

1 # Hierarchical clustering(Dendrogram) Agglomerative clustering

In [11]:

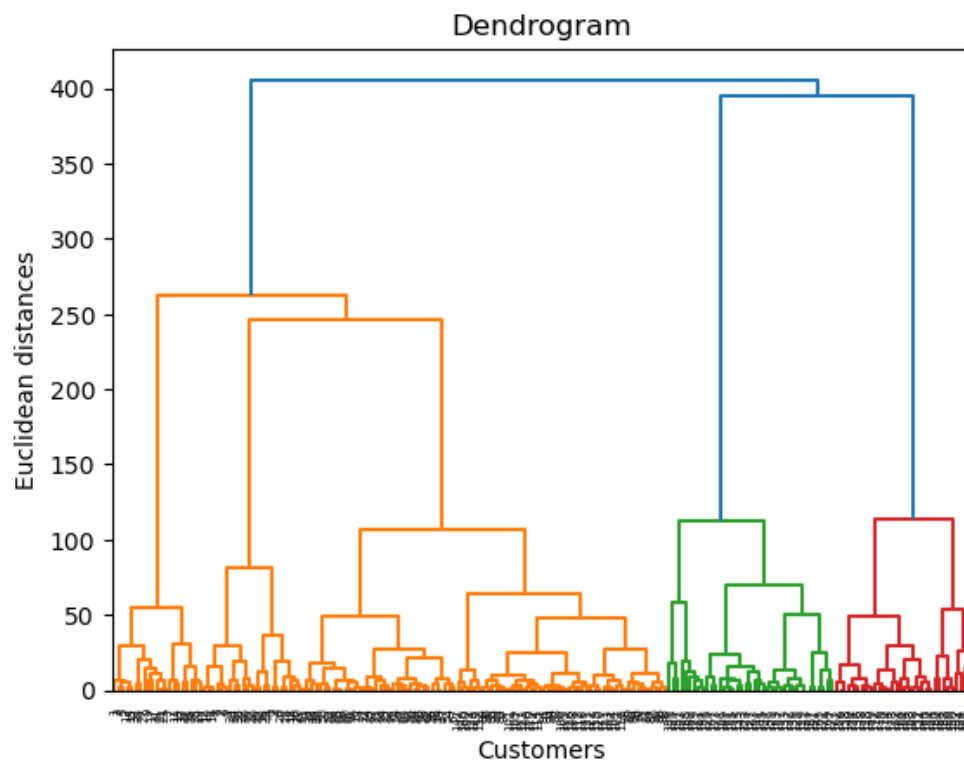
```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
```

In [12]:

```
1 dataset = pd.read_csv('/Users/myyntiimac/Desktop/Mall_Customers.csv')
2 X = dataset.iloc[:, [3, 4]].values
```

In [13]:

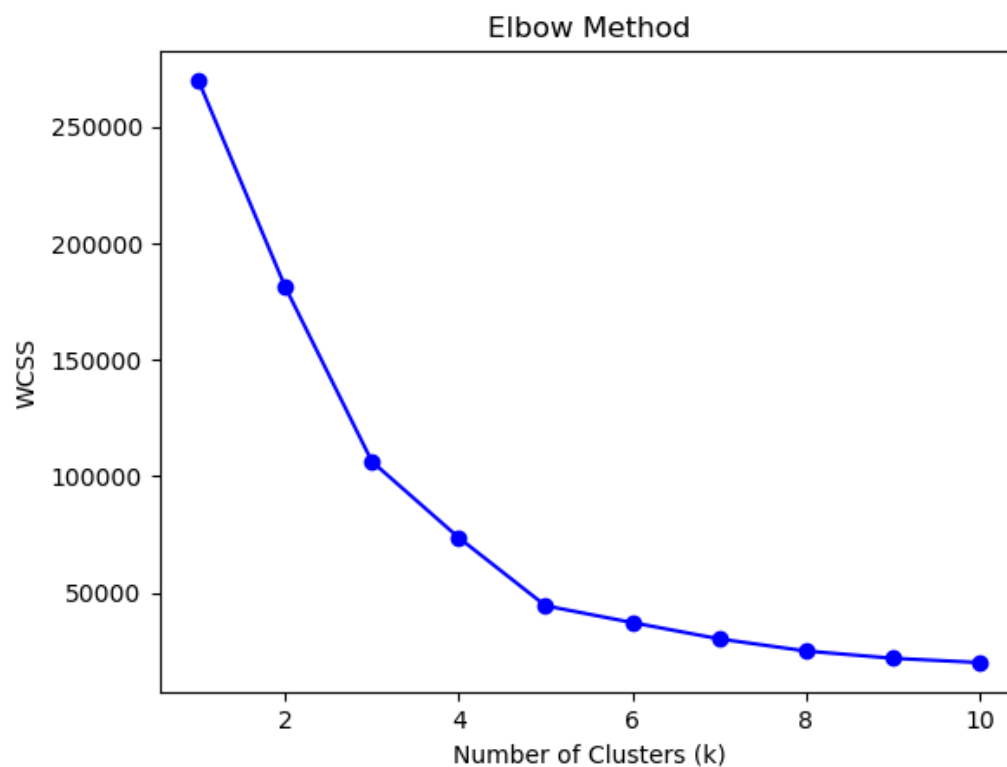
```
1 import scipy.cluster.hierarchy as sch
2 dendrogram = sch.dendrogram(sch.linkage(X, method = 'ward'))
3 plt.title('Dendrogram')
4 plt.xlabel('Customers')
5 plt.ylabel('Euclidean distances')
6 plt.show()
```



```
1 #to determine cut level or threshold , implement elbow method
2 Elbow method: Calculate a measure of dissimilarity (e.g., within-cluster sum of squares,
WCSS) for different cut levels. Plot the dissimilarity measure against the number of
clusters. Look for a point where the rate of decrease in dissimilarity slows down
significantly (forming an "elbow" shape). This point can be considered as a potential cut
level.
```

In [14]:

```
1 import warnings
2
3 # Ignore all warnings
4 warnings.filterwarnings("ignore")
5 from sklearn.cluster import KMeans
6
7 # Assuming you have your data stored in 'X'
8 # X should be a 2D array or matrix with shape (n_samples, n_features)
9
10 # Initialize an empty list to store the WCSS values for different numbers of clusters
11 wcss = []
12
13 # Define the range of cluster numbers to try
14 k_values = range(1, 11) # Try cluster numbers from 1 to 10
15
16 # Calculate WCSS for each cluster number
17 for k in k_values:
18     kmeans = KMeans(n_clusters=k, random_state=42)
19     kmeans.fit(X)
20     wcss.append(kmeans.inertia_) # Inertia is the WCSS value
21
22 # Plot the WCSS values against the number of clusters
23 plt.plot(k_values, wcss, 'bo-')
24 plt.xlabel('Number of Clusters (k)')
25 plt.ylabel('WCSS')
26 plt.title('Elbow Method')
27 plt.show()
28
29
30
31
32
33
```



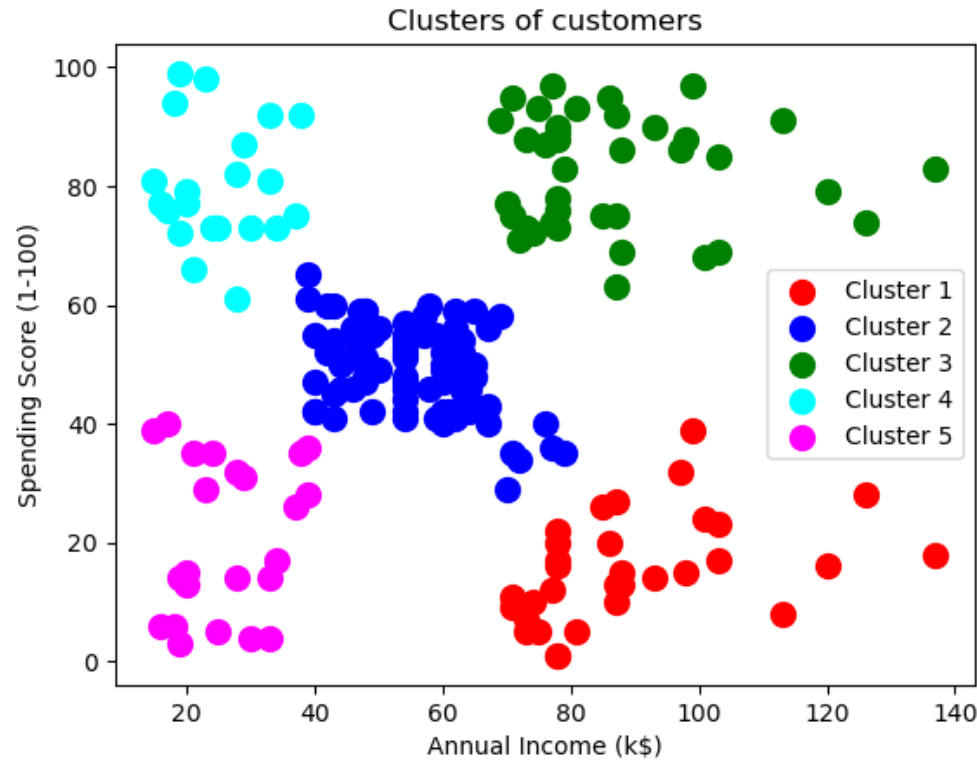
1 insight: from this figure , we assume cluster no is 5

In [15]:

```
1 #training the data with aggluramative clustering
2 from sklearn.cluster import AgglomerativeClustering
3 hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean', linkage = 'ward')
4 y_hc = hc.fit_predict(X)
```

In [17]:

```
1 # Visualizing the clusters
2 plt.scatter(X[y_hc == 0, 0], X[y_hc == 0, 1], s = 100, c = 'red', label = 'Cluster 1')
3 plt.scatter(X[y_hc == 1, 0], X[y_hc == 1, 1], s = 100, c = 'blue', label = 'Cluster 2')
4 plt.scatter(X[y_hc == 2, 0], X[y_hc == 2, 1], s = 100, c = 'green', label = 'Cluster 3')
5 plt.scatter(X[y_hc == 3, 0], X[y_hc == 3, 1], s = 100, c = 'cyan', label = 'Cluster 4')
6 plt.scatter(X[y_hc == 4, 0], X[y_hc == 4, 1], s = 100, c = 'magenta', label = 'Cluster 5')
7 plt.title('Clusters of customers')
8 plt.xlabel('Annual Income (k$)')
9 plt.ylabel('Spending Score (1-100)')
10 plt.legend()
11 plt.show()
```



In [18]:

```
1 df1=pd.read_csv("/Users/myyntiimac/Desktop/modified_dataset.csv")
2 df1.head()
```

Out[18]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	Cluster
0	1	Male	19	15	39	2
1	2	Male	21	15	81	3
2	3	Female	20	16	6	2
3	4	Female	23	16	77	3
4	5	Female	31	17	40	2

In [19]:

```
1 df1["agguluramative"] = y_hc
```

In [20]:

```
1 df1.head()
```

Out[20]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)	Cluster	agguluramative
0	1	Male	19	15	39	2	4
1	2	Male	21	15	81	3	3
2	3	Female	20	16	6	2	4
3	4	Female	23	16	77	3	3
4	5	Female	31	17	40	2	4

In [21]:

```
1 df1.to_csv('new_dataframe.csv', index=False)
```

In [22]:

```
1 import os
2
3 # Get the current working directory
4 current_dir = os.getcwd()
5
6 # Specify the filename
7 filename = 'new_dataframe.csv'
8
9 # Combine the directory and filename to get the full file path
10 file_path = os.path.join(current_dir, filename)
11
12 # Print the file path
13 print("File saved at:", file_path)
```

File saved at: /Users/myyntiimac/new_dataframe.csv

In []:

```
1
```