

```
In [1]: 1 from sklearn.cluster import KMeans
2 import pandas as pd
3 from sklearn.preprocessing import MinMaxScaler
4 from matplotlib import pyplot as plt
5 %matplotlib inline
```

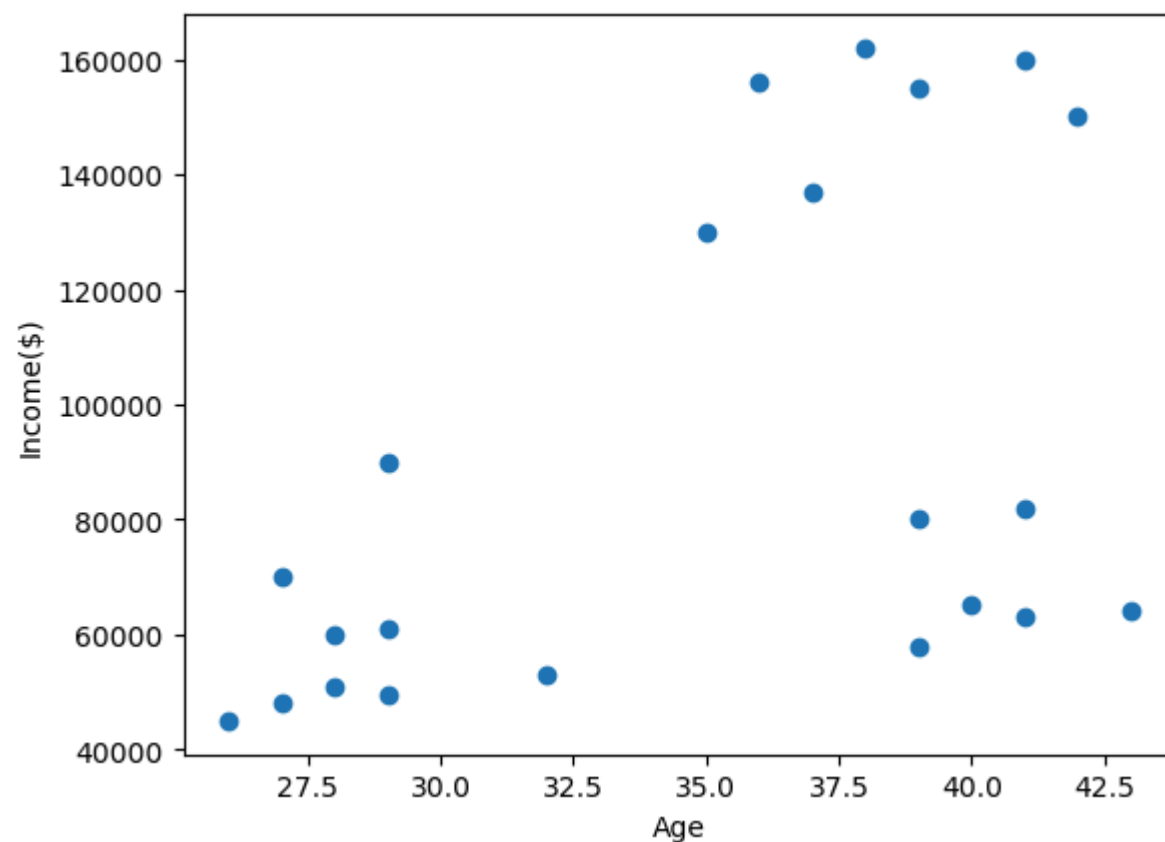
```
In [2]: 1 df = pd.read_csv("income.csv")
2 df.head()
```

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

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



```
In [4]: 1 km = KMeans(n_clusters=3)
2 y_predicted = km.fit_predict(df[['Age', 'Income($)']])
3 y_predicted
```

C:\Users\shikh\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:870: FutureWarning: The default value of `n\_init` will change from 10 to 'auto' in 1.4. Set the value of `n\_init` explicitly to suppress the warning  
warnings.warn(  
C:\Users\shikh\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1382: 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(  
array([0, 0, 2, 2, 1, 1, 1, 1, 1, 1, 1, 2, 2, 2, 2, 2, 2, 2, 0, 0, 2])

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

```
In [5]: 1 df['cluster']=y_predicted
2 df.head()
```

Out[5]:

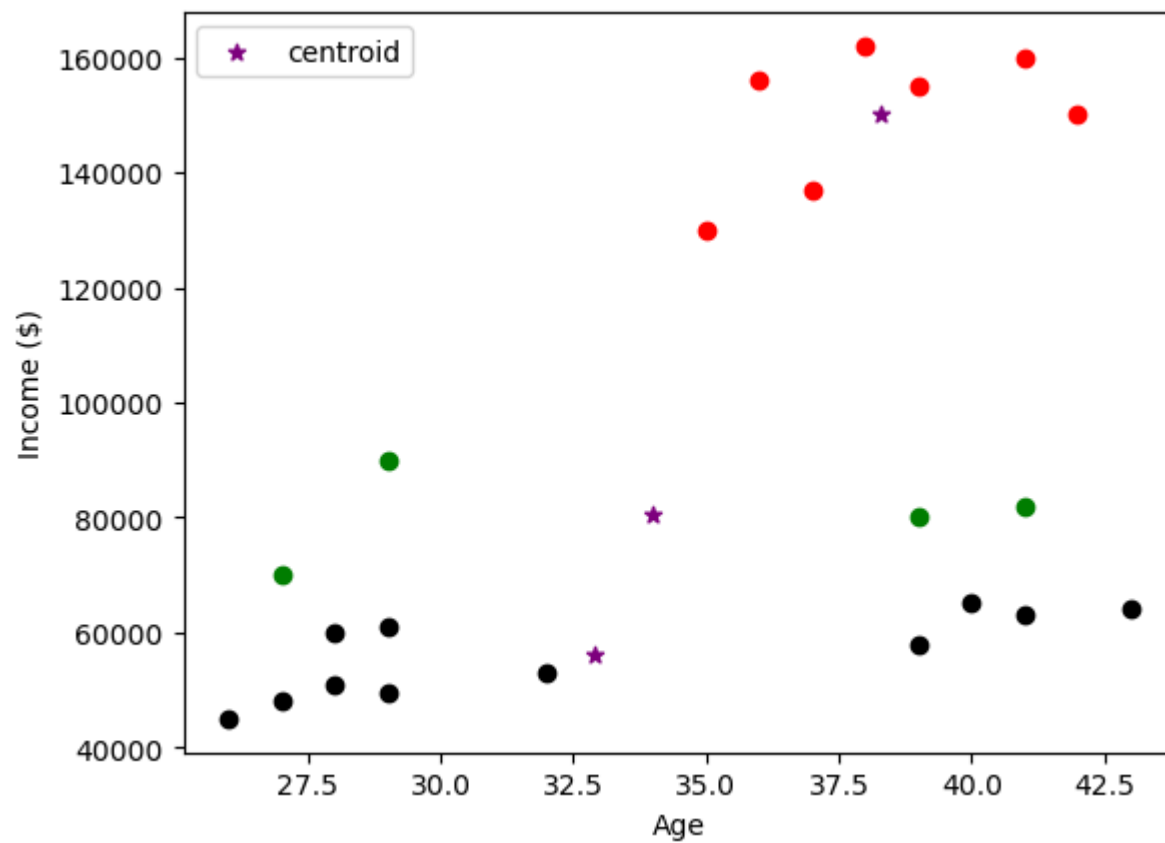
	Name	Age	Income(\$)	cluster
0	Rob	27	70000	0
1	Michael	29	90000	0
2	Mohan	29	61000	2
3	Ismail	28	60000	2
4	Kory	42	150000	1

```
In [6]: 1 km.cluster_centers_
```

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

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

```
Out[7]: <matplotlib.legend.Legend at 0x1a7a325c790>
```



### Preprocessing using min max scaler

```
In [8]: 1 scaler = MinMaxScaler()
2
3 scaler.fit(df[['Income($)']])
4 df['Income($)'] = scaler.transform(df[['Income($)']])
5
6 scaler.fit(df[['Age']])
7 df['Age'] = scaler.transform(df[['Age']])
```

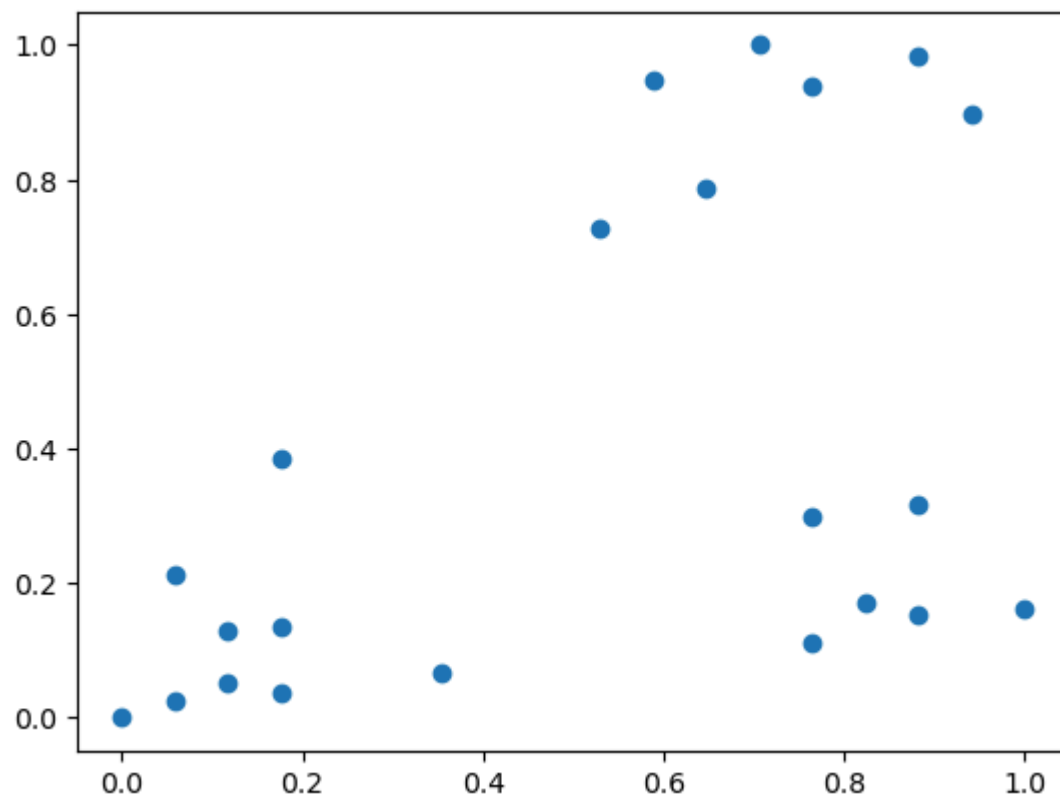
```
In [9]: 1 df.head()
```

```
Out[9]:
```

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

```
In [10]: 1 plt.scatter(df.Age,df['Income($)'])
```

```
Out[10]: <matplotlib.collections.PathCollection at 0x1a7a420ba90>
```



```
In [11]: 1 km = KMeans(n_clusters=3)
2 y_predicted = km.fit_predict(df[['Age', 'Income($)']])
3 y_predicted
```

C:\Users\shikh\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:870: FutureWarning: The default value of `n\_init` will change from 10 to 'auto' in 1.4. Set the value of `n\_init` explicitly to suppress the warning

```
warnings.warn(
```

C:\Users\shikh\anaconda3\Lib\site-packages\sklearn\cluster\\_kmeans.py:1382: 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(
```

```
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])
```

```
In [12]: 1 df['cluster']=y_predicted
2 df.head()
```

```
Out[12]:
```

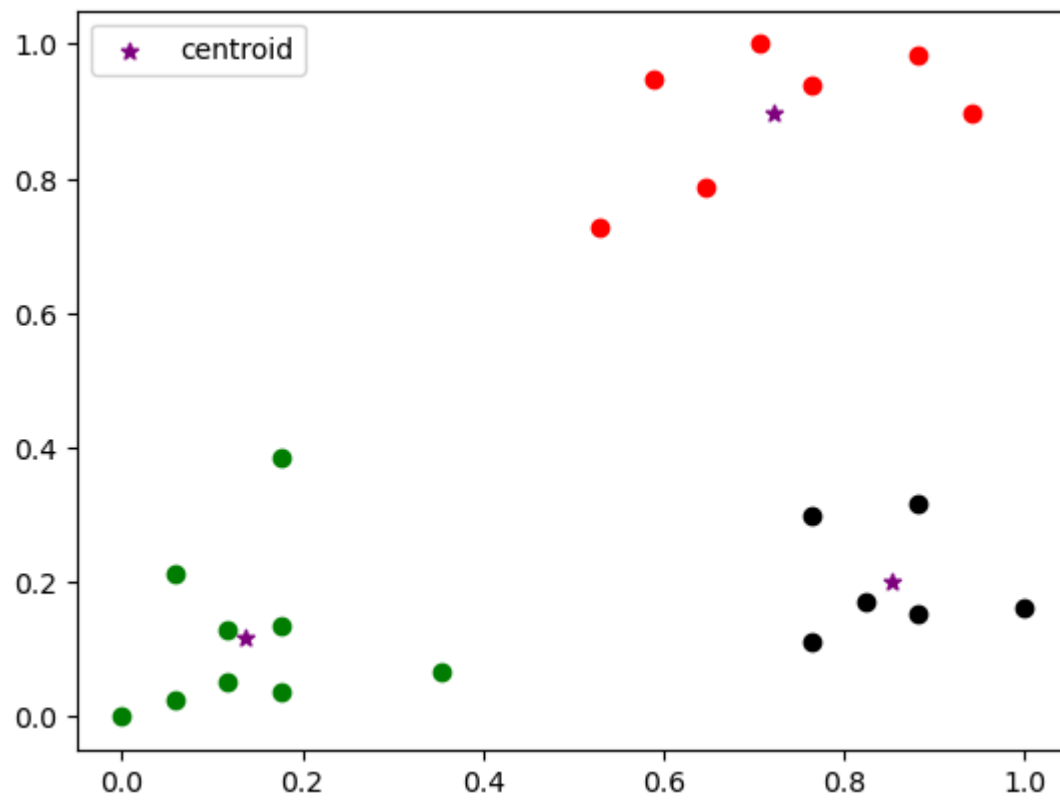
	Name	Age	Income(\$)	cluster
0	Rob	0.058824	0.213675	0
1	Michael	0.176471	0.384615	0
2	Mohan	0.176471	0.136752	0
3	Ismail	0.117647	0.128205	0
4	Kory	0.941176	0.897436	1

```
In [13]: 1 km.cluster_centers_
```

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

```
In [14]: 1 df1 = df[df.cluster==0]
2 df2 = df[df.cluster==1]
3 df3 = df[df.cluster==2]
4 plt.scatter(df1.Age,df1['Income($)'],color='green')
5 plt.scatter(df2.Age,df2['Income($)'],color='red')
6 plt.scatter(df3.Age,df3['Income($)'],color='black')
7 plt.scatter(km.cluster_centers_[0],km.cluster_centers_[1],color='purple',marker='*',label='centroid')
8 plt.legend()
```

Out[14]: <matplotlib.legend.Legend at 0x1a7a41ca590>



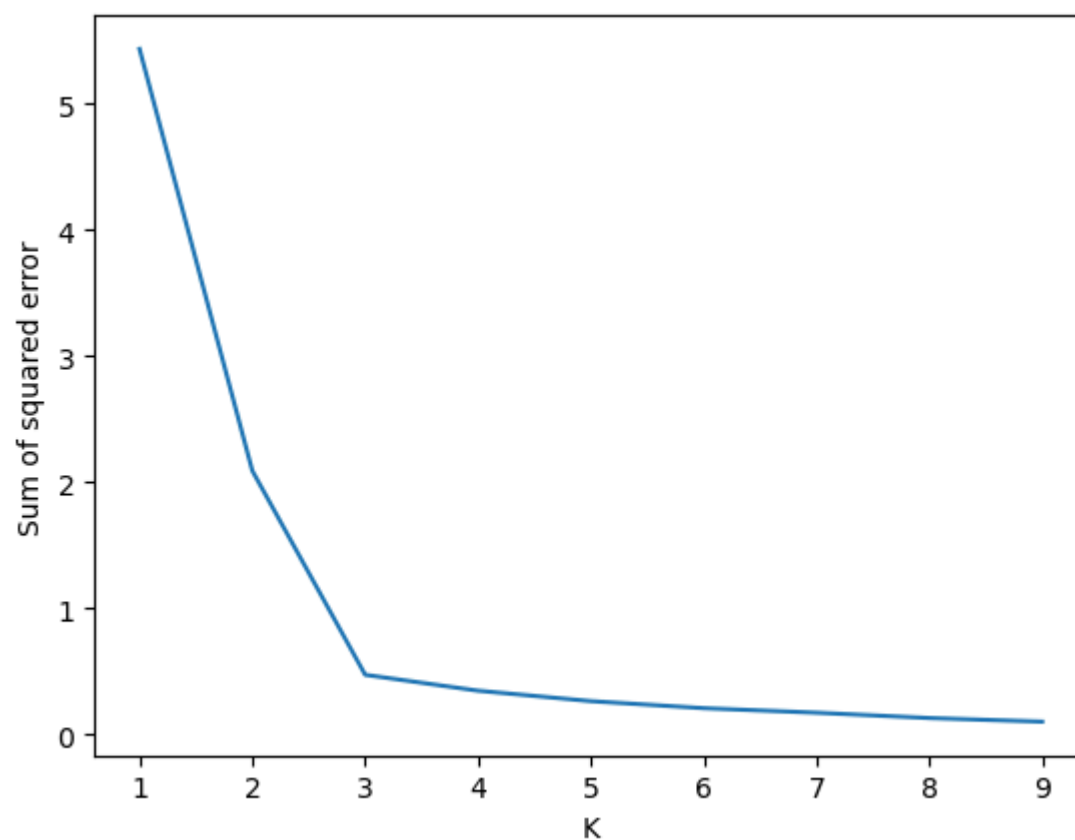
### Elbow Plot

To find the perfect value of k

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

```
In [16]: 1 plt.xlabel('K')
2 plt.ylabel('Sum of squared error')
3 plt.plot(k_rng,sse)
```

Out[16]: [<matplotlib.lines.Line2D at 0x1a7a380aa90>]



DBSCAN

```
In [17]: 1 df = pd.read_csv("income.csv")
        2 df.head()
```

Out[17]:

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

```
In [19]: 1 features = df[['Age', 'Income($)']].values
        2 features
```

Out[19]: array([[ 27, 70000],
 [ 29, 90000],
 [ 29, 61000],
 [ 28, 60000],
 [ 42, 150000],
 [ 39, 155000],
 [ 41, 160000],
 [ 38, 162000],
 [ 36, 156000],
 [ 35, 130000],
 [ 37, 137000],
 [ 26, 45000],
 [ 27, 48000],
 [ 28, 51000],
 [ 29, 49500],
 [ 32, 53000],
 [ 40, 65000],
 [ 41, 63000],
 [ 43, 64000],
 [ 39, 80000],
 [ 41, 82000],
 [ 39, 58000]], dtype=int64)

```
In [21]: 1 from sklearn.cluster import DBSCAN
        2 dbscan = DBSCAN(eps=10000, min_samples=2) # Adjust epsilon and min_samples according to your dataset
        3 clusters = dbscan.fit_predict(features)
```

```
In [22]: 1 plt.scatter(features[:, 0], features[:, 1], c=clusters, cmap='viridis')
        2 plt.xlabel("Age")
        3 plt.ylabel("Income")
        4 plt.title("DBSCAN Clustering")
        5 plt.show()
```

