

```
1 pip install palmerpenguins
```

```
Collecting palmerpenguins
  Downloading palmerpenguins-0.1.4-py3-none-any.whl (17 kB)
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (from palmerpenguins)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from palmerpenguins)
Requirement already satisfied: python-dateutil>=2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil)
Installing collected packages: palmerpenguins
Successfully installed palmerpenguins-0.1.4
```

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from palmerpenguins import load_penguins # For penguins dataset
5 import pandas as pd # For dataframes
6 import matplotlib.pyplot as plt # For plotting functions
7 import seaborn as sns # For additional plotting functions
8 from sklearn.cluster import KMeans # For k-Means
9 from sklearn.model_selection import GridSearchCV # For grid search
10 from sklearn.metrics import silhouette_score # For metrics and scores
11 from sklearn.preprocessing import StandardScaler # For standardizing data
```

```
1 # 1. Import the Vehicle dataset, summarize it and explain the output
2 data = pd.read_csv("Vehicle.csv");
3 print(data.describe())
```

	Year	Selling_Price	Present_Price	Kms_Driven	Owner
count	301.000000	301.000000	301.000000	301.000000	301.000000
mean	2013.627907	4.661296	7.628472	36947.205980	0.043189
std	2.891554	5.082812	8.644115	38886.883882	0.247915
min	2003.000000	0.100000	0.320000	500.000000	0.000000
25%	2012.000000	0.900000	1.200000	15000.000000	0.000000
50%	2014.000000	3.600000	6.400000	32000.000000	0.000000
75%	2016.000000	6.000000	9.900000	48767.000000	0.000000
max	2018.000000	35.000000	92.600000	500000.000000	3.000000

```
1 # 2. Show the structure and dimension of the dataset and explain it
2 print(data.info()) # tells how many null values, datatypes, column names
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 301 entries, 0 to 300
Data columns (total 9 columns):
#   Column          Non-Null Count  Dtype
---  ---
0   Car_Name        301 non-null    object
1   Year            301 non-null    int64
2   Selling_Price   301 non-null    float64
3   Present_Price   301 non-null    float64
4   Kms_Driven      301 non-null    int64
5   Fuel_Type       301 non-null    object
6   Seller_Type     301 non-null    object
7   Transmission    301 non-null    object
8   Owner           301 non-null    int64
dtypes: float64(2), int64(3), object(4)
memory usage: 21.3+ KB
None
```

```
1 # 3. Show the column names of the Vehicle dataset and the first 3 rows and the last 6 rows of i
2 print("\n3. Column names:")
3 print(data.columns)
4 print("\nFirst 3 rows:")
```

```

4 print("\nFirst 3 rows: ")
5 print(data.head(3))
6 print("\nLast 6 rows:")
7 print(data.tail(6))

```

3. Column names:

```

Index(['Car_Name', 'Year', 'Selling_Price', 'Present_Price', 'Kms_Driven',
      'Fuel_Type', 'Seller_Type', 'Transmission', 'Owner'],
      dtype='object')

```

First 3 rows:

	Car_Name	Year	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	\
0	ritz	2014	3.35	5.59	27000	Petrol	
1	sx4	2013	4.75	9.54	43000	Diesel	
2	ciaz	2017	7.25	9.85	6900	Petrol	

	Seller_Type	Transmission	Owner
0	Dealer	Manual	0
1	Dealer	Manual	0
2	Dealer	Manual	0

Last 6 rows:

	Car_Name	Year	Selling_Price	Present_Price	Kms_Driven	Fuel_Type	\
295	city	2015	8.55	13.09	60076	Diesel	
296	city	2016	9.50	11.60	33988	Diesel	
297	brio	2015	4.00	5.90	60000	Petrol	
298	city	2009	3.35	11.00	87934	Petrol	
299	city	2017	11.50	12.50	9000	Diesel	
300	brio	2016	5.30	5.90	5464	Petrol	

	Seller_Type	Transmission	Owner
295	Dealer	Manual	0
296	Dealer	Manual	0
297	Dealer	Manual	0
298	Dealer	Manual	0
299	Dealer	Manual	0
300	Dealer	Manual	0

1 # 4. Show the average Kms_Driven for each type of car (Car_Name)

```

2 average = data.groupby("Car_Name")["Kms_Driven"].mean() #first grouping the Cars and how many kms the
3 print("\n4. Average Kms_Driven for each type of car:")
4 print(average)

```

4. Average Kms_Driven for each type of car:

```

Car_Name
800      127000.000000
Activa 3g  250250.000000
Activa 4g   1300.000000
Bajaj ct 100  35000.000000
Bajaj Avenger 150  7000.000000
...
sx4      50740.000000
verna    42747.285714
vitara brezza  2071.000000
wagon r    40644.750000
xcnt      27448.333333
Name: Kms_Driven, Length: 98, dtype: float64

```

```

1 average_selling_price_by_year = data.groupby("Year")["Selling_Price"].mean()
2 print("\n5. Average Selling_Price of the cars in each year:")
3 print(average_selling_price_by_year)

```

5. Average Selling_Price of the cars in each year:

```

Year
2003    1.300000

```

```

2004    1.500000
2005    2.487500
2006    1.437500
2007    0.160000
2008    1.002857
2009    2.816667
2010    5.262667
2011    2.375263
2012    3.841304
2013    3.540909
2014    4.762105
2015    5.927049
2016    5.213200
2017    6.209143
2018    9.250000

```

Name: Selling_Price, dtype: float64

```

1 n]].drop_duplicates() #drop duplicates removes all the duplcates and keeps only the unique values
2 ssion:")
3

```

6. Unique combinations of Car_Name, Fuel_Type, Seller_Type, and Transmission:

	Car_Name	Fuel_Type	Seller_Type	Transmission
0	ritz	Petrol	Dealer	Manual
1	sx4	Diesel	Dealer	Manual
2	ciaz	Petrol	Dealer	Manual
3	wagon r	Petrol	Dealer	Manual
4	swift	Diesel	Dealer	Manual
..
259	amaze	Petrol	Dealer	Manual
263	jazz	Petrol	Dealer	Manual
275	city	Petrol	Dealer	Automatic
285	jazz	Petrol	Dealer	Automatic
287	amaze	Petrol	Dealer	Automatic

[135 rows x 4 columns]

```

1 # Count the occurrences of unique combinations
2 combination_counts = data.groupby(['Car_Name', 'Fuel_Type', 'Seller_Type', 'Transmission']).size().reset_index(name='Count')
3
4 # Display combinations in ascending order of count
5 ascending_order = combination_counts.sort_values(by='Count', ascending=True)
6 print("Combinations in Ascending Order:")
7 print(ascending_order)
8
9 # Display combinations in descending order of count
10 descending_order = combination_counts.sort_values(by='Count', ascending=False)
11 print("\nCombinations in Descending Order:")
12 print(descending_order)

```

Combinations in Ascending Order:

	Car_Name	Fuel_Type	Seller_Type	Transmission	Count
0	800	Petrol	Individual	Manual	1
86	elantra	Petrol	Dealer	Automatic	1
85	elantra	Diesel	Dealer	Manual	1
82	creta	Petrol	Dealer	Manual	1
77	corolla	Petrol	Dealer	Automatic	1
..
46	Royal Enfield Classic 350	Petrol	Individual	Manual	7
97	fortuner	Diesel	Dealer	Automatic	8
68	brio	Petrol	Dealer	Manual	9
80	corolla altis	Petrol	Dealer	Manual	11
76	city	Petrol	Dealer	Manual	19

[135 rows x 5 columns]

Combinations in Descending Order:

	Car_Name	Fuel_Type	Seller_Type	Transmission	Count
76	city	Petrol	Dealer	Manual	19
80	corolla altis	Petrol	Dealer	Manual	11
68	brio	Petrol	Dealer	Manual	9
97	fortuner	Diesel	Dealer	Automatic	8
130	verna	Petrol	Dealer	Manual	7
..
45	Royal Enfield Bullet 350	Petrol	Individual	Manual	1
44	Mahindra Mojo XT300	Petrol	Individual	Manual	1
43	KTM RC390	Petrol	Individual	Manual	1
41	KTM 390 Duke	Petrol	Individual	Manual	1
67	brio	Petrol	Dealer	Automatic	1

[135 rows x 5 columns]

```
1 # 9. Check for missing values
2 missing_values = data.isnull().sum().sort_values(ascending=False)
3 print("9. Missing values in the dataset:")
4 print(missing_values)
```

9. Missing values in the dataset:

```
Car_Name      0
Year          0
Selling_Price 0
Present_Price 0
Kms_Driven    0
Fuel_Type     0
Seller_Type   0
Transmission  0
Owner         0
dtype: int64
```

```
1 # 10. Replace missing values with the most repeated value in each column
2
3 # Iterate through each column with missing values
4 for column in missing_values.index:
5     # Find the most common (mode) value in the column
6     most_common_value = data[column].mode()[0]
7
8     # Fill missing values in the column with the most common value
9     data[column].fillna(most_common_value, inplace=True)
10
11 # Check if missing values were replaced successfully
12 missing_values_after_replace = data.isnull().sum().sum()
13 print("\n10. Missing values after replacement:", missing_values_after_replace)
```

10. Missing values after replacement: 0

```
1 data.drop_duplicates(inplace=True)
2 print("\n11. Dataset after removing duplicates:")
```

11. Dataset after removing duplicates:

```

1 # 12. Replace values in attributes (Fuel_Type, Seller_Type, Transmission)
2 data['Fuel_Type'].replace({"Petrol": 0, "Diesel": 1, "CNG": 2}, inplace=True)
3 data['Seller_Type'].replace({"Dealer": 0, "Individual": 1}, inplace=True)
4 data['Transmission'].replace({"Manual": 0, "Automatic": 1}, inplace=True)
5
6 # Show the conversion output
7 print("\n12. Dataset after attribute value conversion:")
8 print(data[['Fuel_Type', 'Seller_Type', 'Transmission']])

```

12. Dataset after attribute value conversion:

	Fuel_Type	Seller_Type	Transmission
0	0	0	0
1	1	0	0
2	0	0	0
3	0	0	0
4	1	0	0
..
296	1	0	0
297	0	0	0
298	0	0	0
299	1	0	0
300	0	0	0

[299 rows x 3 columns]

```

1 # 13. Add a new 'Age' field based on the 'Year' column
2 current_year = 2023 # Assuming the current year is 2023
3 data['Age'] = current_year - data['Year']
4
5 # Show the output with the 'Age' field
6 print("\n13. Dataset with the 'Age' field:")
7 print(data[['Year', 'Age']])

```

13. Dataset with the 'Age' field:

	Year	Age
0	2014	9
1	2013	10
2	2017	6
3	2011	12
4	2014	9
..
296	2016	7
297	2015	8
298	2009	14
299	2017	6
300	2016	7

[299 rows x 2 columns]

