

21200889_assignment_2_task2

May 28, 2022

1 NOTE: Some warning messages take too much space. To traverse using ctrl-f and look for the word 'Apple'

Code Based on:

scikit-fuzzy (n.d.) Fuzzy c-means clustering. Available at: https://pythonhosted.org/scikit-fuzzy/auto_examples/plot_cmeans.html Accessed (18 March 2022)

scikit-learn (n.d.) Demo of DBSCAN clustering algorithm. Available at: https://scikit-learn.org/stable/auto_examples/cluster/plot_dbSCAN.html#sphx-glr-auto-examples-cluster-plot-dbscan-py Accessed (19 March 2022)

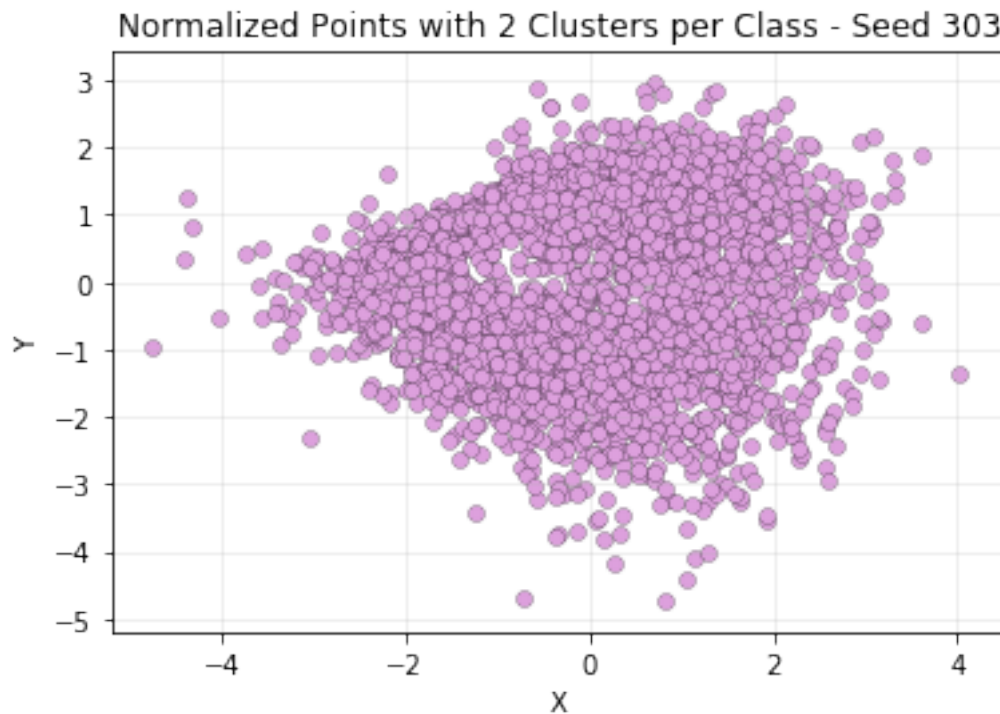
Ruane, E. (2022) assignment_2 Jupyter Notebook.

```
[3]: #Apple
      #Required for installing skfuzzy
      #pip install -U scikit-fuzzy
```

```
[4]: #Apple
      #Imports based on Ruane (2022)
      import numpy as np
      from numpy import unique # can also just use np.unique
      import pandas as pd
      import matplotlib.pyplot as plt
      %matplotlib inline
      from sklearn.cluster import KMeans
      from sklearn.preprocessing import MinMaxScaler
      from sklearn.datasets import make_classification
      from sklearn.cluster import DBSCAN
      import skfuzzy as fuzz
```

```
[5]: #Apple
      # Define dataset. Based on Ruane(2022)
      data_values, class_labels = make_classification(n_samples=3000, n_features=2,
                                                    n_informative=2, n_redundant=0,
                                                    n_classes = 2,
                                                    ↪n_clusters_per_class=2,
                                                    random_state= 303)
```

```
[6]: #Apple
# Scatter plot without showing clusters
# Create scatter plot for samples from each class. Based on Ruane (2022)
for class_value in range(len(unique(class_labels))):
    # get row indexes for samples with this class
    row_ix = np.where(class_labels == class_value)
    # create scatter of these samples
    plt.scatter(data_values[row_ix, 0], data_values[row_ix, 1], c = "plum",
        linewidths = 0.2, edgecolors = "black")
# show the plot
plt.xlabel("X")
plt.ylabel("Y")
plt.title("Normalized Points with 2 Clusters per Class - Seed 303")
plt.grid(True, alpha = 0.25)
plt.show()
```

