

✓ Worksheet 07

Name: Jeong Yong Yang

UID: U95912941

Topics

- Density-Based Clustering

Density-Based Clustering

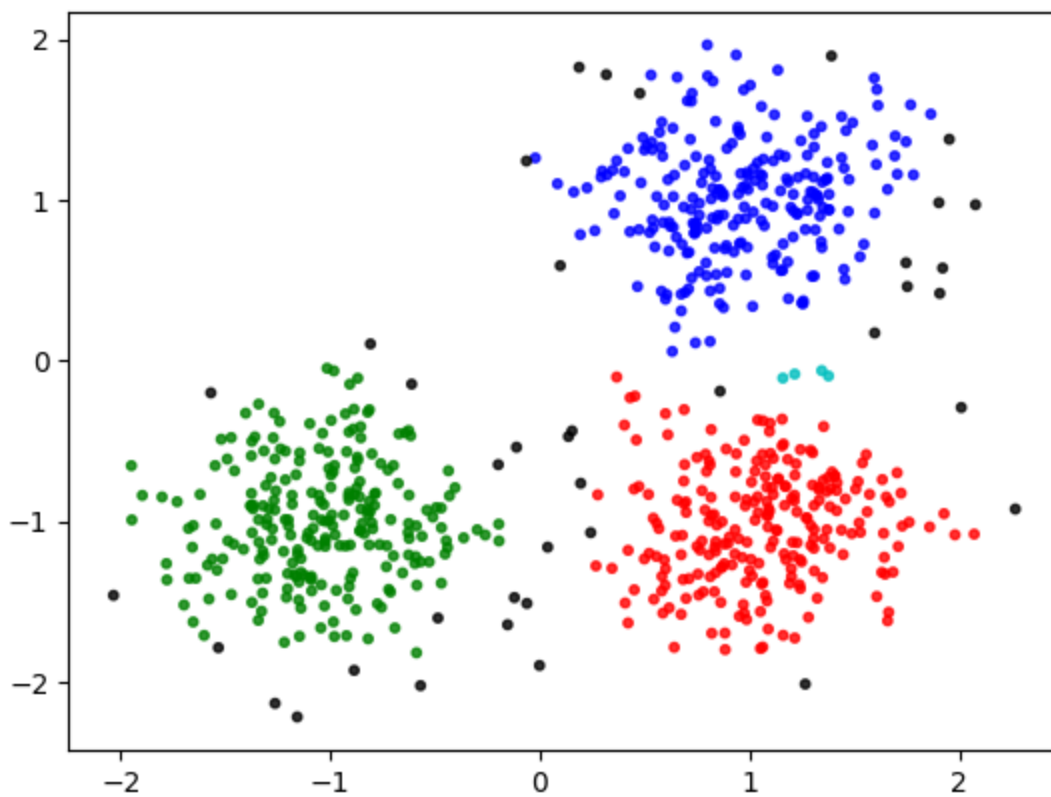
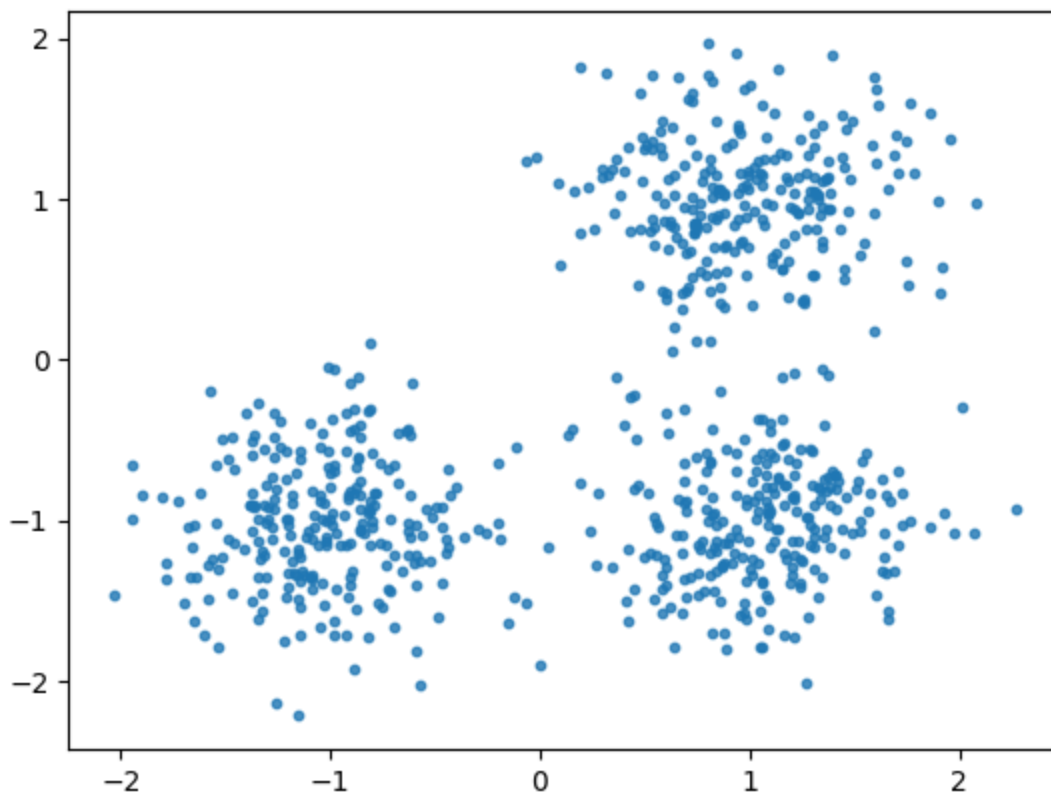
Follow along with the live coding of the DBScan algorithm.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import sklearn.datasets as datasets
4
5 centers = [[1, 1], [-1, -1], [1, -1]]
6 X, _ = datasets.make_blobs(n_samples=750, centers=centers, cluster_std=0.4,
7                             random_state=0)
8 plt.scatter(X[:,0],X[:,1],s=10, alpha=0.8)
9 plt.show()
10
11 class DBC():
12
13     def __init__(self, dataset, min_pts, epsilon):
14         self.dataset = dataset
15         self.min_pts = min_pts
16         self.epsilon = epsilon
17         self.assignments = [-1 for _ in range(len(self.dataset))]
18
19     def is_unassigned(self, i):
20         return self.assignments[i] == -1
21
22     def distance(self, i, j):
23         return np.linalg.norm(self.dataset[i] - self.dataset[j])
24
25     def get_neighborhood(self, i):
26         neighborhood = []
27         for j in range(len(self.dataset)):
28             if i != j and self.distance(i, j) <= self.epsilon:
29                 neighborhood.append(j)
30         return neighborhood
31
32     def get_unassigned_neighborhood(self, i):
33         neighborhood = self.get_neighborhood(i)
34         # return [point for point in neighborhood if self.assignments[point] == -1] # you
35         return [point for point in neighborhood if self.is_unassigned(point)]
36
37     def is_core(self, i):
38         return len(self.get_neighborhood(i)) >= self.min_pts
39
40
41     def make_cluster(self, i, cluster_num):
42         self.assignments[i] = cluster_num
43         neighborhood_queue = self.get_unassigned_neighborhood(i) #TODO: maybe make this :
44
45         while neighborhood_queue:
46             next_pt = neighborhood_queue.pop()
47             if not self.is_unassigned(next_pt): #TODO: make this a function and improve (
48                 continue
49
50             self.assignments[next_pt] = cluster_num #NOTE: border points will be assigned
51             # self.snap
52             if self.is_core(next_pt):
53                 neighborhood_queue += self.get_unassigned_neighborhood(next_pt)
54
55         return

```

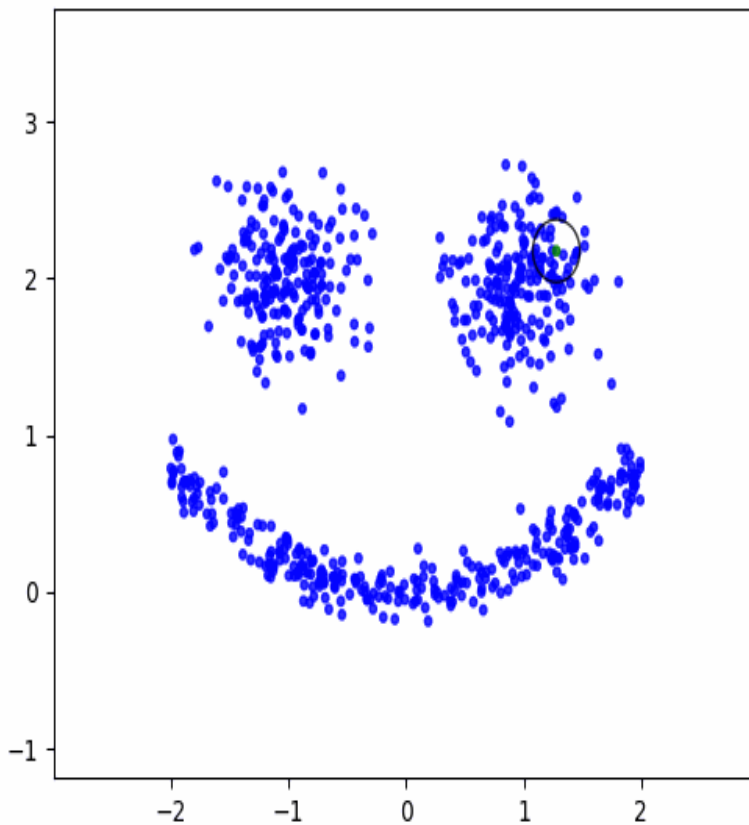
```
56
57     def dbscan(self):
58         """
59         returns a list of assignments. The index of the
60         assignment should match the index of the data point
61         in the dataset.
62         """
63
64         cluster_num = 0
65         for i in range(len(self.dataset)):
66             if self.assignments[i] != -1:
67                 continue
68
69             if self.is_core(i):
70                 # start building a new cluster
71                 self.make_cluster(i, cluster_num)
72                 cluster_num += 1
73
74         return self.assignments
75
76
77 clustering = DBC(X, 3, .2).dbscan()
78 colors = np.array([x for x in 'bgrcmykbgrcmykbgrcmykbgrcmyk'])
79 colors = np.hstack([colors] * 100)
80 plt.scatter(X[:, 0], X[:, 1], color=colors[clustering].tolist(), s=10, alpha=0.8)
```



✓ Challenge Problem

Using the code above and the template provided below, create the animation below of the DBScan algorithm.

```
1 from IPython.display import Image
2 Image(filename="dbscan_2.gif", width=500, height=500)
```



Hints:

- First animate the dbscan algorithm for the dataset used in class (before trying to create the above dataset)
- Take a snapshot of the assignments when the point gets assigned to a cluster
- Confirm that the snapshot works by saving it to a file
- Don't forget to close the matplotlib plot after saving the figure
- Gather the snapshots in a list of images that you can then save as a gif using the code below
- Use `ax.set_aspect('equal')` so that the circles don't appear to be oval shaped
- To create the above dataset you need two blobs for the eyes. For the mouth you can use the following process to generate (x, y) pairs:
 - Pick an x at random in an interval that makes sense given where the eyes are positioned
 - For that x generate y that is $0.2 * x^2$ plus a small amount of randomness
 - zip the x's and y's together and append them to the dataset containing the blobs

```

1 import numpy as np
2 from PIL import Image as im
3 import matplotlib.pyplot as plt
4 import sklearn.datasets as datasets
5
6 TEMPFILE = 'temp.png'
7
8 class DBC():
9
10     def __init__(self, dataset, min_pts, epsilon):
11         self.dataset = dataset
12         self.min_pts = min_pts
13         self.epsilon = epsilon
14         self.snaps = []
15
16
17     def snapshot(self):
18         fig, ax = plt.subplots()
19         num_points = len(self.dataset)
20         colors = np.random.rand(num_points, 3)
21
22         ax.scatter(self.dataset[:, 0], self.dataset[:, 1], s=8, c=colors)
23         cir = plt.Circle((0, 0), 1, edgecolor='black', facecolor='none') # create circle
24         ax.add_patch(cir)
25         ax.set_xlim(-10, 10)
26         ax.set_ylim(-10, 10)
27         ax.set_aspect('equal') # necessary or else the circles appear to be oval shaped
28
29         fig.savefig(TEMPFILE)
30         plt.close()
31
32         return im.fromarray(np.asarray(im.open(TEMPFILE)))
33
34
35     def dbscan(self):
36         from sklearn.cluster import DBSCAN
37         dbscan = DBSCAN(eps=self.epsilon, min_samples=self.min_pts).fit(self.dataset)
38         return dbscan.labels_
39
40
41 centers = [[-5, 0], [5, 0]]
42 eyes, _ = datasets.make_blobs(n_samples=100, centers=centers, cluster_std=.7)
43
44 mouth_x = 3.5 * np.random.random(100)
45 mouth_y = np.abs(mouth_x - 1.75)**2 - 8 + .1 * np.random.randn(100)
46
47 face = np.append(eyes, np.column_stack((mouth_x, mouth_y)), axis=0)
48
49 dbc = DBC(face, min_pts=5, epsilon=3)
50 clustering = dbc.dbscan()
51
52 dbc.snaps.append(dbc.snapshot())
53 # for i in range(len(dbc.dataset)):
54     # dbc.snaps.append(dbc.snapshot())
55

```

```
56
57 dbc.snaps[0].save(
58     'dbscan.gif',
59     optimize=False,
60     save_all=True,
61     append_images=dbc.snaps[1:],
62     loop=0,
63     duration=25
```

```

1 import numpy as np
2 from PIL import Image as im
3 import matplotlib.pyplot as plt
4 import sklearn.datasets as datasets
5
6 TEMPFILE = 'temp.png'
7
8 class DBC():
9
10     def __init__(self, dataset, min_pts, epsilon):
11         self.dataset = dataset
12         self.min_pts = min_pts
13         self.ensilon = ensilon
14
15     import numpy as np
16     from PIL import Image as im
17     import matplotlib.pyplot as plt
18     import sklearn.datasets as datasets
19
20     TEMPFILE = 'temp.png'
21
22     class DBC():
23
24         def __init__(self, dataset, min_pts, epsilon):
25             self.dataset = dataset
26             self.min_pts = min_pts
27             self.epsilon = epsilon
28             self.snaps = []
29
30
31         def snapshot(self, assignments):
32             fig, ax = plt.subplots()
33             colors = np.array([x for x in 'bgrcmymbgrcmymbgrcmymbgrcmymb'])
34             colors = np.hstack([colors] * 100)
35
36             ax.scatter(self.dataset[:, 0], self.dataset[:, 1], color=colors[assignments].tolist())
37             cir = plt.Circle((0, 0), 1, edgecolor='black', facecolor='none')
38
39

```