

```
In [1]: import pandas as pd
```

```
In [2]: import numpy as np
```

```
In [3]: import matplotlib.pyplot as plt
```

```
In [4]: import seaborn as sns
```

```
In [8]: df = pd.read_csv(r'C:\Users\Rutu\Documents\New folder\Mall_Customers.csv')
```

```
In [9]: df.head()
```

Out[9]:

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

```
In [10]: df.tail()
```

Out[10]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

```
In [11]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   CustomerID                            200 non-null    int64
1   Genre                                 200 non-null    object
2   Age                                   200 non-null    int64
3   Annual Income (k$)                   200 non-null    int64
4   Spending Score (1-100)               200 non-null    int64
dtypes: int64(4), object(1)
memory usage: 7.9+ KB
```

```
In [12]: df.isnull()
```

Out[12]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	False	False	False	False	False
1	False	False	False	False	False
2	False	False	False	False	False

3	False	False	False	False	False
4	False	False	False	False	False
...
195	False	False	False	False	False
196	False	False	False	False	False
197	False	False	False	False	False
198	False	False	False	False	False
199	False	False	False	False	False

200 rows × 5 columns

In [13]: `df.dropna()`

Out[13]:

	CustomerID	Genre	Age	Annual Income (k\$)	Spending Score (1-100)
0	1	Male	19	15	39
1	2	Male	21	15	81
2	3	Female	20	16	6
3	4	Female	23	16	77
4	5	Female	31	17	40
...
195	196	Female	35	120	79
196	197	Female	45	126	28
197	198	Male	32	126	74
198	199	Male	32	137	18
199	200	Male	30	137	83

200 rows × 5 columns

In [15]: `df.describe()`

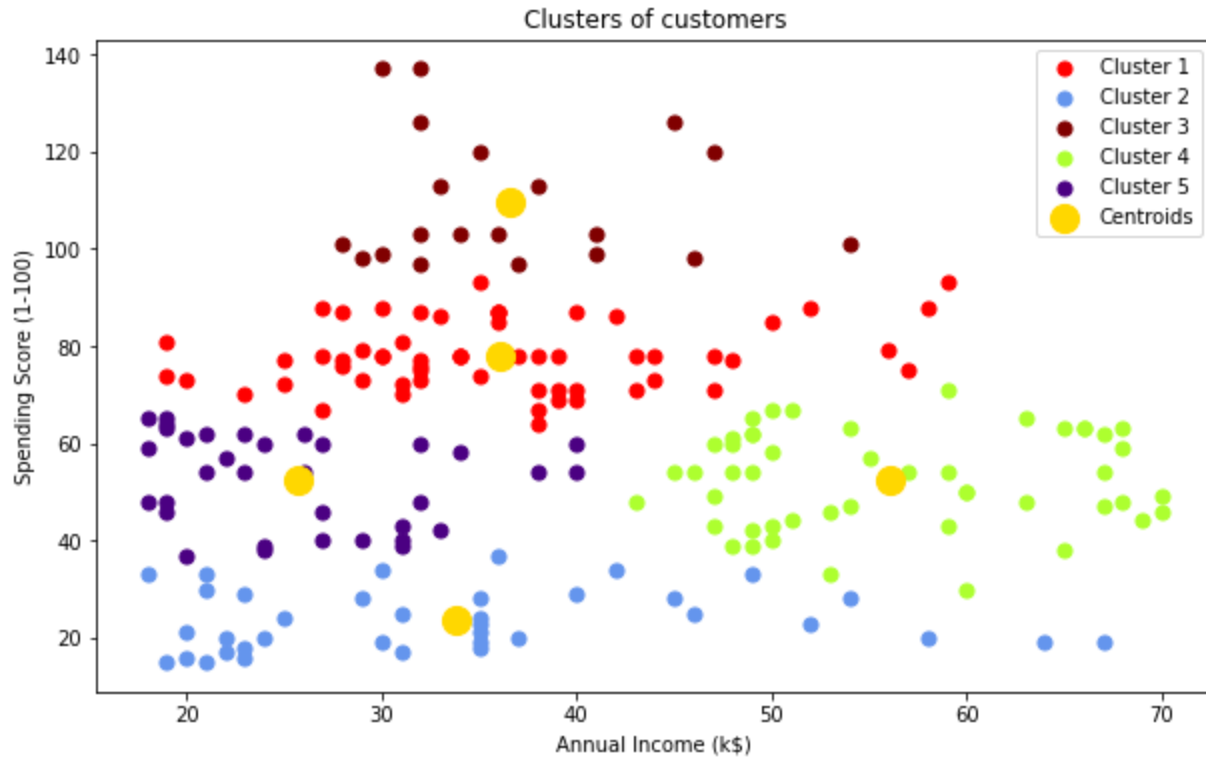
Out[15]:

	CustomerID	Age	Annual Income (k\$)	Spending Score (1-100)
count	200.000000	200.000000	200.000000	200.000000
mean	100.500000	38.850000	60.560000	50.200000
std	57.879185	13.969007	26.264721	25.823522
min	1.000000	18.000000	15.000000	1.000000
25%	50.750000	28.750000	41.500000	34.750000
50%	100.500000	36.000000	61.500000	50.000000
75%	150.250000	49.000000	78.000000	73.000000
max	200.000000	70.000000	137.000000	99.000000

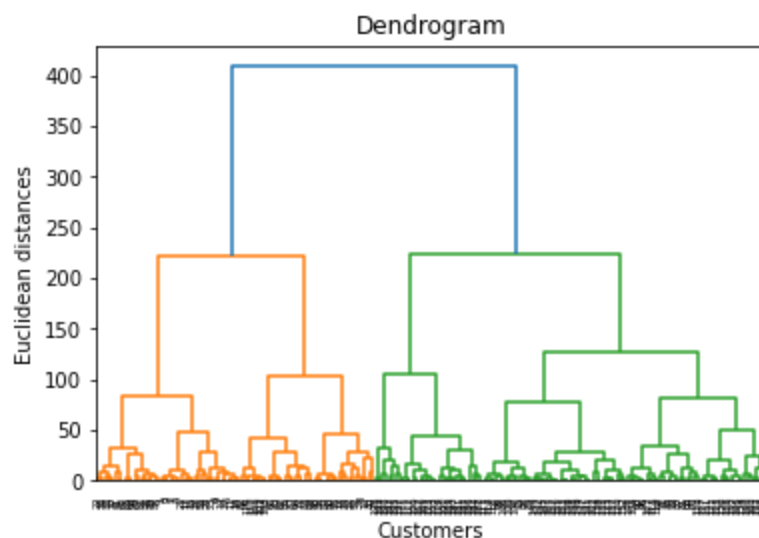
In [16]: `df.shape`

[illegible]

```
In [49]: fig, (ax1) = plt.subplots(1, figsize = (10,6))
ax1.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1], s = 50, c = 'red', label = 'Cluster 1')
ax1.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1], s = 50, c = 'cornflowerblue', label = 'Cluster 2')
ax1.scatter(x[y_kmeans == 2, 0], x[y_kmeans == 2, 1], s = 50, c = 'maroon', label = 'Cluster 3')
ax1.scatter(x[y_kmeans == 3, 0], x[y_kmeans == 3, 1], s = 50, c = 'greenyellow', label = 'Cluster 4')
ax1.scatter(x[y_kmeans == 4, 0], x[y_kmeans == 4, 1], s = 50, c = 'indigo', label = 'Cluster 5')
ax1.scatter(kmeans.cluster_centers[:, 0], kmeans.cluster_centers[:, 1], s = 200, c = 'yellow', label = 'Centroids')
plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```



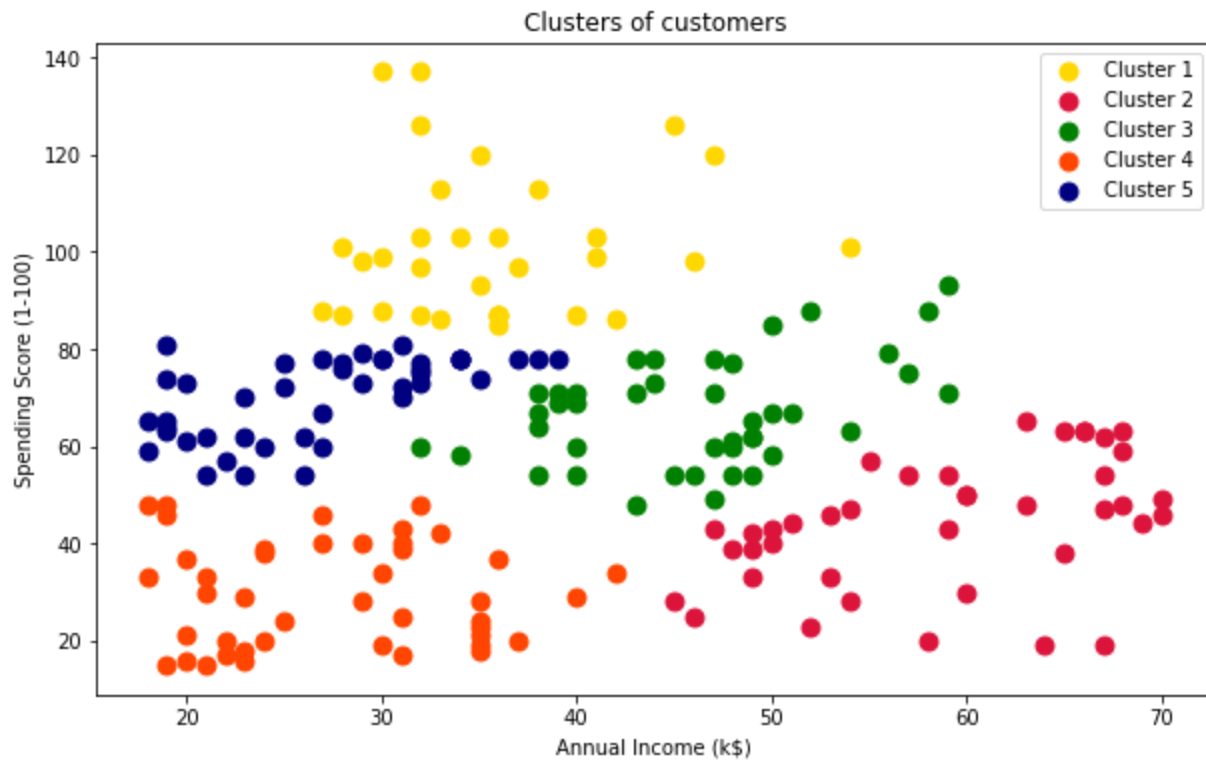
```
In [50]: import scipy.cluster.hierarchy as sch
dendrogram = sch.dendrogram(sch.linkage(x, method = 'ward'))
plt.title('Dendrogram')
plt.xlabel('Customers')
plt.ylabel('Euclidean distances')
plt.show()
```



```
In [51]: from sklearn.cluster import AgglomerativeClustering
```

```
hc = AgglomerativeClustering(n_clusters = 5, affinity = 'euclidean', linkage = 'ward')
y_hc = hc.fit_predict(x)
```

```
In [60]: fig, (ax1) = plt.subplots(1, figsize = (10,6))
ax1.scatter(x[y_hc == 0, 0], x[y_hc == 0, 1], s = 80, c = 'gold', label = 'Cluster 1')
ax1.scatter(x[y_hc == 1, 0], x[y_hc == 1, 1], s = 80, c = 'crimson', label = 'Cluster 2')
ax1.scatter(x[y_hc == 2, 0], x[y_hc == 2, 1], s = 80, c = 'green', label = 'Cluster 3')
ax1.scatter(x[y_hc == 3, 0], x[y_hc == 3, 1], s = 80, c = 'orangered', label = 'Cluster 4')
ax1.scatter(x[y_hc == 4, 0], x[y_hc == 4, 1], s = 80, c = 'navy', label = 'Cluster 5')
plt.title('Clusters of customers')
plt.xlabel('Annual Income (k$)')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()
```



```
In [ ]:
```