# Clustering With K Means - Assignment 2

In [1]:
```python
from sklearn.cluster import KMeans
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from matplotlib import pyplot as plt
%matplotlib inline
```
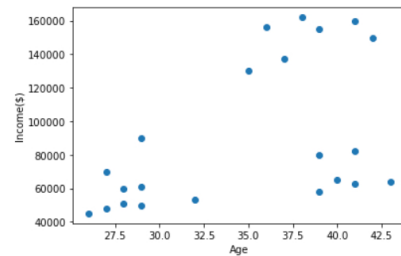
In [2]:
```python
df = pd.read_csv("income.csv")
df.head()
```

Out[2]:

|   | Name | Age | Income($) |
|---|------|-----|-----------|
| 0 | Rob | 27 | 70000 |
| 1 | Michael | 29 | 90000 |
| 2 | Mohan | 29 | 61000 |
| 3 | Ismail | 28 | 60000 |
| 4 | Kory | 42 | 150000 |

In [3]:
```python
plt.scatter(df.Age,df['Income($)'])
plt.xlabel('Age')
plt.ylabel('Income($)')
```

Out[3]: Text(0, 0.5, 'Income($)')



In [4]:
```python
km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df[['Age','Income($)']])
y_predicted
```

Out[4]: array([2, 2, 0, 0, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2, 0],
        dtype=int32)

In [5]:
```python
df['cluster']=y_predicted
df.head()
```
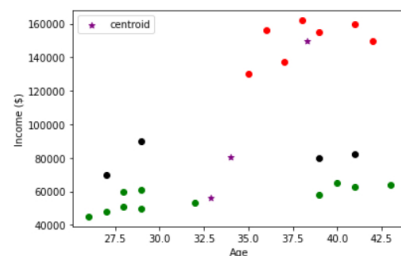
Out[5]:

|   | Name | Age | Income($) | cluster |
|---|------|-----|-----------|---------|
| 0 | Rob | 27 | 70000 | 2 |
| 1 | Michael | 29 | 90000 | 2 |
| 2 | Mohan | 29 | 61000 | 0 |
| 3 | Ismail | 28 | 60000 | 0 |
| 4 | Kory | 42 | 150000 | 1 |

In [6]:
```python
km.cluster_centers_
```

Out[6]: array([[3.29090909e+01, 5.61363636e+04],
        [3.82857143e+01, 1.50000000e+05],
        [3.40000000e+01, 8.05000000e+04]])

In [7]:
```python
df1 = df[df.cluster==0]
df2 = df[df.cluster==1]
df3 = df[df.cluster==2]
plt.scatter(df1.Age,df1['Income($)'],color='green')
plt.scatter(df2.Age,df2['Income($)'],color='red')
plt.scatter(df3.Age,df3['Income($)'],color='black')
plt.scatter(km.cluster_centers_[:,0],km.cluster_centers_[:,1],color='purple',marker='*',label='centroid')
plt.xlabel('Age')
plt.ylabel('Income ($)')
plt.legend()
```

Out[7]: <matplotlib.legend.Legend at 0x7f63e246e668>



### Preprocessing using min max scaler

In [8]:
```python
scaler = MinMaxScaler()

scaler.fit(df[['Income($)']])
```

```
df['Income($)'] = scaler.transform(df[['Income($)']])

scaler.fit(df[['Age']])
df['Age'] = scaler.transform(df[['Age']])
```
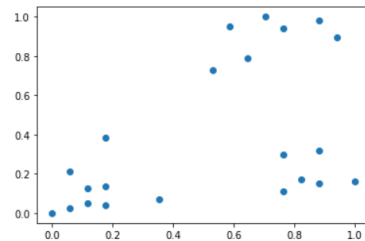
In [9]: `df.head()`

Out[9]:

|   | Name | Age | Income($) | cluster |
|---|------|-----|-----------|---------|
| 0 | Rob | 0.058824 | 0.213675 | 2 |
| 1 | Michael | 0.176471 | 0.384615 | 2 |
| 2 | Mohan | 0.176471 | 0.136752 | 0 |
| 3 | Ismail | 0.117647 | 0.128205 | 0 |
| 4 | Kory | 0.941176 | 0.897436 | 1 |

In [10]: `plt.scatter(df.Age,df['Income($)'])`

Out[10]: `<matplotlib.collections.PathCollection at 0x7f63e23ee048>`



In [11]:
```
km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df[['Age','Income($)']])
y_predicted
```

Out[11]: `array([0, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 2, 2, 2, 2, 2, 2],`
`       dtype=int32)`
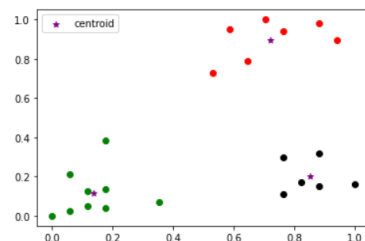
In [12]:
```
df['cluster']=y_predicted
df.head()
```

Out[12]:

|   | Name | Age | Income($) | cluster |
|---|------|-----|-----------|---------|
| 0 | Rob | 0.058824 | 0.213675 | 0 |
| 1 | Michael | 0.176471 | 0.384615 | 0 |
| 2 | Mohan | 0.176471 | 0.136752 | 0 |
| 3 | Ismail | 0.117647 | 0.128205 | 0 |
| 4 | Kory | 0.941176 | 0.897436 | 1 |

In [13]: `km.cluster_centers_`

Out[13]: `array([[0.1372549 , 0.11633428],`
`        [0.72268908, 0.8974359 ],`
`        [0.85294118, 0.2022792 ]])`

In [14]:
```
df1 = df[df.cluster==0]
df2 = df[df.cluster==1]
df3 = df[df.cluster==2]
plt.scatter(df1.Age,df1['Income($)'],color='green')
plt.scatter(df2.Age,df2['Income($)'],color='red')
plt.scatter(df3.Age,df3['Income($)'],color='black')
plt.scatter(km.cluster_centers_[:,0],km.cluster_centers_[:,1],color='purple',marker='*',label='centroid')
plt.legend()
```
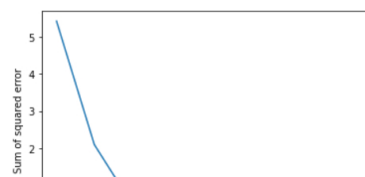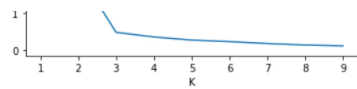
Out[14]: `<matplotlib.legend.Legend at 0x7f63e23d5ef0>`



**Elbow Plot**

In [15]:
```
sse = []
k_rng = range(1,10)
for k in k_rng:
    km = KMeans(n_clusters=k)
    km.fit(df[['Age','Income($)']])
    sse.append(km.inertia_)
```

In [16]:
```
plt.xlabel('K')
plt.ylabel('Sum of squared error')
plt.plot(k_rng,sse)
```

Out[16]: `[<matplotlib.lines.Line2D at 0x7f63e2361390>]`

In [ ]: