

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot
4 import pygad
```

In [4]:

```
1 cluster1_num_samples = 10
2 cluster1_x1_start = 0
3 cluster1_x1_end = 5
4 cluster1_x2_start = 2
5 cluster1_x2_end = 6
6 cluster1_x1 = np.random.random(size=(cluster1_num_samples))
7 cluster1_x1 = cluster1_x1 * (cluster1_x1_end - cluster1_x1_start) + cluster1_x1_start
8 cluster1_x2 = np.random.random(size=(cluster1_num_samples))
9 cluster1_x2 = cluster1_x2 * (cluster1_x2_end - cluster1_x2_start) + cluster1_x2_start
10 cluster2_num_samples = 10
11 cluster2_x1_start = 10
12 cluster2_x1_end = 15
13 cluster2_x2_start = 8
14 cluster2_x2_end = 12
15 cluster2_x1 = np.random.random(size=(cluster2_num_samples))
16 cluster2_x1 = cluster2_x1 * (cluster2_x1_end - cluster2_x1_start) + cluster2_x1_start
17 cluster2_x2 = np.random.random(size=(cluster2_num_samples))
18 cluster2_x2 = cluster2_x2 * (cluster2_x2_end - cluster2_x2_start) + cluster2_x2_start
```

In [5]:

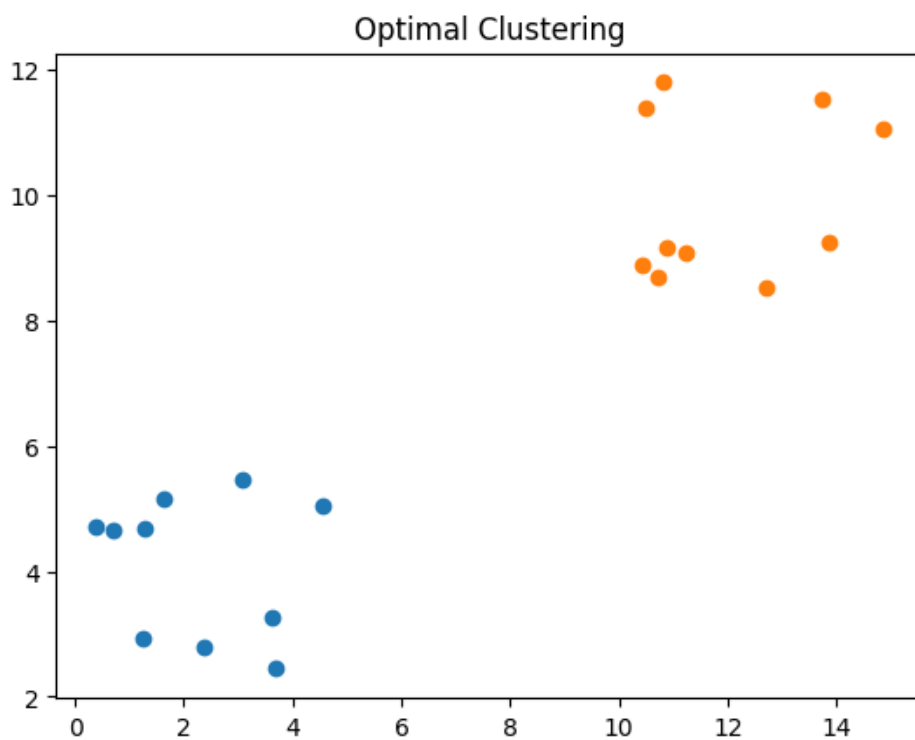
```
1 c1 = np.array([cluster1_x1, cluster1_x2]).T
2 c2 = np.array([cluster2_x1, cluster2_x2]).T
3 data = np.concatenate((c1, c2), axis=0)
4 data
```

Out[5]:

```
array([[ 1.27898768,  4.68548343],
       [ 3.61945701,  3.26190686],
       [ 3.07839194,  5.47084461],
       [ 2.38205502,  2.78383971],
       [ 0.69906865,  4.65059247],
       [ 1.23731864,  2.91766968],
       [ 0.37730323,  4.72014926],
       [ 4.56216652,  5.05882869],
       [ 3.6955702 ,  2.44550353],
       [ 1.63370019,  5.16619157],
       [10.49505539, 11.38343319],
       [10.42020043,  8.87963689],
       [10.71666831,  8.68699844],
       [14.85192003, 11.06524778],
       [12.70516146,  8.53392134],
       [10.88464926,  9.16465559 ],
       [13.73402563, 11.53691185],
       [10.82438817, 11.78955707],
       [11.24654249,  9.07758307],
       [13.86478611,  9.25591543]])
```

In [6]:

```
1 matplotlib.pyplot.scatter(cluster1_x1, cluster1_x2)
2 matplotlib.pyplot.scatter(cluster2_x1, cluster2_x2)
3 matplotlib.pyplot.title("Optimal Clustering")
4 matplotlib.pyplot.show()
5
```



In [7]:

```
1 def euclidean_distance(X, Y):
2     return np.sqrt(np.sum(np.power(X - Y, 2), axis=1))
```

In [25]:

```

1 def cluster_data(solution, solution_idx):
2     global num_cluster, data
3     feature_vector_length = data.shape[1]
4     cluster_centers = []
5     all_clusters_dists = []
6     clusters = []
7     clusters_sum_dist = []
8
9     for clust_idx in range(num_clusters):
10        cluster_centers.append(solution[feature_vector_length*clust_idx:feature_vector_length*(clust_idx+1)])
11        cluster_center_dists = euclidean_distance(data, cluster_centers[clust_idx])
12        all_clusters_dists.append(np.array(cluster_center_dists))
13
14    cluster_centers = np.array(cluster_centers)
15    all_clusters_dists = np.array(all_clusters_dists)
16
17    cluster_indices = np.argmin(all_clusters_dists, axis=0)
18    for clust_idx in range(num_clusters):
19        clusters.append(np.where(cluster_indices == clust_idx)[0])
20        if len(clusters[clust_idx]) == 0:
21            clusters_sum_dist.append(0)
22        else:
23            clusters_sum_dist.append(np.sum(all_clusters_dists[clust_idx, clusters[clust_idx]]))
24    clusters_sum_dist = np.array(clusters_sum_dist)
25
26    return cluster_centers, all_clusters_dists, cluster_indices, clusters, clusters_sum_dist

```

In [26]:

```

1 def fitness_func(ga_instance, solution, solution_idx):
2     _, _, _, _, clusters_sum_dist = cluster_data(solution, solution_idx)
3     fitness = 1.0 / (np.sum(clusters_sum_dist) + 0.00000001)
4     return fitness

```

In [27]:

```

1 num_clusters = 2
2 num_genes = num_clusters * data.shape[1]
3 ga_instance = pygad.GA(num_generations=100,
4                         sol_per_pop=10,
5                         num_parents_mating=5,
6                         init_range_low=-6,
7                         init_range_high=20,
8                         keep_parents=2,
9                         num_genes=num_genes,
10                        fitness_func=fitness_func,
11
12                        suppress_warnings=True)
13 ga_instance.run()
14

```

In [28]:

```

1 best_solution, best_solution_fitness, best_solution_idx = ga_instance.best_solution()
2 print("Best solution is {bs}".format(bs=best_solution))
3 print("Fitness of the best solution is {bsf}".format(bsf=best_solution_fitness))
4 print("Best solution found after {gen} generations".format(gen=ga_instance.best_solution_generation))

```

Best solution is [11.33809485 9.31276481 1.92397849 4.27798275]

Fitness of the best solution is 0.0283305456175161

Best solution found after 67 generations

In [31]:

```

1 cluster_centers, all_clusters_dists, cluster_indices, clusters, clusters_sum_dist= cluster_data(best.

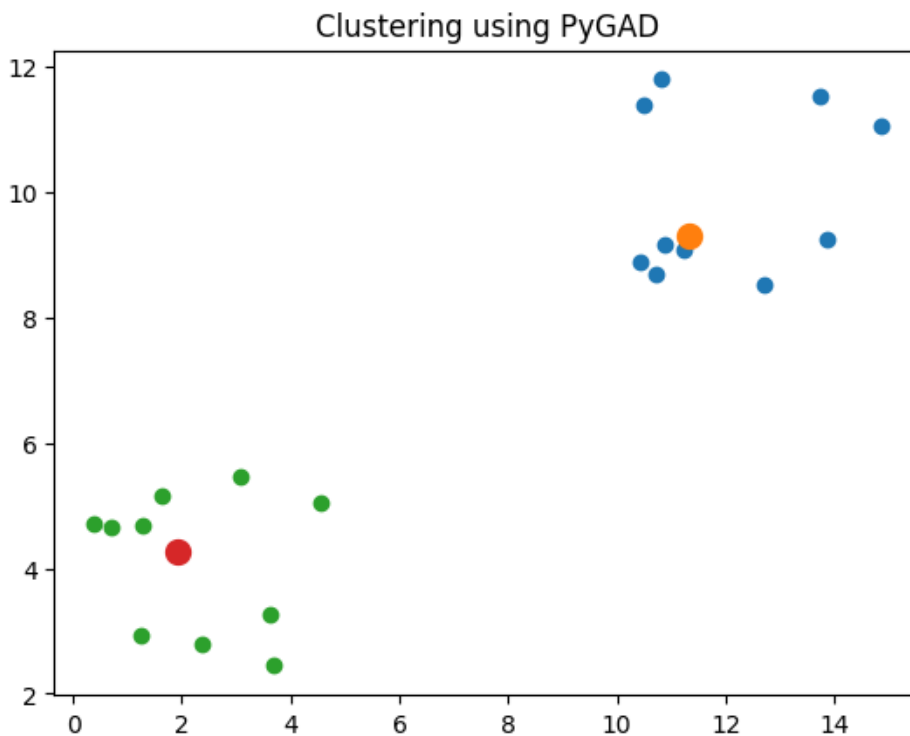
```

In [32]:

```

1 for cluster_idx in range(num_clusters):
2     cluster_x = data[clusters[cluster_idx], 0]
3     cluster_y = data[clusters[cluster_idx], 1]
4     matplotlib.pyplot.scatter(cluster_x, cluster_y)
5     matplotlib.pyplot.scatter(cluster_centers[cluster_idx, 0], cluster_centers[cluster_idx, 1], linewidth=2)
6 matplotlib.pyplot.title("Clustering using PyGAD")
7 matplotlib.pyplot.show()

```



In []:

```

1

```