```

1 new_dataset = data[['Car_Name', 'Selling_Price', 'Present_Price', 'Kms_Driven']]
2
3 # Show the output of the new dataset
4 print("\n14. New dataset:")
5 print(new_dataset)
6
7 # 15. Shuffle the rows of the dataset randomly
8 shuffled_data = data.sample(frac=1, random_state=42) # random_state for reproducibility
9
10 # Show the output of the shuffled dataset
11 print("\n15. Shuffled dataset:")
12 print(shuffled_data)

```

14. New dataset:

	Car_Name	Selling_Price	Present_Price	Kms_Driven
0	ritz	3.35	5.59	27000
1	sx4	4.75	9.54	43000
2	ciaz	7.25	9.85	6900
3	wagon r	2.85	4.15	5200
4	swift	4.60	6.87	42450
..
296	city	9.50	11.60	33988
297	brrio	4.00	5.90	60000
298	city	3.35	11.00	87934
299	city	11.50	12.50	9000
300	brrio	5.30	5.90	5464

[299 rows x 4 columns]

15. Shuffled dataset:

	Car_Name	Year	Selling_Price	Present_Price	\
283	city	2016	8.99	11.80	
267	city	2016	8.35	9.40	
166	Hero Passion Pro	2016	0.45	0.55	
9	ciaz	2015	7.45	8.92	
78	corolla altis	2010	5.25	22.83	
..	
190	Bajaj Pulsar 150	2008	0.20	0.75	
72	corolla altis	2013	7.45	18.61	
108	Royal Enfield Thunder 350	2016	1.20	1.50	
272	city	2015	7.50	10.00	
104	Royal Enfield Classic 350	2017	1.35	1.47	

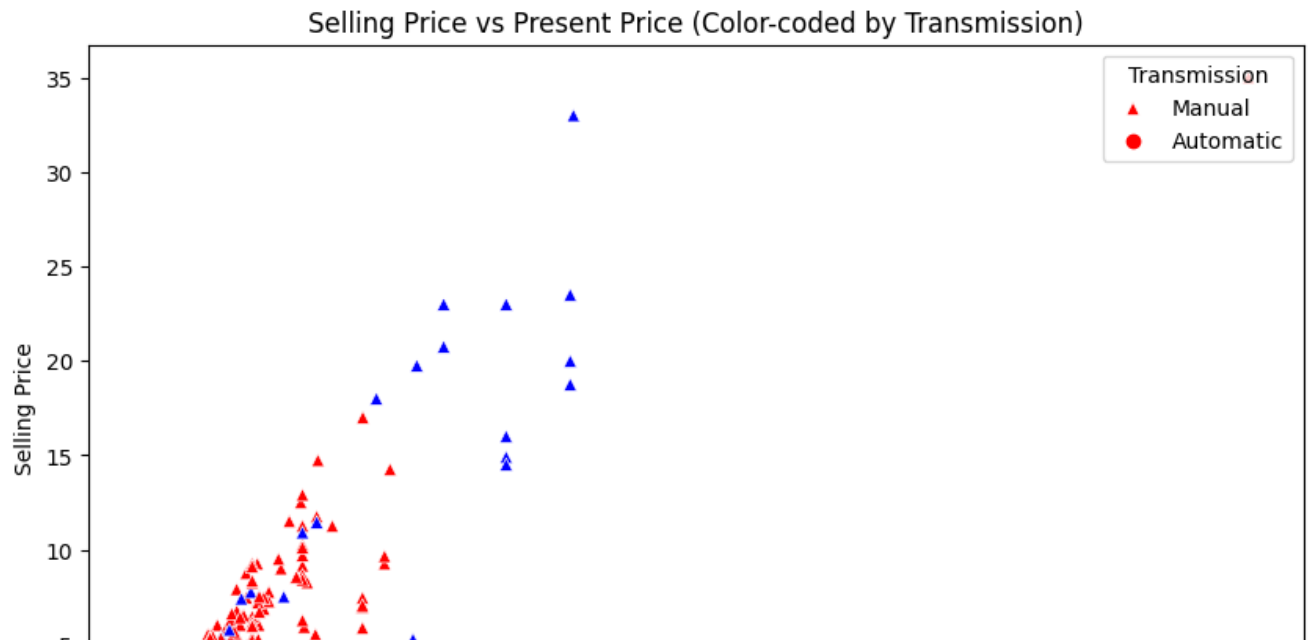
	Kms_Driven	Fuel_Type	Seller_Type	Transmission	Owner	Age
283	9010	0	0	0	0	7
267	19434	1	0	0	0	7
166	1000	0	1	0	0	7
9	42367	1	0	0	0	8
78	80000	0	0	1	0	13
..
190	60000	0	1	0	0	15
72	56001	0	0	0	0	10
108	18000	0	1	0	0	7
272	27600	0	0	0	0	8
104	4100	0	1	0	0	6

[299 rows x 10 columns]

```

1 # a. Create the scatter plot with color-coding
2 plt.figure(figsize=(10, 6))
3 sns.scatterplot(x='Present_Price', y='Selling_Price', data=data, hue='Transmission', palette={0: 'red',
4
5 # Add labels, title, and legend
6 plt.xlabel('Present Price')
7 plt.ylabel('Selling Price')
8 plt.title('Selling Price vs Present Price (Color-coded by Transmission)')
9 plt.legend(title='Transmission', loc='upper right', labels=['Manual', 'Automatic'])
10
11 plt.show()
12
13

```



▼ c. What do you understand from the output:

This scatter plot visualizes the relationship between Selling Price and Present Price of vehicles. The points are color-coded based on the transmission type, with red representing manual transmission (0) and blue representing automatic transmission (1). We can observe that vehicles with manual transmission tend to have a wider range of selling prices, while those with automatic transmission are more clustered in the lower price range.

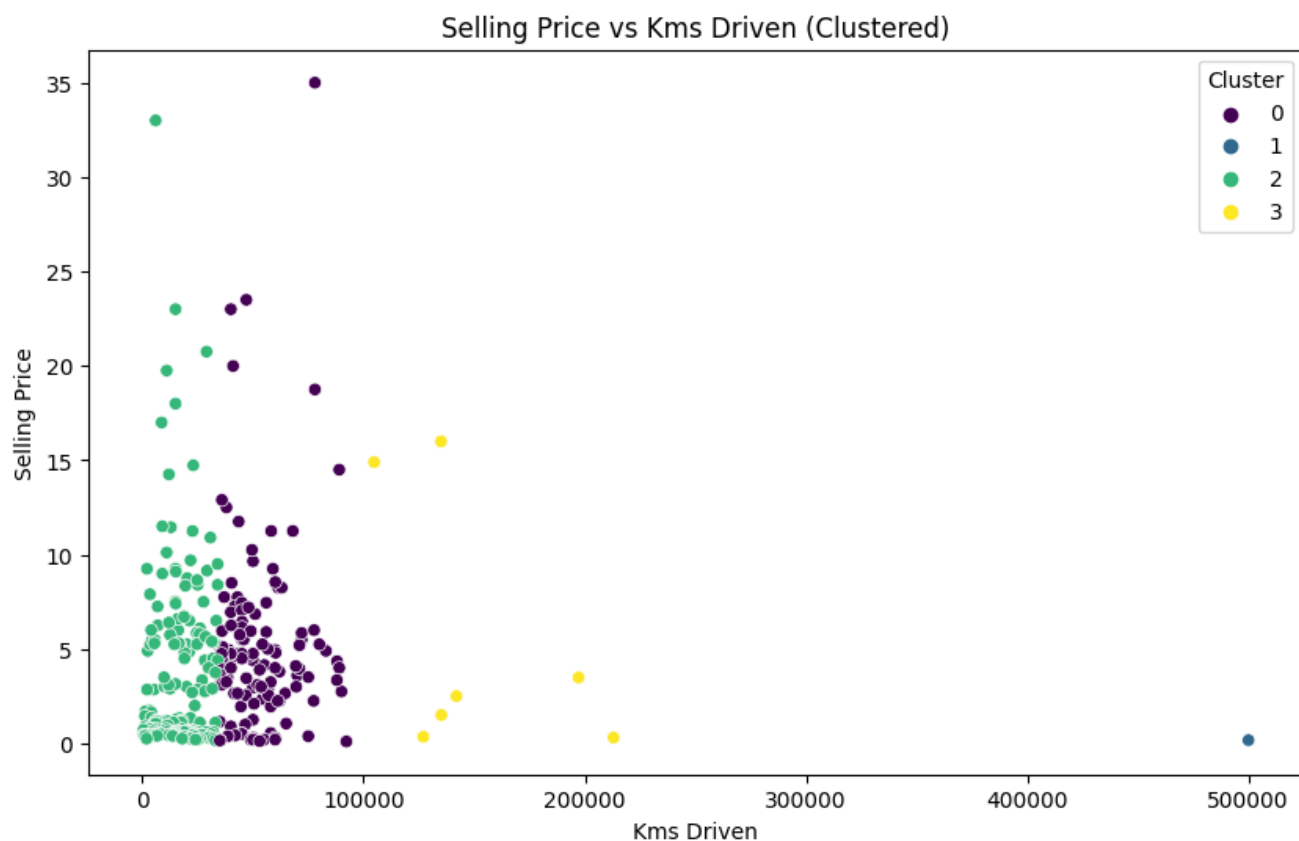
```
1 # 17. Box plot of Selling_Price Vs Transmission and Fuel_Type
2 plt.figure(figsize=(12, 6))
3 sns.boxplot(x='Transmission', y='Selling_Price', hue='Fuel_Type', data=data)
4 plt.xlabel('Transmission')
5 plt.ylabel('Selling Price')
6 plt.title('Box Plot of Selling Price vs Transmission and Fuel Type')
7 plt.legend(title='Fuel Type', loc='upper right')
8 plt.show()
```

```

1 # 18. Scatter plot of Selling_Price Vs Kms_Driven with k-means clustering (4 clusters)
2 from sklearn.cluster import KMeans
3
4 # Select the features for clustering
5 features = data[['Selling_Price', 'Kms_Driven']]
6
7 # Apply k-means clustering
8 kmeans = KMeans(n_clusters=4, random_state=0)
9 data['Cluster'] = kmeans.fit_predict(features)
10
11 # Create a scatter plot with color-coded clusters
12 plt.figure(figsize=(10, 6))
13 sns.scatterplot(x='Kms_Driven', y='Selling_Price', data=data, hue='Cluster', palette='viridis')
14 plt.xlabel('Kms Driven')
15 plt.ylabel('Selling Price')
16 plt.title('Selling Price vs Kms Driven (Clustered)')
17 plt.legend(title='Cluster')
18 plt.show()

```

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v: warnings.warn(



```

1 # 19. Scatter plot of Selling_Price Vs Present_Price with hierarchical clustering (3 clusters)
2 from sklearn.cluster import AgglomerativeClustering
3
4 # Select the features for clustering
5 features = data[['Selling_Price', 'Present_Price']]
6
7 # Apply hierarchical clustering
8 hierarchical_clustering = AgglomerativeClustering(n_clusters=3)
9 data['Hierarchical_Cluster'] = hierarchical_clustering.fit_predict(features)
10
11 # Create a scatter plot with color-coded clusters

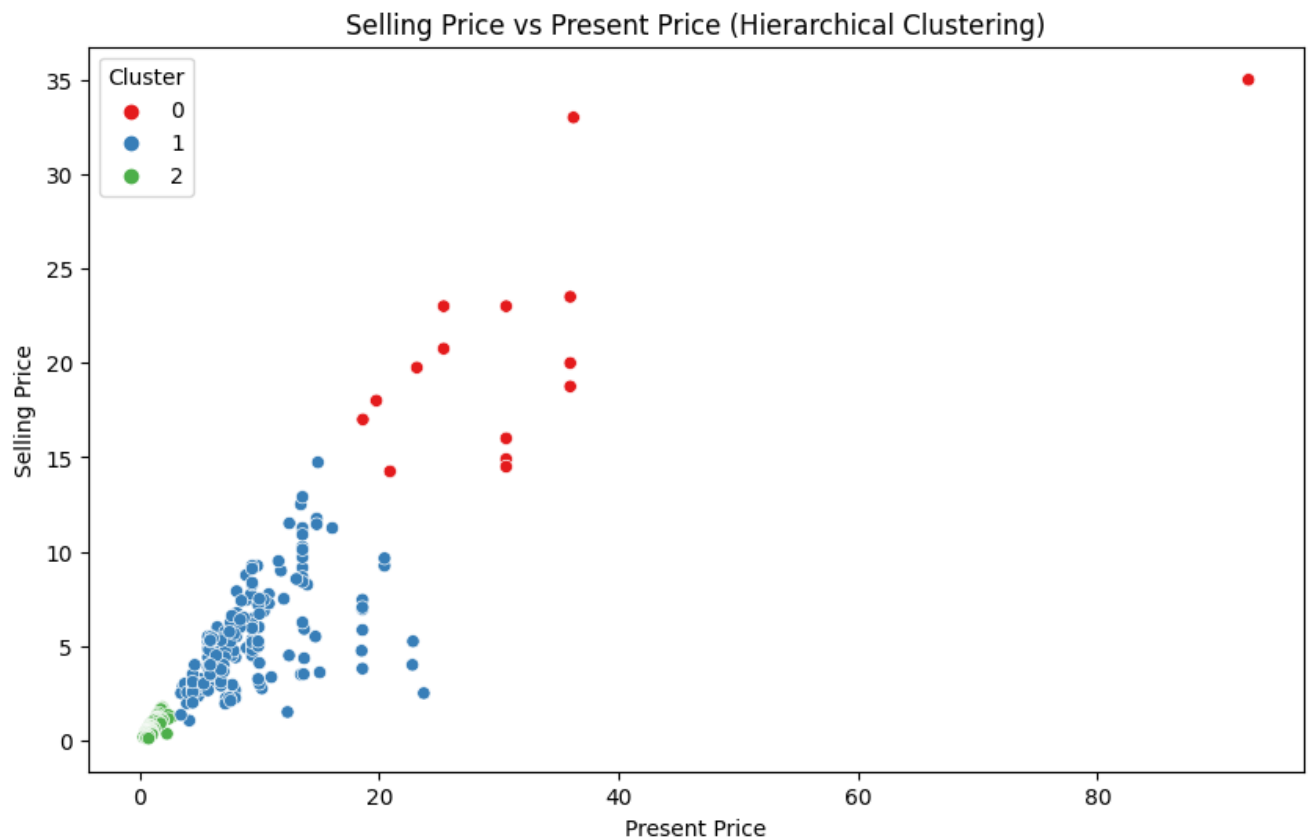
```



```

12 plt.figure(figsize=(10, 6))
13 sns.scatterplot(x='Present_Price', y='Selling_Price', data=data, hue='Hierarchical_Cluster', palette='S
14 plt.xlabel('Present Price')
15 plt.ylabel('Selling Price')
16 plt.title('Selling Price vs Present Price (Hierarchical Clustering)')
17 plt.legend(title='Cluster')
18 plt.show()

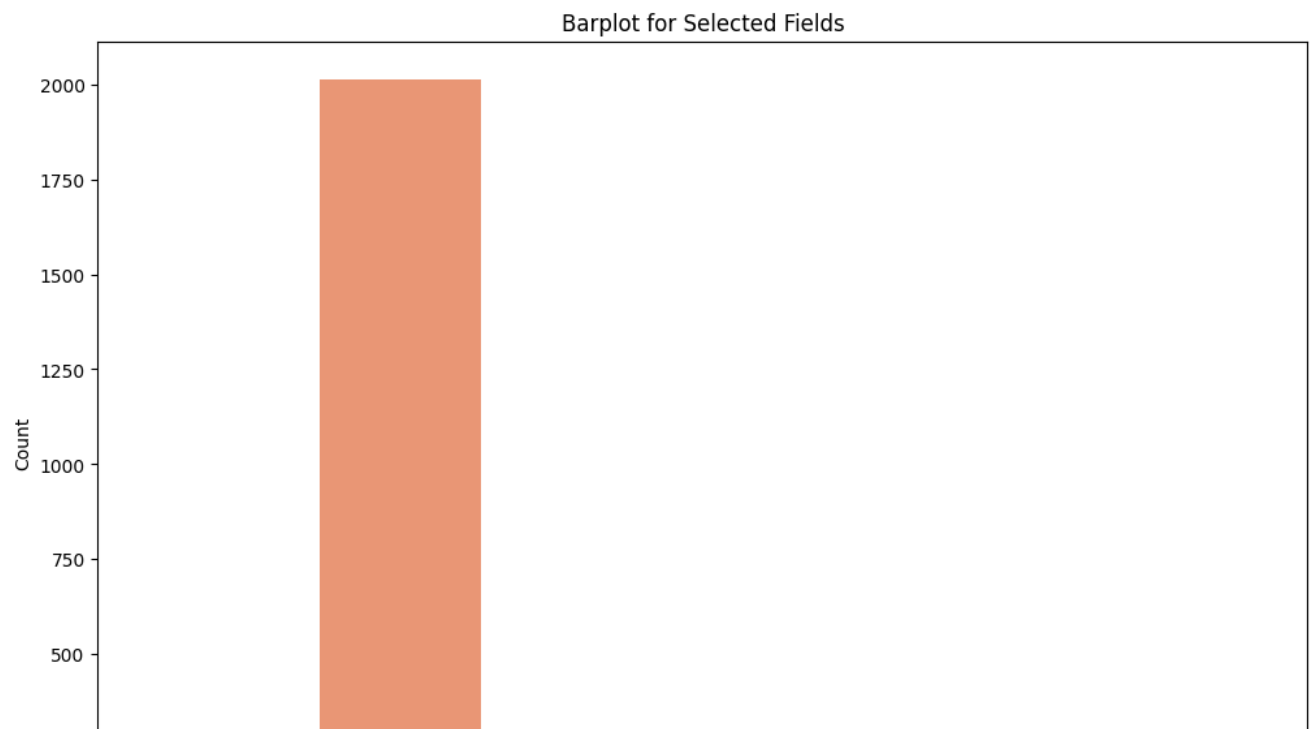
```



```

1 # 20. Create a barplot for selected fields with labels, titles, and colors
2 selected_fields = ['Age', 'Year', 'Transmission', 'Seller_Type', 'Fuel_Type', 'Owner']
3
4 # Create a new DataFrame with selected fields
5 selected_data = data[selected_fields]
6
7 # Plot the barplot
8 plt.figure(figsize=(12, 8))
9 sns.barplot(data=selected_data, palette='Set2')
10 plt.xlabel('Attributes')
11 plt.ylabel('Count')
12 plt.title('Barplot for Selected Fields')
13 plt.xticks(rotation=45)
14 plt.show()

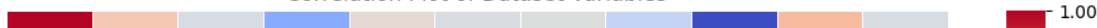
```



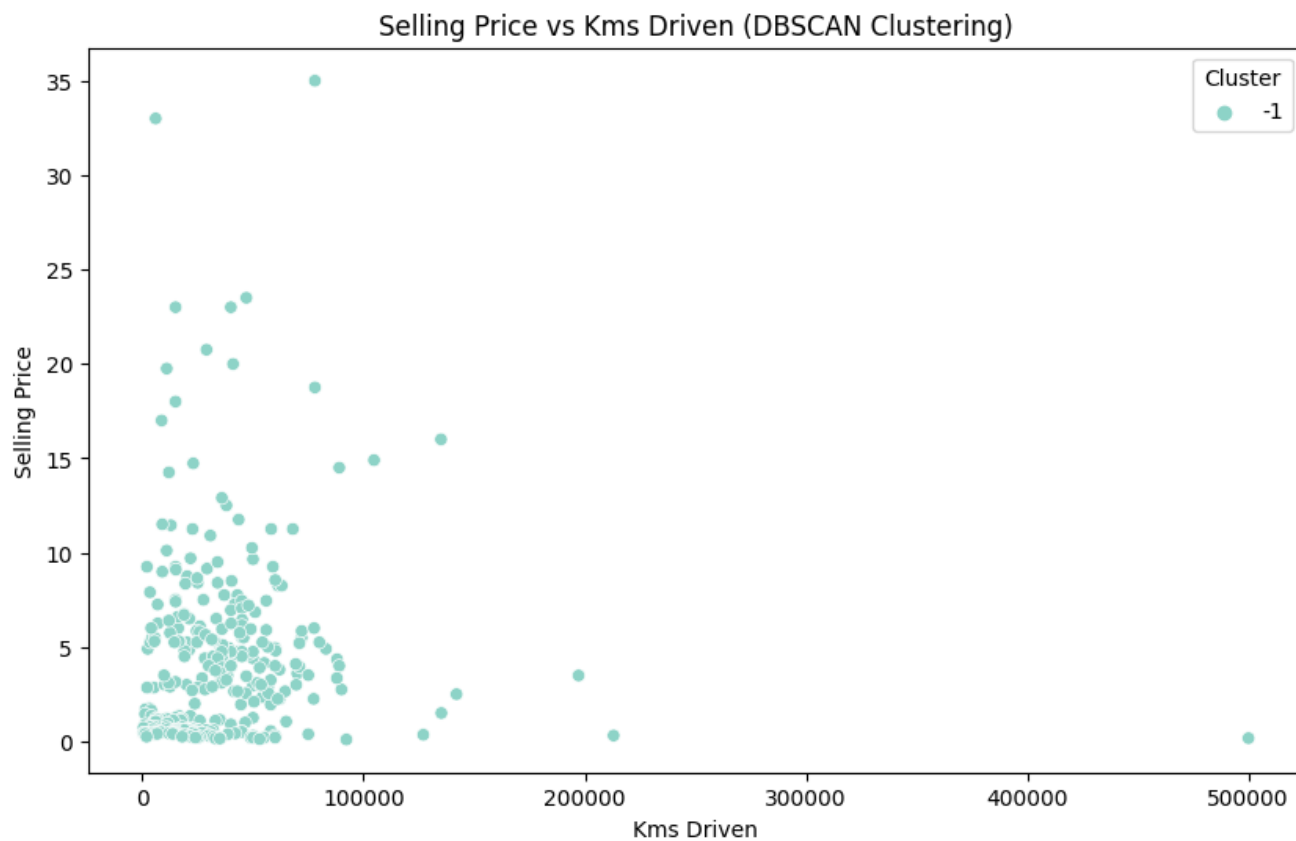
```
1 # 21. Correlation plot of the whole dataset
2 corr_matrix = data.corr()
3 plt.figure(figsize=(12, 8))
4 sns.heatmap(corr_matrix, annot=True, cmap='coolwarm', linewidths=0.5)
5 plt.title('Correlation Plot of Dataset Variables')
6 plt.show()
```

