29 # Separate the data and target variables
30 X = df.drop(columns=['id', 'diagnosis', 'Unnamed: 32'])
31 y = df['diagnosis']
32
33 print("Data:")
34 print(X.shape)
35 print(X.head())
36 print()
37 print("Target:")
38 print(y.shape)
39 print(y.head())
40 X.columns
41
42 # Create a heatmap of the attribute
43 # correlations
44 # High correlation (+ve or -ve) by bright
45 # colors, low correlation by dull
46 # colors.
47 plt.figure(figsize=(12, 8))
48 sns.heatmap(X.corr(), cmap='coolwarm',
49             annot=True, annot_kws={'size': 5})
50 plt.title('Attribute Correlation Heatmap')
51 plt.text(8, -1.2, "'THA076BCT026',",
52         "'THA076BCT027', 'THA076BCT041'",,
53         fontsize=8,color='red')
54 plt.savefig('./plots/Corelation Matrix',
55             bbox_inches='tight')
56 plt.show()
57
58 # **STEP 2: Preprocessing**
59 # Perform label encoding on the 'diagnosis' target variable
60 from sklearn.preprocessing import
61     LabelEncoder
62 label_encoder = LabelEncoder()
63 y_encoded = label_encoder.fit_transform(y)
64 print(y_encoded)
65
66 # change the dataframe's diagnosis column
67 df['diagnosis'] = y_encoded
68 df.head()
69
70 # Standardizing everything to have a
71 # common ground so that magnitude of
72 # some feature's value doesn't affect
73 # the PCA.
74 from sklearn.preprocessing import
75     StandardScaler
76 scaler = StandardScaler()
77 scaler.fit(X)
78 # Scaling of the whole dataframe.
79 df_scaled = pd.DataFrame(scaler.
80     fit_transform(X), columns=X.columns)
81
82 # Define new X and y
83 X_train = df_scaled

```

```

70 y_train = df['diagnosis']
71 print("shape of new input X",X_train.
72     shape)
73 print("shape of new output Y",y_train.
74     shape)
75 df_scaled.head()
76
77 # Assign different colors to the data
78 # points based on the 'diagnosis'
79 # column
80 # Used to visually distinguish 2
81 # different categories in that column
82 colors = ['blue' if t == 0 else 'red' for
83     t in df['diagnosis'].map({0: 'B', 1:
84         'M'})]
85
86 X_train.cov()
87
88 # **STEP 3:CoVariance Matrix Calculation
89 # covariance matrix will be 30*30 since X
90 # =(569,30). [each 30 features has
91 # covariance with each 30. so 30*30
92 # matrix]
93 covariance_matrix = np.cov(X_train,
94     rowvar=False)
95 print("Covariance matrix shape:",,
96     covariance_matrix.shape)
97 attribute_labels = X.columns
98 covariance_matrix = pd.DataFrame(
99     covariance_matrix, columns=
100     attribute_labels, index=
101     attribute_labels)
102 print(covariance_matrix.head())
103
104 # Visualize covariance matrix
105 plt.figure(figsize=(12, 8))
106 sns.heatmap(covariance_matrix, cmap='
107     coolwarm', annot=True, annot_kws={'
108         size': 5})
109 plt.title('Covariance Matrix on Processed
110     Cancer Data')
111 plt.text(8, -1.2, "'THA076BCT026',",
112         "'THA076BCT027', 'THA076BCT041'",,
113         fontsize=8,color='red')
114 plt.savefig('./plots/Covariance Matrix
115     Processed Cancer Data', bbox_inches='
116         tight')
117 plt.show()
118
119 #### **STEP 4: Eigen vectors/values
120 # calculation**
121 eigenvalues, eigenvectors = np.linalg.eig
122     (covariance_matrix)
123 print("shape of eigenvalue:",eigenvalues.
124     shape)
125 print("shape of eigen vectors:",,
126     eigenvectors.shape)
127 print("eigen values are",eigenvalues)
128 print("eigen vectors are",eigenvectors)
129
130 # **STEP 5: See the variance explained by
131 # each eigen values.**
132 # Element wise array division to obtain a

```

