

K Means

68_Adnan Shaikh

```
In [1]: from sklearn.cluster import KMeans
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from matplotlib import pyplot as plt
import warnings
warnings.filterwarnings("ignore")
```

```
In [2]: df=pd.read_excel("./income.xlsx")
```

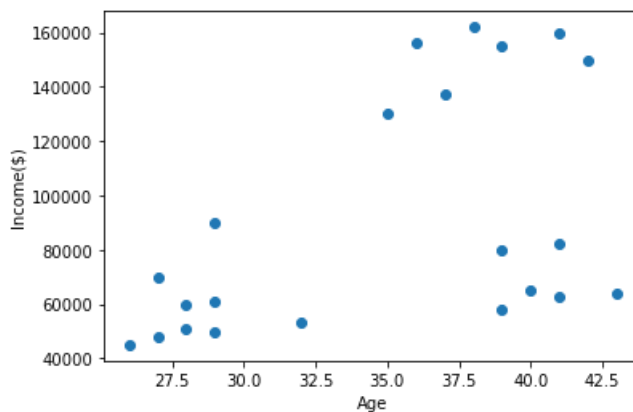
```
In [3]: df.head()
```

Out[3]:

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

```
In [4]: plt.scatter(df.Age,df['Income($)'])
plt.xlabel('Age')
plt.ylabel('Income($)')
```

Out[4]: Text(0, 0.5, 'Income(\$))')



```
In [5]: km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df[['Age','Income($)']])
y_predicted
```

Out[5]: array([2, 2, 0, 0, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 2, 2, 0])

```
In [6]: df['cluster']=y_predicted
df
```

Out[6]:

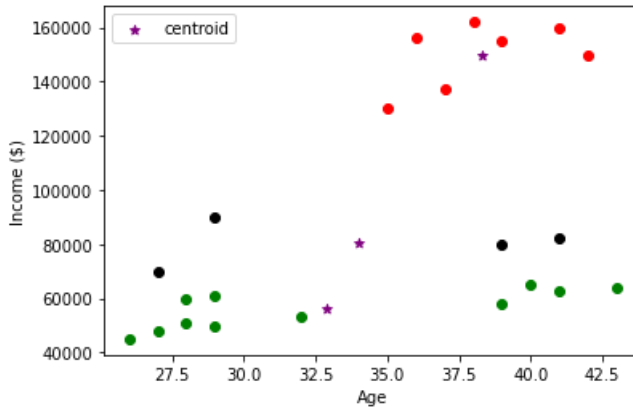
	Name	Age	Income(\$)	Unnamed: 3	cluster
0	Rob	27	70000	NaN	2
1	Michael	29	90000	NaN	2
2	Mohan	29	61000	NaN	0
3	Ismail	28	60000	NaN	0
4	Kory	42	150000	NaN	1
5	Gautam	39	155000	NaN	1
6	David	41	160000	NaN	1
7	Andrea	38	162000	NaN	1
8	Brad	36	156000	NaN	1
9	Angelina	35	130000	NaN	1
10	Donald	37	137000	NaN	1
11	Tom	26	45000	NaN	0
12	Arnold	27	48000	NaN	0
13	Jared	28	51000	NaN	0
14	Stark	29	49500	NaN	0
15	Ranbir	32	53000	NaN	0
16	Dipika	40	65000	NaN	0
17	Priyanka	41	63000	NaN	0
18	Nick	43	64000	NaN	0
19	Alia	39	80000	NaN	2
20	Sid	41	82000	NaN	2
21	Abdul	39	58000	NaN	0

```
In [7]: km.cluster_centers_
```

Out[7]: array([[3.29090909e+01, 5.61363636e+04],
[3.82857143e+01, 1.50000000e+05],
[3.40000000e+01, 8.05000000e+04]])

```
In [8]: 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[8]: <matplotlib.legend.Legend at 0x22d95fdda60>



Preprocessing using min max scaler

```
In [9]: scaler = MinMaxScaler()

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

scaler.fit(df[['Age']])
df['Age'] = scaler.transform(df[['Age']])
```

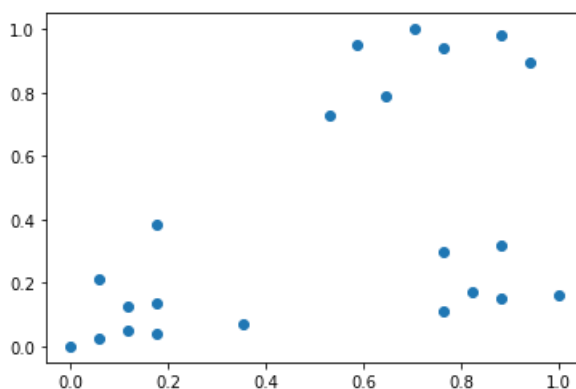
In [10]: df.head()

Out[10]:

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

In [11]: plt.scatter(df.Age,df['Income(\$)'])

Out[11]: <matplotlib.collections.PathCollection at 0x22d96069340>



```
In [12]: km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df[['Age', 'Income($)']])
y_predicted
```

```
Out[12]: array([1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0])
```

```
In [13]: df['cluster']=y_predicted
df.head()
```

```
Out[13]:
```

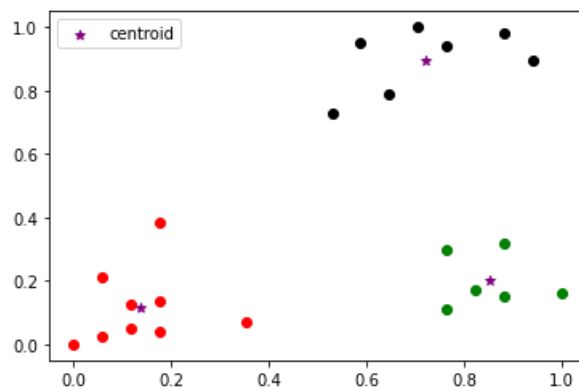
	Name	Age	Income(\$)	Unnamed: 3	cluster
0	Rob	0.058824	0.213675	NaN	1
1	Michael	0.176471	0.384615	NaN	1
2	Mohan	0.176471	0.136752	NaN	1
3	Ismail	0.117647	0.128205	NaN	1
4	Kory	0.941176	0.897436	NaN	2

```
In [14]: km.cluster_centers_
```

```
Out[14]: array([[0.85294118, 0.2022792 ],
                [0.1372549 , 0.11633428],
                [0.72268908, 0.8974359 ]])
```

```
In [15]: 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[15]: <matplotlib.legend.Legend at 0x22d962d5b20>
```



Elbow Plot

```
In [16]: 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 [17]: plt.xlabel('K')  
plt.ylabel('Sum of squared error')  
plt.plot(k_rng,sse)
```

Out[17]: [<matplotlib.lines.Line2D at 0x22d96368610>]

