In [7]:	Clustering  from IPython.display import Image  What is Clustering?
In [8]:	Clustering is a type of unsupervised machine learning, which means the actual "y" (target variable) is not given when traning the model. By seperating the datapoints based on their structures or similarities, we're able to create several groups, called cluster.  Image ("images/clusterpng")
Out[8]:	Cluster 3  Cluster 4  Cluster 4
	(a) Data objects  (b) Clustered data objects  How to implement Clustering in python?  To impletment Clustering, we actually have developed a lot of methods/algorithms. For this project, we will use k-means to achieve our goal. The basic idea is below:
	<ul> <li>Step1: Select k unique points from datapoints as initial centroids (k is the number of clusters we want to generate)</li> <li>Step2: Put each datapoints into the k clusters based on the minimum distance to each centroids</li> <li>Step3: Calculate the mean of each cluster</li> <li>Step4: Repeat Step2 and Step3 so that the mean of each cluster = previously chosen centroids</li> </ul>
In [9]: Out[9]:	Image ("images/cluster.png")  Before K-Means  After K-Means
	Apply "kmeans++" initialization
	When we only use k-means function, we initialize the centroids randomly, which may not generate accurate resulsts when initialization is bad. Thus, in order to improve our model, we can work on a funciton called "select_centroids". This function function as "kmeans++" in scikit-learn. The basic algorightm is below:  • Step1: Randomly select the first centroid from the data points.
	<ul> <li>Step2: For each data point compute its distance from the nearest, previously chosen centroid.</li> <li>Step3: Select the next centroid from the data points such that the probability of choosing a point as centroid is directly proportional to its distance from the nearest, previously chosen centroid. (i.e. the point having maximum distance from the nearest centroid is most likely to be selected next as a centroid)</li> <li>Step4: Repeat steps 2 and 3 until k centroids have been sampled</li> </ul>
n [11]: nt[11]:	<pre>Image("images/++.jpeg")</pre>
	-6 -8
ı [187	-4 -2 0 2 4 6 8
	Try Our k-means on One-Dimensional Data  Suppose we have several grades, and we want to seperate them into several clusters. We want to seperate these groups so that we can have a better understanding of the grades distribution.
[261	<pre>grades = [92.65, 93.87, 74.06, 86.94,92.26, 94.46, 92.94, 80.65, 92.86, 85.94, 91.79, 95.23, 85.37, 87.85, k = 3 grades = np.array(grades).reshape(-1,1) centroids = select_centroids(grades,3) centroids,labels = kmeans(grades,k) print("centroids:", centroids.reshape(1,-1))</pre>
	<pre>print("labels for each x: ", labels) for j in range(k):     print("vector assignments: ",grades[labels==j].reshape(1,-1))  centroids: [[93.23222222 86.762</pre>
[19]: at[19]:	<pre>vector assignments: [[86.94 85.94 85.37 87.85 87.71]] vector assignments: [[74.06 80.65]]  Image("images/one_dim.png")</pre>
	From the graph above, we can see that the "grades" are seperated into three clusters. In a practical perspective, we can use k-
	Try Our k-means on Two-Dimensional Data  Suppose we have a two-dimentional dataset and we want to seperate them into two clusters.
[329	<pre>from sklearn.datasets import make_circles import matplotlib.pyplot as plt  X, _ = make_circles(n_samples=500, noise=0.1, factor=.2) centroids = select_centroids(X,2) #print(centroids) centroids, labels = kmeans(X, 2, centroids=centroids)</pre>
	<pre>colors=np.array(['#4574B4','#A40227']) plt.scatter(X[:,0], X[:,1], c=colors[labels]) plt.show()</pre>
	0.5
	From the graph shown above, kmeans performs poorly on disjoint and nested structures because we separate each point based on
[327	their distance. It's only good at seperating datapoints that are clustered togerther. If the datapoints cluster in a more complex pattern, k-means cannot generate good results on that.  Try Our k-means on Multi-Dimensional Data
	<pre>from sklearn.datasets import load_breast_cancer from sklearn.preprocessing import StandardScaler from sklearn.metrics import confusion_matrix from sklearn.metrics import accuracy_score import seaborn as sns cancer = load_breast_cancer() X = cancer.data y = cancer.target #y=0 means cancer, y=1 benigh</pre>
	<pre>sc = StandardScaler() X = sc.fit_transform(X)  centroids = select_centroids(X, k=2) centroids, labels = kmeans(X, k=2, centroids = centroids, tolerance= 0.01)  if len(X[labels == 1]) &lt; len(X[labels == 0]):</pre>
[14]: t[14]:	<pre>labels = [1-i for i in labels]  Image("images/confusion_matrix.png")  Confusion Matrix</pre>
	-300 -250
	- 200 - 150 - 100 - 50
	False True  Predicted Values  From the confusion matrix above, we can see that k-means is good at clustering binary problem.
. [297	Try Our k-means on Image Processing  Below is a very famous picture - Girl with a Pearl Earring by Dutch Golden Age painter Johannes Vermeer. K-means can generate some interesting artificial effect of the picture by clustering the color of the image. Rather than use millions of colors, we can usuall get away with 256 or even 64 colors by setting different k in our k-means function. Let's start with grey_scaled picture first.
	<pre>import ev2 import warnings warnings.filterwarnings("ignore") image = cv2.imread("pearl_grey.jpg") image = cv2.cvtColor(image, cv2.COLOR_BGR2RGB) plt.imshow(image)</pre>
t[297	100
	200 - 300 - 400 -
[298	<pre>pixel_values = Image.resnape((-1, 4)) pixel_values = np.float32(pixel_values) k=4</pre>
	<pre>centroids = select_centroids(pixel_values, k=k) centroids, labels = kmeans(pixel_values, k=k, centroids=centroids, tolerance=.01) centroids = centroids.astype(np.uint8) segmented_image = centroids[labels.flatten()] segmented_image = segmented_image.reshape(image.shape) plt.imshow(segmented_image) plt.show()</pre>
	200 -
	300 - 400 -
[299	From the graph shown above, we can see the picture got less color, since the original pixels are replaced by the centroids. Original numbers of different colors are reduced to k, which is 4. That results in the change on the picture above.
t[299	<pre>image = cv2.imread("pearl.jpg") image = cv2.cvtColor(image, cv2.COLOR_BGR2RGB) plt.imshow(image)  <matplotlib.image.axesimage 0x7f954c52af40="" at=""></matplotlib.image.axesimage></pre>
	200 - 400 - 600 -
	1200 - 1400 - 0 200 400 600 800 1000
[302	<pre>pixel_values = image.reshape((-1, 30)) pixel_values = np.float32(pixel_values) k=30 centroids = select_centroids(pixel_values, k=k) centroids, labels = kmeans(pixel_values, k=k, centroids=centroids, tolerance=.01) centroids = centroids.astype(np.uint8) segmented image = centroids[labels.flatten()]</pre>
	<pre>segmented_image = segmented_image.reshape(image.shape) plt.imshow(segmented_image) plt.show()</pre>
	400 - 600 - 800 -
	1200 - 1400
	Spectral clustering  As we mentioned above, kmeans performs poorly on disjoint and nested structures because we seperate each point based on their distance. You may wonder, how can we fix it? There's a model in sklearn called SpectralClustering, where clusters are grouped by based on the edges connecting them rather than their distances. The basic steps are shown below:
	<ul> <li>Take our graph and built an adjacency matrix</li> <li>Create the Graph Laplacian by subtracting the adjacency matrix from the degree matrix and calculate the eigenvalues of the Laplacian. The vectors associated with those eigenvalues contain information on how to segment the nodes</li> <li>Performe K-Means on those vectors in order to get the labels for the nodes</li> </ul>
[15]: t[15]:	<pre>Image("images/pic.png")  10 - 05</pre>
	0.5 - 0.0 - 0.5 - 0
[322	from sklearn.cluster import SpectralClustering X, _ = make_circles(n_samples=500, noise=0.1, factor=.2)
	<pre>X, _ = make_circles(n_samples=500, noise=0.1, factor=.2) cluster = SpectralClustering(n_clusters=2, affinity="nearest_neighbors") labels = cluster.fit_predict(X)  # pass X not similarity matrix  #print(labels) colors=np.array(['#4574B4','#A40227']) plt.scatter(X[:,0], X[:,1], c=colors[labels]) plt.show()</pre>
	10 - 0.5 -
	0.00.51.01.0
	Application Clustering is a widely-applied technique in the data science industry. Here are some examples:
[16]:	• Recommendation engines: In this method, the clustering method provided an idea of like-minded users. The computation/estimation as data provided by several users is leveraged for improving the performance of collaborative filtering methods. And this can be implemented for rendering recommendations in diverse applications.  Image ("images/rec.jpeg")
t[16]:	
[17]:	<ul> <li>Market and Customer segmentation: Clustering methods enable us to segment customers into diverse clusters, depending of which companies can consider novel strategies to apply to their customer base.</li> <li>Image ("images/m.png")</li> </ul>
t[17]:	
	Loyal Customers  Potential Customers
	"Shut up, and take my money!"  "I wanna buy something!"
	Middle class, 46-60 years old with small size family, open to online shopping.  "I wanna buy something!"  Upper class, German luxury car lovers, family originated from West German in the old time.
	• Search Result Clustering: Depending on the closest similar objects/properties, the data is assigned to a single cluster, giving the plethora sets of similar results of the users. In simple terms, the search engine attempts to group identical objects in one cluster and non-identical objects in another cluster.
n [18]: nt[18]:	Image ("images/search.jpg")  Google
	<ul> <li>search engine list</li> <li>top 50 search engines</li> <li>what is search engine in computer</li> <li>search engine examples</li> </ul>
	<ul> <li>google search engine</li> <li>best search engine</li> <li>top 20 search engines</li> <li>duckduckgo search engines</li> </ul>
	what is search engine search engine optimization Report inappropriate predictions
	Conclusion  Clustering is a very important method in Machine Learning. This report discussed one types of clustering - k-means and its initialization - kmeans++. K-means is good at group things that are continuous and visually seperated since it divides datapoints based on their distances. This report also mentions Spectral clustering, which group datapoints based on the graph structure.  References
	References  • https://towardsdatascience.com/spectral-clustering-aba2640c0d5b  • https://github.com/parrt/msds689/blob/master/projects/kmeans/kmeans.md#spectral-clustering