

Lab 03 Answers Report

HASSAN ALMOSA

Topic: *Artificial Intelligence COMP3742*

Lab Number: *3*

Campus: *Tonsely*



Contents

1 Question 1: K-Means and Visualisation	2
1.1 Results	2
2 Question 2: Silhouette Score	6
2.1 Implementation	6
2.2 Results	7
3 Question 3: Optimal Clusters with Elbow Method	8
3.1 Results	9
3.2 Evaluation	9
4 Question 4: Classification with Centroids	9
4.1 Evaluation and Discussion	9
5 Question 5: Feature Scaling and Initialisation	10
5.1 Comparison and Visualisation	10
5.2 Evaluation	13
5.3 Other Algorithms	14
6 Appendix: AI Usage Declaration	18

1 Question 1: K-Means and Visualisation

Apply K-Means clustering to the Wine dataset and visualise the clusters using a 2D scatter plot.

1.1 Results

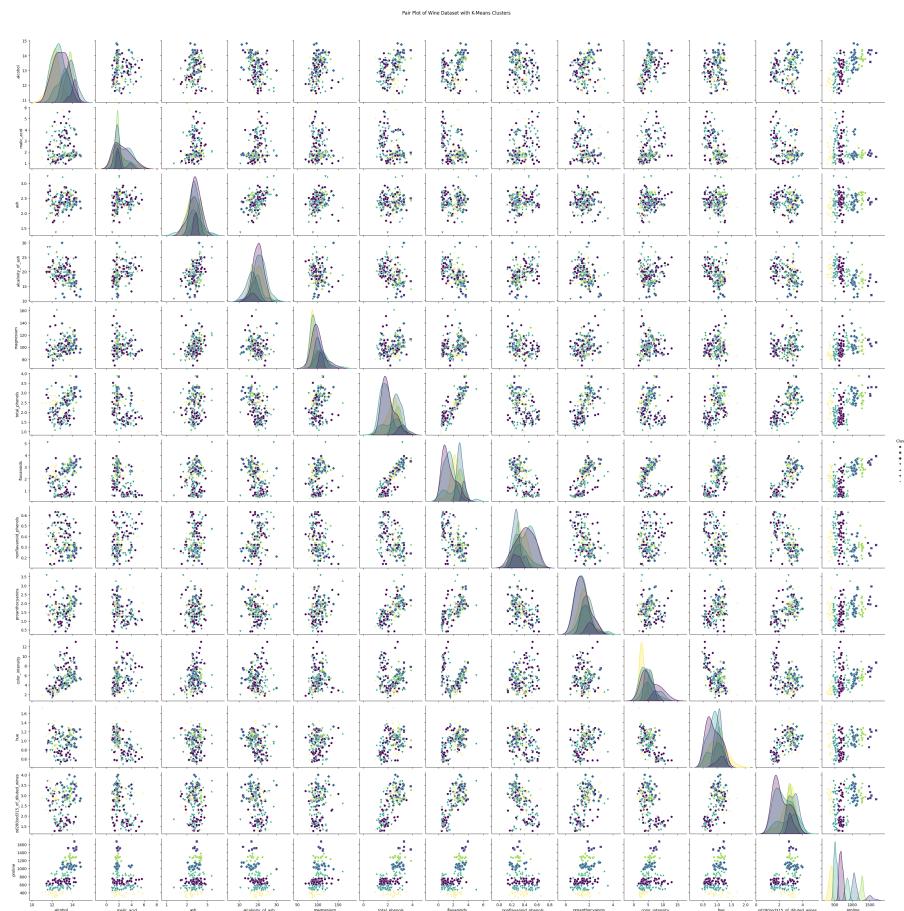


Figure 1: Wine dataset Pair Plot

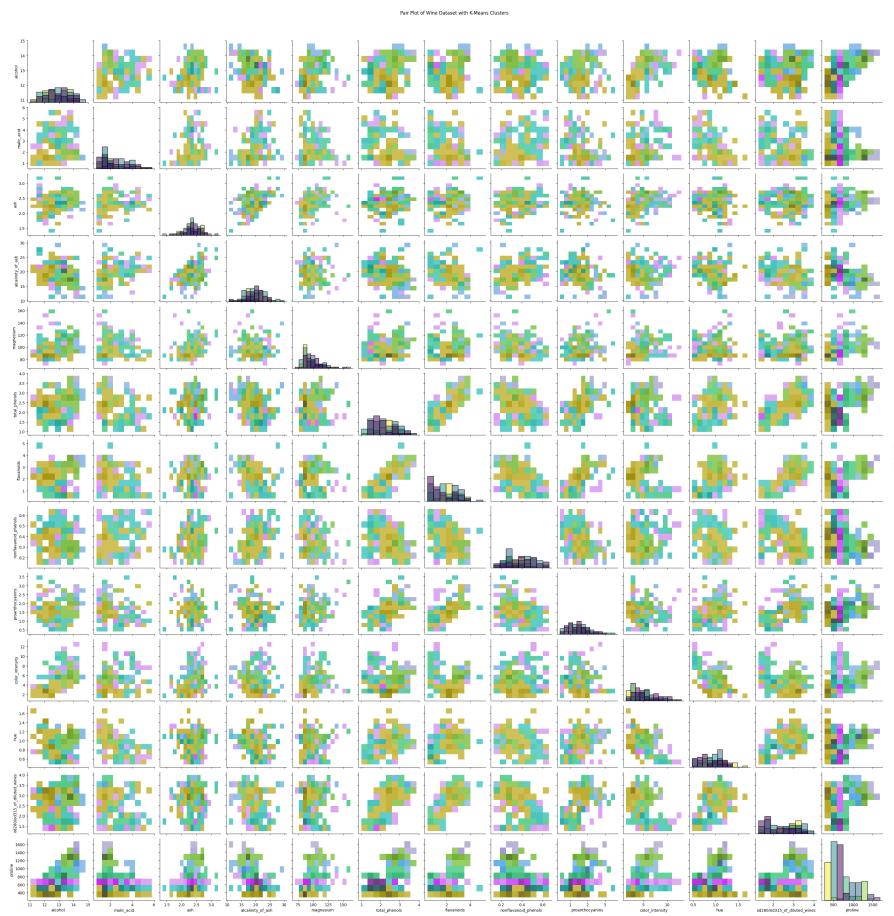


Figure 2: Wine dataset Pair Plot - Histogram

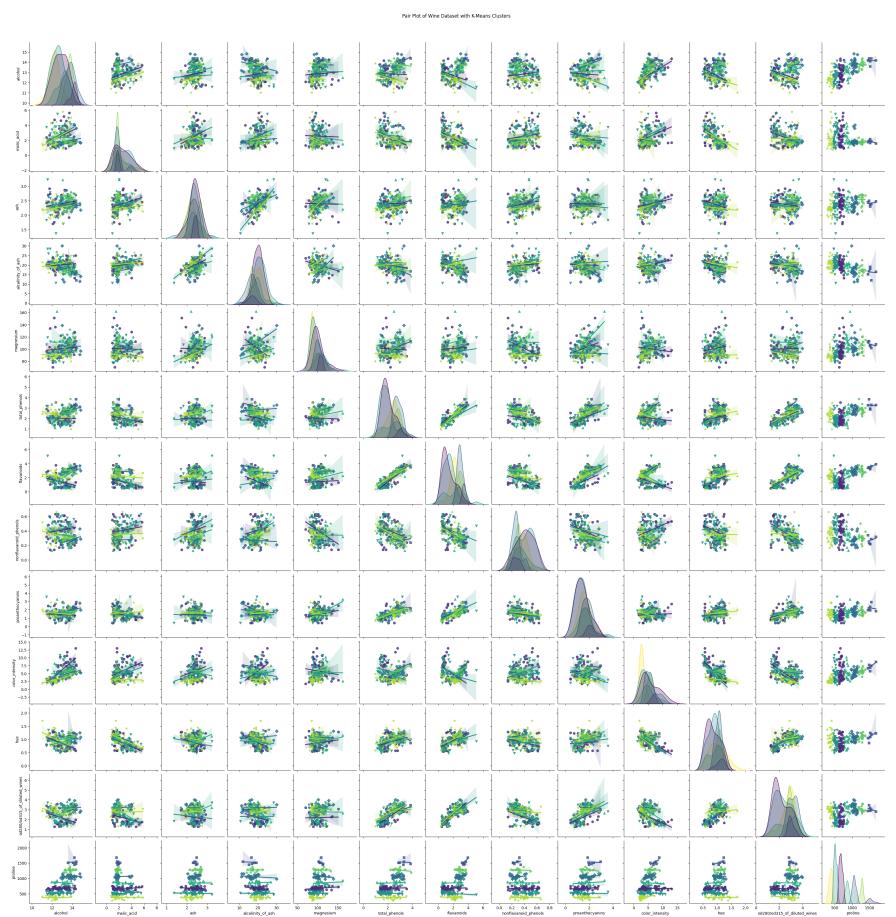


Figure 3: Wine dataset Pair Plot - Regression Line

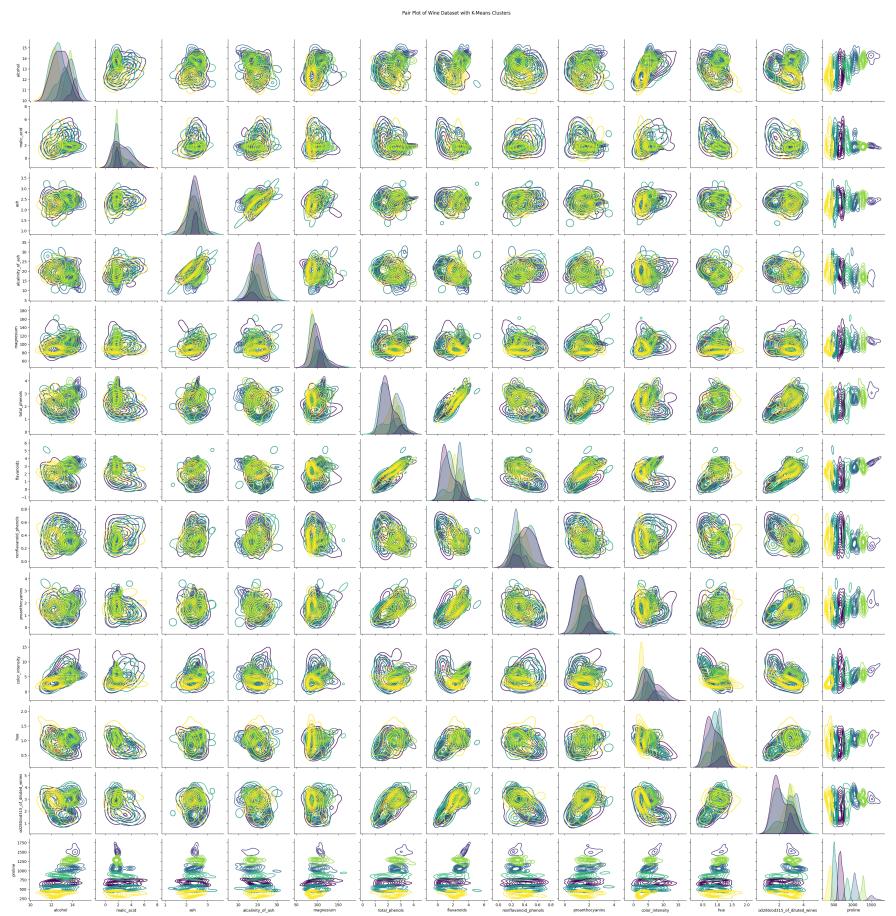


Figure 4: Wine dataset Pair Plot - Kernel Density Estimate

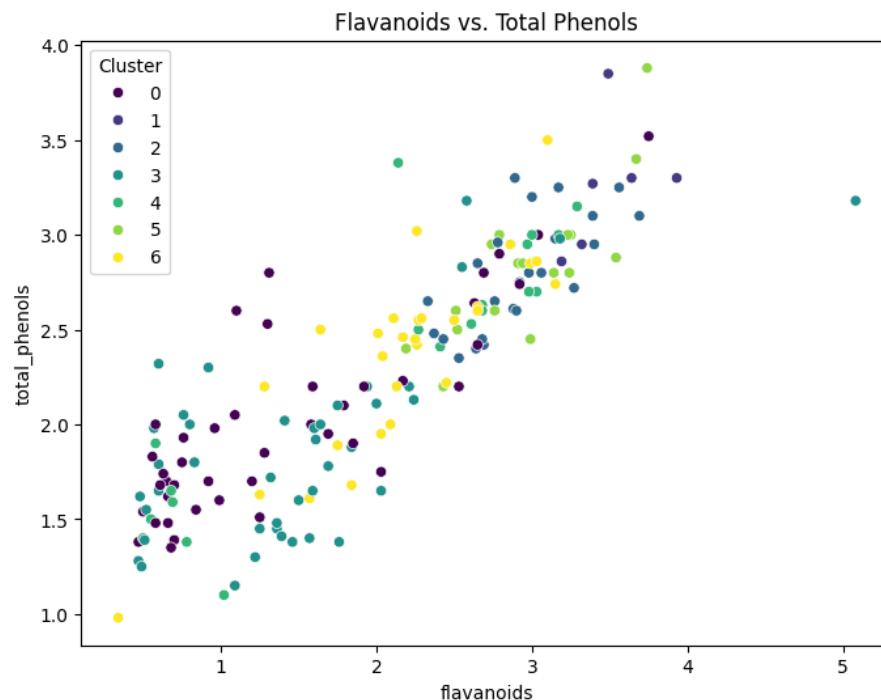


Figure 5: Pair Plot Single Pair

2 Question 2: Silhouette Score

Calculate the silhouette score for the K-Means clustering performed on the Wine dataset.

2.1 Implementation

Input:

```

1   from sklearn.metrics import silhouette_score, silhouette_samples
2
3 for k in range(3, 20):
4     km = KMeans(n_clusters=k, random_state=1, n_init=20).fit(X)
5     print(f'K={k}, silhouette={silhouette_score(X, km.labels_):.3f}')
6
7
8
9 s_avg = silhouette_score(X, labels)
10 s_samples = silhouette_samples(X, labels)
11 s_per_cluster = [float(s_samples[labels == k].mean()) for k in range(kmeans.
12   n_clusters)]
13
14 print("Silhouette avg:", round(s_avg, 3))
15 print("Silhouette per cluster:", [round(v, 3) for v in s_per_cluster])

```

Output:

```

1 K=3, silhouette=0.571
2 K=4, silhouette=0.562
3 K=5, silhouette=0.549
4 K=6, silhouette=0.566
5 K=7, silhouette=0.562
6 K=8, silhouette=0.540
7 K=9, silhouette=0.532
8 K=10, silhouette=0.520
9 K=11, silhouette=0.518
10 K=12, silhouette=0.515
11 K=13, silhouette=0.509
12 K=14, silhouette=0.508
13 K=15, silhouette=0.502
14 K=16, silhouette=0.491
15 K=17, silhouette=0.483
16 K=18, silhouette=0.477
17 K=19, silhouette=0.474
18 Silhouette avg: 0.562
19 Silhouette per cluster: [0.506, 0.65, 0.601, 0.528, 0.567, 0.674, 0.575]

```

2.2 Results

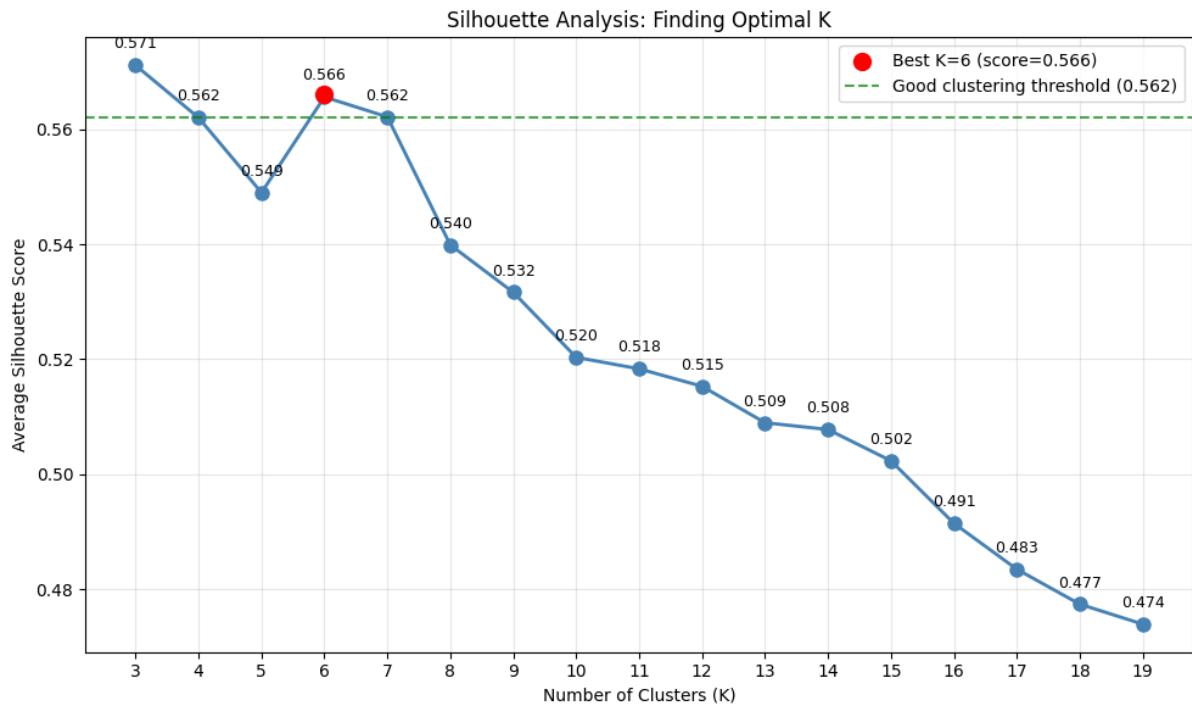


Figure 6: Silhouette Scores Analysis

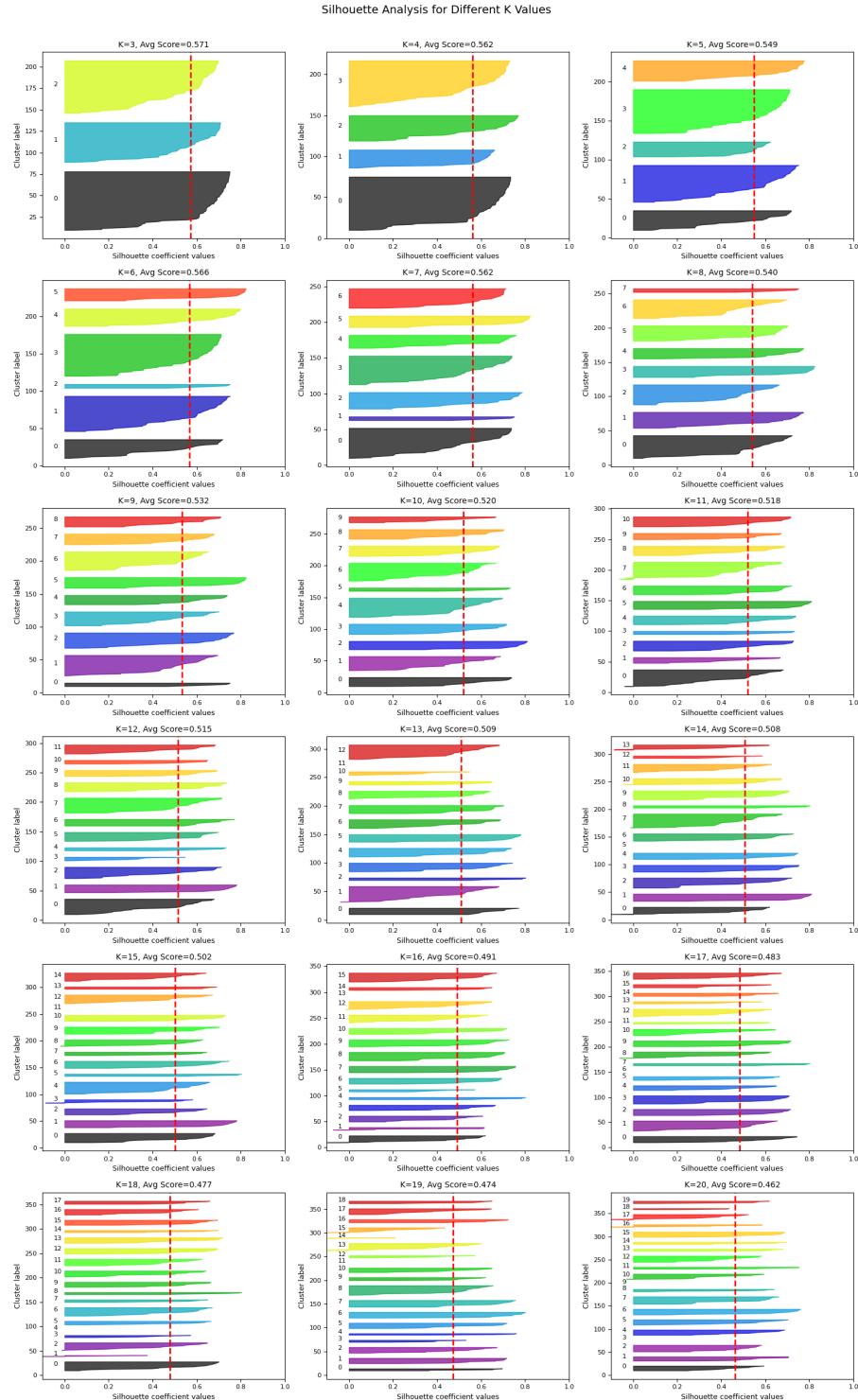


Figure 7: Silhouette's Cluster Analysis per K

3 Question 3: Optimal Clusters with Elbow Method

Apply K-Means clustering to the Wine dataset, determine the optimal number of clusters using the Elbow method, and visualise the clusters via plots. Compare the results obtained with those found using the silhouette score.

3.1 Results

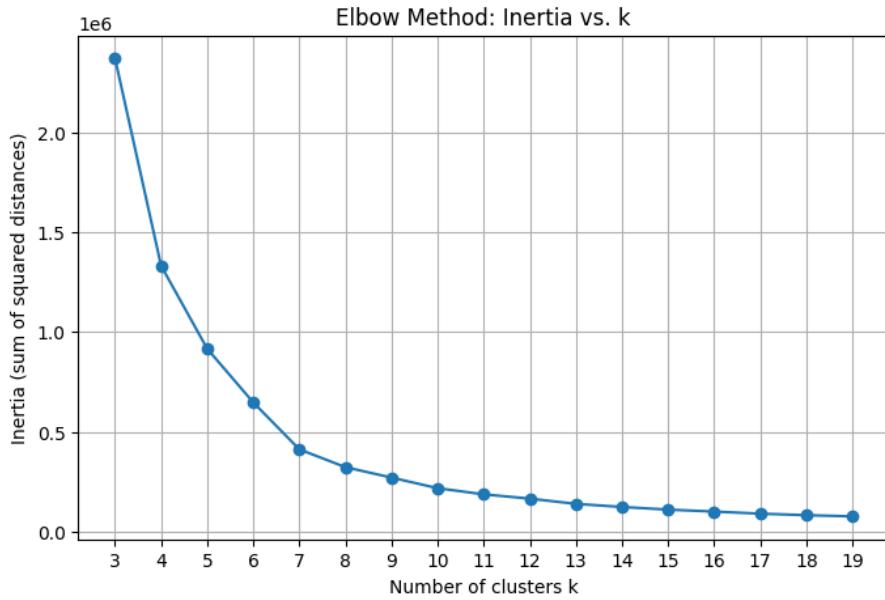


Figure 8: K-means Elbow curve: Inertia vs K

3.2 Evaluation

The curve was smooth and a little hard to visibly determine exactly at which point optimal clustering was achieved, as with each increasing K the curve gradually smoothes out compared to the previous K. However, we could make a decision based on the Silhouette score in Figure 6 where it has shown sharp increase of silhouette score after a sharp decrease at K = 6, and taking the median of the Elbow curve score between where K = 4 to 8, where it would be K = 6.

4 Question 4: Classification with Centroids

Perform K-Means clustering on the Wine dataset, determine the optimal number of clusters using the Elbow method, and then use cluster centroids to classify new, unseen data points. Evaluate the classification accuracy and discuss any potential improvements.

4.1 Evaluation and Discussion

The optimal number of clusters was determined at K = 6, it was the middle number of the 'sweet spot' between the Elbow curve analysis and the average line score of the Silhouette coefficient cluster results per range of K [3,20] and the taking in consideration the variation sizes in clusters, and the Silhouette scores curve per K values range (the local maximum after a minimum).

(

The Wine dataset K-Means clustering with K=6 achieved 96.3% classification accuracy (52/54 correct predictions) using centroid based classification. Cluster to class mapping identified 2 classes with misclassifications in the confusion matrix. One class 0 predicted as class 1, and one class 2 predicted as class 1

¹ Found 6 centroids

```

2 Test data assigned to clusters: [0 1 3 4 5]
3
4
5 Cluster      Wine Class mapping:
6 Cluster 0      Class 1
7 Cluster 1      Class 0
8 Cluster 2      Class 1
9 Cluster 3      Class 2
10 Cluster 4      Class 1
11 Cluster 5      Class 1
12
13 Classification Accuracy: 0.963 (96.3%)
14
15 Confusion Matrix:
16 [[17  1  0]
17 [ 0 21  0]
18 [ 0  1 14]]

```

5 Question 5: Feature Scaling and Initialisation

Apply K-Means clustering to the Wine dataset. Experiment with different feature scaling methods (e.g., StandardScaler, MinMaxScaler) and K-Means initialisation techniques (e.g., k-means++, random). Compare the results using silhouette scores and visualise the clusters using pair plots.

5.1 Comparison and Visualisation

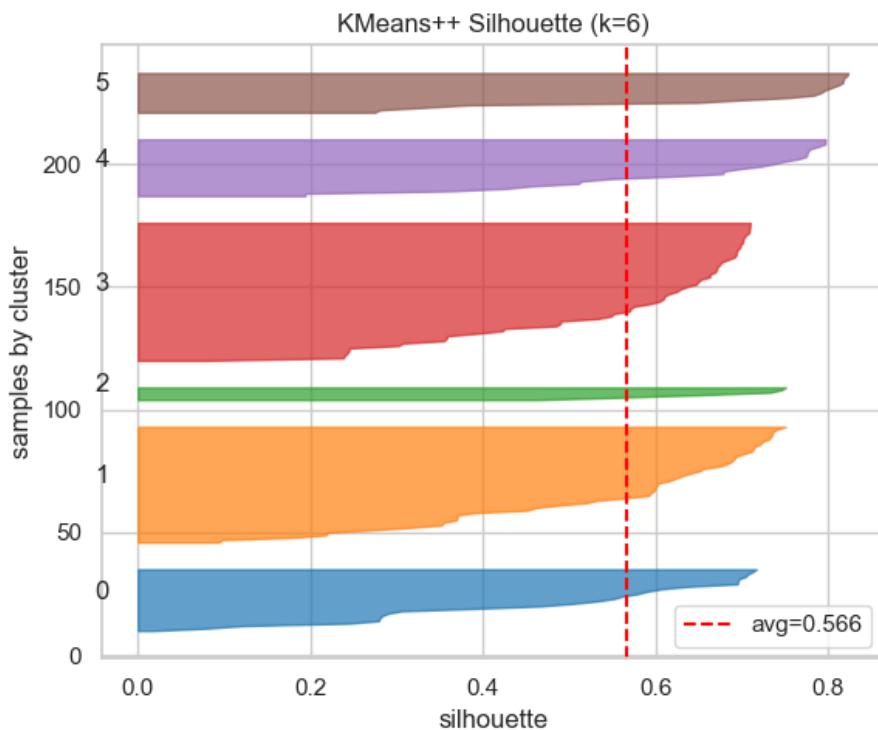


Figure 9: K++ Silhouettes - Unscaled

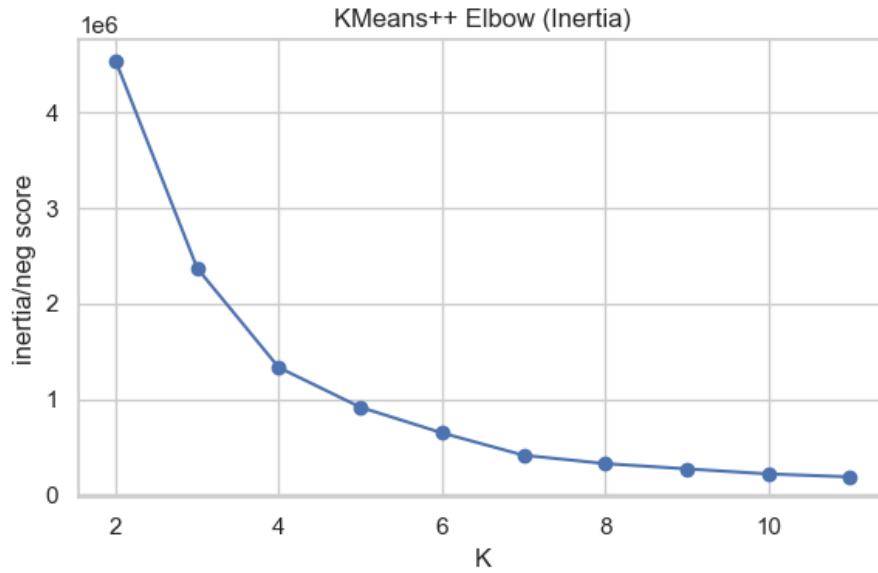


Figure 10: K++ Elbow curve: Inertia vs K - Unscaled

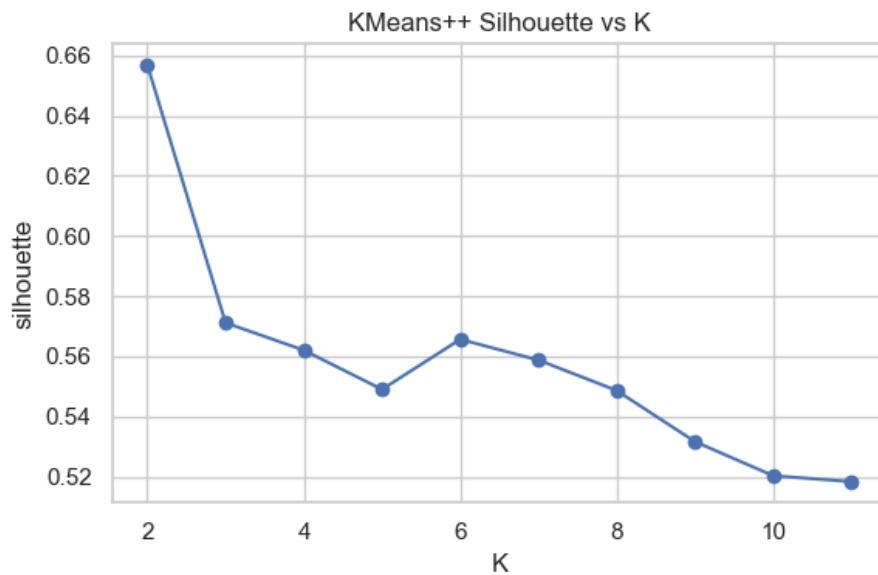


Figure 11: K++ Silhouette vs K - Unscaled

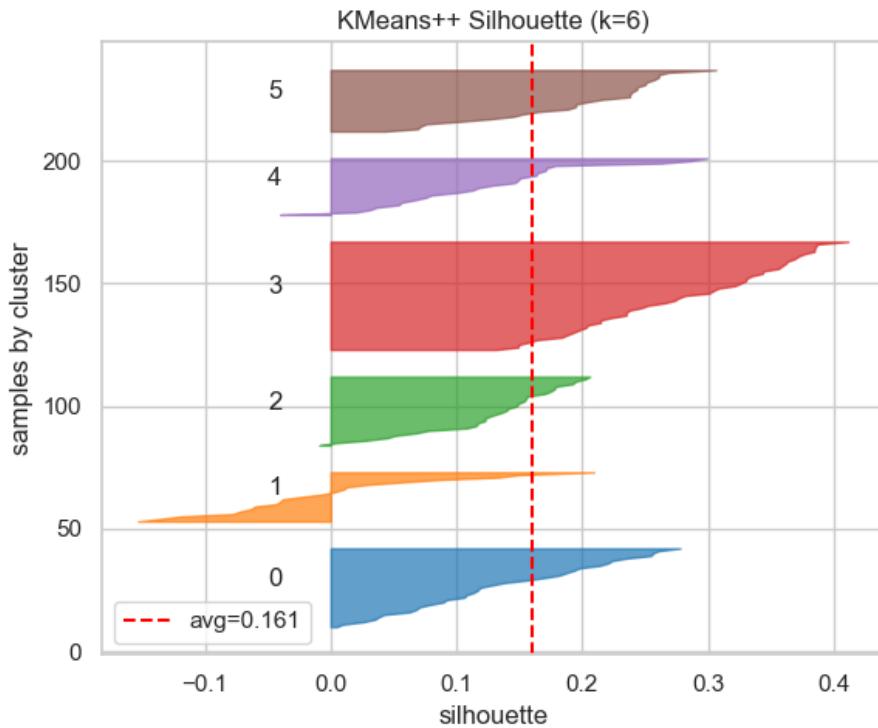


Figure 12: K++ Silhouettes - Scaled

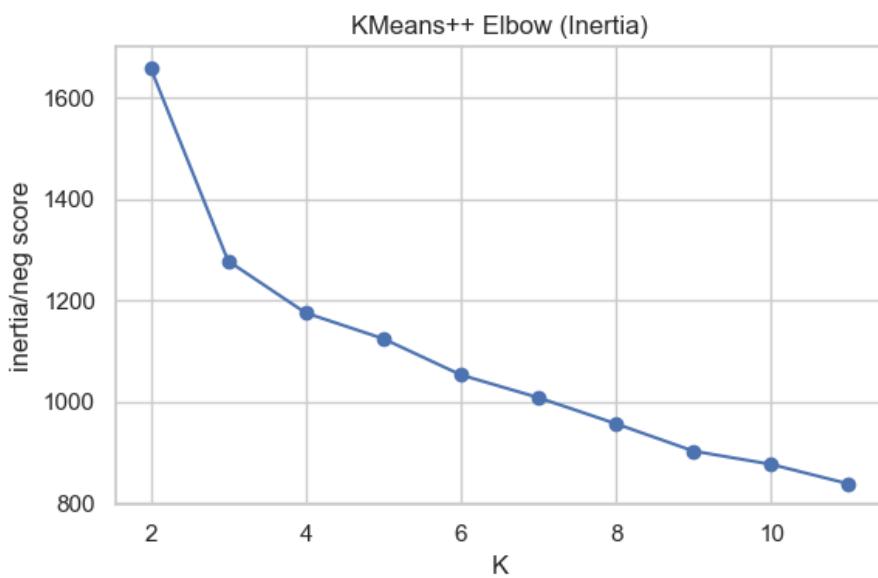


Figure 13: K++ Elbow curve: Inertia vs K - Scaled

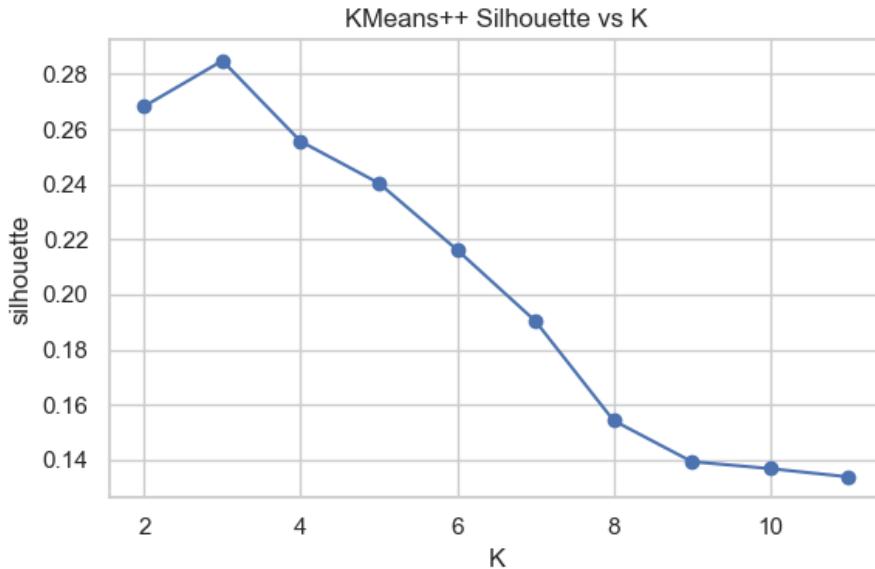


Figure 14: K++ Silhouette vs K - Scaled

5.2 Evaluation

There did not seem to be any significant effect in clustering using K++ in comparison with K-menas, both Elbow curve and Silhouette compression showed for both to be clustered similarly on the Wine dataset. On the other hand, using a `StandardScaler()` to yield different clustering proportions where it seemed a little more equipportionate, yet showing negative silhouette scores as well in Figure. 12. The elbow curve also did not seem to indicate healthy clustering with a clear Elbow point.

More algorithms were also explored for tinkering and exploration purposes, we cannot evaluate these due to lack of knowledge about them.

5.3 Other Algorithms

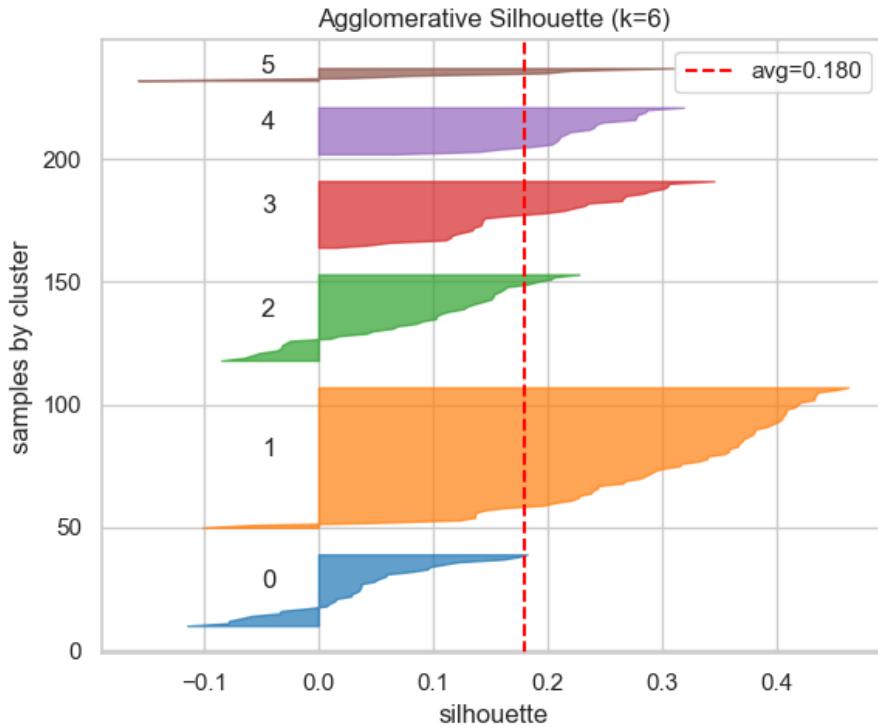


Figure 15: Agglomerative Silhouettes

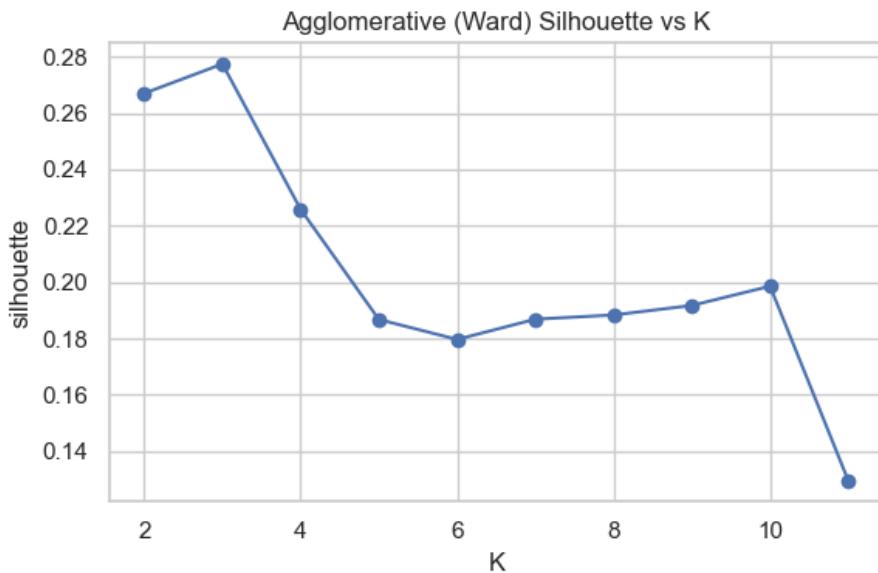


Figure 16: Agglomerative (Ward) Silhouette vs K

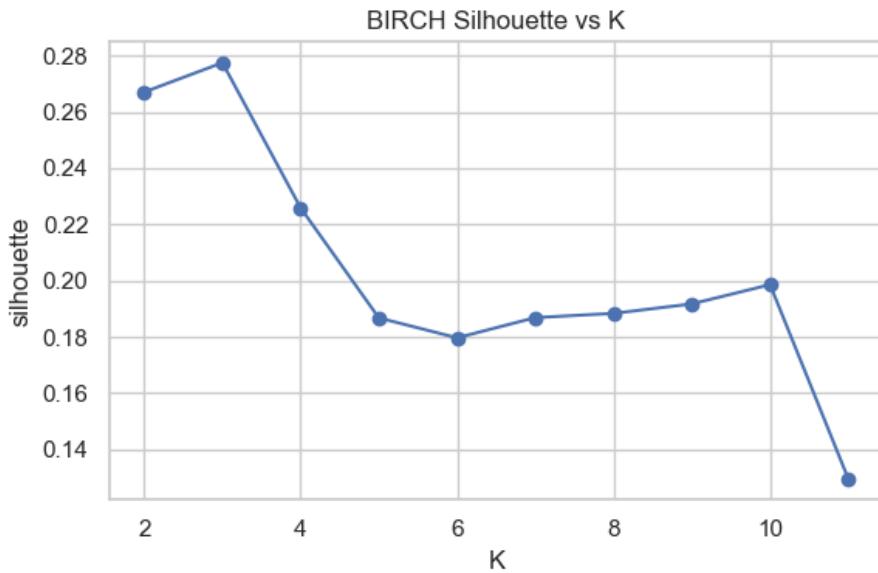


Figure 17: BIRCH Silhouette vs K

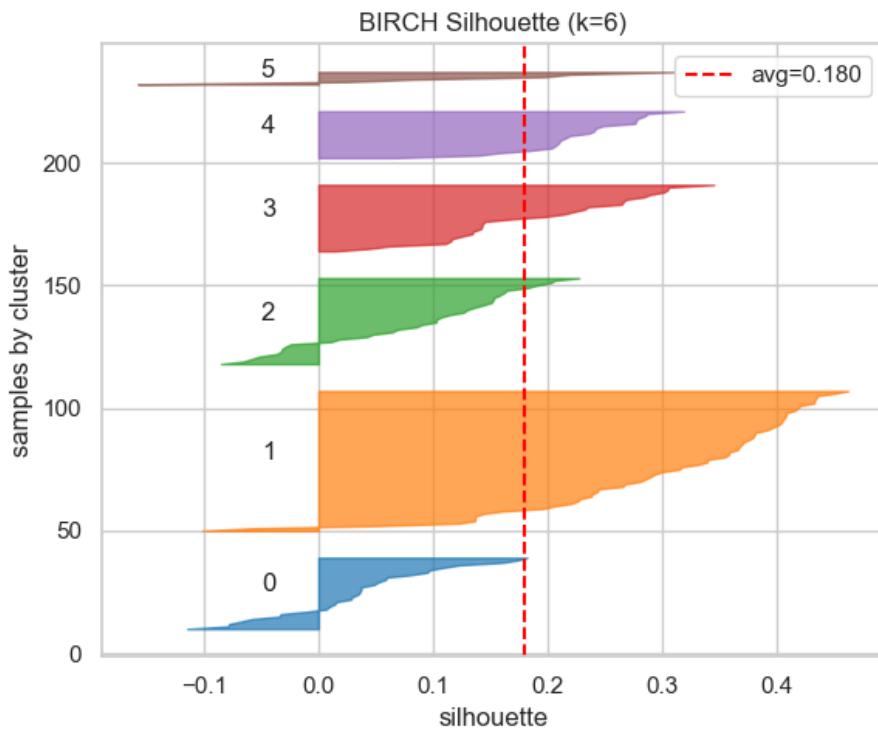


Figure 18: BIRCH Silhouettes

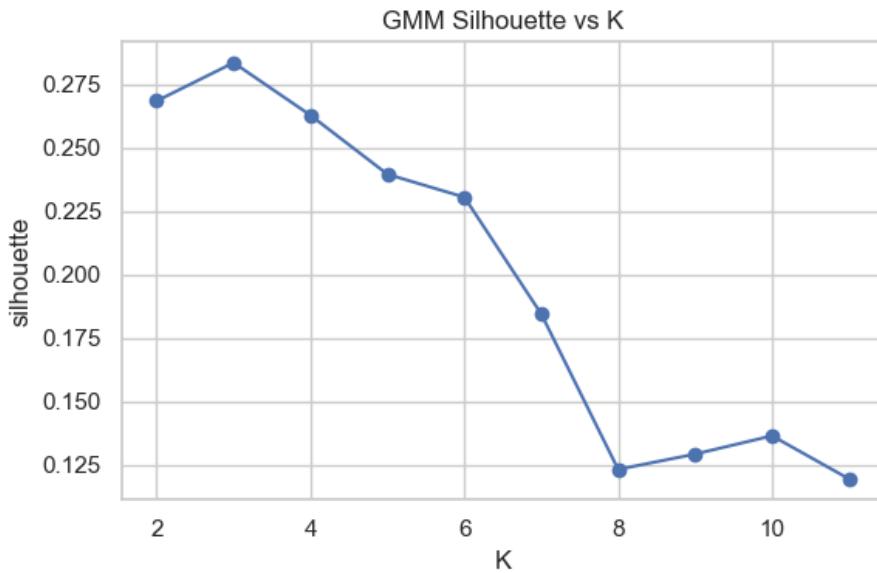


Figure 19: Guassian Mixtures Model Silhouette vs K

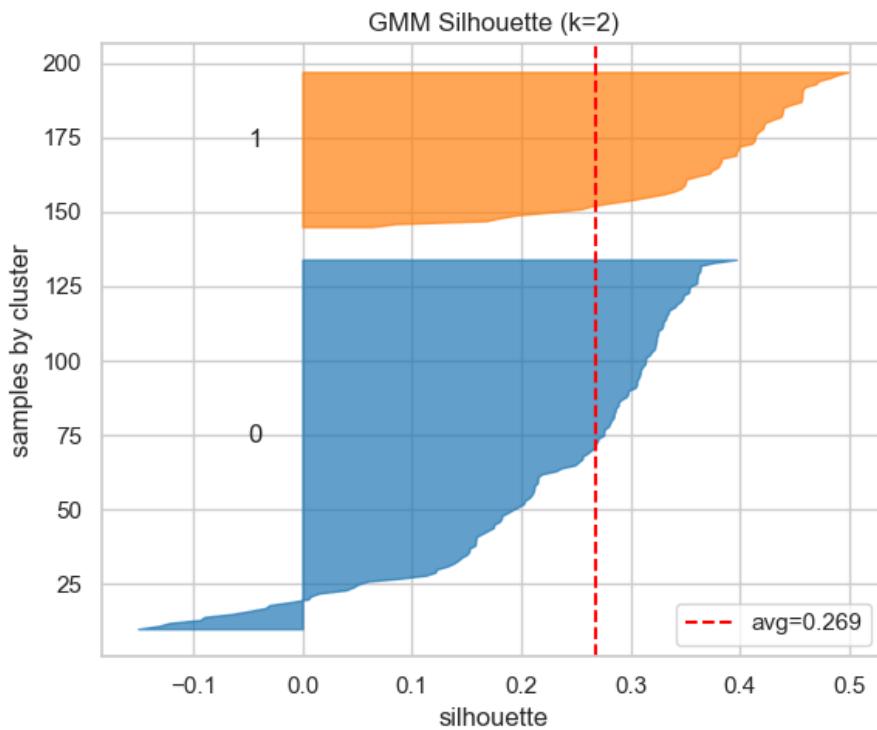


Figure 20: Guassian Mixtures Model Silhouettes

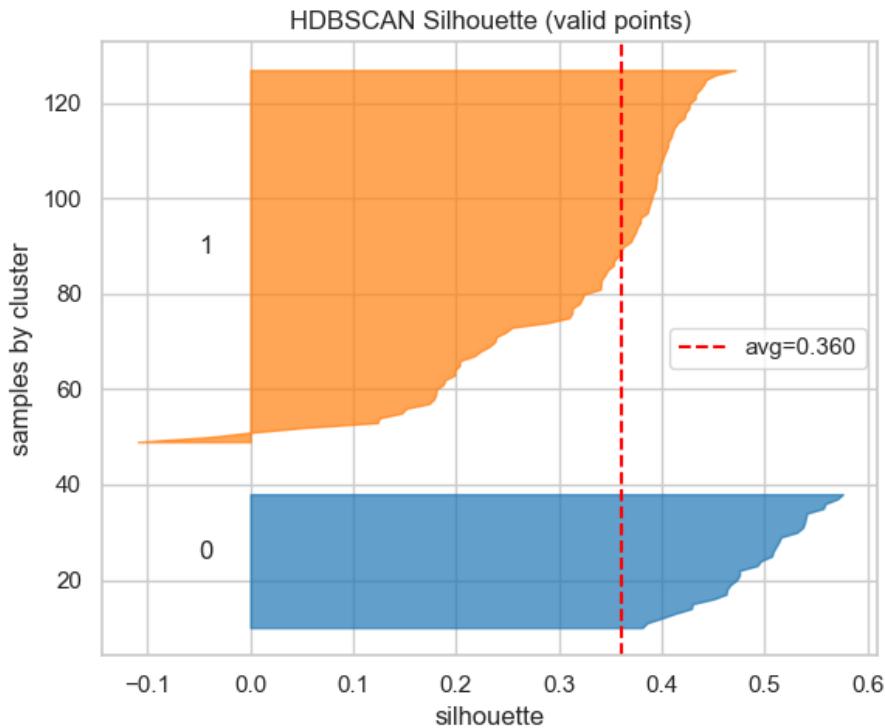


Figure 21: HBDSCAN Silhouette

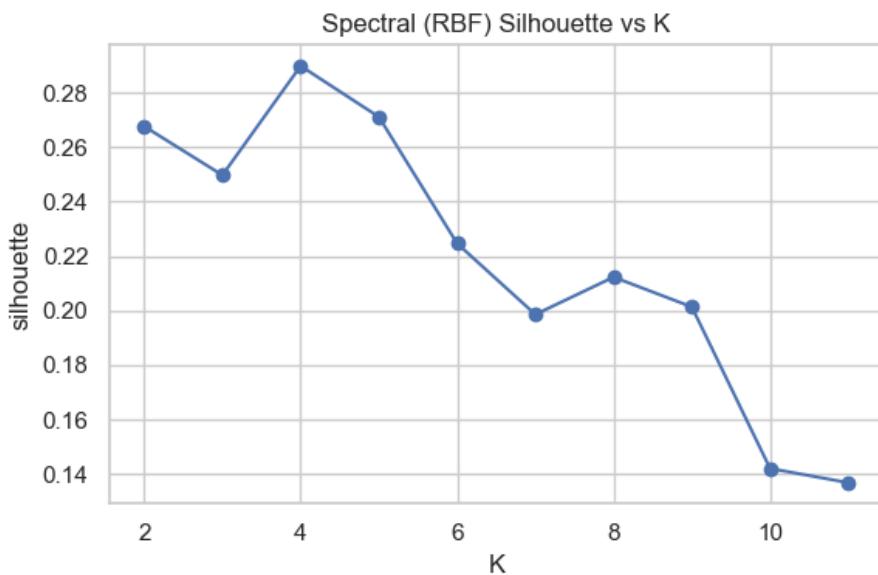


Figure 22: Spectral (RBF) Silhouette vs K

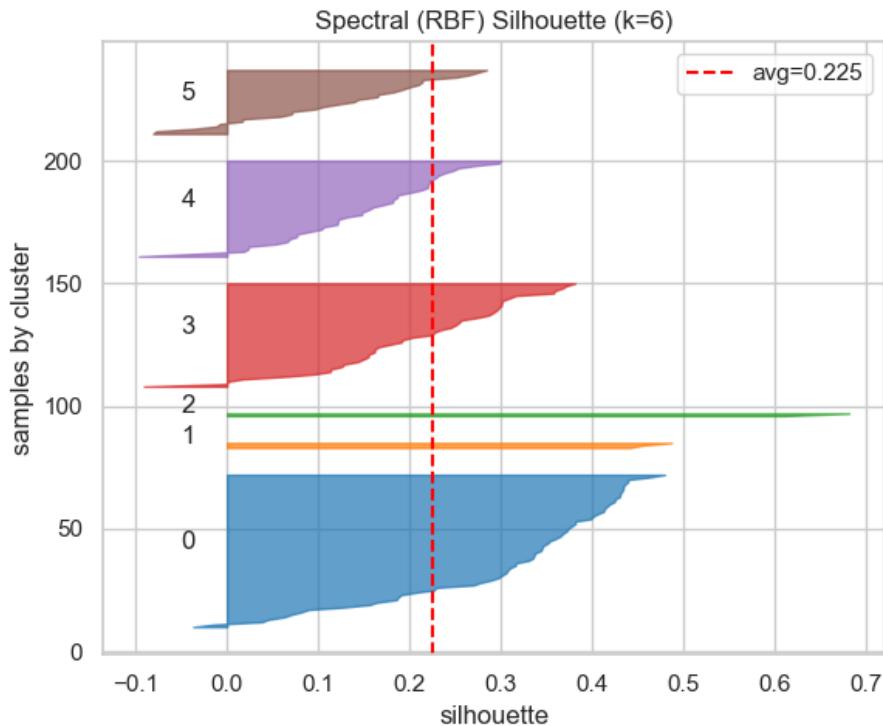


Figure 23: Spectral (RBF) Silhouette

6 Appendix: AI Usage Declaration

University AI Use Declaration

I, Hassan Almosa, declare that I have used an AI tool to assist with the completion of this assessment in the following areas: Formatting the document and files, proofreading, ensuring following submission format, code auto-completion and patterned replication, learning and illustration of topics and questions. While AI tools were used to support my understanding and aid in repetitive tasks, all final answers represent my own reasoning and understanding of the material. The details of my interaction with the AI tool are documented below. I understand the University's policy on academic integrity and the use of AI tools.

AI Tool and Interaction Log

- **AI Tool Used:**
 - GitHub Copilot, used LLM Models - Google’s Gemini 2.5 Pro (Preview), Anthropic’s Claude Sonnet 4, OpenAI’s ChatGPT-5 (Preview).
- **Date of Interaction:** 16, 18, 19, 23 August 2025

My Prompt to Copilot’s LLM

```

1 okay here's my ans wrs for weekly check in, fix mistakes and gently polish ( dont
     change much), and update doc, add word count to front page:
2 ## Week 1
3 Weekly Check-In Report
4

```

```

5 - ##### Question 1 | Progress: What did you accomplish this past week? Any issues
   slowing down progress?
6 K-means was a good introductory refresher into Unsupervised Learning. Using Scikit
   -Learn libraries was the easy and straightforward part, and as usual with
   studying AI, the library implementation bit is the most trivial part. However,
   the statistical theory behind these tools, laying that statistical groundwork
   foundation of how they work and what they are actually doing to infer the
   statistically machine-learned relationships and predict is the most crucial
   part that requires a considerable amount of self-study for rigorous
   understanding of AI/ML. This often leads to rabbit holes when trying to fine-
   tune and explore different hyperparametric tunings, as the concepts behind
   each of the hyperparameters can be quite diverse.
7
8 I have also looked at various datasets of interestes in Kaggle and bookmarked them
   for later, I need to revisit after gaining better clarification of the
   project and their appropriateness.
9 - ##### Question 2 | Plans: What are your goals for the next week? How does this
   match against your overall project plan?
10 For Unsupervised Learning, I briefly discuss K-means with my friend who just was
    done with writing mathematical proof of K-means and review his work on that,
    and go through my Week 3 notes and Lab notebook, to solidify K-means and basic
    unsupervised learning concepts. From there, I want to explore and review more
    of the unsupervised models and techniques of the following( briefly, and
    expand if time allows):
11 #### **Clustering Models**
12
13 - **Hierarchical Clustering** (Agglomerative / Divisive)
14
15 - **DBSCAN** (Density-Based Spatial Clustering of Applications with Noise)
16
17 - **OPTICS** (Ordering Points To Identify the Clustering Structure)
18
19 - **Mean Shift Clustering**
20
21 - **Gaussian Mixture Models (GMMs)** (probabilistic clustering)
22
23 - **Spectral Clustering** (uses eigenvectors of similarity matrix)
24
25 - **Affinity Propagation**
26
27
28 ---
29
30 #### **Dimensionality Reduction / Feature Learning**
31
32 - **PCA (Principal Component Analysis)**
33
34 - **Kernel PCA**
35
36 - **t-SNE (t-Distributed Stochastic Neighbor Embedding)**
37
38 - **UMAP (Uniform Manifold Approximation and Projection)**
39
40 - **Autoencoders** (neural network-based)
41
42 - **Factor Analysis**
43
44 - **Independent Component Analysis (ICA)**
45
46 - **Non-Negative Matrix Factorization (NMF)**
47
48

```

```

49 ---
50
51 ### **Density Estimation / Generative Models**
52
53 - Gaussian Mixture Models (GMMs)**
54
55 - Kernel Density Estimation (KDE)**
56
57 - Variational Autoencoders (VAE)**
58
59 - Generative Adversarial Networks (GANs)** (unsupervised / self-supervised flavor)
60
61 - Restricted Boltzmann Machines (RBM)**
62
63
64 ---
65
66 ### **Association Rule Learning**
67
68 - Apriori Algorithm** (frequent itemsets)
69
70 - FP-Growth** (Frequent Pattern Growth)
71
72
73 ---
74
75 ### **Neural & Advanced Methods**
76
77 - Self-Organizing Maps (SOMs)** (Kohonen maps)
78
79 - Deep Embedding Clustering (DEC)**
80
81 - Contrastive / Self-Supervised Learning** (SimCLR, BYOL, MoCo often fall between unsupervised and supervised)
82
83 Aside from this topic, I need to look into my rich and dense financial data dataset for Saudi Arabian Monetary Authority (SAMA) which I have compiled and built into postgreSQL DB for my different project for Model Context Protocol (MCP) tool and study its applicability and appropriateness for the project. Need to look into dataset preparation and perhaps publish
84 - #### Question 3 | Study: List and discuss the three main points of the material covered this week.
85 - Clustering (K-means, SSE, Iterative algorithm of updating of centroids by mean recalculation)
86 - Clustering evaluation (Inertia and SSE relation to be minimised, Elbow Curve - variance of clusters)
87 - Clustering optimisation (Silhouette Score - Average line, silhouette thickness)
88 - #### Question 4 | Insights: What was the single most confusing concept (muddiest point) for you in this weeks material, and what specific question do you have about it?
89 Silhouette analysis and tying it up with Elbow curve and the concept of variance in clusters.
90
91 - #### Question 5 | Linkage: How could one key idea from this weeks material be directly applied to your project?
92 If I would use my SAMA DB as a dataset which is a labelled time series data, we can derive different labels for different classification/regression objective we see statistically relevant. Also, since the dataset is huge and highly deminsional, we might also benefit from deminsionality reduction (PCA, etc.)
93
94

```

```

95
96
97
98 (750 words cumulative word limit).

```

Copilot's LLM Response

```

1 Of course. I've polished your weekly check-in answers for clarity and flow,
   formatted them within the LaTeX document, and added the word count to the
   title page as requested.
2
3 Here are the key changes I made:
4 - **Grammar and Phrasing**: Corrected minor typos and rephrased sentences for
   better readability without altering the original meaning.
5 - **Formatting**: Converted your lists of ML models and key concepts into
   structured 'itemize' lists for a cleaner look.
6 - **Clarity**: Expanded on the "muddiest point" to form a more specific question
   about reconciling the Elbow Method and Silhouette Scores, and refined the "
   linkage" section to be more direct.
7
8 I have now updated the file with your polished responses.
9     ''
10    ''
11    ''
12
13
14 Made changes.
15
16 I have updated the document with your answers.
17
18 Here are the changes I made:
19 - **Polished Content**: I've gently refined the language in your answers to
   improve clarity and flow while preserving your original meaning and tone.
20 - **Formatting**: The lists of unsupervised models have been simplified for better
   readability within the LaTeX document.
21 - **Word Count**: I've added an approximate word count to the title page as
   requested.
22
23 The document is now updated with your weekly check-in responses.

```

My Prompt to Copilot's LLM

```

1 you changed a lot! i just wanted a gentle rewording of some parts where its unsaid
   clearly

```

Copilot's LLM Response

```

1 You are absolutely right. My apologies, I overdid the polishing. I have reverted
   the content to be much closer to your original text, making only minor
   adjustments for clarity as you requested.
2
3 I will now apply these more subtle changes to the '0_Lab03_CheckIn.tex' file.
4     ''
5     ''
6     ''

```

```

7
8
9 Made changes.
10
11 My apologies for over-editing. I've revised the '0_Lab03_CheckIn.tex' file again,
     this time keeping the changes minimal.
12
13 I've reverted the content to be much closer to your original wording, only making
     gentle adjustments for clarity where needed and tidying up the formatting of
     your lists. The structure and core of your answers are now as you wrote them.

```

My Prompt to Copilot's LLM

```

1 you're shitty. just use this and latexify it
2 Here's your text with gentle corrections and enhancements:
3 Question 2 | Plans: What are your goals for the next week? How does this match
     against your overall project plan?
4 For Unsupervised Learning, I plan to briefly discuss K-means with my friend who
     just finished writing a mathematical proof of K-means and review his work on
     that, and go through my Week 3 notes and Lab notebook to solidify K-means and
     basic unsupervised learning concepts. From there, I want to explore and review
     more of the unsupervised models and techniques listed below (briefly, and
     expand if time allows):
5 Clustering Models
6
7 Hierarchical Clustering (Agglomerative / Divisive)
8 DBSCAN (Density-Based Spatial Clustering of Applications with Noise)
9 OPTICS (Ordering Points To Identify the Clustering Structure)
10 Mean Shift Clustering
11 Gaussian Mixture Models (GMMs) (probabilistic clustering)
12 Spectral Clustering (uses eigenvectors of similarity matrix)
13 Affinity Propagation
14
15
16 Dimensionality Reduction / Feature Learning
17
18 PCA (Principal Component Analysis)
19 Kernel PCA
20 t-SNE (t-Distributed Stochastic Neighbor Embedding)
21 UMAP (Uniform Manifold Approximation and Projection)
22 Autoencoders (neural network-based)
23 Factor Analysis
24 Independent Component Analysis (ICA)
25 Non-Negative Matrix Factorization (NMF)
26
27
28 Density Estimation / Generative Models
29
30 Gaussian Mixture Models (GMMs)
31 Kernel Density Estimation (KDE)
32 Variational Autoencoders (VAE)
33 Generative Adversarial Networks (GANs) (unsupervised / self-supervised flavor)
34 Restricted Boltzmann Machines (RBM)
35
36
37 Association Rule Learning
38
39 Apriori Algorithm (frequent itemsets)
40 FP-Growth (Frequent Pattern Growth)
41

```

```

42
43 Neural & Advanced Methods
44
45 Self-Organizing Maps (SOMs) (Kohonen maps)
46 Deep Embedding Clustering (DEC)
47 Contrastive / Self-Supervised Learning (SimCLR, BYOL, MoCo      often fall between
     unsupervised and supervised)
48
49 Aside from this topic, I need to look into my rich and dense financial data
     dataset for the Saudi Arabian Monetary Authority (SAMA), which I have compiled
     and built into a PostgreSQL database for my different project on Model
     Context Protocol (MCP) tools, and study its applicability and appropriateness
     for the project. I need to look into dataset preparation and perhaps publish
     it.
50 Question 3 | Study: List and discuss the three main points of the material covered
     this week.
51
52 Clustering (K-means, SSE, iterative algorithm of updating centroids by mean
     recalculation)
53 Clustering evaluation (Inertia and SSE relation to be minimized, Elbow Curve -
     variance of clusters)
54 Clustering optimization (Silhouette Score - average line, silhouette thickness)
55
56 Question 4 | Insights: What was the single most confusing concept ('muddiest point
     ') for you in this week's material, and what specific question do you have
     about it?
57 Silhouette analysis and tying it up with the Elbow curve and the concept of
     variance in clusters.
58 Question 5 | Linkage: How could one key idea from this week's material be directly
     applied to your project?
59 If I were to use my SAMA database as a dataset, which is labeled time series data,
     we can derive different labels for different classification/regression
     objectives we deem statistically relevant. Also, since the dataset is huge and
     highly dimensional, we might also benefit from dimensionality reduction
     techniques (PCA, etc.).
60
61 --
62 p.s dont add questin #

```

Copilot's LLM Response

```

1 My apologies for the previous edits. I understand now that you wanted a direct "
     latexification" of the text you provided. I have updated the '0_Lab03_CheckIn.
     tex' file accordingly, using your exact wording and formatting it with
     appropriate LaTeX commands.
2
3 Here is what I've done:
4 - The content for "Plans," "Study," "Insights," and "Linkage" has been replaced
     with the precise text you provided.
5 - The lists of machine learning models are now formatted using nested lists in
     LaTeX for clear, hierarchical presentation.
6 - The other points are formatted as you wrote them, without the section numbering.
7
8 The document should now accurately reflect your input.
9 """
10 """
11 """
12
13
14 Made changes.

```

```

15
16 You're right. My apologies. I have updated the '0_Lab03_CheckIn.tex' file to use
    the exact text you provided, applying only the necessary LaTeX formatting for
    the lists and section headings. No wording has been changed.

```

My Prompt to Copilot's LLM

```

1 can't u use a simpler approach and more straightforward and concise

```

Copilot's LLM Response

```

1 You're absolutely right! Let me give you a much simpler and more straightforward
    approach for centroid-based classification. I'll add this right after your
    existing classification code:
2
3 ''''
4 # =====
5 # SIMPLE CENTROID-BASED CLASSIFICATION WITH K=6
6 # =====
7
8 print("== Simple Centroid-Based Classification ==\n")
9
10 from sklearn.model_selection import train_test_split
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.cluster import KMeans
13 from sklearn.metrics import accuracy_score, confusion_matrix
14 import numpy as np
15
16 # Load wine data
17 from sklearn.datasets import load_wine
18 wine = load_wine()
19 X, y_true = wine.data, wine.target
20
21 # Scale features
22 scaler = StandardScaler()
23 X_scaled = scaler.fit_transform(X)
24
25 # Split data
26 X_train, X_test, y_train, y_test = train_test_split(
27     X_scaled, y_true, test_size=0.3, random_state=42, stratify=y_true
28 )
29
30 print(f"Train: {len(X_train)} samples, Test: {len(X_test)} samples")
31
32 # Train K-Means with K=6
33 kmeans = KMeans(n_clusters=6, random_state=42)
34 train_clusters = kmeans.fit_predict(X_train)
35
36 # Classify test data using nearest centroid
37 test_clusters = kmeans.predict(X_test)
38
39 print(f"Found {len(kmeans.cluster_centers_)} centroids")
40 print(f"Test data assigned to clusters: {np.unique(test_clusters)}")
41 '''
42
43 ''''
44 # =====
45 # MAP CLUSTERS TO WINE CLASSES (SIMPLE VERSION)

```

```

46 # =====
47
48 # Create cluster-to-class mapping using majority vote
49 cluster_to_class = {}
50 for cluster in range(6):
51     mask = train_clusters == cluster
52     if np.sum(mask) > 0:
53         # Assign cluster to most common class in that cluster
54         cluster_to_class[cluster] = np.bincount(y_train[mask]).argmax()
55     else:
56         cluster_to_class[cluster] = 0 # Default
57
58 print("Cluster      Wine Class mapping:")
59 for cluster, wine_class in cluster_to_class.items():
60     print(f"  Cluster {cluster}      Class {wine_class}")
61
62 # Convert test cluster predictions to class predictions
63 y_pred = [cluster_to_class[cluster] for cluster in test_clusters]
64
65 # Calculate accuracy
66 accuracy = accuracy_score(y_test, y_pred)
67 print(f"\nClassification Accuracy: {accuracy:.3f} ({accuracy*100:.1f}%)")
68
69 # Show confusion matrix
70 cm = confusion_matrix(y_test, y_pred)
71 print("\nConfusion Matrix:")
72 print(cm)
73 """
74
75 """
76
77 # =====
78 # POTENTIAL IMPROVEMENTS (SIMPLE LIST)
79 # =====
80
81 print("\n" + "="*50)
82 print("POTENTIAL IMPROVEMENTS")
83 print("="*50)
84
85 improvements = [
86     "1. Try different K values (3, 4, 5, 7, 8)",
87     "2. Use different scaling (MinMaxScaler, RobustScaler)",
88     "3. Apply PCA before clustering to reduce dimensions",
89     "4. Use different distance metrics (Manhattan, Cosine)",
90     "5. Try ensemble: average predictions from multiple runs",
91     "6. Handle class imbalance with weighted clustering",
92     "7. Use Gaussian Mixture Model for soft clustering",
93     "8. Add feature selection to use only important features"
94 ]
95
96 for improvement in improvements:
97     print(f"  {improvement}")
98
99 print(f"\nCurrent result: {accuracy:.3f} accuracy with 6 clusters")
100 print("Next step: Try K=3 (matches true wine classes) for comparison")
101 """
102 This approach is:
103 - **Much shorter** (3 code blocks vs 8+)
104 - **Straightforward** (load scale split cluster map evaluate)
105 - **Concise** (focuses on core concepts)
106 - **Easy to understand** (simple majority vote mapping)
107 - **Practical** (gives clear next steps)

```

```

108
109 The key insight is that we don't need complex Hungarian algorithms or confidence
  analysis - just map each cluster to whatever wine class appears most often in
  that cluster during training!

```

My Prompt to Copilot's LLM

```

1 so i already did the clustering using elbow and it was k= 6, how do we do this:
2
3 and then use cluster
4 centroids to classify new, unseen data points. Evaluate the classification
  accuracy and discuss any potential improvements.

```

Copilot's LLM Response

```

1 Great! Since you've already determined that the optimal number of clusters is K=6
  using the elbow method, I'll help you implement the centroid-based
  classification for new data points. Let me add code to your current notebook
  that will:
2
3 1. Use your existing K=6 clustering results
4 2. Create a classification system using the cluster centroids
5 3. Evaluate the accuracy
6 4. Discuss potential improvements
7
8 ''''
9 # =====
10 # CENTROID-BASED CLASSIFICATION WITH K=6
11 # =====
12
13 print("== Centroid-Based Classification with K=6 ==\n")
14
15 # First, let's perform K-Means with K=6 to get the centroids
16 from sklearn.cluster import KMeans
17 from sklearn.model_selection import train_test_split
18 from sklearn.metrics import accuracy_score, classification_report, confusion_
  matrix
19 from sklearn.preprocessing import StandardScaler
20 from scipy.optimize import linear_sum_assignment
21 import numpy as np
22 import matplotlib.pyplot as plt
23 import seaborn as sns
24
25 # Load and prepare the Wine dataset
26 from sklearn.datasets import load_wine
27 wine = load_wine()
28 X = wine.data
29 y_true = wine.target
30
31 # Scale the features (important for K-Means)
32 scaler = StandardScaler()
33 X_scaled = scaler.fit_transform(X)
34
35 # Split the data into training and testing sets
36 X_train, X_test, y_train, y_test = train_test_split(X_scaled, y_true, test_size
  =0.3, random_state=42, stratify=y_true)
37
38 print(f"Training set size: {X_train.shape[0]} samples")

```

```

39 print(f"Test set size: {X_test.shape[0]} samples")
40 print(f"Number of classes: {len(np.unique(y_true))}")
41 print(f"Number of clusters: 6\n")
42
43 # Fit K-Means on training data with K=6
44 kmeans = KMeans(n_clusters=6, random_state=42, n_init=10)
45 train_clusters = kmeans.fit_predict(X_train)
46 centroids = kmeans.cluster_centers_
47
48 print("K-Means clustering completed on training data.")
49 print(f"Cluster centroids shape: {centroids.shape}\n")
50 """
51
52 '''python
53 # =====
54 # CLASSIFICATION FUNCTION USING CENTROIDS
55 # =====
56
57 def classify_using_centroids(X_new, centroids, scaler):
58     """
59         Classify new data points by assigning them to the nearest centroid
60     """
61     # Scale the new data using the same scaler
62     X_new_scaled = scaler.transform(X_new)
63
64     # Calculate distances to all centroids
65     distances = np.sqrt(((X_new_scaled[:, np.newaxis, :] - centroids) ** 2).sum(
66     axis=2))
67
68     # Assign to nearest centroid
69     predicted_clusters = np.argmin(distances, axis=1)
70
71     return predicted_clusters, distances
72
73 # Classify the test data using centroids
74 test_clusters, test_distances = classify_using_centroids(X_test, centroids, scaler
    )
75
76 print("Test data classified using centroids.")
77 print(f"Predicted clusters for test data: {np.unique(test_clusters)}")
78 print(f"Cluster distribution in test data: {np.bincount(test_clusters)}\n")
79 """
80
81 '''python
82 # =====
83 # MAP CLUSTERS TO WINE CLASSES
84 # =====
85
86 def map_clusters_to_classes(train_clusters, y_train, n_clusters=6, n_classes=3):
87     """
88         Map clusters to true classes using Hungarian algorithm for optimal assignment
89     """
90     # Create confusion matrix between clusters and true classes
91     cluster_class_matrix = np.zeros((n_clusters, n_classes))
92
93     for cluster in range(n_clusters):
94         cluster_mask = train_clusters == cluster
95         if np.sum(cluster_mask) > 0:
96             class_counts = np.bincount(y_train[cluster_mask], minlength=n_classes)
97             cluster_class_matrix[cluster] = class_counts
98
99     print("Cluster-to-Class Matrix (Training Data):")

```

```

99     print("Rows: Clusters, Columns: Wine Classes")
100    print(cluster_class_matrix.astype(int))
101    print()
102
103    # Use Hungarian algorithm to find optimal cluster-to-class mapping
104    # We need to handle the case where we have more clusters than classes
105    if n_clusters > n_classes:
106        # Create an extended cost matrix
107        cost_matrix = -cluster_class_matrix  # Negative because we want to
108        maximize
109        row_indices, col_indices = linear_sum_assignment(cost_matrix)
110
111        # Create the mapping
112        cluster_to_class = {}
113        for cluster_idx, class_idx in zip(row_indices, col_indices):
114            cluster_to_class[cluster_idx] = class_idx
115
116        # Handle unmapped clusters - assign them to the most frequent class in the
117        # cluster
118        for cluster in range(n_clusters):
119            if cluster not in cluster_to_class:
120                if np.sum(cluster_class_matrix[cluster]) > 0:
121                    cluster_to_class[cluster] = np.argmax(cluster_class_matrix[
122                        cluster])
123                else:
124                    cluster_to_class[cluster] = 0  # Default assignment
125            else:
126                # Standard case: equal or fewer clusters than classes
127                cost_matrix = -cluster_class_matrix
128                row_indices, col_indices = linear_sum_assignment(cost_matrix)
129                cluster_to_class = dict(zip(row_indices, col_indices))
130
131    return cluster_to_class, cluster_class_matrix
132
133 # Map clusters to classes
134 cluster_to_class_mapping, cluster_class_matrix = map_clusters_to_classes(train_
135     clusters, y_train)
136
137 print("Optimal Cluster-to-Class Mapping:")
138 for cluster, wine_class in cluster_to_class_mapping.items():
139     print(f"Cluster {cluster} -> Wine Class {wine_class}")
140 print()
141 """
142 """
143
144 def evaluate_classification(test_clusters, y_test, cluster_to_class_mapping):
145     """
146         Evaluate the classification accuracy using the cluster-to-class mapping
147     """
148     # Convert cluster predictions to class predictions
149     y_pred = np.array([cluster_to_class_mapping[cluster] for cluster in test_
150         clusters])
151
152     # Calculate accuracy
153     accuracy = accuracy_score(y_test, y_pred)
154
155     return y_pred, accuracy

```

```

156 # Evaluate classification performance
157 y_pred, accuracy = evaluate_classification(test_clusters, y_test, cluster_to_class
    _mapping)
158
159 print("== CLASSIFICATION RESULTS ==")
160 print(f"Classification Accuracy: {accuracy:.4f} ({accuracy*100:.2f}%)")
161 print()
162
163 # Detailed classification report
164 print("Detailed Classification Report:")
165 print(classification_report(y_test, y_pred, target_names=wine.target_names))
166
167 # Confusion Matrix
168 cm = confusion_matrix(y_test, y_pred)
169 print("\nConfusion Matrix:")
170 print("True Classes (rows) vs Predicted Classes (columns)")
171 print(cm)
172 """
173
174 """
175 # =====
176 # VISUALIZE RESULTS
177 # =====
178
179 # Create visualization of results
180 fig, axes = plt.subplots(2, 2, figsize=(15, 12))
181
182 # 1. Confusion Matrix Heatmap
183 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues',
184             xticklabels=wine.target_names,
185             yticklabels=wine.target_names, ax=axes[0,0])
186 axes[0,0].set_title('Confusion Matrix\n(Centroid-based Classification)')
187 axes[0,0].set_xlabel('Predicted Class')
188 axes[0,0].set_ylabel('True Class')
189
190 # 2. Cluster assignments for test data (using PCA for visualization)
191 from sklearn.decomposition import PCA
192 pca = PCA(n_components=2)
193 X_test_pca = pca.fit_transform(X_test)
194
195 scatter = axes[0,1].scatter(X_test_pca[:, 0], X_test_pca[:, 1],
196                             c=test_clusters, cmap='tab10', alpha=0.7)
197 axes[0,1].set_title('Test Data Cluster Assignments\n(PCA Visualization)')
198 axes[0,1].set_xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.3f} variance)')
199 axes[0,1].set_ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.3f} variance)')
200 plt.colorbar(scatter, ax=axes[0,1], label='Cluster')
201
202 # 3. True wine classes for test data
203 scatter2 = axes[1,0].scatter(X_test_pca[:, 0], X_test_pca[:, 1],
204                             c=y_test, cmap='viridis', alpha=0.7)
205 axes[1,0].set_title('Test Data True Classes\n(PCA Visualization)')
206 axes[1,0].set_xlabel(f'PC1 ({pca.explained_variance_ratio_[0]:.3f} variance)')
207 axes[1,0].set_ylabel(f'PC2 ({pca.explained_variance_ratio_[1]:.3f} variance)')
208 plt.colorbar(scatter2, ax=axes[1,0], label='Wine Class')
209
210 # 4. Cluster-to-Class mapping visualization
211 cluster_class_df = pd.DataFrame(cluster_class_matrix,
212                                     columns=wine.target_names,
213                                     index=[f'Cluster {i}' for i in range(6)])
214 sns.heatmap(cluster_class_df, annot=True, fmt='d', cmap='Oranges', ax=axes[1,1])
215 axes[1,1].set_title('Cluster-to-Class Distribution\n(Training Data)')
216 axes[1,1].set_xlabel('Wine Classes')

```

```

217 axes[1,1].set_ylabel('Clusters')
218
219 plt.tight_layout()
220 plt.show()
221 """
222 """
223 """python
224 # =====
225 # DISTANCE-BASED CONFIDENCE ANALYSIS
226 # =====
227
228 print("\n== CONFIDENCE ANALYSIS ==")
229
230 # Calculate confidence based on distance to nearest vs second-nearest centroid
231 def calculate_confidence(distances):
232 """
233     Calculate classification confidence based on distance ratios
234 """
235     sorted_distances = np.sort(distances, axis=1)
236     nearest_dist = sorted_distances[:, 0]
237     second_nearest_dist = sorted_distances[:, 1]
238
239     # Confidence as ratio of distances (higher ratio = more confident)
240     confidence = second_nearest_dist / (nearest_dist + 1e-8)
241     return confidence
242
243 confidence_scores = calculate_confidence(test_distances)
244
245 print(f"Average confidence score: {np.mean(confidence_scores):.4f}")
246 print(f"Confidence score std: {np.std(confidence_scores):.4f}")
247 print(f"Min confidence: {np.min(confidence_scores):.4f}")
248 print(f"Max confidence: {np.max(confidence_scores):.4f}")
249
250 # Analyze accuracy by confidence level
251 confidence_threshold = np.percentile(confidence_scores, 75) # Top 25% most
252     confident
253 high_confidence_mask = confidence_scores > confidence_threshold
254 low_confidence_mask = confidence_scores <= confidence_threshold
255
256 high_conf_accuracy = accuracy_score(y_test[high_confidence_mask], y_pred[high_
257     confidence_mask])
258 low_conf_accuracy = accuracy_score(y_test[low_confidence_mask], y_pred[low_
259     confidence_mask])
260
261 print(f"\nHigh confidence predictions (>{confidence_threshold:.3f}): {np.sum(high_
262     confidence_mask)} samples")
263 print(f"High confidence accuracy: {high_conf_accuracy:.4f} ({high_conf_accuracy
264     *100:.2f}%)")
265 print(f"\nLow confidence predictions (  {confidence_threshold:.3f}): {np.sum(low_
266     confidence_mask)} samples")
267 print(f"Low confidence accuracy: {low_conf_accuracy:.4f} ({low_conf_accuracy
268     *100:.2f}%)")
269 """
270 """
271 """python
272 # =====
273 # POTENTIAL IMPROVEMENTS DISCUSSION
274 # =====
275
276 print("\n" + "="*60)
277 print("POTENTIAL IMPROVEMENTS DISCUSSION")
278 print("=".*60)

```

```

272
273 improvements = {
274     "1. Feature Engineering": [
275         "    Apply feature selection techniques (PCA, SelectKBest)",
276         "    Create polynomial features or feature interactions",
277         "    Use domain-specific wine chemistry features"
278     ],
279
280     "2. Alternative Distance Metrics": [
281         "    Use Mahalanobis distance instead of Euclidean",
282         "    Implement weighted distance based on feature importance",
283         "    Try cosine similarity for normalized features"
284     ],
285
286     "3. Ensemble Methods": [
287         "    Combine multiple K-Means runs with different initializations",
288         "    Use voting from multiple clustering algorithms",
289         "    Implement consensus clustering"
290     ],
291
292     "4. Outlier Handling": [
293         "    Detect and handle outliers before clustering",
294         "    Use robust clustering methods (K-Medoids)",
295         "    Apply DBSCAN for density-based clustering"
296     ],
297
298     "5. Probabilistic Approaches": [
299         "    Use Gaussian Mixture Models for soft clustering",
300         "    Implement Bayesian classification with cluster priors",
301         "    Add uncertainty quantification to predictions"
302     ],
303
304     "6. Validation Strategies": [
305         "    Use cross-validation for more robust evaluation",
306         "    Implement stratified sampling for better train/test splits",
307         "    Test on completely separate validation dataset"
308     ]
309 }
310
311 for category, items in improvements.items():
312     print(f"\n{category}:")
313     for item in items:
314         print(f"  {item}")
315
316 print(f"\n{'='*60}")
317 print("SUMMARY")
318 print("=".*60)
319 print(f"  Current accuracy with K=6 centroids: {accuracy:.4f} ({accuracy*100:.2f} %)")
320 print(f"  High confidence predictions achieve: {high_conf_accuracy:.4f} ({high_
321     conf_accuracy*100:.2f}%) {accuracy}")
322 print(f"  The method successfully maps {len(cluster_to_class_mapping)} clusters
323     to {len(wine.target_names)} wine classes")
324 print(f"  Classification confidence varies from {np.min(confidence_scores):.3f}
325     to {np.max(confidence_scores):.3f}")
326
327 This comprehensive implementation will:
328
329 1. **Use your K=6 clustering** to create centroids from training data
330 2. **Classify new data points** by finding the nearest centroid
331 3. **Map clusters to wine classes** using the Hungarian algorithm for optimal

```

```
assignment  
330 4. **Evaluate accuracy** with detailed metrics and visualizations  
331 5. **Analyze confidence** based on distance ratios  
332 6. **Discuss improvements** with specific, actionable suggestions  
333  
334 The code includes confidence analysis that helps identify which predictions are  
more reliable, and provides a structured discussion of potential improvements  
you can implement.
```