```

    new array
109 # proportion_of_variance array provides
   insights into relative importance of
   each PC in capturing the overall
   variance in the dataset
110 proportion_of_variance = eigenvalues/sum(
   eigenvalues)
111 proportion_of_variance
112
113 # Plot a scree plot to visualize how much
   variance is captured by which PC
114 cumulative_variance = np.cumsum(
   proportion_of_variance)
115
116 # range(1 to 30) on x axis, cummulative
   pov on y axis
117 plt.plot(range(1, len(
   proportion_of_variance) + 1),
   proportion_of_variance, marker='x',
   color='blue')
118 plt.plot(range(1, len(
   proportion_of_variance) + 1),
   cumulative_variance, marker='o',
   color='orange')
119 plt.xlabel('Principal Components (Eigen
   Vector)')
120 plt.ylabel('Proportion of Variance')
121 plt.title('Scree Plot')
122 legend_class = ["Explained Variance", "
   Cumulative Explained Variance"]
123 plt.legend(labels = legend_class)
124 plt.text(12, 0.9, "'THA076BCT026'",
   'THA076BCT027', 'THA076BCT041'",
   fontsize=8,color='red')
125 plt.savefig('./plots/Scree Plot',
   bbox_inches='tight')
126 plt.show()
127
128
129 # **STEP 6: Select the eigen vectors as
   needed to compute the final output**
   Formula: New data (Y) = Row feature
   vector * Row data
130 # Transpose data into row wise
131 transposed_X_train = np.transpose(X_train
   )
132 print("shape of transposed x train is",
   transposed_X_train.shape)
133 # Transpose eigen vectors into row wise
   for row feature vector
134 eigenvectors_transposed=np.transpose(
   eigenvectors)
135
136 # Dictionary to store the final Y(because
   we want to have multiple Y, each
   obtained by selecting our choice of
   combination of PC).
137 Y = {}
138 selected_components = [[0], [0, 1],
   [0,3], [2,3], [0,29]]
   # best PC, 2
   best PCs, best+medium, medium+medium,
   best+worst
139 count=len(selected_components)
140
141 # Iterate over the selected components
142 for i in range (0,count):

```

```

143     selected_eigenvectors =
   eigenvectors_transposed[
   selected_components[i]]           #
   select the combination of
   components from above list
144     Y[i] = np.transpose(
   selected_eigenvectors @
   transposed_X_train)               #
   projection
145     print(Y[i].shape)
146
147 # checking the nature of Y
148 print(type(Y))
149 print(Y.keys())
150 Y[1].shape
151 print(type(Y[1]))
152
153
154 # **STEP 7: Computing and visualizing
   covariance for the new data.**
155 # Dictionary to store all covariance
   matrices(obtained by selecting our
   choice of combination of PC).
156 new_covariance_matrices_dictionary = {}
157 for i in range(0, count):
158     new_covariance_matrix = np.cov(Y[i],
   rowvar=False)
159     new_covariance_matrices_dictionary[i] =
   new_covariance_matrix
160
161     # Combination:
162     print("*Covariance matrix shape for
       combination",i,"is",
       new_covariance_matrix.shape)
163     # 0=best PC, 1=2 best PCs, 2= best+
       medium PCs
164     print("-->Covariance matrix:",
       new_covariance_matrices_dictionary[
       i])                                # 3=
       medium+medium PCs, 4= best+worst
       PCs
165
166 # Heatmap to visualize covariance matrix
   combination.
167 plt.imshow(
   new_covariance_matrices_dictionary
   [1], cmap='Reds', interpolation='
   nearest') #instead of 'nearest', we
   can use 'bilinear','bicubic','
   spline16', 'spline36', 'hanning', '
   hamming', 'hermite',etc..
168 plt.colorbar()
169 plt.title("New Covariance Matrix
   Combination 1 (2 best PCs)")
170 plt.show()
171
172 #lets try bilinear,pixel coloring is
   smoother(linear interpolation between
   4 nearest data points)
173 plt.imshow(
   new_covariance_matrices_dictionary
   [2], cmap='Reds', interpolation='
   bilinear')
174 plt.colorbar()
175 plt.title("New Covariance Matrix
   Combination 2 (Best and medium PCs)")
176 plt.show()

```

```

174
175 #lets try bicubic,pixel coloring is again
176     smoother(linear interpolation based
177     on 16 nearest data points)
178 plt.imshow(
179     new_covariance_matrices_dictionary
180     [3], cmap='Reds', interpolation='
181     bicubic')
182 plt.colorbar()
183 plt.title("New Covariance Matrix
184     Combination 3 (Two medium PCs)")
185 plt.show()

186 #lets try spline16,pixel coloring is even
187     smoother(16 degree of polynomial
188     used for interpolation, for
189     flexibility)
190 plt.imshow(
191     new_covariance_matrices_dictionary
192     [4], cmap='Reds', interpolation='
193     spline16')
194 plt.colorbar()
195 plt.title("New Covariance Matrix
196     Combination 4 (Best and worst PCs)")
197 plt.show()

