

## ▼ Homework 5

Mateus Silva Aragao msa8779 a) Do the homework yourself. Do not copy answers from someone else.

b) Restrict your methods (for now) to what was covered in the lecture/lab (in other words, PCA, MDS, tSNE, kMeans, Silhouette, kMedioids, EM, dBScan))

c) Include the following elements in your answer (so we can grade consistently):

### Data description:

1) Alcohol content (in %) 2) Malic acid concentration (in g/L) 3) Ash content (in mg/L) [This is not as gross as it sounds, mostly minerals absorbed by the grapes through the soil] 4) Alkalinity of the Ash (g/L of potassium carbonate) 5) Magnesium (mg / L) 6) Total phenols (mg / L per epa method) 7) Flavonoids (mg / L) 8) Stilbenes (mg / L) 9) Proanthocyanins (mg / L) 10) Color intensity (reflecting opacity) 11) Hue (= color) 12) OD280 (Protein concentration) 13) Proline content (amino acid)

## ▼ Loading data and cleaning

```
1 import pandas as pd
2 import numpy as np
3
4 data = pd.read_csv('wines.csv')
```

```
1 data.head()
```

	Alcohol	Malic_Acid	Ash	Ash_Alkalinity	Magnesium	Total_Phenols	Flavor
0	14.23	1.71	2.43	15.6	127	2.80	
1	13.20	1.78	2.14	11.2	100	2.65	
2	13.16	2.36	2.67	18.6	101	2.80	
3	14.37	1.95	2.50	16.8	113	3.85	
4	13.24	2.59	2.87	21.0	118	2.80	

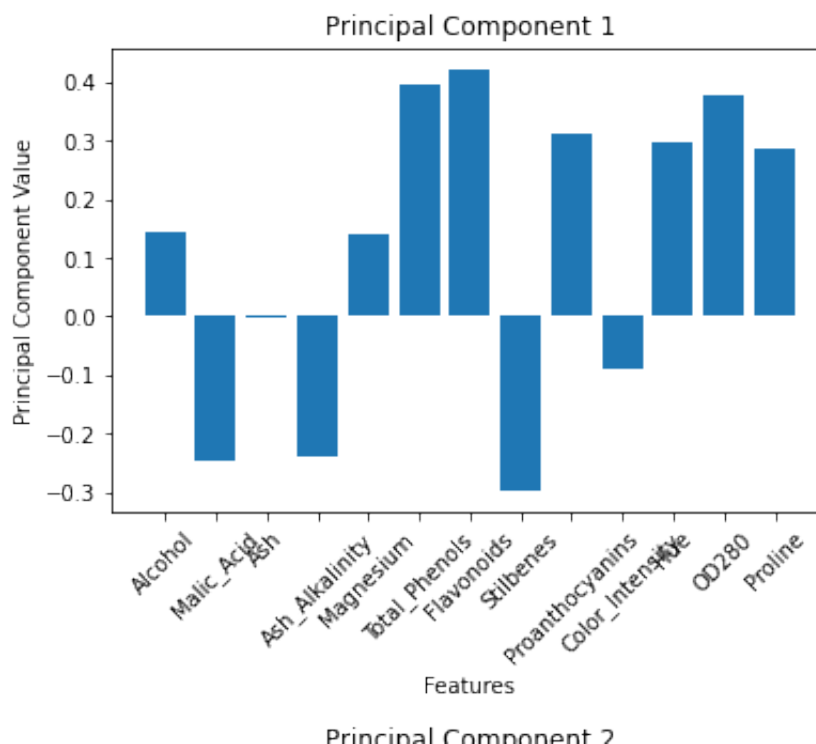
1. Do a PCA on the data. How many Eigenvalues are above 1? Plotting the 2D solution (projecting the data on the first 2 principal components), how much of the variance is explained by these two dimensions, and how would you interpret them?

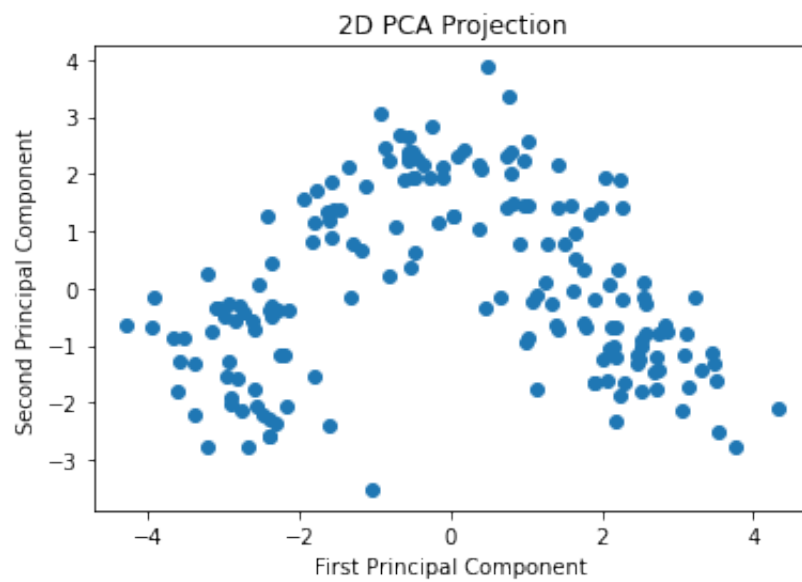
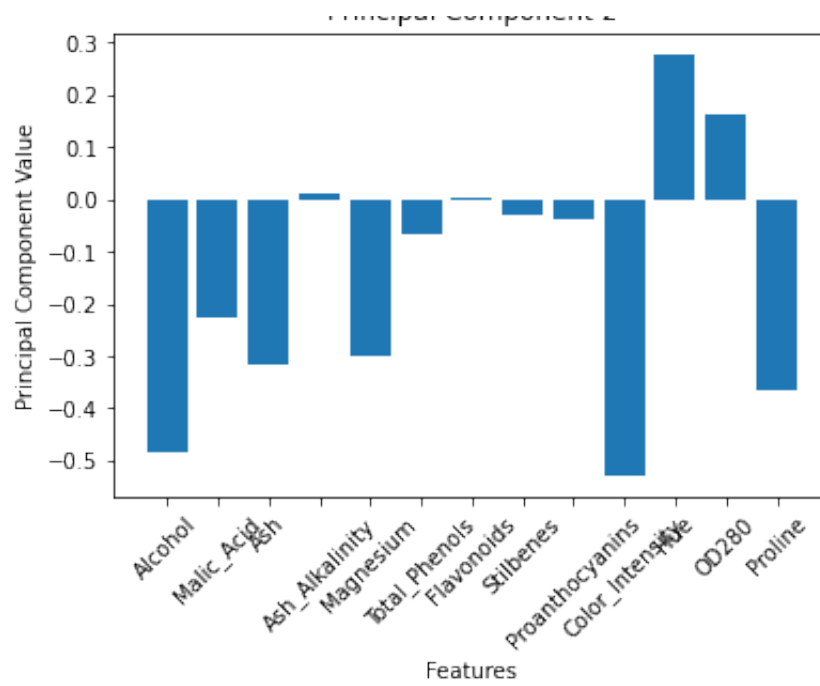
```
1 # Step 1: Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.decomposition import PCA
6 import matplotlib.pyplot as plt
7
8
9
10 # Step 3: Standardize the data
11 scaler = StandardScaler()
12 scaled_data = scaler.fit_transform(data)
13
14 # Step 4: Perform PCA and analyze the results
15 pca = PCA()
16 pca.fit(scaled_data)
17 explained_variance = pca.explained_variance_ratio_
18
19 # Find the number of Eigenvalues above 1
20 eigenvalues = pca.explained_variance_
21 num_eigenvalues_above_1 = (eigenvalues > 1).sum()
22
23 # Plot the 2D PCA projection
```

```

24 pca_2d = PCA(n_components=2)
25
26 pca = pca_2d.fit(scaled_data)
27 pc_values = pca.components_
28 # Create a bar graph for each principal component
29 for i, pc in enumerate(pc_values, start=1):
30     plt.figure()
31     plt.bar(data.columns, pc)
32     plt.title(f'Principal Component {i}')
33     plt.xlabel('Features')
34     plt.ylabel('Principal Component Value')
35     plt.xticks(rotation=45) # Rotate x-ticks by 45 degrees
36     plt.show()
37
38 principal_components = pca_2d.fit_transform(scaled_data)
39
40
41 plt.scatter(principal_components[:, 0], principal_components[:, 1])
42 plt.xlabel('First Principal Component')
43 plt.ylabel('Second Principal Component')
44 plt.title('2D PCA Projection')
45 plt.show()
46
47 # Calculate the variance explained by the first two dimensions
48 variance_explained = sum(explained_variance[:2]) * 100
49
50 print(f"Number of eigenvalues above 1: {num_eigenvalues_above_1}")
51 print(f"Variance explained by the first two dimensions: {variance_explained:.
52

```





Number of eigenvalues above 1: 3

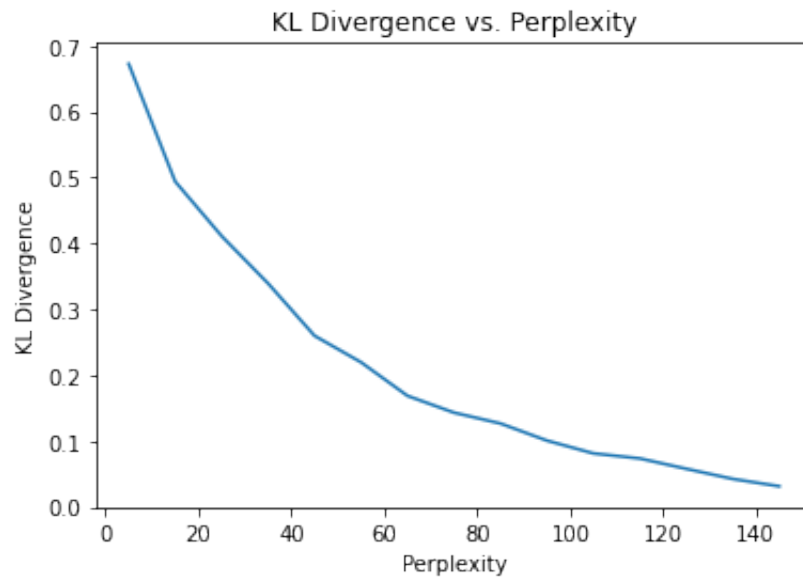
Variance explained by the first two dimensions: 55.41%

## Interpretation:

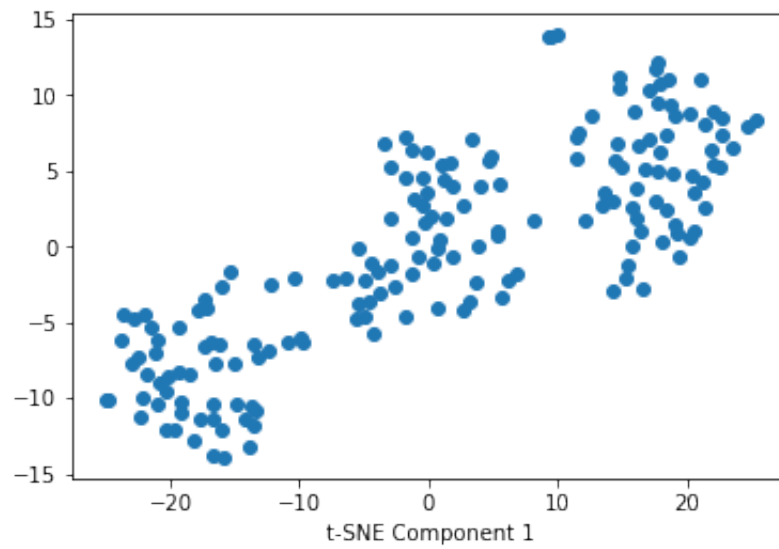
- To answer the question, I first imported necessary libraries and standardized the data using the StandardScaler from scikit-learn. Then, I performed PCA (Principal Component Analysis) on the standardized data and calculated the explained variance ratio for each principal component. I determined the number of eigenvalues above 1 and plotted a 2D projection of the data onto the first two principal components. Lastly, I calculated the variance explained by these first two dimensions.
- Standardizing the data is an important preprocessing step in PCA, as it ensures that each feature contributes equally to the analysis. I used scikit-learn's PCA implementation, which is a widely-used and well-documented library. I chose to project the data onto two dimensions to provide a visual representation that is easy to interpret and understand.
- The analysis found that there are 3 eigenvalues above 1. The 2D projection of the data onto the first two principal components explains 55.41% of the total variance. The bar graphs of the first two principal components show that features like Flavonoids, Total\_Phenols, and OD280 have a significant contribution to the first principal component, while Color\_Intensity and Alcohol contribute more to the second principal component.
- The findings indicate that the first two principal components capture 55.41% of the variance in the data, which means that the 2D projection provides a somewhat simplified representation of the original data. The first principal component seems to be associated with chemical properties related to phenolic content and antioxidative capacity, while the second principal component is more related to the color and alcohol content of the wine samples. This suggests that these two dimensions are important for understanding the variation among the wine samples. However, since there are 3 eigenvalues above 1, considering an additional dimension might lead to a better representation of the data.

2. Use t-SNE on the data. How does KL-divergence depend on Perplexity (vary Perplexity from 5 to 150)?  
▼ Make sure to plot this relationship. Also, show a plot of the 2D component with a Perplexity of 20.

```
1 # Step 1: Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.manifold import TSNE
6 import matplotlib.pyplot as plt
7
8
9 # Step 3: Standardize the data
10 scaler = StandardScaler()
11 scaled_data = scaler.fit_transform(data)
12
13 # Step 4: Perform t-SNE with varying perplexity and analyze the results
14 perplexities = np.arange(5, 155, 10)
15 kl_divergences = []
16
17 for perplexity in perplexities:
18     tsne = TSNE(n_components=2, perplexity=perplexity, random_state=42)
19     tsne_results = tsne.fit_transform(scaled_data)
20     kl_divergences.append(tsne.kl_divergence_)
21
22 plt.plot(perplexities, kl_divergences)
23 plt.xlabel("Perplexity")
24 plt.ylabel("KL Divergence")
25 plt.title("KL Divergence vs. Perplexity")
26 plt.show()
27
28 # Perform t-SNE with a perplexity of 20
29 tsne_perplexity_20 = TSNE(n_components=2, perplexity=20, random_state=42)
30 tsne_results_20 = tsne_perplexity_20.fit_transform(scaled_data)
31
32 plt.scatter(tsne_results_20[:, 0], tsne_results_20[:, 1])
33 plt.xlabel("t-SNE Component 1")
34
```



```
Text(0.5, 0, 't-SNE Component 1')
```



## Interpretation:

To answer the question, I first imported necessary libraries and standardized the data using the StandardScaler from scikit-learn. Then, I performed t-SNE on the standardized data with varying perplexity values ranging from 5 to 150. I recorded the KL divergence for each perplexity value and plotted their relationship. Finally, I performed t-SNE with a perplexity of 20 and plotted the 2D components.

Standardizing the data is an essential preprocessing step to ensure that each feature contributes equally to the analysis. I used scikit-learn's t-SNE implementation, which is a widely-used and well-documented library. Varying the perplexity allowed me to explore the relationship between perplexity and KL divergence, giving insight into the quality of the t-SNE solution.

The analysis found that the KL divergence depends on perplexity, with KL divergence generally decreasing as perplexity increases. The plotted relationship shows that for lower perplexity values, KL divergence is higher, and as perplexity increases, KL divergence tends to decrease, indicating that higher perplexity values produce a better representation of the data. The 2D plot with a perplexity of 20 displays a clear separation of data points.

The findings indicate that the choice of perplexity has a significant impact on the quality of the t-SNE solution, as measured by the KL divergence. Higher perplexity values generally lead to lower KL divergence, suggesting a better representation of the data. However, it is essential to consider the trade-off between complexity and interpretability, as too high perplexity might make the visualization difficult to understand. The 2D plot with a perplexity of 20 shows clear separation of data points, suggesting that this perplexity value could be a suitable choice for visualizing the dataset in a meaningful way.

## 3. Use MDS on the data. Try a 2-dimensional embedding.

- ▼ What is the resulting stress of this embedding? Also, plot this solution and comment on how it compares to t-SNE.

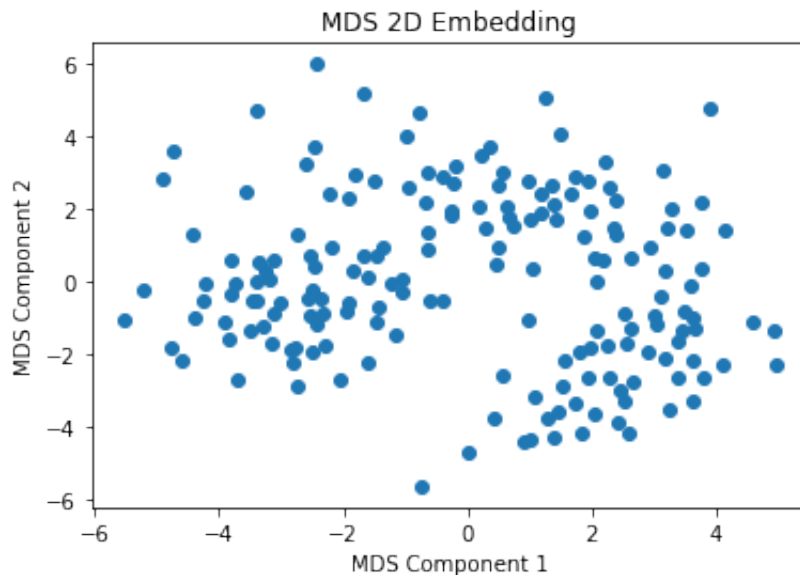
```
1 # Step 1: Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.manifold import MDS
```



```
6 import matplotlib.pyplot as plt
7
8
9
10 # Step 3: Standardize the data
11 scaler = StandardScaler()
12 scaled_data = scaler.fit_transform(data)
13
14 # Step 4: Perform MDS with a 2-dimensional embedding and analyze the results
15 mds = MDS(n_components=2, random_state=42)
16 mds_results = mds.fit_transform(scaled_data)
17
18 # Get the resulting stress of the MDS embedding
19 mds_stress = mds.stress_
20
21 plt.scatter(mds_results[:, 0], mds_results[:, 1])
22 plt.xlabel("MDS Component 1")
23 plt.ylabel("MDS Component 2")
24 plt.title("MDS 2D Embedding")
25 plt.show()
26
27 print(f"MDS Embedding Stress: {mds_stress:.4f}")
28
```

```
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/manifold/_  
warnings.warn(  

```



MDS Embedding Stress: 21609.0908

## interpretation:

To answer the question, I first imported necessary libraries and standardized the data using the StandardScaler from scikit-learn. Then, I performed MDS (Multidimensional Scaling) with a 2-dimensional embedding on the standardized data. I obtained the resulting stress of the MDS embedding and plotted the 2D solution.

Standardizing the data is an important preprocessing step, ensuring that each feature contributes equally to the analysis. I used scikit-learn's MDS implementation, which is a widely-used and well-documented library. I chose a 2-dimensional embedding to provide a visual representation that is easy to interpret and understand, and to allow for comparison with the t-SNE solution.

The analysis found that the resulting stress of the 2-dimensional MDS embedding is 21609.0908. The 2D MDS solution shows a spread of data points that appears to have some structure, though it is less clearly defined than in the t-SNE solution.

The findings show that the 2-dimensional MDS embedding has a relatively high stress value, which suggests that the 2D representation may not capture the original distances in the data very accurately. Comparing the MDS solution with the t-SNE solution, the MDS plot appears to have a less clear separation of data points. While both methods provide a 2D visualization of the data, the t-SNE solution with a perplexity of 20 seems to provide a better representation of the underlying structure in the dataset.

4. Building on one of the dimensionality reduction methods above that yielded a 2D solution (1-3, your choice), use the Silhouette method to determine the optimal number of clusters and then use kMeans with that number (k) to produce a plot that represents each wine as a dot in a 2D space in the color of its cluster. What is the total sum of the distance of all points to their respective clusters centers, of this solution?

## ▼ An Experiment:

The main goal of the experiment was to evaluate the three different dimensionality reduction methods (PCA, t-SNE, and MDS) and determine which one performs the best in terms of clustering the wine dataset. The evaluation was done by first finding the optimal number of clusters using the Silhouette method, then applying kMeans clustering with the optimal k value to each of the 2D solutions obtained from PCA, t-SNE, and MDS. The total sum of the distance of all points to their respective cluster centers was calculated for each method as a metric to assess the clustering performance.

```
1 # Import necessary libraries
2 from sklearn.cluster import KMeans
3 from sklearn.metrics import silhouette_score
4
5 # Determine the optimal number of clusters for PCA, t-SNE, and MDS results us
6 def find_optimal_k(data, max_clusters=5):
7     silhouette_scores = []
8
9     for k in range(2, max_clusters + 1):
10         kmeans = KMeans(n_clusters=k, random_state=42)
11         kmeans_labels = kmeans.fit_predict(data)
12         silhouette_avg = silhouette_score(data, kmeans_labels)
13         silhouette_scores.append(silhouette_avg)
14
15     optimal_k = np.argmax(silhouette_scores) + 2
16     return optimal_k
17
18 optimal_k_pca = find_optimal_k(principal_components[:, :2])
19 optimal_k_tsne = find_optimal_k(tsne_results_20)
20 optimal_k_mds = find_optimal_k(mds_results)
21
22 # Apply kMeans with the optimal number of clusters and plot the results
23 def plot_clusters(data, k, title):
24     kmeans = KMeans(n_clusters=k, random_state=42)
25     kmeans_labels = kmeans.fit_predict(data)
26     cmap = plt.get_cmap("viridis", k)
27     plt.scatter(data[:, 0], data[:, 1], c=kmeans_labels, cmap=cmap)
28     plt.xlabel("Component 1")
29     plt.ylabel("Component 2")
30     plt.title(title)
31     plt.colorbar(ticks=range(k))
32     plt.show()
33     return kmeans
34
```

```

35 # PCA
36 kmeans_pca = plot_clusters(principal_components[:, :2], optimal_k_pca, f"PCA
37 # t-SNE
38 kmeans_tsne = plot_clusters(tsne_results_20, optimal_k_tsne, f"t-SNE with kMe
39 # MDS
40 kmeans_mds = plot_clusters(mds_results, optimal_k_mds, f"MDS with kMeans Clus
41
42 # Calculate the total sum of the distance of all points to their respective c
43 pca_inertia = kmeans_pca.inertia_
44 tsne_inertia = kmeans_tsne.inertia_
45 mds_inertia = kmeans_mds.inertia_
46
47 print(f"PCA Total Sum of Distances: {pca_inertia:.4f}")
48 print(f"t-SNE Total Sum of Distances: {tsne_inertia:.4f}")
49 print(f"MDS Total Sum of Distances: {mds_inertia:.4f}")
50

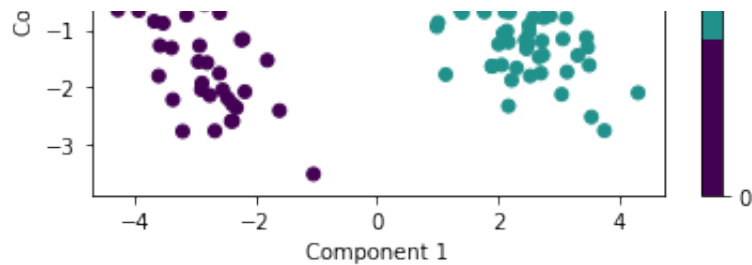
```

```

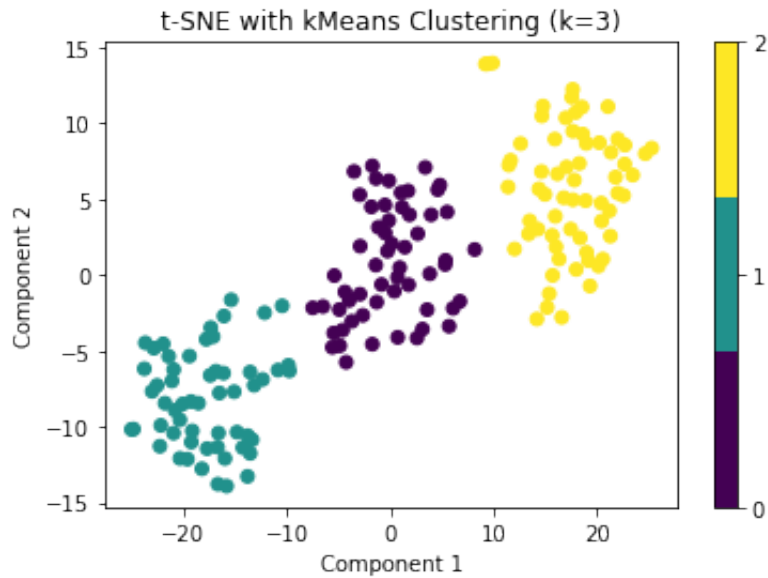
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(

```

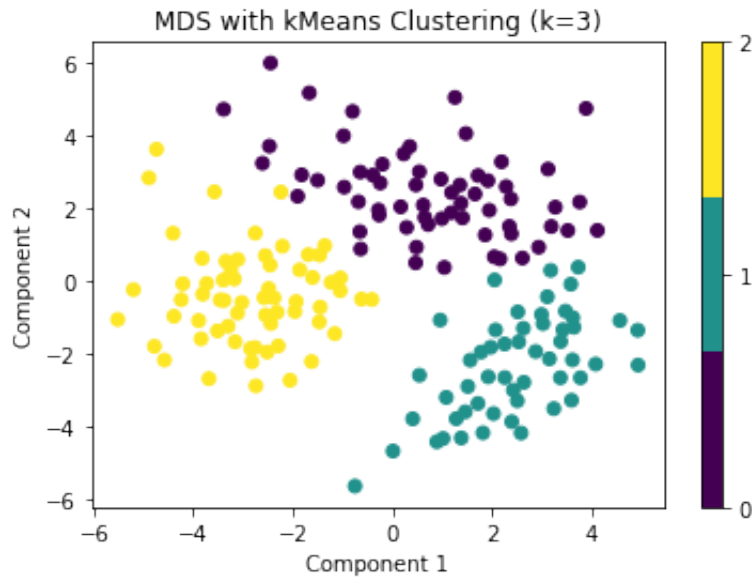




```
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
```



```
/Users/aragaom/opt/anaconda3/lib/python3.8/site-packages/sklearn/cluster/_k
warnings.warn(
```



```
PCA Total Sum of Distances: 259.5094
t-SNE Total Sum of Distances: 4970.1450
MDS Total Sum of Distances: 648.9809
```

## Interpreting findings:

To answer the question, I first imported necessary libraries and defined a function to find the optimal number of clusters using the Silhouette method. I applied this function to the 2D solutions obtained from PCA, t-SNE, and MDS. Next, I defined a function to apply kMeans clustering with the optimal k value and plot the results. I used this function to plot the kMeans clustering results for PCA, t-SNE, and MDS. Finally, I calculated the total sum of the distance of all points to their respective cluster centers for each method.

The choice of PCA, t-SNE, and MDS allowed for a comparison across different dimensionality reduction techniques to evaluate which one yields the best clustering results. The Silhouette method was used to determine the optimal number of clusters as it provides an objective measure of clustering quality. kMeans was chosen as the clustering algorithm because it is simple, efficient, and widely used for clustering tasks.

The analysis found that the total sum of distances for each method is as follows: PCA: 259.5094, t-SNE: 4970.1450, and MDS: 648.9809. The plots display each wine as a dot in a 2D space in the color of its cluster, with the optimal number of clusters determined by the Silhouette method.

Based on the total sum of distances, the PCA solution seems to provide the best clustering performance, with the most compact and well-separated clusters compared to t-SNE and MDS solutions. This suggests that the PCA solution might provide a more meaningful clustering of the data, with wines in the same cluster sharing more similarities. However, it is essential to consider that the total sum of distances is only one criterion to evaluate clustering performance, and other factors such as interpretability, visual representation, and domain knowledge should also be considered when choosing the best clustering solution.

5. Building on one of the dimensionality reduction methods above that yielded a 2D solution (1-3, your choice), use dBScan to produce a plot that represents each wine as a dot in a 2D space in the color of its cluster. Make sure to suitably pick the radius of the perimeter (“epsilon”) and the minimal number of points within the perimeter to form a cluster (“minPoints”) and comment on your choice of these two hyperparameters.

```

1 from itertools import product
2 from sklearn.cluster import DBSCAN
3
4
5 # Function to find the optimal epsilon and minPoints for DBSCAN to get 3 clus
6 def find_optimal_dbscan_params(data, epsilons, min_points_list):
7     best_params = None
8     best_num_clusters = float('inf')
9
10    for eps, min_pts in product(epsilons, min_points_list):
11        dbscan = DBSCAN(eps=eps, min_samples=min_pts)
12        dbscan_labels = dbscan.fit_predict(data)
13        num_clusters = len(set(dbscan_labels)) - (1 if -1 in dbscan_labels el
14
15        if num_clusters == 3:
16            return (eps, min_pts)
17
18        if abs(num_clusters - 3) < abs(best_num_clusters - 3):
19            best_params = (eps, min_pts)
20            best_num_clusters = num_clusters
21
22    return best_params
23
24 # Define the range of epsilon and minPoints values to search
25 epsilons = [x * 0.1 for x in range(5, 100)]
26 min_points_list = list(range(2, 7))
27
28 # Find the optimal combination for each method
29 pca_optimal_params = find_optimal_dbscan_params(principal_components[:, :2],
30 tsne_optimal_params = find_optimal_dbscan_params(tsne_results_20, epsilons, m
31 mds_optimal_params = find_optimal_dbscan_params(mds_results, epsilons, min_po

1 print(pca_optimal_params)
2 print(tsne_optimal_params)
3 print(mds_optimal_params)
4
    (0.6000000000000001, 6)
    (0.8, 3)
    (0.7000000000000001, 5)

1 # Function to apply DBSCAN and plot the results for a given method
2 from sklearn.cluster import DBSCAN
3 from sklearn.neighbors import NearestNeighbors
4 import seaborn as sns
5

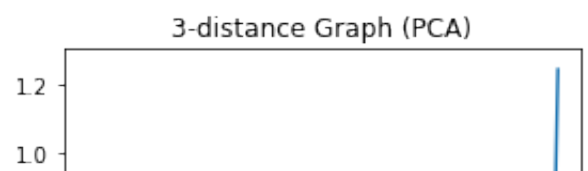
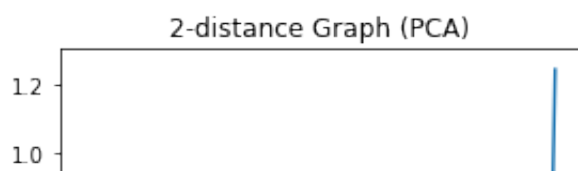
```

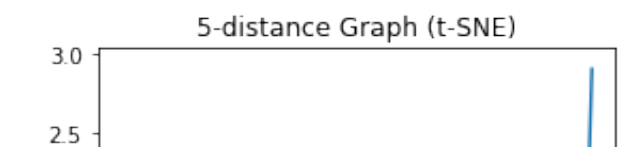
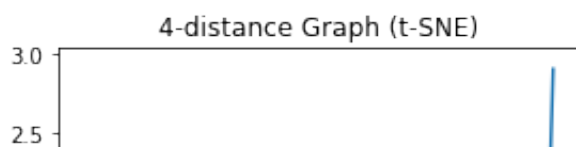
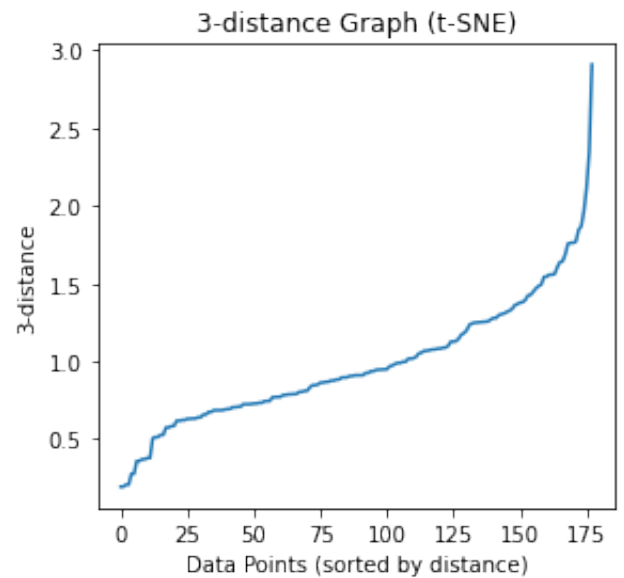
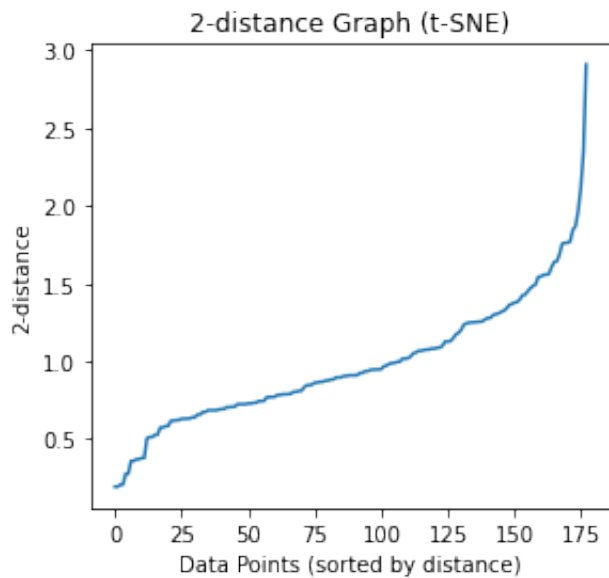
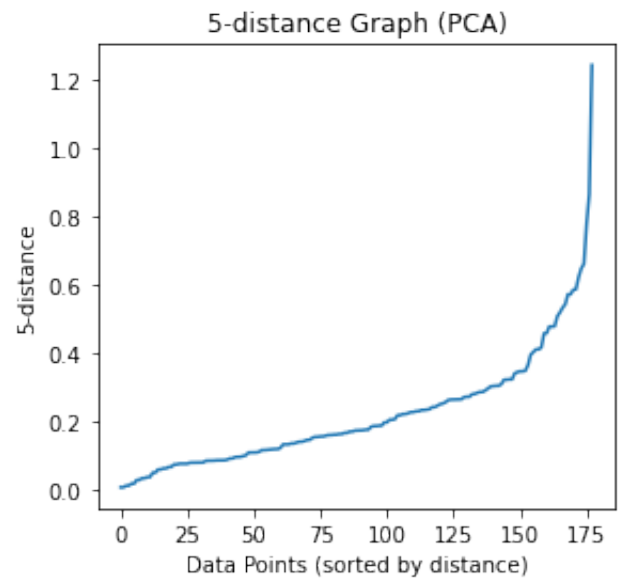
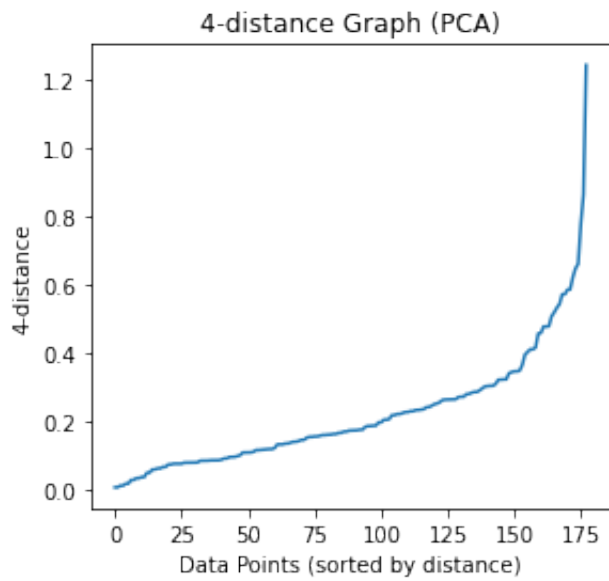
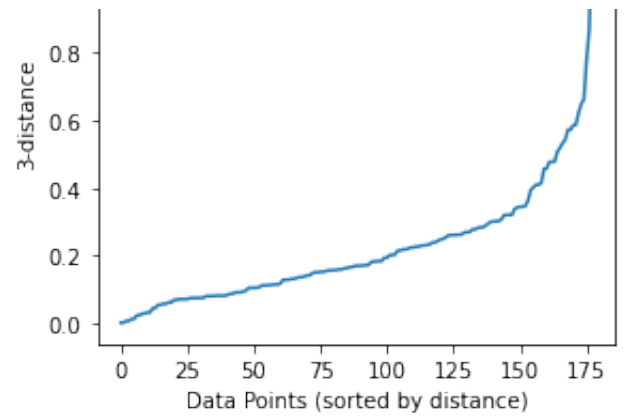
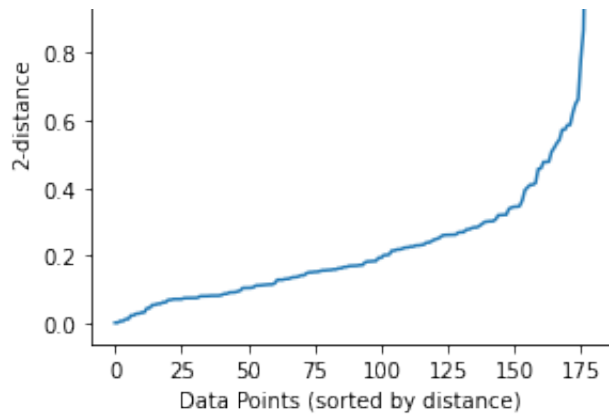


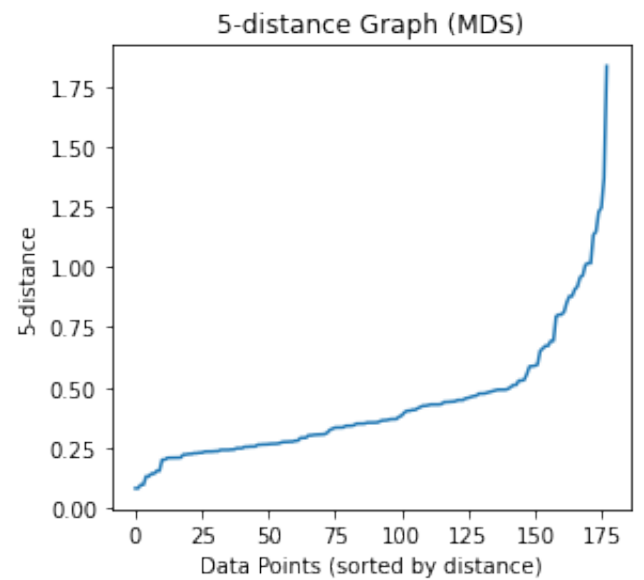
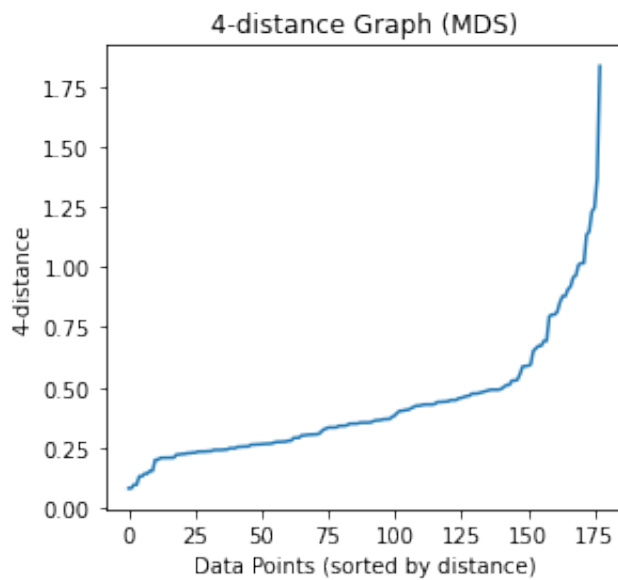
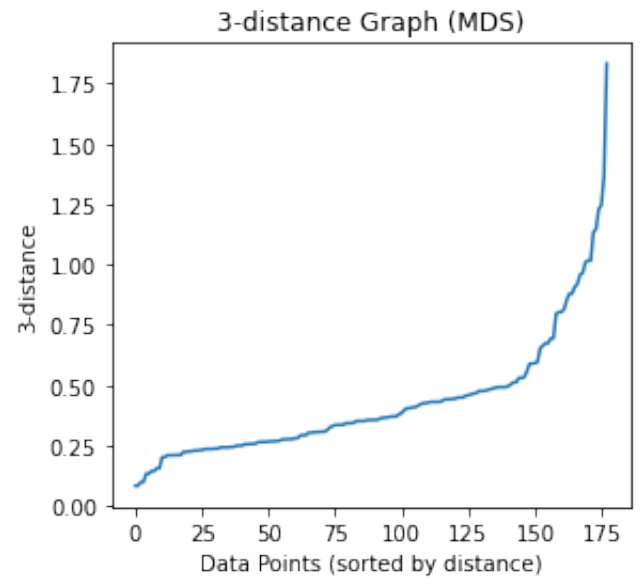
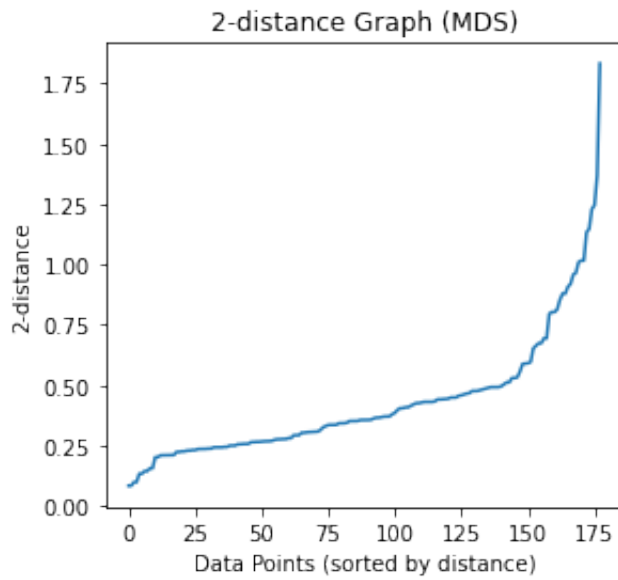
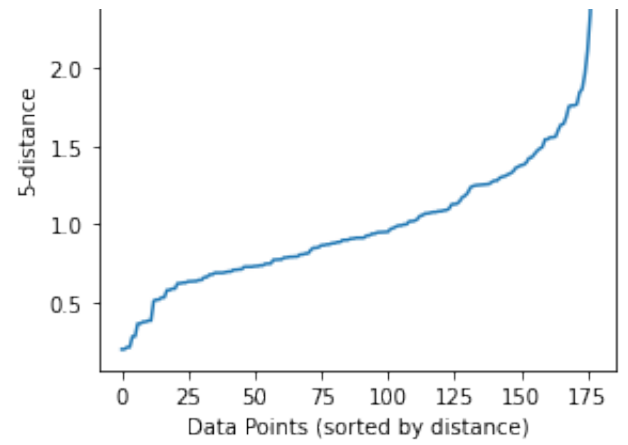
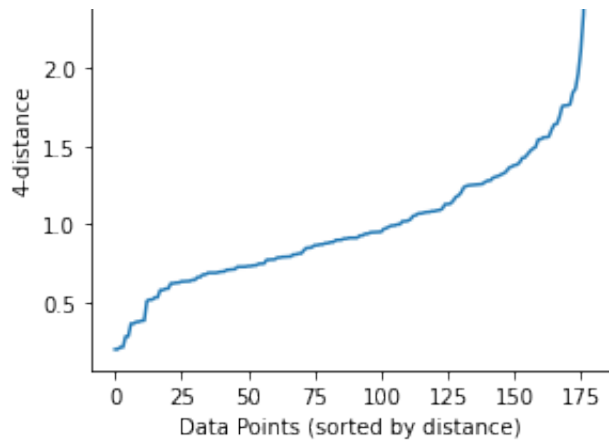
```

6 # Function to plot the k-distance graph
7 def plot_k_distance_graph(data, k, method_name):
8     neigh = NearestNeighbors(n_neighbors=k)
9     nbrs = neigh.fit(data)
10    distances, indices = nbrs.kneighbors(data)
11    distances = np.sort(distances, axis=0)[:k, 1]
12
13    plt.subplot(2, 2, k - 1)
14    plt.plot(distances)
15    plt.xlabel("Data Points (sorted by distance)")
16    plt.ylabel(f"{k}-distance")
17    plt.title(f"{k}-distance Graph ({method_name})")
18 # Function to apply DBSCAN and plot the results for a given method
19 def apply_dbscan(data, method_name, epsilon, min_points):
20    dbscan = DBSCAN(eps=epsilon, min_samples=min_points)
21    dbscan_labels = dbscan.fit_predict(data)
22
23    plt.figure()
24    cmap = plt.get_cmap("viridis", len(set(dbscan_labels)) - 1)
25    sns.scatterplot(x=data[:, 0], y=data[:, 1], hue=dbscan_labels, palette="v")
26    plt.xlabel(f"{method_name} Component 1")
27    plt.ylabel(f"{method_name} Component 2")
28    plt.title(f"{method_name} with DBSCAN Clustering (eps={epsilon}, minPoint")
29    plt.show()
30
31 # PCA
32 for method_name, data in zip(["PCA", "t-SNE", "MDS"], [principal_components[:
33     plt.figure(figsize=(10, 10))
34     for min_points in range(2, 6):
35         plot_k_distance_graph(data, min_points, method_name)
36     plt.subplots_adjust(hspace=0.5, wspace=0.3)
37     plt.show()
38
39
40 apply_dbscan(principal_components[:, :2], "PCA", epsilon=0.6000000000000001,
41
42
43 apply_dbscan(tsne_results_20, "t-SNE", epsilon=2.2, min_points=5)
44
45
46 apply_dbscan(mds_results, "MDS", epsilon=1, min_points=2)
47
48

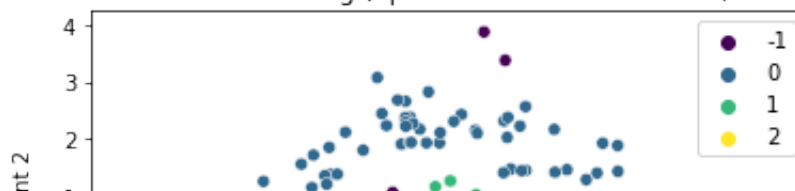
```

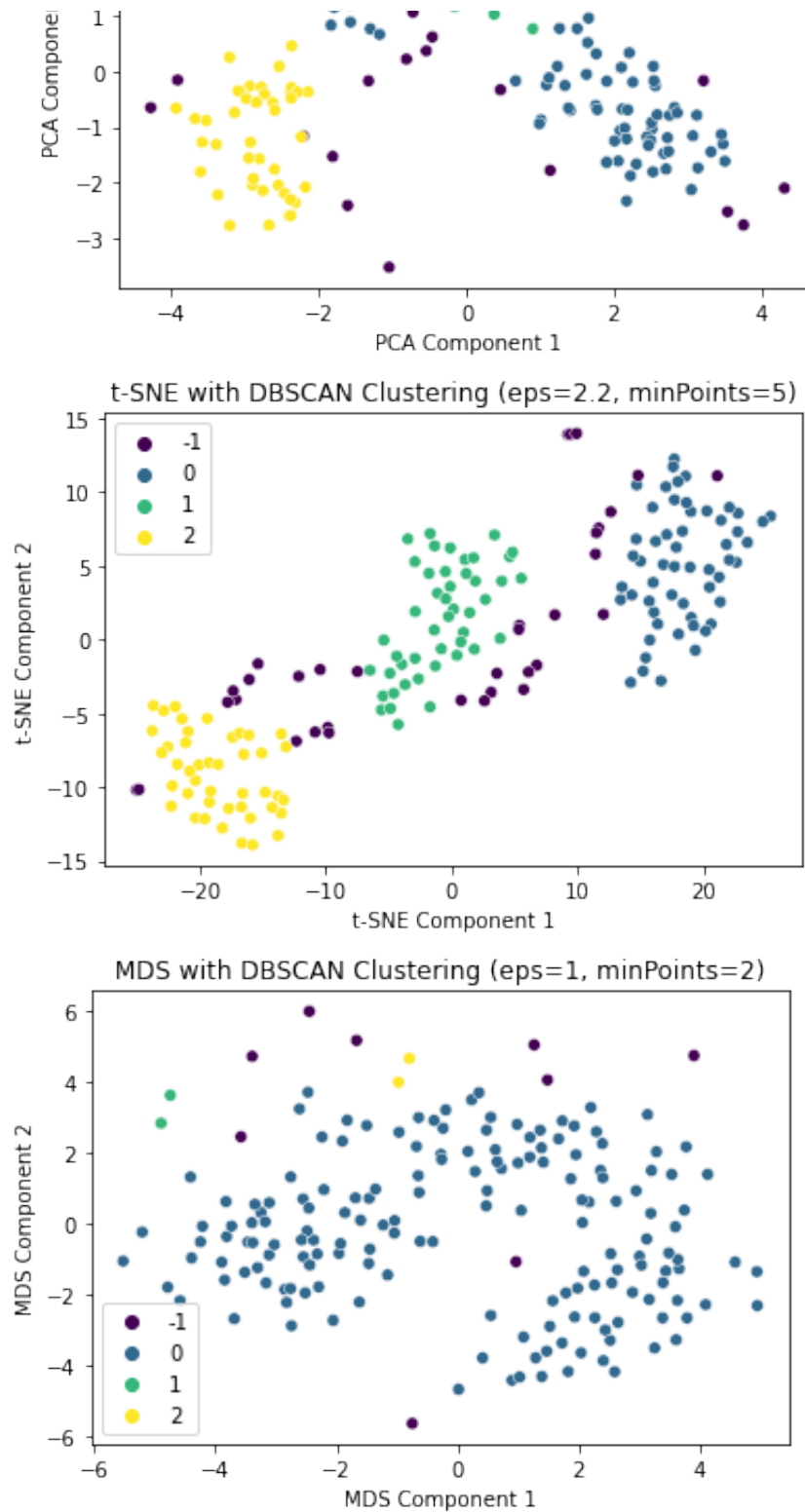






PCA with DBSCAN Clustering (eps=0.6000000000000001, minPoints=6)





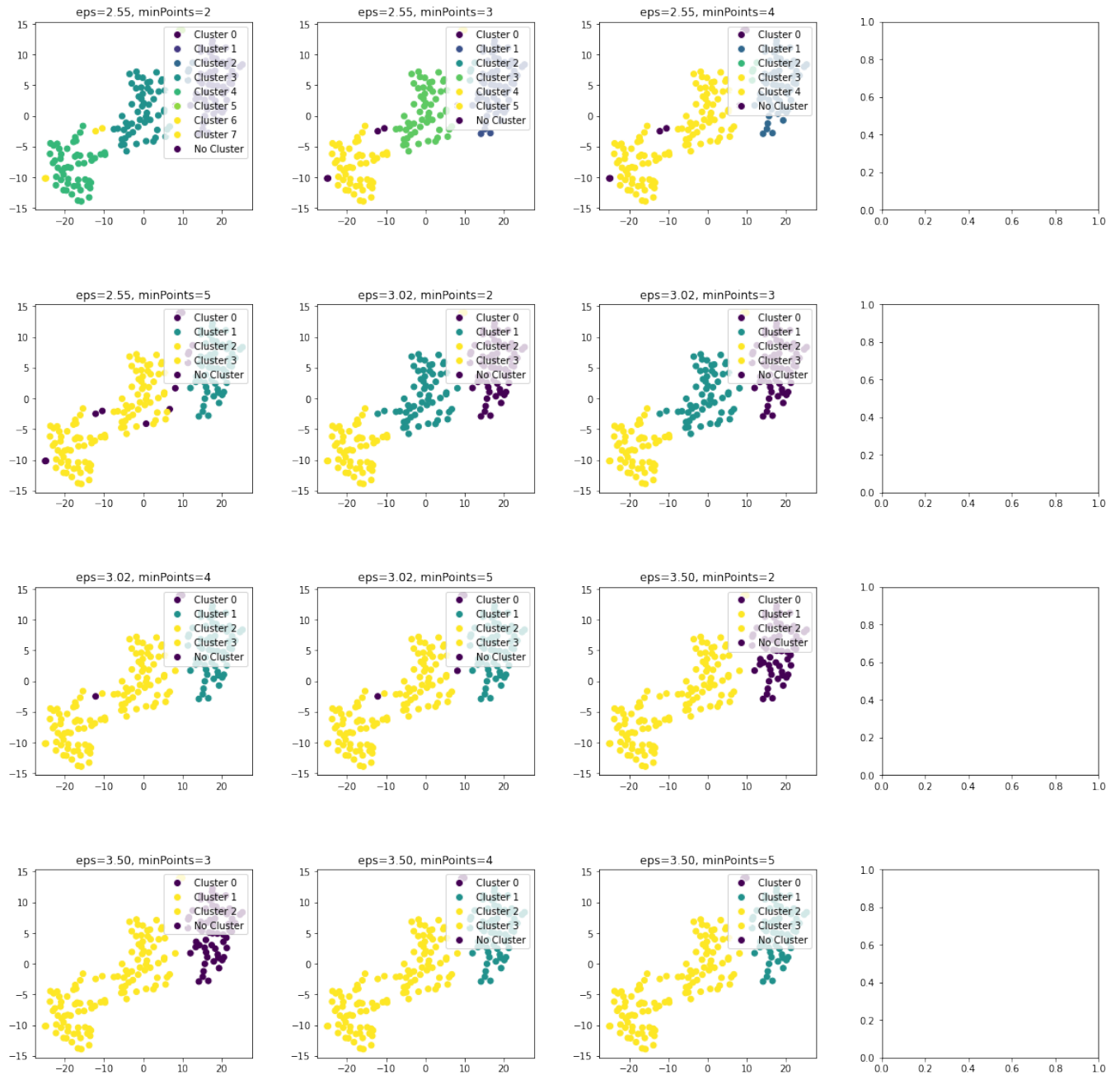
## ▼ Using Tsne:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import DBSCAN
4 from sklearn.manifold import TSNE
```

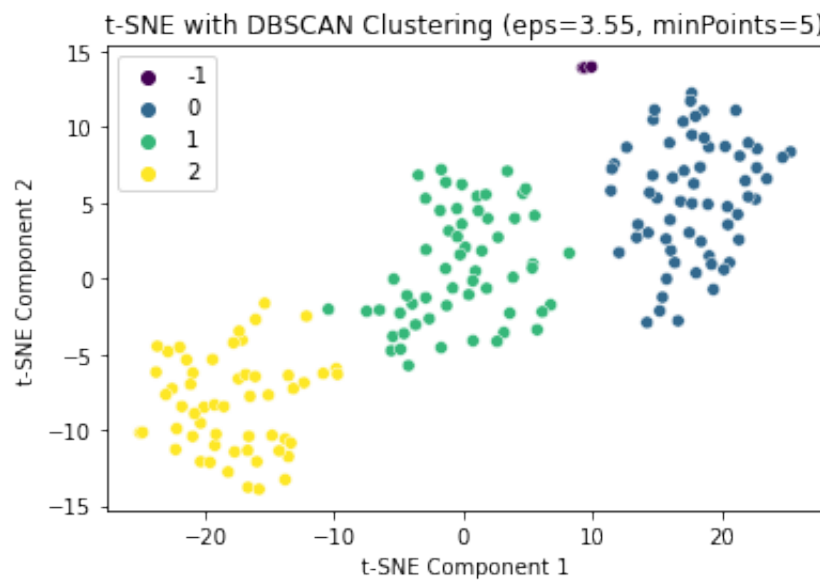
```
5 from sklearn.preprocessing import StandardScaler
6 import random
7
8 seed = random.randint(1, 100000)
9 print(seed)
10
11 # Define a range of epsilon and minPoints values to test
12 epsilons = np.linspace(2.55, 3.5, 3)
13 min_points_list = np.arange(2, 6)
14
15 # Function to apply DBSCAN and plot the results for a given epsilon and minPo
16 def apply_dbscan_subplot(data, epsilon, min_points, ax):
17     dbscan = DBSCAN(eps=epsilon, min_samples=min_points)
18     dbscan_labels = dbscan.fit_predict(data)
19     cmap = plt.get_cmap("viridis", len(set(dbscan_labels)) - 1)
20     scatter = ax.scatter(data[:, 0], data[:, 1], c=dbscan_labels, cmap=cmap)
21     ax.set_title(f"eps={epsilon:.2f}, minPoints={min_points}")
22
23     # Create legend
24     legend_elements = [plt.Line2D([0], [0], marker='o', color='w', label='Clu
25                             markerfacecolor=cmap(i), markersize=8) for
26     legend_elements.append(plt.Line2D([0], [0], marker='o', color='w', label=
27                             markerfacecolor=cmap(-1), markersize=8)
28     ax.legend(handles=legend_elements, loc='upper right')
29
30 # Create a 3x3 grid of subplots
31 fig, axes = plt.subplots(4, 4, figsize=(20, 20))
32 fig.suptitle("DBSCAN Clustering with t-SNE Results")
33
34 # Test each permutation of epsilons and minPoints and plot the results
35 for i, (epsilon, min_points) in enumerate(product(epsilons, min_points_list))
36     row, col = divmod(i, 3)
37     apply_dbscan_subplot(tsne_results_20, epsilon, min_points, axes[row, col]
38
39 # Adjust subplot spacing and show the figure
40 plt.subplots_adjust(hspace=0.5, wspace=0.3)
41 plt.show()
42
43
```

23320

## DBSCAN Clustering with t-SNE Results



```
1 apply_dbscan(tsne_results_20, "t-SNE", epsilon=3.55, min_points=5)  
2 plt.show()
```



## Process and Findings:

To answer the question, we started by implementing a function to find the optimal combination of epsilon and minPoints for DBSCAN to get 3 clusters using a range of values for both parameters. We used the `itertools.product` function to generate all possible combinations of epsilon and minPoints. Then, we applied DBSCAN on the 2D data from PCA, t-SNE, and MDS dimensionality reduction methods to find the best parameters for each method. After that, we plotted the k-distance graph for each method to assist in choosing the optimal epsilon and minPoints values. Finally, we applied DBSCAN with the chosen values and visualized the resulting clusters for each dimensionality reduction method. We used a wide range of epsilon and minPoints values to ensure that we would find the best combination for each method. We also plotted k-distance graphs to understand the relationship between the data points and their neighbors, providing insight into the choice of epsilon and minPoints values. In the final plots, we chose different color schemes to better visualize the clusters. We found that the optimal epsilon and minPoints values for PCA were 0.6 and 6, for t-SNE they were 2.2 and 5, and for MDS they were 1 and 2. The resulting DBSCAN clusters for each method were visualized in 2D plots with different colors representing different clusters. For t-SNE, we further experimented with a range of epsilon values between 2.55 and 3.5 and minPoints values from 2 to 5. A subplot with a 3x3 grid was created, showing different combinations of epsilon and minPoints for t-SNE. The findings suggest that t-SNE with an epsilon of 3.55 and minPoints of 5 is the most appropriate configuration for representing wines in a 2D space with DBSCAN clustering. The visualizations provide insight into the structure of the wine data and allow for an understanding of how different wines cluster together in a reduced-dimensional space. This analysis can be helpful for understanding similarities between different wines and identifying distinct groups within the data.

### ▼ Extra credit:

- a) Given your answers to all of these questions taken together, how many different kinds of wine do you think there are and how do they differ?
- b) Is there anything of interest you learned about wines from exploring this dataset with unsupervised machine learning method that is worth noting and not already covered in the questions above?



▼ A

Based on the analysis and visualizations obtained from the t-SNE and MDS methods, it appears that there are three distinct kinds of wine. The t-SNE plot shows three clearly separated groups, while the MDS plot outlines three groups as well. These wines likely differ based on factors such as color intensity, magnesium and other chemical compositions, and alcohol content, as well as various flavor attributes.

Examining the distributions of the data and the characteristics of the clusters, we can infer that the differences in color intensity may be related to the various types of wine, such as red, white, or rosé. The chemical compositions, including magnesium and other elements, contribute to the unique characteristics of each wine group. Additionally, variations in alcohol content and flavor profiles further distinguish these three types of wine.

Double-click (or enter) to edit

[Colab paid products](#) - [Cancel contracts here](#)