```
<ipython-input-21-a1fb9262dd1e>:2: FutureWarning: The default value of numeric_only in DataFrame.corr
corr_matrix = data.corr()
```

Correlation Plot of Dataset Variables



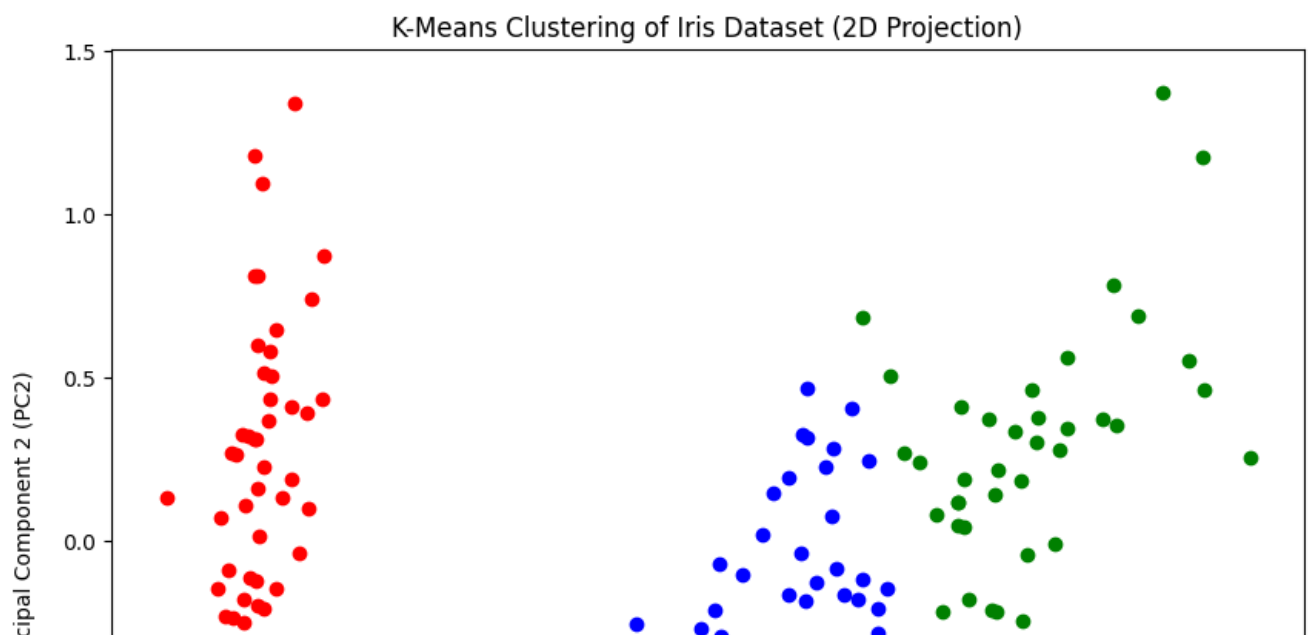
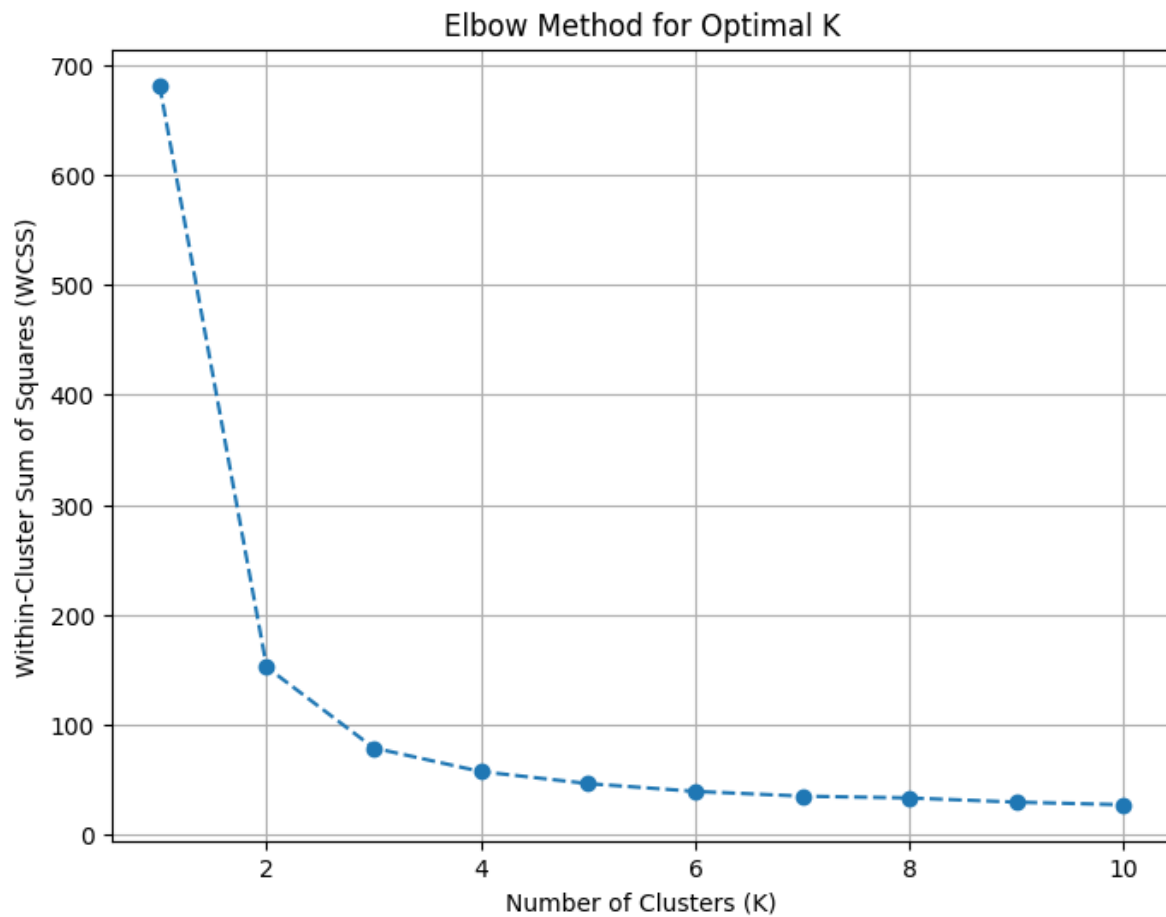
```
1 # 22. Scatter plot of Selling_Price Vs Kms_Driven with DBSCAN clustering (3 clusters)
2 from sklearn.cluster import DBSCAN
3
4 # Select the features for clustering
5 features = data[['Selling_Price', 'Kms_Driven']]
6
7 # Apply DBSCAN clustering
8 dbscan = DBSCAN(eps=0.5, min_samples=5)
9 data['DBSCAN_Cluster'] = dbscan.fit_predict(features)
10
11 # Create a scatter plot with color-coded clusters and add a legend
12 plt.figure(figsize=(10, 6))
13 sns.scatterplot(x='Kms_Driven', y='Selling_Price', data=data, hue='DBSCAN_Cluster', palette='Set3')
14 plt.xlabel('Kms Driven')
15 plt.ylabel('Selling Price')
16 plt.title('Selling Price vs Kms Driven (DBSCAN Clustering)')
17 plt.legend(title='Cluster')
18 plt.show()
```



Part B

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from sklearn.cluster import KMeans
4 from sklearn.datasets import load_iris
5 from sklearn.decomposition import PCA
6 import numpy as np
7
```

```
8 # Load the Iris dataset
9 iris = load_iris()
10 X = iris.data # Features (sepal length, sepal width, petal length, petal width)
11
12 # Find the optimal number of clusters (K) using the elbow method
13 wcss = []
14 for i in range(1, 11):
15     kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init='auto', random_state=0)
16     kmeans.fit(X)
17     wcss.append(kmeans.inertia_)
18
19 # Plot the elbow graph to find the optimal K value
20 plt.figure(figsize=(8, 6))
21 plt.plot(range(1, 11), wcss, marker='o', linestyle='--')
22 plt.title('Elbow Method for Optimal K')
23 plt.xlabel('Number of Clusters (K)')
24 plt.ylabel('Within-Cluster Sum of Squares (WCSS)')
25 plt.grid()
26 plt.show()
27
28 # Based on the elbow method, choose the optimal K value (e.g., K=3)
29 optimal_k = 3
30
31 # Perform K-Means clustering with the optimal K value
32 kmeans = KMeans(n_clusters=optimal_k, init='k-means++', max_iter=300, n_init='auto', random_state=0)
33 kmeans.fit(X)
34
35 # Reduce dimensionality to 2D for visualization
36 pca = PCA(n_components=2)
37 X_2d = pca.fit_transform(X)
38
39 # Add cluster labels to the data
40 iris_df = pd.DataFrame(X_2d, columns=['PC1', 'PC2'])
41 iris_df['Cluster'] = kmeans.labels_
42
43 # Plot the clusters in 2D
44 plt.figure(figsize=(10, 8))
45 colors = ['r', 'g', 'b']
46 for cluster in range(optimal_k):
47     cluster_data = iris_df[iris_df['Cluster'] == cluster]
48     plt.scatter(cluster_data['PC1'], cluster_data['PC2'], c=colors[cluster], label=f'Cluster {cluster}')
49
50 plt.title('K-Means Clustering of Iris Dataset (2D Projection)')
51 plt.xlabel('Principal Component 1 (PC1)')
52 plt.ylabel('Principal Component 2 (PC2)')
53 plt.legend()
54 plt.show()
```



```

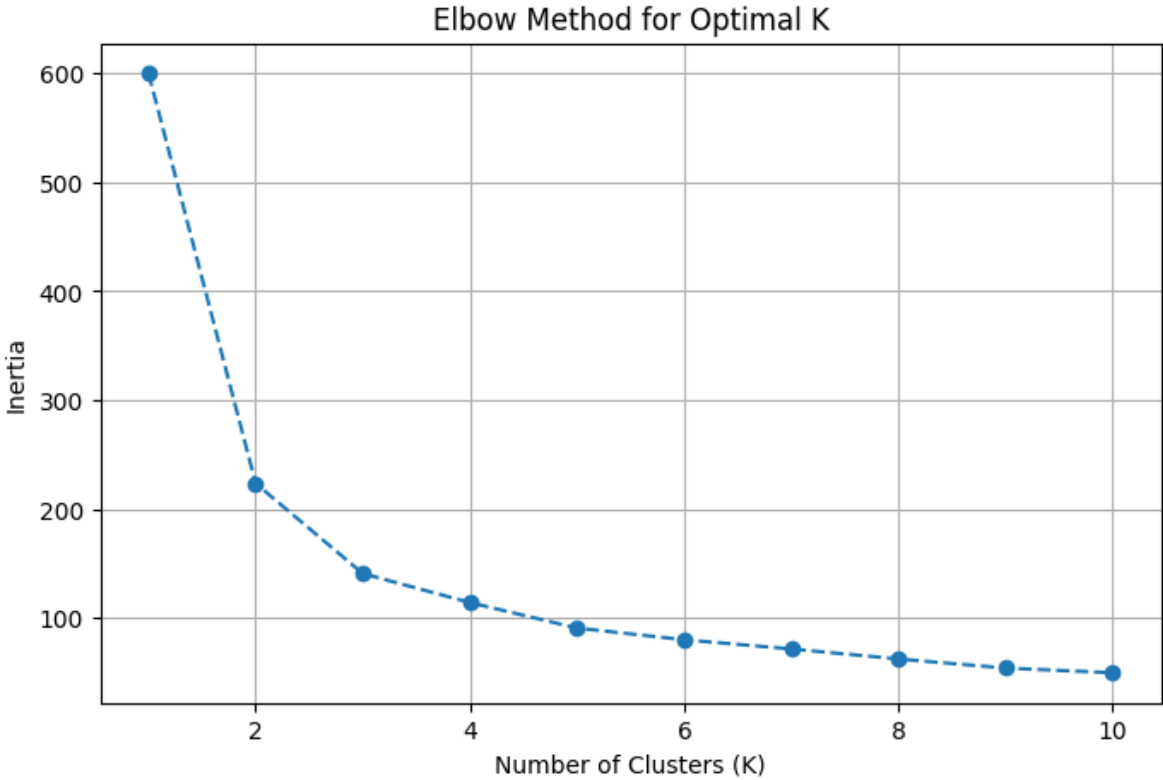
1 # Import necessary libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.cluster import KMeans
7 from sklearn.preprocessing import StandardScaler
8
9 # Load the Iris flower dataset
10 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
11 column_names = ["sepal_length", "sepal_width", "petal_length", "petal_width", "species"]
12 data = pd.read_csv(url, header=None, names=column_names)
13
14 # Explore the data
15 print(data.head())

```

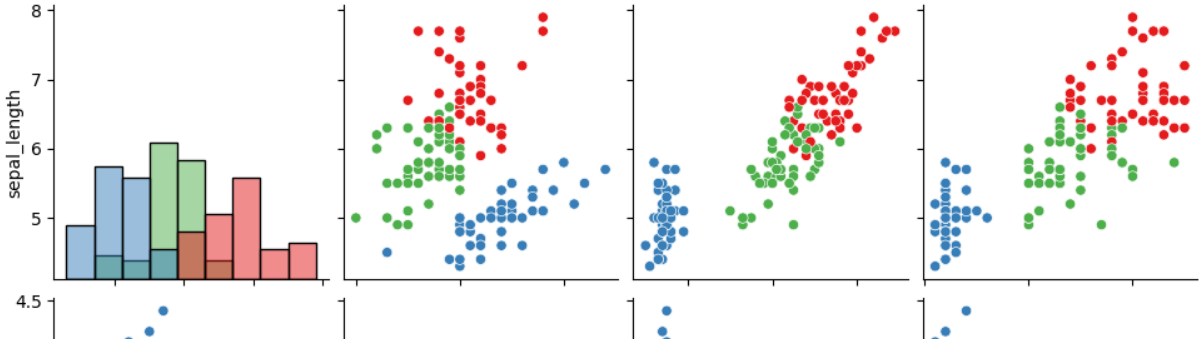
```
16
17 # Remove the 'species' column
18 X = data.drop('species', axis=1)
19
20 # Standardize the data
21 scaler = StandardScaler()
22 X_scaled = scaler.fit_transform(X)
23
24 # Choose the number of clusters (K) using the Elbow method
25 inertia = []
26 for k in range(1, 11):
27     kmeans = KMeans(n_clusters=k, random_state=42)
28     kmeans.fit(X_scaled)
29     inertia.append(kmeans.inertia_)
30
31 # Plot the Elbow method results
32 plt.figure(figsize=(8, 5))
33 plt.plot(range(1, 11), inertia, marker='o', linestyle='--')
34 plt.xlabel('Number of Clusters (K)')
35 plt.ylabel('Inertia')
36 plt.title('Elbow Method for Optimal K')
37 plt.grid(True)
38 plt.show()
39
40 # From the Elbow method, it seems K=3 is a good choice
41 kmeans = KMeans(n_clusters=3, random_state=42)
42 data['cluster'] = kmeans.fit_predict(X_scaled)
43
44 # Visualize the clustered data
45 sns.pairplot(data=data, hue='cluster', palette='Set1', diag_kind='hist')
46 plt.show()
47
```

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
```



```
/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: FutureWarning: The default v
warnings.warn(
```



```
1 # Import necessary libraries
2 import pandas as pd
```

```

3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.model_selection import train_test_split
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
9
10 # Load the breast cancer dataset
11 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/breast-cancer-
12 column_names = ["id", "clump_thickness", "uniformity_cell_size", "uniformity_cell_shape", "marginal_adh
13 data = pd.read_csv(url, header=None, names=column_names)
14
15 # Data preprocessing
16 # Replace missing values represented as '?' with NaN and convert 'bare_nuclei' to numeric
17 data['bare_nuclei'] = pd.to_numeric(data['bare_nuclei'], errors='coerce')
18 data = data.dropna() # Drop rows with missing values
19
20 # Map class labels to 0 (benign) and 1 (malignant)
21 data['class'] = data['class'].map({2: 0, 4: 1})
22
23 # Split the data into features (X) and target (y)
24 X = data.drop(['id', 'class'], axis=1)
25 y = data['class']
26
27 # Split the data into training and testing sets
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
29
30 # Train the KNN classifier
31 knn = KNeighborsClassifier(n_neighbors=5) # You can choose the number of neighbors (K) as needed
32 knn.fit(X_train, y_train)
33
34 # Predict on the test set
35 y_pred = knn.predict(X_test)
36
37 # Evaluate the model
38 accuracy = accuracy_score(y_test, y_pred)
39 confusion = confusion_matrix(y_test, y_pred)
40 classification_report_str = classification_report(y_test, y_pred)
41
42 # Print the evaluation results
43 print("Accuracy:", accuracy)
44 print("Confusion Matrix:\n", confusion)
45 print("Classification Report:\n", classification_report_str)
46

```

Accuracy: 0.9560975609756097

Confusion Matrix:

```
[[125  2]
 [ 7 71]]
```

Classification Report:

	precision	recall	f1-score	support
0	0.95	0.98	0.97	127
1	0.97	0.91	0.94	78
accuracy			0.96	205
macro avg	0.96	0.95	0.95	205
weighted avg	0.96	0.96	0.96	205