

```

1 import matplotlib.pyplot as plt

1 def read_csv(csv_file):
2     csv_lst = []
3     with open(csv_file) as f:
4         lines = f.read().split('\n')[:-1]
5         for line in lines:
6             row = line.split(",")
7             csv_lst.append([int(row[0]), int(row[1])])
8     return csv_lst

1 data = read_csv("/content/drive/MyDrive/Colab Notebooks/COMP5511_AI_Assignment/Q3/I
2 data_x = [x[0] for x in data]
3 data_y = [y[1] for y in data]

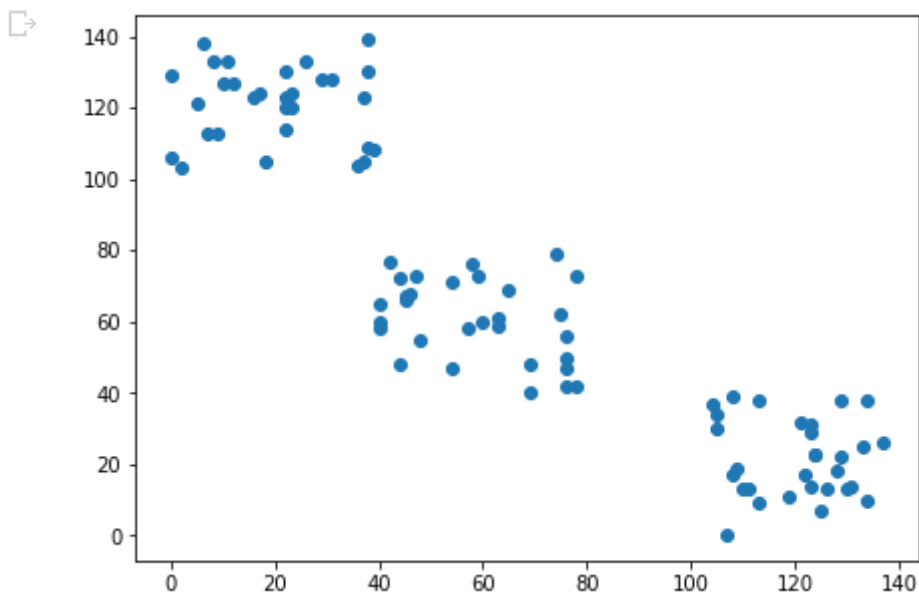
```

From the following graph, we know that there should be 3 clusters.

```

1 fig = plt.figure(figsize=(7,5))
2 plt.scatter(data_x, data_y)
3 fig.show()

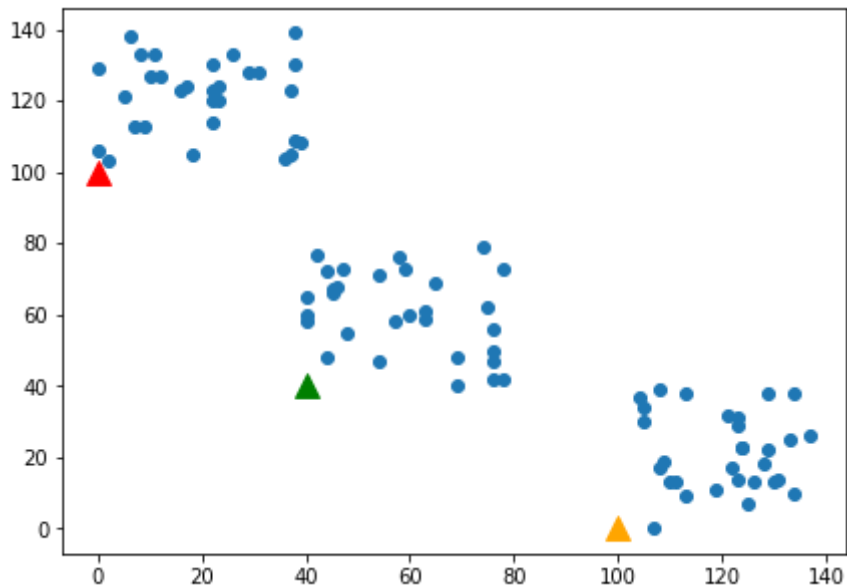
```



```

1 fig = plt.figure(figsize=(7,5))
2 initial_means = [[40, 40], [100, 0], [0, 100]]
3 initial_means_x = [x[0] for x in initial_means]
4 initial_means_y = [y[1] for y in initial_means]
5 plt.scatter(data_x, data_y)
6 for x, y, c in zip(initial_means_x, initial_means_y, ["green", "orange", "red"]):
7     plt.scatter(x, y, marker="^", color=c, s=140)
8 fig.show()

```



```

1 def euclidean_distance(p1_x, p2_x, p1_y, p2_y):
2     x = ((p2_x-p1_x)**2)
3     y = ((p2_y-p1_y)**2)
4     return (x+y)**(1/2)

```

```

1 def cal_mean(points):
2     mean_x = 0
3     mean_y = 0
4     for point in points:
5         mean_x += point[0]
6         mean_y += point[1]
7     mean_x = mean_x/len(points)
8     mean_y = mean_y/len(points)
9     return [mean_x, mean_y]

```

```

1 def k_means(points, k, itr=100):
2     N, M = len(points), len(points[0])
3
4     cluster_num = [0 for i in range(N)]
5
6     means = [[40, 40], [100, 0], [0, 100]]
7
8     for i in range(itr):
9         dist = dict([(k,[]) for k in range(N)])
10
11         for p_idx, p in enumerate(points):
12             for mean in means:
13                 dist[p_idx].append(euclidean_distance(p[0], mean[0], p[1], mean[1]))
14
15         cluster_num = []
16         for p_dist_idx, p_dist in enumerate(dist):
17             min_val = min(dist[p_dist])
18             cluster_num.append(min_val)

```

```

18     min_val_idx = [i for i, v in enumerate(dist[p_dist]) if v == min_val][0]
19     cluster_num.append(min_val_idx)
20
21     for k_class in range(k):
22         curr_point_idx_class = []
23         for cn_idx, cn in enumerate(cluster_num):
24             if cn == k_class:
25                 curr_point_idx_class.append(points[cn_idx])
26         means[k_class] = cal_mean(curr_point_idx_class)
27
28     return means

```

```

1 fig = plt.figure(figsize=(7,5))
2 kmeans = k_means(data, 3)
3 kmeans_x = [x[0] for x in kmeans]
4 kmeans_y = [y[1] for y in kmeans]
5 plt.scatter(data_x, data_y)
6 for x, y, c in zip(kmeans_x, kmeans_y, ["green", "orange", "red"]):
7     plt.scatter(x, y, marker="^", color=c, s=140)
8 fig.show()

```