198 # **STEP 8:Visualize the final output**
199 #Combination 0 = Y[0] = Best PC
200 #plt.ylim(-10,10)
201 #plt.xlim(-10,10)
202 varlegend= "Roll: 26,27,41"
203 plt.text(5,0.04,varlegend,color='red')
204 plt.title("1 best PC")
205 plt.xlabel('PC1')
206 plt.scatter(Y[0], np.zeros_like(Y[0]), c=
207     y_train) # colors = based on
208     y_train

209 # Combination 1 = Y[1] = 2 best PCs
210 #plt.ylim(-3,3)
211 #plt.xlim(-3,3)
212 varlegend= "Roll: 26,27,41"
213 plt.text(5,4,varlegend,color='red')
214 plt.title("2 best PCs")
215 plt.xlabel('PC1')
216 plt.ylabel('PC2')
217 plt.scatter(Y[2].iloc[:, 0], Y[2].iloc[:, 1],
218     c=y_train) # iloc used to
219     access data

220 # Combination 3 = Y[3] = 2 medium PCs
221 #plt.ylim(-3,3)
222 #plt.xlim(-3,3)
223 varlegend= "Roll: 26,27,41"
224 plt.text(2,4.5,varlegend,color='red')
225 plt.title("2 medium PCs")
226 plt.xlabel('PC3')
227 plt.ylabel('PC4')
228 plt.scatter(Y[3].iloc[:,0],Y[3].iloc
229     [:,1],c=y_train)

230 # Combination 4 = Y[4] = Best and worst
231     PCs
232 plt.ylim(-3,3)
233 #plt.xlim(-3,3)
234 varlegend= "Roll: 26,27,41"

235 plt.text(5,2,varlegend,color='red')
236 plt.title("Best and worst PCs")
237 plt.xlabel('PC1')
238 plt.ylabel('PC30')
239 plt.scatter(Y[4].iloc[:,0],Y[4].iloc
240     [:,1],c=y_train)

241 # **STEP 9: PCA using Library**
242 from sklearn.decomposition import PCA
243 # Create instances of PCA with the
244     desired number of components. (pca1 =
245     2components. pca2=so as to retain
246     95% info)
247 pca1 = PCA(n_components=2)
248 pca2 = PCA(0.95)

249 # Fit the PCA model to the data and
250     transform the data
251 reduced_data1 = pca1.fit_transform(
252     X_train)
253 reduced_data2 = pca2.fit_transform(
254     X_train)

255 # Print the shape of the reduced data
256 print("Shape of reduced data:",
257     reduced_data1.shape)
258 print("reduced through first",
259     reduced_data1)
260 print("Shape of reduced data:",
261     reduced_data2.shape)
262 print("reduced through second",
263     reduced_data2)

264 # Here it can be seen that to retain 95%
265     info, 10 principle components are
266     required
267 # Visualize the output for the first PCA
268     done by library
269 #plt.ylim(-3,3)
270 #plt.xlim(-3,3)
271 varlegend= "Roll: 26,27,41"
272 plt.text(10,10,varlegend,color='red')
273 plt.title("Two best PCs-Using Library")
274 plt.xlabel('PC1')
275 plt.ylabel('PC2')
276 plt.scatter(reduced_data1.T[0],
277     reduced_data1.T[1], c= y_train) #
278     After applying transform method of
279     PCA,Each row was represented as a
280     sample. So need to transpose.

281 # For the 2nd PCA done by library(to
282     retain 95% info), 10 components were
283     required! We will try to plot 3d
284     scatter plot to visualize.
285 from mpl_toolkits.mplot3d import Axes3D
286 fig = plt.figure()
287 fig = plt.figure(figsize=(10, 10))
288 ax = fig.add_subplot(111, projection='3d'
289     )
290 varlegend= "Roll: 26,27,41"
291 ax.text(-5,12.5,10,varlegend,color='red')

292 # 3d coordinates required to specify

```

```
the legend position
264 ax.scatter(reduced_data2.T[0],
265     reduced_data2.T[1],reduced_data2.T
266     [2],c=y_train)
267 ax.set_xlabel('PC1')
268 ax.set_ylabel('PC2')
269 ax.set_zlabel('PC3')
270 ax.set_title('3D Scatter Plot of 3 PCs')
271 plt.show()

# For the next step, we can train a model
# and compare the performance to see
# the impact of PCA we did so far.
```

• • •