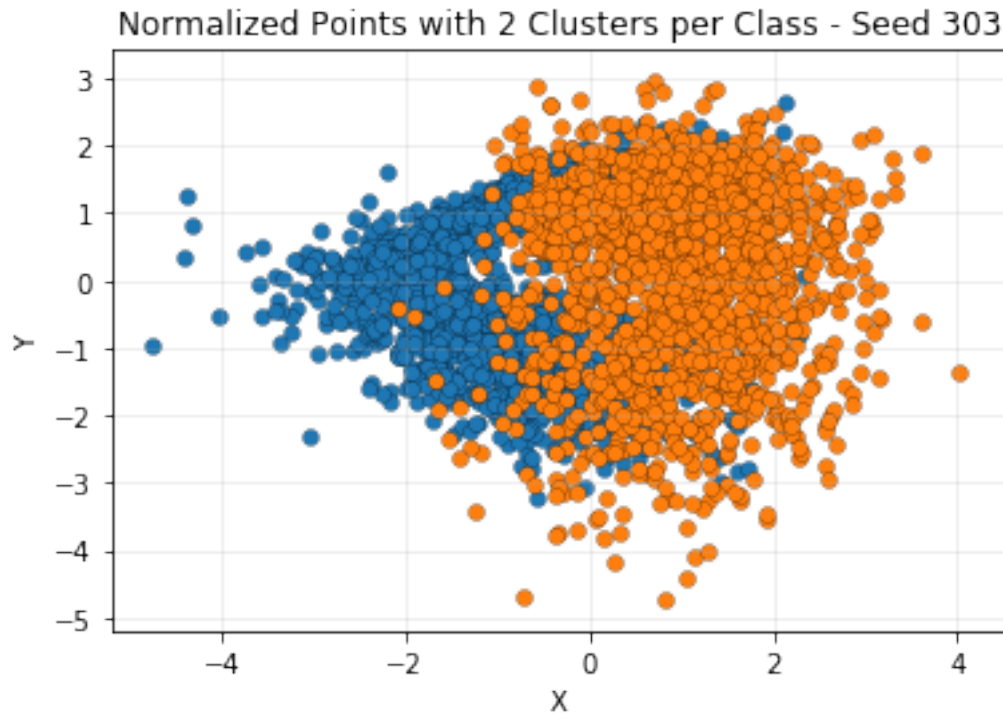


```
[7]: #Apple
# Scatter plot showing clusters
# Create scatter plot for samples from each class. Based on Ruane (2022)
for class_value in range(len(unique(class_labels))):
    # get row indexes for samples with this class
    row_ix = np.where(class_labels == class_value)
    # create scatter of these samples
```

```

plt.scatter(data_values[row_ix, 0], data_values[row_ix, 1], linewidths = 0.
↪2, edgecolors = "black")
# show the plot
plt.xlabel("X")
plt.ylabel("Y")
plt.title("Normalized Points with 2 Clusters per Class - Seed 303")
plt.grid(True, alpha = 0.25)
plt.show()

```



```

[8]: #Apple
      #Creating a DataFrame for later purposes
      k_frame = pd.DataFrame(data_values, columns = ["x", "y"])
      k_frame.sort_values(by = "x")

```

```

[8]:
      x      y
1694 -4.733662 -0.973050
2202 -4.409129  0.328300
2154 -4.358334  1.261389
  58 -4.321743  0.826715
2273 -4.024949 -0.521945
...
  64  3.306552  1.290013
2193  3.323652  1.551523

```

```
2472  3.609009  1.889647
156    3.621020 -0.599586
2428   4.019805 -1.380379
```

```
[3000 rows x 2 columns]
```

K-means++

```
[9]: #Apple
      #Elbow Plot Function from Ruane (2022)
      def get_elbow_plot(dataset, x_col='', y_col='', min_k=1, max_k=10):

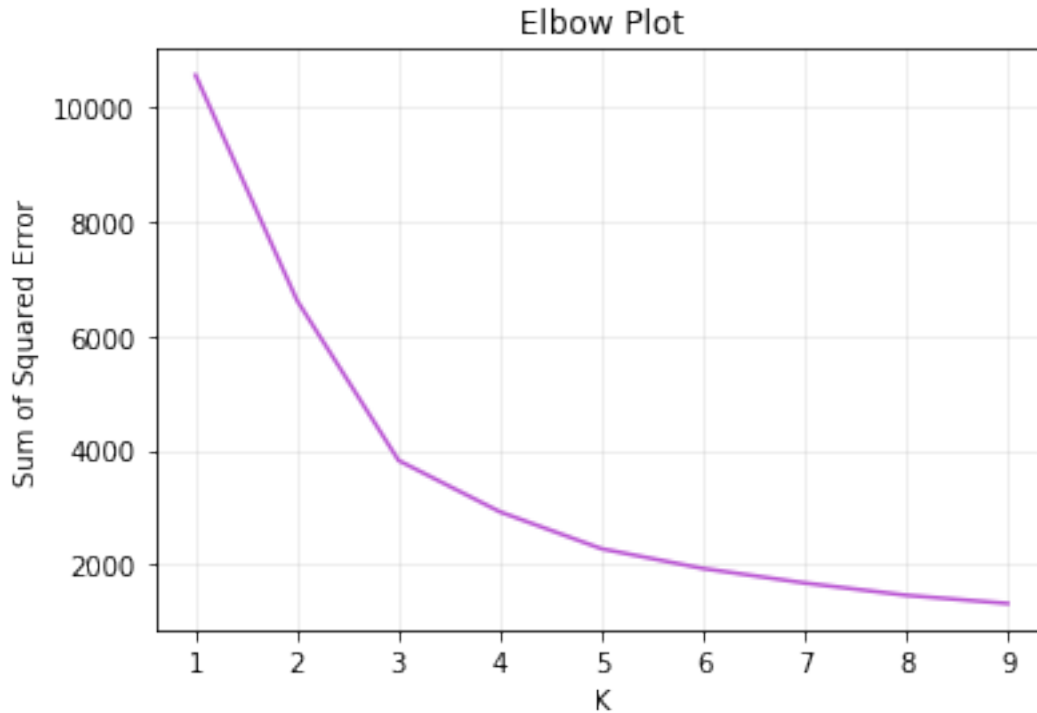
          # for what values of k do we want to run
          k_range = range(min_k, max_k)

          # calculate sum of squared error for each value of k
          sse = [] # initialize empty list

          for k in k_range:
              km = KMeans(n_clusters=k) # set whatever other parameters you want here
              km.fit(dataset)
              sse.append(km.inertia_)

          plt.xlabel('K')
          plt.ylabel('Sum of Squared Error')
          plt.plot(k_range, sse, color = "mediumorchid")
          plt.grid(True, alpha = 0.25)
          plt.title("Elbow Plot")
          plt.show()
```

```
[10]: #Apple
        #Getting the elbow plot for k-means
        get_elbow_plot(k_frame, x_col='x', y_col='y')
```



```
[11]: #Apple
#Running k-means++. Based on Ruane (2022)
k_values = [[2,3],[4,5]]
