

การทดลองจัดกลุ่มข้อมูล 2 มิติ ด้วยวิธี k-Means และ Hierarchical Clustering

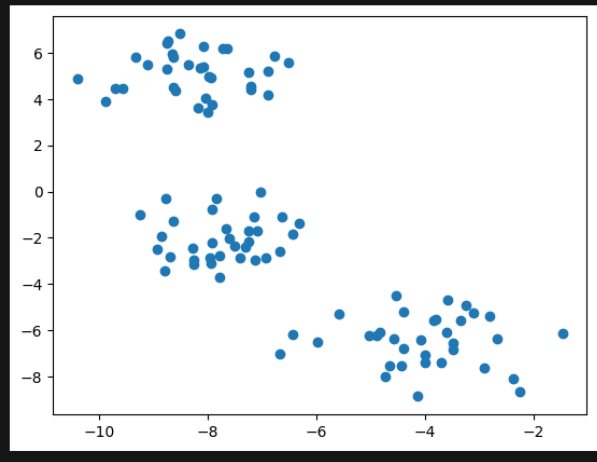
ชุดข้อมูลที่ 1

1. อ่านข้อมูล

```
data = pd.read_csv('data2Dset1.csv', header=None)
```

2. Plot จุดข้อมูล

```
x = data[0]
y = data[1]
plt.scatter(x,y)
plt.show()
```



3. จัดกลุ่มข้อมูลด้วย K-Means

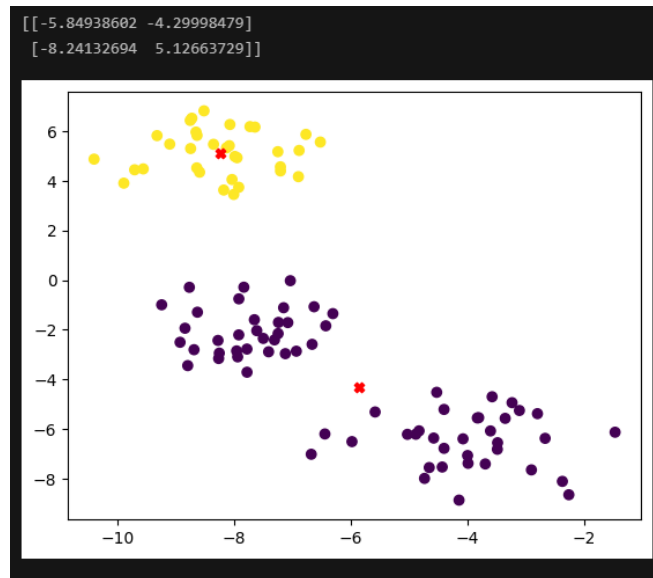
- a. k=2

```
model_kmeans = cluster.KMeans(n_clusters=2, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with k=2

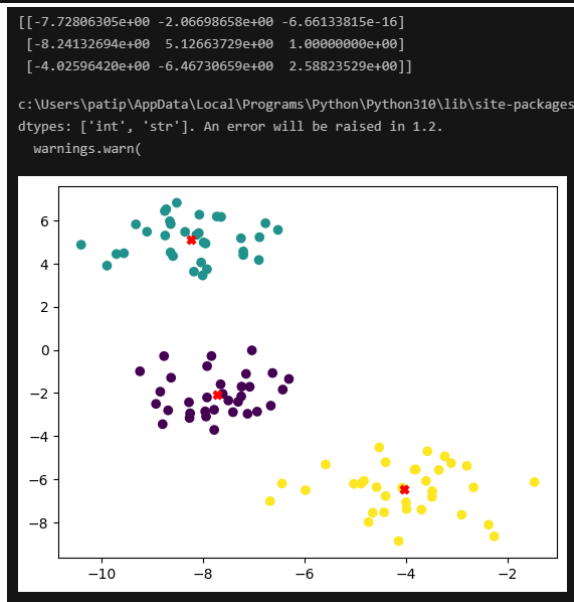
b. k=3

```
model_kmeans = cluster.KMeans(n_clusters=3, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with k=3

c. $k=4$

```

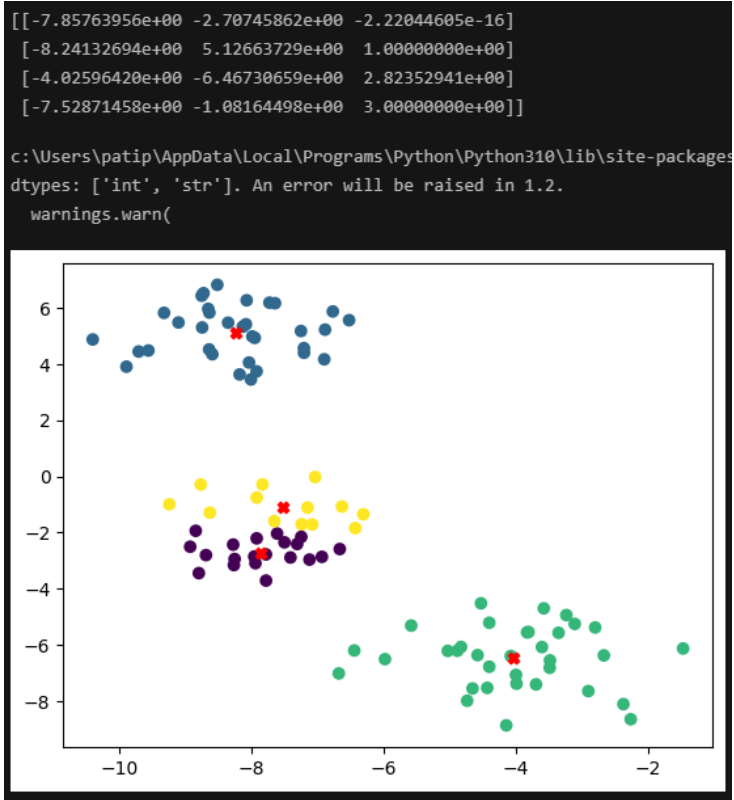
model_kmeans = cluster.KMeans(n_clusters=4, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()

```

Scatter plot with $k=4$

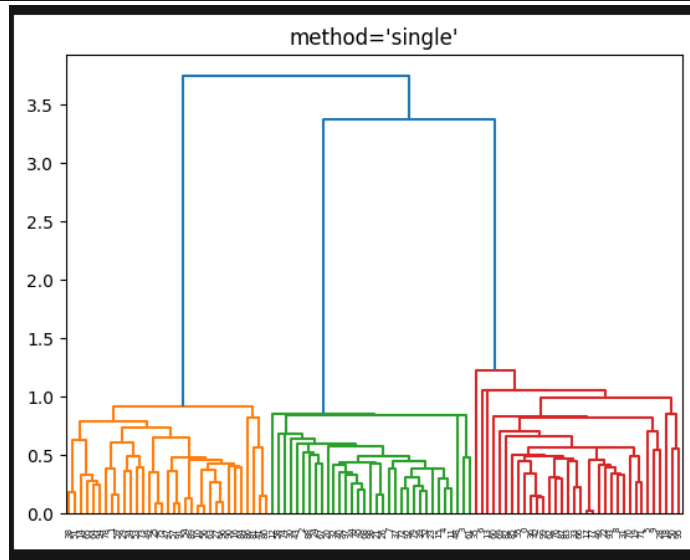
โดยสรุปจาก k ทั้ง 3 ค่าที่ได้ทดลองใช้ พบว่าค่า k ที่ดีที่สุดคือ 3 เพราะจากข้อมูลที่ได้มาเมื่อนำมา plot แล้วจะเห็นได้ว่าข้อมูลมีการแบ่งออกเป็น 3 กลุ่มชัดเจน และเมื่อใช้สีในการแยก cluster ก็จะได้เห็นว่าข้อมูลที่อยู่ในกราฟถูกแย่งแยกอย่างชัดเจนมากขึ้น

4. จัดกลุ่มข้อมูลด้วย Hierarchical Clustering

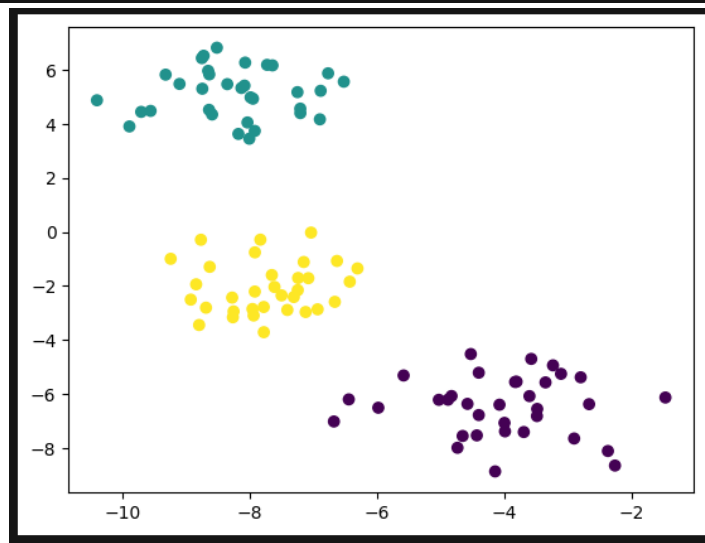
a. Method single

```
linkage_data = linkage(data, method='single', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='single'")
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=1.5, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```

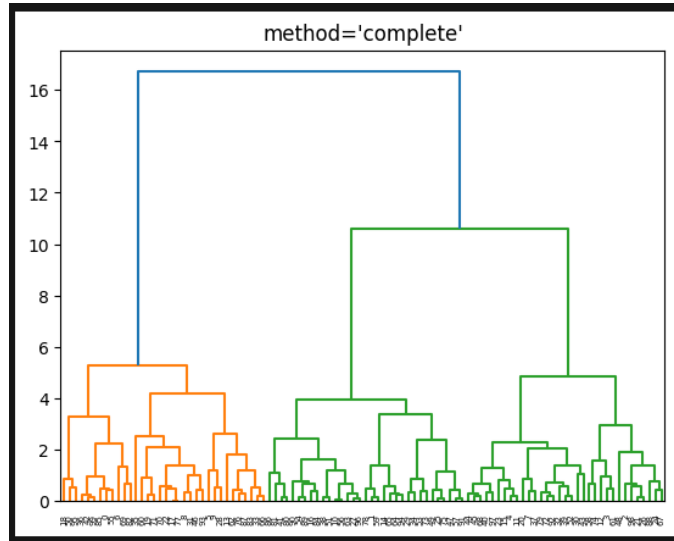


Cutoff โดยกำหนด $t=1.5$

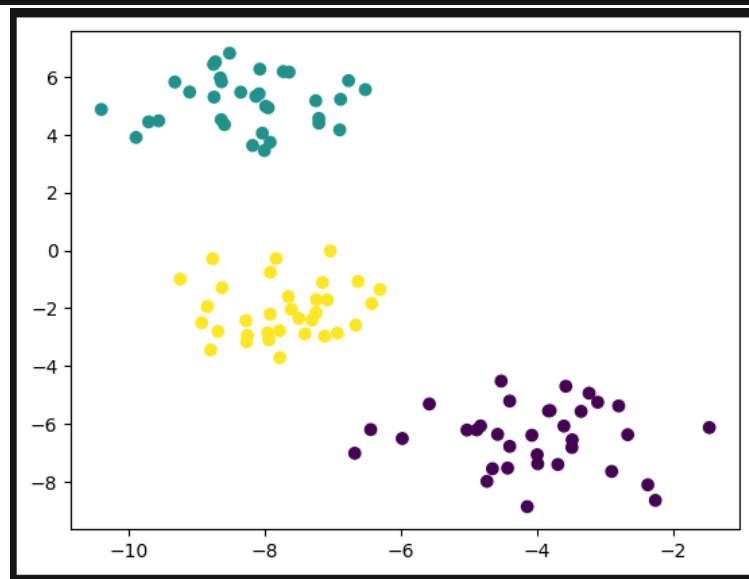
b. Method complete

```
linkage_data = linkage(data, method='complete', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='complete'")
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=6, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```

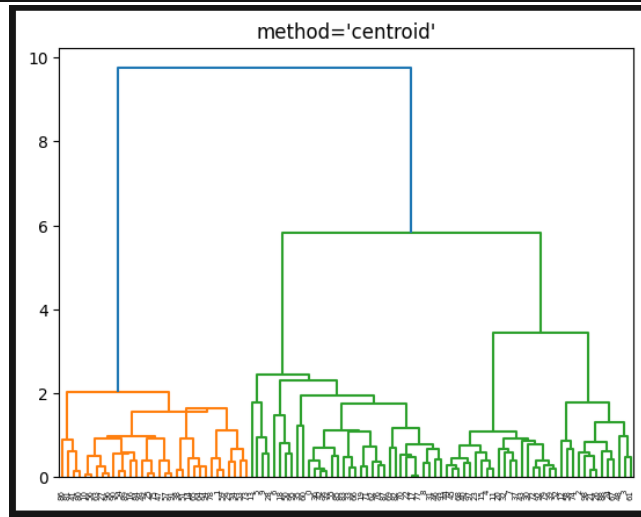


Cutoff โดยกำหนด $t=6$

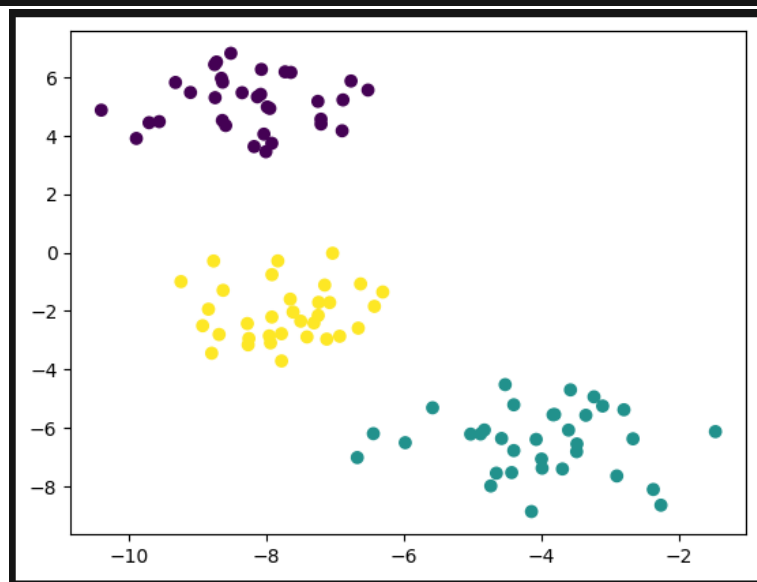
c. Method centroid

```
linkage_data = linkage(data, method='centroid', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='centroid'")
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=4, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```



Cutoff โดยกำหนด $t=4$

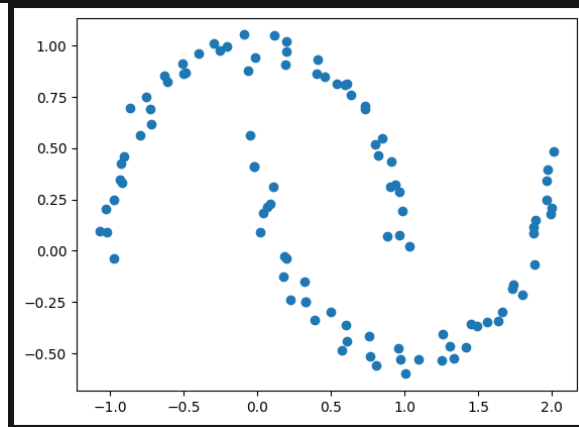
ชุดข้อมูลที่ 2

1. อ่านข้อมูล

```
data = pd.read_csv('data2Dset2.csv', header=None)
```

2. Plot จุดข้อมูล

```
x = data[0]
y = data[1]
plt.scatter(x,y)
plt.show()
```



3. จัดกลุ่มข้อมูลด้วย K-Means

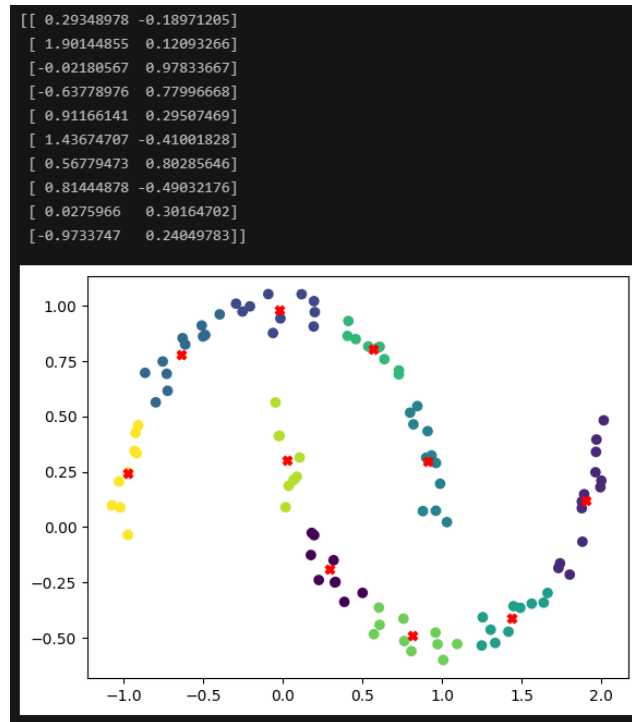
- a. k=10

```
model_kmeans = cluster.KMeans(n_clusters=10, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with $k=10$

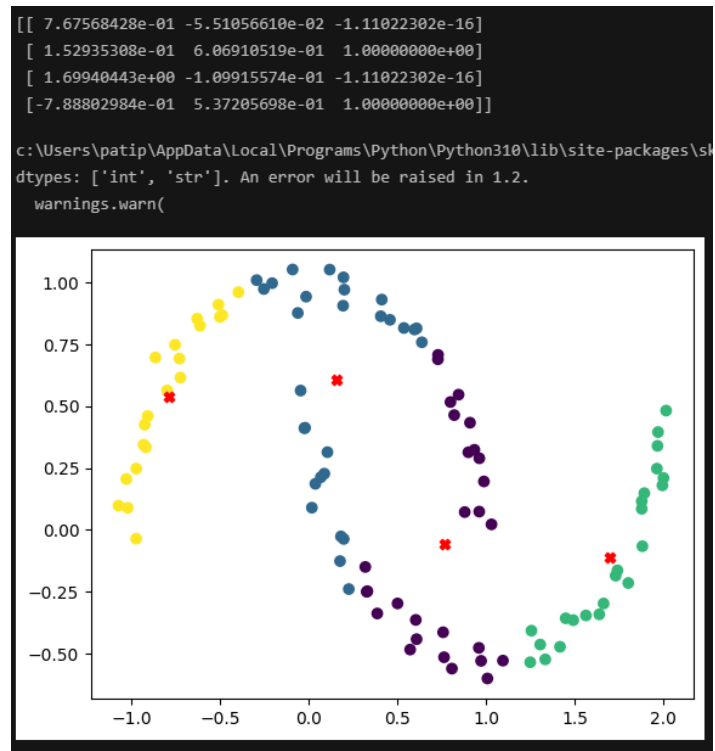
b. $k=4$

```
model_kmeans = cluster.KMeans(n_clusters=4, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```

Scatter plot with $k=4$

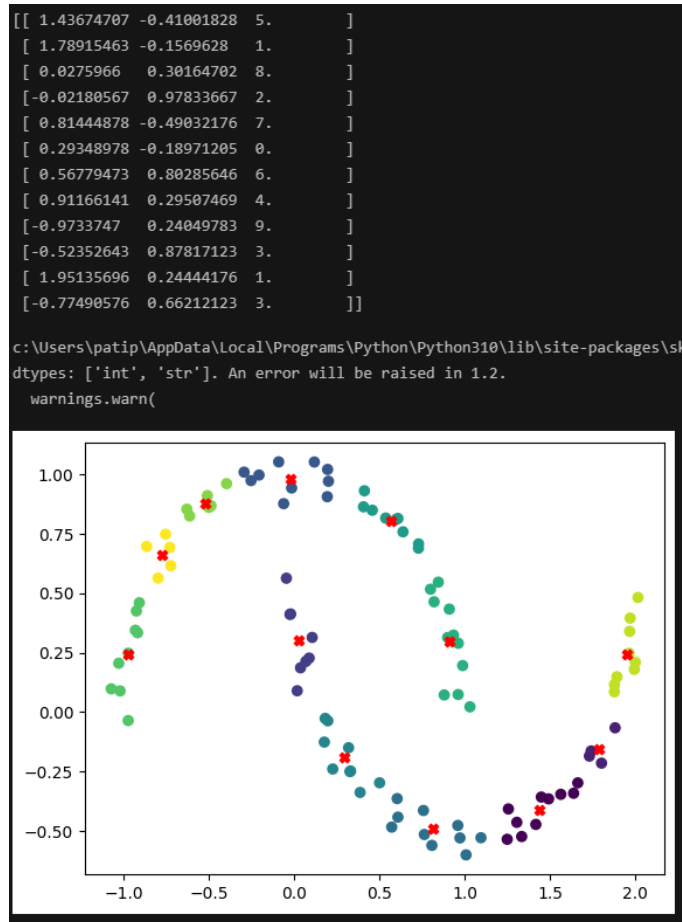
c. $k=12$

```
model_kmeans = cluster.KMeans(n_clusters=12, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with $k=12$

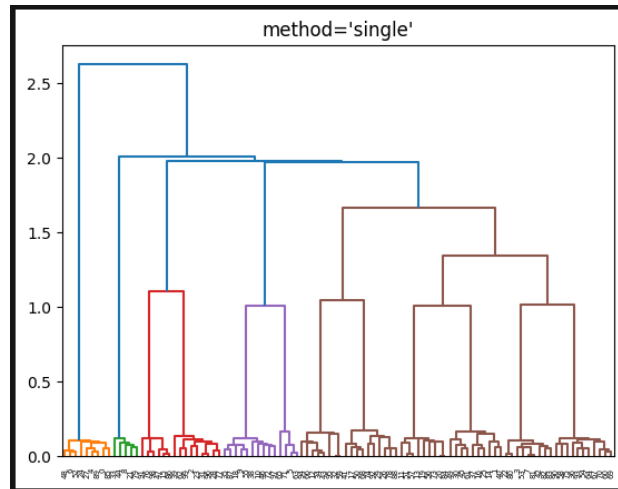
โดยสรุปแล้ว k ที่ทดลองนำมาใช้แล้วคิดว่าดีที่สุดคือ $k=10$ เนื่องจากหากมองที่ $k=4$ จะเห็นว่าถึงแม้จะสามารถแยกข้อมูลตามสีได้ชัด แต่กลุ่มแต่ละกลุ่มนั้นกว้างเกินไป และหากเป็น $k=12$ การแบ่งกลุ่มนั้นละเอียดเกินไปจนมองยาก ดังนั้น $k=10$ จึงเป็นตัวเลือกที่ดีที่สุดเพราะนอกจากจะแยกข้อมูลได้ไม่กว้างเกินไปแล้วก็ไม่ละเอียดเกินไปอีกด้วย

4. จัดกลุ่มข้อมูลด้วย Hierarchical Clustering

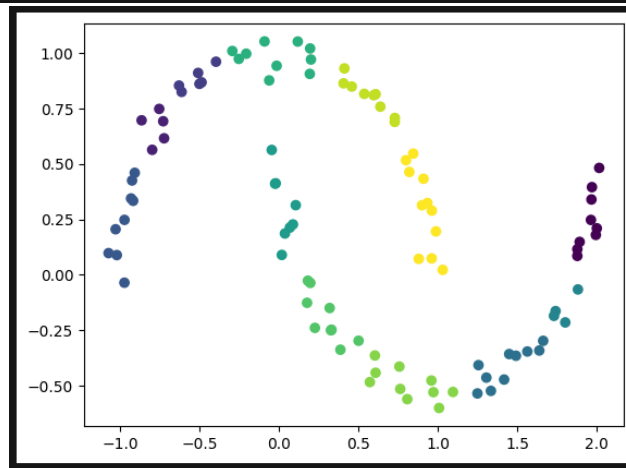
a. Method single

```
linkage_data = linkage(data, method='single', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='single'")
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=1.0, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```

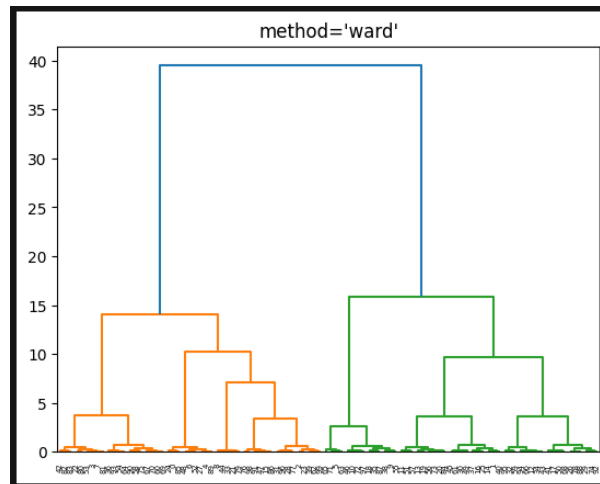


Cutoff ที่ $t=1$

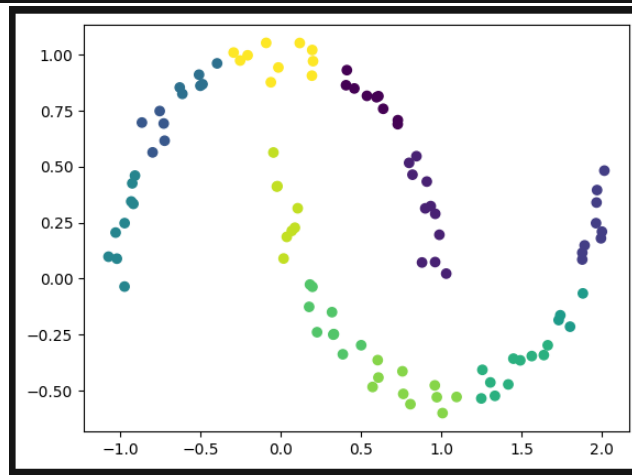
b. Method ward

```
linkage_data = linkage(data, method='ward', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='ward'")
plt.show()
```

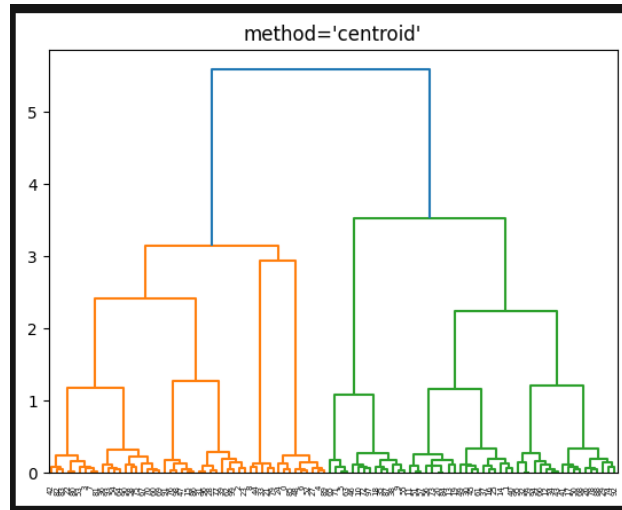


```
cluster_id = fcluster(linkage_data, t=2.5, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```

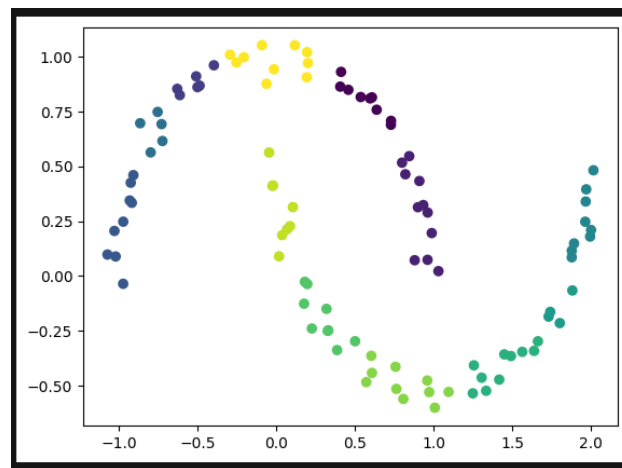


c. Method centroid

```
linkage_data = linkage(data, method='centroid', metric='euclidean')  
dendrogram(linkage_data)  
  
plt.title("method='centroid'")  
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=1, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```



Cutoff ที่ $t=1$

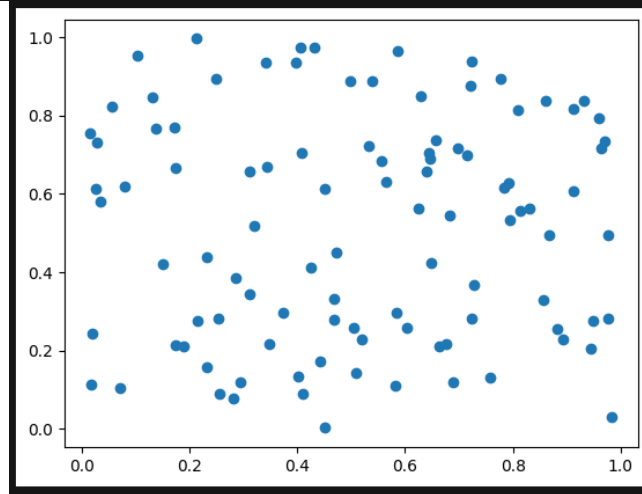
ชุดข้อมูลที่ 3

1. อ่านข้อมูล

```
data = pd.read_csv('data2Dset3.csv', header=None)
```

2. Plot จุดข้อมูล

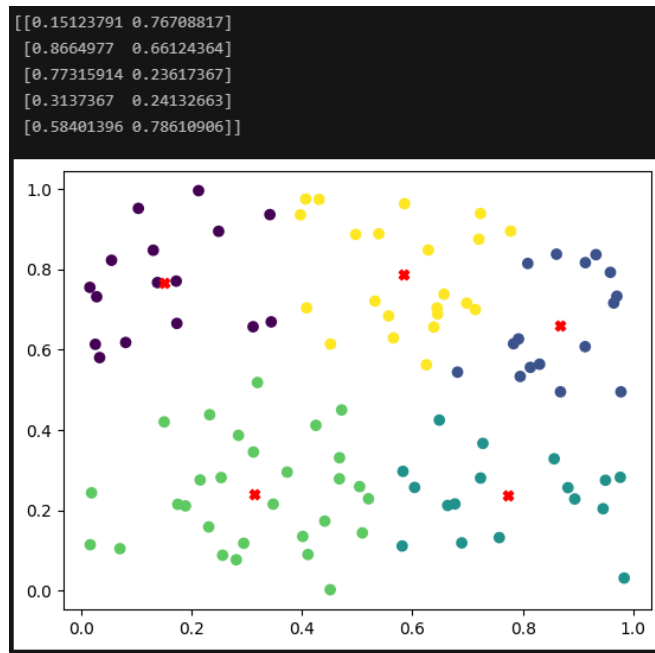
```
x = data[0]  
y = data[1]  
plt.scatter(x,y)  
plt.show()
```



3. จัดกลุ่มข้อมูลด้วย K-Means

- a. k=5

```
model_kmeans = cluster.KMeans(n_clusters=5, max_iter=50, random_state=1)  
model_kmeans.fit(data)  
  
data['cluster_id'] = model_kmeans.labels_  
  
centroids = model_kmeans.cluster_centers_  
print(centroids)  
  
plt.scatter(x,y, c=data['cluster_id'])  
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')  
plt.show()
```



Scatter plot with k=5

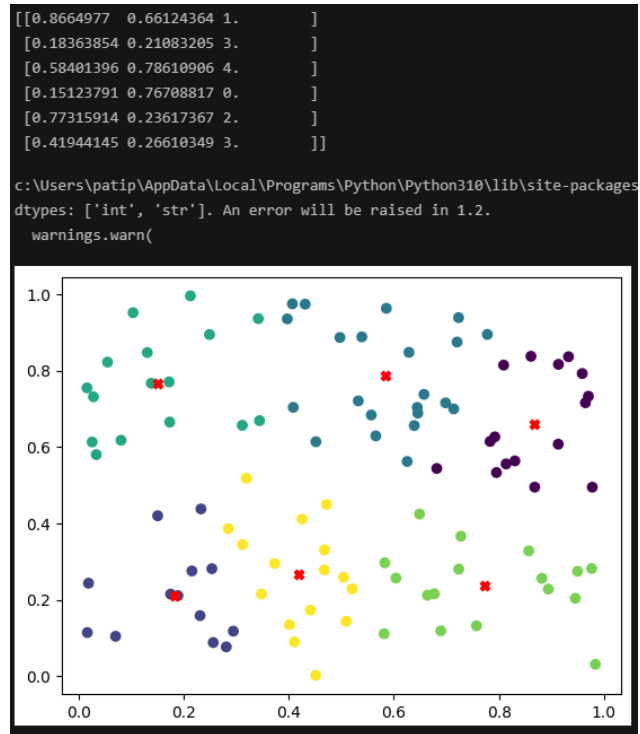
b. k=6

```
model_kmeans = cluster.KMeans(n_clusters=6, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with $k=6$

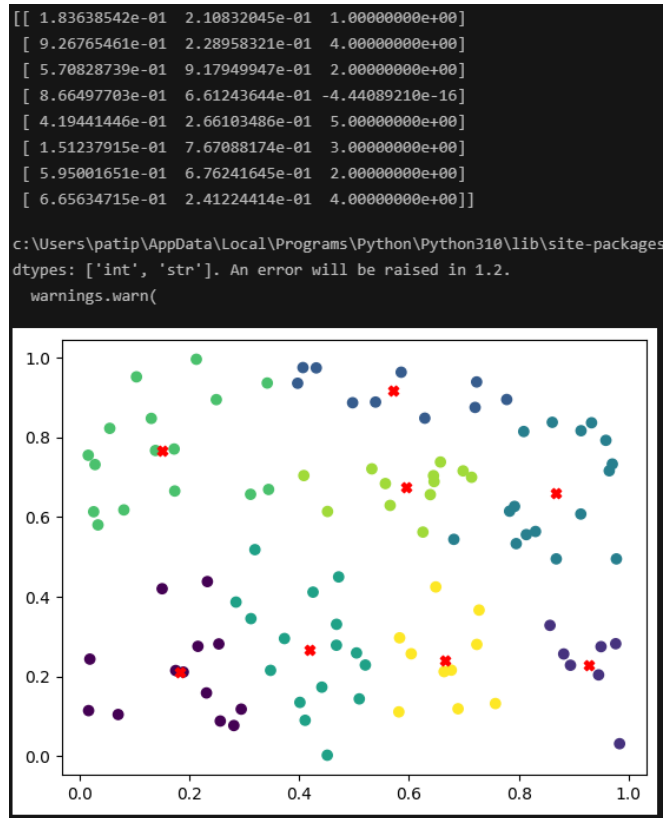
c. $k=8$

```
model_kmeans = cluster.KMeans(n_clusters=8, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```

Scatter plot with $k=8$

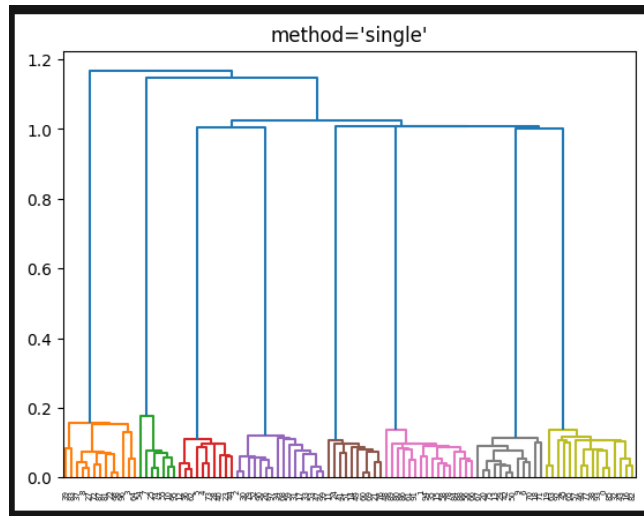
โดยสรุปแล้ว k ที่ดีที่สุดคือ $k=6$ ถึงแม้ข้อมูลที่มีจะดูกระจุกกระจาย แต่เมื่อลองใช้ $k=6$ แล้วจะเห็นข้อมูลที่จัดกลุ่มไว้ได้ชัดเจนและไม่กว้างไม่ละเอียดเกินไป หากเลข k มากกว่านี้จะได้ข้อมูลที่อ่านและวิเคราะห์ได้ยาก

4. จัดกลุ่มข้อมูลด้วย Hierarchical Clustering

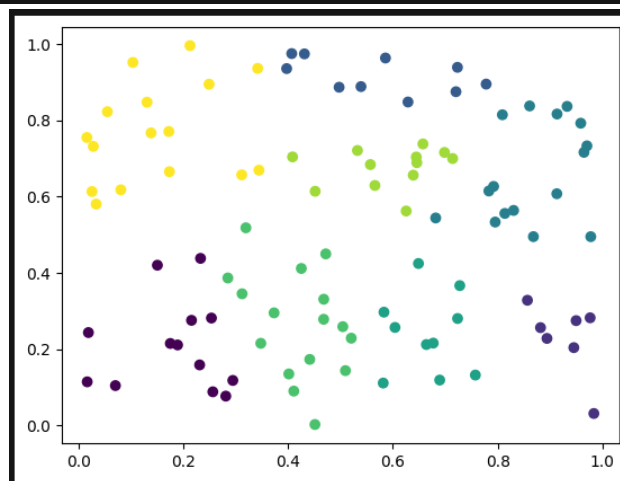
a. Method single

```
linkage_data = linkage(data, method='single', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='single'")
plt.show()
```



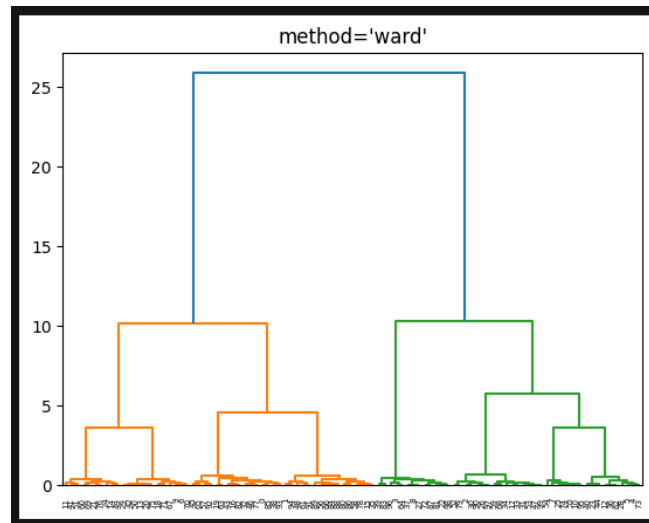
```
cluster_id = fcluster(linkage_data, t=0.2, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```



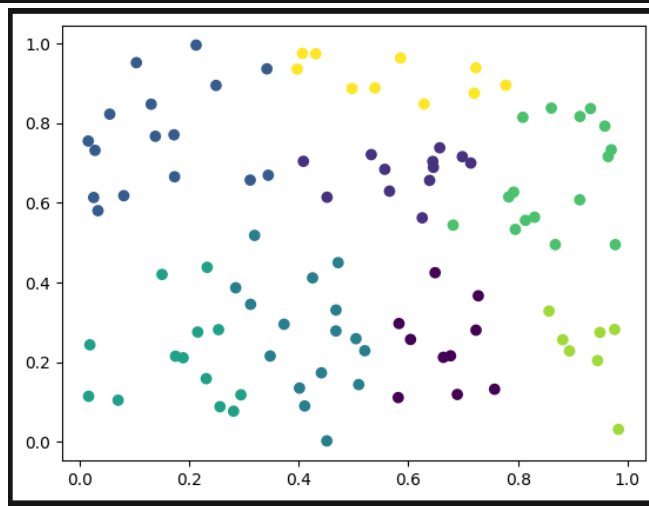
Cutoff ที่ $t=0.2$

b. Method ward

```
linkage_data = linkage(data, method='ward', metric='euclidean')  
dendrogram(linkage_data)  
  
plt.title("method='ward'")  
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=2.5, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```

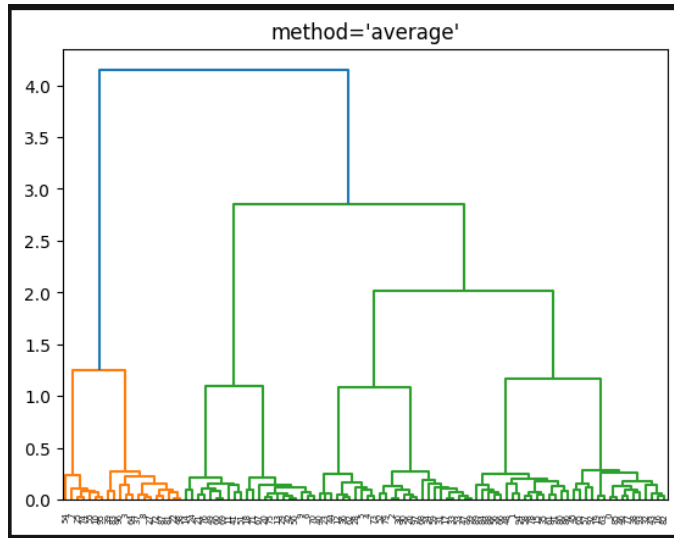


Cutoff ที่ $t=2.5$

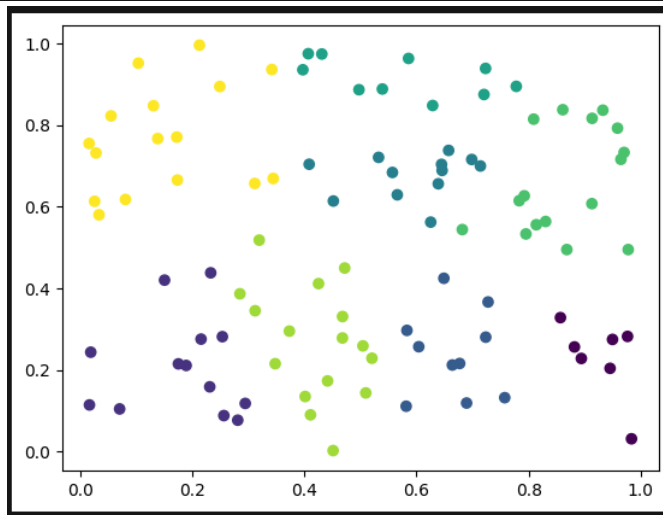
c. Method average

```
linkage_data = linkage(data, method='average', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='average'")
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=0.5, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```



Cutoff ที่ $t=0.5$

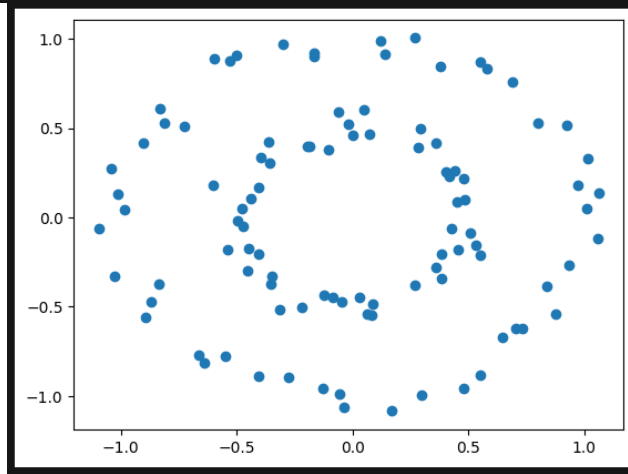
ชุดข้อมูลที่ 4

1. อ่านข้อมูล

```
data = pd.read_csv('data2Dset4.csv', header=None)
```

2. Plot จุดข้อมูล

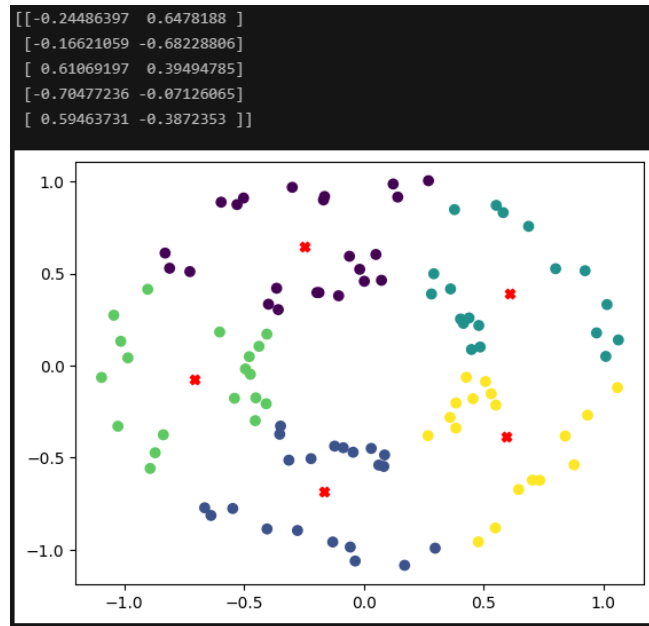
```
x = data[0]  
y = data[1]  
plt.scatter(x,y)  
plt.show()
```



