

```
In [1]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [5]: df=mpd.read_csv(r"C:\Users\aditr\Desktop\project report\1. customer segmentation\Mall_Customers.csv")
df.head()
```

CustomerID	Gender	Age	Annual Income (k\$)	Spending Score (1-100)
0	1 Male	19	15	39
1	2 Male	21	15	81
2	3 Female	20	16	6
3	4 Female	23	16	77
4	5 Female	31	17	40

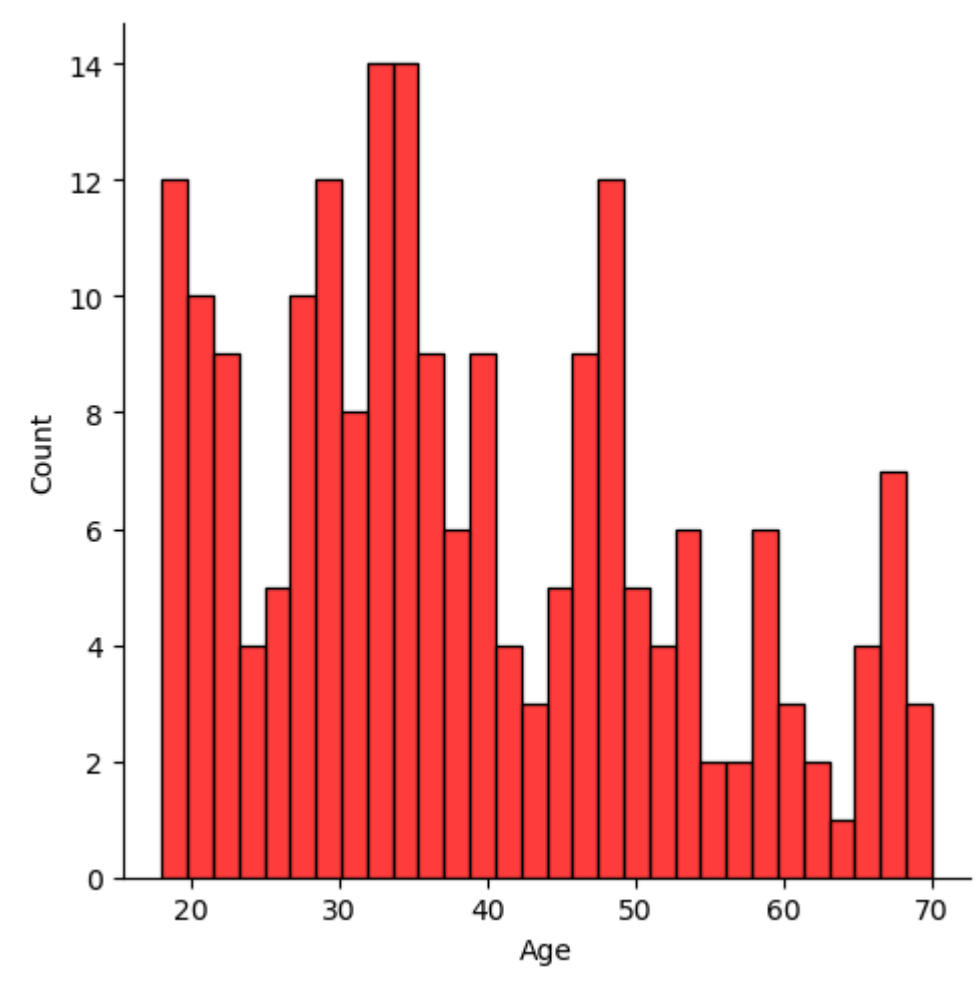
```
In [9]: df.shape
```

```
Out[9]: (200, 5)
```

	df.describe()
count	200.000000
mean	100.500000
std	57.879185
min	1.000000
25%	50.750000
50%	100.500000
75%	150.250000
max	200.000000
CustomerID	
Age	
Annual Income (k\$)	
Spending Score (1-10)	
	200.000000
	50.200000
	25.823525
	1.000000
	34.750000
	50.000000
	73.000000
	99.000000

```
In [17]: sns.displot(df['Age'], kde = False, color = 'red', bins = 30)
```

```
Out[17]: <seaborn.axisgrid.FacetGrid at 0x2a536648b90>
```



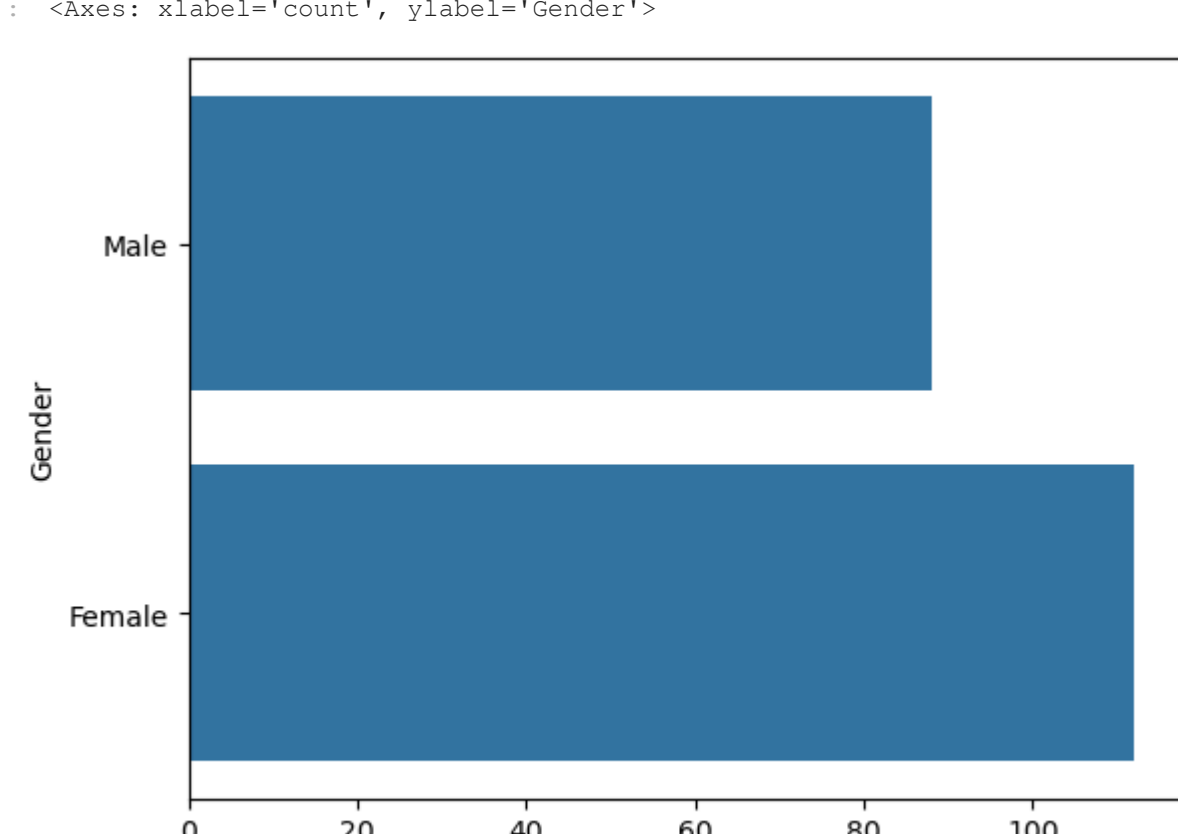
THE AVERAGE CUSTOMER AROUND THE MALL ARE OF GENERALLY AGE 31-34

```
In [22]: sns.countplot(y='Gender', data=df)
```

```

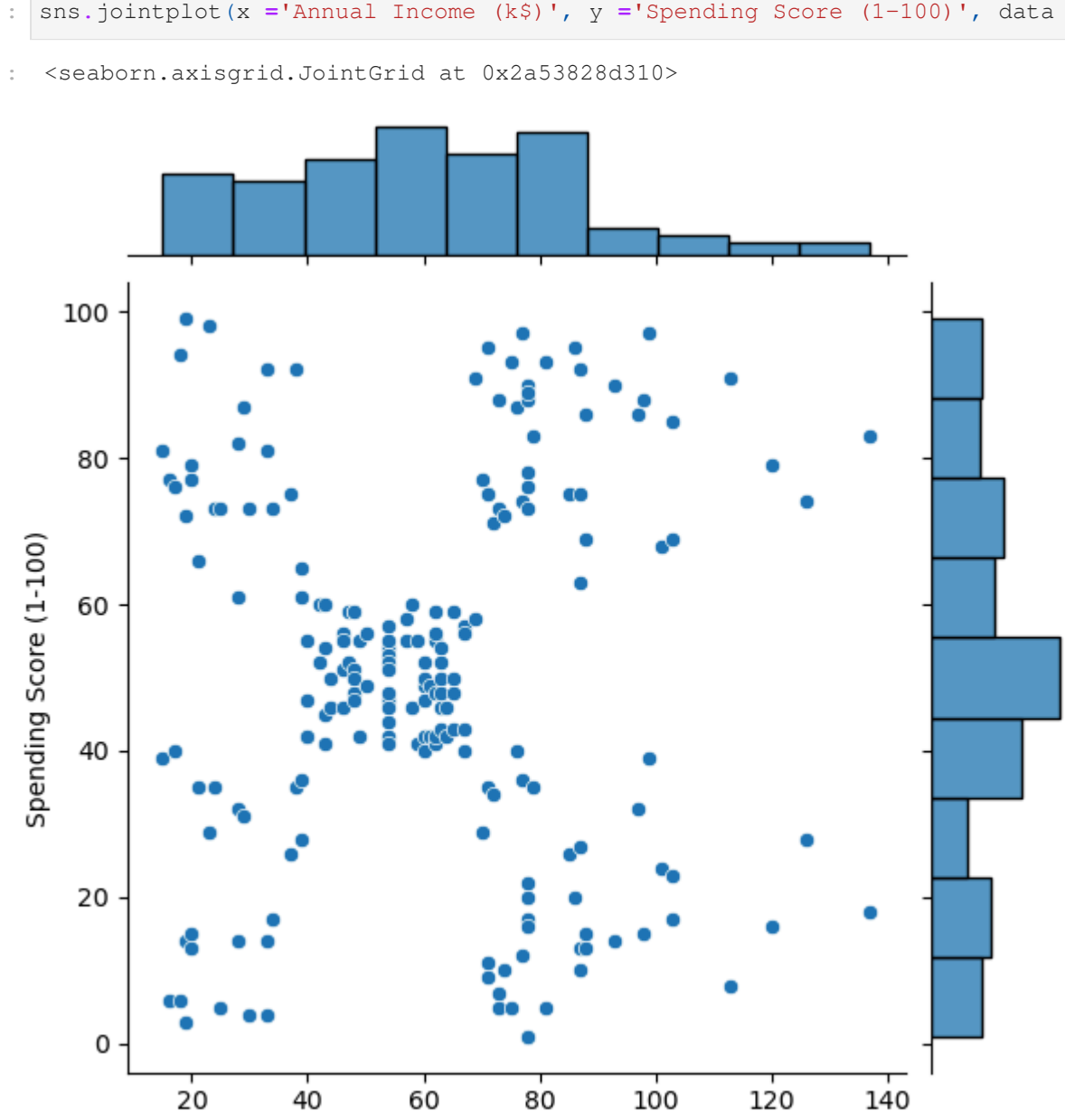
111 [22]: sns.countplot(y = gender , data = df)
112
113

```



count

THIS MALL HA GENERALLY MORE FEMALE CUSTOMER THEN MALE



Annual Income (k\$)

Question : Why choosing K-Clustering over Hierarchical Clustering and Density based Clustering ?

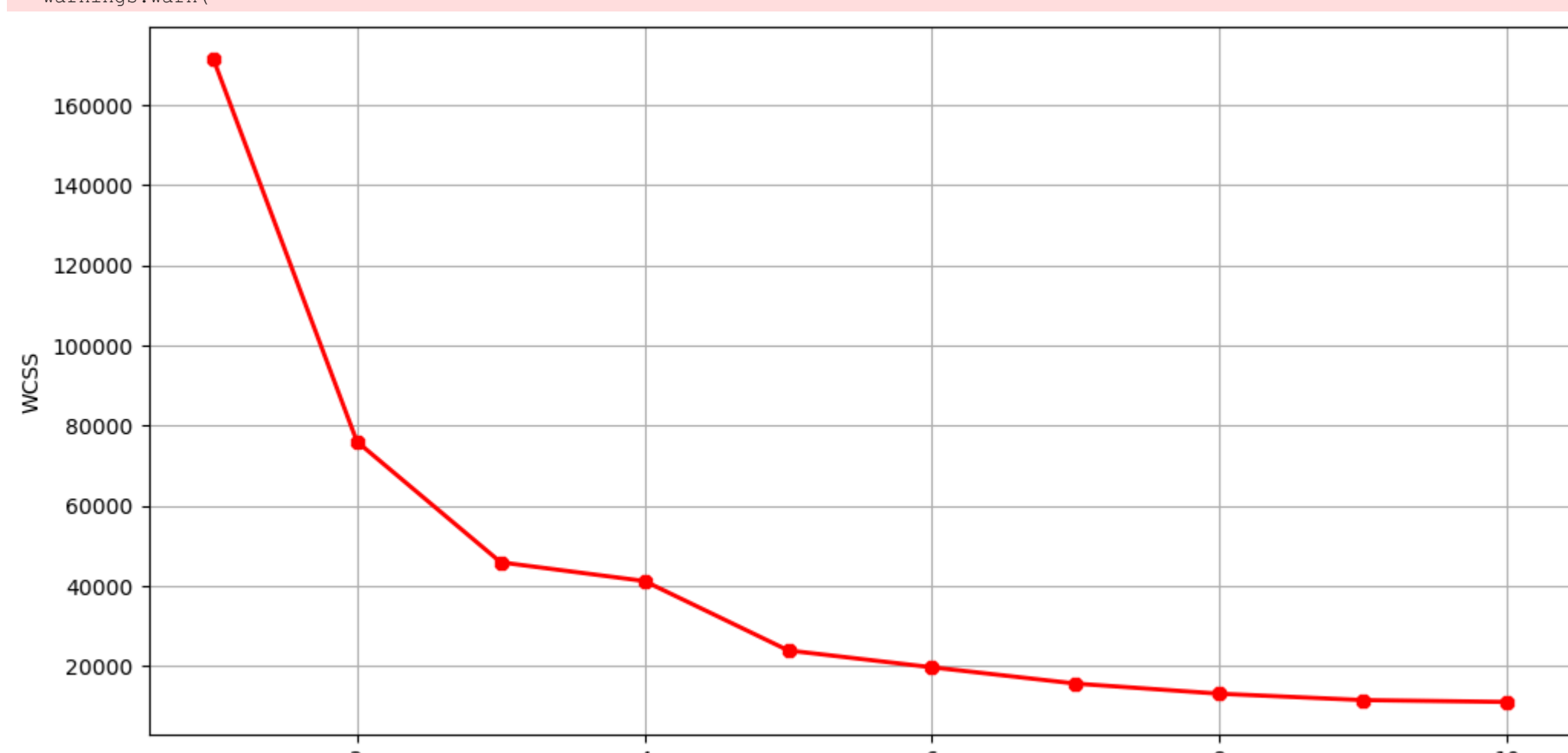
datasets with thousands or even millions of row a

```
In [44]: from sklearn.datasets import make_bl
```

```
In [54]: import os
os.environ["OMP_NUM_THREADS"] = "1"

In [58]: k=def_loop, "kobj">Spending Score [1-100"]].values
from sklearn.cluster import KMeans
wcss = []
for k in range(1,11):
    kmeans = KMeans(n_clusters= k, init="k-means++")
    kmeans.fit(X)
    wcss.append(kmeans.inertia_)

plt.figure(figsize=(12,6))
plt.plot(range(1,11),wcss, linewidth=2, color="red",
plt.xlabel('K-VAlUES')
plt.ylabel('WCSS')
plt.grid()
plt.show()
```

[illegible]

2 4

assuming cluster = 5, as after which the data seems to go constant

```
[73]: kmeans = KMeans(n_clusters=4)
      label = kmeans.fit_predict(k1)
      print(label)
```

```
C:\Users\aditr\anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1446: 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 e
```

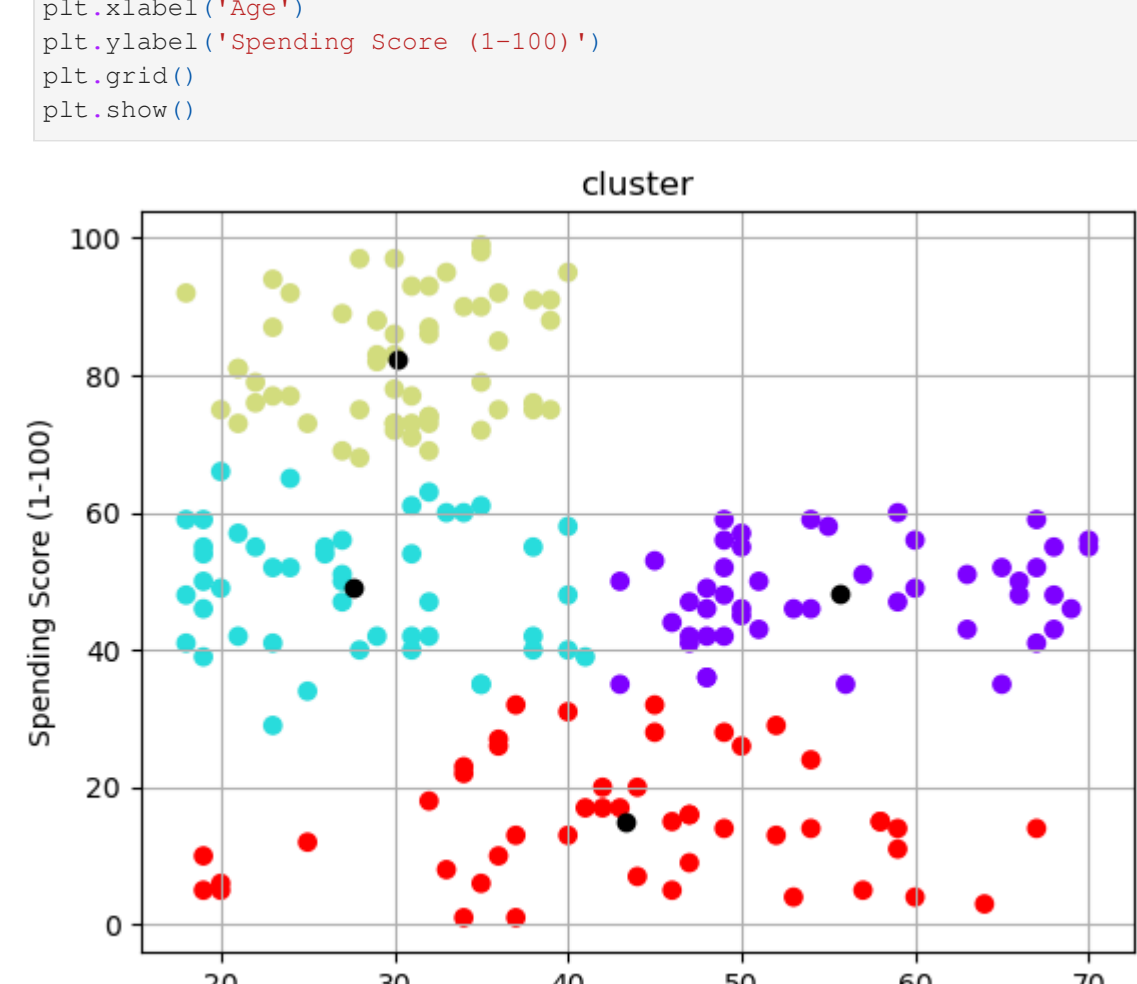
```
environment Variable CRP_NUM_THRE
warnings.warn(
```

centroids

```
In [82]: print(kmeans.cluster_centers_)

[[ 55.70833333  48.22916667]
 [ 127.61702128  49.14893617]
 [ 130.1754386   82.35087719]
 [ 43.29166667  15.02083333]]

In [88]: plt.scatter(kl[:,0],kl[:,1], c=kmeans.labels_, cmap='rainbow')
plt.scatter(kmeans.cluster_centers_[:,0],kmeans.cluster_centers_[:,1],color=
plt.title('cluster')
```



20 30

block data shows the peptide

```
In [92]: k2=df.loc[:,["Age","Annual Income (k$)"]].values
from sklearn.cluster import KMeans
wcss = []
for k in range(1,11):
    kmeans = KMeans(n_clusters=k, init="k-means++",
                    random_state=0)
    kmeans.fit(k2)
    wcss.append(kmeans.inertia_)
plt.figure(figsize=(12,6))
plt.plot(range(1,11),wcss, linewidth=2, color='red')
plt.xlabel('Number of Clusters')
plt.ylabel('Within-Cluster Sum of Squares (WCSS)')
plt.title('Elbow Method for Optimal Number of Clusters')
plt.grid(True)
```

```
plt.xlabel('K-VALUES')
plt.ylabel('WCSS')
plt.grid()
plt.show()
```

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()

C:\Users\adit1\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:1446: UserWarning: KMeans is known to have a memory leak on Windows with MSL, when there are less chunks than available threads. You can avoid it by setting the environment variable OMP\_NUM\_THREADS=1.  
warnings.warn()



Year	Number of people in the workforce (millions)
1990	12,500
1995	14,000
2000	16,000
2005	17,500
2010	19,000

```
In [94]: kmeans = KMeans(n_clusters=6)
label = kmeans.fit_predict(k2)
print(label)
```

0 0 0 0 0 0 0 0 0 2 0 2 0 2 0 0 0 0 0 0 0 0 0 0 0 2 0 0 0 0 0 2 0 2 0 2 0 0  
0 0 5 2 5 2 5 2 5 2

```
C:\Users\adit>anaconda3\Lib\site-packages\sklearn\cluster\_kmeans.py:1446: UserWarning:
nvironment variable OMP_NUM_THREADS=1.
warnings.warn(

In [96]: plt.scatter(k2[:,0],k2[:,1],c=kmeans.labels_, cmap='rainbow')
```

