

```
In [15]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [16]: data = pd.read_csv('Mall_Customers.csv')
```

```
In [11]: data.head()
```

Out[11]:

	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

```
In [17]: data.shape
```

Out[17]: (200, 5)

```
In [18]: data.tail()
```

Out[18]:

	CustomerID	Gender	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 [19]: data.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
```

```
In [20]: data.describe()
```

Out[20]:

	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

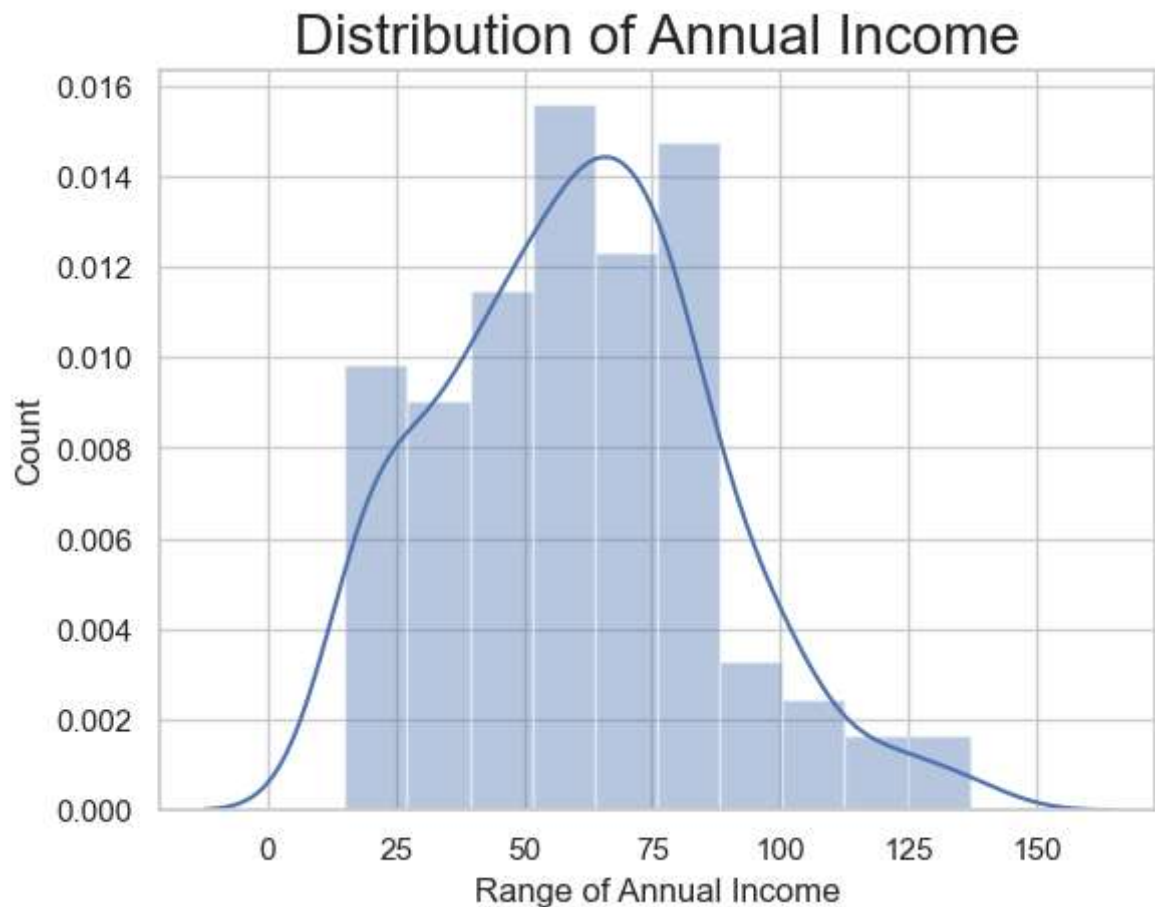
```
In [21]: data.isnull().any()
```

```
Out[21]: CustomerID           False
Gender             False
Age                False
Annual Income (k$) False
Spending Score (1-100) False
dtype: bool
```

```
In [22]: sns.set(style = 'whitegrid')
sns.distplot(data['Annual Income (k$)'])
plt.title('Distribution of Annual Income', fontsize = 20)
plt.xlabel('Range of Annual Income')
plt.ylabel('Count')
plt.show()
```

C:\Users\LENOVO\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

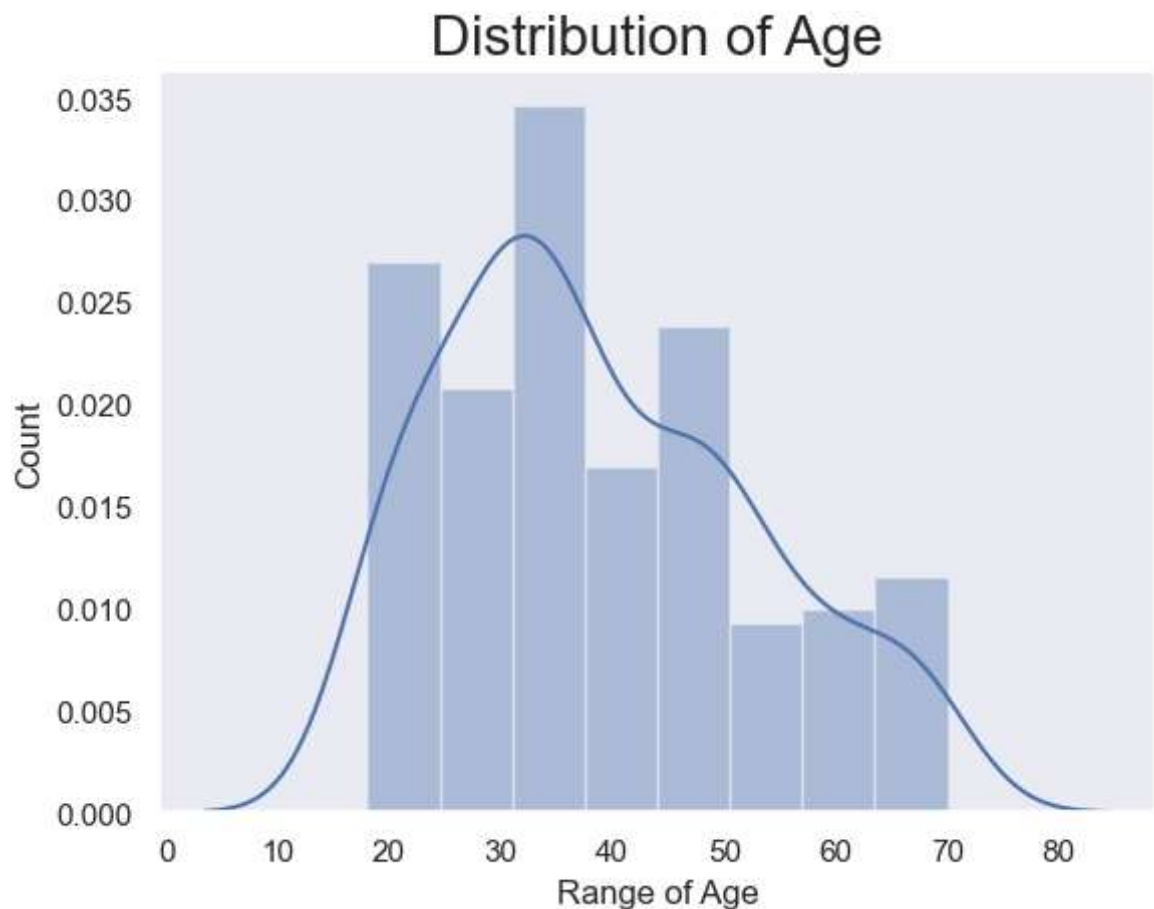
warnings.warn(msg, FutureWarning)



```
In [23]: sns.set(style = 'dark')
sns.distplot(data['Age'])
plt.title('Distribution of Age', fontsize = 20)
plt.xlabel('Range of Age')
plt.ylabel('Count')
plt.show()
```

C:\Users\LENOVO\anaconda3\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: `distplot` is a deprecated function and will be removed in a future version. Please adapt your code to use either `displot` (a figure-level function with similar flexibility) or `histplot` (an axes-level function for histograms).

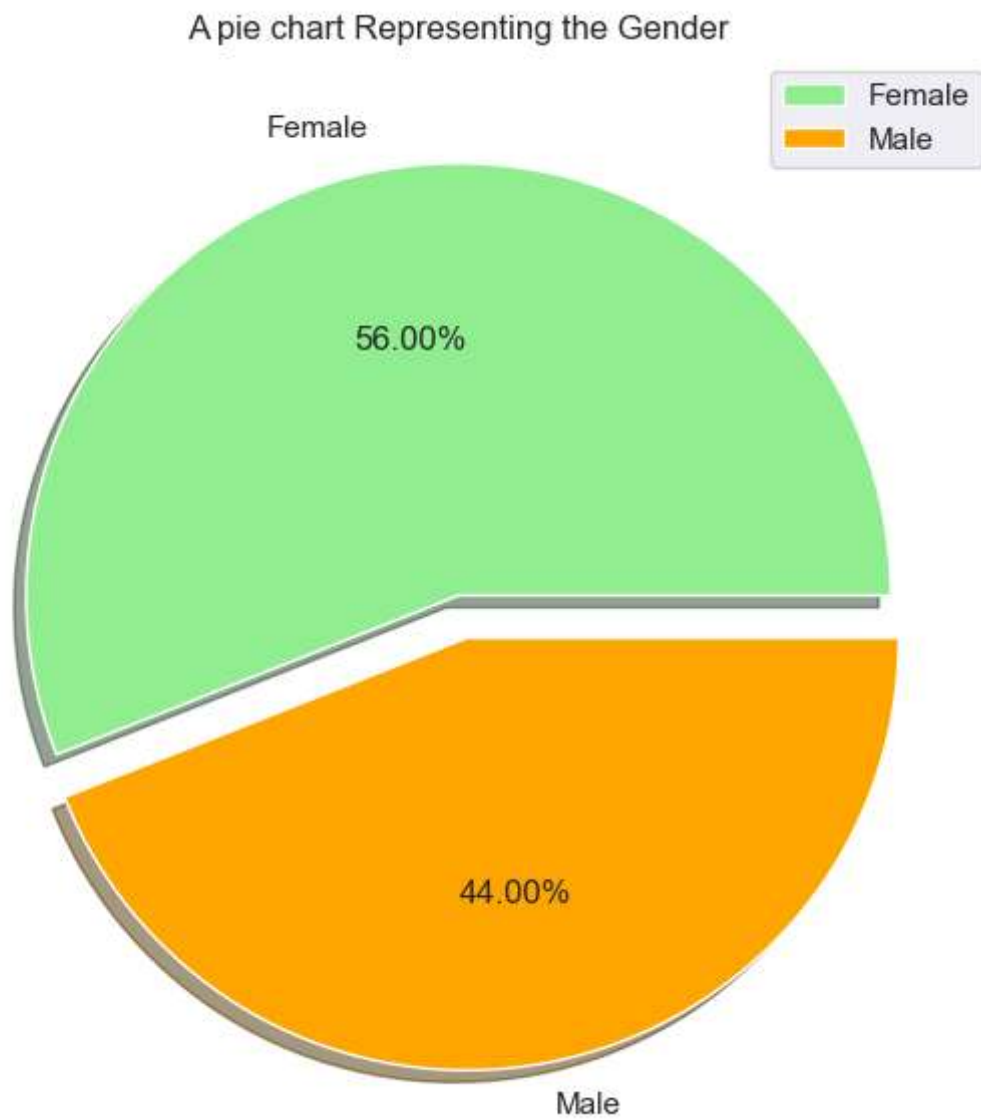
warnings.warn(msg, FutureWarning)



```
In [24]: data['Gender'].value_counts()
```

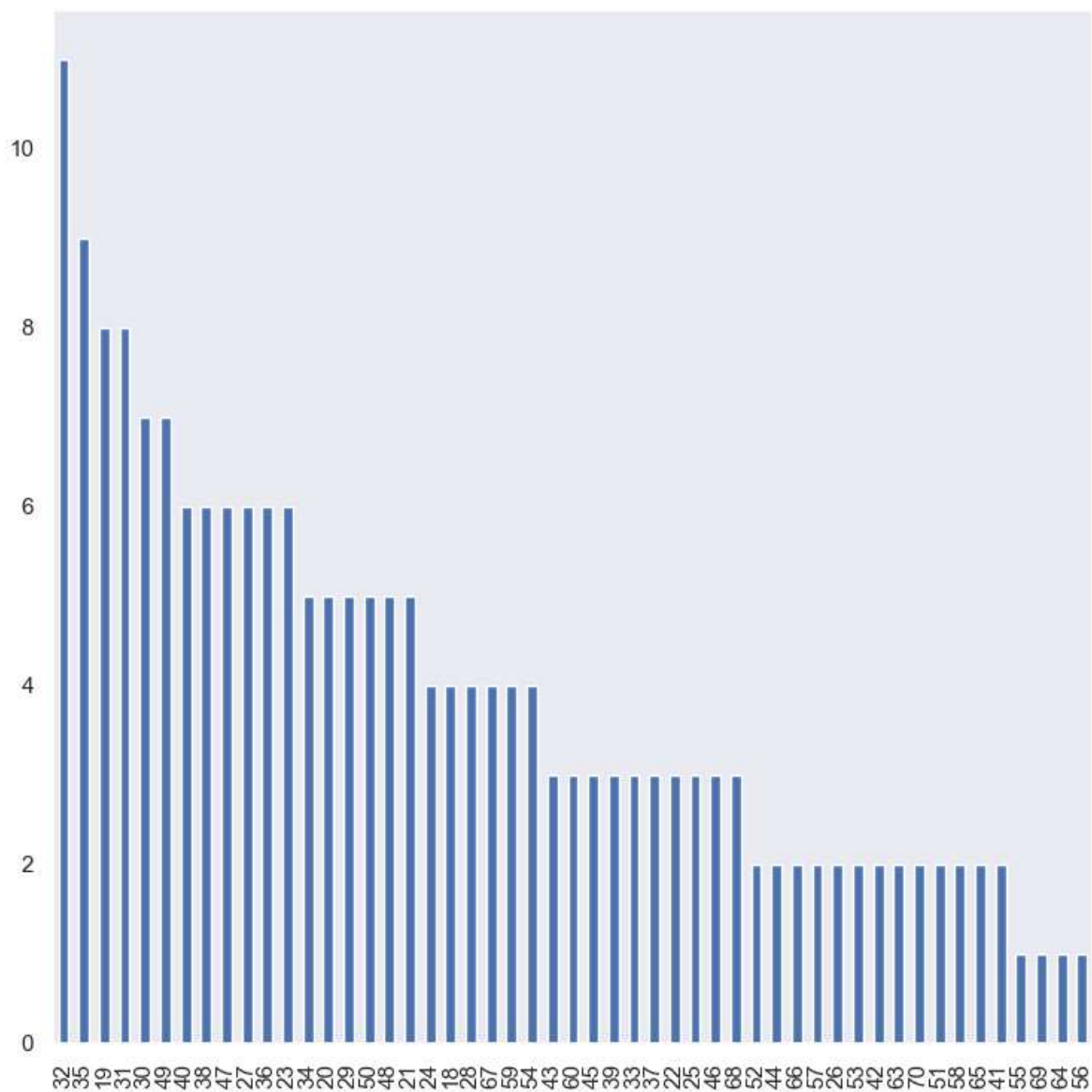
```
Out[24]: Female    112
Male           88
Name: Gender, dtype: int64
```

```
In [25]: labels = ['Female', 'Male']  
size = [112, 88]  
colors = ['lightgreen', 'orange']  
explode = [0, 0.1]  
  
plt.rcParams['figure.figsize'] = (7, 7)  
plt.pie(size, colors = colors, explode = explode, labels = labels, shadow = True)  
plt.title('A pie chart Representing the Gender')  
plt.axis('off')  
plt.legend()  
plt.show()
```



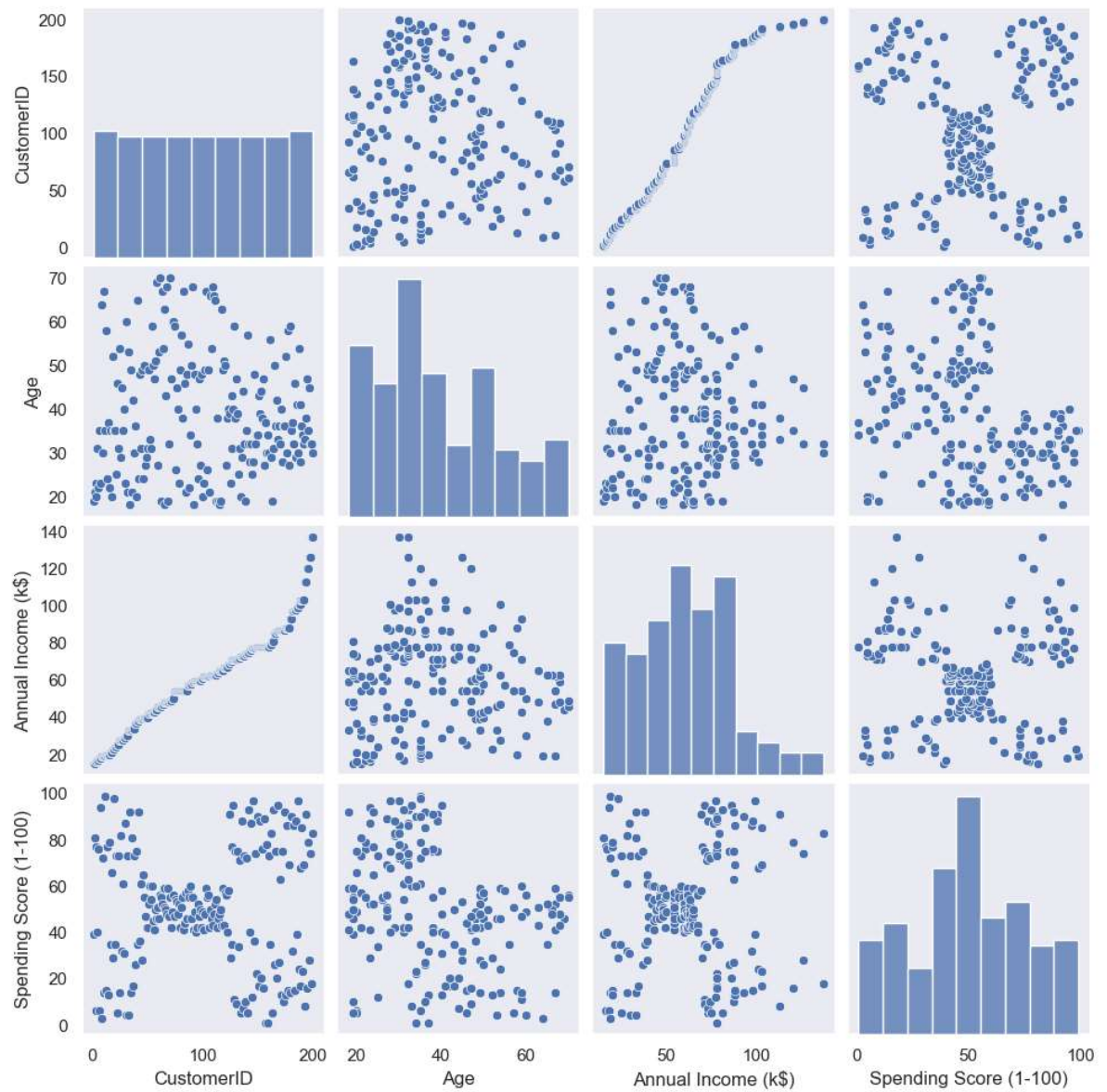
```
In [26]: data['Age'].value_counts().plot.bar(figsize = (9, 9))
```

```
Out[26]: <AxesSubplot:>
```



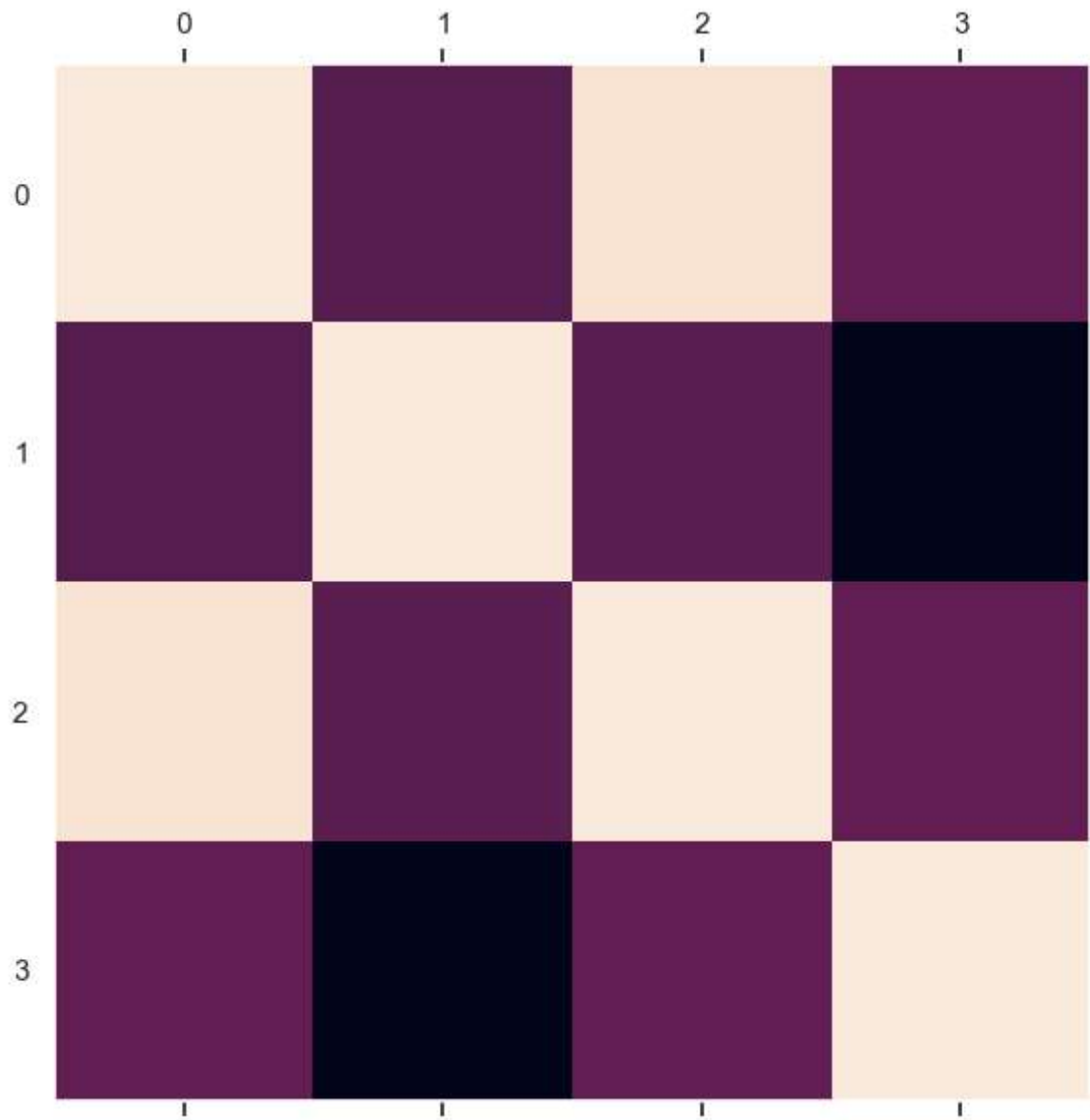
```
In [27]: sns.pairplot(data)
```

```
Out[27]: <seaborn.axisgrid.PairGrid at 0x14a0038cd90>
```



```
In [28]: plt.matshow(data.corr())
```

```
Out[28]: <matplotlib.image.AxesImage at 0x14a00663be0>
```



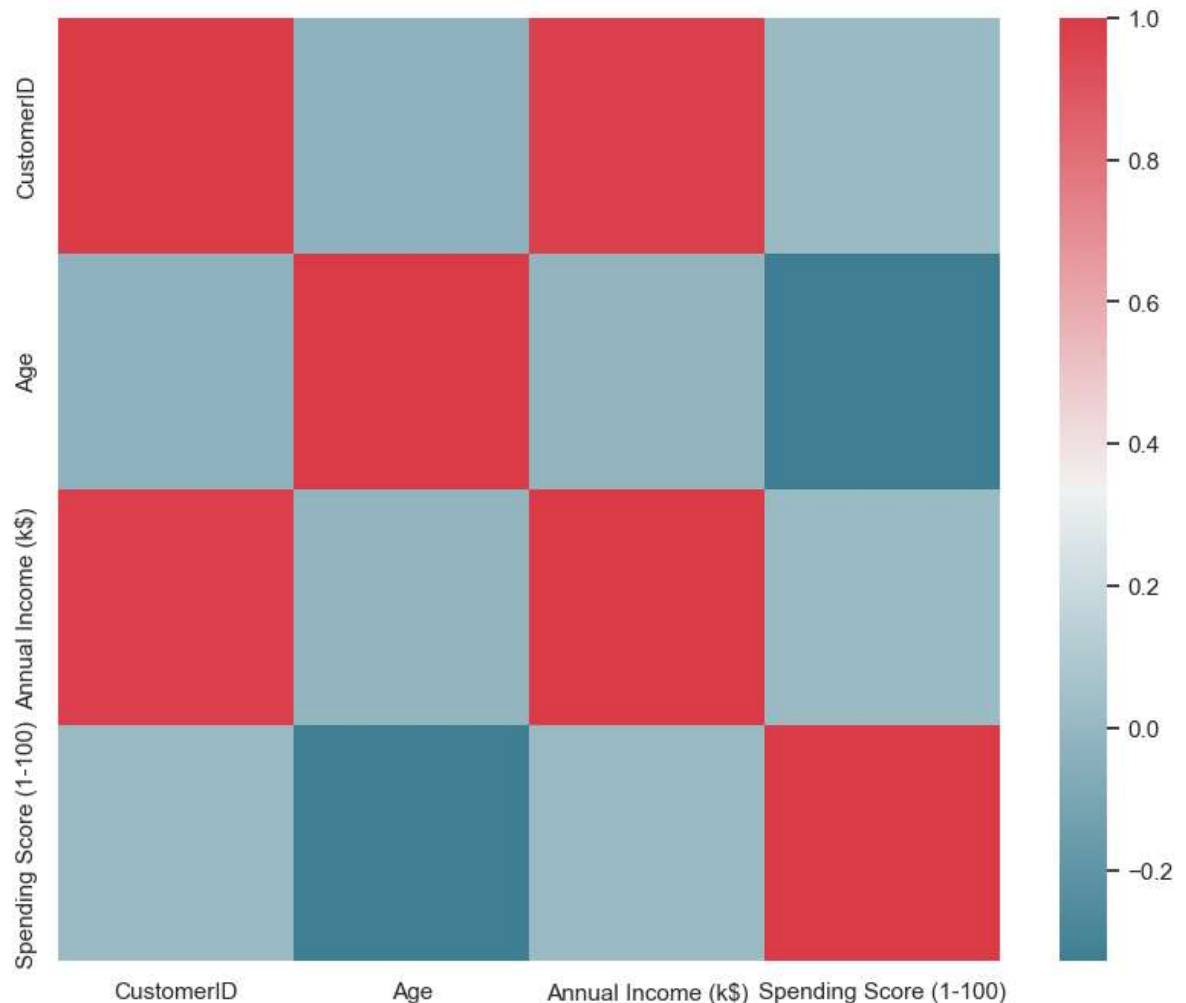

```
In [30]: fig, axis = plt.subplots(figsize=(10, 8))
corr = data.corr()
sns.heatmap(corr, mask = np.zeros_like(corr, dtype = np.bool), cmap = sns.diverging_palette(220, 10, as_cmap = True), ax = axis)
```

C:\Users\LENOVO\AppData\Local\Temp\ipykernel_6436\408083941.py:3: Deprecation Warning: `np.bool` is a deprecated alias for the builtin `bool`. To silence this warning, use `bool` by itself. Doing this will not modify any behavior and is safe. If you specifically wanted the numpy scalar type, use `np.bool_` here.

Deprecated in NumPy 1.20; for more details and guidance: <https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations> (<https://numpy.org/devdocs/release/1.20.0-notes.html#deprecations>)

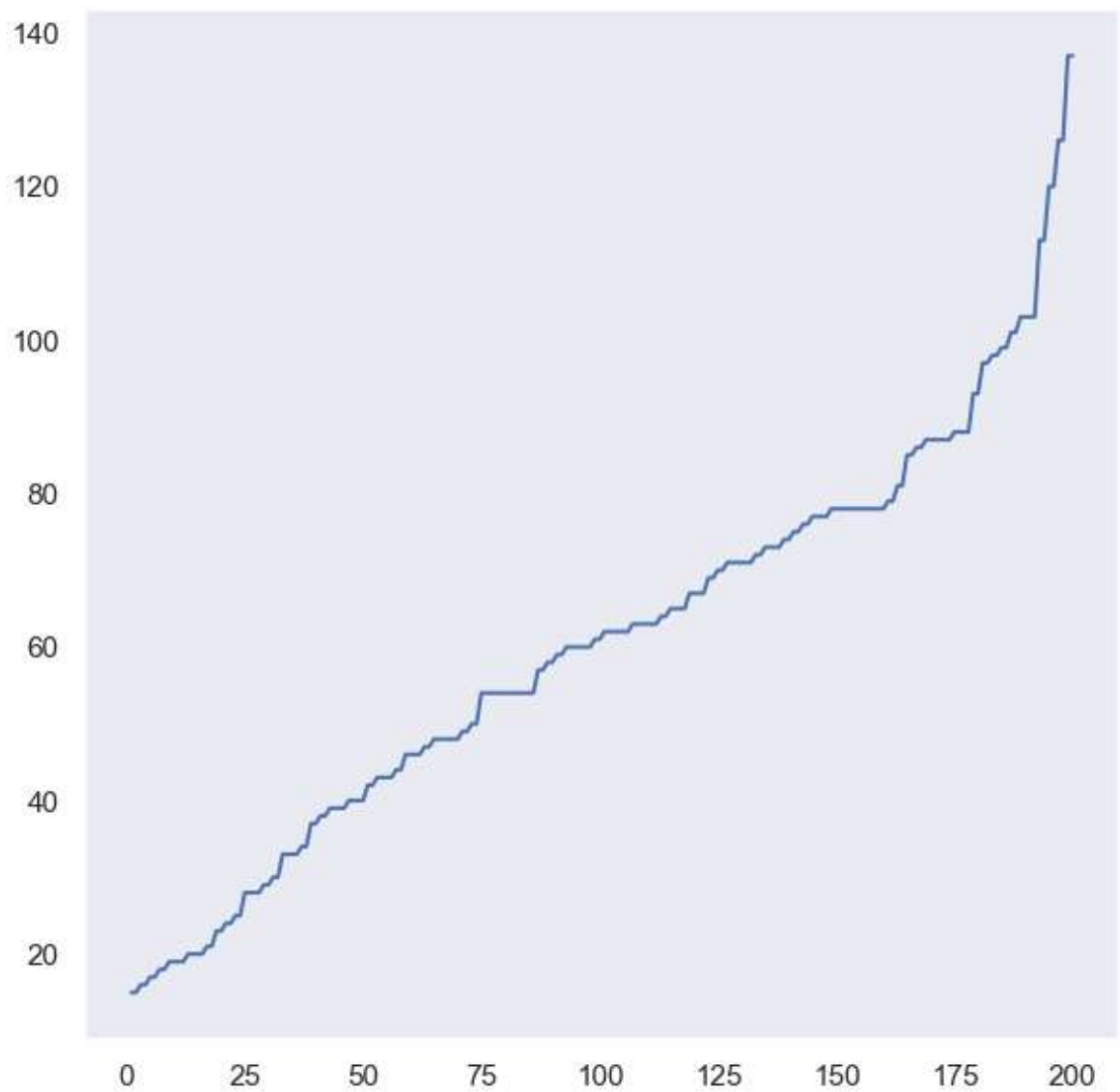
```
sns.heatmap(corr, mask = np.zeros_like(corr, dtype = np.bool), cmap = sns.diverging_palette(220, 10, as_cmap = True),
```

Out[30]: <AxesSubplot:>



```
In [31]: x = data['CustomerID']  
y = data['Annual Income (k$)']  
  
plt.plot(x, y)
```

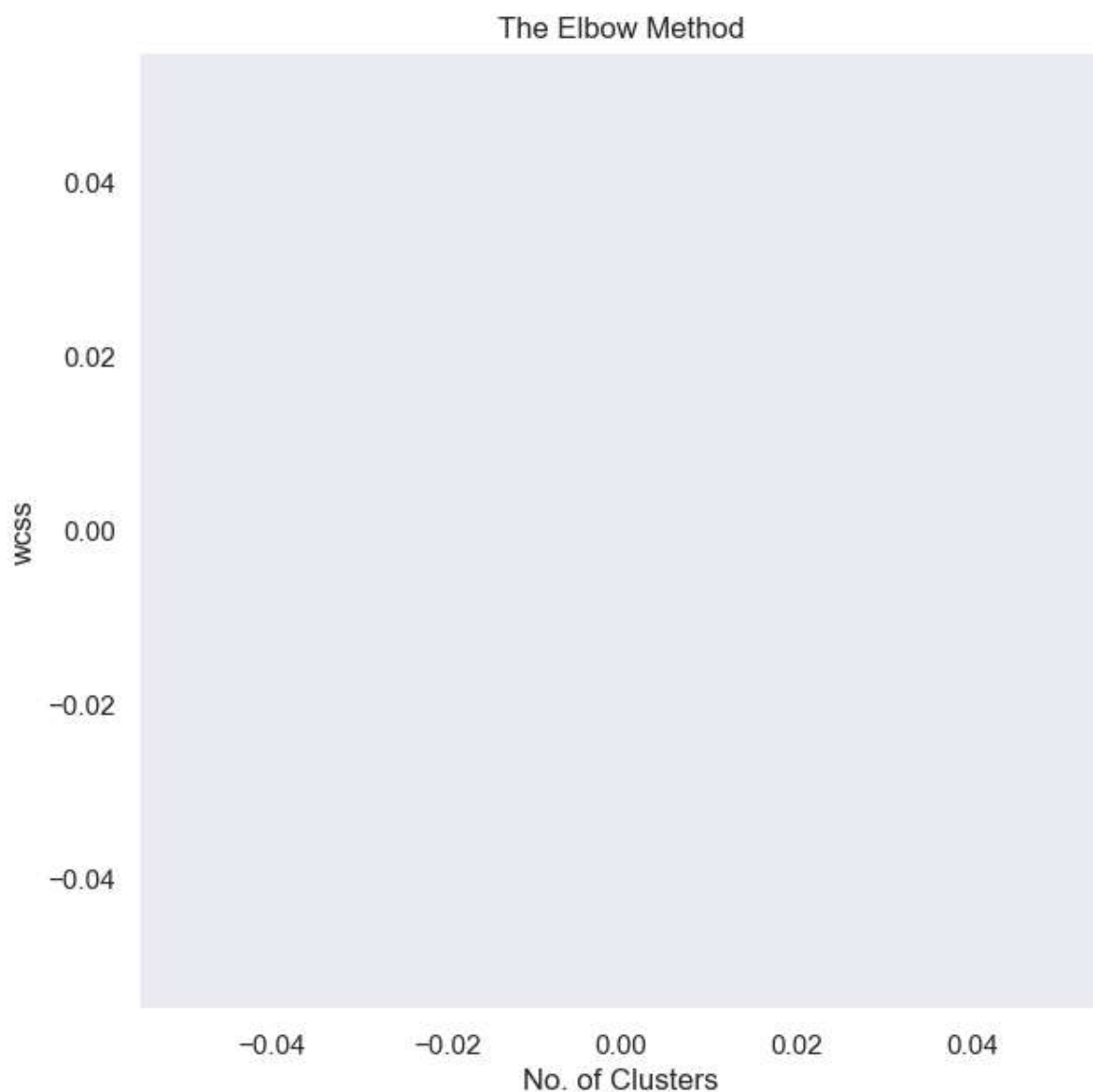
Out[31]: [<matplotlib.lines.Line2D at 0x14a023b61f0>]



```
In [38]: from sklearn.cluster import KMeans

wcss = []
for i in range(1,11):
    km = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10,
    km.fit(x)
    wcss.append(km.inertia_)

plt.plot(range(1,11), wcss)
plt.title('The Elbow Method')
plt.xlabel('No. of Clusters')
plt.ylabel('wcss')
plt.show()
```



```
In [50]: data.columns
```

```
Out[50]: Index(['CustomerID', 'Gender', 'Age', 'Annual Income (k$)',
                'Spending Score (1-100)'],
                dtype='object')
```

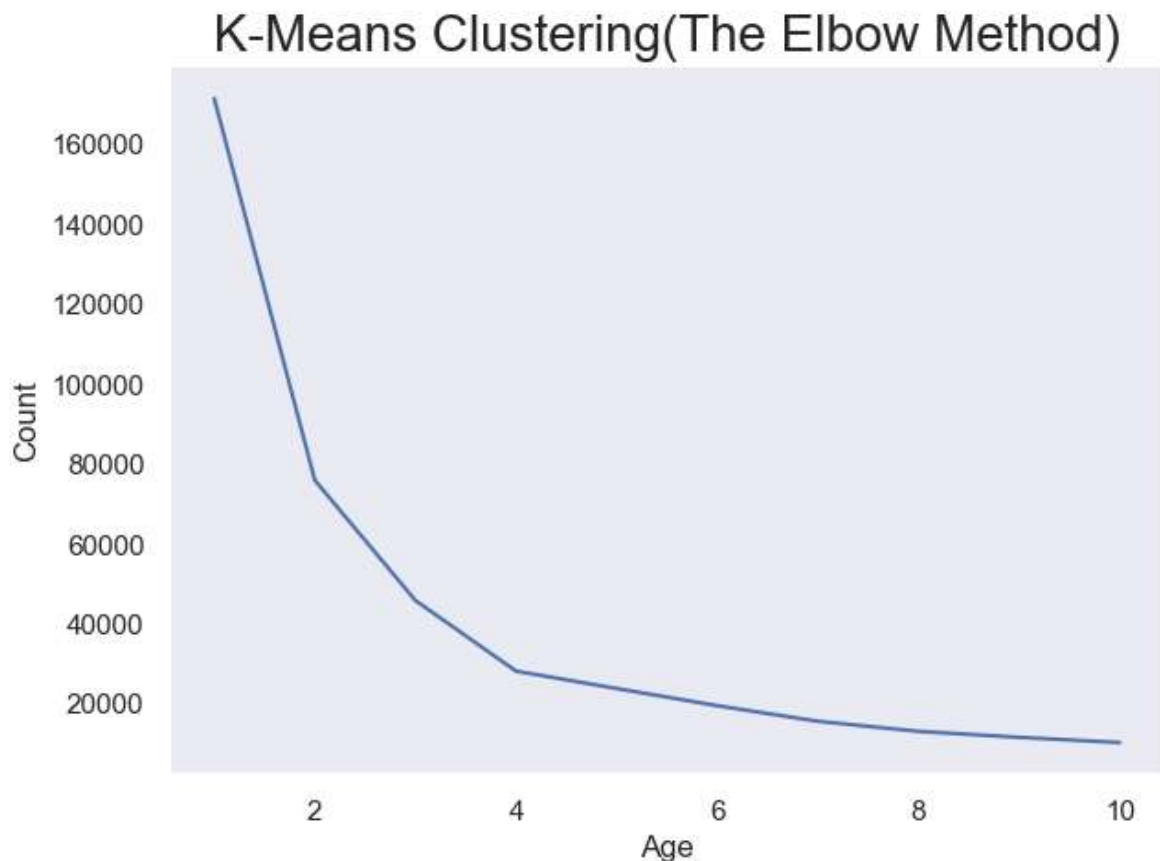
```
In [51]: x = data.iloc[:, [2, 4]].values  
x.shape
```

```
Out[51]: (200, 2)
```

```
In [52]: from sklearn.cluster import KMeans  
  
wcss = []  
for i in range(1, 11):  
    kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10,  
                    random_state = 0)  
    kmeans.fit(x)  
    wcss.append(kmeans.inertia_)  
  
plt.rcParams['figure.figsize'] = (7, 5)  
plt.plot(range(1, 11), wcss)  
plt.title('K-Means Clustering(The Elbow Method)', fontsize = 20)  
plt.xlabel('Age')  
plt.ylabel('Count')  
plt.show()
```

C:\Users\LENOVO\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036:
UserWarning: KMeans is known to have a memory leak on Windows with MKL, when
there are less chunks than available threads. You can avoid it by setting the
environment variable OMP_NUM_THREADS=1.

warnings.warn(



```

In [53]: kmeans = KMeans(n_clusters = 4, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
ymeans = kmeans.fit_predict(x)

plt.rcParams['figure.figsize'] = (10, 10)
plt.title('Cluster of Ages', fontsize = 30)

plt.scatter(x[ymean == 0, 0], x[ymean == 0, 1], s = 100, c = 'pink', label = 'Usual Customers')
plt.scatter(x[ymean == 1, 0], x[ymean == 1, 1], s = 100, c = 'orange', label = 'Priority Customers')
plt.scatter(x[ymean == 2, 0], x[ymean == 2, 1], s = 100, c = 'lightgreen', label = 'Target Customers (Young)')
plt.scatter(x[ymean == 3, 0], x[ymean == 3, 1], s = 100, c = 'red', label = 'Target Customers (Old)')
plt.scatter(kmeans.cluster_centers_[0, 0], kmeans.cluster_centers_[0, 1], s = 100, c = 'black', label = 'Cluster Center 0')
plt.scatter(kmeans.cluster_centers_[1, 0], kmeans.cluster_centers_[1, 1], s = 100, c = 'black', label = 'Cluster Center 1')
plt.scatter(kmeans.cluster_centers_[2, 0], kmeans.cluster_centers_[2, 1], s = 100, c = 'black', label = 'Cluster Center 2')
plt.scatter(kmeans.cluster_centers_[3, 0], kmeans.cluster_centers_[3, 1], s = 100, c = 'black', label = 'Cluster Center 3')

plt.xlabel('Age')
plt.ylabel('Spending Score (1-100)')
plt.legend()
plt.show()

```



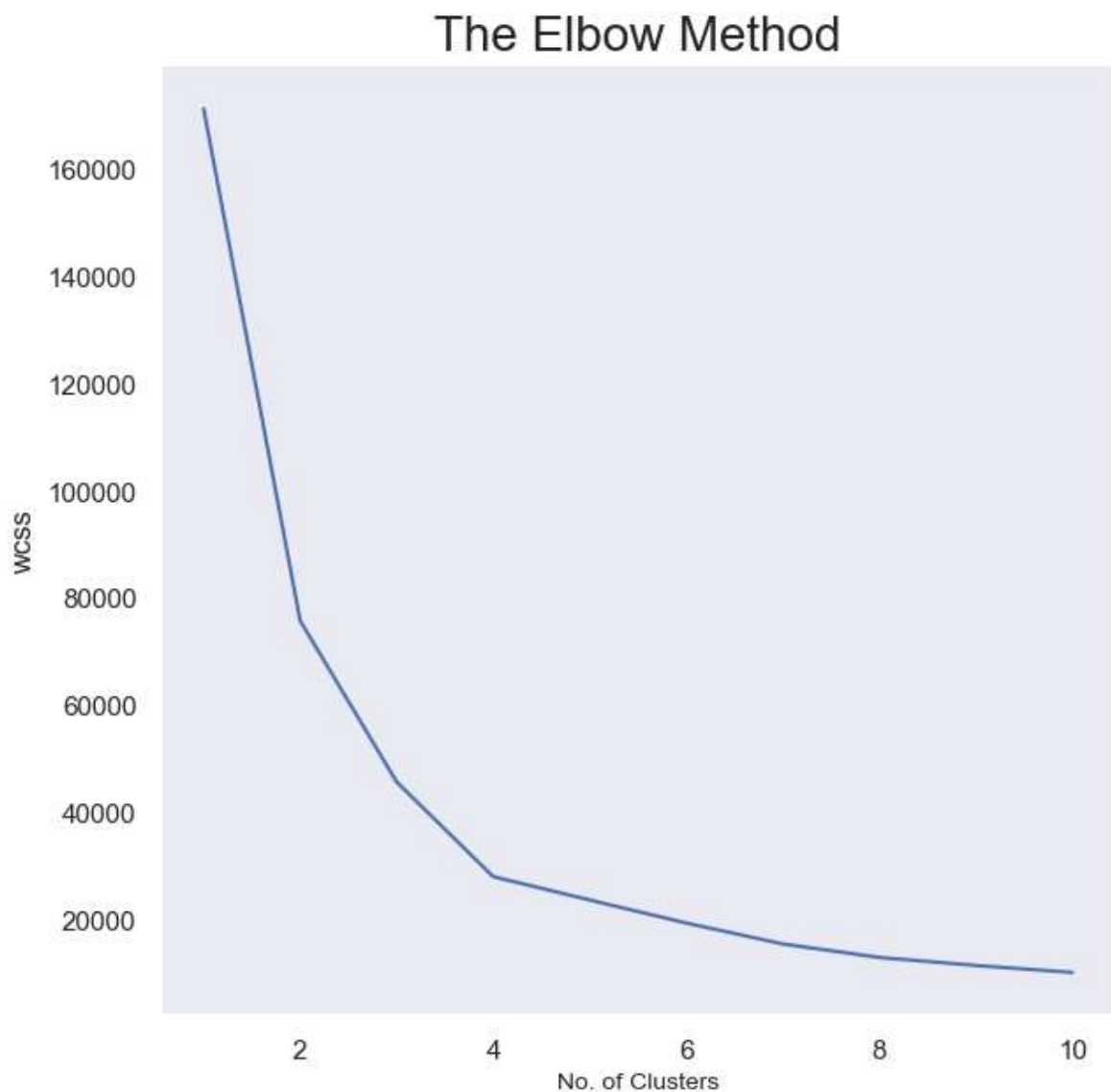
```
In [54]: from sklearn.cluster import KMeans

wcss = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters = i, init = 'k-means++', max_iter = 300, n_init = 10)
    kmeans.fit(x)
    wcss.append(kmeans.inertia_)

plt.rcParams['figure.figsize'] = (7, 7)
plt.title('The Elbow Method', fontsize = 20)
plt.plot(range(1, 11), wcss)
plt.xlabel('No. of Clusters', fontsize = 10)
plt.ylabel('wcss')
plt.show()
```

C:\Users\LENOVO\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036:
UserWarning: KMeans is known to have a memory leak on Windows with MKL, when
there are less chunks than available threads. You can avoid it by setting the
environment variable OMP_NUM_THREADS=1.

warnings.warn(

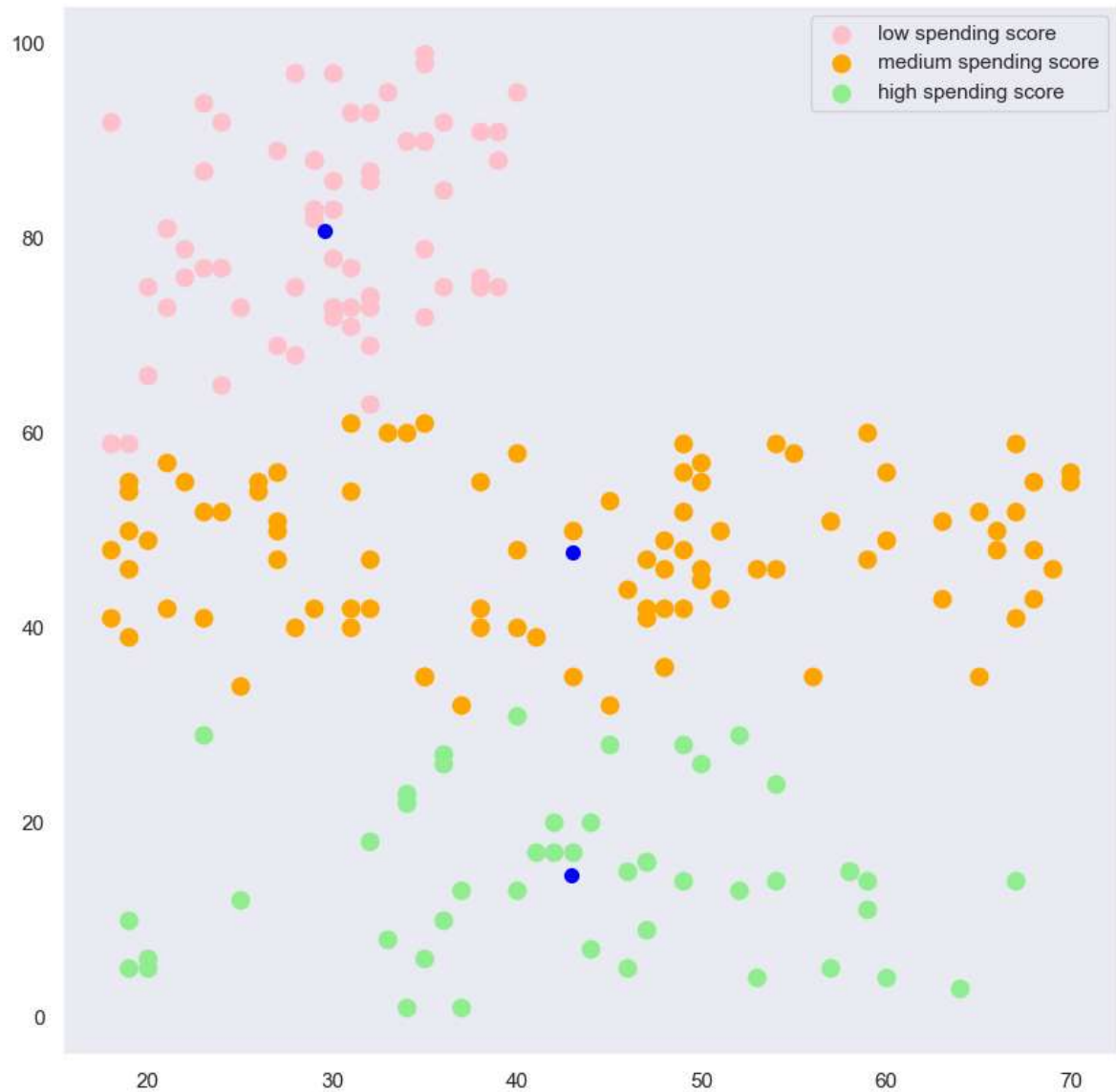


```

In [55]: kmeans = KMeans(n_clusters = 3, max_iter = 300, n_init = 10, random_state = 0)
ymeans = kmeans.fit_predict(x)

plt.rcParams['figure.figsize'] = (10, 10)
plt.scatter(x[ymmeans == 0, 0], x[ymmeans == 0, 1], s = 80, c = 'pink', label =
plt.scatter(x[ymmeans == 1, 0], x[ymmeans == 1, 1], s = 80, c = 'orange', label
plt.scatter(x[ymmeans == 2, 0], x[ymmeans == 2, 1], s = 80, c = 'lightgreen', la
plt.scatter(kmeans.cluster_centers_[0,0], kmeans.cluster_centers_[0, 1], s = 5
plt.legend()
plt.show()

```



In []:

