In [1]:	<pre>import pandas as pd import matplotlib.pyplot as plt import seaborn as sns from sklearn.cluster import KMeans</pre>
In [2]: In [3]: Out[3]:	df.head() CustomerID Gender Age Annual Income (k\$) Spending Score (1-100) 1 Male 19 15 39 Male 21 15 81
In [5]: Out[5]:	(200 5)
In [6]:	df.describe()
	min 1.000000 18.000000 15.000000 1.000000 25% 50.750000 28.750000 34.750000 50% 100.500000 36.00000 50.000000 75% 150.250000 49.000000 73.000000 max 200.000000 70.000000 99.000000
In [7]:	<pre>df.info() <class 'pandas.core.frame.dataframe'=""> RangeIndex: 200 entries, 0 to 199 Data columns (total 5 columns): # Column Non-Null Count Dtype</class></pre>
In [8]:	2 Age 200 non-null int64 3 Annual Income (k\$) 200 non-null int64 4 Spending Score (1-100) 200 non-null int64 dtypes: int64(4), object(1) memory usage: 7.9+ KB #Checking for null values df.isnull().sum()
Out[8]: In [18]:	Age 0 Annual Income (k\$) 0 Spending Score (1-100) 0 dtype: int64 # Perform EDA on the data print(df.describe())
	df.hist(bins=50, figsize=(20,15)) plt.show() CustomerID Age Annual Income (k\$) Spending Score (1-100) count 200.000000 200.000000 200.000000 200.000000 mean 100.500000 38.850000 60.560000 50.200000 std 57.879185 13.969007 26.264721 25.823522 min 1.000000 18.000000 15.000000 15.000000 25% 50.750000 28.750000 41.500000 34.750000 50% 100.500000 36.000000 61.500000 50.000000
	35% 150.35000 49.00000 78.00000 73.0000000 73.000000 73.000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.0000000 73.000000 73.0000000 73.0000000 73.000000 73.0000000 73.000000 73.000000 73.0000000 73.0000000 73.000000 73.000000 7
	3.0 2.5 2.0 1.5
	0 25 50 75 100 125 150 175 200 20 30 40 50 60 70 Annual Income (k\$) Spending Score (1-100)
	15.0 12.5 10.0
	7.5 5.0 2.5
In [19]: Out[19]:	40-0 home and a model Forest Omid at 10,004 45 d55 d0 a Ov
	40 30 75
	10
In [20]: Out[20]:	40-0 h ann an i anni d. Farat Oni d. at. On 04 45 d 0 0 0 0 d O
	35 30 25 10 20
	8 15 10 5
In [21]: Out[21]:	complere evicerid Foodforid at Ov214fbCo0f70>
In [22]: Out[22]:	Spending Score (1-100) sns.catplot(data=df, x="Gender", y="Annual Income (k\$)", kind="bar") sns.catplot(data=df, x="Gender", y="Annual Income (k\$)", kind="bar")
	60 50 60 60 60 60 60 60
	20 10 Male Female Gender
In [23]: Out[23]:	Gender sns.catplot(data=df, x="Gender", y="Spending Score (1-100)", kind="bar")
	00 Soore (1-100)
	10 Male Female Gender
In [24]: Out[24]:	sns.scatterplot(data=df, x="Age", y="Spending Score (1-100)", hue="Gender") <pre> </pre> <pre> <axessubplot:xlabel='age', ylabel="Spending Score (1-100)"> </axessubplot:xlabel='age',></pre> <pre> Gender Male Male Formule Formule </pre>
	So Spending Soor (1-10) 60 20 0
In [25]: Out[25]:	20 30 40 50 60 70 Age sns.scatterplot(data=df, x="Age", y="Annual Income (k\$)", hue="Gender") <axessubplot:xlabel='age', ylabel="Annual Income (k\$)"> 140 Gender</axessubplot:xlabel='age',>
	Gender Male Female
In [26]: Out[26]:	20 30 40 50 60 70 sns.scatterplot(data=df, x="Annual Income (k\$)", y="Spending Score (1-100)", hue="Gender")
Out[26]:	100 (001-1) 00 (001-1)
To [40].	20 40 60 80 100 120 140 Annual Income (k\$)
In [10]:	<pre>#Now through dataset as we can see that customer's Annual income and spending score decides if they are our customer or not X = df.iloc[:,[3,4]].values print(X) [[15 39] [15 81] [16 6] [16 77] [17 40] [17 70] [17 70] [18 70] [19 70] [19 70] [10</pre>
	[17 76] [18 6] [18 94] [19 3] [19 72] [19 14] [19 99] [20 15] [20 77] [20 13]
	[20 79] [21 35] [21 66] [23 29] [23 98] [24 35] [24 73] [25 73]
	[28 14] [28 82] [28 32] [28 61] [29 31] [29 87] [30 4] [30 73] [33 4]
	[33 92] [33 14] [33 81] [34 17] [34 73] [37 26] [37 75] [38 35] [38 92]
	[39
	[42 52] [42 60] [43 54] [43 60] [43 45] [43 41] [44 50] [44 46] [46 51]
	[46 46] [46 56] [46 55] [47 52] [47 59] [48 51] [48 59] [48 50] [48 48]
	[48 59] [48 47] [49 55] [49 42] [50 49] [54 47] [54 54] [54 53]
	<pre>[54 48] [54 52] [54 42] [54 51] [54 55] [54 41] [54 44] [54 57] [54 46]</pre>
	[57 58] [57 55] [58 60] [58 46] [59 55] [59 41] [60 49] [60 42]
	[60 52] [60 47] [60 50] [61 42] [61 49] [62 41] [62 48] [62 59]
	 [62 56] [62 42] [63 50] [63 46] [63 48] [63 52] [63 54] [64 42] [64 46]
	[65
	<pre>[69 91] [70 29] [70 77] [71 35] [71 11] [71 75] [71 9]</pre>
	[71 75] [72 34] [72 71] [73 5] [73 88] [73 7] [73 73] [74 10] [74 72]
	[74 72] [75 5] [75 93] [76 40] [76 87] [77 12] [77 97] [77 36] [77 4] [78 22] [78 90]
	[78
	[78 73] [79 35] [79 83] [81 5] [81 93] [85 26] [85 75] [86 20] [86 95]
	[87 27] [87 63] [87 13] [87 75] [87 10] [88 13] [88 86] [88 15]
	[88 69] [93 14] [93 90] [97 32] [97 86] [98 15] [98 88] [99 39] [99 97]
	[101 24] [101 68] [103 17] [103 85] [103 23] [103 69] [113 8] [113 91] [120 16]
In [11]:	[120 79] [126 28] [126 74] [137 18] [137 83]] : #Choosing the number of Clusters #Clustering with respect to Spending score and annual income wcss = []
	##wcsss = Within Clusters Sum of Squares for i in range(1,11): kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans.fit(X) wcss.append(kmeans.inertia_) C:\Users\Atharva Yadav\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available thre ads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1. warnings.warn(
In [12]:	<pre>sns.set() plt.plot(range(1,11) , wcss) plt.title('The Elbow Point Graph') plt.xlabel('Number of Clusters') plt.ylabel('WCSS') plt.show()</pre>
	The Elbow Point Graph 250000 200000 88 150000
	10000 2
In [15]:	#Training the Kmeans Clustering Model kmeans = KMeans(n_clusters = 5 , init = 'k-means++' , random_state=0) Y = kmeans.fit_predict(X) print(Y) [4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3 4 3
In [16]:	0 2 0 2 0 2 0 2 0 2 0 2 0 2 0 2 0 2 0 2
	<pre>plt.scatter(X[Y==3,0], X[Y==3,1], s=50, c='violet', label='Cluster 4') plt.scatter(X[Y==4,0], X[Y==4,1], s=50, c='blue', label='Cluster 5') # plot the centroids plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], s=100, c='cyan', label='Centroids') plt.title('Customer type groups') plt.xlabel('Annual Income') plt.ylabel('Spending Score')</pre>
	plt.ylabel('Spending Score') plt.show() Customer type groups 80
	nding Score
	20
In [27]:	20 40 60 80 100 120 140 #clustering using age and spending score X1 = df.iloc[:,[2,4]].values wcss1 = []
T⊬	<pre>for i in range(1,11): kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans.fit(X1) wcss1.append(kmeans.inertia_) C:\Users\Atharva Yadav\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available thre ads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1. warnings.warn(#plotting an elbow graph</pre>
ыг [29]:	<pre>sns.set() plt.plot(range(1,11) , wcss1) plt.title('The Elbow Point Graph') plt.xlabel('Number of Clusters') plt.ylabel('WCSS') plt.show()</pre> The Elbow Point Graph
	160000 140000 120000 80000
Tr	60000 20000 2 4 6 8 10 Number of Clusters #Training the Kmeans Clustering Model
_п [31]:	<pre>kmeans = KMeans(n_clusters = 5 , init = 'k-means++' , random_state=0) Y1 = kmeans.fit_predict(X1) print(Y1) [0 1 3 1 0 1 3 1 3 1 3 1 3 1 3 1 4 0 4 1 4 1 3 1 3 1 4 0 4 1 3 1 3 1 3 1 3 1 3 1 3 1 2 1 4 0 4 0 2 0 0 0 0 2 0 0 2 4 4 2 2 0 2 2 0 2 2 2 0 4 2 0 0 2 2 4 2 2 2 0 4 4 0 4 2 0 2 4 0 4 2 0 0 4 2 0 4 4 0 0 4 0 0 4 2 0 2 0</pre>
In [33]:	3 1 3 1 3 1 3 1 3 1 4 1 3 1 4 1 3 1 4 0 3 1 3 1 3 1 3 1 3 1 4 1 3 1 4 1 3 1 3 1 3 1 3 1 3 1 4 1 3 1] #Visualizing the Clusters plt.figure(figsize=(8,8)) plt.scatter(X1[Y1==0,0], X1[Y1==0,1], s=50, c='green', label='Cluster 1') plt.scatter(X1[Y1==1,0], X1[Y1==1,1], s=50, c='red', label='Cluster 2') plt.scatter(X1[Y1==2,0], X1[Y1==2,1], s=50, c='yellow', label='Cluster 3') plt.scatter(X1[Y1==3,0], X1[Y1==3,1], s=50, c='violet', label='Cluster 4')
	<pre>plt.scatter(X1[Y1==3,0], X1[Y1==3,1], s=50, c='violet', label='Cluster 4') plt.scatter(X1[Y1==4,0], X1[Y1==4,1], s=50, c='blue', label='Cluster 5') # plot the centroids plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], s=100, c='cyan', label='Centroids') plt.title('Customer type groups') plt.xlabel('Age') plt.ylabel('Age') plt.ylabel('Spending Score') plt.show()</pre>
	Customer type groups 80
	80 Source 40
	20
In [47]:	20 30 40 50 60 70 Age
	<pre>#clustering using age and spending score X3 = df.iloc[:,[2,3]].values wcss2 = [] for i in range(1,11):</pre>
т.	X3 = df.iloc[:,[2,3]].values wcss2 = []
In [51]:	X3 = df.iloc[;,[2,3]].values wcss2 = [] for i in range(1,11): kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans = KMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans =
In [51]:	X3 = df.iloc[;[2,3]].values Wcss2 = [] for i in range(1,11): kmeans = kMeans(n_clusters= i , init = 'k-means++' , random_state=0) kmeans.sit(X3) wcss2.append(kmeans.inertia_) C:\Users\Atharva Yadav\anaconda3\lib\site-packages\sklearn\cluster_kmeans.py:1036: UserWarning: KMeans is known to have a memory leak on Windows with MKL, when there are less chunks than available thre ads. You can avoid it by setting the environment variable OMP_NUM_THREADS=1. #plotting an elbow graph sns.set() plt.plot(range(1,11) , wcss2) plt.title('The Elbow Point Graph') plt.ylabel('Wcss')
	X3 = df.iloc(; [2,3]) values vesS2 = [] for i in range(d,11):
	X3 = df. lloc[:[P,3]] values vex.V = [Ver vex.V = [Ver vex.V = [Vex.V = Vex
In [52]:	A
In [52]:	State of the (the (the (the (the (the (the (the
In [52]:	Has eff time: (1.5.) (1.5.) (value) of the results
In [52]:	More and control of the first control of the contro
In [52]:	The first control (1) file to the control (1) file to
In [52]:	So and Land (1921) whose If you can improve the control of the co