fig1, axes1 = plt.subplots(2, 2, figsize = (10, 10))
row = 0

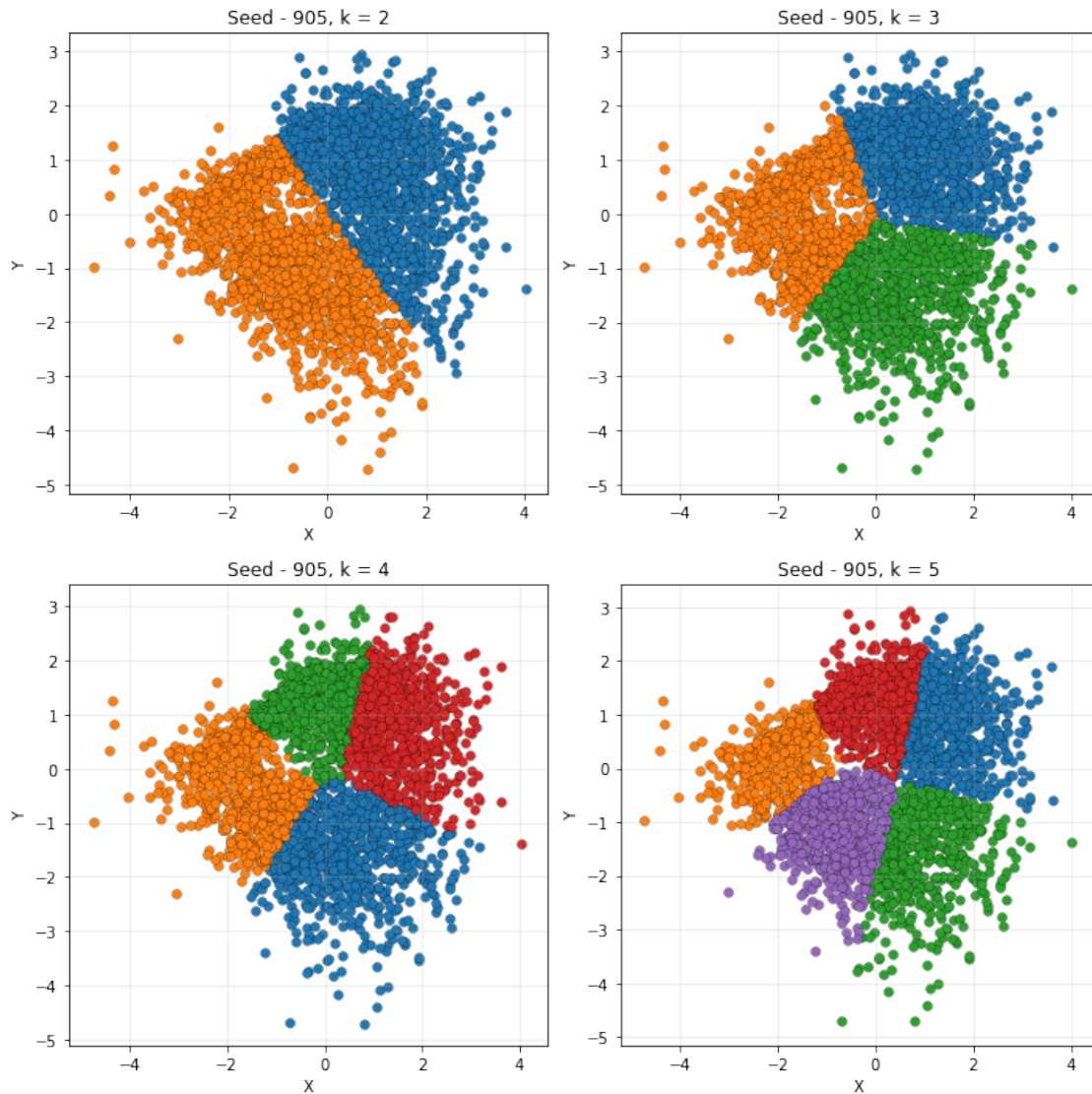
for i in range(2):
    for j in range(2):
        kmeans_model = KMeans(n_clusters= k_values[i][j], init='k-means++',
        ↪n_init=10,
                                random_state = 905) # set k to whatever makes
        ↪sense for your data
        kmeans_model.fit(data_values)
        kmeans_clusters = kmeans_model.predict(data_values)

# create scatter plot for samples from each cluster.
    for cluster in unique(kmeans_clusters):
        # get row indexes for samples with this cluster
        row_ix = np.where(kmeans_clusters == cluster)
        # create scatter of these samples.
        axes1[i,j].scatter(data_values[row_ix, 0], data_values[row_ix, 1],
        ↪linewidths = 0.2, edgecolors = "black")
        axes1[i,j].set_xlabel("X")
```

```

axes1[i,j].set_ylabel("Y")
axes1[i,j].set_title("Seed - 905, k = {}".format(k_values[i][j]))
axes1[i,j].grid(True, alpha = 0.25)
# show the plot
fig1.tight_layout()
plt.show()

