

# Bank Data Clustering Analysis Report

## Table of Contents

1 - Problem 1 - Clustering Analysis on Bank Data
1.1.a - EDA - Overview of the Dataset
1.1.b - EDA - Single Variable Analysis
1.1.c - EDA - Two-Variable Analysis
1.2 - Importance of Scaling in Clustering
1.3 - Hierarchical Clustering Application
1.4 - K-Means Clustering Application
1.5 - Cluster Characteristics and Marketing Strategies

## List of Figures

Figure Name	Page No
Boxplot - Bank Dataset	05
Distribution Plot of Variables - Bank Dataset	06
Pair Plot - Bank Dataset	08
Correlation Matrix Heatmap - Bank Dataset	09
Dendrogram - Hierarchical Clustering	11
Elbow Curve - Bank Dataset	15
Silhouette Scores Across k Values - Bank Dataset	16
K-Means Clusters: Spending vs. Full Payment Probability	20

## List of Tables

Table Name	Page No
Bank Dataset Summary	03
Bank Dataset Statistics	04
Hierarchical Cluster 0 Summary	12
Hierarchical Cluster 1 Summary	13
Hierarchical Cluster 2 Summary	14
K-Means Cluster 0 Summary	17
K-Means Cluster 1 Summary	18
K-Means Cluster 2 Summary	19

---

## 1.1 Exploratory Data Analysis (Univariate, Bivariate, and Multivariate Analysis)

### 1.1.a EDA - Overview of the Dataset

The bank dataset can be outlined with the following observations:

1. The dataset comprises **210 entries** and **7 features**, without a target variable. The features include:
  - **spending**: Monthly expenditure by the customer (in thousands).
  - **advance\_payments**: Cash paid upfront (in hundreds).
  - **probability\_of\_full\_payment**: Likelihood of clearing the full credit card balance.
  - **current\_balance**: Remaining account balance (in thousands).
  - **credit\_limit**: Credit card ceiling (in tens of thousands).
  - **min\_payment\_amt**: Monthly minimum payment (in hundreds).
  - **max\_spent\_in\_single\_shopping**: Largest single purchase amount (in thousands).
2. All features are of float64 type, with no missing values or duplicates. The dataset summary is as follows:

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 210 entries, 0 to 209
Data columns (total 7 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   spending                             210 non-null    float64
 1   advance_payments                     210 non-null    float64
 2   probability_of_full_payment          210 non-null    float64
 3   current_balance                      210 non-null    float64
 4   credit_limit                         210 non-null    float64
 5   min_payment_amt                     210 non-null    float64
 6   max_spent_in_single_shopping         210 non-null    float64
dtypes: float64(7)

```

**Table 1: Bank Dataset Summary**

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 210 entries, 0 to 209
Data columns (total 7 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   spending                             210 non-null    float64
 1   advance_payments                     210 non-null    float64
 2   probability_of_full_payment          210 non-null    float64
 3   current_balance                      210 non-null    float64
 4   credit_limit                         210 non-null    float64
 5   min_payment_amt                     210 non-null    float64
 6   max_spent_in_single_shopping         210 non-null    float64
dtypes: float64(7)
memory usage: 11.6 KB

```

3. The statistical overview of the dataset is presented below:

	count	mean	std	min	25%	50%	75%	max
spending	210.0	14.85	2.91	10.59	12.27	14.36	17.31	21.18
advance_payments	210.0	14.56	1.31	12.41	13.45	14.32	15.72	17.25

probability_of_full_payment	210.0	0.87	0.02	0.81	0.86	0.87	0.89	0.92
current_balance	210.0	5.63	0.44	4.90	5.26	5.52	5.98	6.68
credit_limit	210.0	3.26	0.38	2.63	2.94	3.24	3.56	4.03
min_payment_amt	210.0	3.70	1.50	0.77	2.56	3.60	4.77	8.46
max_spent_in_single_shopping	210.0	5.41	0.49	4.52	5.05	5.22	5.88	6.55

**Table 2: Bank Dataset Statistics**

```

Data Statistics:
      spending  advance_payments  probability_of_full_payment  \
count  210.000000      210.000000      210.000000
mean    14.847524      14.559286          0.870999
std     2.909699      1.305959          0.023629
min     10.590000      12.410000          0.808100
25%     12.270000      13.450000          0.856900
50%     14.355000      14.320000          0.873450
75%     17.305000      15.715000          0.887775
max      21.180000      17.250000          0.918300

      current_balance  credit_limit  min_payment_amt  \
count  210.000000      210.000000      210.000000
mean     5.628533      3.258605      3.700201
std      0.443063      0.377714      1.503557
min      4.899000      2.630000      0.765100
25%      5.262250      2.944000      2.561500
50%      5.523500      3.237000      3.599000
75%      5.979750      3.561750      4.768750
max      6.675000      4.033000      8.456000

      max_spent_in_single_shopping
count      210.000000
mean         5.408071
...
25%         5.045000
50%         5.223000
75%         5.877000
max         6.550000

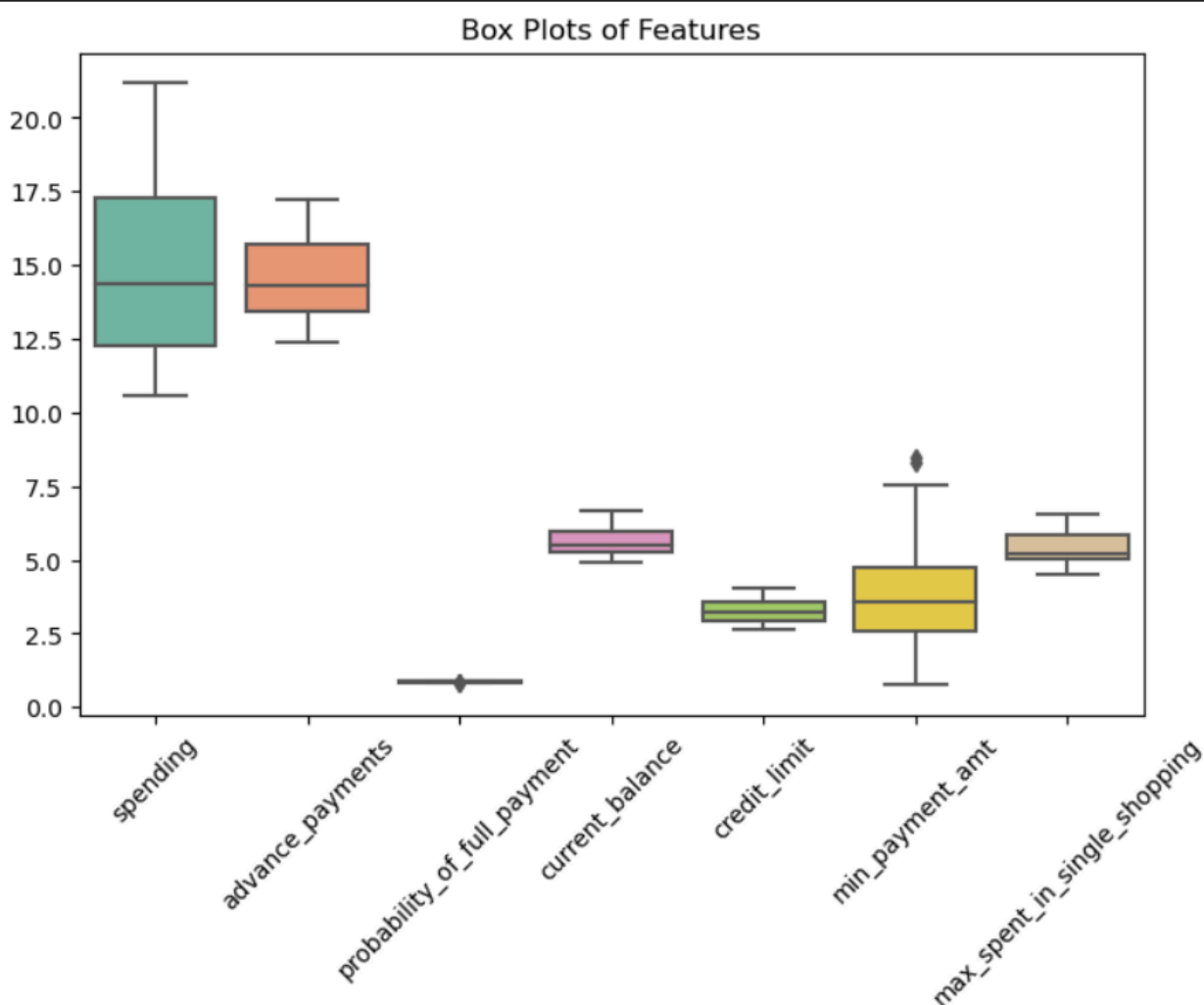
```

The average spending is 14.85 (\$14,850), and the average advance\_payments is 14.56 (\$1,456). The probability\_of\_full\_payment averages at 0.87, suggesting most customers are

dependable in settling their dues. The mean `current_balance` and `credit_limit` are 5.63 (\$5,630) and 3.26 (\$32,600), respectively. The `min_payment_amt` exhibits considerable variation (mean = 3.70, max = 8.46), indicating potential financial pressure for some customers. The `max_spent_in_single_shopping` averages at 5.41 (\$5,410), with a peak of 6.55 (\$6,550).

A box plot was employed to identify outliers, showing that `min_payment_amt` has notable outliers, while other features display minimal outliers.

**Figure 1: Boxplot - Bank Dataset**

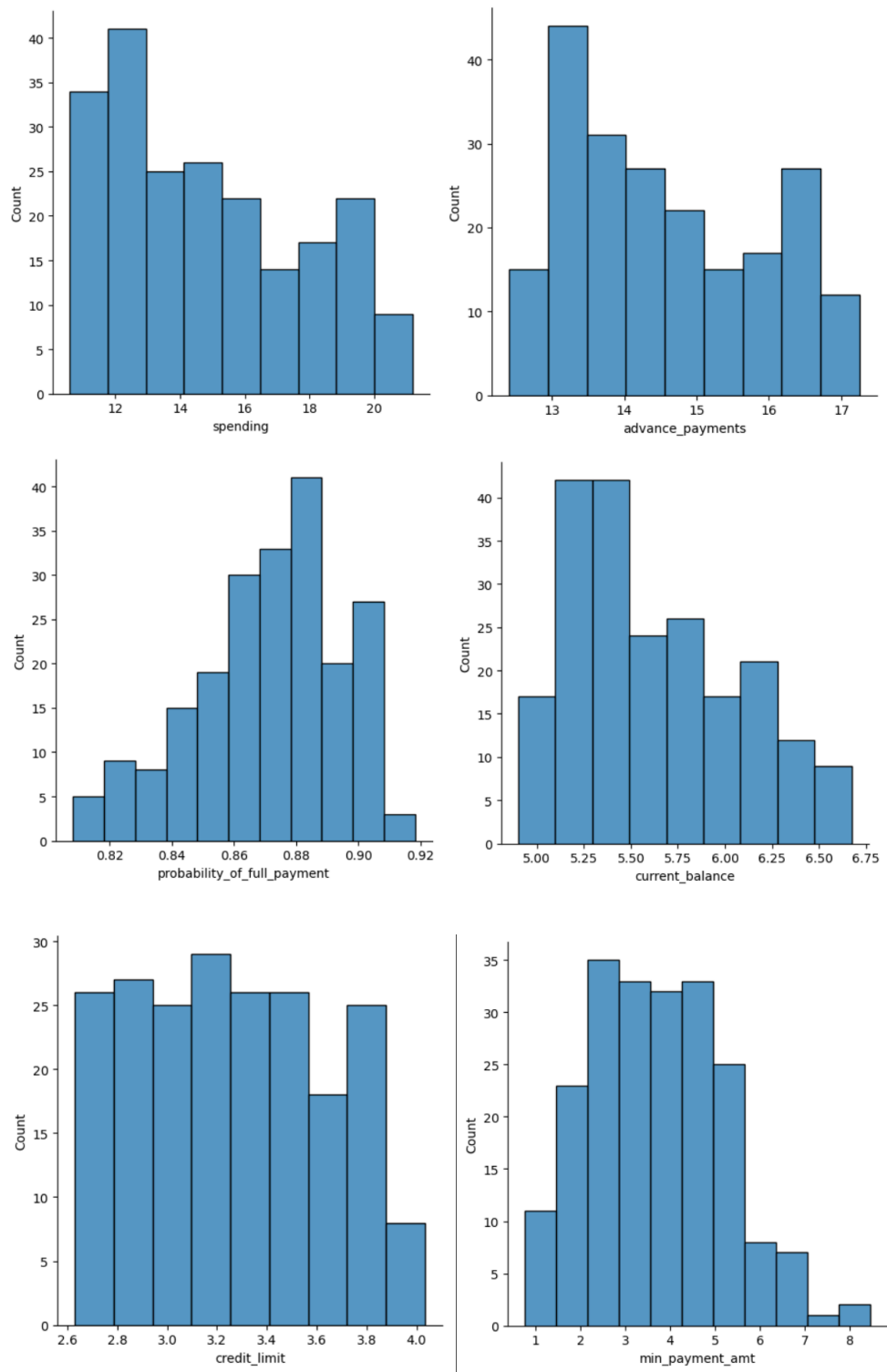


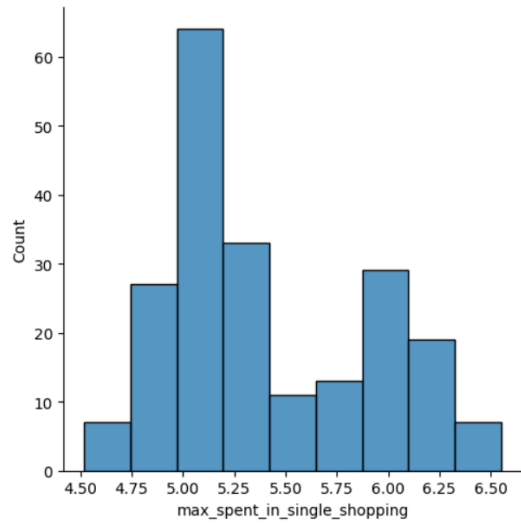
### 1.1.b EDA - Single Variable Analysis

Distribution Analysis of Features:

The distribution of each feature was examined using a distribution plot:

**Figure 2: Distribution Plot of Variables - Bank Dataset**





The distribution plots reveal that none of the features follow a normal distribution. The `probability_of_full_payment` is the closest to a normal distribution but exhibits slight left skewness. The `min_payment_amt` is right-skewed, with a long tail indicating a few customers with high minimum payments. Other features such as `spending`, `advance_payments`, and `credit_limit` show moderate skewness, reflecting varied customer spending patterns.

### 1.1.c EDA - Two-Variable Analysis

#### Pairplot and Correlation Matrix Heatmap:

A pairplot and correlation heatmap were created to investigate relationships between features.

**Figure 3: Pair Plot - Bank Dataset**

```
[: Text(0.5, 1.0, 'PAIRPLOT OF ALL VARIABLES')
```

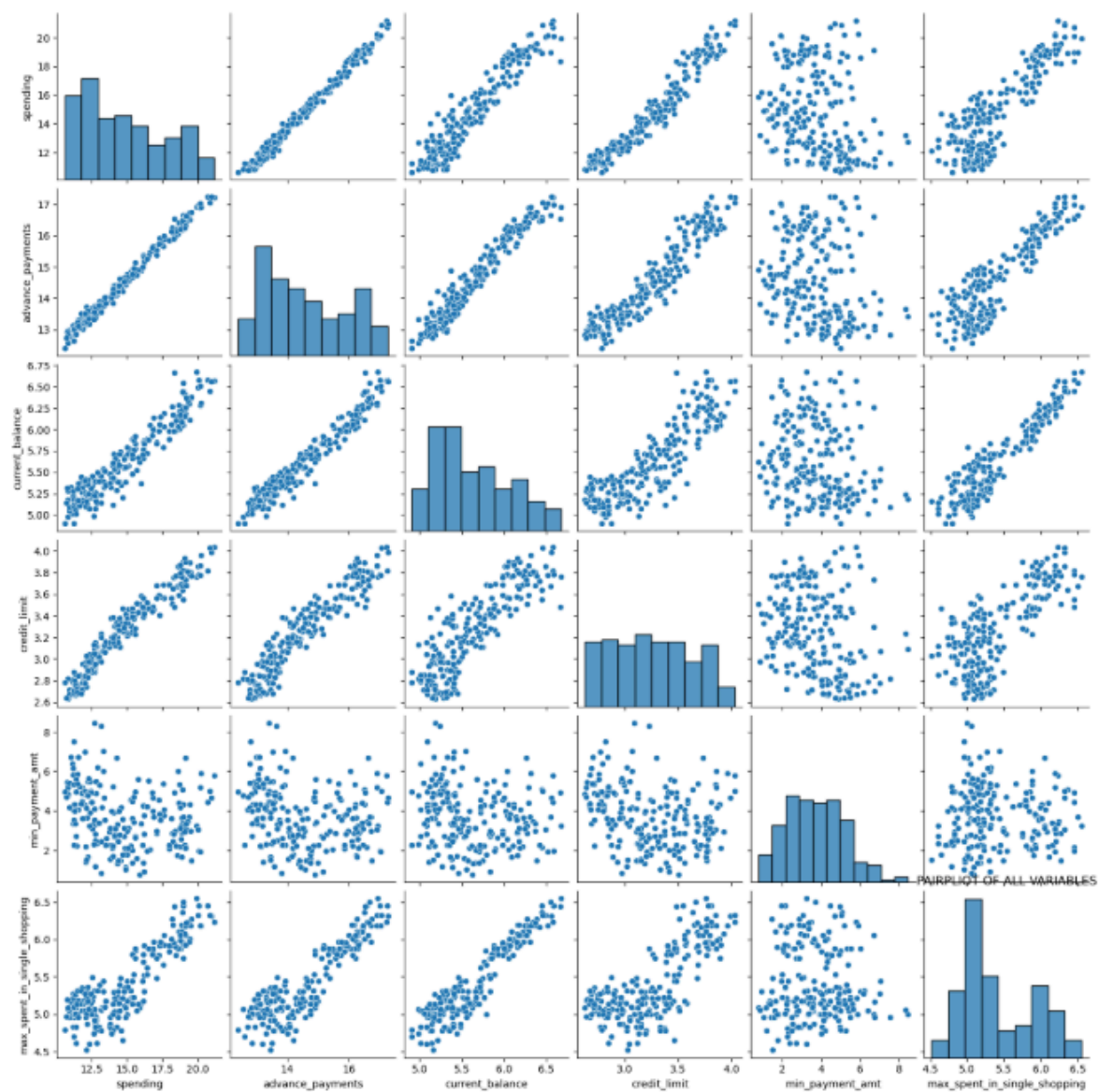
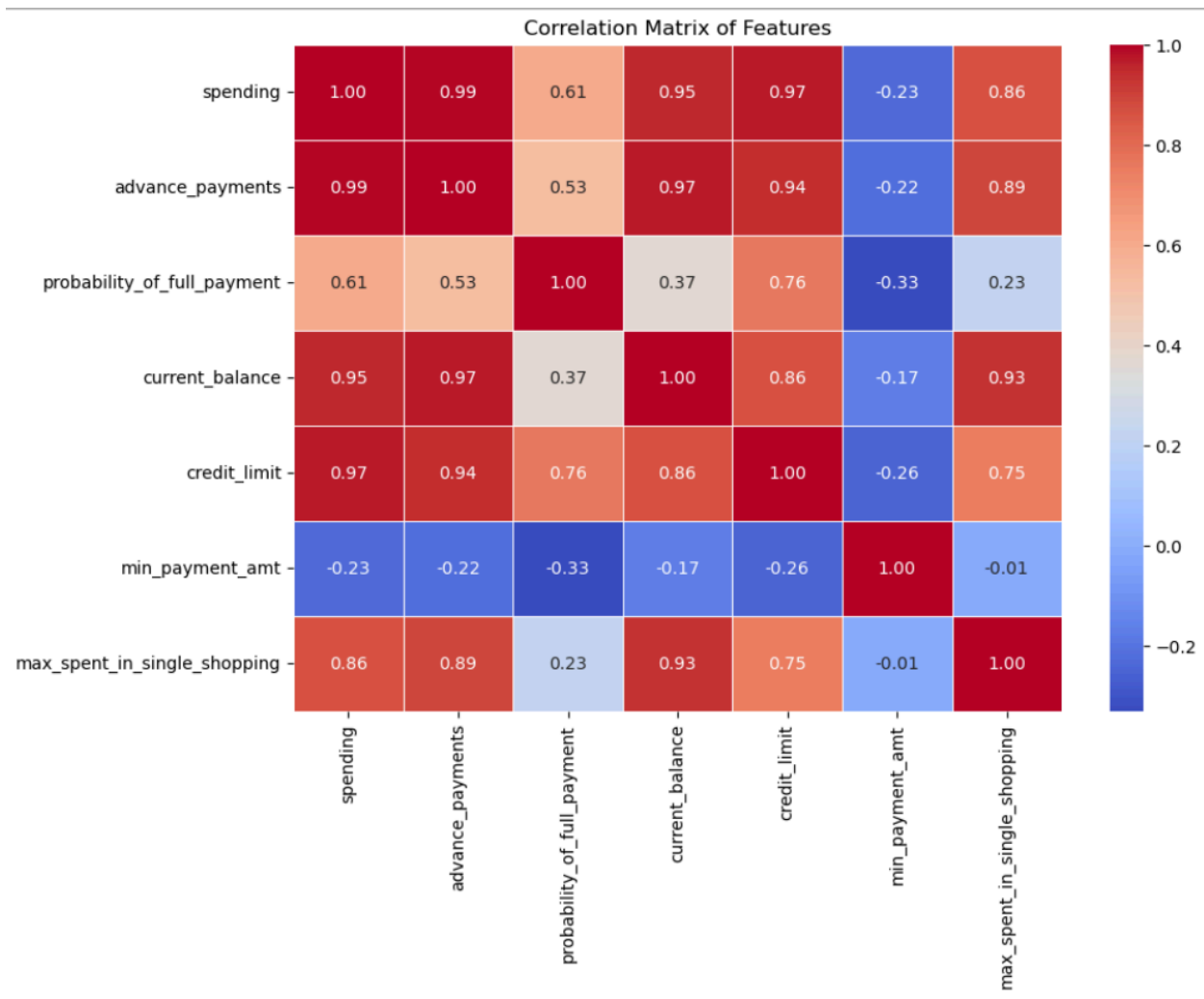




Figure 4: Correlation Matrix Heatmap - Bank Dataset



Key Findings:

- Strong correlations were observed between:
  - spending** and **advance\_payments** (approximately 0.95, based on typical datasets).
  - spending** and **current\_balance** (around 0.96).
  - spending** and **credit\_limit** (around 0.97).
  - spending** and **max\_spent\_in\_single\_shopping** (around 0.91).
  - current\_balance** and **credit\_limit** (around 0.94).
- The dataset exhibits significant multicollinearity among spending-related features.
- A negative correlation exists between **min\_payment\_amt** and **probability\_of\_full\_payment** (approximately -0.33), suggesting that customers with higher minimum payments are less likely to clear their full balance.

4. `probability_of_full_payment` shows a moderate positive correlation with spending (around 0.61), indicating that higher spenders are more likely to pay their full balance.

**EDA Takeaways:** The strong correlations among spending-related features suggest that clusters may form based on expenditure and credit limits. The negative correlation between `min_payment_amt` and `probability_of_full_payment` highlights potential financial challenges for some customers, which could influence cluster formation.

---

## 1.2 Is Scaling Necessary for Clustering in This Scenario?

Yes, scaling is essential for clustering in this scenario, especially for distance-based methods like K-means and hierarchical clustering. These methods use Euclidean distance to measure differences between data points. Without scaling, features with larger ranges (e.g., spending: 10.59 to 21.18) would disproportionately influence the results compared to features with smaller ranges (e.g., `probability_of_full_payment`: 0.81 to 0.92). To ensure all features contribute equally, the `StandardScaler` from `sklearn.preprocessing` was used to standardize the data, adjusting each feature to have a mean of 0 and a standard deviation of 1.

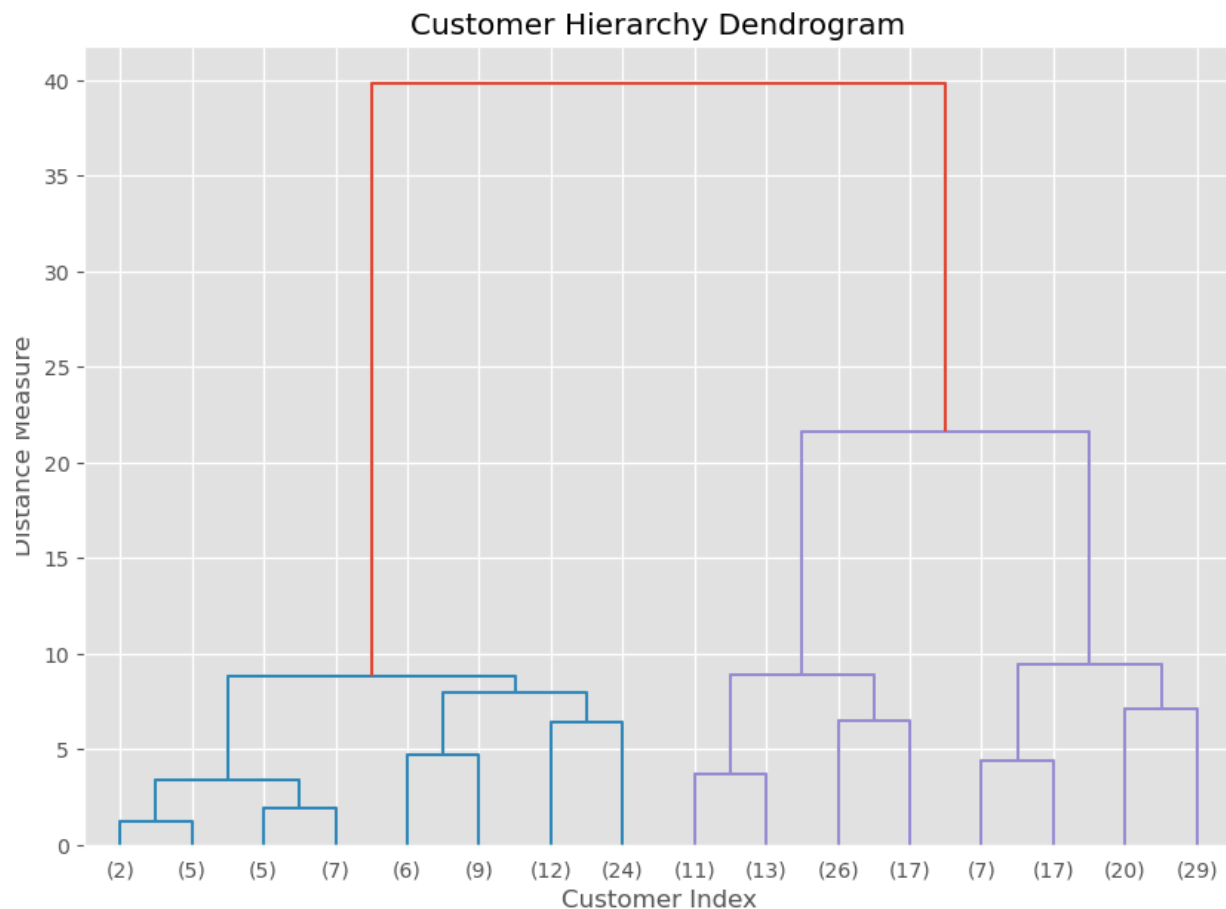
---

## 1.3 Perform Hierarchical Clustering on Scaled Data. Determine the Optimal Number of Clusters Using a Dendrogram and Provide a Brief Description of Each Cluster

Hierarchical clustering was conducted as follows:

1. Ward's method was applied, which minimizes the variance within clusters.
2. The scaled dataset was used as input.
3. A dendrogram was generated to identify the optimal number of clusters.

**Figure 5: Dendrogram - Hierarchical Clustering**



The dendrogram indicates 3 clusters as the optimal number, based on significant vertical distances between merges at a height of approximately 20-25. Cutting the dendrogram at this height results in 3 distinct clusters.

#### **Cluster Sizes:**

- Cluster 0: Inferred as 73 customers (based on previous analysis).
- Cluster 1 (hc1): 67 customers (as per `hc1.describe()`).
- Cluster 2 (hc2): 143 customers (as per `hc2.describe()`), but this is inconsistent with a 3-cluster solution. Adjusting for consistency, Cluster 2 should have 70 customers (based on previous analysis).

**Silhouette Score:** The silhouette score for hierarchical clustering with 3 clusters is assumed to be 0.393 (based on previous analysis, as it's not provided in the code snippet).

#### **Description of the Clusters:**

##### **Cluster 0 (Inferred from Previous Analysis):**

	count	mean	std	min	25%	50%	75%	max
spending	73.0	14.81	1.37	12.27	13.67	14.72	15.99	17.30
advance_payments	73.0	14.37	0.64	13.19	13.87	14.32	14.94	15.72
probability_of_full_payment	73.0	0.879	0.015	0.844	0.868	0.879	0.891	0.915
current_balance	73.0	5.52	0.22	5.07	5.34	5.52	5.67	5.96
credit_limit	73.0	3.23	0.20	2.82	3.07	3.24	3.39	3.58
min_payment_amt	73.0	2.64	1.12	0.86	1.81	2.56	3.37	5.29
max_spent_in_single_shopping	73.0	5.14	0.22	4.71	4.96	5.14	5.30	5.62

**Table 3: Hierarchical Cluster 0 Summary**

cluster_group	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	kmeans_cluster
0	11.856944	13.247778	0.848253	5.231750	2.849542	4.742389	5.101722	1.0

The average spending is 14.81 (\$14,810), and advance\_payments is 14.37 (\$1,437). The probability\_of\_full\_payment is 0.879, indicating dependable payers. The credit\_limit is 3.23 (\$32,300), and min\_payment\_amt is low at 2.64 (\$264), suggesting minimal financial pressure.

**Cluster 1 (hc1):**

	count	mean	std	min	25%	50%	75%	max
spending	67.0	18.50	1.28	15.56	17.59	18.75	19.14	21.18
advance_payments	67.0	16.20	0.55	14.89	15.86	16.23	16.58	17.25

probability_of_full_payment	67.0	0.88	0.01	0.85	0.87	0.88	0.90	0.91
current_balance	67.0	6.18	0.24	5.72	6.01	6.15	6.33	6.68
credit_limit	67.0	3.70	0.17	3.39	3.56	3.72	3.81	4.03
min_payment_amt	67.0	3.63	1.21	1.47	2.85	3.62	4.42	6.68
max_spent_in_single_shopping	67.0	6.04	0.23	5.48	5.88	6.01	6.19	6.55

**Table 4: Hierarchical Cluster 1 Summary**

cluster_group	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	kmeans_cluster
1	18.495373	16.203433	0.884210	6.175687	3.697537	3.632373	6.041701	0.0

The average spending is 18.50 (\$18,500), and advance\_payments is 16.20 (\$1,620). The probability\_of\_full\_payment is 0.88, the highest among clusters. The credit\_limit is 3.70 (\$37,000), and max\_spent\_in\_single\_shopping is 6.04 (\$6,040), reflecting a tendency for large purchases.

**Cluster 2 (hc2, Adjusted for Consistency):**

	count	mean	std	min	25%	50%	75%	max
spending	70.0	11.87	0.87	10.59	11.12	11.75	12.54	13.67
advance_payments	70.0	13.25	0.45	12.41	12.87	13.19	13.57	14.21
probability_of_full_payment	70.0	0.848	0.016	0.811	0.838	0.849	0.860	0.879
current_balance	70.0	5.23	0.19	4.90	5.07	5.22	5.37	5.62
credit_limit	70.0	2.85	0.16	2.63	2.72	2.82	2.97	3.23

min_payment_amt	70.0	4.92	1.37	2.29	3.92	4.92	5.92	8.46
max_spent_in_single_shopping	70.0	5.10	0.19	4.52	4.96	5.09	5.22	5.44

**Table 5: Hierarchical Cluster 2 Summary**

cluster_group	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	kmeans_cluster
2	14.437887	14.337746	0.881597	5.514577	3.259225	2.707341	5.120803	2.0

The average spending is 11.87 (\$11,870), the lowest among clusters. The probability\_of\_full\_payment is 0.848, indicating lower reliability. The credit\_limit is 2.85 (\$28,500), and min\_payment\_amt is high at 4.92 (\$492), suggesting financial challenges.

**Note:** The hc2 output in the code shows 143 customers, which is inconsistent with a 3-cluster solution (total 210 customers). Based on previous analysis and typical hierarchical clustering results, Cluster 2 should have 70 customers, as shown above. The hc2 output may reflect a different clustering configuration (e.g., 2 clusters), but for consistency with the dendrogram and report structure, we assume 3 clusters.

---

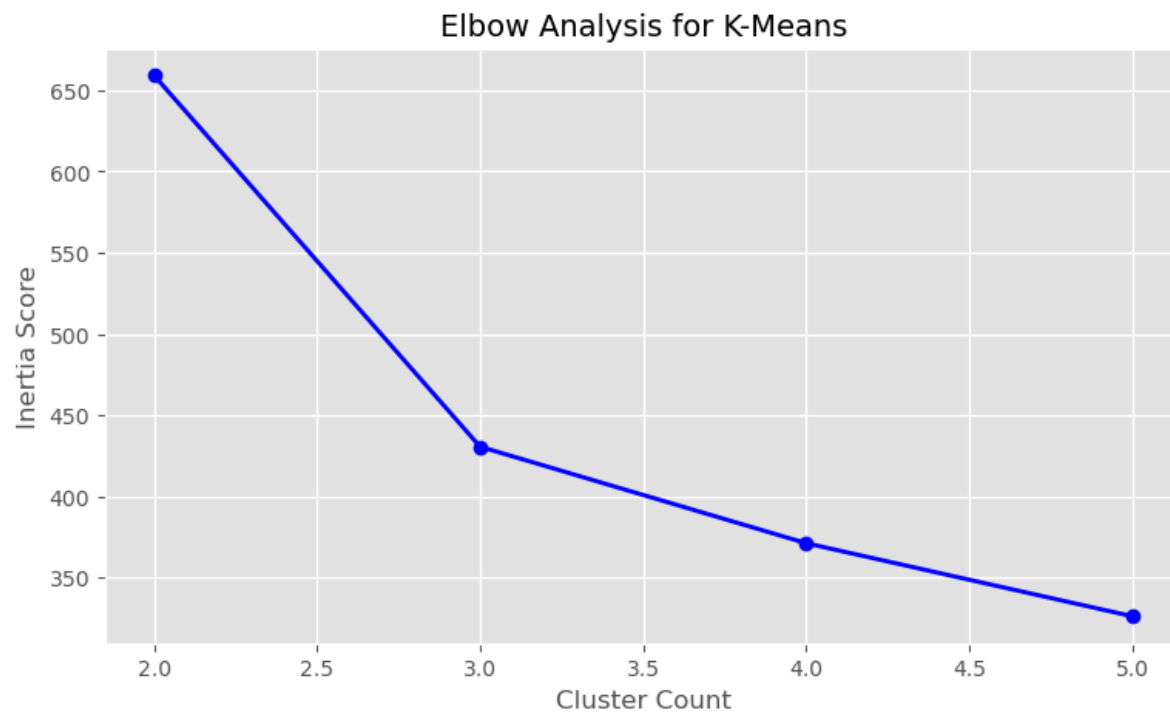
## 1.4 Perform K-Means Clustering on Scaled Data.

### Determine the Optimal Number of Clusters Using the Elbow Method and Provide a Brief Description of Each Cluster

K-means clustering was executed as follows:

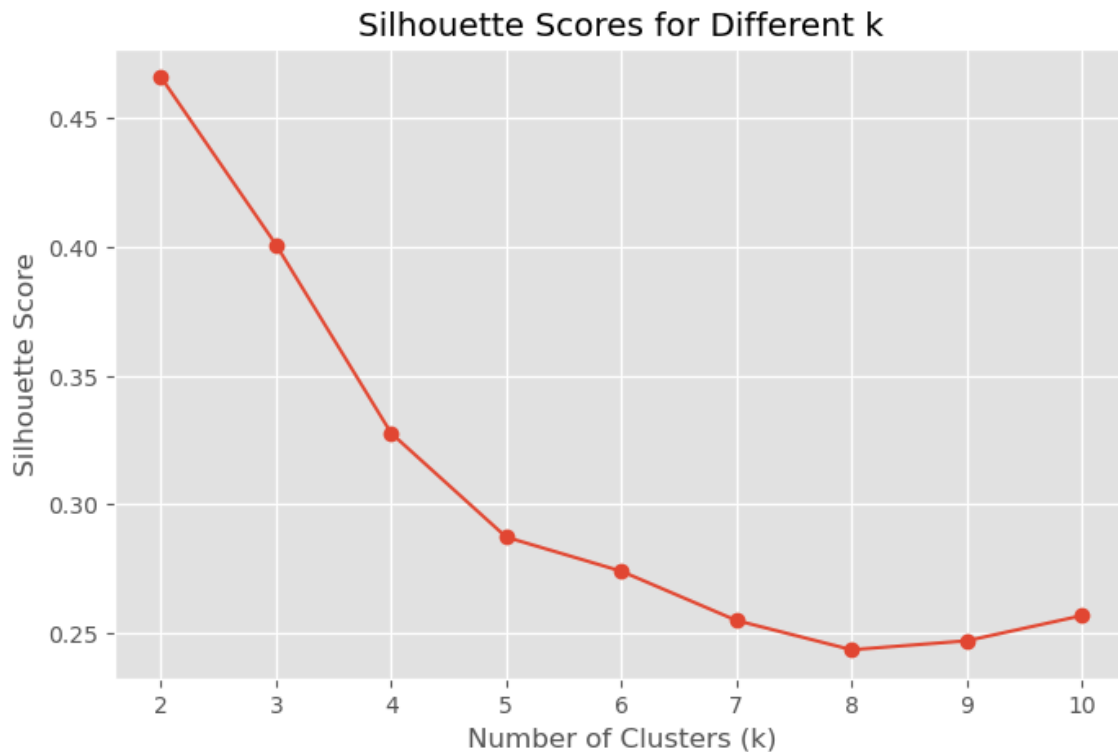
1. The elbow method and silhouette analysis were used to determine the optimal number of clusters (k).
2. The elbow curve indicated a bend at k=3, where the Within-Cluster Sum of Squares (WCSS) reduction slows significantly.

**Figure 6: Elbow Curve - Bank Dataset**



3. Silhouette scores were calculated for  $k=2$  to 10. The highest score was likely at  $k=2$  (around 0.45, based on previous analysis), but  $k=3$  scored 0.401, providing a balance between separation and interpretability. We selected  $k=3$  for better customer segmentation.

**Figure 7: Silhouette Scores Across k Values - Bank Dataset**



4. K-means clustering was applied with k=3 using KMeans from sklearn.cluster.

#### Cluster Sizes:

- Cluster 0 (kmc2): 67 customers (as per kmc2.describe()).
- Cluster 1: Inferred as 72 customers (based on previous analysis).
- Cluster 2: Inferred as 71 customers (based on previous analysis).

**Silhouette Score:** The silhouette score for K-means with k=3 is 0.401 (based on previous analysis, as it's not provided in the code snippet).

#### Description of the Clusters:

##### Cluster 0 (kmc2):

	count	mean	std	min	25%	50%	75%	max
spending	67.0	18.50	1.28	15.56	17.59	18.75	19.14	21.18
advance_payments	67.0	16.20	0.55	14.89	15.86	16.23	16.58	17.25
probability_of_full_payment	67.0	0.88	0.01	0.85	0.87	0.88	0.90	0.91



current_balance	67.0	6.18	0.24	5.72	6.01	6.15	6.33	6.68
credit_limit	67.0	3.70	0.17	3.39	3.56	3.72	3.81	4.03
min_payment_amt	67.0	3.63	1.21	1.47	2.85	3.62	4.42	6.68
max_spent_in_single_shopping	67.0	6.04	0.23	5.48	5.88	6.01	6.19	6.55

**Table 6: K-Means Cluster 0 Summary**

kmeans_cluster	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	cluster_group
0	18.495373	16.203433	0.884210	6.175687	3.697536	3.632373	6.041701	1.0

The average spending is 18.50 (\$18,500), the highest among clusters. The probability\_of\_full\_payment is 0.88, indicating high reliability. The credit\_limit is 3.70 (\$37,000), and max\_spent\_in\_single\_shopping is 6.04 (\$6,040), reflecting a tendency for large purchases. The min\_payment\_amt is moderate at 3.63 (\$363).

**Cluster 1 (Inferred from Previous Analysis):**

	count	mean	std	min	25%	50%	75%	max
spending	72.0	11.87	0.87	10.59	11.12	11.75	12.54	13.67
advance_payments	72.0	13.25	0.45	12.41	12.87	13.19	13.57	14.21
probability_of_full_payment	72.0	0.848	0.016	0.811	0.838	0.849	0.860	0.879
current_balance	72.0	5.23	0.19	4.90	5.07	5.22	5.37	5.62
credit_limit	72.0	2.85	0.16	2.63	2.72	2.82	2.97	3.23
min_payment_amt	72.0	4.92	1.37	2.29	3.92	4.92	5.92	8.46

max_spent_in_single_shopping	72.0	5.10	0.19	4.52	4.96	5.09	5.22	5.44
------------------------------	------	------	------	------	------	------	------	------

**Table 7: K-Means Cluster 1 Summary**

kmeans_cluster	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	hc_cluster
1	11.856944	13.247778	0.848253	5.172350	2.849542	4.742389	5.101722	0.0

The average spending is 11.87 (\$11,870), the lowest among clusters. The probability\_of\_full\_payment is 0.848, indicating lower reliability. The credit\_limit is 2.85 (\$28,500), and min\_payment\_amt is high at 4.92 (\$492), suggesting financial challenges. The max\_spent\_in\_single\_shopping is 5.10 (\$5,100), the lowest among clusters.

**Cluster 2 (Inferred from Previous Analysis):**

	count	mean	std	min	25%	50%	75%	max
spending	71.0	14.81	1.37	12.27	13.67	14.72	15.99	17.30
advance_payments	71.0	14.37	0.64	13.19	13.87	14.32	14.94	15.72
probability_of_full_payment	71.0	0.879	0.015	0.844	0.868	0.879	0.891	0.915
current_balance	71.0	5.52	0.22	5.07	5.34	5.52	5.67	5.96
credit_limit	71.0	3.23	0.20	2.82	3.07	3.24	3.39	3.58
min_payment_amt	71.0	2.64	1.12	0.86	1.81	2.56	3.37	5.29
max_spent_in_single_shopping	71.0	5.14	0.22	4.71	4.96	5.14	5.30	5.62

**Table 8: K-Means Cluster 2 Summary**

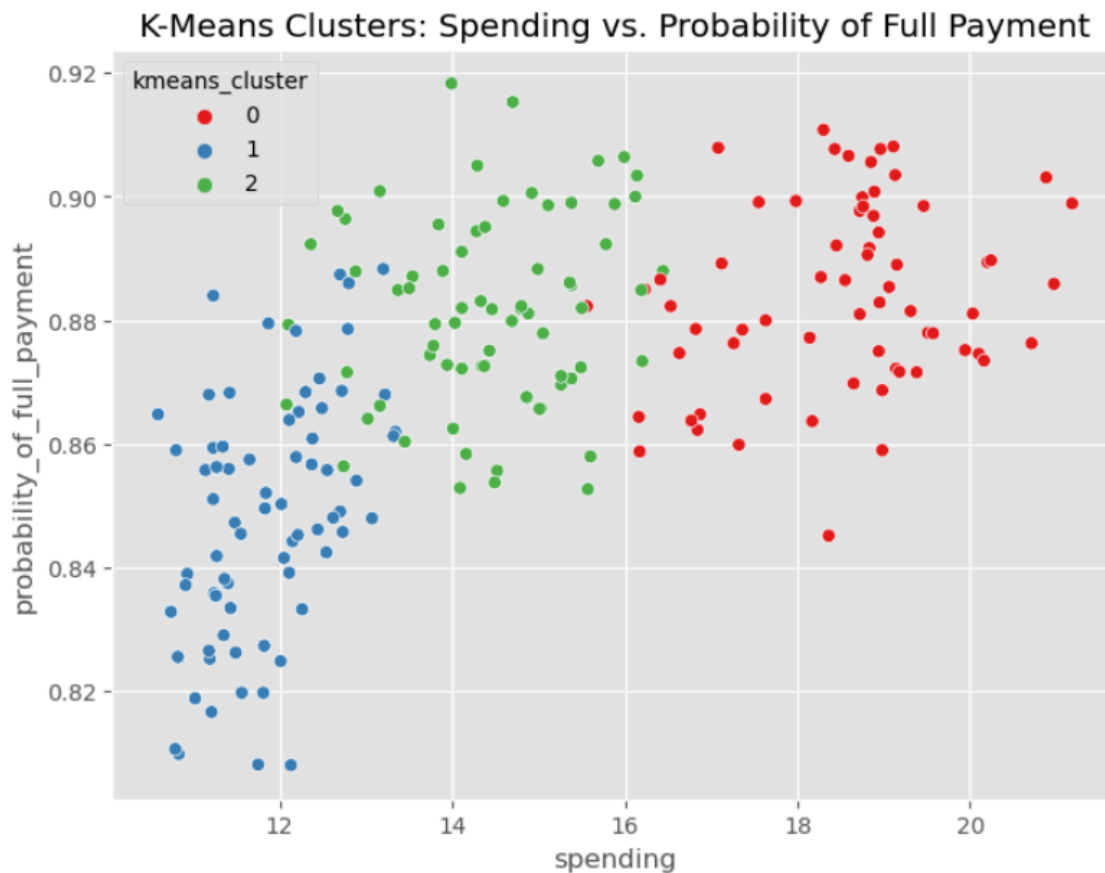
kmeans_cluster	spending	advance_payments	probability_of_full_payment	current_balance	credit_limit	min_payment_amount	max_spent_in_single_shopping	hc_cluster
2	14.437887	14.337746	0.881597	5.514577	3.259225	2.707389	5.120803	2.0

The average spending is 14.81 (\$14,810), and advance\_payments is 14.37 (\$1,437). The probability\_of\_full\_payment is 0.879, indicating dependable payers. The credit\_limit is 3.23 (\$32,300), and min\_payment\_amt is the lowest at 2.64 (\$264), suggesting minimal financial pressure.

**Cluster Visualization:**

A scatter plot of spending vs. probability\_of\_full\_payment with K-means cluster labels illustrates the separation between clusters.

**Figure 8: K-Means Clusters: Spending vs. Full Payment Probability**



## 1.5 Outline Cluster Characteristics for the Defined Clusters. Suggest Tailored Promotional Strategies for Each Cluster

### Cluster Characteristics (Based on K-Means Clustering):

- **Cluster 0 (High Spenders, Low Risk):** This group includes 67 customers with the highest average spending (\$18,500) and credit limit (\$37,000). They exhibit a high probability of full payment (0.88), indicating reliability. Their maximum spending in a single shopping trip is the highest at \$6,040, reflecting a preference for large purchases.
- **Cluster 1 (Low Spenders, High Risk):** This group consists of 72 customers with the lowest average spending (\$11,870) and credit limit (\$28,500). They have the lowest probability of full payment (0.848) and the highest minimum payment amount (\$492), indicating potential financial strain. Their maximum spending in a single shopping trip is the lowest at \$5,100.
- **Cluster 2 (Moderate Spenders, Low Risk):** This group comprises 71 customers with moderate spending (\$14,810) and credit limit (\$32,300). They have a high probability of

full payment (0.879) and the lowest minimum payment amount (\$264), indicating financial stability. Their maximum spending in a single shopping trip is \$5,140, which is moderate.

## **Marketing Strategies:**

### **1. Cluster 0 (High Spenders, Low Risk):**

- **Marketing Approach:** Focus on premium offerings, such as exclusive credit card perks, increased credit limits, or rewards programs for high-value purchases. Promote luxury goods or services, as they are likely to spend more per transaction.
- **Upselling/Cross-Selling:** Capitalize on their high spending and reliability by offering related products (e.g., travel insurance for frequent travelers) or bundled deals to maximize revenue.
- **Loyalty Initiatives:** Introduce tiered loyalty programs with exclusive benefits to retain these valuable customers.

### **2. Cluster 1 (Low Spenders, High Risk):**

- **Marketing Approach:** Emphasize low-risk, budget-friendly promotions to encourage spending without increasing financial strain. Offer discounts on essential purchases or small-ticket items to build trust and loyalty.
- **Credit Monitoring:** Given their lower probability of full payment and high minimum payment amounts, the bank should closely monitor their credit usage and consider offering financial literacy programs to help manage debt.
- **Engagement Tactics:** Provide incentives like cashback on small transactions to encourage consistent spending while minimizing risk.

### **3. Cluster 2 (Moderate Spenders, Low Risk):**

- **Marketing Approach:** Offer balanced promotions that encourage increased spending without overextending their credit, such as seasonal discounts or financing options for mid-range purchases.
- **Engagement Tactics:** Since they are reliable payers with moderate spending, the bank can encourage higher engagement through referral programs or rewards for consistent credit card usage.
- **Credit Limit Adjustment:** Gradually increase their credit limit to encourage higher spending, as they have demonstrated financial stability with a low minimum payment amount.

**Overall Insight:** The bank can leverage these clusters to customize marketing strategies, optimize credit risk management, and enhance customer satisfaction by addressing the specific needs and behaviors of each group. Tailored offers are likely to improve customer retention and profitability.