

# 1.KMeans (SKLearn)

November 13, 2021

K-Means (With SKLearn)

```
[1]: from sklearn.cluster import KMeans
import pandas as pd
from sklearn.preprocessing import MinMaxScaler
from matplotlib import pyplot as plt
%matplotlib inline
```

## 0.0.1 Step 1: Load the dataset

```
[2]: df = pd.read_csv("E:\\MY LECTURES\\8.2021-09-03 DATA SCIENCE (KNU)\\3.
↳Programs\\dataset\\income.csv")
df.head()
```

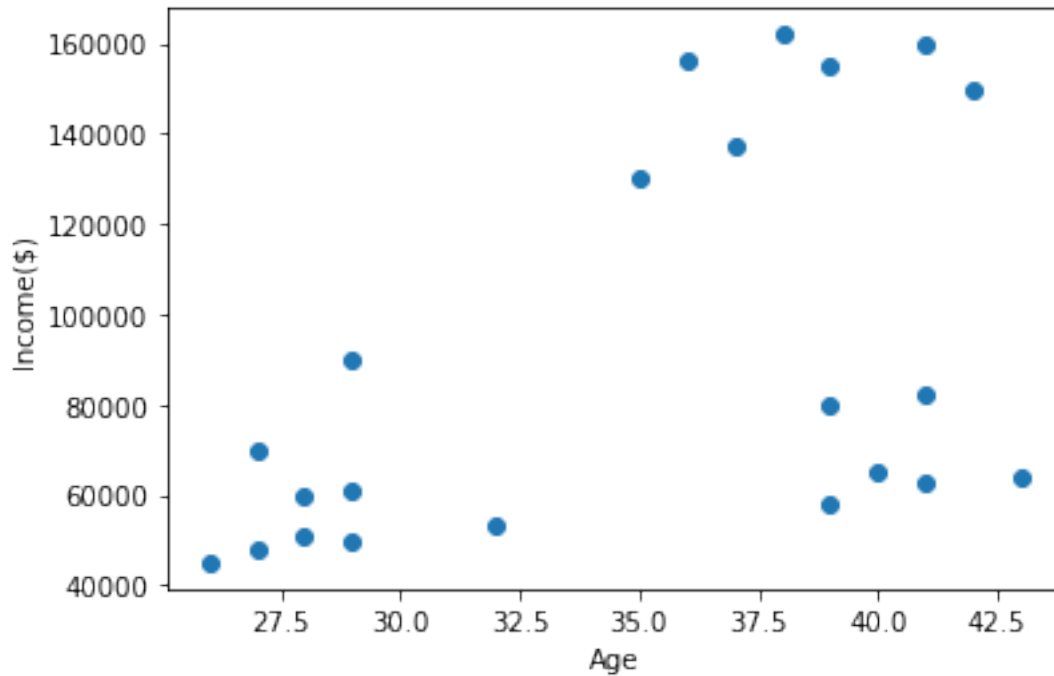
```
[2]:
```

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

## 0.0.2 Step 2: EDA

```
[3]: # Scatter plot to find the pattern
plt.scatter(df.Age,df['Income($)'])
plt.xlabel('Age')
plt.ylabel('Income($)')
```

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



### 0.0.3 Step 4: Training the model

```
[4]: # Find the initial
km = KMeans(n_clusters=3)
y_predicted = km.fit_predict(df[['Age', 'Income($)']])
y_predicted
```

```
[4]: array([2, 2, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2, 0])
```

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

```
[5]:
```

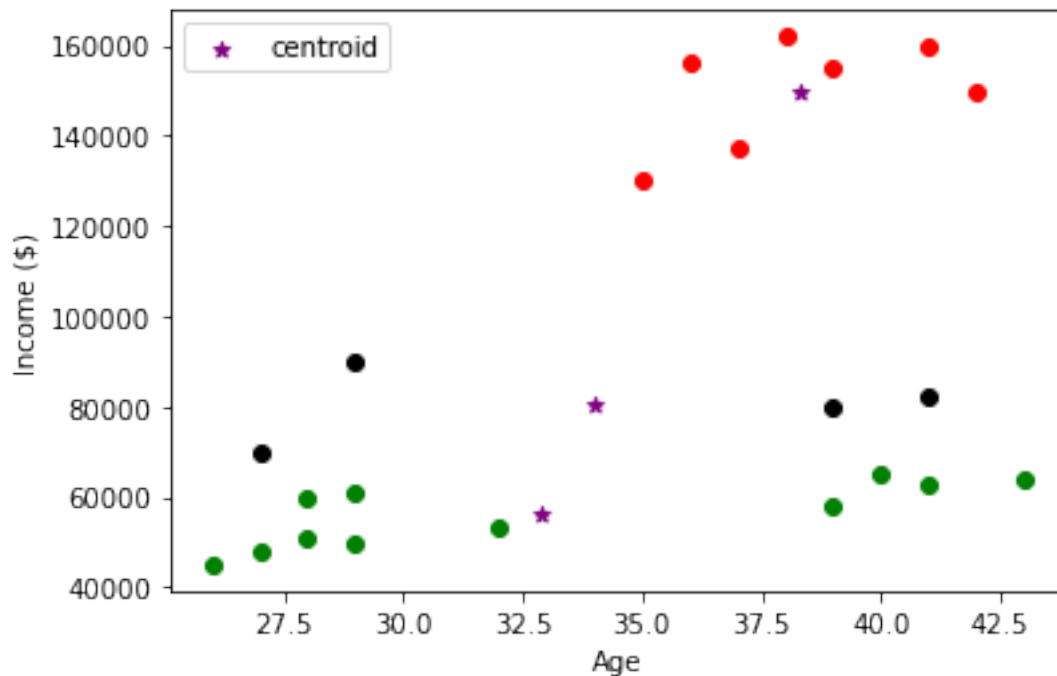
	Name	Age	Income(\$)	cluster
0	Rob	27	70000	2
1	Michael	29	90000	2
2	Mohan	29	61000	0
3	Ismail	28	60000	0
4	Kory	42	150000	1

```
[6]: # cluster centres
km.cluster_centers_
```

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

```
[7]: # Plot the output
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()
```

```
[7]: <matplotlib.legend.Legend at 0x25a88751d90>
```



#### 0.0.4 Step 3: Pre-processing

We do this after fitting, because we found cluster is not in place as the scale of tick values are highly varying on x and y axis

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

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

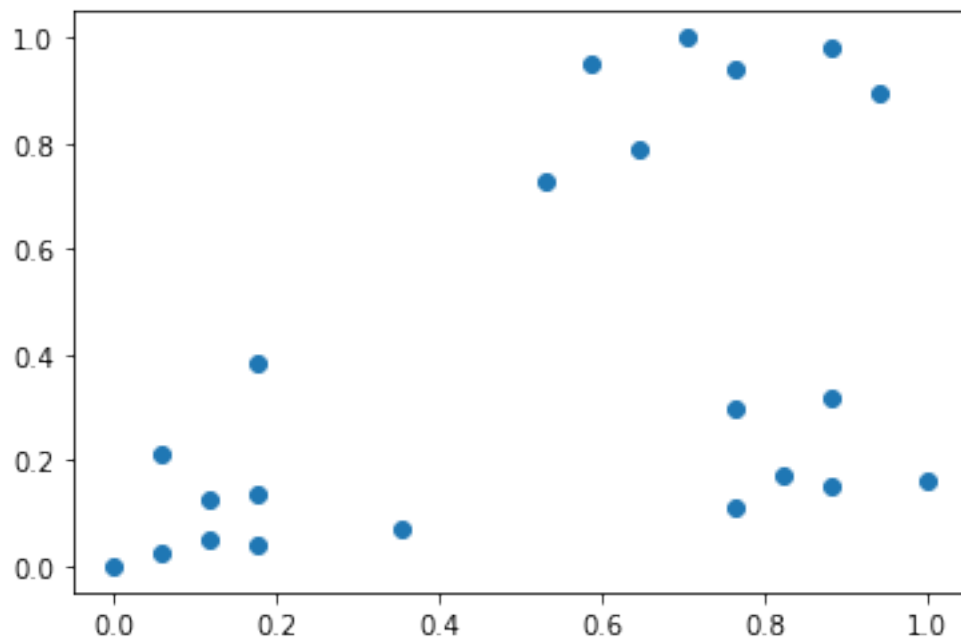
scaler.fit(df[['Age']])
df['Age'] = scaler.transform(df[['Age']])
```

```
[9]: df.head()
```

```
[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
```

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

```
[10]: <matplotlib.collections.PathCollection at 0x25a887e0220>
```



### 0.0.5 Step 4: Training the model again

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

```
[11]: array([1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2])
```

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

```
[12]:
```

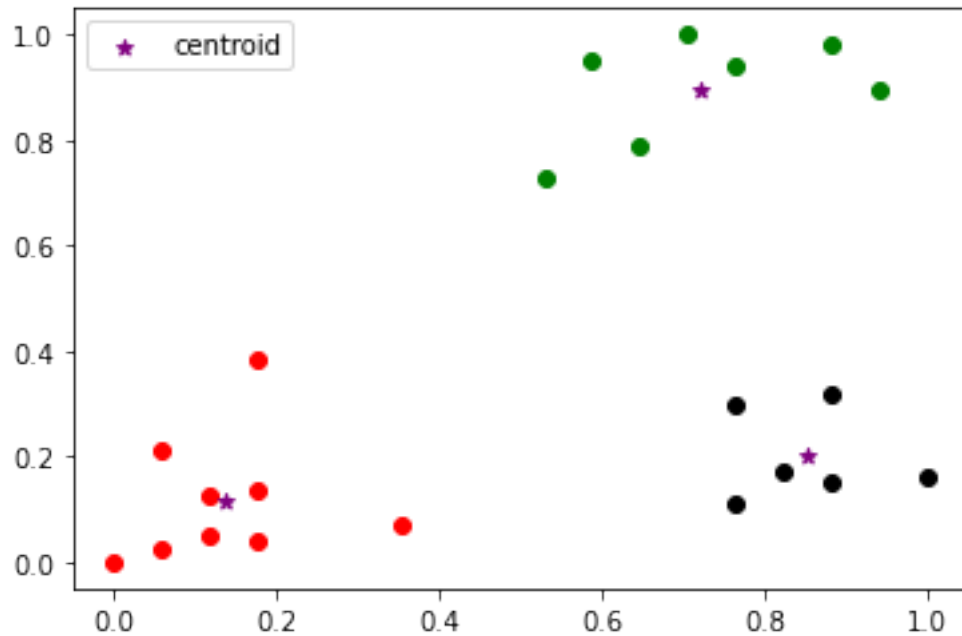
	Name	Age	Income(\$)	cluster
0	Rob	0.058824	0.213675	1
1	Michael	0.176471	0.384615	1
2	Mohan	0.176471	0.136752	1
3	Ismail	0.117647	0.128205	1
4	Kory	0.941176	0.897436	0

```
[13]: km.cluster_centers_
```

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

```
[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()
```

```
[14]: <matplotlib.legend.Legend at 0x25a8883ba30>
```



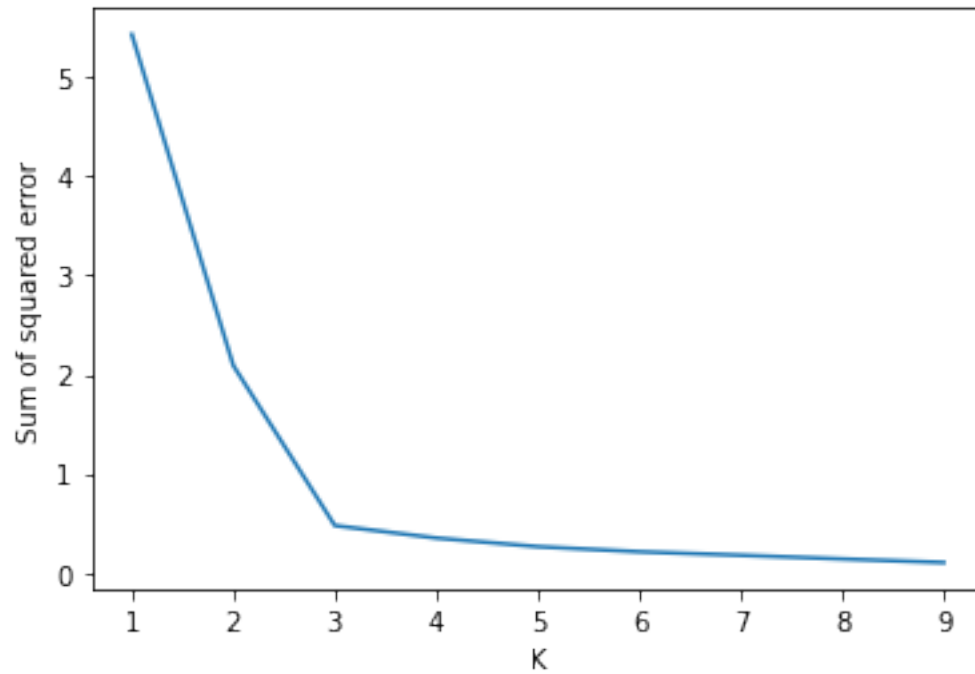
### 0.0.6 Step 5: Elbow plot to find

To find the right number of cluster

```
[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_)
```

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

```
[16]: [<matplotlib.lines.Line2D at 0x25a888c2130>]
```



### 0.0.7 Exercise

1. Use iris flower dataset from sklearn library and try to form clusters of flowers using petal width and length features. Drop other two features for simplicity.
2. Figure out if any preprocessing such as scaling would help here
3. Draw elbow plot and from that figure out optimal value of k