```



DBSCAN

```

[15]: #Apple
#Plotting DBSCAN. Based on Ruane (2022) and sci-kit(n.d.). No color repetition
epsilon = [0.05, 0.15, 0.30]
samples = [5, 8, 10]

```

```

fig2, axes2 = plt.subplots(3, 3, figsize = (12, 12))
row = 0

for e in epsilon:
    col = 0
    for s in samples:
        dbscan_model = DBSCAN(eps= e, min_samples= s)
        dbscan_clusters = dbscan_model.fit_predict(data_values)

        colors = [plt.cm.Spectral(each) for each in np.linspace(0, 1,
→len(unique(dbscan_clusters)))]

        #for k, color in zip(unique_labels, colors):
        #    if k == -1:
        #        # Black used for noise.
        #        color = [0, 0, 0, 1]

        # create scatter plot for samples from each cluster
        for cluster in unique(dbscan_clusters):

            # get row indexes for samples with this cluster
            row_ix = np.where(dbscan_clusters == cluster)
            # create scatter of these samples
            axes2[row,col].scatter(data_values[row_ix, 0], data_values[row_ix,
→1], c = colors[cluster], linewidths = 0.2, edgecolors = "black")
            axes2[row,col].set_xlabel("X")
            axes2[row,col].set_ylabel("Y")
            axes2[row,col].set_title("No. Clusters: {}, : {}, MinPts: {}".
→format(len(unique(dbscan_clusters))-1,e,s))
            axes2[row,col].grid(True, alpha = 0.25)

            col += 1
        row += 1
# show the plot
fig2.tight_layout()
plt.show()

