•	 1.Supervised learning 2.Unsupervised Learning 3.Semi-supervised learning. In the supervised machine learning the target variable or label is defined. In the supervised machine learning the label is not defined so for getting the label on the group of the data we have to do the clustring of sample data and label each cluster with some numerical.like 0,1,2 and 3 so on.
→	B
,	In this kind of unsupervied machine learning technique we will get the similar kind of observation in individual cluster whereas the data point's behaviour or characteristics from one cluster to another cluster are getting changes. It means each and every cluster data having different characteristics and nature. If we plot the scatter plot in that we make the custering(grouping) of similar kind of object or attributes those are differ from the other attributes. Moreover we can say here that It is an amount that reflect the straigth of the relationship between two data object or two classes.
,	 So,the clustering mainly use for the exploratory data mining and it used in many field such as machine learning,pattern recognition,image analysis,information retrival,bio informatics, data compression, computer graphics. Moreover the k-means algorithm identified that k number of centroids for example if we having the 3 grouping that means we will have the 3 numbers of centroid. Centroid is the each group (cluster) is the mean of that group (cluster). It start from the random selection of the datapoints and then it will start to conversing the final clustered one by one iteration.
•	 In the K-Means clustering we do have data space inside that we have so many data points. While,First iteration the it will select random centroid like if we want to make the binary classification that time the centroid will be 2. Then it measure the distance from the that random centroid that distance might be euclidean or manhattan distance and whichever the datapoints come close to the centroid which are belonging to that particular centroid. Then again based on the nearest data points It will calculate the mean and consider than mean as next centroid. It is going happen for the some iteration until and unless the we will get the fix number of group and no point movement will there during mean is getting shited inside the data space at each iteration.
H	ypes Of Clustering. ard Clustring: Each data object or classes or point either belongs to cluster completly or not. In this kind of clustering points are clustered in such manner so that all attributes will belongs to the only same cluster and others. It means that both group will be the diffrent. It will never consider as part of any other cluser. Means they are perfectly classified into their classes. Oft Clustering:-
O:	 But in soft clustering the classification it will so strongs some points will be associated with one cluster or many cluster with some probability those neibour it or nearer to it. ther Types Of the Clustring Algorithm. Connectivity Based clustering. Centroid Based clustering. Distribution Based clustering. Density Based clustering.
	Clustering • Centroid clustering (Partitioning algorithms): Construct various partitions and then evaluate them by some criterion
	 Connectivity-based clustering (Hierarchical algorithms): Create a hierarchical decomposition of the set of data (or objects) using some criterion Density-based clustering: based on connectivity and density functions Model-based clustering: A model is hypothesized for each of the clusters and the idea is to find the
	onnectivity Based Clustering: The main idea behind this clustering is that data points that are the closer in the data space are more related (similar) than the to data points farther away. They are not very roubust towards outliers, which might show up as additional clusters or even cause other cluster to the merge. Here cluser form according to their diistances and datapoints will being connected to each are away from the centroid.
	 The both side of the dataspee there could be formation of cluster and these kind of cluster connection is known as dindogram. This is also commonly called the Heirarchical clustering. This method do not produce unique partitioning of data rather heirarchi from which the user still needs to choose appropriate cluster by using the label where they want to cluster. Agglomerative: This is a "bottom-up" approach: each observation starts in its own cluster, and pairs of clusters are merged as one moves up the hierarchy. Divisive: This is a "top-down" approach: all observations start in one cluster, and splits are performed recursively as one moves down the hierarchy. In order to decide which clusters should be combined (for agglomerative), or where a cluster should be split (for divisive), a measure of dissimilarity between sets of observations is required. In most methods of hierarchical clustering, this is achieved by use of an appropriate metric (a measure of distance between pairs of observations), and a linkage criterion which specifies the dissimilarity of sets as a function of the pairwise distances of observations in the sets.
Di	 In this kind of clustering, the cluster are represented by central vector or a centroid. This centroid might not necessarily be a member of datasets. This is an ittrerative cluster algorithm which is the motion of the similarity derive by how close data point is to be centroid of cluster. istribution Based Clustring. This model having the strong therotical foundation, however they ofter suffer for the overfitting. Gaussian mixture models, using expectation-maximization algorithm is the fomous distribution based clustring method. This is clustring is closed to the stastical distribution modeling. Datapoints are belongs the same distribution and model have strong therotical foundation however they often suffer from overfitting.
In [8]:	 Density based methods search the data space for the area of varied density of data points. Cluster will define the area of highly density within the dataspace to region. Vorking With K-Means Clustering. import numpy as np import pandas as pd import matplotlib.pyplot as plt import seaborn as sns
In [11]: Out[11]:	#for indexing the columns used index_col whenever uname is avalible that time use the index_col =0 data = pd.read_csv('data.csv',index_col=0) data.head() x y cluster
In [12]: Out[12]: 1 0 2	-11.999447 -7.606734 2 -1.736810 10.478015 1 #to know about the how many cluster avalible data['cluster'].value_counts() 67 67 66 ame: cluster, dtype: int64
In [13]:	Please check this to get more about the cmap:-https://matplotlib.org/3.1.0/tutorials/colors/colormaps.html Here if we can see that there is no any perfect the clustering has been done due to some points these are the yellow and voilet #c based on the clustering #cmap (colormap) in the python plt.scatter(data['x'], data['y'], c=data['cluster'], cmap='viridis') plt.show()
-	5- 051012.5 -10.0 -7.5 -5.0 -2.5 0.0 2.5 5.0
In [15]:	<pre>from sklearn.cluster import KMeans from sklearn.preprocessing import StandardScaler • We had created the vector of x and y • We had taken out the two columns x and y from the data datasets. • In unsupervised machine learning we don't need any labeling but for comparision we will using this final output y as well. #We had created the vector of x and y x = data[['x','y']] #we had taken out the two columns x and y from the data datasets. y = data['cluster'] #In unsupervised machine learning we don't need any labeling but for comparision we will to the comparision we will to th</pre>
In [16]:	we see here x is neither standardized or normalized #If we see here x is neither standardized or normalized x y 0 -8.482852 -5.603349 1 -7.751632 -8.405334 2 -10.967098 -9.032782
1 1	3 -11.999447 -7.606734 4 -1.736810 10.478015 95 -8.820126 -9.479259 96 -1.573419 -6.650994 97 -2.619581 8.269253 98 -2.634418 6.697531
20 Fin In [18]:	99 -11.951634 -7.121327 00 rows × 2 columns rst we need to stardardized and normalized the data by using Standardscaler. scaler = StandardScaler() x = scaler.fit_transform(x) x rray([[-1.01200363, -0.60606415],
	[-1.5097118 , -1.14041707], [-1.71653856, -0.91821912], [0.33953731,
	[1.5381542 , -0.52935538], [0.27161133, -0.71651135], [-1.35289151, -0.93641563], [0.79308339, -0.35212834], [0.14633873, 2.31330628], [1.14228061, -0.30354302], [1.12503491, -0.04177872], [-1.17977302, -1.09578279], [-2.01238779, -0.86540408], [0.19152374, 1.57228952], [0.83261025, 1.80799546], [0.66508288, -0.6110459], [-1.12827817, -0.15250042], [0.52241815, -0.65585313],
	[0.32241613, -0.63383313], [1.08583214,
	[0.65096239, -0.19455043], [-1.39325586, -1.1523753], [-1.49484828, -0.92704231], [-1.45695253, -0.70502951], [0.5305949,
	[-2.03584663, -0.69245567], [0.59854512,
	[-1.97542172, -0.86113934], [0.55709844,
	[1.59514484, -0.57810326], [0.15529192, -0.46145428], [1.07078756, -0.73342709], [-1.84390846, -0.81992791], [0.90990583, 1.40916629], [0.3576306, -0.73136536], [1.29042424, -0.89780574], [-1.36084266, -0.77007109], [-1.0236826, -1.24868072], [0.7946996, -1.12120567], [0.79884711, 1.66676256], [0.99941321, 1.13889426], [-1.36837982, -1.15198567], [-0.3156759, -1.46559894],
	[0.70861868, -0.75627388], [-0.21082427, -0.41430294], [0.32589705,
	[0.24621536, -0.32010977], [0.53736804, 1.0345916], [1.21425236, 1.07790672], [-0.00967353, 1.27374229], [0.53207638, 1.10912864], [-2.00277598, -1.07059812], [1.33013807, -0.90971954], [-1.25463745, -1.18101253], [0.53775645, 1.29669311], [0.27952731, -0.01116382], [0.32201814, 1.62673893], [0.4329941, -0.57783171], [0.57365299, 1.3368478], [0.65082523, 1.46367753],
	[0.65445327, 1.46194122], [0.30583921, 0.15038392], [0.58866226, -0.19466067], [0.77058657, -0.65992081], [-1.73473402, -0.85951427], [1.21530606, -0.27942699], [0.37146484, 1.57464784], [0.13101932, 1.17557222], [0.55048135, -0.073215], [0.65268228, -0.76005108], [-0.12524786, 1.16172152], [0.9416221, -0.8284188], [0.72962605, -0.78012446], [-0.7595301, -1.02429527],
	[0.42537401,
	[0.23509561,
	[-1.77472219, -0.84542156], [0.8071665 ,
	[-1.39989229, -0.77754158], [0.94719192, -0.41984737], [-2.06564945, -0.83849631], [-0.71883787, -0.62819834], [-1.17765783, -1.01506935], [0.70341832, -0.73729614], [-1.26359243, -0.85907062], [-1.08899436, -0.88801288], [1.132047, -0.55152432], [0.52620981, 1.43321291], [0.77533462, -0.4080321], [0.77533462, -0.82445019], [1.43342856, -0.2291756], [0.72870678, 1.06563353],
In [20]: W	[-1.07957503, -1.20998434], [0.37227198, -0.76930174], [0.16267793,
Out[21]: K In [22]:	<pre>#we already fitted our data it's means training is done. kmeans = KMeans(n_clusters =2,random_state =42) kmeans.fit(x) Means(algorithm='auto', copy_x=True, init='k-means++', max_iter=300,</pre>
In [23]:	0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
In [24]:	Now, do the scatter plot and check the how the data being clustered. In the above the code we had used the 2 cluster. If we use he scatter with the 3 cluster that time we will get the 3 cluster in the sactter plot. #Now, do the scatter plot and check the how the data being clustered #In the above the code we had used the 2 cluster .If we use he scatter with the 3 cluster that time we will get plt.scatter(data['x'],data['y'],c=kmeans.labels_,cmap='viridis') plt.show()
-	2.0 - 1.5 - 1.0 - 0.5 0.0 0.5 1.0 1.5 - 1.0 0.0 0.5 1.0 1.5 - 1.0 0.0 0.5 1.0 1.5 - 1.0 0.0 0.5 1.0 1.5 - 1.0 0.0 0.5 1.0 1.0 1.5 - 1.0 0.0 1.
	<pre>kmeans = KMeans(n_clusters = 3, random_state = 42) kmeans.fit(x) kmeans.labels_ kmeans.cluster_centers_ #In the above the code we had used the 2 cluster .If we use he scatter with the 3 cluster that time we will get plt.scatter(data['x'],data['y'],c=kmeans.labels_,cmap='viridis') plt.show()</pre> 25 20 15
-	Now if you see the cluster are classified according to the our convenience of perfect classification.
- Yo	 Now if you see the cluster are classified according to the our convenience of perfect classification. But how do I cloose the right value of the n_cluster value. Then we can use the n_value whatever we want it it will not always give us the better clustering over the scatter plot so So,we can get the n_value of the cluster so that the make the data clustring precisely. There is a techinique for this one. How do choose the right value of K? Assuming we having input X_1,X_2,X_3,,X_nx
R	 Assuming we having input X_1,X_2,X_3,,X_nx Pick k-random points are cluster centers called centroid and at the meantime calculate the mean square error or sum of square error (SSE). For every different value of k we need to calculate the sum of squared error. Assign the each x_i to the nearest cluster by calculating its distance to each centroid. Finally new cluster by taking the avarage of the assigned points. Repeat the step 2 and 3 until none of cluster assignments change. elationship between the sum of squared error and k-mean value Let's consider that we having the graph of the k values against SSE (sum of squared error).
F	 When the value of k is increasing that time the SSE(sum of squared error) start Decreasing. After the some point the The decresing of the SSE is becoming the low and It's slop is become is very low and horizontal line having the 0 slop. This is konw as elbow method. Because this graph having the shape of elbow. This graph can able to help to get the perfect value of the k at this condition the value of the sse(sum of squared error) is become low and at the same points the k value is to be consider as perfect value of k. or that we need to find the value of the SSE (sum of squared error).
	SSE = Sum Of Squared Error empty list has been created. Make the variable as index in that we had passed the range from the 1 to 10.Beacsue we want to iterate the k value from the 1 to 10. Create the kmeans variable in that we had put attributes that are the n_clusters and random state with the help of KMeans. Make the fit on the x to train it. Append the values of kmeansineria in the empty list i.e.SSE kmeansineria = It's helps to get the values of the sum squared errors inside the list of the kmeans from the from the range 1 to 10. Then the print the values of the kmeans.inertia It will help us give it back the sum of square error.
In [26]:	 If we kindly obesrved that the sum of square errors are decresing from the 400 to 156 then to 44 but after the 44 it will not decresing if you observe the previous decreasing order. The decrement is going foreword not going the highly decrement. Decrement strat from the 400 to 156. Decrement done succefully upto 44 but after that the there so much diffrences in the declination SSE (sum of square error).
1 2 3 4 5 6	<pre>for i in index: kmeans = KMeans(n_clusters=i,random_state=42) kmeans.fit(x) SSE.append(kmeans.inertia_) print(i,kmeans.inertia_) 400.00000000000001 156.41078579574986 44.05704845329278 36.72638711866608 31.016427615314644 25.39959047192452 22.547184727743293</pre>
In [30]:	19.923438178477447 17.295836404824975 • We have to plot the k vs SSE. • If we can check it then we will get to know that the 3 is the k value over the declination SSE is lowest declination point over the SSE slop is being low and giving the less diffrences after the k=3. • So,We have to consider the out appropriate value of k as 3. • plt.plot(index, SSE) plt.show
4 3 3 2 2	<pre><matplotlib.lines.line2d 0xf81e410="" at="">] 000 - 00</matplotlib.lines.line2d></pre>
	50 1 2 3 4 5 6 7 8 9

K-Means Clustering Algorithm.

Machine learning can braodly be classified three types :-