	The Spark Foundation-GRIP- Data Science and Business Analytics-Aug21 Task by: Deepti Taram¶ Language used: Python 3 IDE: Jupyter Notebook Task-2 Prediction using unsupervised Machine Learning • For the given Iris Dataset, predict the optimum number of clusters and present it visually We will using Kmean clustering for this Task.
In [1]	 K-Means¶ K-means is a centroid-based algorithm, or a distance-based algorithm, where we calculate the distances to assign a point to a cluster. In K-Means, each cluster is associated with a centroid # Importing the libraries import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns from sklearn import datasets from sklearn.cluster import KMeans # Load the iris dataset df = pd.read_csv('Iris.csv')
Out[2]	Id SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm Species
In [3]	148 149 6.2 3.4 5.4 2.3 Iris-virginica 149 150 5.9 3.0 5.1 1.8 Iris-virginica 150 rows × 6 columns df.drop('Id', axis=1, inplace= True) # As ID column will not help in clustering so we drop it. Explore the data df.shape
In [5]: Out[5]:	
Out[6]	SepalLengthCm SepalWidthCm PetalLengthCm Species 145 6.7 3.0 5.2 2.3 Iris-virginica 146 6.3 2.5 5.0 1.9 Iris-virginica 147 6.5 3.0 5.2 2.0 Iris-virginica 148 6.2 3.4 5.4 2.3 Iris-virginica 149 5.9 3.0 5.1 1.8 Iris-virginica
In [8]	Data columns (total 5 columns): # Column Non-Null Count Dtype
In [9]:	: SepalLengthCm 0
In [10]: Out[10]: In [11]:	PetalLengthCm
Out[11]: In [12]: Out[12]: In [13]: Out[13]: Out[14]:	4.5, 5.3, 7., 6.4, 6.9, 6.5, 6.3, 6.6, 5.9, 6., 6.1, 5.6, 6.7, 6.2, 6.8, 7.1, 7.6, 7.3, 7.2, 7.7, 7.4, 7.9]) : df['SepalWidthCm'].unique() : array([3.5, 3., 3.2, 3.1, 3.6, 3.9, 3.4, 2.9, 3.7, 4., 4.4, 3.8, 3.3, 4.1, 4.2, 2.3, 2.8, 2.4, 2.7, 2., 2.2, 2.5, 2.6]) : df['PetalWidthCm'].unique() : array([0.2, 0.4, 0.3, 0.1, 0.5, 0.6, 1.4, 1.5, 1.3, 1.6, 1., 1.1, 1.8, 1.2, 1.7, 2.5, 1.9, 2.1, 2.2, 2., 2.4, 2.3])
In [15]	4.6, 3.3, 3.9, 3.5, 4.2, 3.6, 4.4, 4.1, 4.8, 4.3, 5., 3.8, 3.7, 5.1, 3., 6., 5.9, 5.6, 5.8, 6.6, 6.3, 6.1, 5.3, 5.5, 6.7, 6.9, 5.7, 6.4, 5.4, 5.2]) : df['Species'].unique() : array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object) : df.corr() : SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm SepalLengthCm 1.00000 -0.109369 0.871754 0.817954 SepalWidthCm -0.109369 1.000000 -0.420516 -0.356544
In [17]:	sns.heatmap(df.corr(),annot=True) <pre> </pre> <pre> <pre> <pre> <pre></pre></pre></pre></pre>
	SepalWidthCm0.11
In [18]	<pre>df.plot(subplots=True)</pre>
	Finding the optimum number of clusters Before clustering the data using Kmeans, we need to specify the number of clusters in order to find the optimum number of clusters, there are various method avilable like silhouette coefficient and elbow method. Here the elbow method is used. Brief about Elbow method In cluster analysis, the elbow method is a heuristic used in determining the number of clusters in a data set. The method consists of plotting the explained variation as a function of the number of clusters, and picking the elbow of the curve as the number of clusters to use.
In [19]	<pre># Finding the optimum number of clusters for k-means classification x = df.iloc[:, [0, 1, 2, 3]].values from sklearn.cluster import KMeans k_meansclus = range(1,10) wcss = [] for k in k_meansclus: km = KMeans(n_clusters = k) km.fit(x) wcss.append(km.inertia_)</pre>
Out[20]	152.36870647733915, 78.94084142614601, 57.317873214285726, 46.535582051282034, 38.930963049671746, 34.190687924796634, 30.00336433172303, 27.977297315406027]
	plt.grid() plt.show() The elbow method 700 400 200 100
In [22]	You can clearly see why it is called 'The elbow method' from the above graph, the optimum clusters is where the elbow occurs. This is when the within cluster sum of squares (WCSS) doesn't decrease significantly with every iteration CONCLUSION - OPTIMUM NUMBER OF CLUSTERS IS 3 CREATING THE KMEANS CLASSIFIER
In [23]	<pre>y_kmeans = kmeans.fit_predict(x) y_kmeans array([1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1</pre>
In [24]	x array([[5.1, 3.5, 1.4, 0.2],
	[5.8, 4. , 1.2, 0.2], [5.7, 4.4, 1.5, 0.4], [5.4, 3.9, 1.3, 0.4], [5.1, 3.5, 1.4, 0.3], [5.7, 3.8, 1.7, 0.3], [5.1, 3.8, 1.5, 0.3], [5.1, 3.4, 1.7, 0.2], [5.1, 3.7, 1.5, 0.4], [4.6, 3.6, 1. , 0.2], [5.1, 3.3, 1.7, 0.5], [4.8, 3.4, 1.9, 0.2], [5. , 3. , 1.6, 0.2], [5. , 3.4, 1.6, 0.4], [5.2, 3.5, 1.5, 0.2], [5.2, 3.4, 1.6, 0.4], [5.2, 3.5, 1.5, 0.2], [5.2, 3.4, 1.6, 0.2], [5.3, 3.4, 1.6, 0.2], [5.4, 3.4, 1.6, 0.2], [5.4, 3.4, 1.6, 0.2], [5.4, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.2], [5.5, 3.4, 1.6, 0.4],
	[5.2, 4.1, 1.5, 0.1], [5.5, 4.2, 1.4, 0.2], [4.9, 3.1, 1.5, 0.1], [5. , 3.2, 1.2, 0.2], [5.5, 3.5, 1.3, 0.2], [4.9, 3.1, 1.5, 0.1], [4.4, 3. , 1.3, 0.2], [5.1, 3.4, 1.5, 0.2], [5.1, 3.4, 1.5, 0.2], [5. , 3.5, 1.3, 0.3], [4.5, 2.3, 1.3, 0.3], [4.4, 3.2, 1.3, 0.2], [5. , 3.5, 1.6, 0.6], [5. 1, 3.8, 1.9, 0.4], [4.8, 3. , 1.4, 0.3], [4.8, 3. , 1.4, 0.3], [5.1, 3.8, 1.6, 0.2], [4.6, 3.2, 1.4, 0.2], [5.3, 3.7, 1.5, 0.2],
	[5., 3.3, 1.4, 0.2], [7., 3.2, 4.7, 1.4], [6.4, 3.2, 4.5, 1.5], [6.9, 3.1, 4.9, 1.5], [5.5, 2.3, 4., 1.3], [6.5, 2.8, 4.6, 1.5], [5.7, 2.8, 4.5, 1.3], [6.3, 3.3, 4.7, 1.6], [4.9, 2.4, 3.3, 1.], [6.6, 2.9, 4.6, 1.3], [5.2, 2.7, 3.9, 1.4], [5., 2., 3.5, 1.], [5.9, 3., 4.2, 1.5], [6., 2.2, 4., 1.], [6.1, 2.9, 4.7, 1.4], [6.1, 2.9, 4.7, 1.4], [6.7, 3.1, 4.4, 1.4],
	[5.6, 3. , 4.5, 1.5], [5.8, 2.7, 4.1, 1.], [6.2, 2.2, 4.5, 1.5], [5.6, 2.5, 3.9, 1.1], [5.9, 3.2, 4.8, 1.8], [6.1, 2.8, 4., 1.3], [6.3, 2.5, 4.9, 1.5], [6.1, 2.8, 4.7, 1.2], [6.4, 2.9, 4.3, 1.3], [6.6, 3. , 4.4, 1.4], [6.8, 2.8, 4.8, 1.4], [6.7, 3. , 5. , 1.7], [6. , 2.9, 4.5, 1.5], [5.7, 2.6, 3.5, 1.], [5.5, 2.4, 3.8, 1.1], [5.5, 2.4, 3.8, 1.1], [5.5, 2.4, 3.7, 1.], [5.5, 2.4, 3.7, 1.],
	[6., 2.7, 5.1, 1.6], [5.4, 3., 4.5, 1.5], [6., 3.4, 4.5, 1.6], [6.7, 3.1, 4.7, 1.5], [6.3, 2.3, 4.4, 1.3], [5.6, 3., 4.1, 1.3], [5.5, 2.5, 4., 1.3], [5.5, 2.6, 4.4, 1.2], [6.1, 3., 4.6, 1.4], [5.8, 2.6, 4., 1.2], [5.8, 2.6, 4., 1.2], [5.7, 3., 4.2, 1.3], [5.7, 3., 4.2, 1.3], [5.7, 3., 4.2, 1.3], [5.7, 2.9, 4.2, 1.3], [6.2, 2.9, 4.3, 1.3], [6.2, 2.9, 4.3, 1.3], [6.2, 2.9, 4.3, 1.3], [5.7, 2.8, 4.1, 1.3],
