

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.preprocessing import MinMaxScaler
6 from sklearn.metrics import silhouette_score
7 from sklearn.cluster import KMeans
8 pd.set_option('display.max_columns', None)
9 pd.set_option('display.max_rows', None)
10 pd.set_option('display.width', None)
11
12
13 class DataSet:
14     def __init__(self, df):
15         self.df = df
16
17     def compute_statistical_summary(self):
18         numerical_columns = self.df.select_dtypes(
19             include=[np.number])
20         if not numerical_columns.empty:
21             numerical_stats = numerical_columns.agg(
22                 ['mean', 'median', 'std'])
23             print("Statistics for numerical columns")
24             print(numerical_stats)
25         else:
26             print("No numerical columns found.")
27
28         categorical_columns = self.df.select_dtypes(
29             include=['object', 'category'])
30         if not categorical_columns.empty:
31             print("\nFrequency distribution for")
32             print("categorical columns:")
33             for column in categorical_columns:
34                 freq_dist = ((self.df[column].
35                     value_counts() / len(df)) * 100).astype(str) + ' %'
36                 print(f"\nColumn: {column}")
37                 print(freq_dist)
38
39     def visualize_data(self):
40         sns.set_palette('bright')
```

```

36         font_title = {'fontweight': 'bold', '
    fontsize': 14}
37         font_label = {'fontweight': 'bold', '
    fontsize': 12}
38         numerical_columns = self.df.select_dtypes(
    include=[np.number])
39         for column in numerical_columns:
40             plt.figure()
41             sns.histplot(self.df[column])
42             plt.title(f"Histogram of {column}", **
    font_title)
43             plt.xlabel(column, **font_label)
44             plt.ylabel("Frequency", **font_label)
45             plt.show()
46         for column in numerical_columns:
47             plt.figure()
48             sns.boxplot(x=self.df[column])
49             plt.title(f"Box plot of {column}", **
    font_title)
50             plt.xlabel(column, **font_label)
51             plt.ylabel("Value", **font_label)
52             median_value = df[column].median()
53             plt.text(0.5, 0.5, f"Median: {
    median_value:.2f}", transform=plt.gca().transAxes,
54                     horizontalalignment='center',
    verticalalignment='center', **font_label)
55             plt.show()
56             sns.set(font_scale=1.2)
57             sns.pairplot(self.df[numerical_columns.
    columns])
58             plt.title("Scatter Plot Matrix", **
    font_title)
59             plt.show()
60             categorical_columns = df.select_dtypes(
    include=['object', 'category'])
61             for column in categorical_columns:
62                 plt.figure()
63                 sns.countplot(data=df, x=column)
64                 plt.title(f"Bar Graph of {column}", **
    font_title)
65                 plt.xlabel(column, **font_label)

```

```

66         plt.ylabel("Count", **font_label)
67         value_counts = df[column].value_counts
68         ()
69         for i, count in enumerate(value_counts
70 ):
71             plt.text(i, count, str(count), ha=
72 'center', va='bottom', **font_label)
73             plt.show()
74
75     def preprocess_data(self):
76         numerical_columns = self.df.select_dtypes(
77 include=[np.number])
78         scaler = MinMaxScaler()
79         self.df[numerical_columns.columns] =
80 scaler.fit_transform(self.df[numerical_columns.
81 columns])
82         categorical_columns = self.df.
83 select_dtypes(include=['object', 'category']).
84 columns.tolist()
85         self.df = pd.get_dummies(self.df, columns=
86 categorical_columns)
87
88     def find_optimal_clusters(self):
89         X = self.df.drop('Customer_Id', axis=1)
90         silhouette_scores = []
91         k_range = range(2, 21)
92         max_k = 2
93         max_score = -1
94         for k in k_range:
95             print(f"Running for k = {k}")
96             kmeans = KMeans(n_clusters=k)
97             kmeans.fit(X)
98             labels = kmeans.predict(X)
99             score = silhouette_score(X, labels)
100            silhouette_scores.append(score)
101            if score > max_score:
102                max_k = k
103                max_score = score
104        print("Silhouette scores: ",
105 silhouette_scores)
106        return max_k

```

```

97
98     def perform_cluster_and_analyze(self,
    optimal_clusters):
99         X = self.df.drop('Customer_Id', axis=1)
100         kmeans = KMeans(n_clusters=
    optimal_clusters)
101         kmeans.fit(X)
102         print("Number of customers in each cluster
    :", kmeans.labels_.size)
103         cluster_labels = kmeans.labels_
104         unique_labels, counts = np.unique(
    cluster_labels, return_counts=True)
105         for cluster_label, count in zip(
    unique_labels, counts):
106             print(f"Cluster {cluster_label}: {
    count} samples")
107             X['Cluster'] = cluster_labels
108             cluster_means = X.groupby('Cluster').mean
    ()
109             print(cluster_means)
110             cluster_means.to_csv('cluster_means.csv')
111
112
113 if __name__ == '__main__':
114     df = pd.read_csv("JustBuy_data.csv")
115     data_set = DataSet(df)
116     data_set.compute_statistical_summary()
117     data_set.visualize_data()
118     data_set.preprocess_data()
119     optimal_clusters = data_set.
    find_optimal_clusters()
120     print("Optimal number of clusters: ",
    optimal_clusters)
121     data_set.perform_cluster_and_analyze(
    optimal_clusters)
122
123

```