

▼ Import Libraries

```
from sklearn.preprocessing import StandardScaler
```

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
```

```
import os
import warnings
```

```
warnings.filterwarnings('ignore')
```

```
df = pd.read_csv('/content/Mall_Customers.csv')
df.head()
```

	CustomerID	Gender	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

Next steps: [Generate code with df](#) [View recommended plots](#)

```
df.info()
```

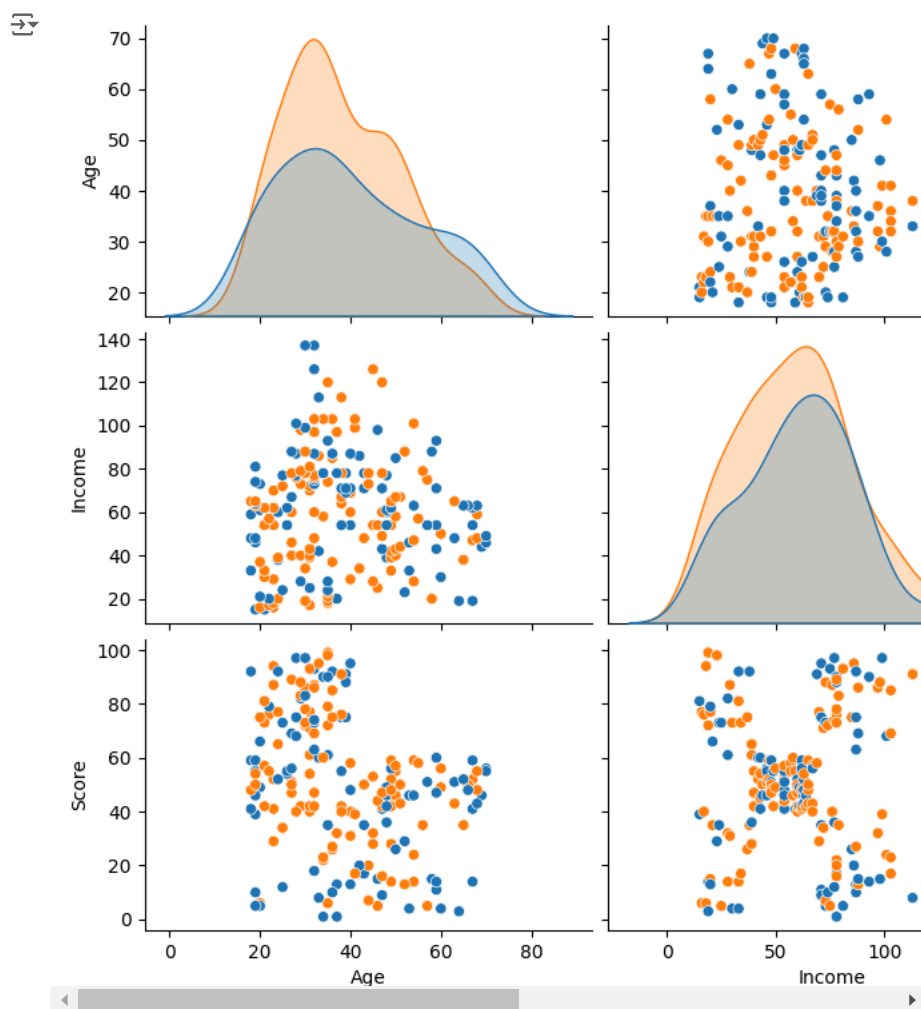
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column                Non-Null Count  Dtype
---  -
0   CustomerID            200 non-null   int64
1   Gender                200 non-null   object
2   Age                   200 non-null   int64
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
```

```
df.rename(index=str, columns={'Annual Income (k$)': 'Income',
                              'Spending Score (1-100)': 'Score'}, inplace=True)
df.head()
```

	CustomerID	Gender	Age	Income	Score
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

Next steps: [Generate code with df](#) [View recommended plots](#)

```
# Let's see our data in a detailed way with pairplot
X = df.drop(['CustomerID', 'Gender'], axis=1)
sns.pairplot(df.drop('CustomerID', axis=1), hue='Gender', aspect=1.5)
plt.show()
```



From the above plot we see that gender has no direct relation to segmenting customers. That's why we can drop it and move on with other features which is why we will X parameter from now on.

```
from sklearn.cluster import KMeans

clusters = []

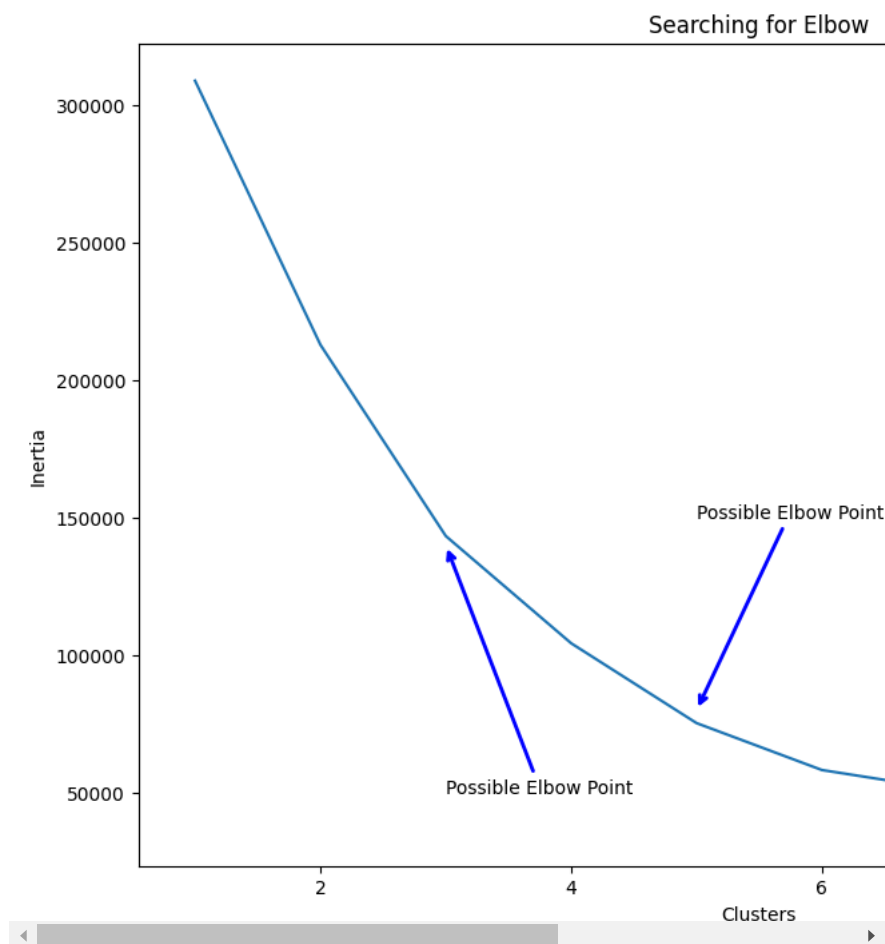
for i in range(1, 11):
    km = KMeans(n_clusters=i).fit(X)
    clusters.append(km.inertia_)

fig, ax = plt.subplots(figsize=(12, 8))
sns.lineplot(x=list(range(1, 11)), y=clusters, ax=ax)
ax.set_title('Searching for Elbow')
ax.set_xlabel('Clusters')
ax.set_ylabel('Inertia')

# Annotate arrow
ax.annotate('Possible Elbow Point', xy=(3, 140000), xytext=(3, 50000), xycoords='data',
           arrowprops=dict(arrowstyle='->', connectionstyle='arc3', color='blue', lw=2))

ax.annotate('Possible Elbow Point', xy=(5, 80000), xytext=(5, 150000), xycoords='data',
           arrowprops=dict(arrowstyle='->', connectionstyle='arc3', color='blue', lw=2))

plt.show()
```



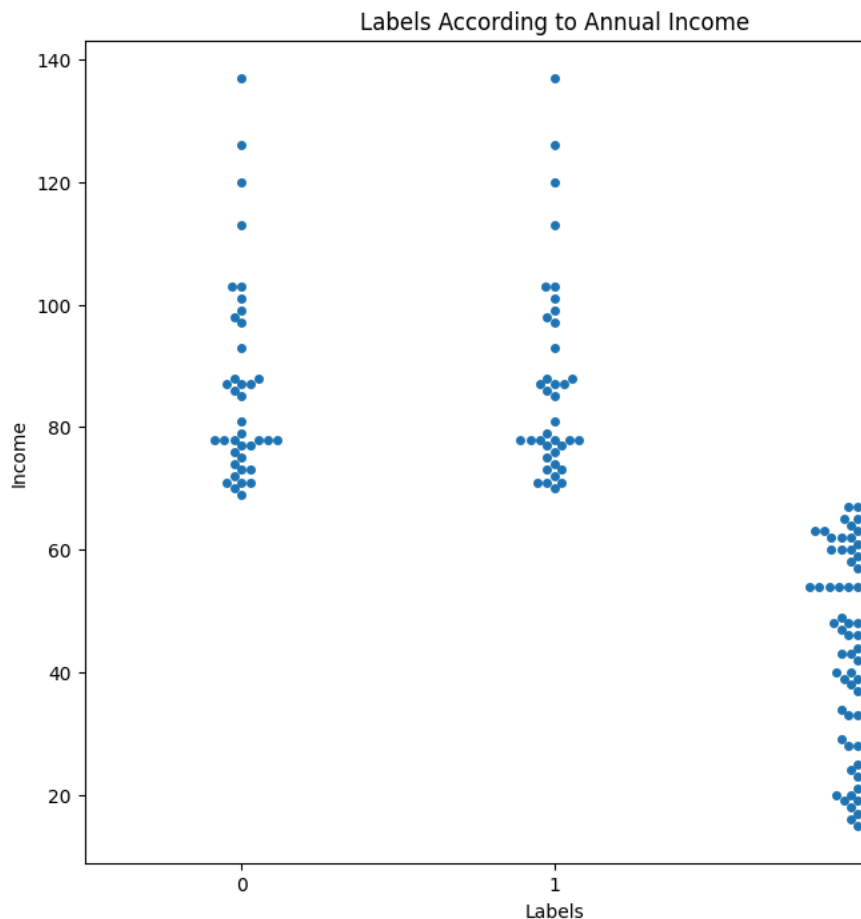
Elbow method tells us to select the cluster when there is a significant change in inertia. As we can see from the graph, we can say this may be either 3 or 5. Let's see both results in graph and decide.

Creating the Visual Plots

```
fig = plt.figure(figsize=(20,8))
ax = fig.add_subplot(121)
sns.swarmplot(x='Labels', y='Income', data=X, ax=ax)
ax.set_title('Labels According to Annual Income')

ax = fig.add_subplot(122)
sns.swarmplot(x='Labels', y='Score', data=X, ax=ax)
ax.set_title('Labels According to Scoring History')

plt.show()
```



✓ Hierarchical Clustering

Agglomerative We will be looking at a clustering technique, which is Agglomerative Hierarchical Clustering. Agglomerative is the bottom up approach which is more popular than Divisive clustering.

We will also be using Complete Linkage as the Linkage Criteria.

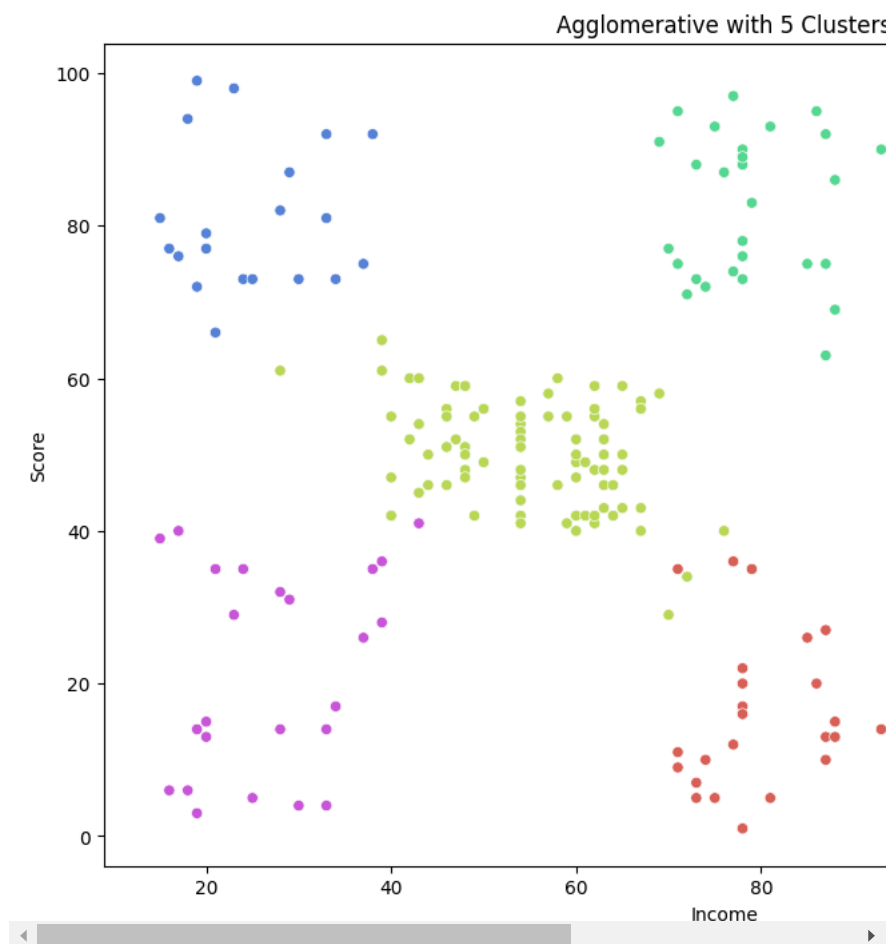
The Agglomerative Clustering class will require two inputs:

`n_clusters`: The number of clusters to form as well as the number of centroids to generate. `linkage`: Which linkage criterion to use. The linkage criterion determines which distance to use between sets of observation. The algorithm will merge the pairs of cluster that minimize this criterion. Value will be: 'complete' Note: It is recommended that try everything with 'average' as well

```
import seaborn as sns
```

```
# ...
```

```
plt.figure(figsize=(12, 8))
sns.scatterplot(x="Income", y="Score", hue="Labels", data=X,
               palette=sns.color_palette("hls", 5))
plt.title("Agglomerative with 5 Clusters")
plt.show()
```



Dendrogram Associated for the Agglomerative Hierarchical Clustering Remember that a distance matrix contains the distance from each point to every other point of a dataset . We can use the function `distance_matrix`, which requires two inputs. Remember that the distance values are symmetric, with a diagonal of 0's. This is one way of making sure your matrix is correct.

```
from scipy.cluster import hierarchy
from scipy.spatial import distance_matrix
```

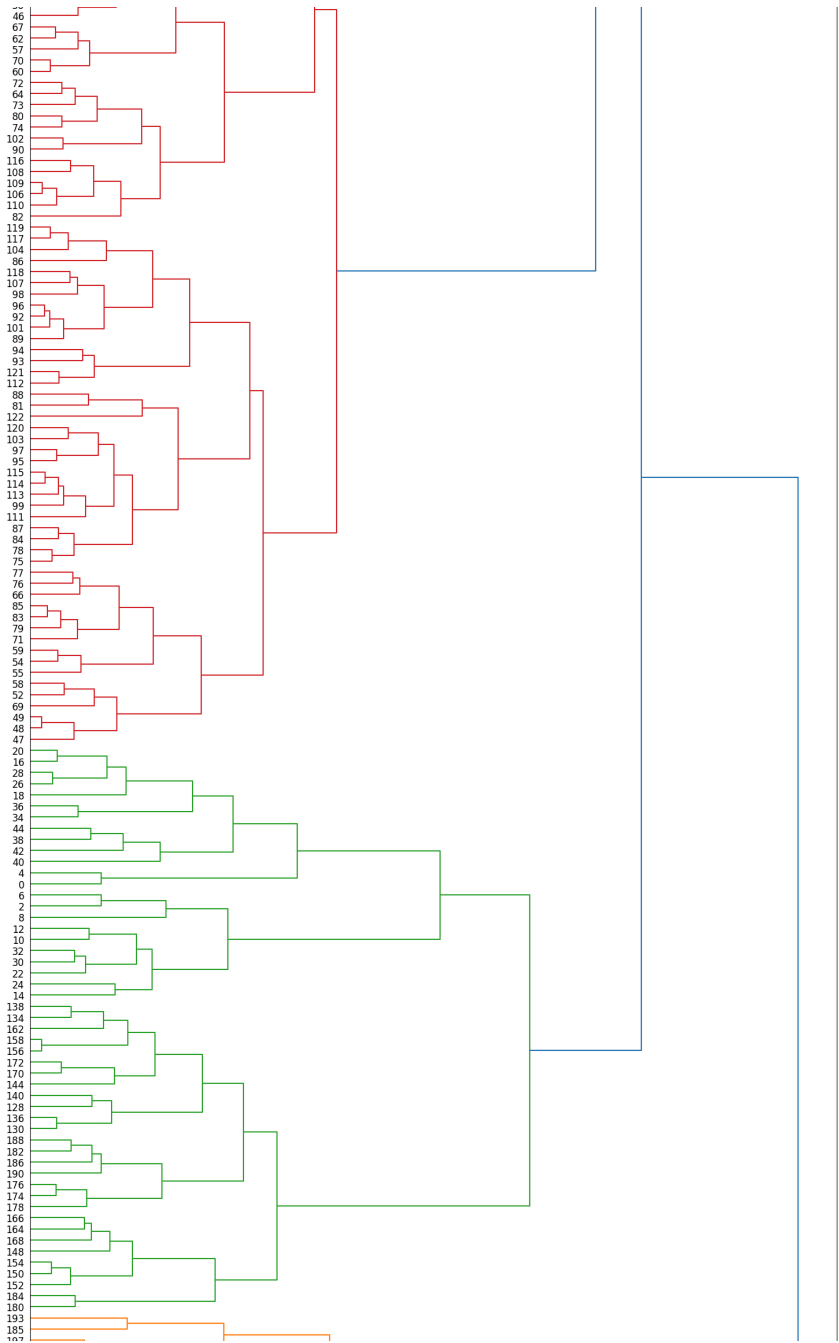
```
dist = distance_matrix(X, X)
print(dist)
```

```
[[ 0.          42.05948169  33.03028913 ... 117.12813496 124.53915047
 130.17296186]
 [ 42.05948169  0.          75.01999733 ... 111.76761606 137.77880824
 122.35195135]
 [ 33.03028913  75.01999733  0.          ... 129.89226305 122.24974438
 143.78456106]
 ...
 [117.12813496 111.76761606 129.89226305 ... 0.          57.10516614
 14.35270009]
 [124.53915047 137.77880824 122.24974438 ... 57.10516614 0.
 65.06150936]
 [130.17296186 122.35195135 143.78456106 ... 14.35270009 65.06150936
 0.          ]]
```

```
Z = hierarchy.linkage(dist, 'complete')
```

A Hierarchical clustering is typically visualized as a dendrogram as shown in the following cell. Each merge is represented by a horizontal line. The y-coordinate of the horizontal line is the similarity of the two clusters that were merged, where cities are viewed as singleton clusters. By moving up from the bottom layer to the top node, a dendrogram allows us to reconstruct the history of merges that resulted in the depicted clustering.

```
plt.figure(figsize=(18, 50))
dendro = hierarchy.dendrogram(Z, leaf_rotation=0, leaf_font_size=12, orientation='right')
```



We used complete linkage for our case, let's change it to average linkage to see how the dendrogram changes

```
Z = hierarchy.linkage(dist, 'average')
plt.figure(figsize=(18, 50))
dendro = hierarchy.dendrogram(Z, leaf_rotation=0, leaf_font_size =12, orientation = 'right')
```