3. จัดกลุ่มข้อมูลด้วย K-Means

- a. k=5

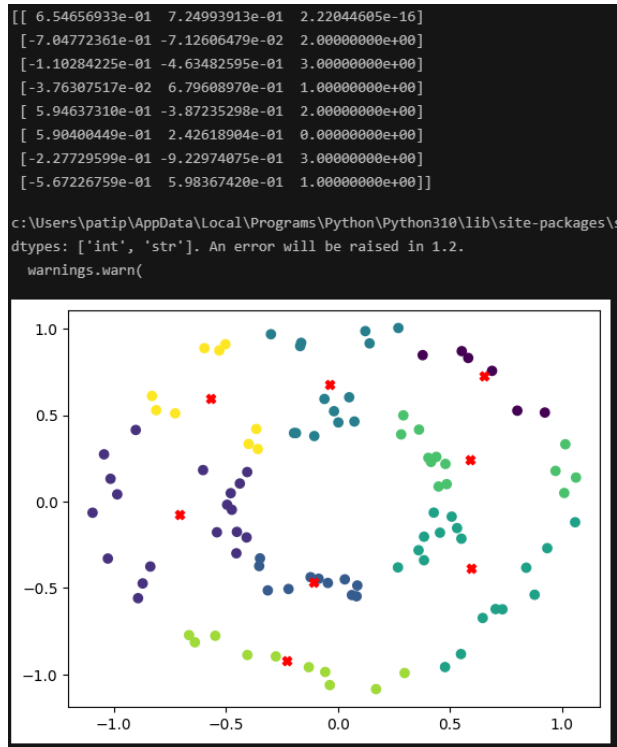
```
model_kmeans = cluster.KMeans(n_clusters=5, max_iter=50, random_state=1)  
model_kmeans.fit(data)  
  
data['cluster_id'] = model_kmeans.labels_  
  
centroids = model_kmeans.cluster_centers_  
print(centroids)  
  
plt.scatter(x,y, c=data['cluster_id'])  
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')  
plt.show()
```



Scatter plot with $k=5$

b. $k=8$

```
model_kmeans = cluster.KMeans(n_clusters=8, max_iter=50, random_state=1)  
model_kmeans.fit(data)  
  
data['cluster_id'] = model_kmeans.labels_  
  
centroids = model_kmeans.cluster_centers_  
print(centroids)  
  
plt.scatter(x,y, c=data['cluster_id'])  
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')  
plt.show()
```



Scatter plot with k=8

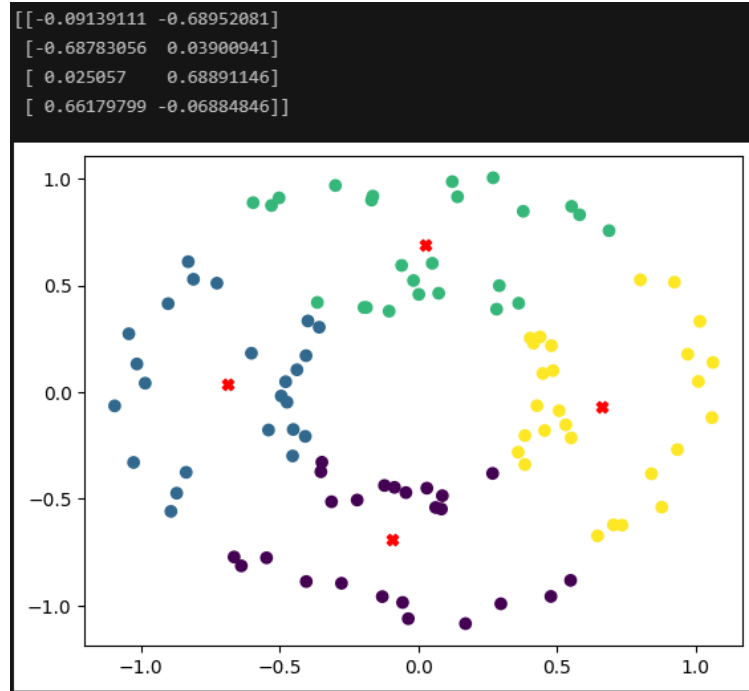
c. k=4

```
model_kmeans = cluster.KMeans(n_clusters=4, max_iter=50, random_state=1)
model_kmeans.fit(data)

data['cluster_id'] = model_kmeans.labels_

centroids = model_kmeans.cluster_centers_
print(centroids)

plt.scatter(x,y, c=data['cluster_id'])
plt.scatter(centroids[:,0],centroids[:,1],marker='X', c='r')
plt.show()
```



Scatter plot with $k=4$

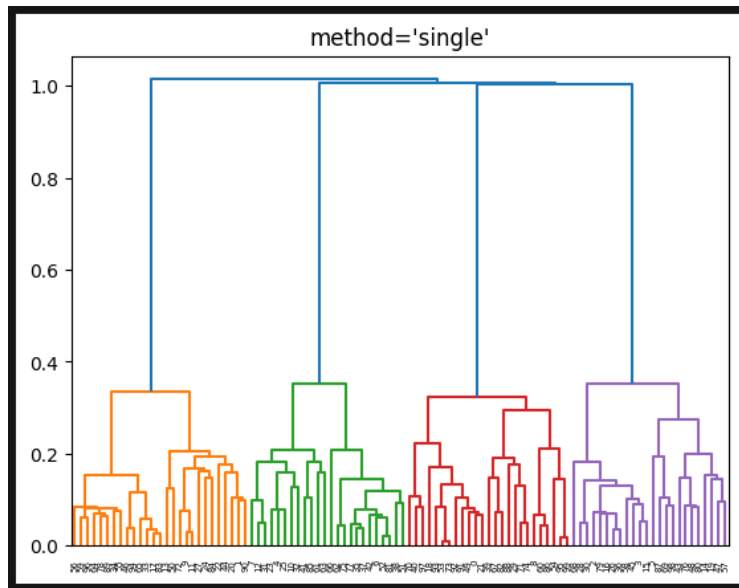
โดยสรุป k ที่ดีที่สุดคือ 4 เพราะข้อมูลแบ่งสัดส่วนกันสวยงามชัดเจน ดูง่าย แต่ถ้าหากมองว่าชุดข้อมูลกว้างเกินไปก็สามารถเลือก $k=5$ ได้เช่นกัน

4. จัดกลุ่มข้อมูลด้วย Hierarchical Clustering

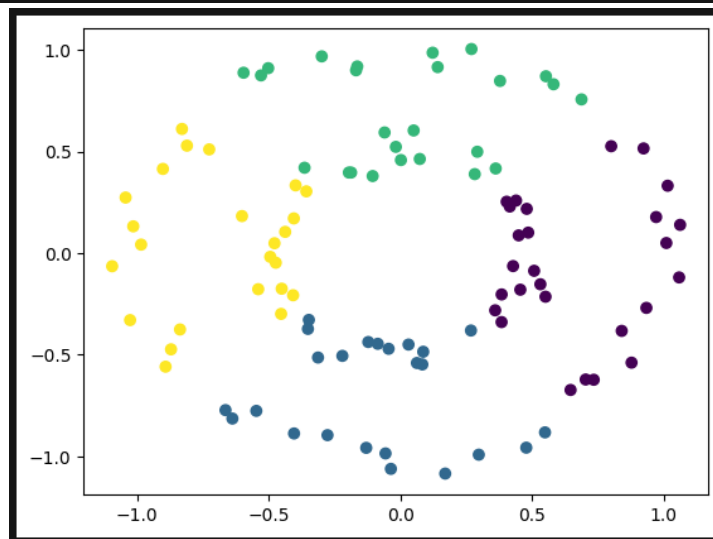
a. Method single

```
linkage_data = linkage(data, method='single', metric='euclidean')
dendrogram(linkage_data)

plt.title("method='single'")
plt.show()
```

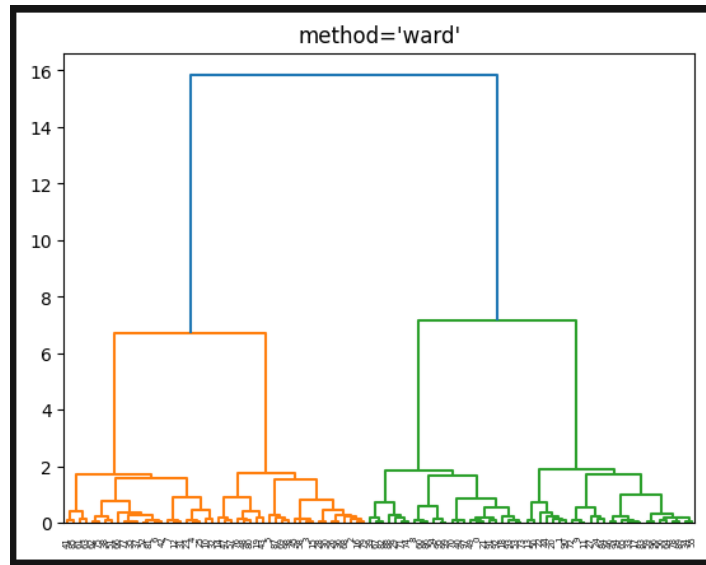
```
cluster_id = fcluster(linkage_data, t=0.4, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```



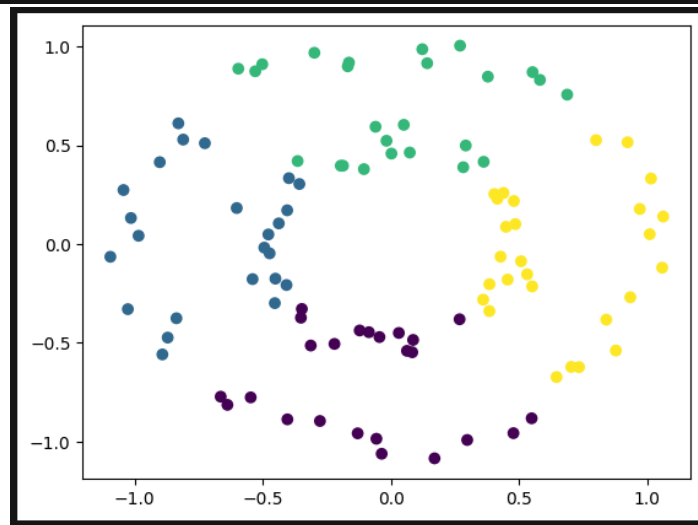
Cutoff ที่ $t=0.4$

b. Method ward

```
linkage_data = linkage(data, method='ward', metric='euclidean')  
dendrogram(linkage_data)  
  
plt.title("method='ward'")  
plt.show()
```



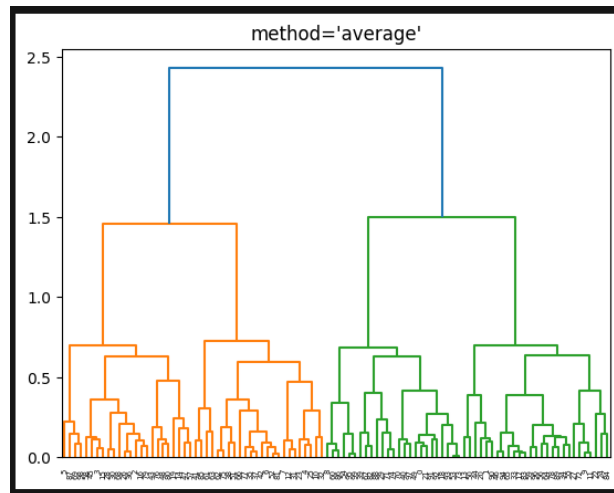
```
cluster_id = fcluster(linkage_data, t=2, criterion='distance')  
plt.scatter(x,y, c=cluster_id)  
plt.show()
```



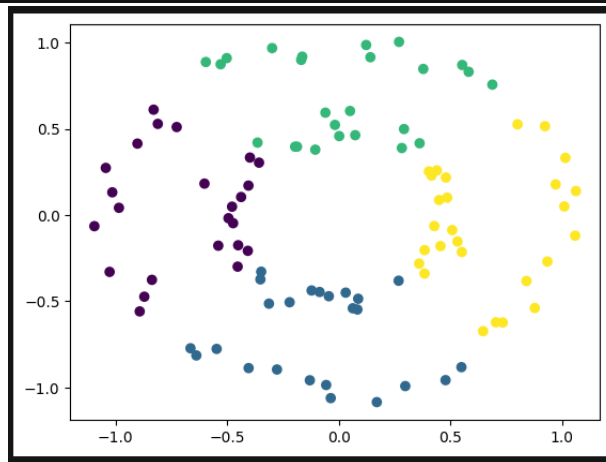
Cutoff ที่ $t=2$

c. Method average

```
linkage_data = linkage(data, method='average', metric='euclidean')  
dendrogram(linkage_data)  
  
plt.title("method='average'")  
plt.show()
```



```
cluster_id = fcluster(linkage_data, t=1, criterion='distance')
plt.scatter(x,y, c=cluster_id)
plt.show()
```



Cutoff ที่ $t=1$

สรุปผลการทดลอง

K-means นั้นเป็นการระบุจำนวน cluster ที่เราต้องการจะแยกหมวดหมู่ ซึ่งวิธีนี้ควรใช้กับชุดข้อมูลที่เราเห็น pattern ชัดเจน สามารถระบุจำนวน cluster หรือ k ได้ง่าย ซึ่งการที่เราสามารถระบุ k ที่เราต้องการลงไปได้เลยนั้นก็คือข้อดีของ K-means แต่ข้อเสียคือค่อนข้างเสีย performance ในการทำงาน ใช้เวลามาก

ส่วน Hierarchical Clustering จากการทดลอง จะเป็นการลองใช้ method ที่ละแบบ ซึ่งแต่ละแบบอาจจะได้ผลลัพธ์หรือการทำ dendrogram ได้ผลต่างหรือใกล้เคียงกันขึ้นอยู่กับข้อมูล แต่ทุกแบบจะทำให้เราเห็น hierarchy ที่ชัดเจนเห็น pattern ได้ดีขึ้น และเมื่อเราเห็นข้อมูลที่ถูกจัดอยู่ในรูปแบบของ dendrogram แล้ว เราก็สามารถที่จะวิเคราะห์เพื่อเลือกค่า t ที่จะนำไปใช้ cutoff เพื่อทำให้เกิด cluster ที่เหมาะสมได้ ซึ่งหมายความว่า

62130500048 นางสาวปณิญา ทองอ่วม

วิธีนี้จะใช้ได้ดีก็ต่อเมื่อข้อมูลมีจำนวนไม่เยอะมาก หากข้อมูลมีเยอะเกินไป จะทำให้วิเคราะห์ค่า t ที่เหมาะสมได้ยาก