

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 Decision tree analysis is a popular and widely used machine learning technique for classification tasks. This project applied decision tree analysis to the Dry Bean dataset, aiming to develop an accurate classification model for different dry bean varieties based on their physical properties and shape characteristics. The dataset was pre-processed, and a decision tree algorithm was trained and optimized using cross-validation. Performance evaluation metrics, including accuracy, precision, recall, and F1-score, were used to assess the model. The decision tree model achieved high accuracy and favourable results across multiple metrics. Visualizing the decision tree structure provided insights into important features and decision paths. Overall, this project highlights the effectiveness of decision tree analysis in classifying dry bean varieties based on their physical properties and shape characteristics. The developed decision tree model can be a useful tool for automatic classification of dry bean varieties in agricultural and food industry applications, facilitating quality control and product differentiation.

INDEX TERMS Decision Tree, cross-validation

I. INTRODUCTION

Decision tree analysis is a visual representation that illustrates possible outcomes resulting from a sequence of interconnected choices. It provides a framework for evaluating different factors and decisions to determine the most favourable outcome. Typically, a decision tree analysis begins with a central decision and expands into branches representing additional choices or potential results. It provides a structured and intuitive approach to analyse complex datasets. This report focuses on applying decision tree analysis to the Dry Bean dataset obtained from the UCI Machine Learning Repository. The dataset comprises physical properties and shape characteristics of various dry bean varieties consumed globally. The primary goal is to develop a decision tree model capable of accurately classifying dry bean varieties based on their attributes. To achieve this, we perform crucial pre-processing steps, including handling missing values, normalizing numerical features, and encoding categorical variables. Subsequently, the dataset is divided into training and testing sets to evaluate the model's performance.

The report then proceeds with training the decision tree

algorithm and optimizing its hyperparameters through cross-validation. Various evaluation metrics such as accuracy, precision, recall, and F1-score are utilized to assess the effectiveness of the decision tree analysis on the Dry Bean dataset. Additionally, the decision tree structure is visualized to gain valuable insights into the significant attributes and decision paths contributing to the classification process. By showcasing the applicability of decision tree analysis, this report emphasizes its potential in accurately classifying dry bean varieties based on their physical properties and shape characteristics. It highlights the relevance of decision tree analysis in agricultural and food industry applications, including quality control and product differentiation.

II. METHODOLOGY

A. THEORY

1) Background

The Decision tree is an intuitive non-parametric supervised machine-learning model that can be used for classification as well as regression tasks by creating a tree-like flowchart structure to make predictions based on input features. It mimics the human decision-making process as the decision

rules and splits in the tree structure make it clear how the input features influence the final predictions or classifications. Each internal node in the decision tree represents a test on an attribute/feature, each leaf node denotes the class label, and each branch represents the outcome of the test leading to the class label in the leaf node. The classification rule can be obtained by tracing the paths from root to leaf. Decision trees offer a measure of feature importance, which can help us identify the most influential factors driving the predictions or classifications.

A decision tree algorithm called Iterative Dichotomizer (ID3) was developed in the late 1970s. A successor of ID3 called C4.5 was presented later which became a standard. Likewise, CART (Classification and Regression Trees) was put forward which described the generation of binary trees. These algorithms are called greedy algorithms where the top-down approach of recursively dividing and conquering is done to construct the tree.

2) Decision Tree Induction

Decision tree induction starts with a dataset containing training examples with their corresponding class labels. The greedy algorithms mentioned identify the best feature to split the dataset based on certain criteria like information gain. The selected feature is used as the root (topmost node) of the decision tree, and the dataset is split into subsets based on the feature values. The process is repeated recursively for each subset, creating branches and sub-trees until a stopping condition is met. The resulting decision tree can be visualized as a flowchart-like structure, with each node representing a test on a specific feature. To make predictions on unseen data, the algorithm traverses the decision tree by following the path based on the feature values of the input instance until it reaches a leaf node. During the induction of the decision trees, there might be various stopping criteria which helps the tree become less prone to being complex and avoiding overfitting.

3) Attribute Selection Criterion

It's a heuristic to select the splitting criteria that can separate the given data into individual classes in the best way possible. Different attribute selection measures have their own calculation methods and criteria for evaluating the importance of attributes. Some of them include Information gain, Gini Index, Gain ratio, Chi-Square, Relieff, etc.

a: Information Gain

Information gain is a widely used attribute selection measure in decision tree algorithms such as ID3. It quantifies the amount of information gained about the target variable by splitting the data based on a specific attribute. It is based on the concept of entropy, which measures the impurity or uncertainty of a dataset.

$$E = - \sum_{i=1}^N p_i \log_2(p_i) \quad (1)$$

Considering a dataset with N classes, 'pi' is the probability of randomly selecting an example in class 'i'. The formula is used to calculate the initial entropy before split (Ebefore) and the entropy after the split (Espli). Then the information gain is calculated using the formula:

$$IGain = E_{\text{before}} - E_{\text{split}} \quad (2)$$

The more the entropy removed, the greater the information gain. The higher the information gain, the better the split. The attribute with the highest information gain is selected as the best attribute for splitting.

b: Gini Index

The Gini index is a measure of impurity or inequality often used in decision tree algorithms such as CART. The formula to calculate the Gini index is:

$$Gini = 1 - \sum_{i=1}^n (p_i)^2 \quad (3)$$

The Gini index ranges from 0 to 1, where 0 indicates a perfectly pure node (all instances belong to the same class) and 1 indicates maximum impurity (instances are evenly distributed across classes). To calculate the Gini index for a specific attribute, the dataset is split based on the attribute's values, and the Gini index is calculated for each resulting subset. The attribute with the lowest Gini index (i.e., highest purity) is selected as the best attribute for splitting. Unlike information gain, the Gini index does not explicitly consider the information content of attributes but focuses solely on the distribution of class labels.

c: Gain Ratio

The Gain ratio is a modification of the information gain measure that addresses the bias towards attributes with many distinct values. The C4.5 algorithm utilizes the Gain ratio as the attribute selection measure. The formula for calculating the Gain ratio is as follows:

$$\text{Gain_ratio} = \frac{\text{Gain(Attribute)}}{\text{Split_info(Attribute)}} \quad (4)$$

d: Chi-Square

It measures the statistical significance of the relationship between the attribute and the class variable. The formula for calculating Chi-square is:

$$\chi^2 = \sum \left(\frac{(O_i - E_i)^2}{E_i} \right) \quad (5)$$

where, Oi is the observed value and Ei is the expected value.

The Chi-Squared statistic follows a Chi-Squared distribution, and its value indicates the strength of association between the attribute and the class variable. The higher the value of chi-square more the sub-nodes are different from the parent node and hence the homogeneity is more. Thus, the attribute with the highest Chi-Squared value is selected as

the best attribute for splitting the data. This test is particularly useful for datasets with categorical attributes and categorical class variables.

e: Relieff

This is an algorithm rather than a mathematical formula such as information gain. Relieff, unlike other greedy algorithms, can detect attribute dependencies. This algorithm uses a sampling-based approach to estimate attribute weights, making it computationally efficient for large datasets. It can handle noisy or incomplete data and is robust against irrelevant attributes. Relieff assigns higher weights to attributes that have consistent values among the nearest neighbors of the same class and different values among the nearest neighbors of different classes. The weights assigned to attributes represent their importance in discriminating between different classes. The main purpose of including Relieff in this section is to mention the use of this algorithm to delete the irrelevant attributes for the classification so that we can obtain an improved decision tree.

4) Decision Tree Induction combinations

The attributes can be discrete-valued or continuous-valued. Likewise, depending on the target, the problem at hand might be classification or regression. The following points clarify the induction process for each of the cases that may arise depending on the dataset chosen:

a: Classification for Discrete-valued attributes

In the context of classification problems in decision trees, categorical attributes and discrete attributes are handled similarly, and thus used interchangeably although they are not exactly the same in meaning. Examples of categorical attributes include "color" (with categories like red, blue, green) or "gender" (with categories like male, female). The best discrete-valued attribute is chosen from the available attributes as the splitting attribute. Then the dataset is partitioned into subsets. The previously chosen attribute need not be considered for any future partitioning, that's why it is removed from the attribute list. The recursive partitioning is carried out for each child node.

b: Classification for Continuous-valued attributes

Examples of continuous (or numeric) attributes include "age," "income," or "temperature." The continuous-valued attributes can be discretized using buckets, intervals, percentiles, or clusters to divide the range of attribute values. Instances falling within each bucket are assigned to the corresponding subset. An impurity measure, such as the Gini index or entropy, for each potential split is calculated to maximize the information gain (minimize impurity). Besides discretizing using buckets, the threshold-based approach can be used to select a threshold value for the attribute. Instances with attribute values greater than or equal to the threshold are assigned to one subset, while instances with values below the threshold are assigned to the other subset. Similarly, the

splitting process for each subset created by the split is done recursively, considering only the remaining attributes.

c: Classification for a mixture of Categorical and numeric attributes

The approach is like handling discrete and continuous attributes separately. For numerical attributes, discretization techniques can be applied to convert them into categorical attributes. This can be done through methods such as binning, percentile-based discretization, or clustering-based discretization. Alternatively, binary decisions can be made by selecting a threshold value to separate the data based on whether the attribute value is above or below the threshold. The obtained categorical data is partitioned based on the attribute values, and the process is repeated recursively for each partition.

d: Regression for numeric data

Numerical data for regression refers to the case where the attribute values are numeric, and the target variable is also numeric. Both the attribute and the target are continuous variables. For the regression problem, we can't use the splitting criteria like calculating entropy like we do in classification problems. Thus, we use an appropriate attribute selection measure, such as mean squared error or mean absolute error, to determine the best attribute for splitting the data.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y}_i)^2 \quad (6)$$

For the regression decision tree, we try to minimize the MSE (Mean Square Error) instead of entropy as we do in the Classification trees. For numerical attributes, binary decisions are made based on a threshold value to determine the direction of the split. The attribute values are compared to the threshold, and the data is partitioned accordingly. Splitting of the data is based on the selected attribute and continued recursively until a stopping condition is met or a desired level of accuracy is achieved. Once the decision tree is built, it can be used to make predictions on new instances by traversing the tree and assigning the predicted numeric value based on the path followed. To reduce the risk of overfitting, we can specify the minimum number of children per leaf node in advance and use cross-validation to find this value.

5) Tree Pruning

Tree pruning is a technique used in decision tree algorithms to reduce the complexity and overfitting of the tree model. Overfitting usually occurs when a decision tree captures the noise and irregularities in the training data, leading to poor generalization on unseen data. The goal of tree pruning is to find the right balance between complexity and accuracy by removing unnecessary branches and nodes from the tree. Two common approaches to tree pruning include: pre-pruning and post-pruning. Pre-pruning is a pruning method that stops the tree growth early based on predefined criteria, such as

maximum depth or minimum number of instances in a leaf. Post-pruning, also known as error-based pruning, involves growing the tree to its entirety and then removing branches based on their estimated error rate. Another could be Rule Post-Pruning, where the decision tree is converted into a set of if-then rules and their accuracy is evaluated on the validation data. Rules that do not contribute significantly to the model's performance are pruned. CART uses Cost complexity pruning, also called minimal cost complexity pruning or alpha-beta pruning, employing a parameter called the complexity parameter to control the pruning process. The pruning process involves evaluating the error rate or other metrics on a validation set or using cross-validation to estimate the performance of the pruned tree. Pruning reduces the risk of overfitting by simplifying the decision tree model and improving its generalization on unseen data. Thus, Tree pruning is an essential technique in decision tree learning, contributing to improved accuracy, reduced complexity, and enhanced model interpretability.

B. SYSTEM BLOCK DIAGRAM

The system architecture in Figure 1 demonstrates the sequence of blocks involved in the decision tree analysis. Initially, diverse datasets were collected, and data visualization techniques such as scatter plots, histograms, and heatmaps were utilized to gain insights into the data's distribution and relationships. For decision tree analysis, the preprocessing stage involves steps such as handling missing values, normalizing numerical features, and encoding categorical variables. This ensures that the data is brought to a common scale, enabling all features to contribute equally to the analysis. The dataset is split into training and testing sets, and the decision tree algorithm is trained using the training set. Hyperparameters are optimized, and the algorithm's performance is evaluated using metrics such as accuracy, precision, recall, and F1-score. Once the decision tree model is trained, it can be visualized to gain insights into the significant attributes and decision paths that contribute to the classification process. The visualization provides a clear representation of the decision tree structure, aiding in understanding the model's behavior and interpreting its results. In summary, the decision tree analysis follows a similar sequence to the architecture in Figure 1, but with distinct steps focusing on feature selection, model training, hyperparameter optimization, and visualization of the decision tree structure.

C. INSTRUMENTATION

The implementation of Decision Tree in our study involved the use of several libraries and tools to facilitate data visualization and mathematical operations.

1) Pandas

Pandas, a powerful Python library for data manipulation, plays a crucial role in the decision tree pipeline. It provides essential functions and data structures for data handling, exploration, and analysis. Pandas is used to load and read data,

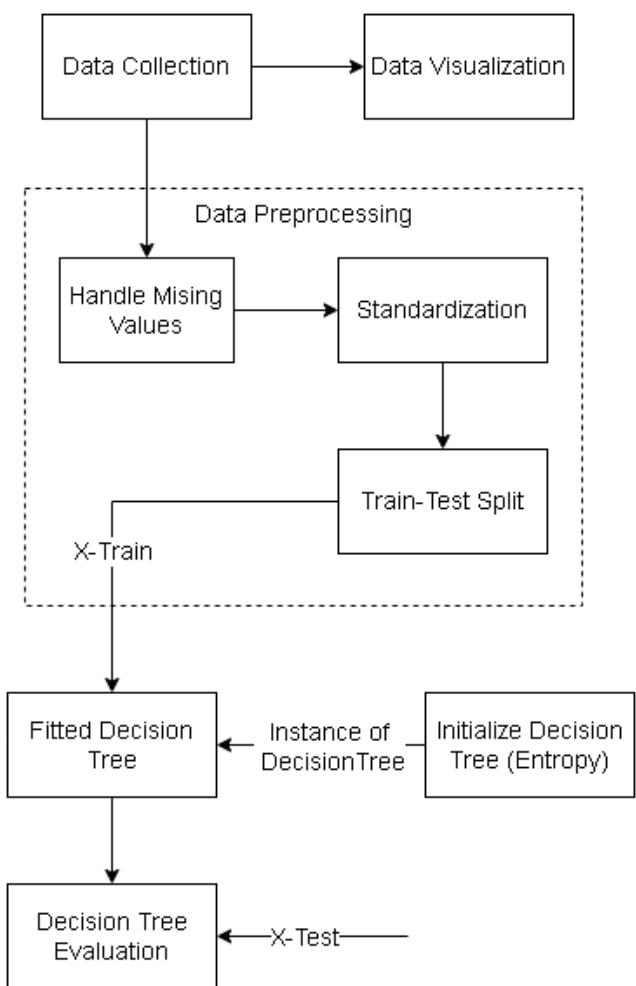


FIGURE 1: System Architecture of Decision Tree

inspect its structure, handle missing values, encode categorical variables, scale numerical data, and perform descriptive statistics. It enables efficient data pre-processing and feature extraction, making it an invaluable tool for decision tree analysis.

2) Matplotlib

Matplotlib is a powerful data visualization library in Python. It contributes to the instrumentation of decision tree models by offering essential plotting capabilities. While it doesn't directly construct decision trees, matplotlib enables the visualization of decision tree structures as hierarchical diagrams or tree plots, aiding in understanding the decision paths. Additionally, matplotlib facilitates the creation of feature importance plots, allowing users to identify influential attributes. Moreover, it enables the visualization of decision boundaries, aiding in comprehending how the model classifies different regions of the feature space. By leveraging matplotlib's visualization capabilities, users can gain valuable insights into decision tree models and effectively communicate their findings.

3) NumPy

NumPy, also known as Numerical Python, is essential in the implementation of decision tree models, offering efficient array operations and mathematical functions. With its multidimensional array object and comprehensive set of functions, NumPy enables high-performance computations, including mathematical operations, shape manipulation, and linear algebra operations. In the context of decision tree analysis, NumPy is instrumental in performing array operations such as element-wise calculations, matrix multiplication, transposition, and the computation of eigenvalues and eigenvectors. Key NumPy functions like `np.dot()`, `np.transpose()`, and `np.linalg.eig()` are utilized for these tasks. Specifically, `np.linalg.eig()` is particularly valuable in decision tree analysis as it determines the direction and magnitude of components crucial for constructing the decision tree. By leveraging NumPy's capabilities, the implementation of decision tree models becomes more efficient, allowing for seamless manipulation and analysis of arrays required for accurate classification.

4) Seaborn

Seaborn, a Python data visualization library, enriches the analysis and interpretation of decision tree models by providing powerful visualization capabilities. While not involved in constructing decision trees directly, seaborn's plotting functions allow for visualizing the decision tree structure through tree plots and dendograms, aiding in understanding the decision paths. Seaborn further facilitates the visualization of feature importance, enabling the identification of influential attributes and their impact on classification. Additionally, seaborn offers a variety of plotting options, such as pair plots, scatter plots, and heatmaps, which assist in exploring feature relationships and their effects on the target variable. By leveraging seaborn's visualization capabilities, users can gain deeper insights into decision tree models, effectively communicate their findings, and enhance the overall analysis and interpretation process.

5) Dtreetviz

Dtreetviz is a Python library that enhances decision tree analysis through interactive and visually appealing visualizations. It allows users to explore the decision tree structure interactively, gaining insights into the decision paths and splits. Dtreetviz provides feature importance plots to identify influential features and offers other visualizations such as class prediction histograms and sample-based explanations. By leveraging dtreetviz, users can gain a deeper understanding of decision tree models, improving interpretation and decision-making.

By utilizing these libraries, we were able to seamlessly integrate data visualization and mathematical operations into our Decision Tree Analysis. These libraries enabled effective data visualization, allowing us to analyse datasets, visualize decision tree structures, and gain valuable insights. By lever-

aging these libraries, we combined data analysis and visualization to improve decision tree analysis and interpretation.

D. WORKING PRINCIPLE

The decision tree Analysis consists a the pipeline of steps involving data collection, preprocessing, training the algorithm, hyperparameter optimization, performance evaluation, and decision tree visualization. The dataset is collected, pre-processed to handle missing values and normalize features, and split into training and testing sets. The decision tree algorithm is trained using the training set, creating a structure that recursively splits the data based on informative attributes. Hyperparameters are optimized to enhance the model's performance. The model's accuracy, precision, recall, and F1-score are evaluated using the testing set. Decision tree visualization provides insights into significant features and decision paths. Overall, this pipeline enables accurate classification of dry bean varieties based on their physical properties and shape characteristics, benefiting applications in agriculture and the food industry.

1) Data Pre-processing

The input dataset is pre-processed to ensure its suitability for decision tree analysis. This includes handling missing values, scaling the data, and addressing outliers. In our implementation, specific data pre-processing techniques were applied to each dataset, such as standardization or normalization. Categorical target variables were encoded to convert them into a numerical format. For instance, label encoding was used to represent the seven classes (e.g., 'SEKER', 'BARBUNYA', 'BOMBAY', 'CALI', 'HOROZ', 'SIRA', 'DERMASON') with numeric values. The processed dataset can be visualized using the provided figures. This pre-processing step is essential to ensure that the decision tree algorithm can effectively analyse the data and make accurate classifications based on the attributes. By pre-processing the data, we address any inconsistencies or anomalies that could affect the performance of the decision tree model. The visualization of the processed dataset provides insights into the transformed data and allows for a better understanding of its distribution and relationships. Overall, data pre-processing is a crucial step in preparing the input dataset for decision tree analysis, ensuring that it meets the necessary requirements for accurate classification.

2) Train-Test Split

After pre-processing, the dataset is partitioned into two subsets: a training set and a testing set. The training set is utilized to train the decision tree algorithm. By leveraging the provided features and their associated labels, the algorithm learns to construct a decision tree structure. This structure is formed through a recursive process of splitting the data based on informative attributes, creating branches and leaf nodes that represent different classification outcomes. The decision tree algorithm learns from the training set to capture patterns and relationships between the features and labels, enabling it to make informed decisions during classification.

3) Model Initialization

The initialization of a decision tree model involves setting up the parameters and configurations necessary to create the decision tree structure. This includes specifying the splitting criterion (e.g., entropy or Gini index), setting the maximum depth of the tree, and defining other hyperparameters like the minimum number of samples required for a node split. By initializing the decision tree model with appropriate parameters, we lay the foundation for the subsequent training and optimization processes.

4) Hyperparameter tuning

Hyperparameter tuning optimizes the decision tree model by fine-tuning parameters like tree depth, minimum sample split, and splitting criterion. Techniques like cross-validation determine the optimal parameter values. This iterative process aims to strike a balance between model complexity and generalization ability, ensuring accurate predictions on unseen data. Through hyperparameter optimization, the decision tree model achieves enhanced performance and optimal results.

5) Model Evaluation

The performance of the decision tree model is assessed using the testing set, where multiple metrics, including accuracy, precision, recall, and F1-score, are computed to evaluate its classification capabilities. These metrics provide insights into the model's accuracy in correctly predicting the class labels of the dry bean varieties. By analysing these performance measures, we can assess the effectiveness of the decision tree model in accurately classifying the dry bean varieties based on their attributes.

6) Model Visualization

Visualizing the decision tree structure allows to gain valuable insights into the important features and decision paths within the model. By employing visualization techniques, we can represent the decision tree's decision-making process graphically, enabling easier interpretation and comprehension of the key factors contributing to the classification outcomes. These visual representations aid in understanding how the decision tree makes decisions based on different attributes, shedding light on the underlying logic and helping us identify the most influential features in the classification process.

To summarize, the decision tree pipeline involves collecting the data, pre-processing it, training the algorithm, optimizing the hyperparameters, evaluating the model's performance, and visualizing the decision tree. This comprehensive process enables the accurate classification of dry bean varieties by considering their physical properties and shape characteristics. The application of this pipeline holds significant potential in agricultural and food industry domains, providing valuable insights for quality control and differentiation purposes. By following these steps, the decision tree analysis becomes a powerful tool in efficiently and effectively classifying dry bean varieties, contributing to advancements in the agricultural sector.

III. RESULTS

A. DISTRIBUTION OF DATA

The distribution of the dataset feature was analysed using a histogram plot. The histogram provided insights into the frequency distribution of the feature values, allowing us to understand the underlying distribution pattern and identify any potential outliers or data skewness.

B. PRE-PROCESSING

The original dataset contained the class labels as 'SEKER', 'BARBUNYA', 'BOMBAY', 'CALI', 'HOROZ', 'SIRA', 'DERMASON' class. As part of the pre-processing phase, label encoding was applied to transform categorical data into numerical values. These labels were encoded using a range from 0 to 6, assigning a unique numerical representation to each category. This transformation enabled us to effectively utilize the categorical feature in subsequent analysis or modelling tasks, as many machine learning algorithms require numerical input.

C. MODEL PERFORMANCE

1) General model performance

This is the initial performance of the model without doing any processing or hyperparameter tuning. The general decision tree model achieved an accuracy of 0.888. The classification report showed precision, recall, and F1-score values for each class, along with the support count. The macro-average of 0.9 and the weighted-average scores of 0.89 indicate the overall performance of the model. Similarly, the confusion matrix shown is useful in visualizing the performance of the classification model for given set of test data.

[insert classification report figure here]

[insert confusion matrix report figure here]

2) Model performance after PCA

The model is expected to perform better after the application of PCA (Principal Component Analysis). It is observed that after applying PCA with 7 components covering 99.60

[insert Feature importance (having both cumulative and individual importance curve) figure here]

[insert scree plot here]

Likewise, the classification report for this model shows a macro average of 0.91 and a weighted average of 0.90. The confusion matrix shown can further help in visualizing the performance of this model for given test data.

[insert classification report here]

[insert confusion matrix here]

3) Model performance after Pre-Pruning

Similarly, the decision tree model's performance was evaluated after applying pre-pruning by specifying the maximum depth of the tree. The accuracy, precision, recall, and F1-score values reflect the model's performance after pre-pruning. When a maximum depth of 3 was chosen, the model gave an accuracy of 83.54%. It was also seen that the macro

average was 0.73 and the weighted average was 0.82. The confusion matrix for this model shows a different kind of performance in the test data.

[insert classification report/ confusion matrix figure here]

Another pre-pruning technique of minimum sample split was done as well to evaluate the model performance. Here we specified the minimum number of samples required to do a split. Keeping the criteria as Entropy and a minimum sample split of 1000 gave an accuracy of 87.64%. Likewise, the macro average and weighted average of 0.89 and 0.88 was obtained respectively. It is further explained from the obtained classification report and confusion matrix.

[insert report/confusion matrix]

4) Model performance after Post-Pruning

Likewise, the decision tree model's performance was evaluated after applying post-pruning techniques, such as Grid SearchCV and Cost Complexity pruning path. GridSearchCV is a technique for finding the optimal parameter values from a given set of parameters in a grid., which is essentially a cross-validation technique. The best parameters were obtained which showed that the optimal maximum depth was 5 and the optimal criterion was log-loss. Doing so gave an accuracy of 85.90

[insert report/confusion matrix here]

Likewise, another post-pruning technique called Cost Complexity Pruning Path helps in finding the optimal pruning parameter by exploring the trade-off between tree complexity and accuracy. The accuracy, precision, recall, and F1-score values reflect the improved performance of the model after Cost Complexity Pruning Path.

[insert effective alpha and impurity figure here]

IV. DISCUSSION AND ANALYSIS

The dry bean dataset contains a collection of 13611 instances belonging to 7 classes of beans. The data were normally distributed with respect to each feature, and some were skewed to the left or right, as seen in the histogram. The successful application of the Decision Tree Classifier to the dry bean dataset demonstrated its effectiveness in accurately classifying the different varieties of dry beans. By constructing a hierarchical structure of decision rules based on the dataset's features, the decision tree model achieved distinct separations between the dry bean varieties, providing reliable class predictions.

A noteworthy finding from the analysis is the impact of the decision tree's depth on its performance. Deeper trees generally exhibit higher accuracy in classifying the dry bean varieties. However, it is crucial to find the right balance to avoid overfitting, where the model captures irrelevant or noisy features. Therefore, selecting an optimal depth for the decision tree is essential to achieve optimal performance without sacrificing generalization. Methods like pruning, Increasing the minimum sample split, fixed depth, etc can be used to limit the depth of the Decision Tree.

The change in the decision tree model's performance after Grid SearchCV post-pruning can be attributed to the optimization of hyperparameters. By fine-tuning the pruning parameters and exploring different combinations of hyperparameters to find the best configuration Grid SearchCV helps to reduce overfitting and enhance the generalization capabilities of the model. Likewise, the Cost Complexity Pruning Path allows for the identification of the optimal pruning parameter by striking a balance between simplicity and accuracy. Overall, both Grid SearchCV and Cost Complexity Pruning Path techniques have shown to be effective in improving the performance of the decision tree model through appropriate post-pruning. The various curve in relation to alpha showed different insights, where alpha refers to the pruning parameter used in the decision tree pruning process. It determines the level of pruning applied to the tree. Initially, as alpha increases, the accuracy of the model shows a gradual improvement. However, after a certain threshold, further increasing alpha leads to a plateau in accuracy, indicating that additional pruning has a limited impact on enhancing model performance. Similarly, as alpha increases, the depth of the tree decreases sharply, resulting in a simpler and less complex model. However, beyond a certain alpha value, the depth curve levels off, indicating that additional pruning has minimal impact on reducing the tree's depth. Also, as alpha increases initially, the total impurity sharply increases, indicating a reduction in model performance and an increase in impurity. However, beyond a certain alpha value, the impurity curve levels off, suggesting that further pruning does not significantly affect the overall impurity. Thus, a higher alpha value results in more aggressive pruning, leading to a simpler model with potentially reduced overfitting. The curves show the effect of different alpha values on various aspects of the decision tree model, providing insights into the trade-off between complexity and performance.

Furthermore, the Decision Tree Classifier allows for the assessment of feature importance. By examining the top-splitting nodes and their associated features, we can identify the most significant attributes in classifying the dry bean varieties. This insight provides valuable information for further analysis and understanding of the factors influencing the classification process. Its hierarchical decision-making structure simplifies the classification process, while the selection of an appropriate tree depth and consideration of feature importance contribute to achieving accurate and interpretable results. By employing pre-pruning and post-pruning techniques, we have successfully optimized the model's complexity and achieved a balance between accuracy and generalization.

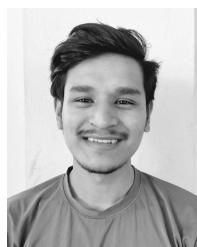
V. CONCLUSION

In conclusion, the implementation of the decision tree algorithm on the dry bean dataset has provided us with valuable insights into its classification capabilities. The objective of implementing the decision tree algorithm using Information gain as the attribute selection measure on the dry bean dataset

was successfully achieved. By utilizing the decision tree classifier with the entropy criterion and a maximum depth of five, we obtained a commendable accuracy score of 89.36% on the test data. The decision tree algorithm demonstrated its ability to effectively handle the classification task, enabling it to identify the most informative features for classification, and therefore showcasing its power in distinguishing between different classes based on the dataset's attributes. The implementation of pre-pruning techniques, such as limiting the maximum depth, helped prevent overfitting and ensured the model's generalization capability. Thus, Decision tree algorithms offer a powerful and interpretable approach to solving classification and regression problems by creating a hierarchical structure of decision rules based on input features. With their ability to handle both discrete and continuous attributes, decision trees can provide flexibility in modeling a wide range of datasets and can effectively capture complex patterns and relationships. With further exploration and fine-tuning of hyperparameters, decision tree models hold great potential for solving a wide range of classification problems in various domains.

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. Patnaik, "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-eeab7cb73b72>.



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
10 # x1 and x2 represent the 2 axes
11 x1 = np.random.randn(50)
12 x2 = np.random.randn(50)
13 print(x1)
14 print(x2)
15 plt.scatter(x1,x2)
16 plt.xlabel("X1")
17 plt.ylabel("X2")
18 plt.text(-2, 2.8, "'THA076BCT026',",
19           'THA076BCT027','THA076BCT041'",
20           fontsize=8,color='red')
21 plt.title('Scatter Plot for Random Value')
22 plt.ylim(-3,3)
23 plt.xlim(-3,3)
24 plt.savefig('Scatter Plot for Random
25             Values', bbox_inches='tight')
26
27
28 # Mean and standard deviation
29 mean1 = np.mean(x1)
30 mean2 = np.mean(x2)
31 std1 = np.std(x1)
32 std2 = np.std(x2)
33 print(mean1, std1)
34 print(mean2, std2)
35
36
37 # 2x2 random matrix with elements sampled
38 # from a uniform distribution (range
39 # [0, 1]).
40 random_matrix = np.random.rand(2,2)
41 print(random_matrix)
42
43
44 # Converts the lists x1 and x2 into NumPy
45 # arrays, then creates a data matrix
46 # by column stacking the arrays.
47 x1_array = np.array(x1)
48 x2_array = np.array(x2)
49 data_matrix = np.column_stack((x1_array,
50                               x2_array))
51 print(data_matrix.shape)
52 print(data_matrix)
53
54
55
56
57
58
59
60 # **STEP 1: Reshape the data by
61 # multiplication with a random matrix
62 # .** To visualize the linear
63 # transformation, click on the link :
64 # https://www.geogebra.org/m/YCZA8TAH?
65 # fbclid=IwAR3LK2VOQooxb-qzH-
66 # LlhlyzVEhi8Qsb8Of_VlzeiqWb3FcjYibVPb
67
68 transformed_data_by_random_matrix =
69     data_matrix @ random_matrix
70 print(transformed_data_by_random_matrix.
71       shape)

```

```

44 # first column [ :, 0 ] as x-axis value .
        similarly 2nd column as y-axis
45 plt.scatter(
    transformed_data_by_random_matrix
    [ :, 0 ],
    transformed_data_by_random_matrix
    [ :, 1 ])
46 plt.xlabel("X1-transformed")
47 plt.ylabel("X2-transformed")
48 plt.ylim(-3,3)
49 plt.xlim(-3,3)
50 plt.text(-1.5, 2.8, "'THA076BCT026', "
      "'THA076BCT027', 'THA076BCT041'",
      fontsize=8,color='red')
51 plt.title('Scatter Plot for transformed
Random Values')
52 plt.savefig('transformed Random Values',
bbox_inches='tight')
53
54
55 # **STEP 2: Calculate covariance matrix**
56     The covariance matrix provides
      insights into the relationships
      between different dimensions of the
      transformed_data_by_random_matrix.
57 # See if elements in main diagonal are
      max and other diagonal are min. If
      not, we have to proceed further .
58 covariance_matrix = np.cov(
      transformed_data_by_random_matrix.T)
59 covariance_matrix
#same as above
60 covariance_matrix = np.cov(np.transpose(
      transformed_data_by_random_matrix))
61 print(covariance_matrix)

62
63
64 # **STEP 3: Eigen values and vector
      calculation**
65 # The eigenvalues represent the variance
      explained by each principal component
      (magnitude of new feature space)
66 # The eigenvectors define the directions
      of the principal components. ((
      direction of new feature space)
67 eigen_values, eigen_vectors = np.linalg.
      eig(covariance_matrix)
68 print(eigen_values)
69 print(eigen_vectors)

70
71 # Just To check if the eigenvalues really
      represent the variance, lets find
      out the proportion of variance for
      each components
72 proportion1 = eigen_values[0] / np.sum(
      eigen_values)
73 print(proportion1*100)
74 proportion2 = eigen_values[1] / np.sum(
      eigen_values)
75 print(proportion2*100)

76
77
78 # **STEP 4: Transform the data with the
      selected eigen vector.**
79 #first write the eigenvectors row wise (
      do the transpose, from column wise to
      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
# Set the size of the figure
54 plt.figure(figsize=(12, 7))
55 df['target'] = df['target'].replace({0: 'Setosa', 1: 'Versicolor', 2: 'Virginica'})
56
57 # Create a pair plot with hue based on
# the 'target' column in the dataframe
58 g=sns.pairplot(df, hue='target',plot_kws
# ={'legend': False})
59 g._legend.set_title('')
60 # Set the title of the plot
61 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             )
```

```

    medium medium best components)
new_covariance_matrices_dictionary[5]

251 #Lets access combination 7, which is
252     eigen[0],eigen[1],eigen[2] and eigen
253         [3]          (all components combined)
new_covariance_matrices_dictionary[7]

254 # Heatmap plot, darker shade indicate
255     higher values. (combination starts
256         from 0 to 7)
257 plt.imshow(
    new_covariance_matrices_dictionary
    [1], cmap='Reds', interpolation='
        nearest') #for smoother
    interpolation of color on pixel
    values, instead of 'nearest', we can
    use 'bilinear','bicubic','spline16',
    'spline36', 'hanning', 'hamming', '
    hermite',etc..
258 plt.colorbar()
259 plt.title("New Covariance Matrix
    Combination 1 (2 best PC)")
260 plt.text(-0.3, -0.66, "'THA076BCT026', ''
    THA076BCT027', 'THA076BCT041'",'
    fontsize=8,color='red')
261 plt.savefig('Covariance Matrix for 2 best
    PC', bbox_inches='tight')
262 plt.show()

263
264 plt.imshow(
    new_covariance_matrices_dictionary
    [3], cmap='Reds', interpolation='
        nearest')
265 plt.colorbar()
266 plt.title("New Covariance Matrix
    Combination 3 (best and worst PC)")
267 plt.text(-0.3, -0.66, "'THA076BCT026', ''
    THA076BCT027', 'THA076BCT041'",'
    fontsize=8,color='red')
268 plt.savefig('Covariance Matrix for best
    and worst PC', bbox_inches='tight')
269 plt.show()

270
271 plt.imshow(
    new_covariance_matrices_dictionary
    [5], cmap='Reds', interpolation='
        nearest')
272 plt.colorbar()
273 plt.title("Covariance Matrix for two
    medium PCs")
274 plt.text(-0.3, -0.66, "'THA076BCT026', ''
    THA076BCT027', 'THA076BCT041'",'
    fontsize=8,color='red')
275 plt.savefig('Covariance Matrix for two
    medium PCs', bbox_inches='tight')
276 plt.show()

277
278 plt.imshow(
    new_covariance_matrices_dictionary
    [7], cmap='Reds', interpolation='
        nearest')
279 plt.colorbar()
280 plt.title("New Covariance Matrix
    Combination 7 (All Components)")
281 plt.show()

```

```

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={ 'size': 5})
108 plt.title('Covariance Matrix on Processed
109     Cancer Data')
110 plt.text(8, -1.2, "'THA076BCT026',",
111         "'THA076BCT027', 'THA076BCT041'",,
112         fontsize=8,color='red')
113 plt.savefig('./plots/Covariance Matrix
114     Processed Cancer Data', bbox_inches='
115     tight')
116 plt.show()
117
118 #### **STEP 4: Eigen vectors/values
119 # calculation**
120 eigenvalues, eigenvectors = np.linalg.eig
121     (covariance_matrix)
122 print("shape of eigenvalue:",eigenvalues.
123     shape)
124 print("shape of eigen vectors:",,
125     eigenvectors.shape)
126 print("eigen values are",eigenvalues)
127 print("eigen vectors are",eigenvectors)
128
129 # **STEP 5: See the variance explained by
130 # each eigen values.**
131 # 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
   #
   Combination:
160     print("*Covariance matrix shape for
       combination",i,"is",
       new_covariance_matrix.shape)
   # 0=best PC, 1=2 best PCs, 2= best+
   medium PCs
161     print("-->Covariance matrix:",
       new_covariance_matrices_dictionary[
       i])                                # 3=
   medium+medium PCs, 4= best+worst
   PCs
162
163 # Heatmap to visualize covariance matrix
   combination.
164 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..
165 plt.colorbar()
166 plt.title("New Covariance Matrix
   Combination 1 (2 best PCs)")
167 plt.show()
168
169 #lets try bilinear,pixel coloring is
   smoother(linear interpolation between
   4 nearest data points)
170 plt.imshow(
   new_covariance_matrices_dictionary
   [2], cmap='Reds', interpolation='
   bilinear')
171 plt.colorbar()
172 plt.title("New Covariance Matrix
   Combination 2 (Best and medium PCs)")
173 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.
```

• • •