


```
In [35]: # Counting the num of customers in each segment
segment_counts = dfm_data['Segment'].value_counts()
segment_counts
Out[35]:
Standard    1215
Silver      134
Gold         68
Premium     51
Name: Segment, dtype: int64
```

```
In [36]: # Creating a bar plot to visualize the distribution of segments
plt.figure(figsize=(10, 6))
segment_counts.plot(kind='bar', color='skyblue')
plt.title('Customer Segmentation')
plt.xlabel('Segments')
plt.ylabel('Number of Customers')
plt.xticks(rotation=45)
plt.show()
```



Strategy for each segment

- Premium:** We must focus on retaining these high-value customers with personalized offers and excellent service. Encourage them to refer friends and family and give them discounts in return.
- Gold:** We must come out with incentives for increased spending and engagement. It is a must to offer them loyalty rewards to maintain their loyalty.
- Silver:** we must keep these customers engaged by sending them regular promotions and updates. We must encourage them to upgrade to higher segments.
- Standard:** Target these customers with promotions to increase their spending and engagement. Focus on converting them into higher segments.

Part 2: K-means(Scientific) Segmentation

This is the second type of segmentation we are going to apply to categorize customers into optimal number of clusters. The numeric variables will be scaled and compressed to Principal components and then clusters would be made on the basis of euclidean distance between points.

```
In [85]: # Calculating days since last transaction
max_date = merged_data['Transaction_Date'].max()
merged_data['Recency'] = (max_date - merged_data['Transaction_Date']).dt.days
merged_data['Recency']

Out[85]:
0      364
1      364
2      364
3      364
4      364
...
630683    0
630684    0
630685    0
630686    0
630687    0
Name: Recency, Length: 630688, dtype: int64
```

```
In [86]: # Dropping the unnecessary columns
merged_data=merged_data.drop(['Transaction_Date','CustomerID','Transaction_ID'], axis=1)
merged_data

Out[86]:
```

	Gender	Location	Tenure_Months	Invoice_Value	totalmarketingspend	Recency
0	M	Chicago	12	1515.229	6924.50	364
1	M	Chicago	12	3206.039	6924.50	364
2	M	Chicago	12	4896.849	6924.50	364
3	M	Chicago	12	1515.229	6924.50	364
4	M	Chicago	12	3206.039	6924.50	364
...
630683	F	California	7	6711.482	6058.75	0
630684	F	California	7	10254.362	6058.75	0
630685	F	California	7	3168.602	6058.75	0
630686	F	California	7	6711.482	6058.75	0
630687	F	California	7	10254.362	6058.75	0

630688 rows x 6 columns

```
In [88]: # Perform one-hot encoding
merged_data = pd.get_dummies(merged_data, columns= ('Gender', 'Location'))
merged_data

Out[88]:
```

	Tenure_Months	Invoice_Value	totalmarketingspend	Recency	Gender_F	Gender_M	Location_California	Location_Chicago	Location_New_Jersey
0	12	1515.229	6924.50	364	0	1	0	0	1
1	12	3206.039	6924.50	364	0	1	0	0	1
2	12	4896.849	6924.50	364	0	1	0	0	1
3	12	1515.229	6924.50	364	0	1	0	0	1
4	12	3206.039	6924.50	364	0	1	0	0	1
...
630683	7	6711.482	6058.75	0	1	0	1	0	0
630684	7	10254.362	6058.75	0	1	0	1	0	0
630685	7	3168.602	6058.75	0	1	0	1	0	0
630686	7	6711.482	6058.75	0	1	0	1	0	0
630687	7	10254.362	6058.75	0	1	0	1	0	0

630688 rows x 11 columns

```
In [145]: # Standardizing the data
scaler = StandardScaler()
merged_data_scaled = scaler.fit_transform(merged_data)
merged_data_scaled

Out[145]:
array([[ -1.04939813, -0.08891677,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.35842471,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.79760219,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.35842471,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.42048942,  0.34452883,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ],
       [ -1.42048942,  1.27332267,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ],
       [ -1.42048942,  2.20211744,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ]])
```

```
In [146]: # Create an imputer object
imputer = SimpleImputer(strategy='mean')

# Fit and transform the imputer on your data
merged_data_scaled = imputer.fit_transform(merged_data_scaled)
merged_data_scaled

Out[146]:
array([[ -1.04939813, -0.08891677,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.35842471,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.79760219,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.04939813,  0.35842471,  1.54086686, ..., -0.51753946,
        -0.23320675, -0.77371162],
       [ -1.42048942,  0.34452883,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ],
       [ -1.42048942,  1.27332267,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ],
       [ -1.42048942,  2.20211744,  0.93396959, ..., -0.51753946,
        -0.23320675,  1.7131431 ]])
```

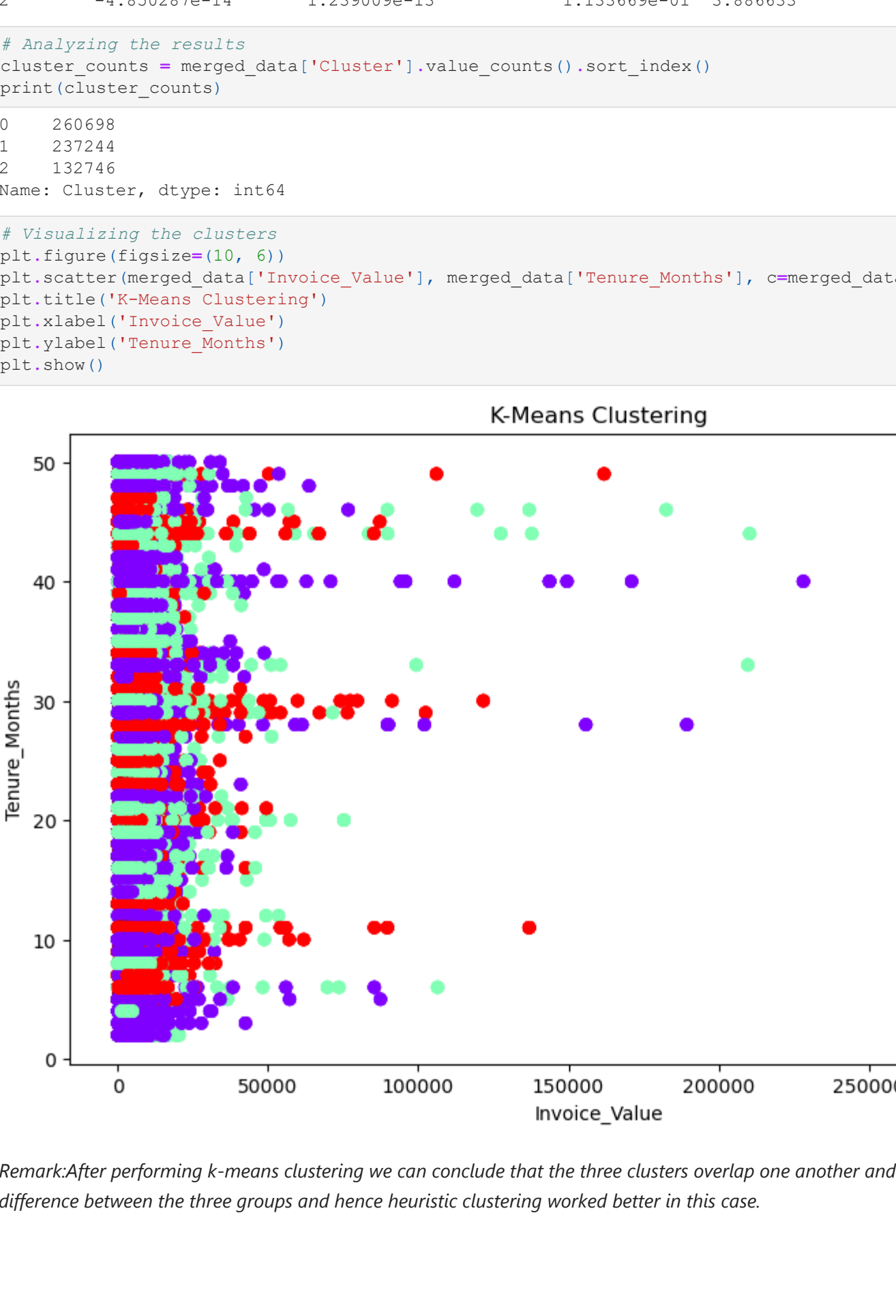
```
In [149]: # Determining the optimal number of clusters using the Elbow method
wcss = []
for i in range(1, 7):
    kmeans = KMeans(n_clusters=i, init='k-means++', random_state=42)
    kmeans.fit(merged_data_scaled)
    wcss.append(kmeans.inertia_)
```

```
In [150]: # Plotting the optimal clusters graph
plt.figure(figsize=(8, 8))
plt.plot(range(1, 7), wcss, marker='o', linestyle='--')
plt.title('Elbow Method')
plt.xlabel('Number of Clusters')
plt.ylabel('WCSS (Within-Cluster Sum of Squares)')

diff = [wcss[i] - wcss[i - 1] for i in range(1, len(wcss))]
optimal_num_clusters = diff.index(max(diff)) + 1

plt.axvline(x=optimal_num_clusters, color='red', linestyle='--')

plt.grid()
plt.show()
```



```
In [151]: # Finding the "elbow" point by looking for the steepest slope
diff = [wcss[i] - wcss[i - 1] for i in range(1, len(wcss))]
optimal_num_clusters = diff.index(max(diff)) + 1
optimal_num_clusters

Out[151]:
3
```

```
In [152]: # Plotting the K-Means clustering results with the optimal number of clusters
kmeans_optimal = KMeans(n_clusters=optimal_num_clusters, init='k-means++', random_state=42)
kmeans_optimal.fit(merged_data_scaled)
labels = kmeans_optimal.labels_
```

```
In [153]: # Based on the Elbow method the optimal number has been set
optimal_clusters = 3
```

```
In [154]: # Performing K-Means clustering with the optimal number of clusters
kmeans = KMeans(n_clusters=optimal_clusters, init='k-means++', random_state=42)
merged_data['Cluster'] = kmeans.fit_predict(merged_data_scaled)
```

```
In [157]: # Viewing the cluster centers
cluster_centers = pd.DataFrame(
    scaler.inverse_transform(kmeans.cluster_centers_),
    columns=merged_data.columns
)
print(cluster_centers)
```

	Tenure_Months	Invoice_Value	totalmarketingspend	Recency
0	25.949355	1794.781643	4703.756100	175.987027
1	26.878421	1863.326819	4663.678146	184.980189
2	25.191622	1955.536061	4883.087213	179.271255

	Gender_F	Gender_M	Location_California	Location_Chicago
0	0.000000e+00	4.040657e-13	-7.494005e-14	5.250251e-01
1	-6.936675e-13	1.000000e+00	3.141365e-01	3.462553e-01
2	1.000000e+00	1.859646e-13	8.865331e-01	1.394468e-13

	Location_New_Jersey	Location_New_York	Location_Washington_DC	Cluster
0	1.400778e-01	3.348971e-01	-3.868433e-14	1.719844
1	7.230952e-02	1.936108e-01	7.368785e-02	1.075066
2	-4.850287e-14	1.239009e-13	1.133669e-01	3.886633

```
In [158]: # Analyzing the results
cluster_counts = merged_data['Cluster'].value_counts().sort_index()
print(cluster_counts)
```

	Cluster
0	260658
1	237244
2	132746

Name: Cluster, dtype: int64

```
In [159]: # Visualizing the clusters
plt.figure(figsize=(10, 6))
plt.scatter(merged_data['Invoice_Value'], merged_data['Tenure_Months'], c=merged_data['Cluster'], cmap='rainbow')
plt.title('K-Means Clustering')
plt.xlabel('Invoice_Value')
plt.ylabel('Tenure_Months')
plt.show()
```