```

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

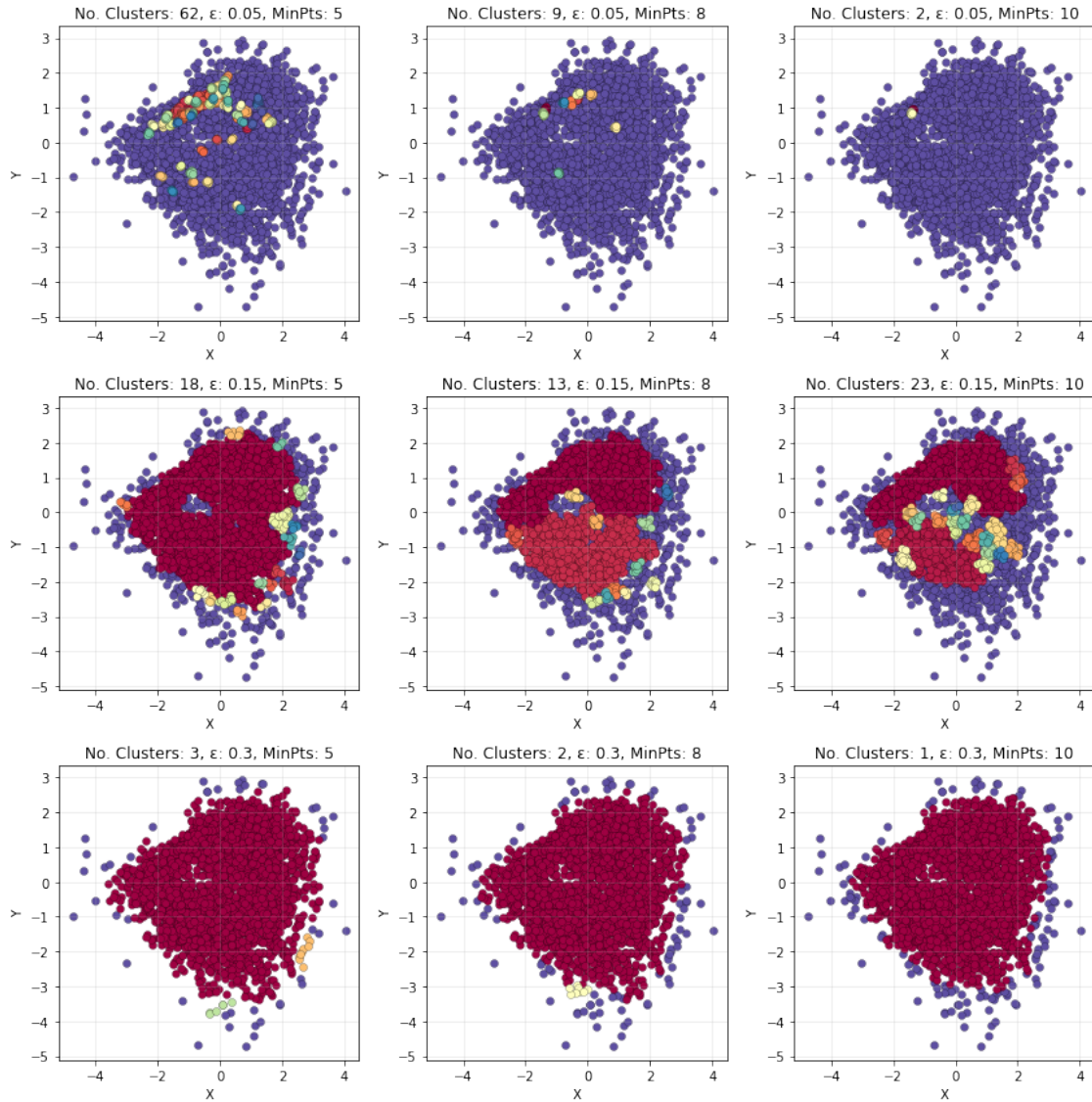
'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.

'c' argument looks like a single numeric RGB or RGBA sequence, which should be avoided as value-mapping will have precedence in case its length matches with 'x' & 'y'. Please use a 2-D array with a single row if you really want to specify the same RGB or RGBA value for all points.



```
[16]: #Apple
#Plotting DBSCAN. Based on Ruane (2022). With color repetition for comparison.
↪sake
epsilon = [0.05, 0.15, 0.30]
samples = [5, 8, 10]
fig2, axes2 = plt.subplots(3, 3, figsize = (12, 12))
row = 0

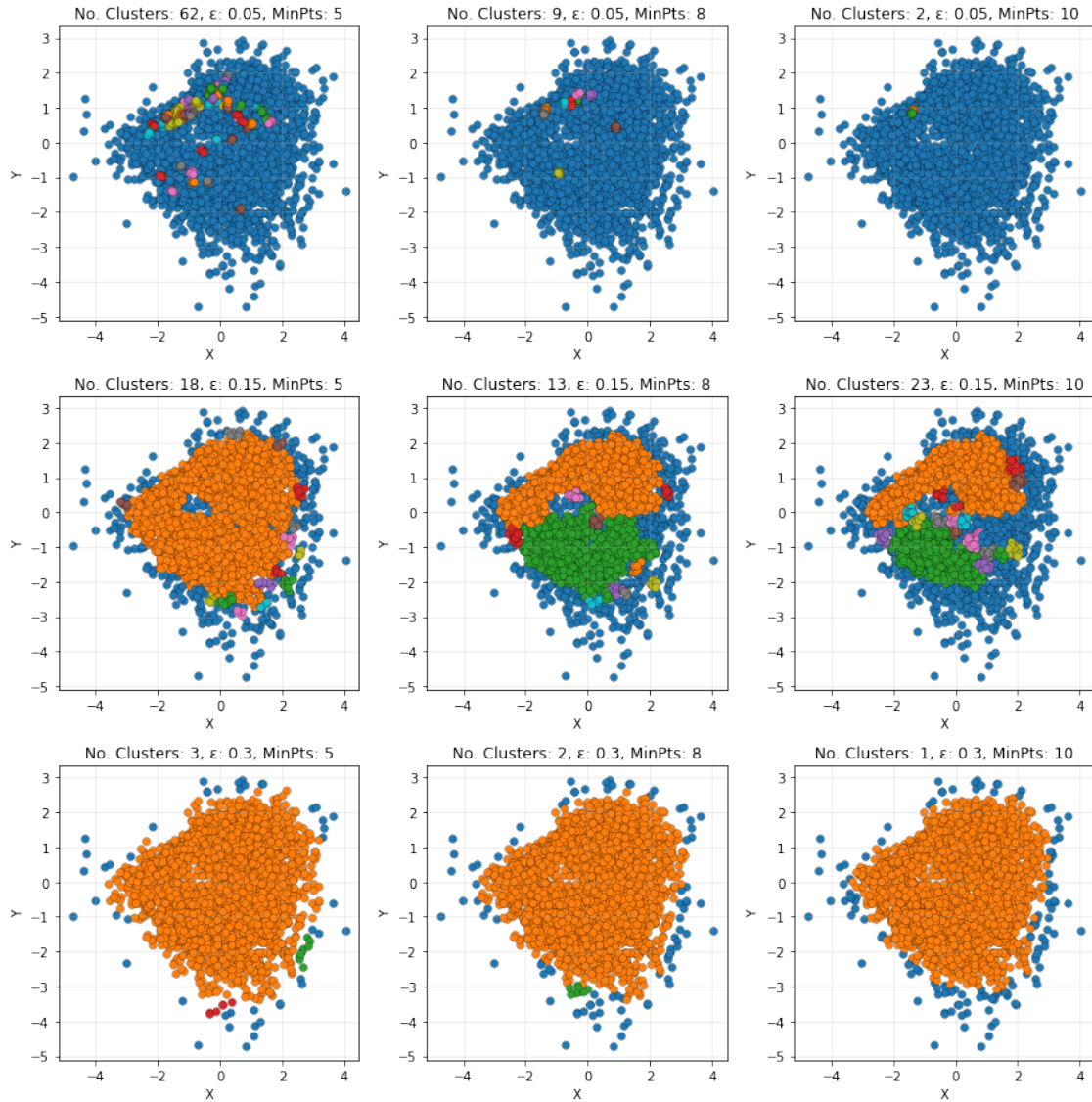
for e in epsilon:
    col = 0
    for s in samples:
        dbscan_model = DBSCAN(eps= e, min_samples= s)
```

```

dbscan_clusters = dbscan_model.fit_predict(data_values)

# create scatter plot for samples from each cluster
for cluster in unique(dbscan_clusters):
    # get row indexes for samples with this cluster
    row_ix = np.where(dbscan_clusters == cluster)
    # create scatter of these samples
    axes2[row,col].scatter(data_values[row_ix, 0], data_values[row_ix,1],linewidths = 0.2, edgecolors = "black")
    axes2[row,col].set_xlabel("X")
    axes2[row,col].set_ylabel("Y")
    axes2[row,col].set_title("No. Clusters: {}, : {}, MinPts: {}".format(len(unique(dbscan_clusters))-1,e,s))
    axes2[row,col].grid(True, alpha = 0.25)
    col += 1
    row += 1
# show the plot
fig2.tight_layout()
plt.show()

