# Machine Learning HW5

0556562 陳鴻君

(All programs are implemented in Python 2.7 in Linux.)

(Videos is made by "ffmpeg".

The usage is "ffmpeg –y –framerate 3 –i $sourgeImages $outputpath")

## 1. K-means clustering

Here's import part.

"deepcopy" is used to copy array elements in detail.

"numpy" is a powerful matrix operation library.

"matplotlib" is used to visualize data.

"argv" is input argument list.

```
1 from copy import deepcopy
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sys import argv
```

A function named "dist" is used to calculate norm from all elements in matrix a to all elements in matrix b.

```
6 def dist(a, b, ax=1):
7     return np.linalg.norm(a-b, axis=ax)
```

In the first part of main function, save initial in following steps:
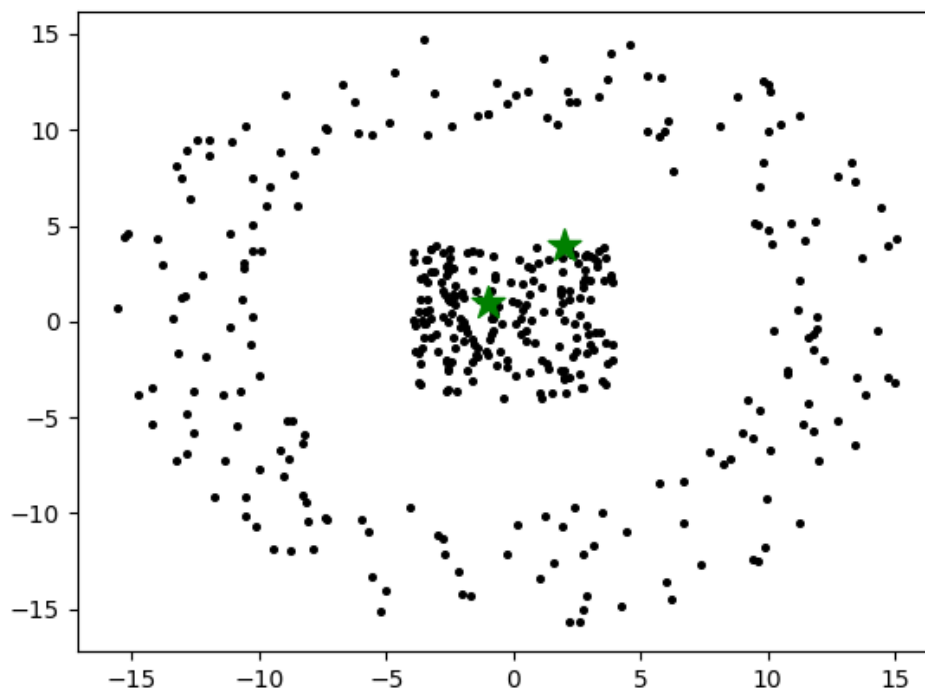
● Read data from test file, and k from arguments.

● Initial centroid matrix C, plot and save them in a figure.

```
10 if __name__ == "__main__":
11     path = '/home/jordan/Documents/ml/hw5/homework#5/'
12     filename = argv[1]
13     k = int(argv[2])
14     output = '/home/jordan/Documents/ml/hw5/output/'
15     colors = ['r', 'g', 'b', 'y', 'c', 'm']
16
17     x = []
18     y = []
19
20     with open(path+filename, 'r') as f:
21         for line in f:
22             tmp = line.split()
23             x.append(float(tmp[0]))
24             y.append(float(tmp[1]))
25     X = np.array(list(zip(x, y)))
26     plt.scatter(x, y, c="black", s=7)
27     init_x = np.random.randint(np.min(X)+2, np.max(X)-2, size=k)
28     init_y = np.random.randint(np.min(X)+2, np.max(X)-2, size=k)
29     C = np.array(list(zip(init_x, init_y)), dtype=np.float32)
30     plt.scatter(init_x, init_y, marker='*', s=200, c='g')
31     plt.savefig(output+'0.png')
32
```

For example:



In the second part of main function, keep clustering data into k clusters until centroid fixed.

- Define "error" is the converge condition, which is the difference from this turn's centroid to previous one.
- The data are clustered into k clusters by their distance to centroids. They will be clustered into the shortest one.
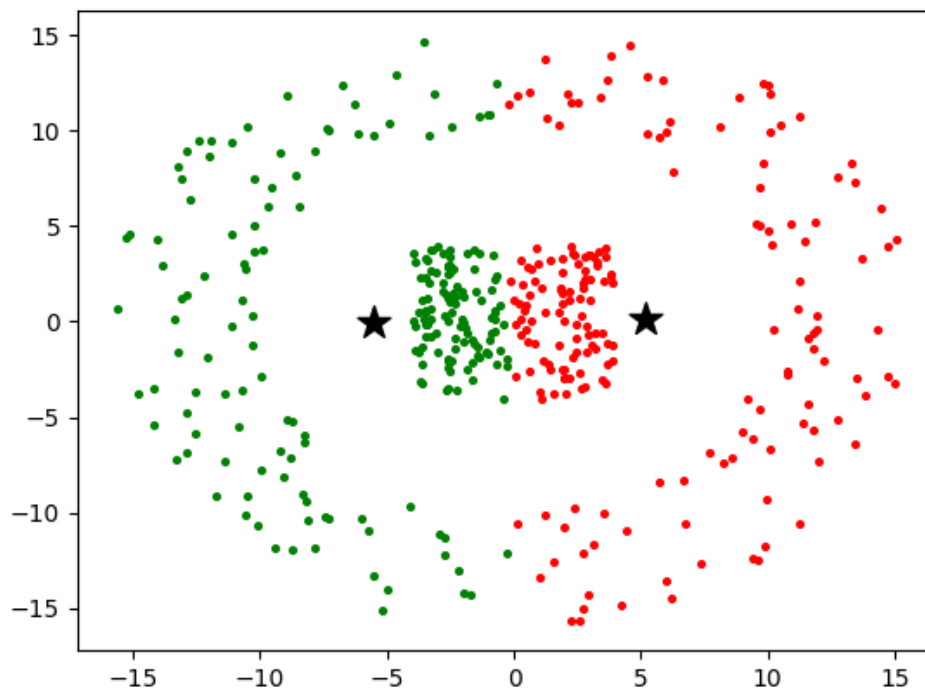
- Update centroids, by calculating the mean of all data belong to this cluster.
- Save clustering result figure in each iterator.

```
33      C_old = np.zeros(C.shape)
34      clusters = np.zeros(len(X))
35      error = dist(C, C_old, None)
36
37      iterator = 0
38      while error != 0:
39          iterator += 1
40          for i in range(len(X)):
41              distances = dist(X[i], C)
42              cluster = np.argmin(distances)
43              clusters[i] = cluster
44
45          C_old = deepcopy(C)
46          fig, ax = plt.subplots()
47          for i in range(k):
48              points = [X[j] for j in range(len(X)) if clusters[j] == i]
49              C[i] = np.mean(points, axis=0)
50              points = np.array(points)
51              ax.scatter(points[:, 0], points[:, 1], s=7, c=colors[i])
52          error = dist(C, C_old, None)
53          ax.scatter(C[:, 0], C[:, 1], marker='*', s=200, c='black')
54          plt.savefig(output+str(iterator)+'.png')
55
```

The result of test 2:



## 2. Kernel K-means clustering with RBF kernel

The import is similar to K-means.

"exp" is the exponent of e. "sqrt" is square root.

```
1 from copy import deepcopy
2 import numpy as np
3 import matplotlib.pyplot as plt
4 from sys import argv
5 from math import exp, sqrt
```

"dist" function is the same of k-means.

"dist2" function is used to calculate norm between two data points.

"rbf" function is used to calculate gram matrix G. Sigma default equals to 3.

```
 7 def dist(a, b, ax=1):
 8     return np.linalg.norm(a-b, axis=ax)
 9
10 def dist2(a, b):
11     return (a[0] - b[0])**2 + (a[1] - b[1])**2
12
13 def rbf(a, sigma=3):
14     n = len(a)
15     G = np.zeros((n, n))
16     for i in range(n):
17         for j in range(n):
18             G[i][j] = exp(-dist2(a[i], a[j])/(2*sigma**2))
19     return G
20
```

The first part of main function is same as k-means. Load and plot data.

```
21 if __name__ == "__main__":
22     path = '/home/jordan/Documents/ml/hw5/homework#5/'
23     filename = argv[1]
24     k = int(argv[2])
25     output = '/home/jordan/Documents/ml/hw5/output2/'
26     colors = ['r', 'g', 'b', 'y', 'c', 'm']
27
28     x = []
29     y = []
30
31     with open(path+filename, 'r') as f:
32         for line in f:
33             tmp = line.split()
34             x.append(float(tmp[0]))
35             y.append(float(tmp[1]))
36     X = np.array(list(zip(x, y)))
37     plt.scatter(x, y, c="black", s=7)
38     plt.savefig(output+'0.png')
39
```

In the second part of main function, plotting data points into high dimensional space, and doing k-means with high dimension data.

- Project data into high dimensional space by RBF kernel, and randomly select k rows as initial centroids C. (C_old is last round's C, and it is used to calculate error. C_low is the mean of corresponding original data points in each cluster.)
- For clustering phase, meaningful, it's using high dimensional data to do k-means clustering. Therefore, the progress is the same.
- Plotting corresponding original data points and their centroids.

```
40      X_kernel = rbf(X)
41      idx = np.random.randint(len(X_kernel), size=k)
42      C = X_kernel[idx, :]
43      C_old = np.zeros(C.shape)
44      C_low = np.zeros((k, 2))
45      clusters = np.zeros(len(X_kernel))
46      error = dist(C, C_old, None)
47
48      iterator = 0
49      while error != 0:
50          iterator += 1
51          for i in range(len(X_kernel)):
52              distances = dist(X_kernel[i], C)
53              cluster = np.argmin(distances)
54              clusters[i] = cluster
55
56          C_old = deepcopy(C)
57          fig, ax = plt.subplots()
58          for i in range(k):
59              points = [X[j] for j in range(len(X)) if clusters[j] == i]
60              points_kernel = [X_kernel[j] for j in range(len(X_kernel)) if cl
   usters[j] == i]
61              C[i] = np.mean(points_kernel, axis=0)
62              points = np.array(points)
63              ax.scatter(points[:, 0], points[:, 1], s=7, c=colors[i])
64              C_low = np.mean(points, axis=0)
65              ax.scatter(C_low[0], C_low[1], marker='*', s=200, c='black')
66          error = dist(C, C_old, None)
67          plt.savefig(output+str(iterator)+'.png')
```

The result of test 2:



# 3. Spectral clustering

Import, "dist" function, and "rbf" function is skipped because they are almost the same with K-means and kernel K-means.

- Project data into high dimensional space with RBF kernel.
- Calculate normalized Laplacian matrix L.
- Eigenvalue matrix of L is called v, and select first k columns as U.
- T is the row normalized matrix of U.
- Do k-means with matrix T

```
37      X_kernel = rbf(X)
38      L = laplacian(X_kernel, normed=True)
39      w, v = np.linalg.eig(L)
40      U = v[:,:k]
41      T = np.zeros(U.shape)
42      for i in range(len(U)):
43          for j in range(k):
44              T[i][j] = U[i][j]/sqrt(reduce(lambda a, b: a**2 + b**2, U[i]))
45
46      kmeans = KMeans(n_clusters=k).fit(T)
47      pred = kmeans.labels_
48
49      for i in range(k):
50          points = [X[j] for j in range(len(X)) if pred[j] == i]
51          points = np.array(points)
52          plt.scatter(points[:, 0], points[:, 1], s=7, c=colors[i])
53      plt.savefig(output+'result.png')
```

The result of test 2:

# 4. More clusters result

K-means:
Data = test_1
Cluster = 2



Cluster = 3

Cluster = 4



Cluster = 5

Data = test_2
Cluster = 2

Cluster = 3



Cluster = 4

Cluster = 5



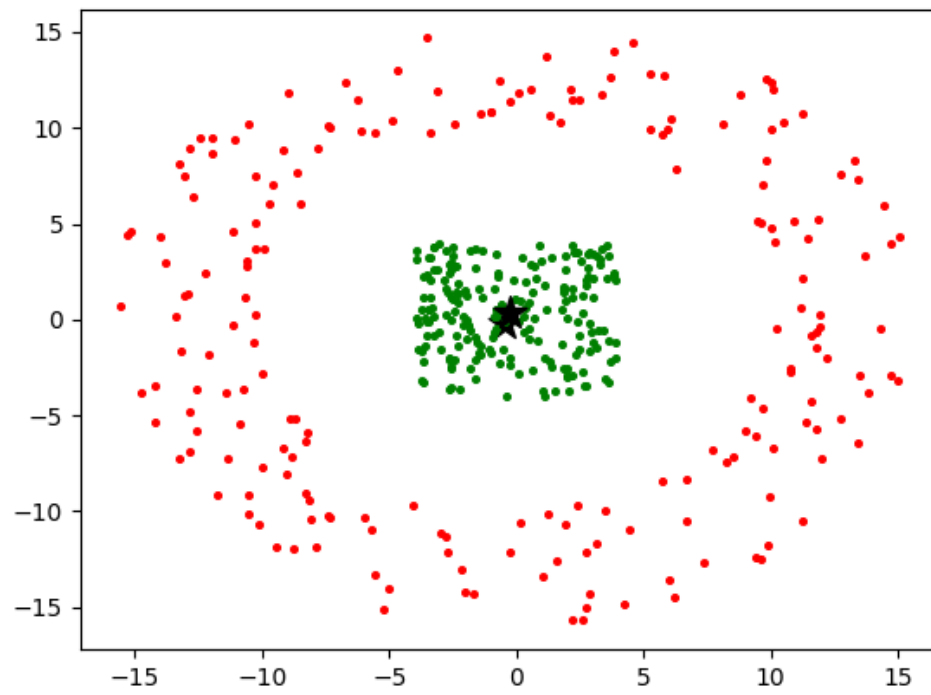Kernel K-means:

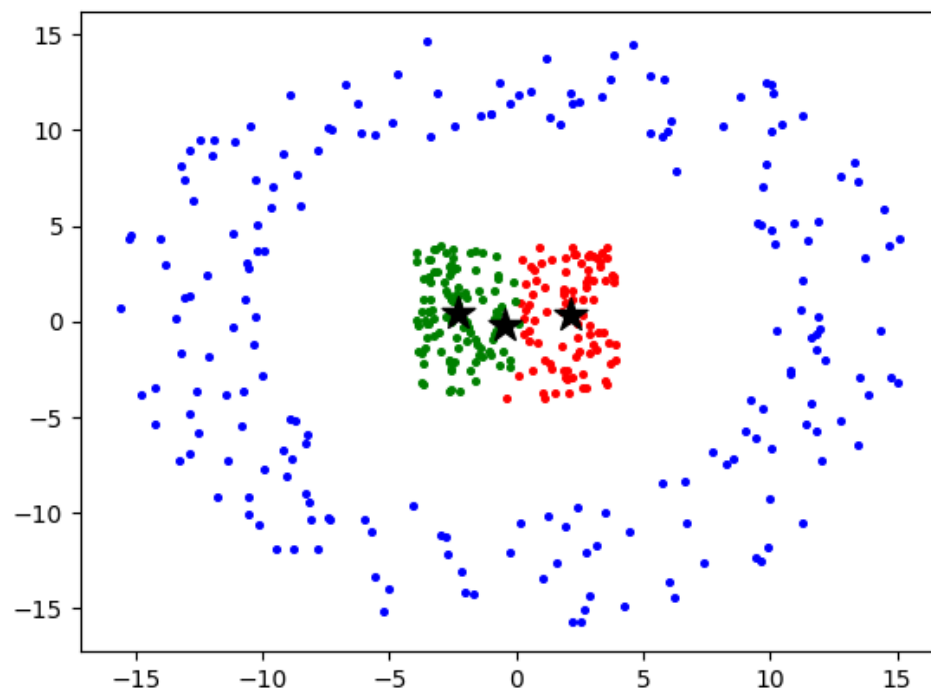Data = test_1

Cluster = 2

Cluster = 3



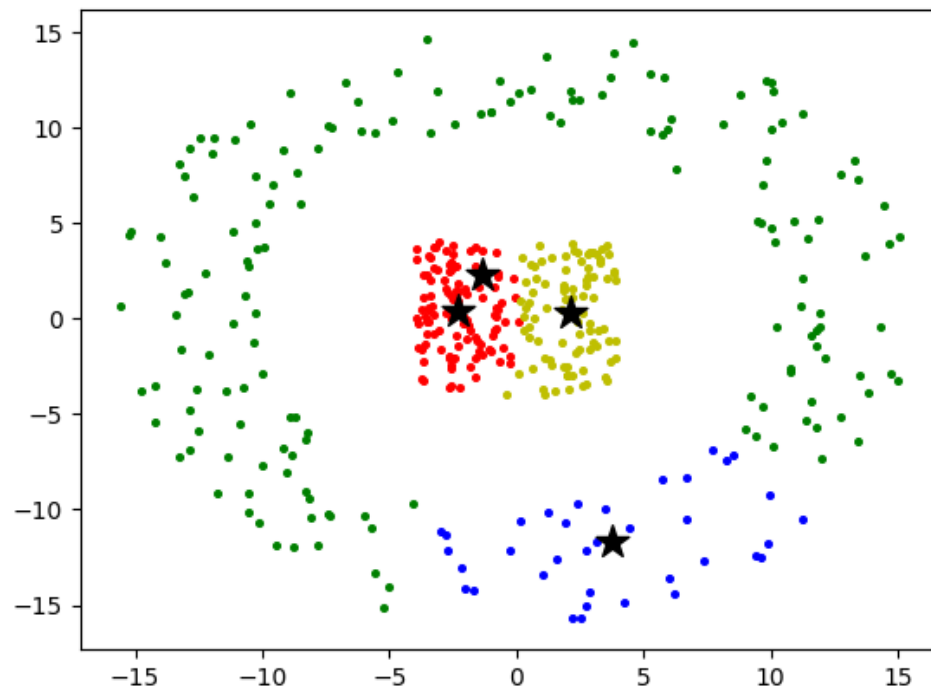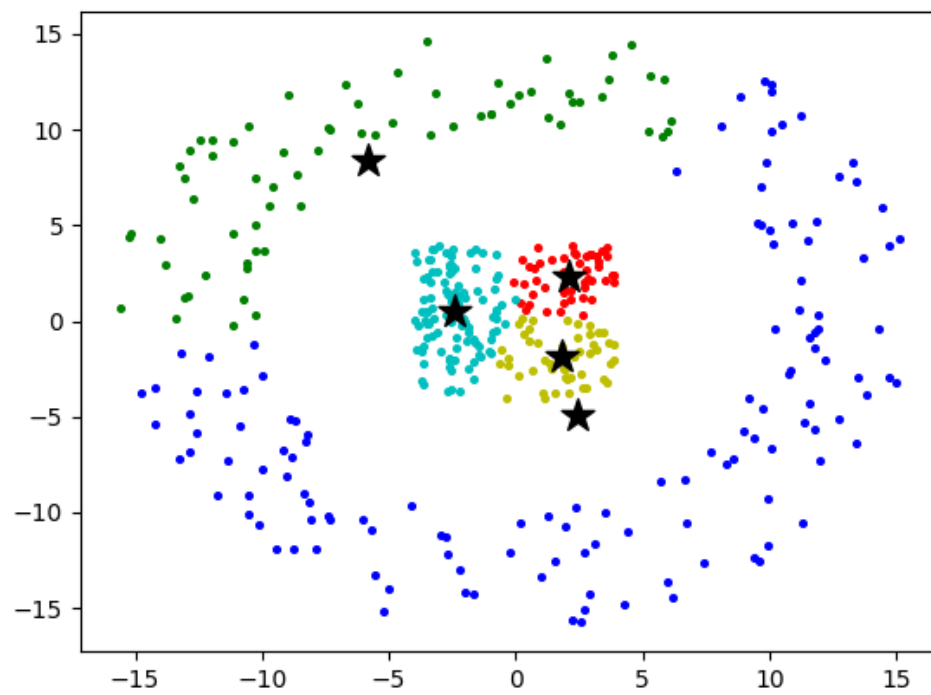Cluster = 4

Cluster = 5

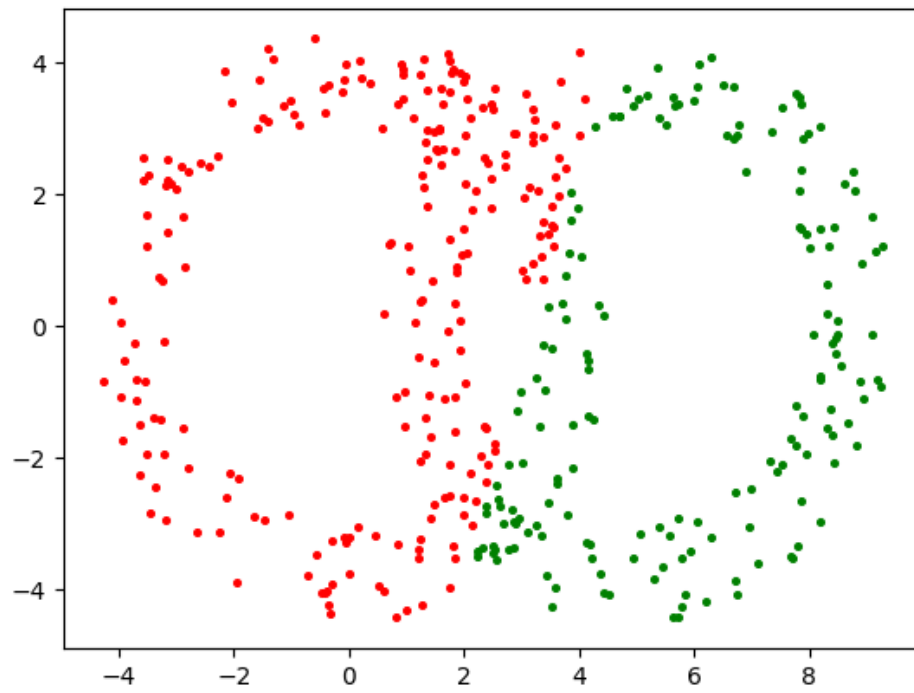

Data = test_2

Cluster = 2
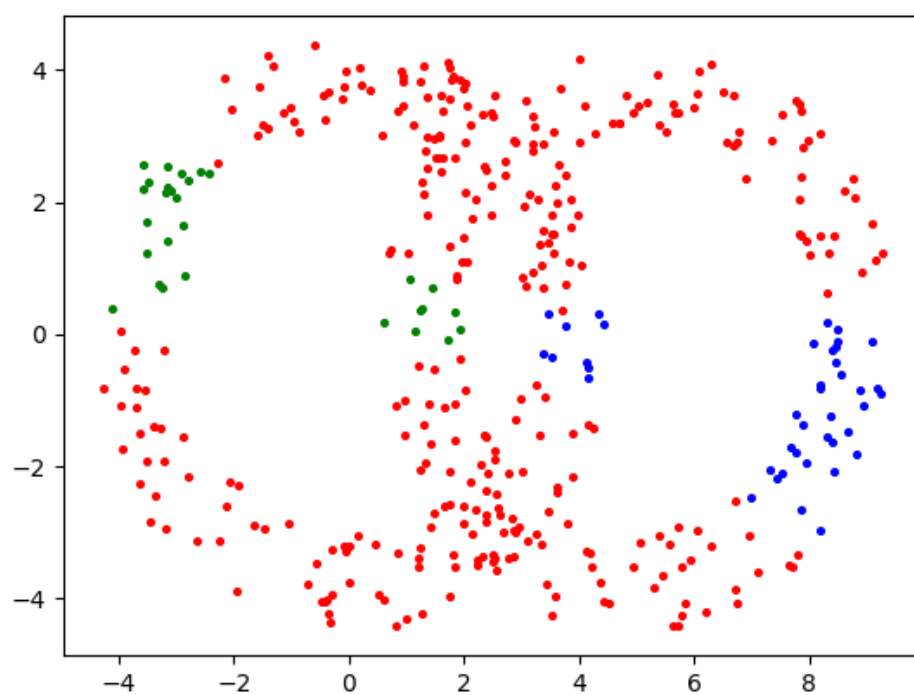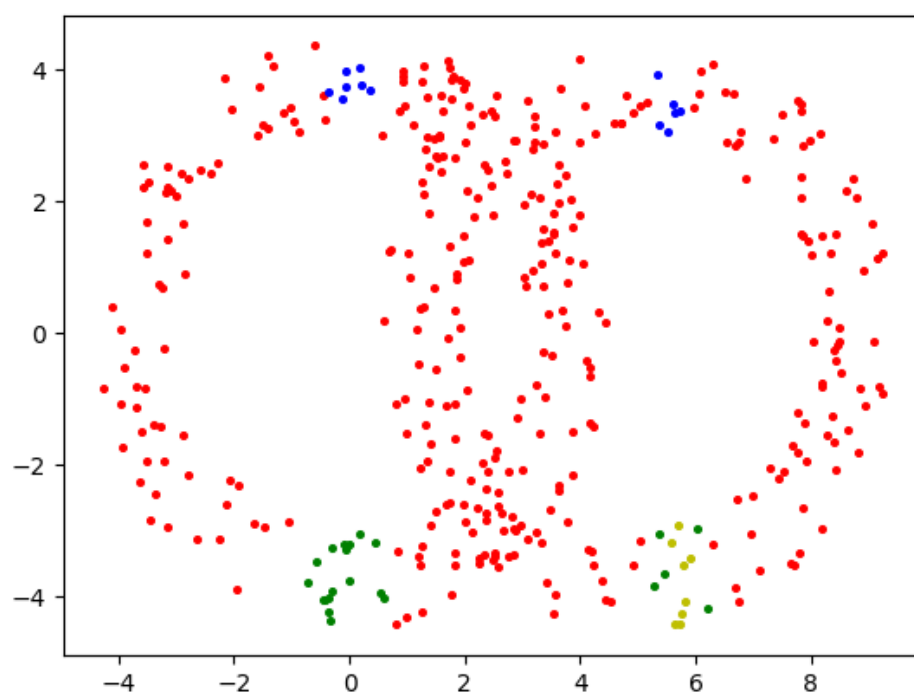


Cluster = 3

Cluster = 4



Cluster =5
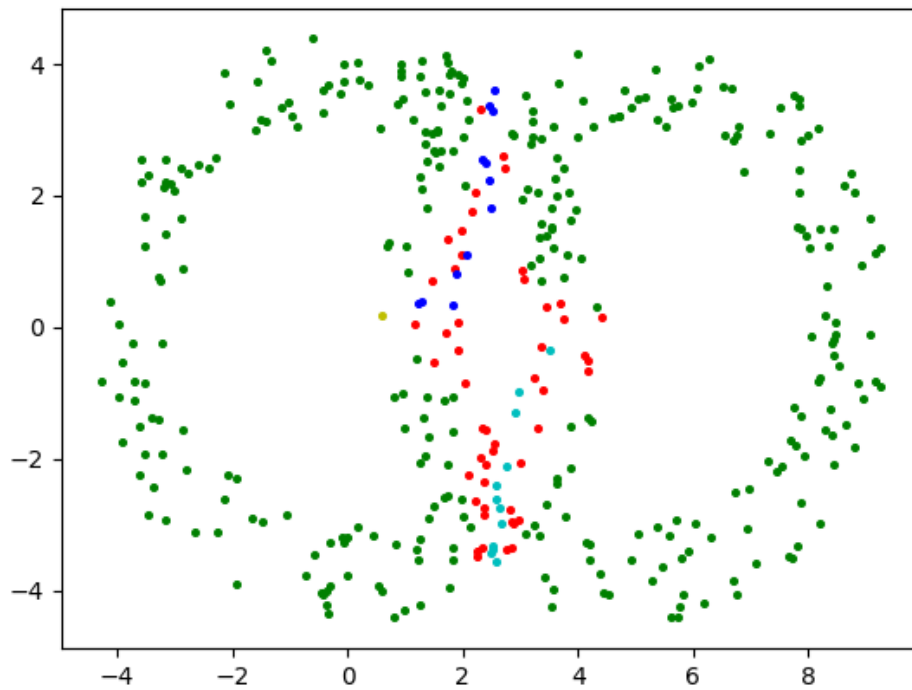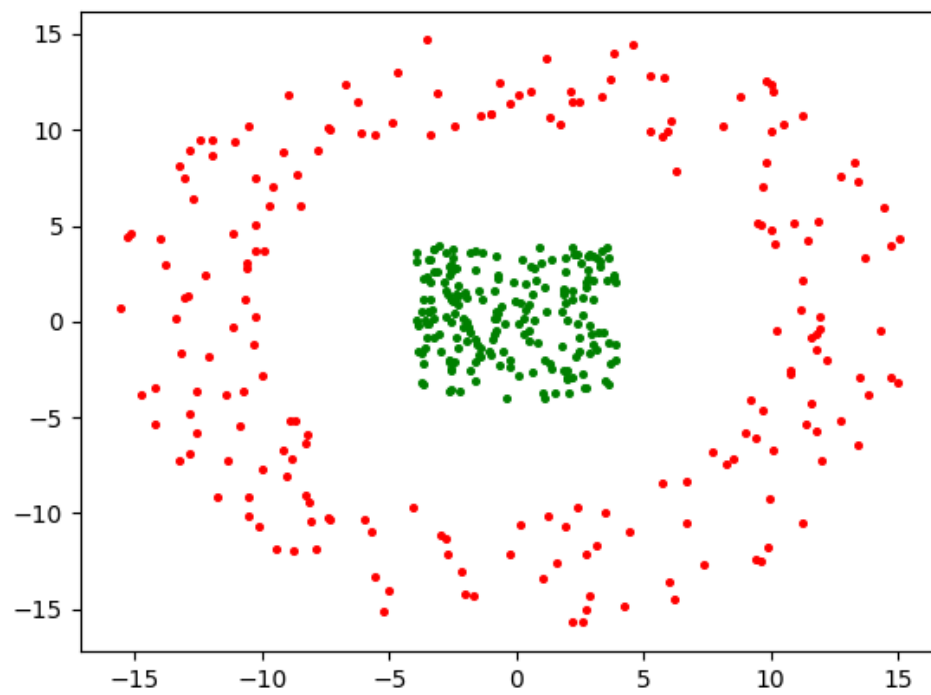
Spectral clustering:

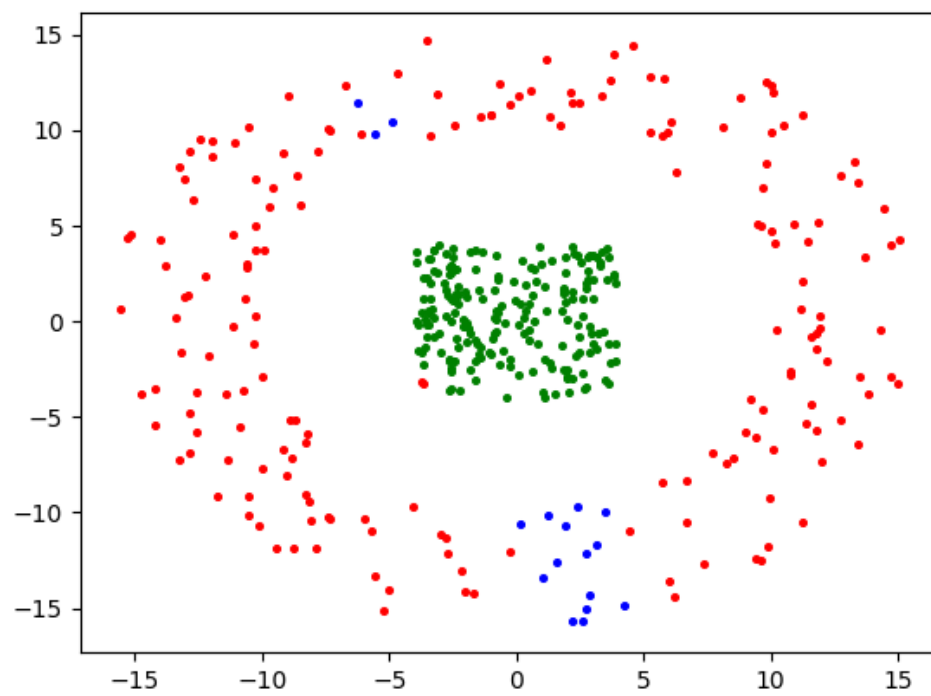Data = test_1

Cluster = 2


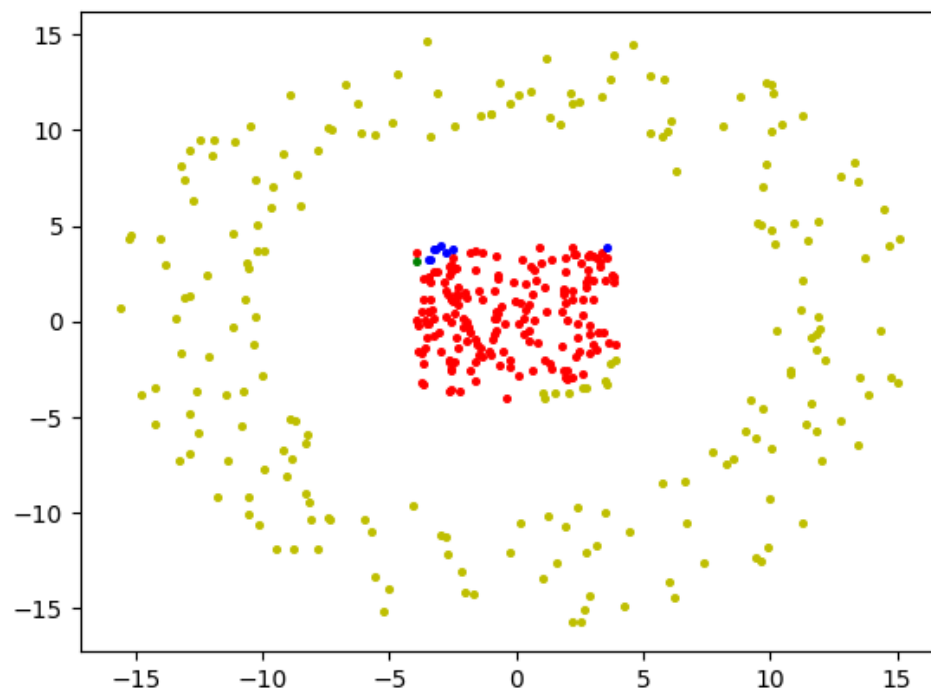
Cluster = 3

Cluster = 4



Cluster = 5

Data = test_2

Cluster = 2

Cluster = 3



Cluster = 4

Cluster = 5