

SECTION A

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

```
In [2]: 1 #Q1-Display the number of attributes available in the dataset
        2
        3 data = pd.read_csv('exam_dataset.csv')
```

```
In [3]: 1 data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 24 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Age                                   1470 non-null   int64
1   BusinessTravel                       1470 non-null   object
2   MonthlyIncome                       1470 non-null   int64
3   JobSatisfaction                     1470 non-null   int64
4   Bonus                               1470 non-null   int64
5   Department                           1470 non-null   object
6   DistanceFromHome                   1470 non-null   int64
7   Education                           1470 non-null   int64
8   EducationField                      1470 non-null   object
9   EmployeeCount                      1470 non-null   int64
10  EmployeeNumber                     1470 non-null   int64
11  EnvironmentSatisfaction             1470 non-null   int64
12  Gender                             1470 non-null   object
13  JobLevel                           1470 non-null   int64
14  JobRole                             1470 non-null   object
15  MaritalStatus                      1470 non-null   object
16  PerformanceRating                  1470 non-null   int64
17  StockOptionLevel                   1470 non-null   int64
18  TrainingTimesLastYear              1470 non-null   int64
19  WorkLifeBalance                    1470 non-null   int64
20  YearsAtCompany                     1470 non-null   int64
21  YearsSinceLastPromotion            1470 non-null   int64
22  OverTime                           1470 non-null   object
23  Attrition                          1470 non-null   object
dtypes: int64(16), object(8)
memory usage: 275.8+ KB
```

In [4]: 1 data.head()

Out[4]:

	Age	BusinessTravel	MonthlyIncome	JobSatisfaction	Bonus	Department	DistanceFromHome	Education	EducationField	EmployeeCount	...	JobRole	MaritalStatus	PerformanceRating	StockOptionLevel	Trainir
0	41	Travel_Rarely	5993	4	17979	Sales	1	2	Life Sciences	1	...	Sales Executive	Single	3	0	
1	49	Travel_Frequently	5130	2	20520	Research & Development	8	1	Life Sciences	1	...	Research Scientist	Married	4	1	
2	37	Travel_Rarely	2090	3	6270	Research & Development	2	2	Other	1	...	Laboratory Technician	Single	3	0	
3	33	Travel_Frequently	2909	3	8727	Research & Development	3	4	Life Sciences	1	...	Research Scientist	Married	3	0	
4	27	Travel_Rarely	3468	2	10404	Research & Development	2	1	Medical	1	...	Laboratory Technician	Married	3	1	

5 rows × 24 columns



In [5]: 1 #Q2-Find the dimension number of this dataset
2
3 data.shape

Out[5]: (1470, 24)

In [6]: 1 print("The dimension of dataset is: ",data.ndim)

The dimension of dataset is: 2

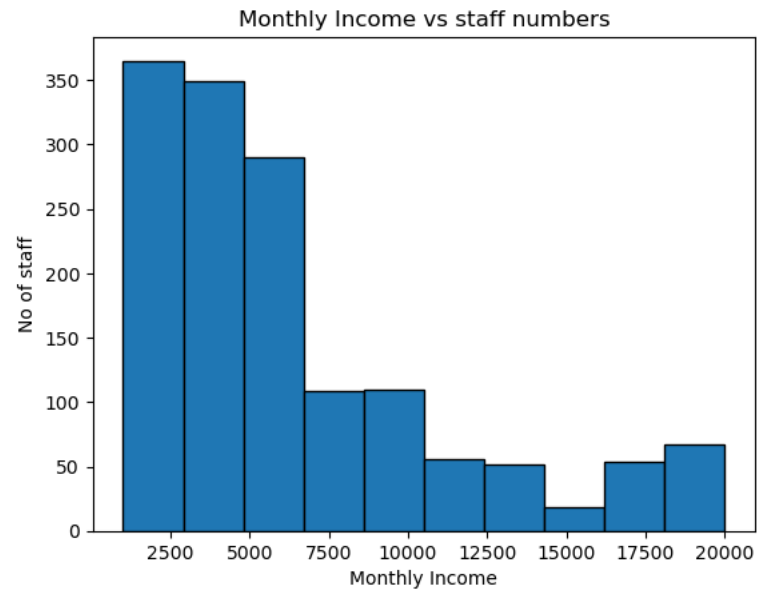
In [7]: 1 #Q3-Display the average of these attributes: 'Age', 'Monthly Income'
2 #and 'Years at Company'. Be sure to round your answer to 2 decimal places
3
4 age = np.average(data['Age'])
5 income = np.average(data['MonthlyIncome'])
6 years = np.average(data['YearsAtCompany'])
7
8 print("Average age is: ", round(age,2))
9 print("Average Monthly income is: ", round(income,2))
10 print("Average Years at company is: ", round(years,2))

Average age is: 36.92
Average Monthly income is: 6502.93
Average Years at company is: 7.01

In [8]: 1 #Q4-Find the minimum and maximum 'Monthly Income'
2
3 minMI = min(data['MonthlyIncome'])
4 print("Minimum income is: ", minMI)
5
6 maxMI = max(data['MonthlyIncome'])
7 print("Maximum income is: ", maxMI)

Minimum income is: 1009
Maximum income is: 19999

```
In [9]: 1 #Q5-Histogram of 'Monthly Income vs staff numbers'
2
3 plt.hist(data['MonthlyIncome'], edgecolor='black')
4 plt.title('Monthly Income vs staff numbers')
5 plt.xlabel('Monthly Income')
6 plt.ylabel('No of staff')
7 plt.show()
```



```
In [10]: 1 #Q6-Provide graphical Visualization of the distribution between
2 #‘Year at Company’ and ‘Monthly Income’ using the scatter plot
3
4 x = data.iloc[:, 20].values
5 y = data.iloc[:, 2].values
6
7 plt.scatter(x,y)
8 plt.title('Years at Company vs Monthly Income')
9 plt.xlabel('Years at Company')
10 plt.ylabel('Monthly Income')
11 plt.show()
```



```
In [11]: 1 #Q7-Find the correlation between ‘Years at Company’ and ‘Monthly Income’
2
3 cor = cor = data[['YearsAtCompany', 'MonthlyIncome']].corr()
4 cor
```

Out[11]:

	YearsAtCompany	MonthlyIncome
YearsAtCompany	1.000000	0.514285
MonthlyIncome	0.514285	1.000000

#Q8-

a) Range of Monthly Income at Company A is from 1009 to 19999

```
In [12]: 1 data.MonthlyIncome.value_counts()
```

```
Out[12]: 2342      4
        6142      3
        2741      3
        2559      3
        2610      3
        ..
        7104      1
        2773      1
        19513     1
        3447      1
        4404      1
        Name: MonthlyIncome, Length: 1349, dtype: int64
```

b) From the histogram, monthly income values at Company A: -Most Frequent monthly income values at Company A around +-2500 ;more specifically 2342 with frequency of 4 -Least Frequent monthly income values at Company A around +-15000

c) Histogram of Monthly Income shows a distribution which peaks in the lower income(up until 5000) , with average income of 6502.93.The monthly income is quite concentrated in the first quartile, and scattered in forth quartile.These observations suggest that most frequent monthly income is first quartile.

d) There are positive linear relationship with correlation of 0.51 correlation between them. However, these correlation is not significant in determine the relationship between income and years of service at Company A

SECTION B

```
In [13]: 1 #Q1-Import necessary libraries
        2
        3 import numpy as np
        4 import pandas as pd
        5 import matplotlib.pyplot as plt
```

```
In [14]: 1 #Q2-Import dataset
        2
        3 data = pd.read_csv('exam_dataset.csv')
```

```
In [15]: 1 data.head()
```

```
Out[15]:
```

	Age	BusinessTravel	MonthlyIncome	JobSatisfaction	Bonus	Department	DistanceFromHome	Education	EducationField	EmployeeCount	...	JobRole	MaritalStatus	PerformanceRating	StockOptionLevel	Trainir
0	41	Travel_Rarely	5993	4	17979	Sales	1	2	Life Sciences	1	...	Sales Executive	Single	3	0	
1	49	Travel_Frequently	5130	2	20520	Research & Development	8	1	Life Sciences	1	...	Research Scientist	Married	4	1	
2	37	Travel_Rarely	2090	3	6270	Research & Development	2	2	Other	1	...	Laboratory Technician	Single	3	0	
3	33	Travel_Frequently	2909	3	8727	Research & Development	3	4	Life Sciences	1	...	Research Scientist	Married	3	0	
4	27	Travel_Rarely	3468	2	10404	Research & Development	2	1	Medical	1	...	Laboratory Technician	Married	3	1	

5 rows × 24 columns

In [16]:

```
1 data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1470 entries, 0 to 1469
Data columns (total 24 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Age                   1470 non-null   int64
 1   BusinessTravel        1470 non-null   object
 2   MonthlyIncome         1470 non-null   int64
 3   JobSatisfaction       1470 non-null   int64
 4   Bonus                1470 non-null   int64
 5   Department            1470 non-null   object
 6   DistanceFromHome     1470 non-null   int64
 7   Education             1470 non-null   int64
 8   EducationField        1470 non-null   object
 9   EmployeeCount         1470 non-null   int64
10   EmployeeNumber        1470 non-null   int64
11   EnvironmentSatisfaction 1470 non-null   int64
12   Gender                1470 non-null   object
13   JobLevel              1470 non-null   int64
14   JobRole               1470 non-null   object
15   MaritalStatus         1470 non-null   object
16   PerformanceRating     1470 non-null   int64
17   StockOptionLevel      1470 non-null   int64
18   TrainingTimesLastYear 1470 non-null   int64
19   WorkLifeBalance       1470 non-null   int64
20   YearsAtCompany        1470 non-null   int64
21   YearsSinceLastPromotion 1470 non-null   int64
22   OverTime              1470 non-null   object
23   Attrition             1470 non-null   object
dtypes: int64(16), object(8)
memory usage: 275.8+ KB
```

In [17]:

```
1 #Q3-Allocate the relevant attributes as input and output
2
3 x = data.iloc[:, [0,1,2,3]].values
4 y = data.iloc[:, 23].values
5
6 print(x)
7 print(y)

[[41 'Travel_Rarely' 5993 4]
 [49 'Travel_Frequently' 5130 2]
 [37 'Travel_Rarely' 2090 3]
 ...
 [27 'Travel_Rarely' 6142 2]
 [49 'Travel_Frequently' 5390 2]
 [34 'Travel_Rarely' 4404 3]]
['Yes' 'No' 'Yes' ... 'No' 'No' 'No']
```

In [18]:

```
1 np.unique(x[:,1])
```

Out[18]: array(['Non-Travel', 'Travel_Frequently', 'Travel_Rarely'], dtype=object)

In []:

```
1
```

```
In [19]: 1 #Q4-Use LabelEncoder to encode categorical data
2
3 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
4 lc_x = LabelEncoder()
5 x[:,1] = lc_x.fit_transform(x[:,1])
6
7 lc_y = LabelEncoder()
8 y = lc_y.fit_transform(y)
```

```
In [20]: 1 print(x)
2 print(y)

[[41 2 5993 4]
 [49 1 5130 2]
 [37 2 2090 3]
 ...
 [27 2 6142 2]
 [49 1 5390 2]
 [34 2 4404 3]]
[1 0 1 ... 0 0 0]
```

```
In [21]: 1 from sklearn.compose import ColumnTransformer
2
3 ct = ColumnTransformer([('BusinessTravel', OneHotEncoder(), [1])], remainder='passthrough')
4 x = ct.fit_transform(x)
5
6 print(x)

[[0.0 0.0 1.0 41 5993 4]
 [0.0 1.0 0.0 49 5130 2]
 [0.0 0.0 1.0 37 2090 3]
 ...
 [0.0 0.0 1.0 27 6142 2]
 [0.0 1.0 0.0 49 5390 2]
 [0.0 0.0 1.0 34 4404 3]]
```

```
In [22]: 1 x = x[:, 2:]
2 print(x)

[[1.0 41 5993 4]
 [0.0 49 5130 2]
 [1.0 37 2090 3]
 ...
 [1.0 27 6142 2]
 [0.0 49 5390 2]
 [1.0 34 4404 3]]
```

```
In [23]: 1 #Q5-Split your data into training and test sets
2
3 from sklearn.model_selection import train_test_split
4 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.2,
5                                               random_state=0)
```

```
In [24]: 1 #Q6-Normalized your data using StandardScaler
2
3 from sklearn.preprocessing import StandardScaler
4 sc = StandardScaler()
5 x_train = sc.fit_transform(x_train)
6 x_test = sc.transform(x_test)
7
8 print(x_train)
```

```
[ [ 0.64565275  2.3389367  2.41725694  1.14972558]
[ 0.64565275  0.9043263  -0.91612115  1.14972558]
[ 0.64565275  0.35255307  0.41020041 -1.57257768]
...
[ 0.64565275  0.68361701  0.29395671  1.14972558]
[ 0.64565275  0.13184377  -0.72026428  0.24229116]
[ 0.64565275  0.35255307  0.68736435 -0.66514326]
```

```
In [25]: 1 #Q7-Fit the and predict results using the Naïve Bayes Classifier
2
3 from sklearn.naive_bayes import GaussianNB
4 classifier = GaussianNB()
5 classifier.fit(x_train,y_train)
```

```
Out[25]: GaussianNB()
```

```
In [26]: 1 y_pred = classifier.predict(x_test)
```

```
In [27]: 1 print(y_test)
          2 print(y_pred)
```

[illegible]

```
In [28]: 1 #Q8-Evaluate your results using confusion matrix and calculate
          2 #the prediction accuracy
          3
          4 from sklearn.metrics import confusion_matrix, accuracy_score
          5 cm = confusion_matrix(y_test,y_pred)
          6 cm
```

```
Out[28]: array([[242,  3],
                [ 47,  2]], dtype=int64)
```

```
In [29]: 1 score = accuracy_score(y_test,y_pred)
          2 score
```

Out[29]: 0.8299319727891157


```
In [30]: 1 #Q9-Discuss your results and findings
2
3 from sklearn.metrics import classification_report
4 accuracy = round(accuracy_score(y_test, y_pred),4)*100
5 error = round(100 - accuracy,4)
6 print("Accuracy:",accuracy,'%')
7 print("Error rate:",error,'%')
8 print(classification_report(y_test, y_pred))
```

Accuracy: 82.99 %

Error rate: 17.01 %

	precision	recall	f1-score	support
0	0.84	0.99	0.91	245
1	0.40	0.04	0.07	49
accuracy			0.83	294
macro avg	0.62	0.51	0.49	294
weighted avg	0.76	0.83	0.77	294

The model has an accuracy of 82.99%. This indicate that it is a good model.

Based on the result, we could say that "Age", "BusinessTravel", "MonthlyIncome", "JobSatisfaction" are important features to predict staff attrition.

SECTION C

```
In [31]: 1 #Q1-Perform K-Means clustering (use WCSS to help find best K value)
2
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7
8 import warnings
9 warnings.filterwarnings("ignore")
```

```
In [32]: 1 data = pd.read_csv('clustering.csv')
```

```
In [33]: 1 data.head()
```

```
Out[33]:
```

	Unnamed: 0	A	B
0	0	0.329241	0.841783
1	1	1.697407	-0.236075
2	2	-0.831460	0.584743
3	3	1.825271	-0.297894
4	4	1.236577	0.121528

In [34]:

```
1 data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 300 entries, 0 to 299
Data columns (total 3 columns):
#   Column      Non-Null Count  Dtype
---  -
0   Unnamed: 0    300 non-null    int64
1   A             300 non-null    float64
2   B             300 non-null    float64
dtypes: float64(2), int64(1)
memory usage: 7.2 KB
```

In [35]:

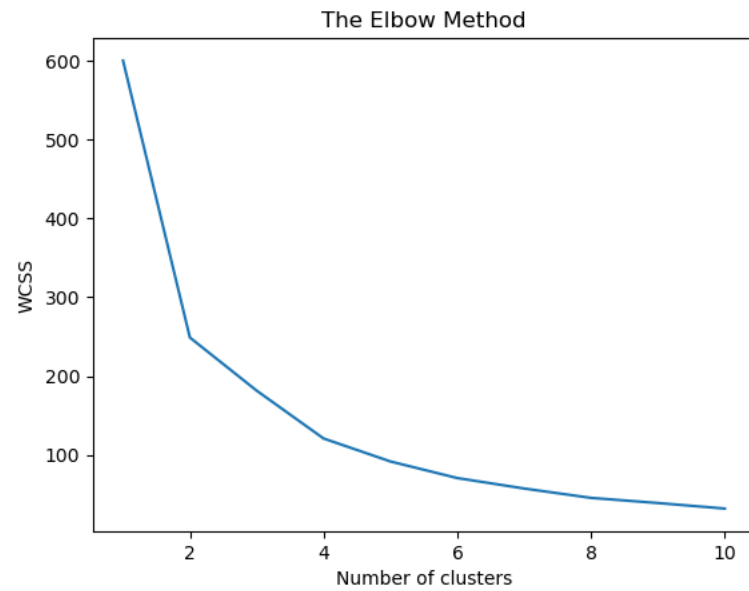
```
1 x = data.iloc[:, [1,2]].values
2
3 print(x)

[-3.16613472e-01  9.12637022e-01]
[-7.44288620e-01  5.33883584e-01]
[ 8.82870511e-01 -9.39191030e-02]
[ 4.70862998e-01 -4.48137064e-01]
[ 3.30129133e-01  9.48687978e-01]
[ 1.85021110e+00  5.27278518e-01]
[ 5.10160050e-02  2.08282719e-01]
[-3.04212080e-02  4.55428711e-01]
[ 2.65729662e-01  9.58333858e-01]
[ 1.14466723e+00 -4.72098220e-01]
[ 2.11063283e+00  3.15386651e-01]
[ 7.47395558e-01  7.27057584e-01]
[-1.01618313e+00  9.48843610e-02]
[ 9.60881972e-01 -3.70942319e-01]
[ 2.11616125e+00  1.28657950e-02]
[ 1.65823396e+00 -4.74901130e-02]
[-9.64387544e-01  4.18223971e-01]
[-5.88863100e-02  8.52369728e-01]
[ 1.31991089e+00 -4.68789511e-01]
[ 7.06101557e-01  5.84692524e-01]
```

In [36]:

```
1 #Apply StandardScaler
2
3 from sklearn.preprocessing import StandardScaler
4 sc = StandardScaler()
5 x_scaled = sc.fit_transform(x)
```

```
In [37]: 1 #Apply Elbow Method
2
3 from sklearn.cluster import KMeans
4 wcss = []
5 for i in range(1,11):
6     kmeans = KMeans(n_clusters = i, init = 'k-means++', random_state=42)
7     kmeans.fit(x_scaled)
8     wcss.append(kmeans.inertia_)
9
10 plt.plot(range(1,11),wcss)
11 plt.title("The Elbow Method")
12 plt.xlabel("Number of clusters")
13 plt.ylabel("WCSS")
14 plt.show()
```

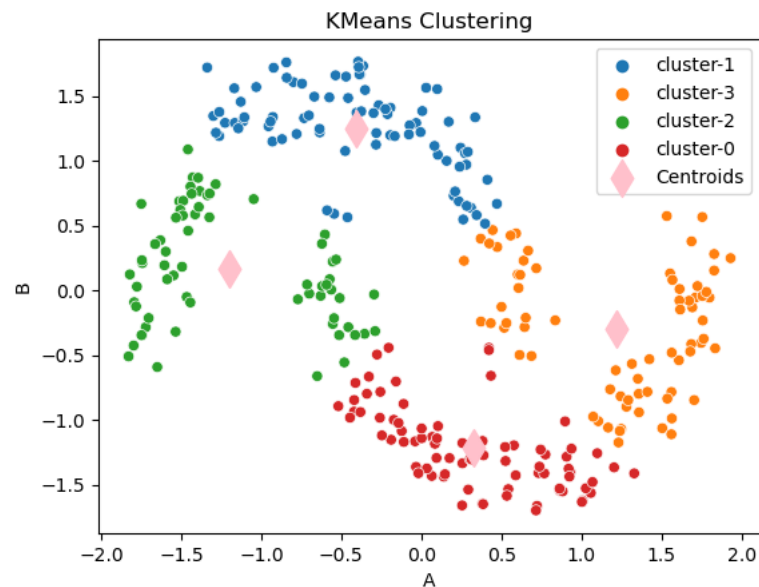


```
In [38]: 1 #Fitting training sets to KMeans
2
3 kmeans = KMeans(n_clusters = 4, init = 'k-means++', random_state=42)
4 kmeans.fit(x_scaled)
5
6 y_kmeans = kmeans.predict(x_scaled)
```

In [39]: 1 y_kmeans

Out[39]: array([[1, 3, 2, 3, 3, 3, 2, 2, 0, 0, 2, 3, 0, 2, 2, 2, 3, 0, 1, 3, 3, 0,
1, 1, 1, 0, 1, 0, 2, 0, 2, 2, 0, 1, 3, 3, 2, 0, 1, 1, 0, 0, 1, 0,
3, 1, 1, 2, 0, 1, 3, 2, 0, 2, 1, 0, 1, 3, 1, 3, 3, 1, 3, 0, 3, 3,
3, 2, 3, 3, 3, 0, 2, 0, 1, 1, 2, 3, 1, 3, 0, 2, 1, 2, 0, 0, 1, 3,
2, 2, 1, 0, 3, 1, 2, 0, 3, 3, 2, 1, 0, 1, 0, 0, 2, 1, 3, 0, 1, 3,
0, 1, 1, 1, 0, 1, 2, 1, 0, 3, 2, 2, 1, 0, 2, 3, 0, 0, 1, 1, 0, 3,
3, 0, 1, 0, 1, 0, 0, 0, 1, 3, 0, 1, 3, 0, 0, 3, 1, 3, 1, 0, 3, 1,
1, 3, 1, 1, 1, 0, 1, 1, 0, 3, 1, 1, 3, 3, 2, 2, 0, 0, 3, 3, 1, 2,
2, 0, 3, 1, 1, 0, 1, 3, 2, 0, 1, 1, 0, 1, 0, 1, 2, 2, 2, 0, 2, 1,
2, 3, 1, 1, 0, 1, 3, 0, 0, 3, 2, 0, 3, 1, 2, 2, 3, 2, 1, 2, 2, 2,
3, 0, 1, 2, 2, 2, 0, 0, 0, 3, 1, 0, 3, 3, 0, 0, 1, 1, 1, 0, 1,
3, 1, 2, 3, 0, 3, 0, 2, 3, 2, 2, 2, 3, 2, 3, 2, 1, 3, 2, 2, 3, 1,
0, 1, 3, 3, 2, 0, 3, 0, 3, 3, 0, 2, 1, 1, 0, 1, 2, 1, 1, 1, 0, 1, 2,
2, 1, 2, 2, 0, 2, 1, 0, 0, 3, 0, 3, 1, 3]])

```
In [40]: 1 #Plot the clusters
2
3 sns.scatterplot(x_scaled[:,0],x_scaled[:,1], hue=
4               ['cluster-{}'.format(x) for x in y_kmeans])
5 plt.scatter(kmeans.cluster_centers_[0,0],kmeans.cluster_centers_[0,1],
6             marker='d',s=200, c='pink', label="Centroids")
7
8 plt.title('KMeans Clustering')
9 plt.legend()
10 plt.xlabel('A')
11 plt.ylabel('B')
12 plt.show()
```

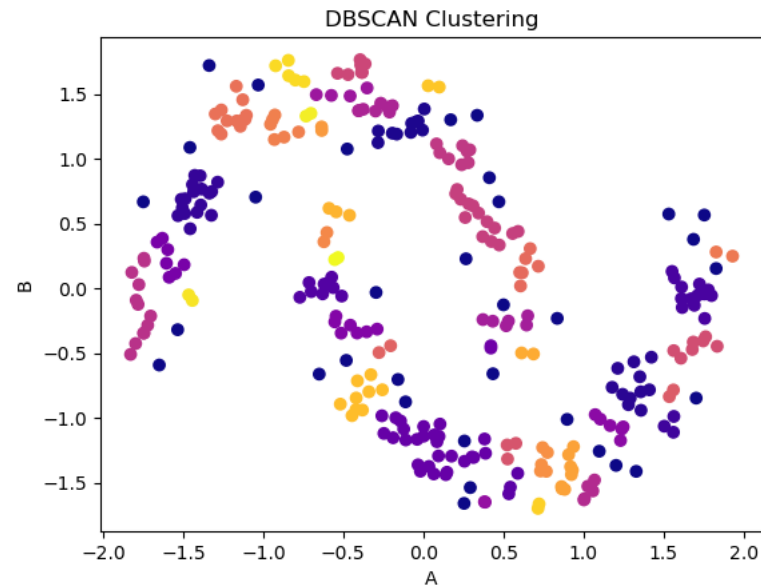


```
In [41]: 1 #Silhouette Score
2
3 from sklearn.metrics import silhouette_score
4 score = silhouette_score(x_scaled, kmeans.labels_, metric='euclidean')
5 score
```

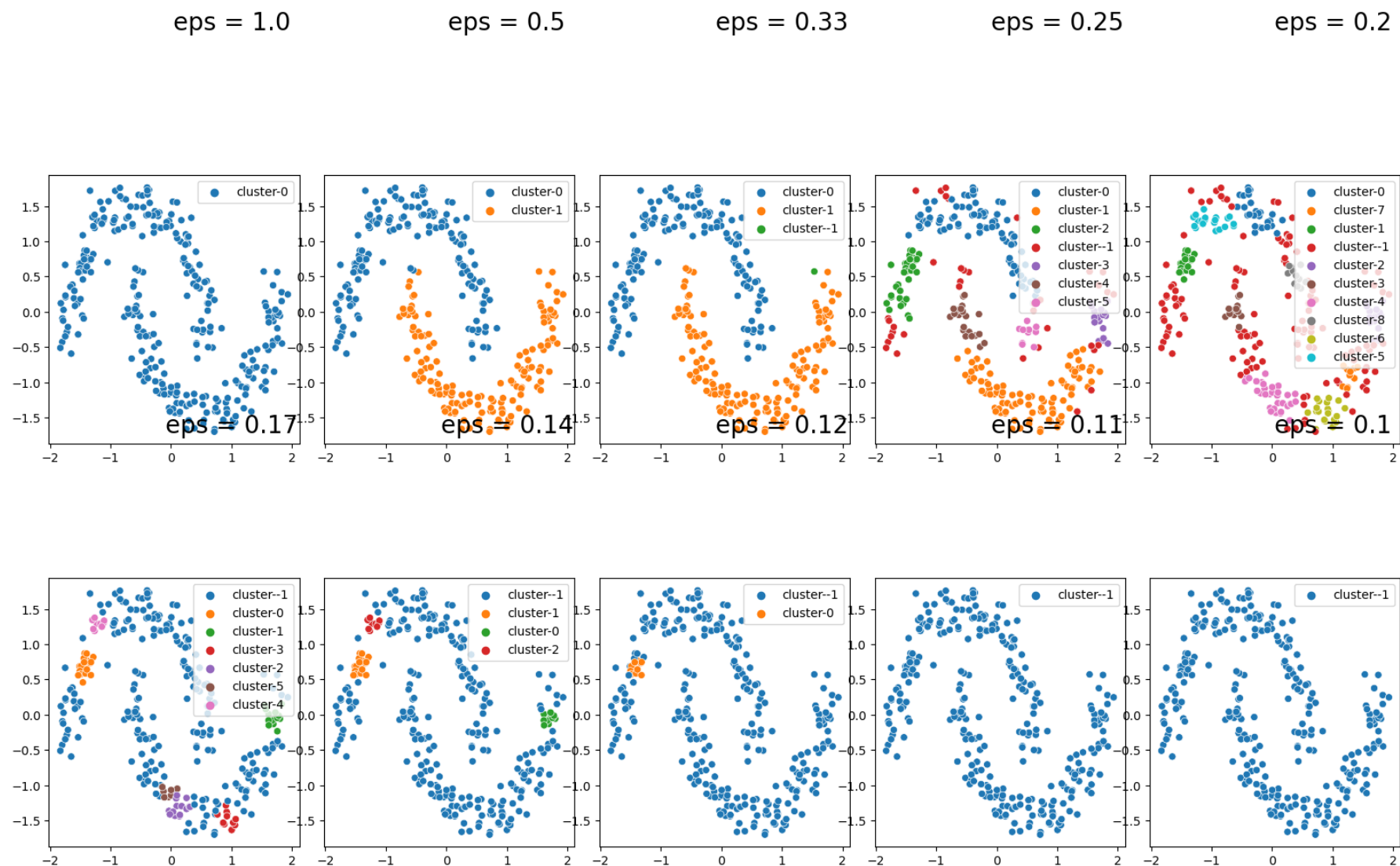
Out[41]: 0.4329118241119466

```
In [42]: 1 #Q2-Perform DBSCAN clustering (use knee locator to help find optimal
2 #parameter) on the given dataset
3
4 from sklearn.cluster import DBSCAN
5 dbs = DBSCAN(eps=0.123, min_samples=2)
6 clusters = dbs.fit_predict(x_scaled)
```

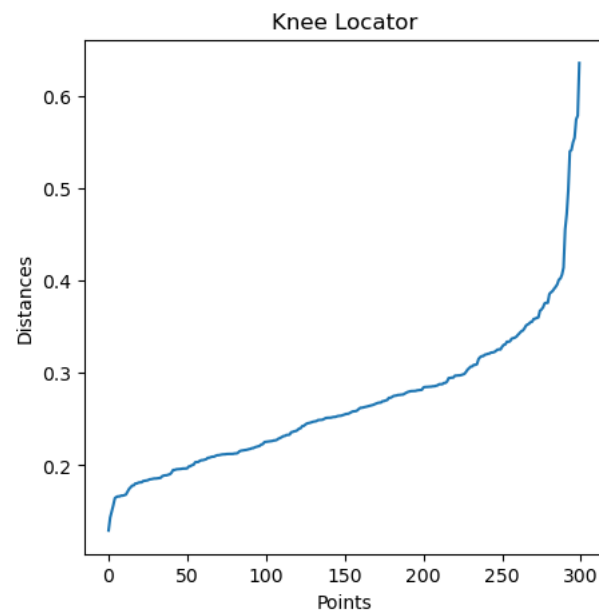
```
In [43]: 1 #Plot DBSCAN Cluster
2
3 plt.scatter(x_scaled[:,0],x_scaled[:,1], c=clusters, cmap='plasma')
4 plt.title("DBSCAN Clustering")
5 plt.xlabel('A')
6 plt.ylabel('B')
7 plt.show()
```



```
In [44]: 1 #DBSCAN Fine Tuning with Varied eps ( $\epsilon$ ) (Varied  $\epsilon = 0.1 \rightarrow 1$ )
2
3 fig = plt.figure(figsize=(20,10))
4 fig.subplots_adjust(hspace=0.5, wspace=.1)
5
6 i=1
7 for x in range(10,0,-1):
8     eps = 1/(11-x)
9     db = DBSCAN(eps=eps, min_samples=10)
10    db.fit(x_scaled)
11
12    core_samples_mask = np.zeros_like(db.labels_, dtype=bool)
13    core_samples_mask[db.core_sample_indices_] = True
14    clusters = db.labels_
15
16    ax = fig.add_subplot(2, 5, i)
17    ax.text(1, 4, "eps = {}".format(round(eps, 2)), fontsize=20, ha="center")
18    sns.scatterplot(x_scaled[:,0], x_scaled[:,1], hue=["cluster-{}".format(x) for x in clusters])
19    i += 1
```



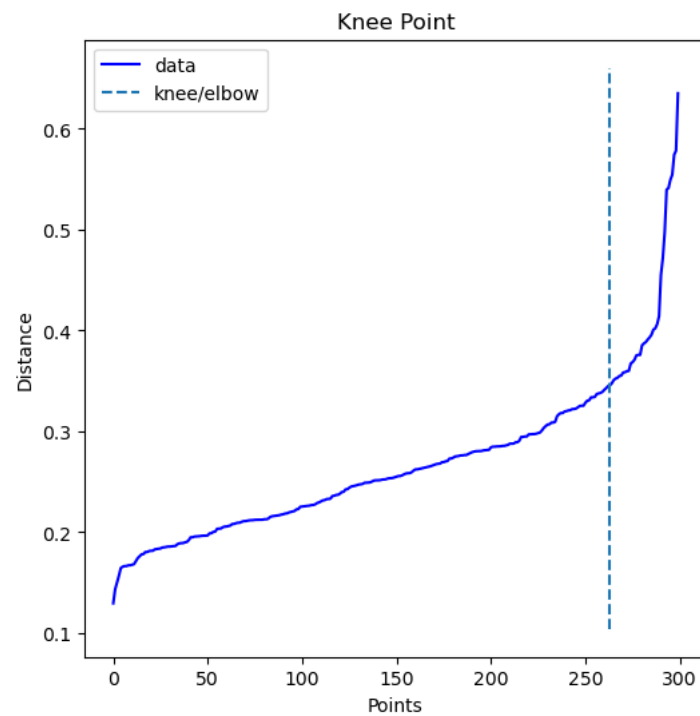
```
In [45]: 1 #DBSCAN fine tuning with knee locator
2
3 from sklearn.neighbors import NearestNeighbors
4 from kneed import KneeLocator
5
6 NN = NearestNeighbors(n_neighbors=11)
7 neighbors = NN.fit(x_scaled)
8 distances, indices = neighbors.kneighbors(x_scaled)
9
10 distances = np.sort(distances[:,10],axis=0)
11 i = np.arange(len(distances))
12 knee = KneeLocator(i, distances, S=1, curve='convex',
13                   direction = 'increasing', interp_method = 'polynomial')
14
15 fig = plt.figure(figsize=(5,5))
16 plt.plot(distances)
17 plt.title('Knee Locator')
18 plt.xlabel('Points')
19 plt.ylabel('Distances')
20 plt.show()
```




```
In [46]: 1 #Find optimum epsilon
2
3 fig = plt.figure(figsize=(5, 5))
4 knee.plot_knee()
5 plt.xlabel("Points")
6 plt.ylabel("Distance")
7 #plt.savefig("knee.png", dpi=300)
8 print(distances[knee.knee])
```

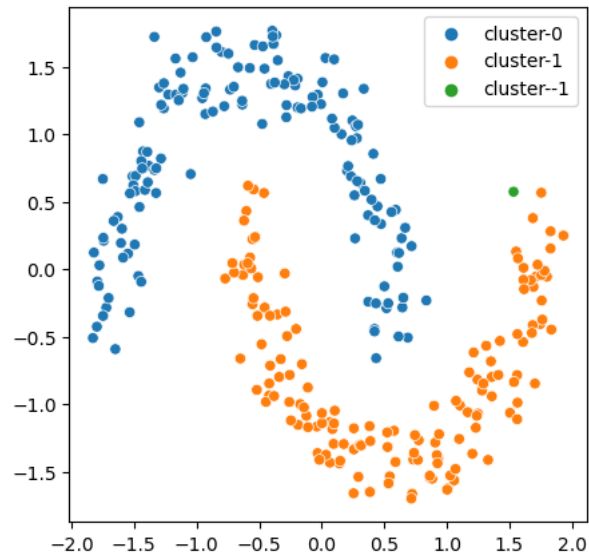
0.34567786876134754

<Figure size 500x500 with 0 Axes>



```
In [47]: 1 db = DBSCAN(eps=distances[knee.knee], min_samples=10).fit(x_scaled)
2 labels = db.labels_
3
4 fig = plt.figure(figsize=(5, 5))
5 sns.scatterplot(x_scaled[:,0], x_scaled[:,1],
6               hue=["cluster-{}".format(x) for x in labels])
```

Out[47]: <AxesSubplot:>



```
In [49]: 1 #Score for DBSCAN
2
3 from sklearn.metrics import silhouette_score
4 score = silhouette_score(x_scaled, db.labels_, metric = 'euclidean')
5 score
```

Out[49]: 0.2144392865052958

```
1 #Q3
2
3 Silhouette Score:
4 Kmeans = 0.43
5 DBSCAN = 0.21
6 K-Means model gives better silhouette score compared to DBSCAN,
7 so this means Kmeans is a better model than DBSCAN
8
9 Kmeans-4 clusters and no outlier
10 DBSCAN-2 clusters + outlier
```