```



Fuzzy c-means

```
[17]: #Apple
#Plotting Fuzzy c-means. Based on scikit-fuzzy (n.d.) and Ruane (2022)

cluster_matrix = [[2,3,4], [5,6,7], [8,9,10]]

# Set up the loop and plot
fig3, axes3 = plt.subplots(3, 3, figsize=(15, 15))
#mi data es data_values
fpc_values = []
```

```

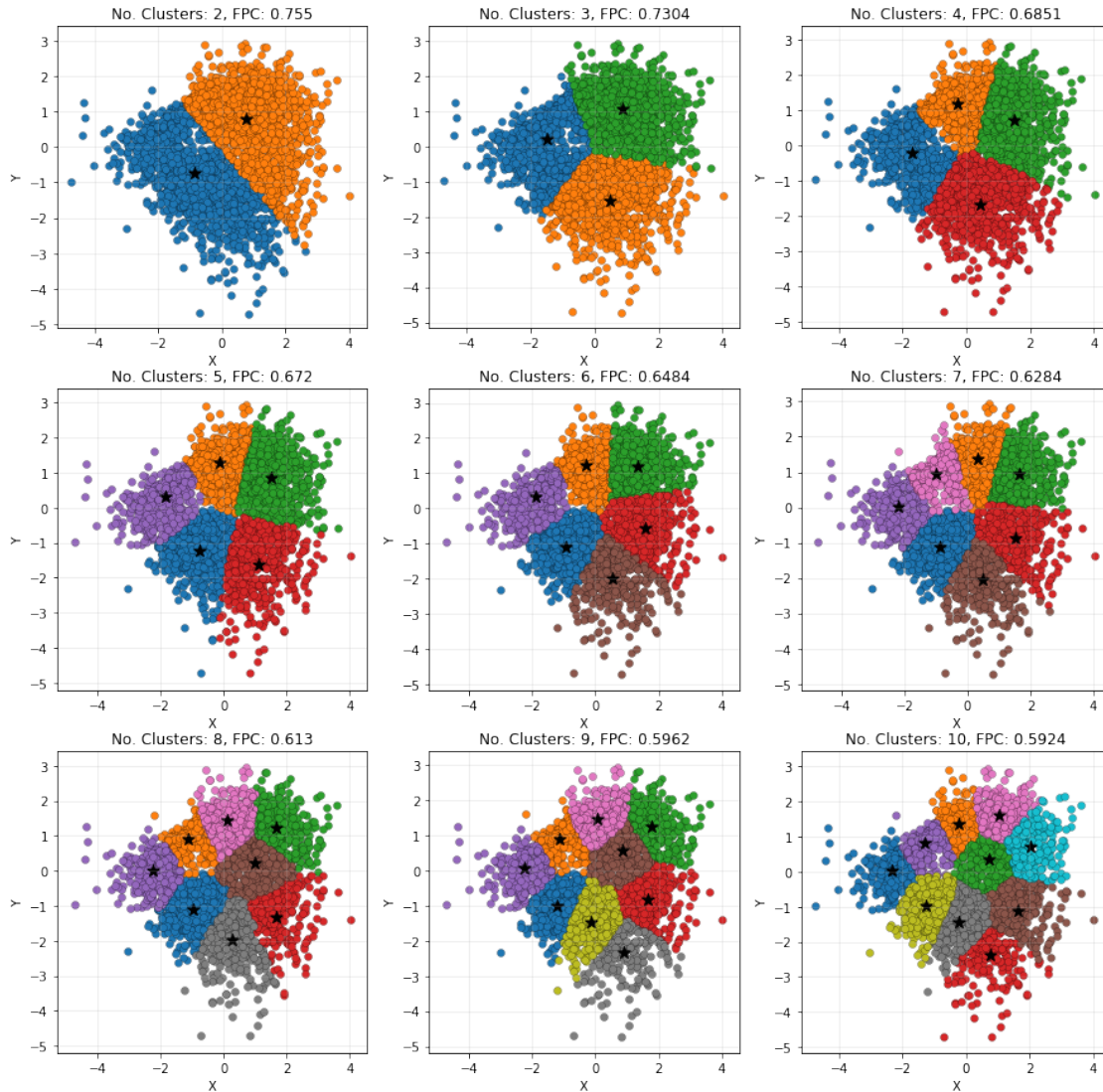
row = 0
for r in range(3):
    col = 0
    for c in range(3):
        cntr, u, u0, d, jm, p, fpc = fuzz.cluster.cmeans(
            k_frame.T, cluster_matrix[row][col], 1.75, error=0.005, maxiter=1000,
            ↪init=None, seed = 905)

        # Plot assigned clusters, for each data point in training set
        cluster_membership = np.argmax(u, axis=0)
        fpc_values.append(fpc)

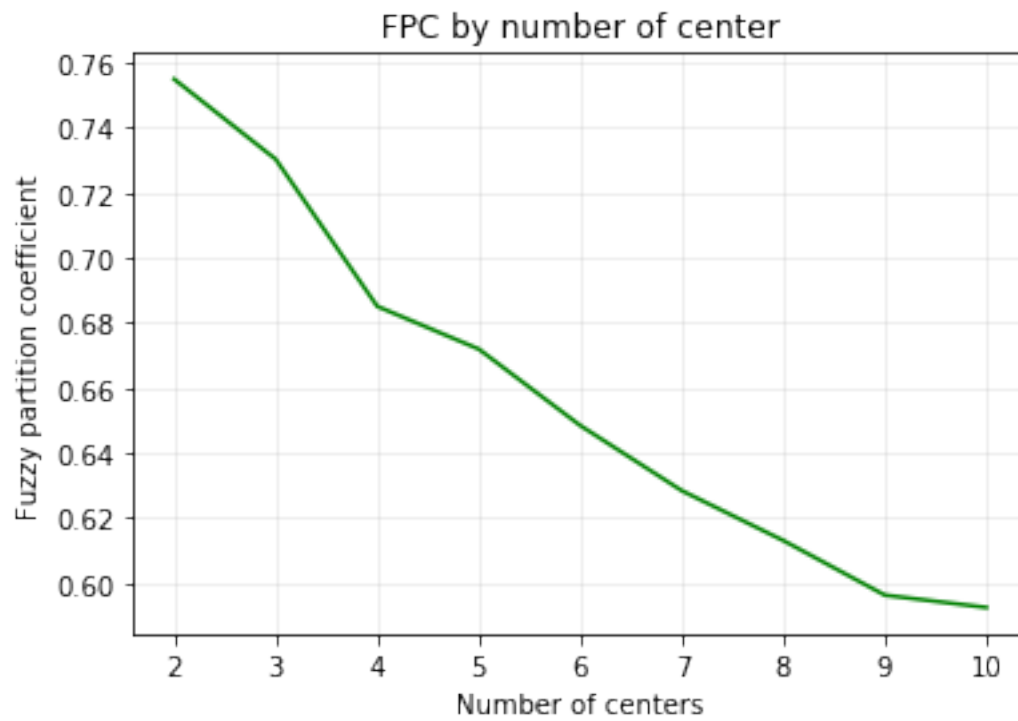
        for j in range(cluster_matrix[row][col]):
            row_ix = np.where(cluster_membership == j )
            axes3[row,col].scatter(
                k_frame["x"].loc[row_ix[0]], k_frame["y"].
            ↪loc[row_ix[0]], linewidths = 0.2, edgecolors = "black")
            for c in range(cluster_matrix[row][col]):
                axes3[row,col].scatter(
                    x = cntr[c][0], y = cntr[c][1], marker = "*", color = "black",
            ↪linewidth = 1, edgecolors = "black", s = 100)
                axes3[row,col].set_xlabel("X")
                axes3[row,col].set_ylabel("Y")
                axes3[row,col].set_title("No. Clusters: {}, FPC: {}".
            ↪format(cluster_matrix[row][col],round(fpc,4)))
                axes3[row,col].grid(True, alpha = 0.25)
                col += 1

    row += 1

```



```
[18]: #Apple
#Plotting FPC Coefficient. Based on scikit-fuzzy (n.d.)
fig5, ax5 = plt.subplots()
ax5.plot(np.r_[2:11], fpc_values, color = "green")
ax5.set_xlabel("Number of centers")
ax5.set_ylabel("Fuzzy partition coefficient")
ax5.set_title("FPC by number of center")
ax5.grid(True, alpha = 0.25)
```

[]: