

DOCUMENTATION OF CLUSTERING CREDIT CARD

```
54 # Data Visualization
55 sns.histplot(data=df, x='BALANCE')
56 plt.title('Distribution of Balance')
57 plt.show()
58
59 sns.histplot(data=df, x='PURCHASES')
60 plt.title('Distribution of Purchases')
61 plt.show()
62
63 sns.histplot(data=df, x='CASH_ADVANCE')
64 plt.title('Distribution of Cash Advances')
65 plt.show()
66
67 # Based on Income Levels
68 b = [0, 50000, 100000, 150000, 200000, np.inf]
69 lab = ['<50K', '50K-100K', '100K-150K', '150K-200K', '>200K']
70 df['INCOME_LEVEL'] = pd.cut(df['PURCHASES']/df['TENURE']*12, bins=b, labels=lab)
71
72 plt.figure(figsize=(12,6))
73 sns.histplot(x='BALANCE', hue='INCOME_LEVEL', data=df, kde=True)
74 plt.title('Balance Distribution for different Income Levels')
75 plt.show()
76
77 plt.figure(figsize=(12,6))
78 sns.histplot(x='PURCHASES', hue='INCOME_LEVEL', data=df, kde=True)
79 plt.title('Purchases Distribution for different Income Levels')
80 plt.show()
81
82 plt.figure(figsize=(12,6))
83 sns.boxplot(x='INCOME_LEVEL', y='BALANCE', data=df)
84 plt.title('Balance Distribution for different Income Levels')
85 plt.show()
86
87 plt.figure(figsize=(12,6))
88 sns.boxplot(x='INCOME_LEVEL', y='PURCHASES', data=df)
89 plt.title('Purchases Distribution for different Income Levels')
90 plt.show()
91
92 plt.hist(df['CREDIT_LIMIT'], bins=20)
93 plt.title('Credit Limit Distribution')
94 plt.xlabel('Credit Limit')
95 plt.ylabel('Frequency')
96 plt.show()
```

Basic analysis obtained from the given graphs:

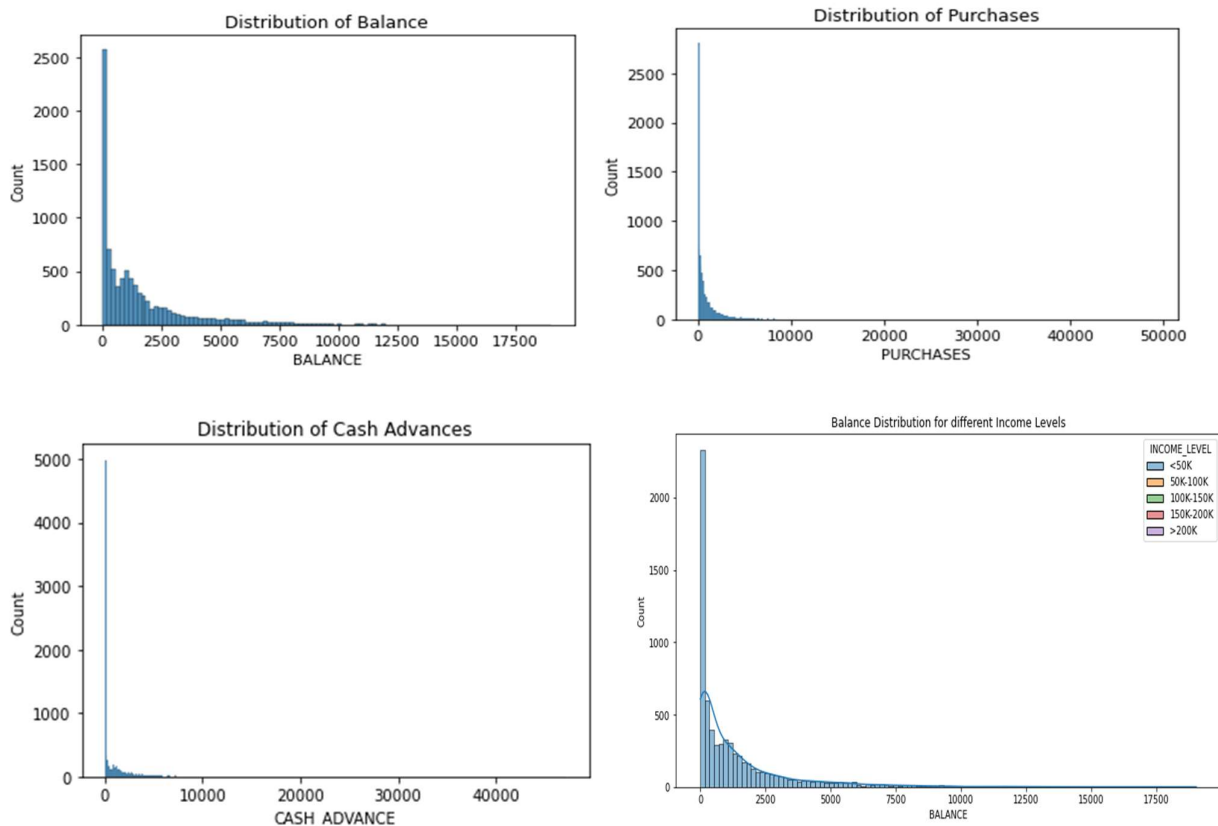
- **Distribution of Balance:** The histogram shows the distribution of balance values in the dataset, indicating that most customers have a balance between 0 and 5000, and there are some customers with very high balances (over 15000).
- **Distribution of Purchases:** The histogram shows the distribution of purchase values in the dataset, indicating that most customers have purchases between 0 and 2500, and there are some customers with very high purchase values (over 10000).
- **Distribution of Cash Advances:** The histogram shows the distribution of cash advance values in the dataset, indicating that most customers have cash

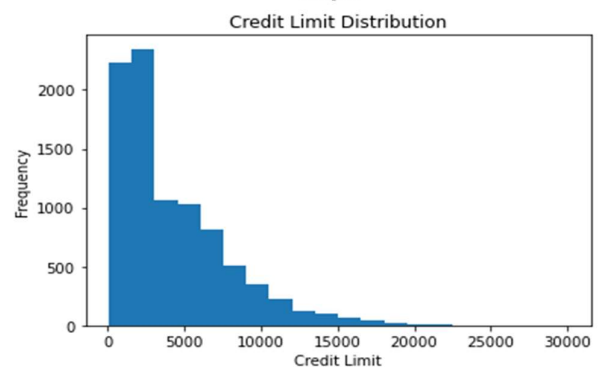
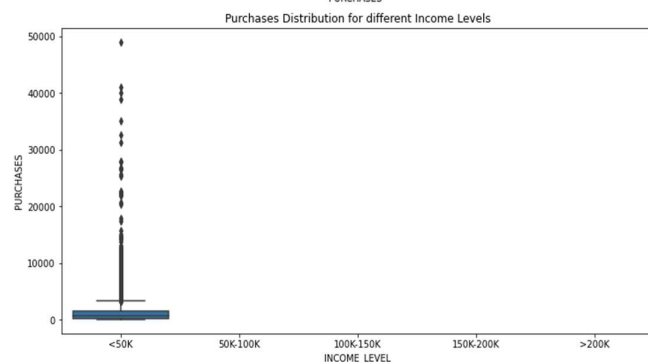
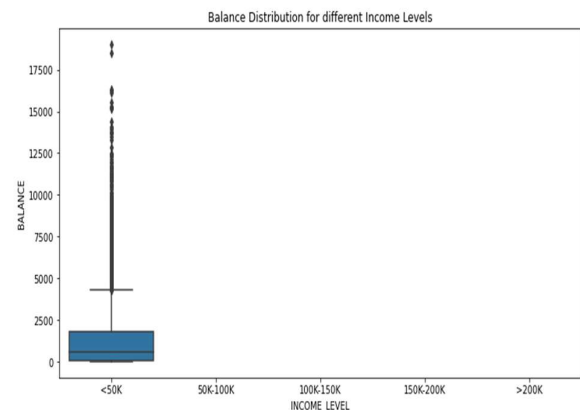
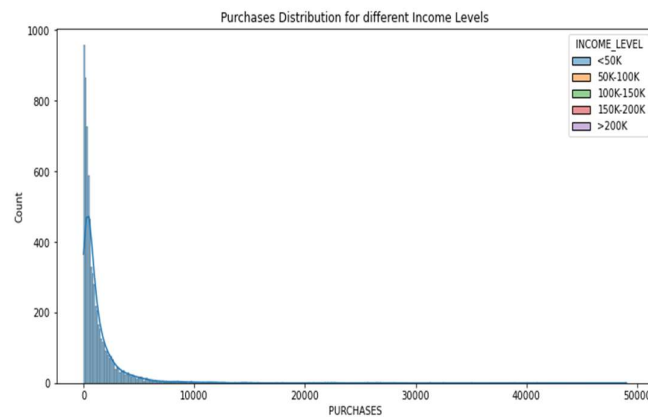
advances between 0 and 1000, and there are some customers with very high cash advance values (over 5000).

- **Balance and Purchases Distribution by Income Level:** The histograms show the distribution of balance and purchases values in the dataset, grouped by income level. The plots indicate that customers with higher income levels tend to have higher balances and make more purchases than customers with lower income levels.
- **Balance and Purchases Boxplot by Income Level:** The boxplots show the distribution of balance and purchases values in the dataset, grouped by income level. The plots indicate that customers with higher income levels tend to have higher median balances and make more purchases than customers with lower income levels.
- **Credit Limit Distribution:** The histogram shows the distribution of credit limit values in the dataset, indicating that most customers have a credit limit between 0 and 10000, and there are some customers with very high credit limits (over 20000).

This can be useful in clustering analysis as it can help in identifying the relevant features for clustering and in defining the clusters. For example, the balance and purchases distribution for different income levels suggest that these features could be important in clustering customers. The credit limit distribution could also be useful in identifying different groups of customers based on their creditworthiness.

Graphs:





```

98 # Clustering
99
100 wcss = []
101 for i in range(1, 11):
102     kmeans = KMeans(n_clusters=i, init='k-means++', max_iter=300, n_init=10, random_state=0)
103     kmeans.fit(df_sc)
104     wcss.append(kmeans.inertia_)
105
106 plt.plot(range(1, 11), wcss)
107 plt.title('Elbow Method')
108 plt.xlabel('Number of Clusters')
109 plt.ylabel('WCSS')
110 plt.show()
111
112 km = KMeans(n_clusters=3, init='k-means++', max_iter=300, n_init=10, random_state=0)
113 py = kmeans.fit_predict(df_sc)
114
115 df['Cluster'] = py
116 # Analyzing
117 df_cluster_0 = df[df['Cluster']==0]
118 df_cluster_1 = df[df['Cluster']==1]
119 df_cluster_2 = df[df['Cluster']==2]

```

This code is performing K-means clustering on the credit card dataset. The purpose of K-means clustering is to group similar customers together based on their purchasing behavior and other characteristics.

This code is using the elbow method to determine the optimal number of clusters to use for the K-means algorithm. This involves running the algorithm for a range of cluster numbers and plotting the within-cluster sum of squares (WCSS)* against the number of clusters. The "elbow" point on the graph is where the rate of decrease in WCSS slows down, indicating that additional clusters do not significantly improve the clustering performance. This elbow point is used to determine the optimal number of clusters for the data.

*(WCSS stands for "within-cluster sum of squares". In K-means clustering, WCSS is used as a metric to evaluate how well the data points within each cluster are grouped together.)

```
150 # Visualization
151 # 1. Hierarchical cluster
152 Z = linkage(df_sc, method='ward')
153 plt.figure(figsize=(12, 6))
154 dendrogram(Z)
155 plt.show()
156
157 # 2. DBSCAN cluster
158 db = DBSCAN(eps=0.5, min_samples=5)
159 clust = db.fit_predict(df_sc)
160 plt.scatter(df_sc[:,0], df_sc[:,1], c=clust, cmap='rainbow')
161 plt.xlabel('PCA Component 1')
162 plt.ylabel('PCA Component 2')
163 plt.show()
164
165 # 3. Spectral cluster
166 sp_c = SpectralClustering(n_clusters=3, affinity='nearest_neighbors', assign_labels='kmeans')
167 sp_c = sp_c.fit_predict(df_sc)
168 fig = plt.figure(figsize=(10,10))
169 ax = fig.add_subplot(111, projection='3d')
170 ax.scatter(df_sc[:, 0], df_sc[:, 1], df_sc[:, 2], c=sp_c, cmap='rainbow')
171 plt.show()
172
173 from sklearn.metrics import silhouette_score, silhouette_samples
174 import matplotlib.cm as cm
175
176 # silhouette
177 km = KMeans(n_clusters=3)
178 km.fit(df_sc)
179 labels = km.labels_
180 sil = silhouette_samples(df_sc, km.labels_)
181 avg = silhouette_score(df_sc, km.labels_)
182 fig, ax = plt.subplots(figsize=(8,6))
183 y_lower, y_upper = 0, 0
184 for i, cl in enumerate(np.unique(km.labels_)):
185     cl_vals = sil[km.labels_ == cl]
186     cl_vals.sort()
187     y_upper += len(cl_vals)
188     col = cm.nipy_spectral(float(i) / len(np.unique(km.labels_)))
189     ax.barh(range(y_lower, y_upper), cl_vals, height=1.0,
190            edgecolor='none', color=col)
191     y_lower += len(cl_vals)
192 ax.axvline(avg, color="red", linestyle="--")
193 ax.set_yticks([])
194 ax.set_xlim([-0.1, 1.0])
195 ax.set_xlabel("Silhouette coefficient values")
196 ax.set_ylabel("Cluster label")
197
198 plt.show()
```

From the graphs, following conclusions are drawn :

- **Elbow Method:** The elbow method is used to determine the optimal number of clusters for k-means clustering. In this case, the graph shows a sharp decrease in the Within-Cluster-Sum-of-Squares (WCSS) until 3 clusters, after which the decrease is less steep. Therefore, the optimal number of clusters for this dataset could be 3.
- **Hierarchical Clustering:** The dendrogram shows the relationship between the data points and how they are grouped based on their similarity. The height of the dendrogram represents the distance between the clusters. From the dendrogram, we can see that there are three main clusters.
- **DBSCAN:** Density-Based Spatial Clustering of Applications with Noise (DBSCAN) is a clustering algorithm that groups together points that are close to each other. In this case,

we can see that there are three main clusters and some noise points that are not part of any cluster.

- **Spectral Clustering:** Spectral clustering is a clustering algorithm that uses the eigenvectors of a similarity matrix to perform dimensionality reduction before clustering. The 3D scatter plot shows the clusters formed by the algorithm. From the plot, we can see that there are three distinct clusters.
- **Silhouette Analysis:** The silhouette analysis measures how well each data point fits into its assigned cluster. The graph shows the silhouette coefficient for each cluster, and a red line represents the average silhouette coefficient for all clusters. In this case, all clusters have a silhouette coefficient above the average, indicating that the clustering is effective.

GRAPHS:

