

In [5]:

```
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns
```

In [6]:

```
data = pd.read_csv('Mall_Customers.csv')
print(data.shape)
```

(200, 5)

In [7]:

```
data = data.drop('CustomerID',axis=1)
data
```

Out[7]:

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

200 rows × 4 columns

In [8]:

```
data.head()
```

Out[8]:

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

In [9]:

```
data.describe()
```

Out[9]:

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

In [66]:

```

from sklearn.preprocessing import LabelEncoder

# label_encoder object knows how to understand word labels.
label_encoder = LabelEncoder()

# Encode labels in column 'species'.
data['Gender'] = label_encoder.fit_transform(data['Gender'])
data

```

Out[66]:

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

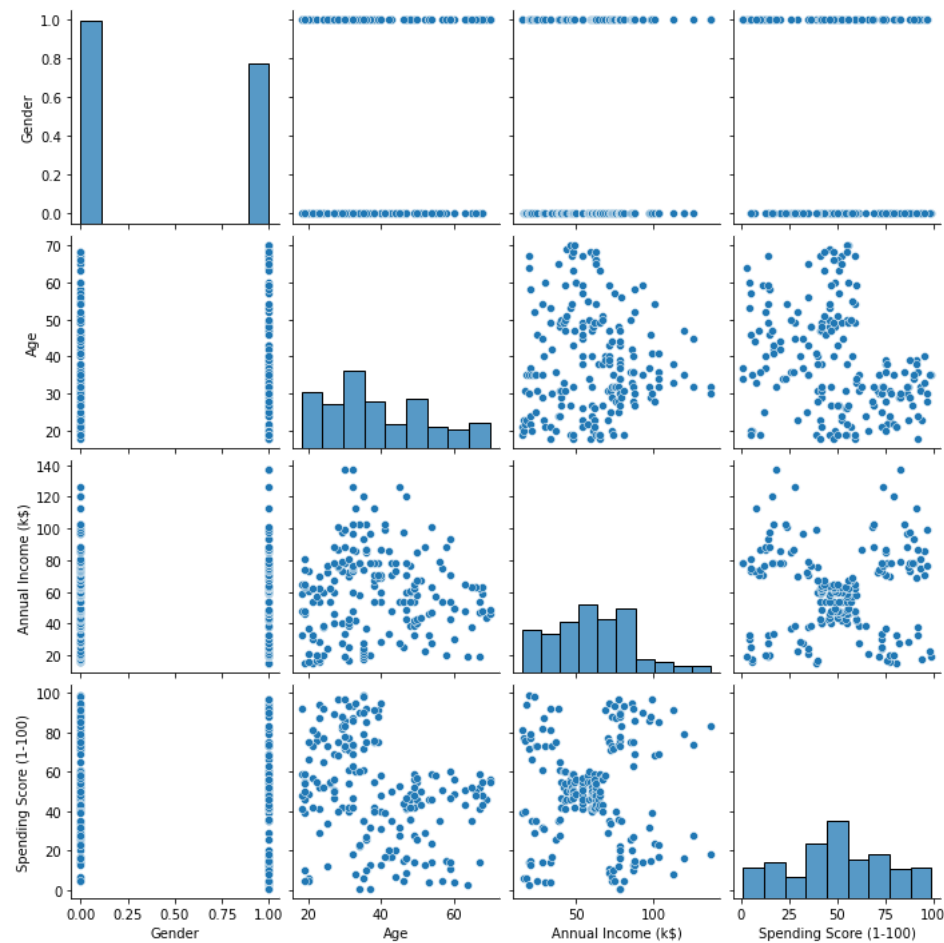
200 rows × 5 columns

In [68]:

```
sns.pairplot(data)
```

Out[68]:

<seaborn.axisgrid.PairGrid at 0x1cbea560970>



In [69]:

```

x = data.iloc[:, 2:4].values
print(x.shape)

```

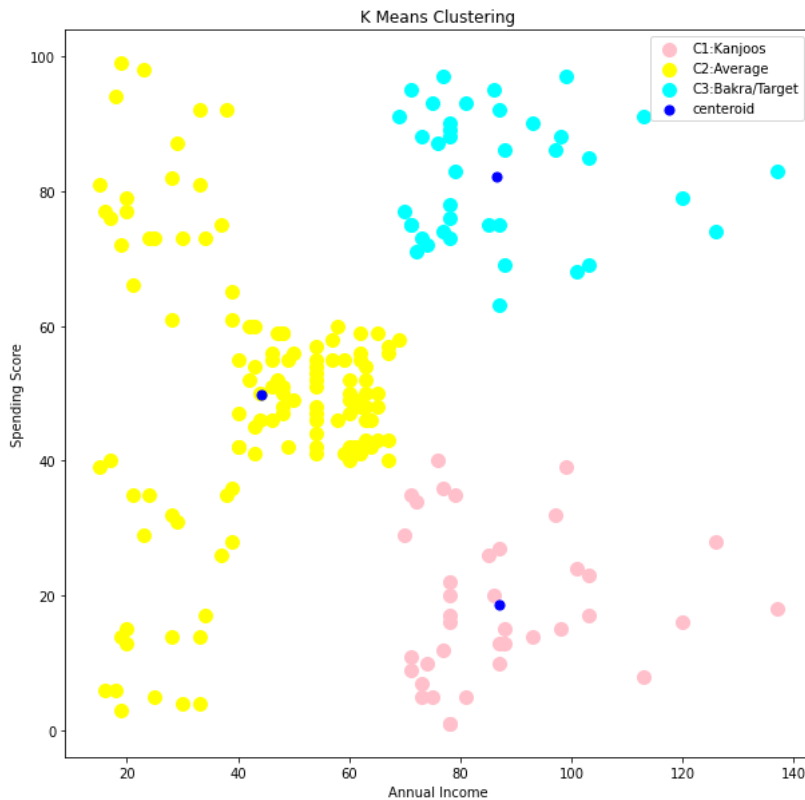
(200, 2)

In [74]:

```
plt.scatter(x[y_means == 0, 0], x[y_means == 0, 1], s = 100, c = 'pink', label = 'C1:Kanjooos')
plt.scatter(x[y_means == 1, 0], x[y_means == 1, 1], s = 100, c = 'yellow', label = 'C2:Average')
plt.scatter(x[y_means == 2, 0], x[y_means == 2, 1], s = 100, c = 'cyan', label = 'C3:Bakra/Target')
plt.scatter(x[y_means == 3, 0], x[y_means == 3, 1], s = 100, c = 'magenta', label = 'C4:Pokiri')
plt.scatter(x[y_means == 4, 0], x[y_means == 4, 1], s = 100, c = 'orange', label = 'C5:Intelligent')

plt.scatter(km1.cluster_centers[:,0], km1.cluster_centers[:, 1], s = 50, c = 'blue' , label = 'centroid')

plt.title('K Means Clustering')
plt.xlabel('Annual Income')
plt.ylabel('Spending Score')
plt.legend()
plt.show()
```



In [22]:

```
x[y_means == 1, 0]
```

Out[22]:

```
array([39, 40, 40, 40, 40, 42, 42, 43, 43, 43, 43, 44, 44, 46, 46, 46, 46,
       47, 47, 48, 48, 48, 48, 48, 48, 49, 49, 50, 50, 54, 54, 54, 54, 54,
       54, 54, 54, 54, 54, 54, 54, 57, 57, 58, 58, 59, 59, 60, 60, 60, 60,
       60, 60, 61, 61, 62, 62, 62, 62, 62, 62, 63, 63, 63, 63, 63, 63, 64,
       64, 65, 65, 65, 65, 67, 67, 67, 67, 69, 71, 72, 76], dtype=int64)
```

In []:

```
km1.inertia_
```

In []:

```
km1.cluster_centers_
```

Clusters of Customers Based on their Ages

In [38]:

```
data.columns
```

Out[38]:

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

In [39]:

```
x = data.iloc[:, [1, 3]].values  
x.shape
```

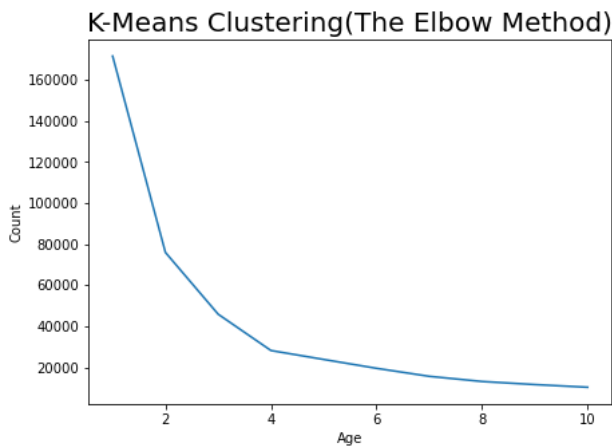
Out[39]:

(200, 2)

In [40]:

```
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\YASH\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 [41]:

```

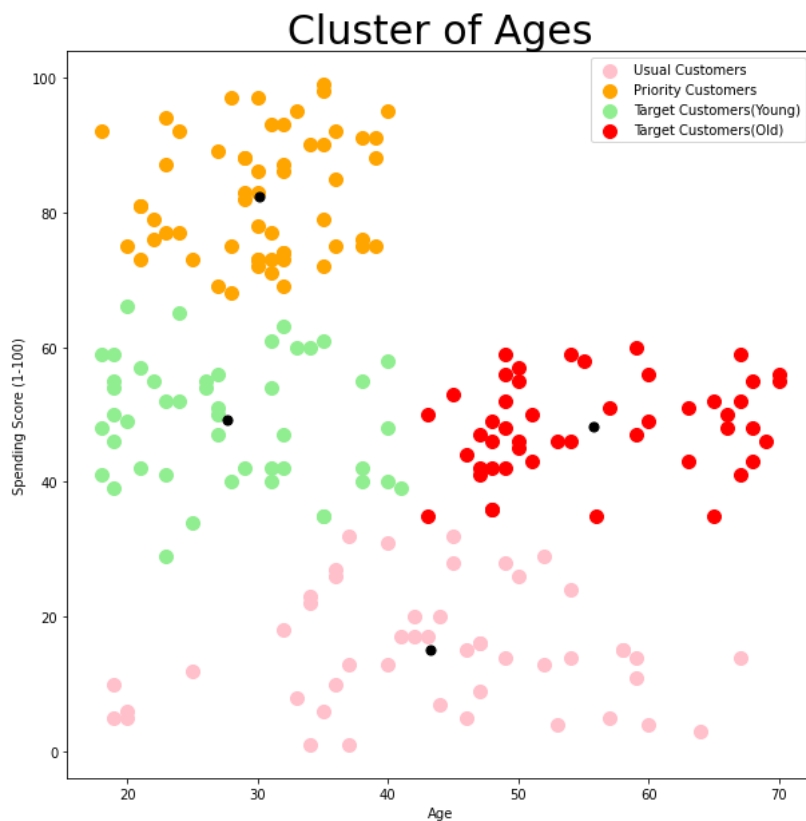
km2 = KMeans(n_clusters = 4, init = 'k-means++', max_iter = 300, n_init = 10, random_state = 0)
ymeans = km2.fit_predict(x)

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

plt.scatter(x[ymmeans == 0, 0], x[ymmeans == 0, 1], s = 100, c = 'pink', label = 'Usual Customers' )
plt.scatter(x[ymmeans == 1, 0], x[ymmeans == 1, 1], s = 100, c = 'orange', label = 'Priority Customers')
plt.scatter(x[ymmeans == 2, 0], x[ymmeans == 2, 1], s = 100, c = 'lightgreen', label = 'Target Customers(Young)')
plt.scatter(x[ymmeans == 3, 0], x[ymmeans == 3, 1], s = 100, c = 'red', label = 'Target Customers(Old)')
plt.scatter(km2.cluster_centers_[0, 0], km2.cluster_centers_[0, 1], s = 50, c = 'black')

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

```



Clustering based on gender

In [42]:

```

x = data.iloc[:, [0, 3]].values
x.shape

```

Out[42]:

(200, 2)

In [43]:

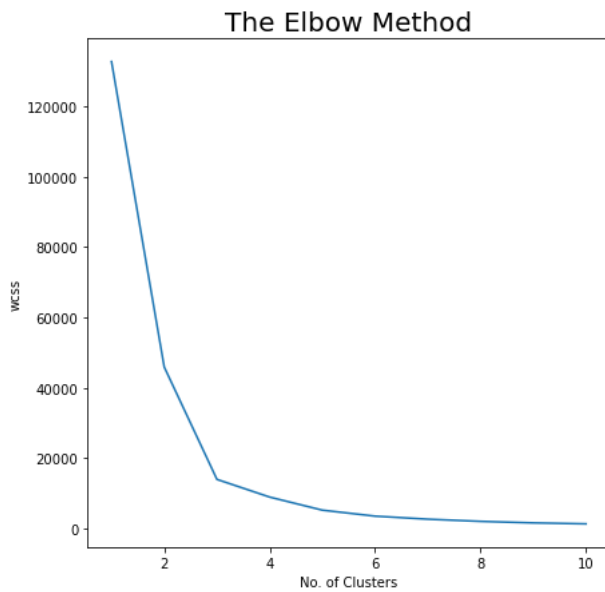
```
from sklearn.cluster import KMeans

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, 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\YASH\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(



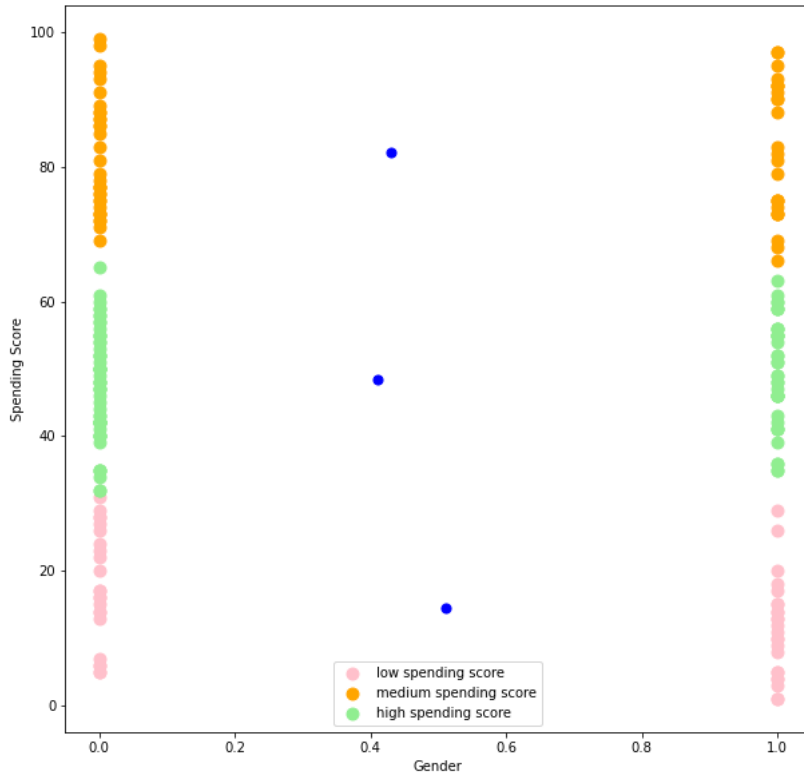
In [44]:

```

km3 = KMeans(n_clusters = 3, max_iter = 300, n_init = 10, random_state = 0)
ymeans = km3.fit_predict(x)

plt.rcParams['figure.figsize'] = (10, 10)
plt.scatter(x[ymean == 0, 0], x[ymean == 0, 1], s = 80, c = 'pink', label = 'low spending score')
plt.scatter(x[ymean == 1, 0], x[ymean == 1, 1], s = 80, c = 'orange', label = 'medium spending score')
plt.scatter(x[ymean == 2, 0], x[ymean == 2, 1], s = 80, c = 'lightgreen', label = 'high spending score')
plt.scatter(km3.cluster_centers[:,0], km3.cluster_centers[:, 1], s = 50, color = 'blue')
plt.legend()
plt.xlabel('Gender')
plt.ylabel('Spending Score')
plt.show()

```



From Above cluster plot we can clearly see that males and females are in all the category that is high low and medium spending score category