	[6.3, 3.3, 6. , 2.5], [5.8, 2.7, 5.1, 1.9], [7.1, 3. , 5.9, 2.1], [6.3, 2.9, 5.6, 1.8], [6.5, 3. , 5.8, 2.2], [7.6, 3. , 6.6, 2.1], [4.9, 2.5, 4.5, 1.7], [7.3, 2.9, 6.3, 1.8], [6.7, 2.5, 5.8, 1.8], [7.2, 3.6, 6.1, 2.5], [6.5, 3.2, 5.1, 2.], [6.4, 2.7, 5.3, 1.9], [6.8, 3. , 5.5, 2.1], [5.7, 2.5, 5. , 2.], [5.8, 2.8, 5.1, 2.4], [6.4, 3.2, 5.3, 2.3], [6.5, 3. , 5.5, 1.8],
	[7.7, 3.8, 6.7, 2.2], [7.7, 2.6, 6.9, 2.3], [6., 2.2, 5., 1.5], [6.9, 3.2, 5.7, 2.3], [5.6, 2.8, 4.9, 2.], [7.7, 2.8, 6.7, 2.], [6.3, 2.7, 4.9, 1.8], [6.7, 3.3, 5.7, 2.1], [7.2, 3.2, 6., 1.8], [6.2, 2.8, 4.8, 1.8], [6.2, 2.8, 4.8, 1.8], [6.4, 2.8, 5.6, 2.1], [7.2, 3., 5.8, 1.6], [7.4, 2.8, 6.1, 1.9], [7.9, 3.8, 6.4, 2.], [6.4, 2.8, 5.6, 2.2],
	[6.3, 2.8, 5.1, 1.5], [6.1, 2.6, 5.6, 1.4], [7.7, 3., 6.1, 2.3], [6.3, 3.4, 5.6, 2.4], [6.4, 3.1, 5.5, 1.8], [6.9, 3.1, 5.4, 2.1], [6.7, 3.1, 5.6, 2.4], [6.9, 3.1, 5.4, 2.1], [6.9, 3.1, 5.1, 2.3], [6.9, 3.1, 5.1, 2.3], [6.8, 2.7, 5.1, 1.9], [6.8, 3.2, 5.9, 2.3], [6.7, 3.3, 5.7, 2.5], [6.7, 3., 5.2, 2.3], [6.7, 3., 5.2, 2.3], [6.3, 2.5, 5., 1.9], [6.3, 3.5, 5.2, 2.3], [6.3, 3.5, 5.2, 2.3], [6.3, 3.5, 5.4, 2.3], [5.9, 3.7, 5.1, 1.8]])
In [25]	<pre>plt.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1], s = 100, c = 'red', label = 'Iris-setosa') plt.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1], s = 100, c = 'blue', label = 'Iris-versicolour') plt.scatter(x[y_kmeans == 2, 0], x[y_kmeans == 2, 1], s = 100, c = 'green', label = 'Iris-virginica') # Plotting the centroids of the clusters plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:,1], s = 100, c = 'yellow', label = 'Centroids') plt.legend()</pre>
	4.0 - 3.5 - 3.0 - 2.5 - 2.0 - 2.5 - 4.5 5.0 5.5 6.0 6.5 7.0 7.5 8.0
In [26]	3D Scatterplot # 3D scatterplot using matplotlib fig = plt.figure(figsize = (15,15)) ax = fig.add_subplot(111, projection='3d') plt.scatter(x(y_kmeans == 0, 0], x(y_kmeans == 0, 1], s = 50, c = 'blue', label = 'Iris-setosa') plt.scatter(x(y_kmeans == 1, 0], x(y_kmeans == 1, 1], s = 50, c = 'orange', label = 'Iris-versicolour') plt.scatter(x(y_kmeans == 2, 0], x(y_kmeans == 2, 1], s = 50, c = 'green', label = 'Iris-virginica') #Plotting the centroids of the clusters plt.scatter(kmeans.cluster_centers_[:, 0], kmeans.cluster_centers_[:,1], s = 50, c = 'red', label = 'Centroids') plt.show()
	0.04
	0.02
	45 50 5.5 6.0 6.5 7.0 7.5 8.0 2.0
In [27]:	#1 to 'Iris-versicolour' #2 to 'Iris-virginica' y_kmeans = np.where(y_kmeans==0, 'Iris-setosa', y_kmeans) y_kmeans = np.where(y_kmeans=='1', 'Iris-versicolour', y_kmeans) y_kmeans = np.where(y_kmeans=='2', 'Iris-virginica', y_kmeans) Adding the prediction to the dataset
	data_with_clusters["Cluster"] = y_kmeans print(data_with_clusters.head(5)) SepalLengthCm
In [29]:	BARPLOT- CLUSTER DISTRIBUTION : #Bar plot sns.set_style('darkgrid') sns.barplot(x = data_with_clusters["Cluster"] .unique(),
	Bar Plot Inference There are around 62 iris-versicolour, 50 Iris-virginica and roughly 38 Iris-setosa samples in the dataset as predicted