

T.R.
GEBZE TECHNICAL UNIVERSITY
FACULTY OF ENGINEERING
DEPARTMENT OF COMPUTER ENGINEERING

CSE 555 - STATISTICAL DATA ANALYSIS
FINAL PROJECT

EMIRKAN BURAK YILMAZ

SUPERVISOR
HABIL KALKAN

GEBZE
2025

T.R.
GEBZE TECHNICAL UNIVERSITY
FACULTY OF ENGINEERING
COMPUTER ENGINEERING DEPARTMENT

CSE 555 - STATISTICAL DATA ANALYSIS
FINAL PROJECT

EMIRKAN BURAK YILMAZ

SUPERVISOR
HABIL KALKAN

2025
GEBZE

ABSTRACT

This study provides an in-depth analysis of the Dry Bean dataset, applying a combination of statistical methods, machine learning, and dimensionality reduction techniques to investigate feature discriminability, clustering behavior, and the impact of outliers. The analysis begins with exploratory data examination, identifying key geometric features such as EquivDiameter, Perimeter, and MinorAxisLength as the most discriminative, based on Fisher Score assessments. Dimensionality reduction comparisons show that Principal Component Analysis (PCA) outperforms Linear Discriminant Analysis (LDA) in separating classes (silhouette scores of 0.3347 vs. 0.2105), indicating better preservation of non-linear patterns. Four clustering methods (K-Means, DBSCAN, t-SNE, and Self-Organizing Maps) are evaluated, with t-SNE yielding the most distinct visual cluster separation. Additionally, applying Isolation Forest for outlier removal improves both feature discriminability (e.g., EquivDiameter Fisher Score rising from 3.2778 to 3.6325) and clustering performance (silhouette score increasing from 0.2565 to 0.2781). Overall, the results emphasize the critical role of geometric features in classifying bean types, the strengths of unsupervised learning approaches for this dataset, and the benefits of outlier detection during preprocessing.

Keywords: Machine Learning, Dry Bean Dataset, Feature Selection, Dimensionality Reduction, Clustering Analysis, Outlier Detection.

LIST OF SYMBOLS AND ABBREVIATIONS

Symbol or

Abbreviation : Explanation

DBSCAN : Density-based Spatial Clustering of Applications with Noise

LDA : Linear Discriminant Analysis

PCA : Principal Component Analysis

SOM : Self Organizing Map

t-SNE : t-distributed Stochastic Neighbor Embedding

CONTENTS

Abstract	iii
List of Symbols and Abbreviations	iv
Contents	v
List of Figures	vi
List of Tables	vii
1 Introduction	1
2 Analysis and Results	2
2.1 Boxplot Analysis of Attributes	2
2.2 Z-Score Normalization and Fisher Distance	3
2.3 Principal Component Analysis (PCA)	4
2.4 Scatter Plot of Principal Components	6
2.5 Linear Discriminant Analysis (LDA)	7
2.6 Clustering Algorithms	8
2.6.1 K-Means	8
2.6.2 DBSCAN	9
2.6.3 t-SNE	10
2.6.4 SOM	11
2.7 Outlier Detection	14
3 Conclusion	17
Bibliography	19

LIST OF FIGURES

2.1	Boxplot of the attributes	2
2.2	Scatter plot using top 2 and bottom 2 PCA components	6
2.3	LDA vs PCA	7
2.4	Clustering result using K-Means (PCA-reduced data)	9
2.5	Clustering result using DBSCAN (PCA-reduced data)	10
2.6	t-SNE 2D visualization	11
2.7	SOM clusters	12
2.8	SOM neuron frequencies	13
2.9	PCA plot after outlier removal	16

LIST OF TABLES

2.1	Fisher Scores of Features (Descending Order)	4
2.2	Fisher Scores and Explained Variance of PCA Components	5
2.3	Average Silhouette Scores for PCA and LDA Projections	8
2.4	Fisher Scores Before and After Outlier Removal Using Isolation Forest	15
2.5	Silhouette Scores Before and After Outlier Removal	15

1. INTRODUCTION

This project explores the application of statistical and machine learning techniques to the Dry Bean dataset, which contains samples from seven distinct types of beans, each described by sixteen measurable features.

The main goal of this work is to investigate and analyze the dataset through a structured pipeline that includes data normalization, dimensionality reduction, clustering, and outlier detection. Techniques such as Principal Component Analysis (PCA), Linear Discriminant Analysis (LDA), and various clustering algorithms (including K-Means, DBSCAN, t-SNE, and Self-Organizing Maps) are applied to assess how effectively they can cluster the Dry Bean dataset.

The report is organized as follows: Chapter 2 presents the methods and results, structured step-by-step according to the project tasks. Chapter 3 summarizes the key findings and conclusions. All code used throughout the experiments is provided in the appendices.

2. ANALYSIS AND RESULTS

This chapter presents the statistical analysis steps applied to the Dry Bean dataset. Each task specified in the project instructions is addressed in a separate section, with corresponding visualizations and interpretations.

2.1. Boxplot Analysis of Attributes

The distribution of each feature was visualized using boxplots, as illustrated in Figure 2.1. This analysis allowed us to identify skewness, variability, and potential outliers.

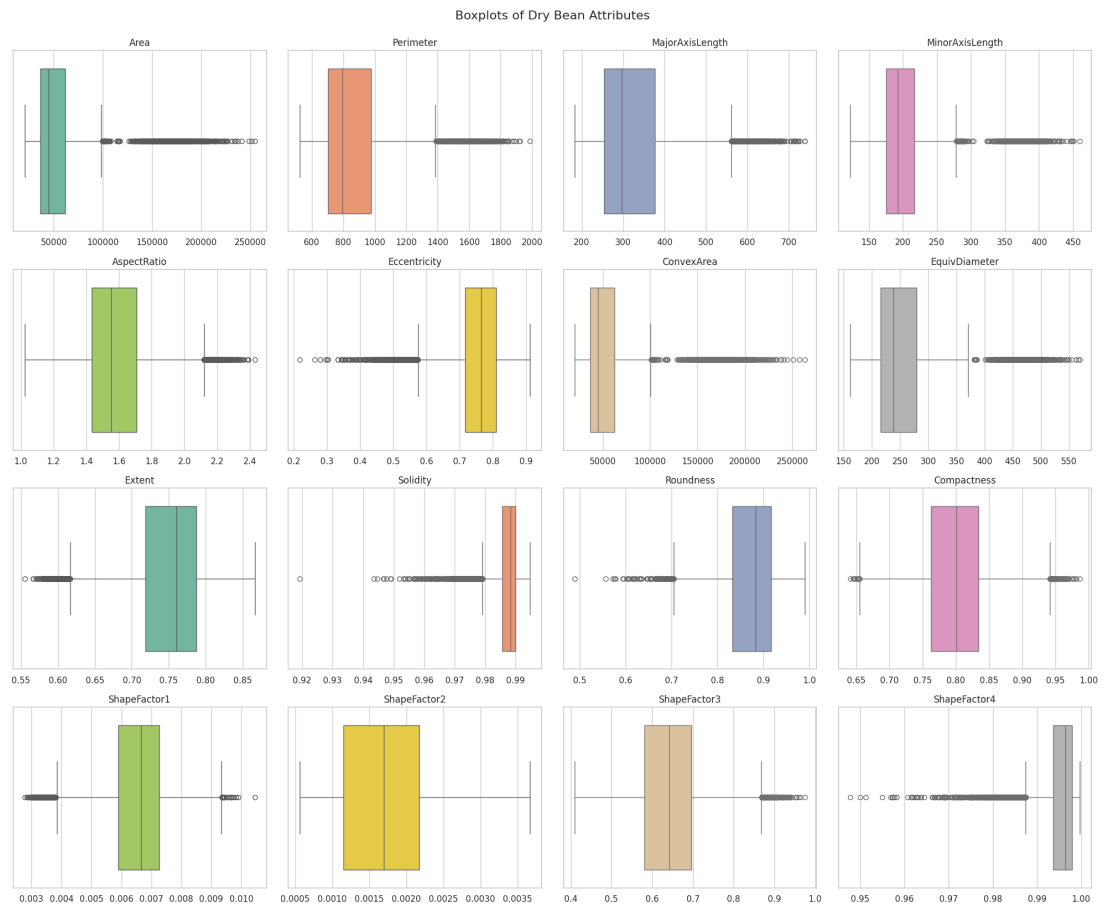


Figure 2.1: Boxplot of the attributes

At first glance, all features except ShapeFactor2 contain numerous outliers. These

outliers are mostly located at the upper end of the range, particularly in area-related features such as AspectRatio, ConvexArea, and EquivDiameter, while for features like Eccentricity, Extent, Solidity, and Roundness, the outliers are more pronounced at the lower end.

All features exhibit a narrow interquartile range (IQR), indicating that the majority of the data points are densely clustered, with relatively few extreme values.

Regarding the distribution patterns, many features display approximately normal distributions but with varying means and variances. Some features show noticeable skewness either to the left or right, suggesting asymmetric data spread and the presence of tails on one side of the distribution.

These insights provide an initial understanding of the dataset's variability, symmetry, and potential anomalies, which are important when selecting and applying further preprocessing and modeling techniques.

2.2. Z-Score Normalization and Fisher Distance

To begin, the data were normalized using Z-score normalization, ensuring that each feature has a mean of zero and a variance of one. This standardization step is essential because the Fisher Distance metric is sensitive to the scale of the input attributes. Without such normalization, features with larger scales could disproportionately influence the distance calculations.

Given that the dataset comprises seven distinct classes, the Fisher Distance was adapted to handle this multi-class setting by employing the generalized Fisher Distance approach proposed by Li et al. (2017) [1]. While the classical Fisher Distance is designed for two-class problem, the generalized method extends this framework to multi-class problems by computing the average of all pairwise Fisher Distances across class pairs. This approach allows for a comprehensive assessment of class separability in multi-class datasets. The generalized Fisher Distance FD_j for feature j is computed as:

$$FD_j = \frac{2}{C(C-1)} \sum_{i=1}^C \sum_{k=i+1}^C \frac{|\mu_{ij} - \mu_{kj}|}{\sqrt{\sigma_{ij}^2 + \sigma_{kj}^2}}$$

where:

- C : number of classes
- μ_{ij} : mean of feature j in class i

- σ_{ij}^2 : variance of feature j in class i

Table 2.1: Fisher Scores of Features (Descending Order)

Feature	Fisher Score
EquivDiameter	3.2778
Perimeter	3.2119
MinorAxisLength	3.1344
ShapeFactor1	3.1215
MajorAxisLength	2.9482
Area	2.9035
ConvexArea	2.9027
ShapeFactor2	2.3974
Compactness	1.8473
ShapeFactor3	1.8223
AspectRatio	1.8027
Eccentricity	1.5937
Roundness	1.5048
ShapeFactor4	0.7290
Solidity	0.5141
Extent	0.3750

The results in Table 2.1 indicate that EquivDiameter, Perimeter, and MinorAxisLength have the highest Fisher Scores, suggesting they are the most discriminative features for class separation in this dataset. These geometric features likely capture significant shape or size variations across classes, which contribute to their high class separability. In contrast, features like ShapeFactor4, Solidity, and Extent exhibit very low Fisher Scores, implying they contribute minimally to class discrimination.

2.3. Principal Component Analysis (PCA)

This section presents the application of PCA to the dataset [2]. Each principal component was treated as a new feature, and its Fisher Distance was calculated to assess class separability. The results, summarized in Table 2.2, reveal the relationship between the explained variance ratio and the Fisher Score for each component.

Table 2.2: Fisher Scores and Explained Variance of PCA Components

Principal Component	Explained Variance Ratio	Fisher Score
PC1	0.5547	2.8558
PC2	0.2643	2.4563
PC3	0.0801	1.1114
PC4	0.0511	0.7795
PC5	0.0274	0.7703
PC6	0.0115	0.4922
PC7	0.0070	0.4819
PC8	0.0033	0.4768
PC9	0.0005	0.3897
PC10	0.0001	0.3297
PC11	0.0001	0.2757
PC12	0.0000	0.2565
PC13	0.0000	0.2350
PC14	0.0000	0.1734
PC15	0.0000	0.1302
PC16	0.0000	0.0754

The first principal component (PC1), which accounts for approximately 55.47% of the total variance, achieved the highest Fisher Score of 2.86, indicating that it captures both substantial variability in the data and a meaningful level of class-discriminative information. Similarly, PC2, which explains 26.43% of the variance, attained a Fisher Score of 2.46, further underscoring the discriminative relevance of the leading components.

Beyond PC2, a sharp decline is observed in both explained variance and Fisher Scores. For instance, PC3 and PC4 account for 8.01% and 5.11% of the variance, respectively, with corresponding Fisher Scores of 1.11 and 0.78. This trend continues, with later components contributing progressively less to both variance and class separation. Notably, components PC10 through PC16 contribute no variance and exhibit very low Fisher Scores suggesting negligible utility in distinguishing between classes.

The results in Table 2.2 further reveal that the Fisher Scores for PCA-transformed features are generally lower than those computed for the original features, as shown in Table 2.1. While PCA is effective for dimensionality reduction and data compression, it does not inherently optimize for class separability. As an unsupervised technique, PCA seeks to maximize variance without accounting for class labels, which can result

in components that do not capture class-specific structure effectively.

2.4. Scatter Plot of Principal Components

To further assess the effectiveness of the principal components in capturing data structure and enabling class discrimination, two scatter plots were analyzed: one displaying the projection of the dataset onto the first two principal components (PC1 and PC2) and the other onto the last two components (PC15 and PC16). These visualizations are presented in Figure 2.2.

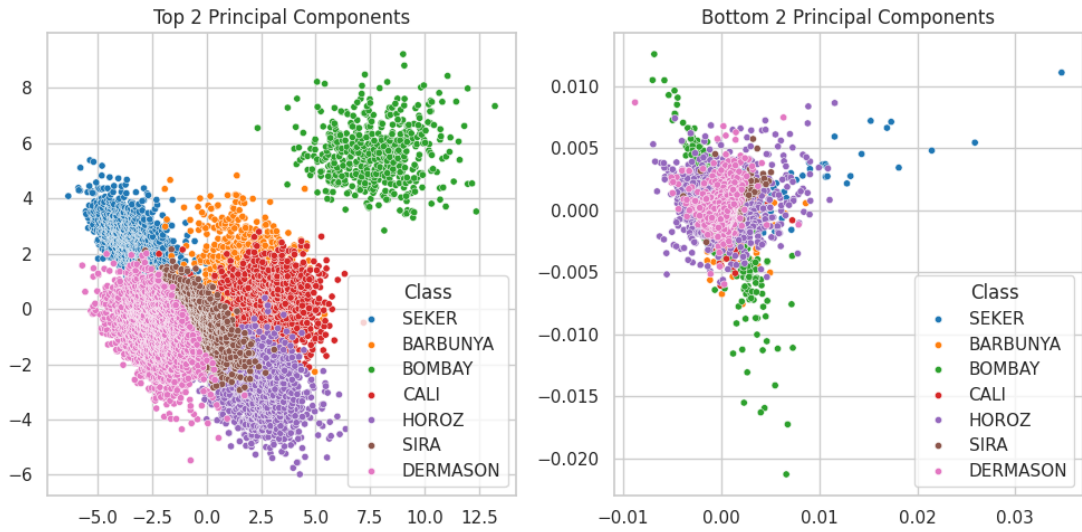


Figure 2.2: Scatter plot using top 2 and bottom 2 PCA components

The projection onto PC1 and PC2 captures approximately 80.9% of the total variance (55.47% and 26.43%, respectively), indicating that these components retain the majority of the original data's structural information. The scatter plot reveals distinct clustering and directional separation among different bean classes, suggesting a high degree of class separability in this reduced two-dimensional space.

These visual observations are consistent with the numerical results reported in Table 2.2, where PC1 and PC2 achieved the highest Fisher Scores of 2.86 and 2.46, respectively. The combination of high variance capture and strong class discrimination confirms the value of these components for both data representation and supervised learning tasks.

In contrast, the projection onto PC15 and PC16 explains virtually no variance in the dataset. Correspondingly, the Fisher Scores for these components are minimal (0.13 for

PC15 and 0.075 for PC16) indicating that these dimensions carry negligible discriminative information. The scatter plot illustrates this, displaying a highly overlapping and noisy distribution with no discernible class structure, thereby reinforcing the conclusion that these components contribute little to classification performance.

2.5. Linear Discriminant Analysis (LDA)

Linear Discriminant Analysis (LDA) is a supervised dimensionality reduction technique designed to maximize class separability by projecting data along directions that increase between-class variance and reduce within-class variance [3]. In contrast, PCA is an unsupervised method that identifies directions of maximum overall variance without considering class labels. LDA was applied to reduce the dataset to two dimensions, and the resulting projection is visualized in Figure 2.3, alongside the PCA projection from Section 2.4 for comparison.

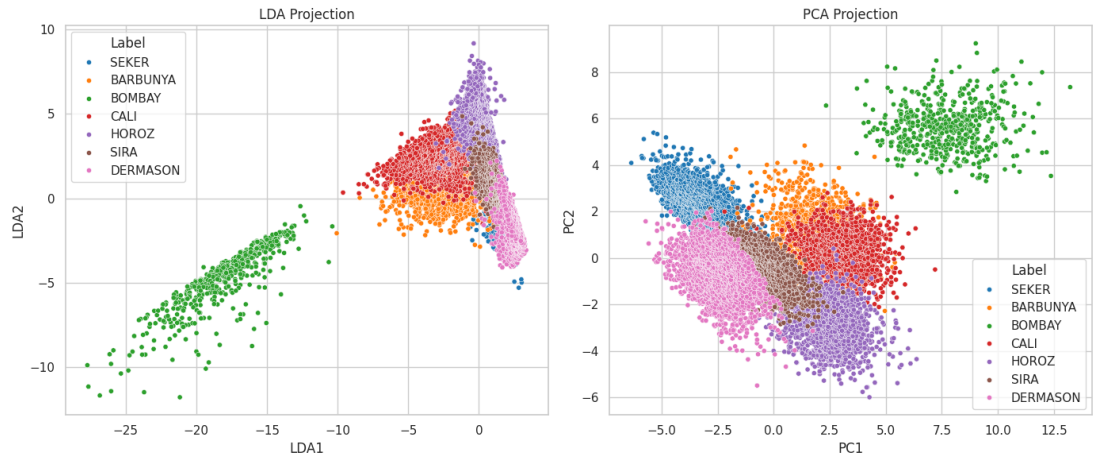


Figure 2.3: LDA vs PCA

Visually, the BOMBAY class (shown in green in the plot) is clearly distinguishable in both projections. However, the PCA projection (right plot) shows better separation among several major classes compared to the LDA projection (left plot), where multiple classes remain tightly clustered with significant overlap.

To quantify the quality of class separation, the silhouette score was used [4]. This score measures how similar each point is to its own cluster compared to other clusters, with values ranging from -1 (poor separation) to $+1$ (well-separated clusters). As shown in Table 2.3, the PCA projection achieved a higher average silhouette score (0.3347) compared to the LDA projection (0.2105), indicating better class separation

with PCA.

Table 2.3: Average Silhouette Scores for PCA and LDA Projections

Projection	Silhouette Score
PCA	0.3347
LDA	0.2105

The higher silhouette score for PCA (0.3347) compared to LDA (0.2105) can be explained by several factors. First, the Dry Bean dataset likely contains non-linear relationships between features. LDA, which is optimized for linear class separation, may struggle when classes are not perfectly linearly separable. In contrast, PCA is an unsupervised method that focuses on maximizing variance, capturing the overall data structure, even if it's not perfectly aligned with class boundaries.

Additionally, LDA is more sensitive to noise. If the dataset has overlapping classes or irrelevant features, this can hinder LDA's performance. In contrast, PCA is less sensitive to noise and focuses on retaining the overall variance in the data, which may contribute to better class separation in some cases.

In conclusion, while LDA is designed for class separability, PCA performed better in this dataset, likely due to the non-linear nature of the data, the dimensionality constraints of LDA, and its sensitivity to noise.

2.6. Clustering Algorithms

In this section, the Dry Bean dataset is grouped using four different clustering algorithms: K-Means, DBSCAN, t-SNE, and SOM. The goal is to explore the structure of the data through unsupervised learning and visualize the results. For the K-Means and DBSCAN clustering algorithms, the first two principal components obtained from the PCA transformation are used as the reduced feature set for clustering. On the other hand, t-SNE and SOM are applied directly to the original dataset without any dimensionality reduction.

2.6.1. K-Means

The K-Means algorithm partitions the data into k clusters by minimizing the within-cluster variance [5]. In this analysis, the number of clusters was set to $k = 7$, based

on prior knowledge of the dataset, specifically reflecting the known number of distinct classes in the Dry Bean dataset. Figure 2.4 illustrates the clustering results obtained using the K-Means algorithm with a comparison to the original classes, applied to the first two principal components extracted through PCA.

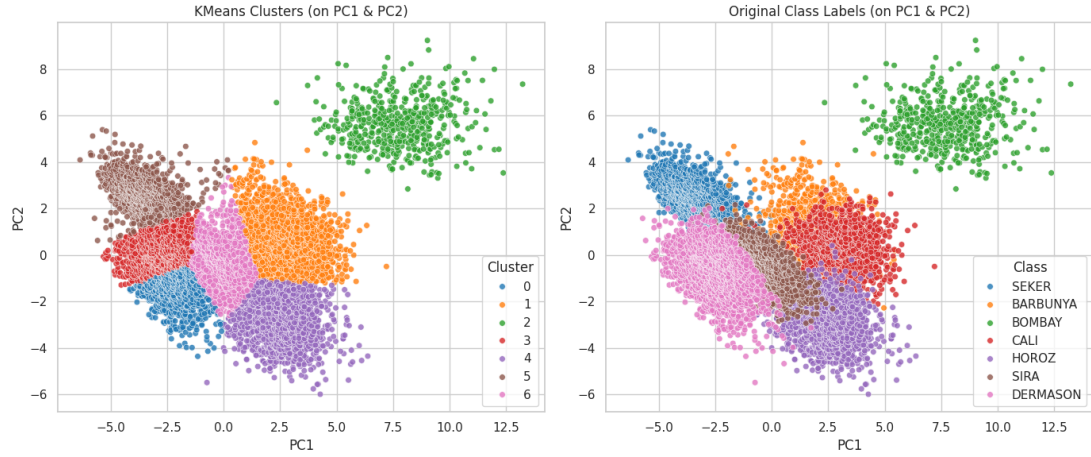


Figure 2.4: Clustering result using K-Means (PCA-reduced data)

The figure demonstrates that K-Means separates the data into distinct clusters; however, the proximity of several cluster boundaries suggests that the algorithm struggles to differentiate certain classes, particularly in our case where the clusters are not well-separated in the reduced PCA space. This observation highlights the sensitivity of the algorithm to the chosen value of k . Without prior knowledge of the actual number of classes, selecting a smaller k such as 2 or 3 may have been more appropriate, as the current configuration shows clusters positioned so closely that minor noise could easily cause a data point to shift from one cluster to another.

2.6.2. DBSCAN

Unlike K-Means, DBSCAN (Density-Based Spatial Clustering of Applications with Noise) does not require a predefined number of clusters and can detect clusters of varying shapes and densities [6]. In Figure 2.5, the results of the DBSCAN algorithm illustrates the identified two primary clusters, as indicated by the clear separation in the plot. This result suggests that DBSCAN is more flexible than K-Means, as it can identify clusters without prior knowledge of number of classes and also detect noise. DBSCAN's ability to find dense regions and mark the rest as noise shows its advantage when dealing with unevenly distributed data.

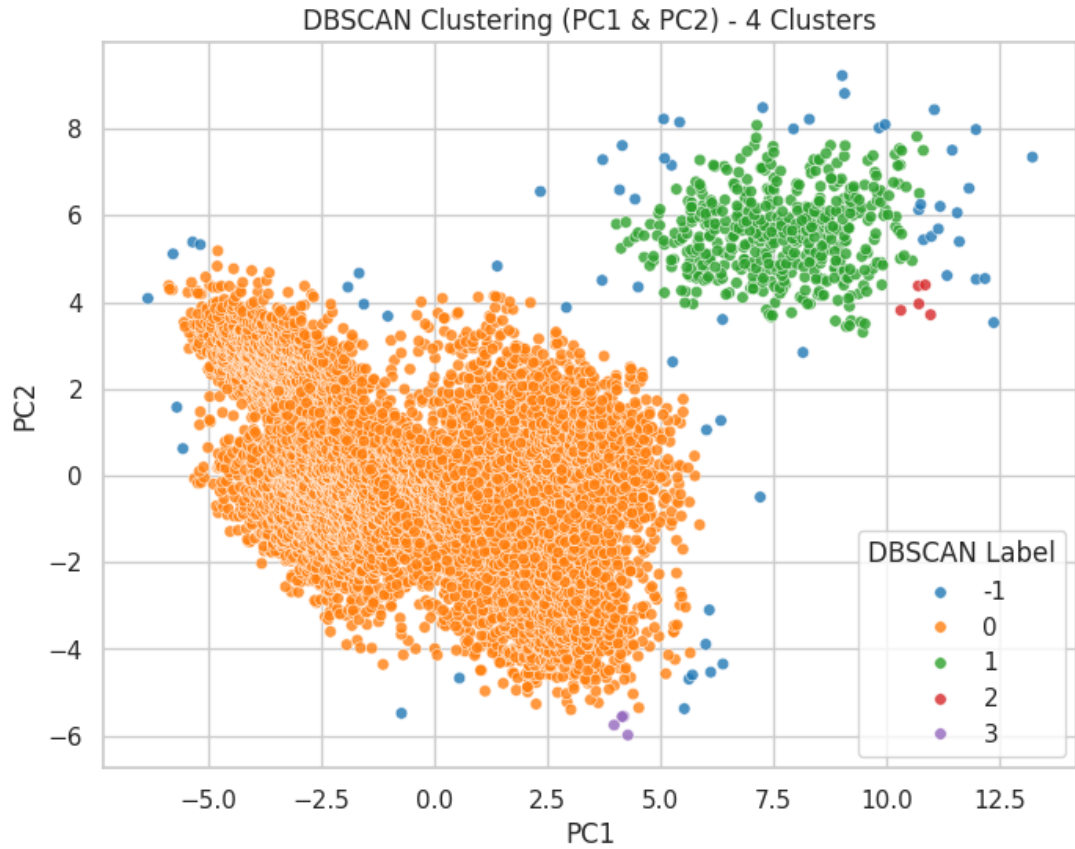


Figure 2.5: Clustering result using DBSCAN (PCA-reduced data)

2.6.3. t-SNE

t-SNE (t-distributed Stochastic Neighbor Embedding) is a dimensionality reduction technique that is particularly useful for visualizing high-dimensional data in lower dimensions [7]. The algorithm works by preserving local data structures, making it easier to identify clusters within the data. In this figure, the clusters are better separated compared to the PCA-based result. Figure 2.6 shows clearer boundaries between most of the classes compared to PCA, suggesting that t-SNE is effective in uncovering the structure of the data that may not be as apparent in higher dimensions.

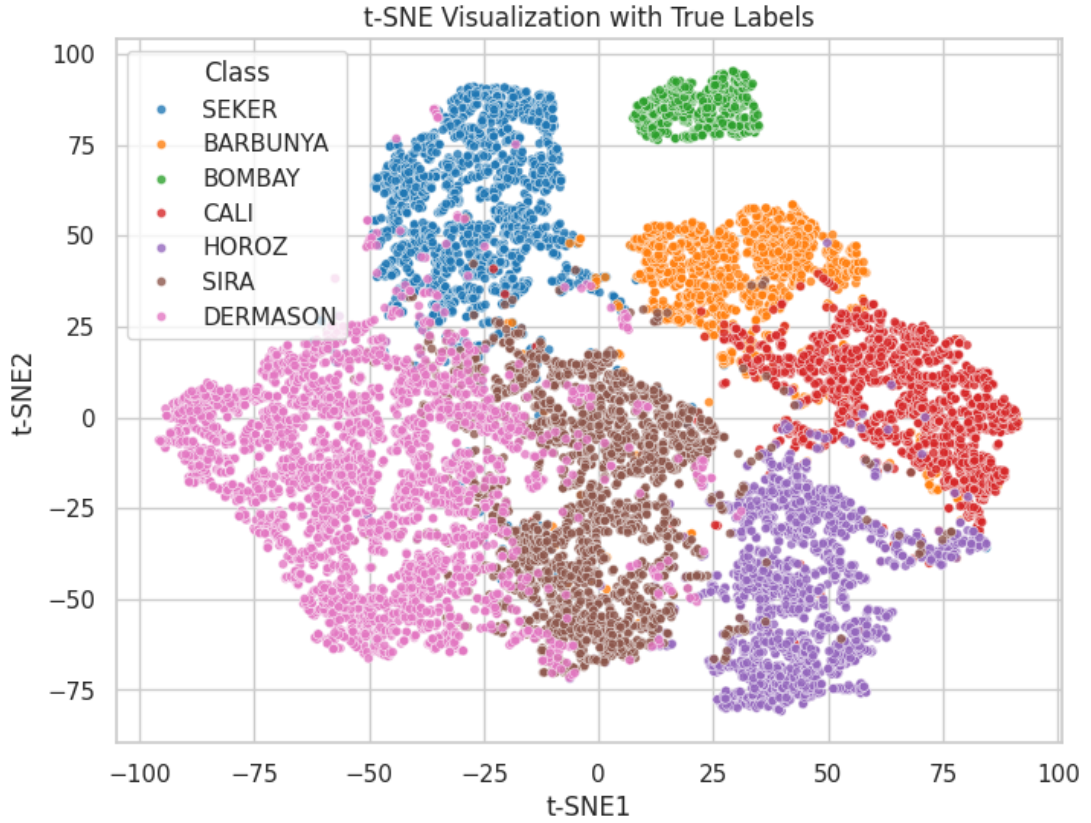


Figure 2.6: t-SNE 2D visualization

2.6.4. SOM

SOM (Self-Organizing Map) is an unsupervised artificial neural network that projects high-dimensional data onto a typically two-dimensional grid, preserving the topological relationships among data points [8]. During training, the SOM adjusts its neurons to represent the input space, effectively grouping similar instances into neighboring grid cells. In this case, a 10×10 grid was used to map the Dry Bean dataset.

Figure 2.7 reveals that similar data points are clustered around certain grid regions, with some neurons capturing the boundaries between different classes. While the SOM effectively identifies general class groupings, some grid cells contain data points from multiple classes, indicating classification or boundary errors. This is partly due to the competitive and neighborhood-based learning process of SOMs, where neighboring neurons influence each other's updates.

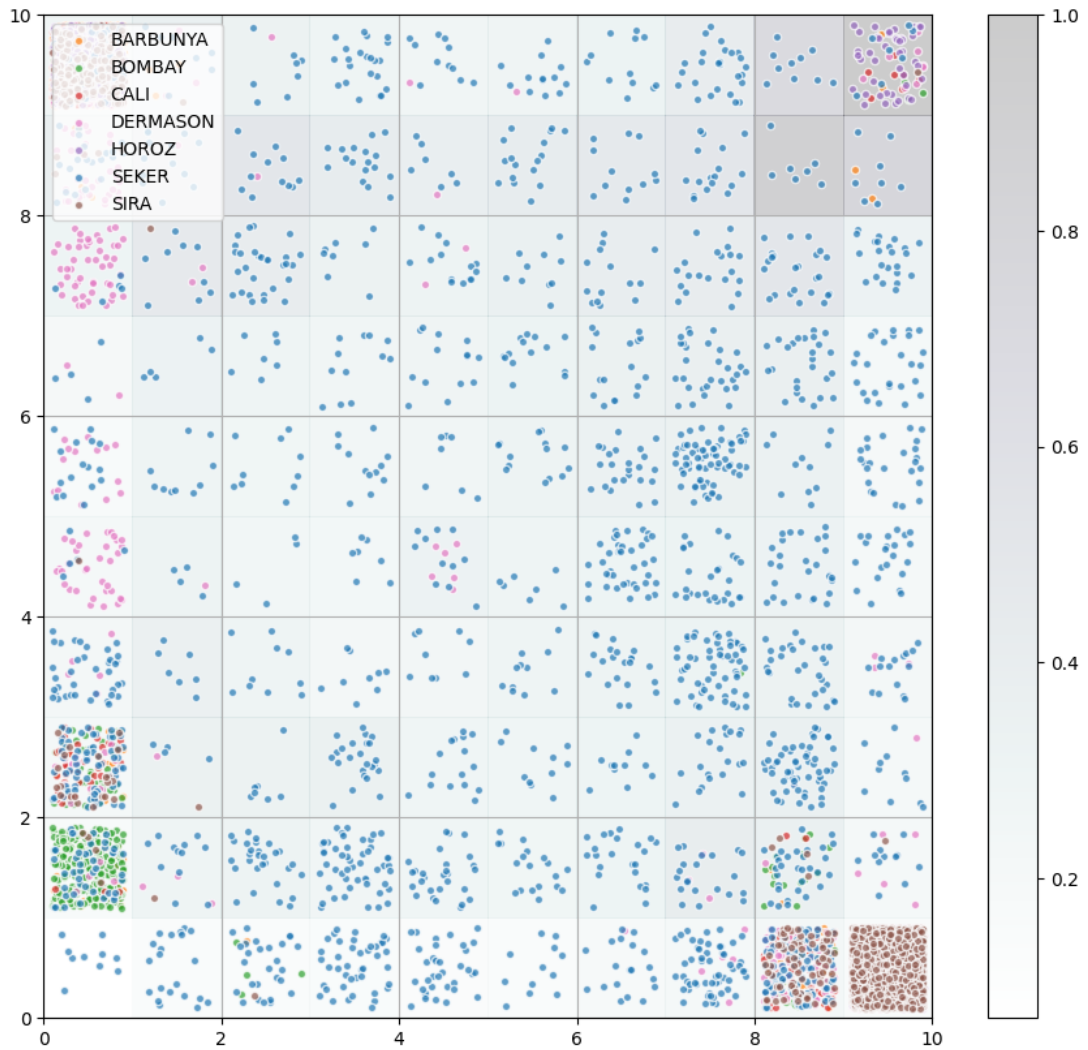


Figure 2.7: SOM clusters

The SOM algorithm is particularly valuable for high-dimensional datasets because it provides a non-linear, topology-preserving mapping, allowing to visualize complex patterns and relationships that may not be easily captured through linear techniques like PCA. However, it is important to note that the interpretability of SOM results can depend on factors such as grid size, training iterations, and neighborhood radius, which influence how well the map captures the true class structure.

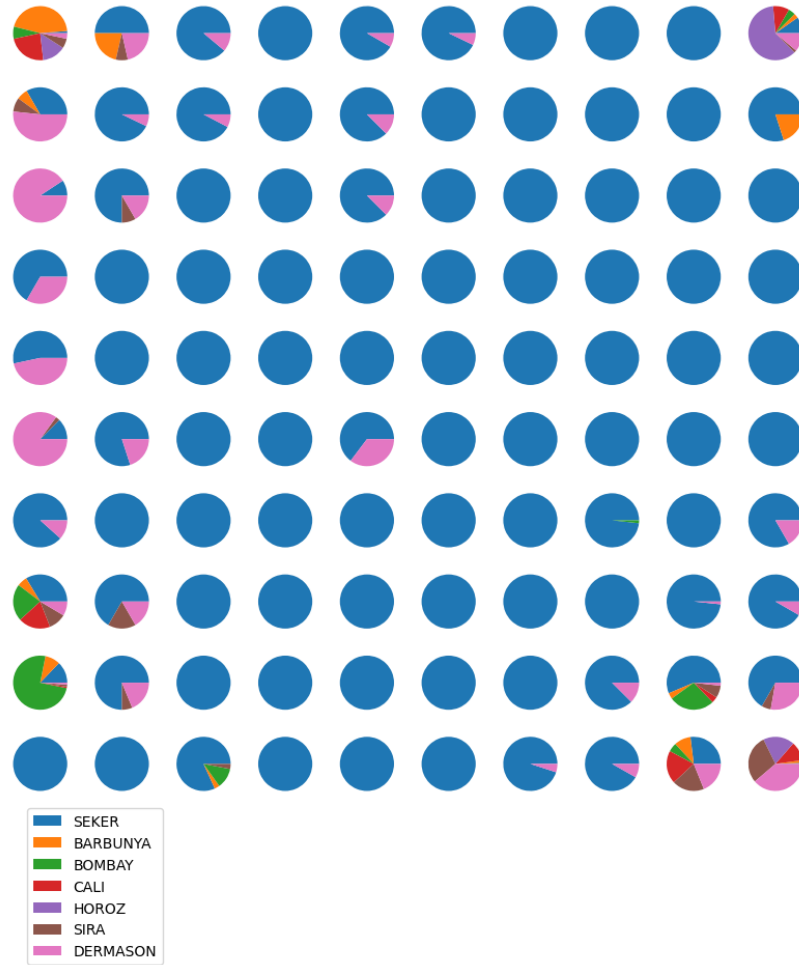


Figure 2.8: SOM neuron frequencies

Figure 2.8 presents SOM visualization that depicts the distribution of classes across the grid's neurons. In this map, each neuron is visualized using a pie-chart-like color scheme, where the proportion of colors reflects the percentage of different classes assigned to that neuron. The visualization reveals that many neurons are entirely dominated by the SEKER class (shown in blue), while several other neurons exhibit mixed-class compositions, highlighting regions where the algorithm clusters similar but not perfectly separable samples. This outcome is consistent with observations from the PCA and t-SNE analyses, where certain class boundaries were also found to be ambiguous or overlapping. The SOM's clustering behavior here reflects the inherent complexity of the dataset and the presence of confusing or borderline examples, underscoring a common challenge faced across multiple clustering techniques.

2.7. Outlier Detection

For outlier detection and noise cleaning in the Dry Bean dataset, the Isolation Forest algorithm was chosen. Isolation Forest is a tree-based ensemble method designed specifically for anomaly detection. It isolates anomalies instead of profiling normal data points, which makes it efficient for high-dimensional datasets like ours with 16 features [9]. Additionally, Isolation Forest does not assume any underlying data distribution, making it suitable for datasets with complex, non-linear feature interactions. Its ability to handle large datasets efficiently and its effectiveness in identifying anomalies with few hyperparameters make it a practical choice for our cleaning task.

The Isolation Forest algorithm was applied to the original dataset with a contamination parameter set to 0.01, meaning that approximately 1% of the data points were flagged as outliers. These detected outliers were then removed from the dataset before further analysis. The goal was to reduce noise and improve the quality of clustering and class separability without removing too many valid samples.

After removing the identified outliers, the Fisher scores were recalculated to assess the effect on feature discriminability. Table 2.4 summarizes the Fisher scores before and after outlier removal.

Table 2.4: Fisher Scores Before and After Outlier Removal Using Isolation Forest

Feature	Fisher Score Before	Fisher Score After
EquivDiameter	3.2778	3.6325
Perimeter	3.2119	3.5138
MinorAxisLength	3.1344	3.4840
ConvexArea	2.9027	3.3320
Area	2.9035	3.3265
ShapeFactor1	3.1215	3.1560
MajorAxisLength	2.9482	3.1193
ShapeFactor2	2.3974	2.3655
Compactness	1.8473	1.8491
ShapeFactor3	1.8223	1.8237
AspectRatio	1.8027	1.8069
Eccentricity	1.5937	1.5960
Roundness	1.5048	1.5145
ShapeFactor4	0.7290	0.7310
Solidity	0.5141	0.5235
Extent	0.3750	0.3847

Notably, the most discriminative features, such as EquivDiameter, Perimeter, and MinorAxisLength, showed a marked increase in Fisher scores after cleaning, indicating improved class separability. This suggests that outlier removal reduced the within-class variance or clarified the differences between classes.

Table 2.5: Silhouette Scores Before and After Outlier Removal

Condition	Silhouette Score
Before Outlier Removal	0.2565
After Outlier Removal	0.2781

Furthermore, the silhouette score for clustering increased from 0.2565 before cleaning to 0.2781 after cleaning, as presented in Table 2.5. Although the increase is moderate, it indicates an improvement in clustering quality, attributed to the reduction of noise and anomalous samples that previously obscured cluster boundaries.

To visualize the overall structure of the dataset and the effect of outlier removal, PCA was performed, reducing the 16-dimensional feature space to two principal components.

Figure 2.9 shows the 2D PCA scatter plot of the cleaned dataset, revealing clearer class groupings and less noise compared to the original data.

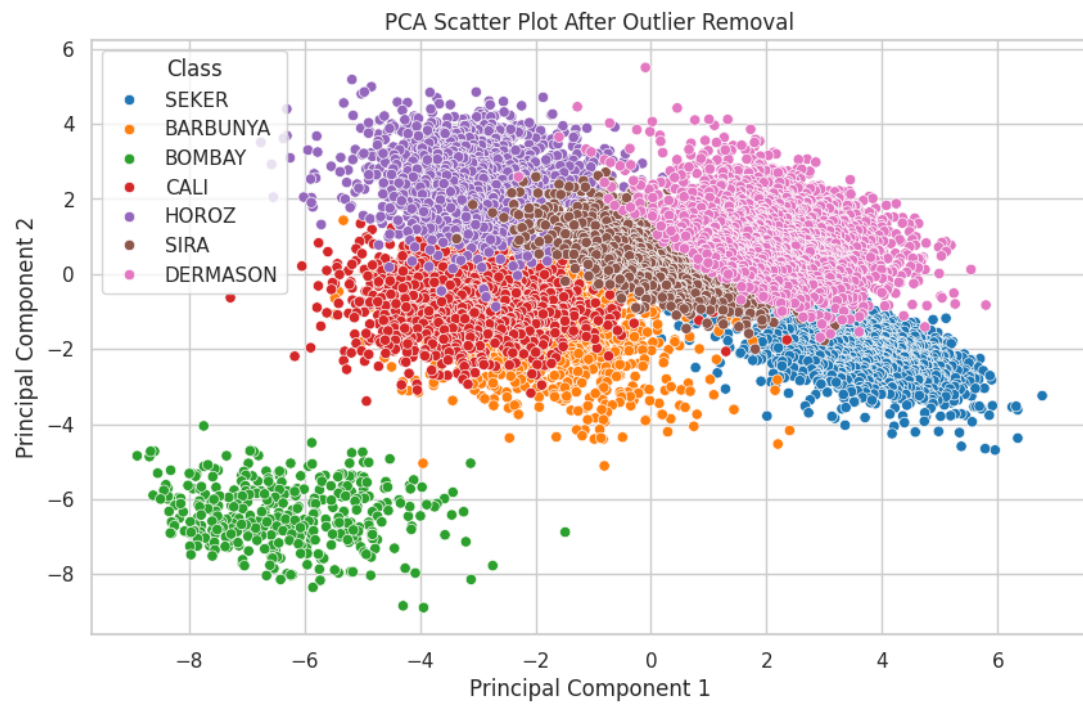


Figure 2.9: PCA plot after outlier removal

Overall, these improvements highlight the importance of outlier detection and removal in enhancing the quality and interpretability of statistical analyses.

3. CONCLUSION

This project explored the Dry Bean dataset through a comprehensive set of statistical analyses, dimensionality reduction techniques, clustering algorithms, and outlier detection methods, each revealing key insights into the structure and separability of the data.

The boxplot analysis highlighted that, although the dataset contains numerous outliers particularly in area-related features, the interquartile ranges were narrow, suggesting a compact distribution of most data points. Z-score normalization enabled fair comparison across features, and the generalized Fisher Distance identified `EquivDiameter`, `Perimeter`, and `MinorAxisLength` as the most discriminative features for class separation.

Applying Principal Component Analysis (PCA) revealed that the first two principal components captured over 80% of the dataset's variance while retaining meaningful class-discriminative information, as confirmed by both Fisher Scores and scatter plots. However, PCA, being an unsupervised technique, did not explicitly optimize for class separation. When compared to Linear Discriminant Analysis (LDA), PCA surprisingly achieved a higher silhouette score, suggesting that the underlying class structure may not be perfectly aligned with linear boundaries and that variance-preserving projections can sometimes outperform class-focused methods in complex datasets.

Clustering experiments using K-Means, DBSCAN, t-SNE, and SOM further illustrated the dataset's rich internal structure. While K-Means and DBSCAN, applied in the PCA-reduced space, provided reasonable but imperfect clusters, t-SNE and SOM, applied to the full dataset, offered superior visual separation of classes, highlighting their strength in capturing non-linear patterns and topological relationships.

Finally, applying the Isolation Forest algorithm for outlier detection and noise removal resulted in notable improvements in feature Fisher Scores and silhouette scores, confirming that removing even a small proportion of anomalous points can enhance class separability and clustering performance.

In summary, this project demonstrated the importance of combining multiple analytical approaches ranging from statistical summaries and supervised methods to unsupervised dimensionality reduction, clustering, and outlier detection to comprehensively understand and improve the performance of classification and clustering tasks on real-world datasets. The findings suggest that for datasets like Dry Bean, where non-linear re-

relationships and noise are present, unsupervised methods such as PCA and t-SNE, complemented by robust cleaning techniques, can be highly effective, sometimes even outperforming supervised alternatives.

BIBLIOGRAPHY

- [1] J. Li, K. Cheng, S. Wang, *et al.*, “Feature selection: A data perspective,” *ACM computing surveys (CSUR)*, vol. 50, no. 6, pp. 1–45, 2017.
- [2] K. Pearson, “Liii. on lines and planes of closest fit to systems of points in space,” *The London, Edinburgh, and Dublin philosophical magazine and journal of science*, vol. 2, no. 11, pp. 559–572, 1901.
- [3] R. A. Fisher, “The use of multiple measurements in taxonomic problems,” *Annals of eugenics*, vol. 7, no. 2, pp. 179–188, 1936.
- [4] P. J. Rousseeuw, “Silhouettes: A graphical aid to the interpretation and validation of cluster analysis,” *Journal of computational and applied mathematics*, vol. 20, pp. 53–65, 1987.
- [5] J. MacQueen, “Some methods for classification and analysis of multivariate observations,” in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Statistics*, University of California press, vol. 5, 1967, pp. 281–298.
- [6] M. Ester, H.-P. Kriegel, J. Sander, X. Xu, *et al.*, “A density-based algorithm for discovering clusters in large spatial databases with noise,” in *kdd*, vol. 96, 1996, pp. 226–231.
- [7] L. Van der Maaten and G. Hinton, “Visualizing data using t-sne.,” *Journal of machine learning research*, vol. 9, no. 11, 2008.
- [8] T. Kohonen, “Self-organized formation of topologically correct feature maps,” *Biological cybernetics*, vol. 43, no. 1, pp. 59–69, 1982.
- [9] F. T. Liu, K. M. Ting, and Z.-H. Zhou, “Isolation forest,” in *2008 eighth ieee international conference on data mining*, IEEE, 2008, pp. 413–422.

APPENDICES

The following Python code is an automatic conversion from a Jupyter notebook to a ‘.py’ script, used for reproducing the experiments and analysis described in this work. To reproduce the results, and better understanding of the code, please check out the given Jupyter notebook.

```
1 # -*- coding: utf-8 -*-
2 """cse555_final_project.ipynb
3
4 Automatically generated by Colab.
5
6 # 0. Setup
7 """
8
9 !pip install ucimlrepo
10 !pip install minisom
11
12 import pandas as pd
13 import numpy as np
14 import matplotlib.pyplot as plt
15 import seaborn as sns
16 from sklearn.preprocessing import StandardScaler
17 from sklearn.decomposition import PCA
18 from sklearn.discriminant_analysis import
19     LinearDiscriminantAnalysis as LDA
20 from sklearn.cluster import KMeans, DBSCAN
21 from sklearn.manifold import TSNE
22 from sklearn.metrics import silhouette_score
23 from sklearn.ensemble import IsolationForest
24 from minisom import MiniSom
25 from ucimlrepo import fetch_ucirepo
26
27 # Fetch and load the Dry Bean dataset using ucimlrepo
28 dry_bean = fetch_ucirepo(id=602)
29 X = dry_bean.data.features
30 y = dry_bean.data.targets.iloc[:, 0]
31 class_labels = y.unique()
```

```

32 dry_bean.variables
33
34 """# 1. Plot and Interpret Boxplots"""
35
36 import seaborn as sns
37 import matplotlib.pyplot as plt
38
39 sns.set_theme(style="whitegrid")
40 fig, axes = plt.subplots(nrows=4, ncols=4, figsize=(20, 16))
41 axes = axes.flatten()
42 palette = sns.color_palette("Set2")
43
44 # Plot each attribute's boxplot horizontally
45 for i, column in enumerate(X.columns):
46     sns.boxplot(x=X[column], ax=axes[i], color=palette[i % len
47         (palette)])
48     axes[i].set_title(column, fontsize=12)
49     axes[i].set_xlabel("")
50
51 plt.tight_layout()
52 plt.suptitle("Boxplots of Dry Bean Attributes", fontsize=16, y
53     =1.02)
54 plt.show()
55
56 """# 2. Z-Score Normalization and Fisher Distance Analysis"""
57
58 def calculate_fisher_distance(X, y):
59     class_labels = y.unique()
60     C = len(class_labels)
61     fisher_scores = {}
62
63     for i in range(X.shape[1]):
64         total_fd = 0
65         valid_pairs = 0
66         for j in range(C):
67             for k in range(j + 1, C):
68                 group_j = X.loc[y == class_labels[j], X.
69                     columns[i]]
70                 group_k = X.loc[y == class_labels[k], X.
71                     columns[i]]

```

```

69         # Skip if either group is empty
70         if group_j.empty or group_k.empty:
71             continue
72
73         mu_j = group_j.mean()
74         mu_k = group_k.mean()
75         var_j = group_j.var()
76         var_k = group_k.var()
77
78         # Skip if variance sum is zero (avoid division
79         # by zero)
80         std_sum = np.sqrt(var_j + var_k)
81         if std_sum == 0 or np.isnan(std_sum):
82             continue
83
84         mu_diff = abs(mu_j - mu_k)
85         total_fd += mu_diff / std_sum
86         valid_pairs += 1
87
88         # Avoid division by zero if no valid pairs found
89         if valid_pairs == 0:
90             fisher_scores[X.columns[i]] = np.nan
91         else:
92             fisher_scores[X.columns[i]] = total_fd /
93             valid_pairs
94
95         # Sort scores in descending order, ignoring NaNs
96         return dict(sorted(
97             ((k, v) for k, v in fisher_scores.items() if not np.
98              isnan(v)),
99             key=lambda item: item[1], reverse=True))
100
101     scaler = StandardScaler()
102     X_scaled = scaler.fit_transform(X)
103     X_scaled_df = pd.DataFrame(X_scaled, columns=X.columns)
104     fisher_scores = calculate_fisher_distance(X_scaled_df, y)
105
106     print("Fisher_Scores:")
107     for k, v in fisher_scores.items():
108         print(f"{k}: {v:.4f}")

```

```

107 """# 3. PCA Transformation and Fisher Distance Comparison"""
108
109 from sklearn.decomposition import PCA
110
111 # PCA transformation
112 pca = PCA()
113 X_pca = pca.fit_transform(X_scaled)
114
115 # Fisher Distance on PCA-transformed features
116 fisher_pca = calculate_fisher_distance(pd.DataFrame(X_pca), y)
117
118 # Compare Fisher Scores with Eigenvalues
119 explained_var = pca.explained_variance_ratio_
120
121 # Print both for comparison
122 print(f"{'PC':<6}_{'Explained_Var':<15}_{'Fisher_Score':<15}")
123 for i, (ev, fs) in enumerate(zip(explained_var, fisher_pca.
    values())):
124     print(f"PC{i+1:<3}_{'ev':<15.4f}_{'fs':<15.4f}")
125
126 """# 4. Scatter Plot Visualization"""
127
128 # Top 2 components
129 plt.figure(figsize=(10,5))
130 plt.subplot(1,2,1)
131 sns.scatterplot(x=X_pca[:,0], y=X_pca[:,1], hue=y, palette="
    tab10", s=20)
132 plt.title("Top_2_Principal_Components")
133
134 # Bottom 2 components
135 plt.subplot(1,2,2)
136 sns.scatterplot(x=X_pca[:, -2], y=X_pca[:, -1], hue=y, palette="
    tab10", s=20)
137 plt.title("Bottom_2_Principal_Components")
138 plt.tight_layout()
139 plt.show()
140
141 """# 5. Dimensionality Reduction with LDA"""
142
143 import seaborn as sns
144 import pandas as pd

```

```

145 import matplotlib.pyplot as plt
146 from sklearn.discriminant_analysis import
    LinearDiscriminantAnalysis as LDA
147
148 lda = LDA(n_components=2)
149 X_lda = lda.fit_transform(X_scaled, y)
150
151 lda_df = pd.DataFrame(X_lda, columns=["LDA1", "LDA2"])
152 lda_df["Label"] = y
153
154 plt.figure(figsize=(8, 6))
155 sns.scatterplot(data=lda_df, x="LDA1", y="LDA2", hue="Label",
    palette="tab10", s=20)
156
157 plt.title("LDA_Projection_of_Dry_Bean_Dataset_(2D)")
158 plt.xlabel("LDA1")
159 plt.ylabel("LDA2")
160 plt.legend(title="Bean_Type")
161 plt.grid(True)
162 plt.tight_layout()
163 plt.show()
164
165 import seaborn as sns
166 import matplotlib.pyplot as plt
167 import pandas as pd
168 from sklearn.metrics import silhouette_score
169
170 # Prepare DataFrames for Seaborn
171 lda_df = pd.DataFrame(X_lda, columns=["LDA1", "LDA2"])
172 lda_df["Label"] = y
173
174 pca_df = pd.DataFrame(X_pca[:, :2], columns=["PC1", "PC2"])
175 pca_df["Label"] = y
176
177 # Set up side-by-side subplots
178 fig, axs = plt.subplots(1, 2, figsize=(14, 6))
179
180 sns.scatterplot(data=lda_df, x="LDA1", y="LDA2", hue="Label",
    palette="tab10", s=20, ax=axs[0])
181 axs[0].set_title("LDA_Projection")
182 axs[0].grid(True)

```

```

183
184 sns.scatterplot(data=pca_df, x="PC1", y="PC2", hue="Label",
185                 palette="tab10", s=20, ax=axes[1])
186 axes[1].grid(True)
187
188 plt.tight_layout()
189 plt.show()
190
191 sil_lda = silhouette_score(X_lda, y)
192 sil_pca = silhouette_score(X_pca[:, :2], y)
193
194 print(f"Silhouette_Score_(LDA):_{sil_lda:.4f}")
195 print(f"Silhouette_Score_(PCA):_{sil_pca:.4f}")
196
197 """# 6. Clustering with K-Means, DBSCAN, t-SNE, and SOM
198
199 ## K-Means
200 """
201
202 import numpy as np
203 import matplotlib.pyplot as plt
204 import seaborn as sns
205 from sklearn.cluster import KMeans
206
207 # Fit KMeans on PC1 and PC2
208 kmeans = KMeans(n_clusters=7, random_state=0)
209 k_labels = kmeans.fit_predict(X_pca[:, :2])
210
211 # Create side-by-side plots
212 fig, axes = plt.subplots(1, 2, figsize=(14, 6))
213
214 sns.scatterplot(
215     x=X_pca[:, 0], y=X_pca[:, 1], hue=k_labels, palette="tab10",
216     s=30, edgecolor='w', alpha=0.8, ax=axes[0]
217 )
218 axes[0].set_title("KMeans_Clusters_(on_PC1_&_PC2)")
219 axes[0].set_xlabel("PC1")
220 axes[0].set_ylabel("PC2")
221 axes[0].legend(title="Cluster")

```



```

222
223 sns.scatterplot(
224     x=X_pca[:, 0], y=X_pca[:, 1], hue=y, palette="tab10",
225     s=30, edgecolor='w', alpha=0.8, ax=axes[1]
226 )
227 axes[1].set_title("Original_Class_Labels_(on_PC1_&_PC2)")
228 axes[1].set_xlabel("PC1")
229 axes[1].set_ylabel("PC2")
230 axes[1].legend(title="Class")
231
232 plt.tight_layout()
233 plt.show()
234
235 """## DBSCAN"""
236
237 from sklearn.cluster import DBSCAN
238
239 # eps is the maximum distance for a point to be considered as
    part of a neighborhood
240 dbscan = DBSCAN(eps=0.5, min_samples=5)
241 y_db = dbscan.fit_predict(X_pca[:, :2])
242
243 df_dbscan = pd.DataFrame({
244     'PC1': X_pca[:, 0],
245     'PC2': X_pca[:, 1],
246     'Cluster': y_db
247 })
248
249 plt.figure(figsize=(8, 6))
250 sns.scatterplot(
251     data=df_dbscan,
252     x='PC1',
253     y='PC2',
254     hue='Cluster',
255     palette=sns.color_palette("tab10", len(set(y_db))),
256     s=30,
257     edgecolor='w',
258     alpha=0.8
259 )
260
261 plt.title(f"DBSCAN_Clustering_(PC1_&_PC2)_-_{len(set(y_db))}_-_{

```

```

        (1_if_1_in_y_db_else_0)}_Clusters")
262 plt.xlabel("PC1")
263 plt.ylabel("PC2")
264 plt.legend(title="DBSCAN_Label")
265 plt.show()
266
267 sil_dbscan = silhouette_score(X_pca[:, :2], y_db)
268 print(f"Silhouette_Score_(DBSCAN):_{sil_dbscan:.4f}")
269
270 """## t-SNE"""
271
272 from sklearn.manifold import TSNE
273
274 # Perform t-SNE
275 tsne = TSNE(n_components=2, random_state=0)
276 X_tsne = tsne.fit_transform(X_scaled)
277
278 plt.figure(figsize=(8,6))
279 sns.scatterplot(x=X_tsne[:,0], y=X_tsne[:,1], hue=y, palette="
        tab10", alpha=0.8, s=20)
280 plt.title("t-SNE_Visualization_with_True_Labels")
281 plt.xlabel("t-SNE1")
282 plt.ylabel("t-SNE2")
283 plt.show()
284
285 sil_tsne = silhouette_score(X_tsne, y)
286 print(f"Silhouette_Score_(t-SNE):_{sil_tsne:.4f}")
287
288 """## SOM"""
289
290 from minisom import MiniSom
291
292 # Define SOM parameters
293 som_shape = (10, 10) # grid size
294
295 # Initialize SOM
296 som = MiniSom(som_shape[0], som_shape[1], X_scaled.shape[1],
        sigma=1.5, learning_rate=.5)
297
298 # Initialize weights randomly
299 som.random_weights_init(X_scaled)

```

```

300 # som.pca_weights_init(X_scaled)
301
302 # Train the SOM
303 som.train_random(X_scaled, 1000) # 1000 iterations
304
305 w_x, w_y = zip(*[som.winner(x) for x in X_scaled])
306 w_x = np.array(w_x)
307 w_y = np.array(w_y)
308
309 # Define the color palette
310 unique_classes = y.unique()
311 palette = sns.color_palette("tab10", n_colors=len(
    unique_classes))
312 colors_in_order = [palette[i] for i, cls in enumerate(
    unique_classes)]
313
314 # Create a mapping of classes to colors
315 class_colors = {cls: palette[i] for i, cls in enumerate(
    unique_classes)}
316
317 plt.figure(figsize=(11, 10))
318 plt.pcolor(som.distance_map().T, cmap='bone_r', alpha=.2)
319 plt.colorbar()
320
321 for c in np.unique(y):
322     idx_target = y == c
323     plt.scatter(w_x[idx_target] + 0.5 + (np.random.rand(np.sum
324         (idx_target)) - 0.5) * 0.8,
325                 w_y[idx_target] + 0.5 + (np.random.rand(np.sum
326         (idx_target)) - 0.5) * 0.8,
327                 s=20, label=c, alpha=0.7, edgecolor='w', color
328                 =class_colors[c])
329
330 plt.legend(loc='upper_left')
331 plt.grid()
332 plt.show()
333
334 import matplotlib.gridspec as gridspec
335
336 labels_map = som.labels_map(X_scaled, y)
337

```

```

335 fig = plt.figure(figsize=(10, 10))
336 the_grid = gridspec.GridSpec(som_shape[0], som_shape[1], fig)
337 for position in labels_map.keys():
338     label_fracs = [labels_map[position][y] for y in y.unique()
339                    ]
340     plt.subplot(the_grid[som_shape[1] - 1 - position[1],
341                        position[0]], aspect=1)
342     patches, texts = plt.pie(label_fracs, colors=
343                             colors_in_order)
344
345 plt.legend(labels=y.unique(), loc="best", bbox_to_anchor=(0,
346                -0.05), ncol=1)
347 plt.show()
348
349 """"# 7. Outlier Detection and Data Cleaning""""
350
351 from sklearn.ensemble import IsolationForest
352
353 # Normalize the features using z-score
354 scaler = StandardScaler()
355 X_scaled = scaler.fit_transform(X)
356
357 # Initialize and fit the Isolation Forest
358 iso_forest = IsolationForest(contamination=0.01, random_state
359                             =42)
360 outlier_preds = iso_forest.fit_predict(X_scaled)
361
362 # Outlier mask: -1 indicates outlier, 1 indicates inlier
363 outliers = outlier_preds == -1
364 inliers = outlier_preds == 1
365 mask = outlier_preds != -1
366 X_clean = X_scaled[inliers]
367 y_clean = y[inliers]
368
369 print(f"Original_sample_size:_{X.shape[0]}")
370 print(f"Cleaned_sample_size:_{X_clean.shape[0]}")
371 print(f"Number_of_outliers_detected_and_removed:_{X.shape[0]}_-
372       _X_clean.shape[0]}")
373
374 # Re-run PCA on cleaned data
375 X_pca_clean = pca.fit_transform(X_clean)

```

```

370
371 plt.figure(figsize=(10,6))
372 sns.scatterplot(x=X_pca_clean[:,0], y=X_pca_clean[:,1], hue=
    y_clean, palette='tab10', legend='full')
373 plt.title('PCA_Scatter_Plot_After_Outlier_Removal')
374 plt.xlabel('Principal_Component_1')
375 plt.ylabel('Principal_Component_2')
376 plt.show()
377
378 X_clean_df = pd.DataFrame(X_clean, columns=X.columns).
    reset_index(drop=True)
379 y_clean_reset = y_clean.reset_index(drop=True)
380
381 # Recalculate Fisher Scores on cleaned data
382 fisher_clean = calculate_fisher_distance(X_clean_df,
    y_clean_reset)
383
384 print("Fisher_Scores_AFTER_Outlier_Removal:")
385 for k, v in fisher_clean.items():
386     print(f"{k}: {v:.4f}")
387
388 from sklearn.metrics import silhouette_score
389 from sklearn.cluster import KMeans
390
391 # Before cleaning
392 kmeans_before = KMeans(n_clusters=7, random_state=0).fit(
    X_scaled)
393 sil_before = silhouette_score(X_scaled, kmeans_before.labels_)
394
395 # After cleaning
396 kmeans_after = KMeans(n_clusters=7, random_state=0).fit(
    X_clean)
397 sil_after = silhouette_score(X_clean, kmeans_after.labels_)
398
399 print(f"Silhouette_Score_BEFORE_Cleaning: {sil_before:.4f}")
400 print(f"Silhouette_Score_AFTER_Cleaning: {sil_after:.4f}")

```

Listing 3.1: Project Python Code