

Clustering

Speaker: Hong-Han Shuai
ECE, NCTU

What is Cluster Analysis?

- **Cluster**: a collection of data objects

Similar to one another within the same cluster

Dissimilar to the objects in other clusters

- **Cluster Analysis**

Grouping a set of data objects into clusters

- Typical applications:

As a stand-alone tool to get insight into data distribution

As a preprocessing step for other algorithms

Unsupervised learning

General Applications of Clustering

➤ Spatial Data Analysis

Create thematic maps in GIS by clustering feature spaces.

Detect spatial clusters and explain them in spatial data mining.

➤ Image Processing

➤ Pattern Recognition

➤ Economic Science (especially market research)

➤ WWW

Document classification

Cluster Web-log data to discover groups of similar access patterns

Examples of Clustering Applications

Marketing: Help marketers discover distinct groups in their customer bases, and then use this knowledge to develop targeted marketing programs.

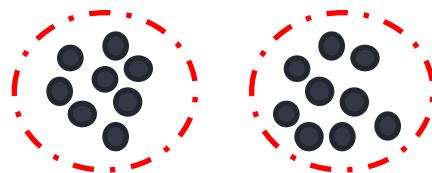
Land use: Identification of areas of similar land use in an earth observation database.

Insurance: Identifying groups of motor insurance policy holders with a high average claim cost.

City-planning: Identifying groups of houses according to their house type, value, and geographical location.

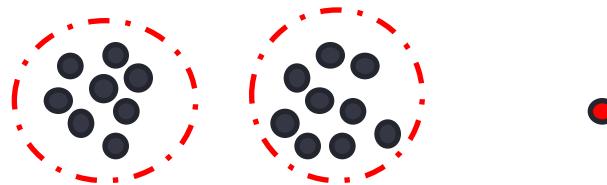
What is Good Clustering?

- A good clustering method will produce high quality clusters with
 - High intra-class similarity
 - Low inter-class similarity
- The quality of a clustering result depends on both the **similarity measure used by the method** and its **implementation**.
- The quality of a clustering method is also measured by its ability to discover hidden patterns.



Requirements of Clustering in Data Mining (1/2)

- Scalability
- Ability to deal with different types of attributes
- Discovery of clusters with arbitrary shape
- Minimal requirements of domain knowledge for input
- Able to deal with outliers



Requirements of Clustering in Data Mining (2/2)



- Inensitive to order of input records
- High dimensionality
- Curse of dimensionality
- Incorporation of user-specified constraints
- Interpretability and usability

Clustering Methods (I)

Partitioning Method

Construct various partitions and then evaluate them by some criterion, e.g., minimizing the sum of square errors

K-means, k-medoids, CLARANS

Hierarchical Method

Create a hierarchical decomposition of the set of data (or objects) using some criterion

Diana, Agnes, BIRCH, ROCK, CHAMELEON

Density-based Method

Based on connectivity and density functions

Typical methods: DBSACN, OPTICS, DenClue

Clustering Methods (II)

➤ Grid-based approach

based on a multiple grid

Typical methods: SNN

➤ Model-based approach

A model is hypothesized for each cluster

model to each other

Typical methods: EM

➤ Frequent pattern-based approach

Based on the analysis of frequent patterns

Typical methods: p-Means

➤ User-guided clustering

Clustering by considering user-specified or application-specific constraints

Typical methods: cluster-on-demand, constrained clustering

	C_0	C_1	C_2	C_3	C_4
O_0	995	942	827	436	391
O_1	359	604	902	963	465
O_2	379	153	292	983	485
O_3	382	421	716	718	895
O_4	259	447	726	863	365

best fit of that

Calculate the Distance between **Clusters**

- **Single link**: smallest distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \min(t_{ip}, t_{jq})$
- **Complete link**: largest distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \max(t_{ip}, t_{jq})$
- **Average**: average distance between an element in one cluster and an element in the other, i.e., $\text{dis}(K_i, K_j) = \text{avg}(t_{ip}, t_{jq})$
- **Centroid**: distance between the centroids of two clusters,
i.e., $\text{dis}(K_i, K_j) = \text{dis}(C_i, C_j)$
- **Medoid**: distance between the medoids of two clusters,
i.e., $\text{dis}(K_i, K_j) = \text{dis}(M_i, M_j)$
- Medoid: one chosen, **centrally located object** in the cluster

Centroid, Radius and Diameter of a Cluster (for numerical data sets)

Centroid: the “middle” of a cluster

$$C_m = \frac{\sum_{i=1}^N (t_{ip})}{N}$$

Radius: square root of average mean squared distance from **any point** of the cluster to its centroid

$$R_m = \sqrt{\frac{\sum_{i=1}^N (t_{ip} - c_m)^2}{N}}$$

Diameter: square root of average mean squared distance between **all pairs** of points in the cluster

$$D_m = \sqrt{\frac{\sum_{i=1}^N \sum_{j=1}^N (t_{ip} - t_{jq})^2}{N(N-1)}}$$

diameter != 2 * radius

Partitioning Algorithms: Basic Concept

- **Partitioning method:** construct a partition of a database D of n objects into a set of k clusters.
Given a number k , find a partition of k clusters that optimizes the chosen partitioning criterion.
- **Global optimal:** exhaustively enumerate all partitions.
- **Heuristic methods:** k-means, k-medoids
 - k-means (MacQueen'67)
 - k-medoids or PAM, partition around medoids (Kaufman & Rousseeuw'87)

The K-Means Clustering Method

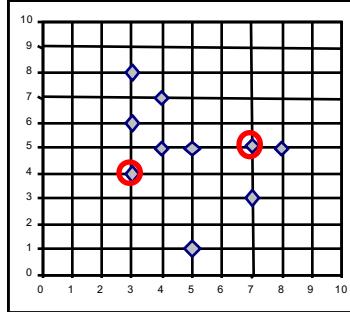
Given k , the k-means algorithm is implemented in four steps:

- 1. Arbitrarily choose k points as initial cluster centroids.
- 2. **Update Means (Centroids)**: Compute seed points as the center of the clusters of the current partition. (center: mean point of the cluster)
- 3. **Re-assign Points**: Assign each object to the cluster with the nearest seed point.
- 4. Go back to Step 2, stop when no more new assignment.

loop



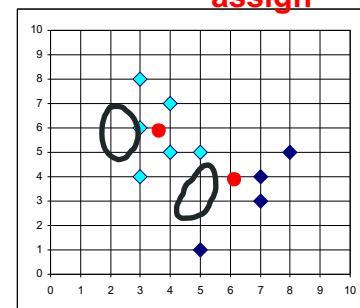
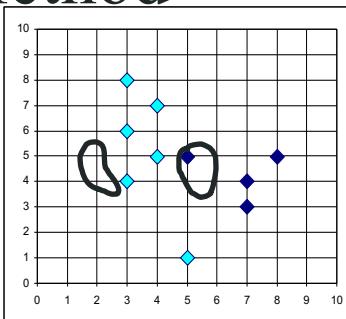
Example of the K-Means Clustering Method



Given $k = 2$:

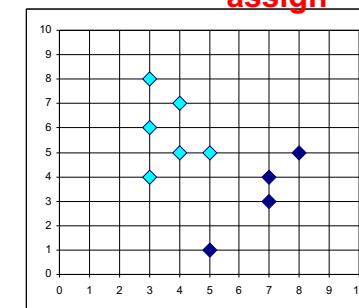
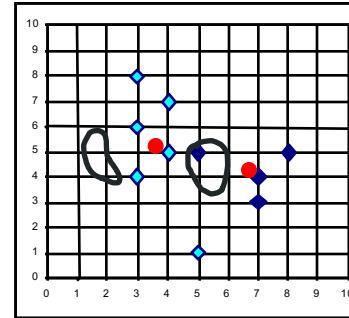
Arbitrarily choose k objects as initial cluster centroids

Assign each object to the most similar centroid



Re-assign

Update the cluster means



Re-assign

Update the cluster means

Comments on the K-Means Clustering

Time Complexity: $O(tkn)$, where n is # of objects, k is # of clusters, and t is # of iterations. Normally, $k,t \ll n$.

Often terminates at a local optimum.

(The global optimum may be found using techniques such as: deterministic annealing and genetic algorithms)

Weakness:

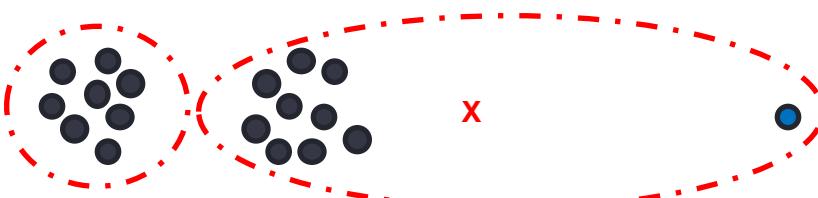
Applicable only when mean is defined, how about categorical data?

Need to specify k , the number of clusters, in advance
Unable to handle noisy data and outliers

Why is K-Means Unable to Handle Outliers?

The k-means algorithm is sensitive to outliers

Since an object with an extremely large value may substantially distort the distribution of the data.



K-Medoids: Instead of taking the mean value of the object in a cluster as a reference point, medoids can be used, which is the most centrally located object in a cluster.

PAM: The K-Medoids Method

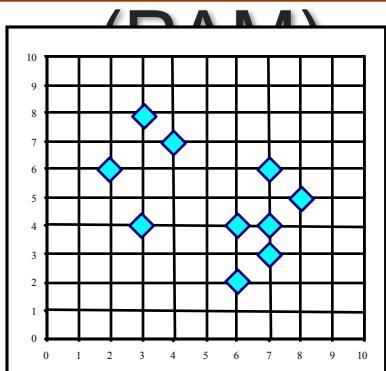
PAM: Partition Around Medoids

Use real object to represent the cluster

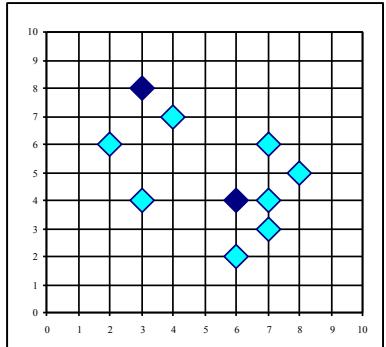
loop

- 1. Randomly select k representative objects as medoids.
- 2. Assign each data point to the closest medoid.
- 3. For each medoid m ,
 - a. For each non-medoid data point o
 - b. Swap m and o , and compute the total cost of the configuration.
- 4. **Select the configuration with the lowest cost.**
- 5. Repeat steps 2 to 5 until there is no change in the medoid.

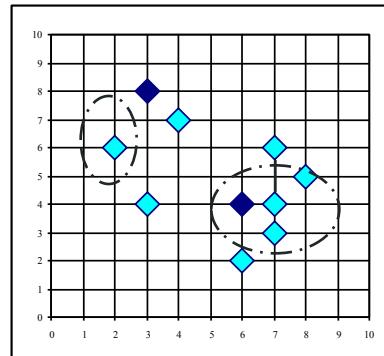
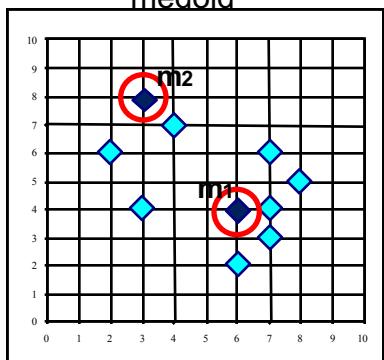
A Typical K-Medoids Algorithm



Arbitrarily choose k object as initial medoid



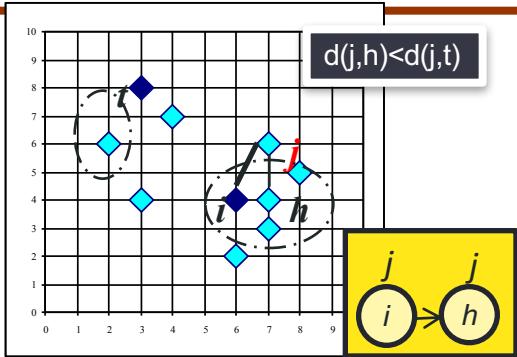
Assign each remaining object to the nearest medoid



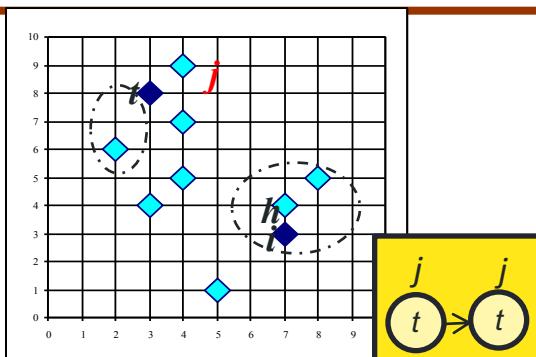
Swap each medoid and each data point, and compute the total cost of the configuration

PAM Clustering: Total swapping cost $TC_{ih} = \sum_j C_{jih}$

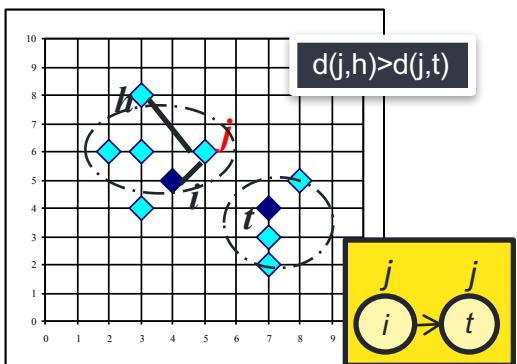
- Original medoid: t, i
- h: swap with i**
- j: any non-selected object



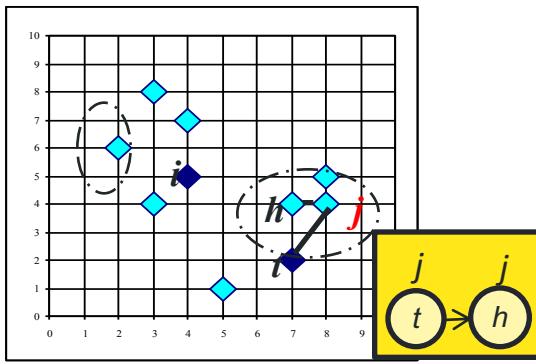
$$C_{jih} = d(j, h) - d(j, i)$$



$$C_{jih} = 0$$



$$C_{jih} = d(j, t) - d(j, i)$$



$$C_{jih} = d(j, h) - d(j, t)$$

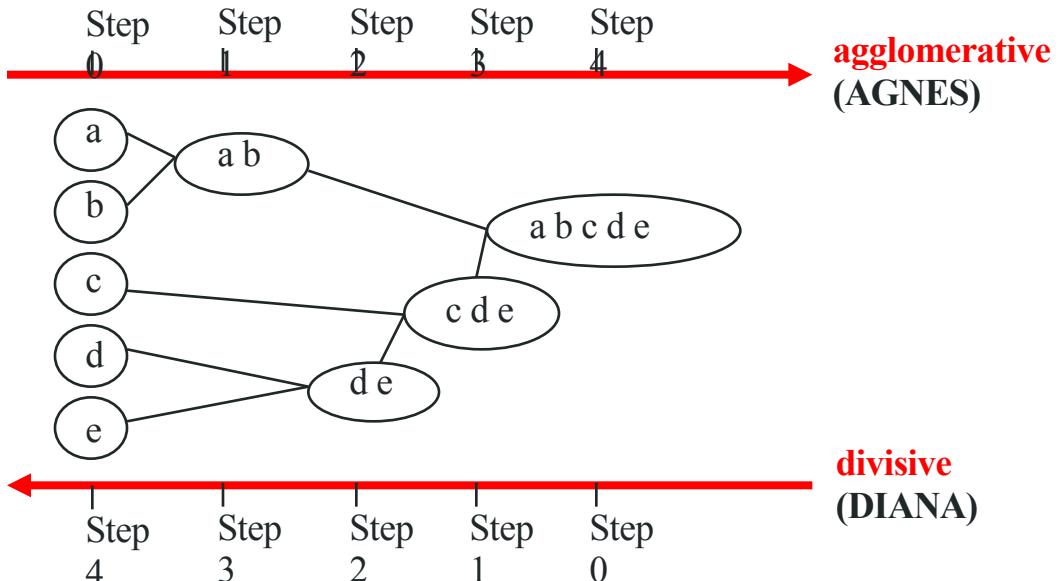
What is the Problem with PAM?

- PAM is more robust than k-means in the presence of noise and outliers because a medoid is less influenced by outliers or other extreme values than a mean.
- PAM works efficiently for small data sets but does not scale well for large data sets.
- $O(k(n-k)(n-k))$ for each iteration, where n is # of data, k is # of clusters
- Improvements: CLARA (uses a sampled set to determine medoids), CLARANS

Hierarchical Clustering

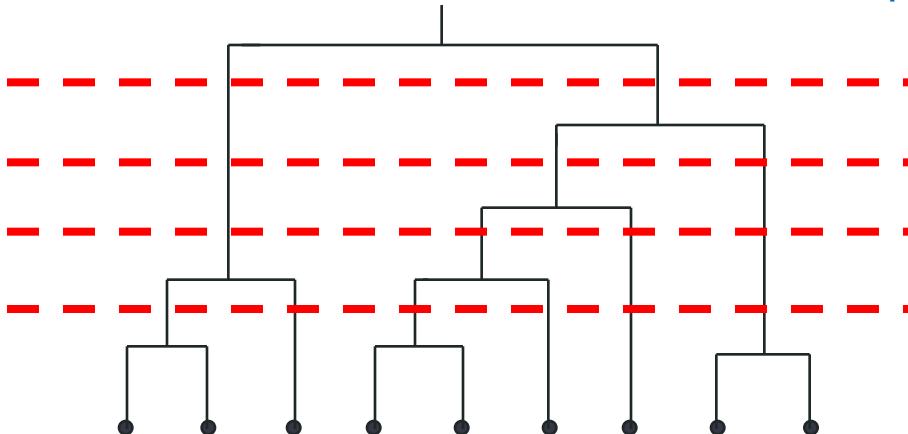
Use distance matrix as clustering criteria.

This method does not require the number of clusters k as an input, but needs a termination condition.



Dendrogram: Shows How the Clusters are Merged

- Decompose data objects into a several levels of nested partitioning (tree of clusters), called a dendrogram.
- A clustering of the data objects is obtained by cutting the dendrogram at the desired level, then each connected component forms a cluster.



More on Hierarchical Clustering

Major weakness:

Do not scale well: time complexity is at least $O(n^2)$, where n is the number of total objects.

Can never undo what was done previously.

Integration of hierarchical with distance-based clustering

BIRCH(1996): uses CF-tree data structure and incrementally adjusts the quality of sub-clusters.

CURE(1998): selects well-scattered points from the cluster and then shrinks them towards the center of the cluster by a specified fraction.

Density-Based Clustering Methods

Clustering based on **density (local cluster criterion)**, such as density-connected points

Major features:

Discover clusters of arbitrary shape

Handle noise

One scan

Need density parameters as termination condition

Several interesting studies:

DBSCAN: Ester, et al. (KDD'96)

OPTICS: Ankerst, et al (SIGMOD'99).

DENCLUE: Hinneburg & D. Keim (KDD'98)

CLIQUE: Agrawal, et al. (SIGMOD'98) (more grid-based)

DBSCAN: Density Based Spatial Clustering of Applications with Noise

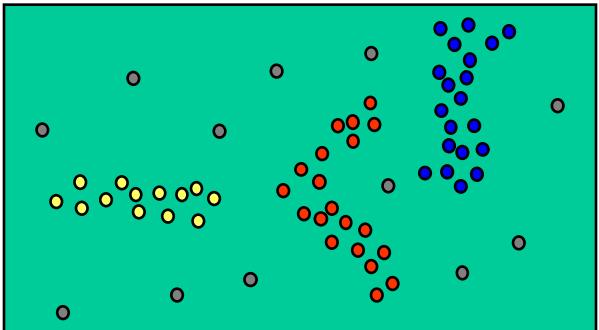


- Proposed by Ester, Kriegel, Sander, and Xu (KDD96)
- Relies on a density-based notion of cluster: A cluster is defined as a maximal set of density-connected points.
- Discovers clusters of arbitrary shape in spatial databases with noise

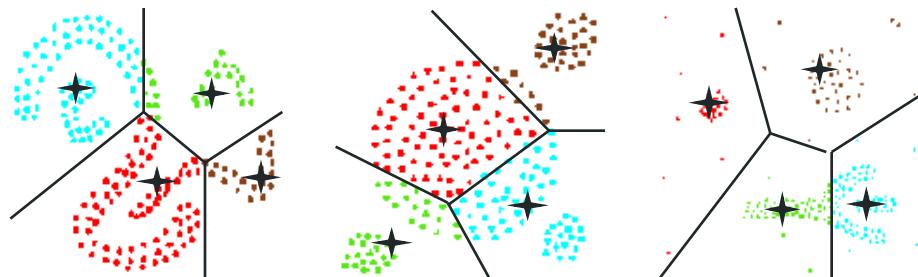
Density-Based Clustering

★ Basic Idea:

Clusters are dense regions in the data space, separated by regions of lower object density



Why Density-Based Clustering?



Results of a k -medoid algorithm for $k=4$

Different density-based approaches exist (see Textbook & Papers)
Here we discuss the ideas underlying the DBSCAN algorithm

Density Based Clustering: Basic Concept

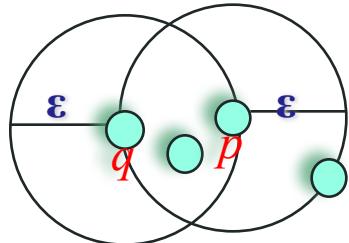
- Intuition for the formalization of the basic idea
- For any point in a cluster, the local point density around that point has to exceed some threshold
- The set of points from one cluster is spatially connected
- Local point density at a point p defined by two parameters
- ε – radius for the neighborhood of point p :
$$N_\varepsilon(p) := \{q \text{ in data set } D \mid dist(p, q) \leq \varepsilon\}$$
- $MinPts$ – minimum number of points in the given neighbourhood $N(p)$

ε -Neighborhood

- ε -Neighborhood – Objects within a radius of ε from an object.

$$N_\varepsilon(p) : \{q \mid d(p, q) \leq \varepsilon\}$$

- “High density” - ε -Neighborhood of an object contains at least *MinPts* of objects.



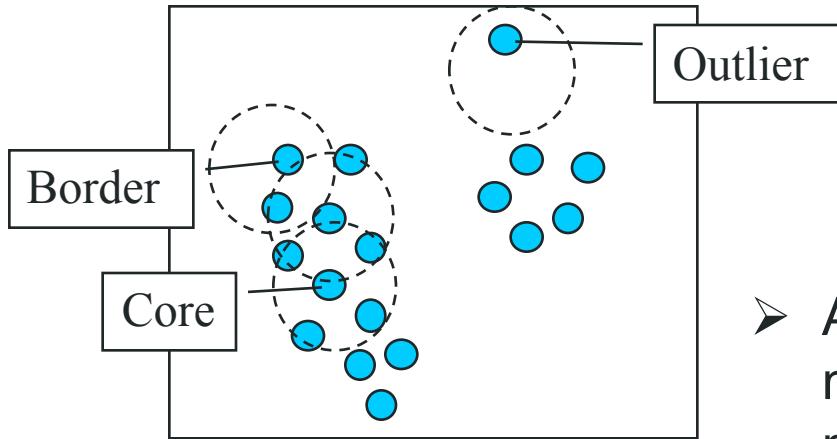
ε -Neighborhood of p

ε -Neighborhood of q

Density of p is “high” (*MinPts* = 4)

Density of q is “low” (*MinPts* = 4)

Core, Border & Outlier



$$\epsilon = 1 \text{ unit}, \text{MinPts} = 5$$

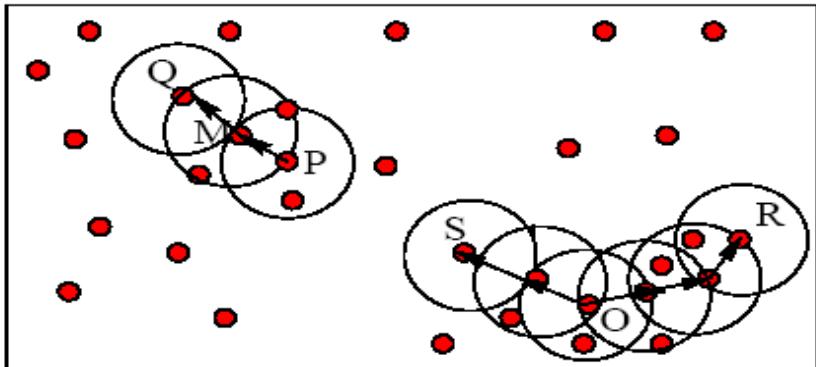
- A **border point** has fewer than MinPts within Eps , but is in the neighborhood of a core point.
- A **noise point** is any point that is not a core point nor a border point.

Given ϵ and MinPts , categorize the objects into three exclusive groups.

- A point is a **core point** if it has more than a specified number of points (MinPts) within Eps . These are points that are at the interior of a cluster.

Example

M, P, O, and R are core objects since each is in an Eps neighborhood containing at least 3 points



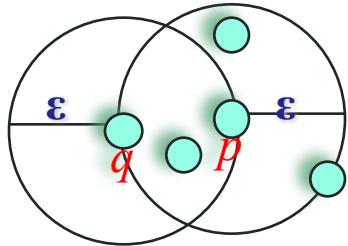
Minpts = 3

Eps=radius
of the circles

Density-Reachability

■ Directly density-reachable

- An object q is directly density-reachable from object p if p is a core object and q is in p 's ε -neighborhood.

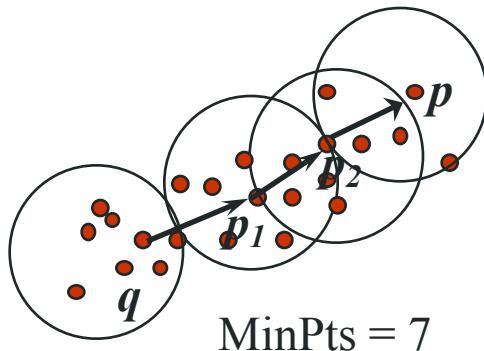


MinPts = 4

- q is directly density-reachable from p
- p is not directly density-reachable from q ?
- Density-reachability is asymmetric.

Density-reachability

- Density-Reachable (directly and indirectly):
- A point p is directly density-reachable from p_2 ;
- p_2 is directly density-reachable from p_1 ;
- p_1 is directly density-reachable from q ;
- $p \leftarrow p_2 \leftarrow p_1 \leftarrow q$ form a chain.



- p is (indirectly) density-reachable from q
- q is not density- reachable from p ?

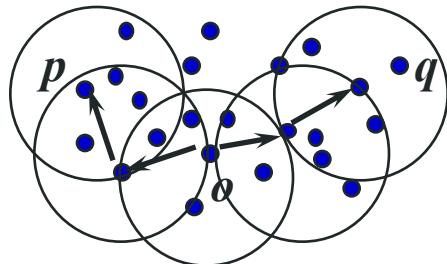
Density-Connectivity

■ Density-reachable is not symmetric

- not good enough to describe clusters

■ Density-Connected

- A pair of points p and q are density-connected if they are commonly density-reachable from a point o .



■ Density-connectivity is symmetric

Formal Description of Cluster

Given a data set D, parameter ε and threshold MinPts.

A cluster C is a subset of objects satisfying two criteria:

- **Connected:** for all p,q in C: p and q are density-connected.
- **Maximal:** for p,q: if p in C and q is density-reachable from p, then q in C. (avoid redundancy)



P is a core object.

DBSCAN Algorithm

Input: The data set D

Parameter: ε , MinPts

For each object p in D

 if p is a core object and not processed then

 C = retrieve all objects density-reachable from p

 mark all objects in C as processed

 report C as a cluster

 else mark p as outlier

 end if

End For

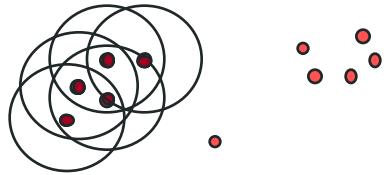
DBScan Algorithm

DBSCAN Algorithm: Example

Parameter

$\varepsilon = 2 \text{ cm}$

$MinPts = 3$



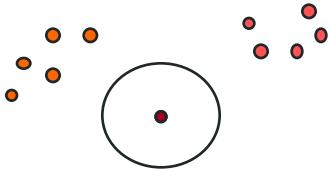
```
for each  $o \in D$  do
    if  $o$  is not yet classified then
        if  $o$  is a core-object then
            collect all objects density-reachable from  $o$ 
            and assign them to a new cluster.
        else
            assign  $o$  to NOISE
```

DBSCAN Algorithm: Example

Parameter

$\varepsilon = 2 \text{ cm}$

$MinPts = 3$



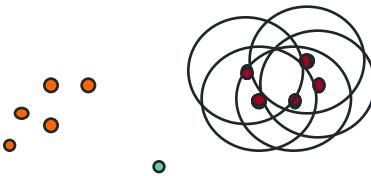
```
for each  $o \in D$  do
    if  $o$  is not yet classified then
        if  $o$  is a core-object then
            collect all objects density-reachable from  $o$ 
            and assign them to a new cluster.
        else
            assign  $o$  to NOISE
```

DBSCAN Algorithm: Example

Parameter

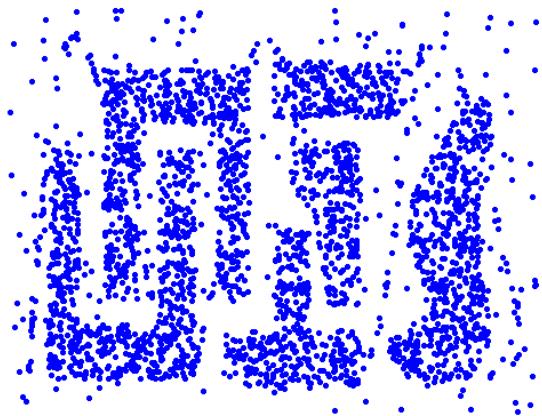
$\varepsilon = 2 \text{ cm}$

$MinPts = 3$

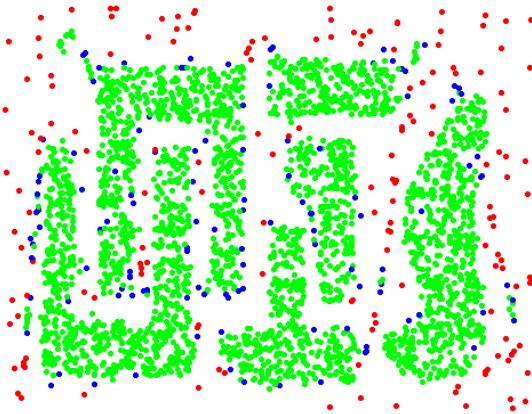


```
for each  $o \in D$  do
    if  $o$  is not yet classified then
        if  $o$  is a core-object then
            collect all objects density-reachable from  $o$ 
            and assign them to a new cluster.
        else
            assign  $o$  to NOISE
```

Example



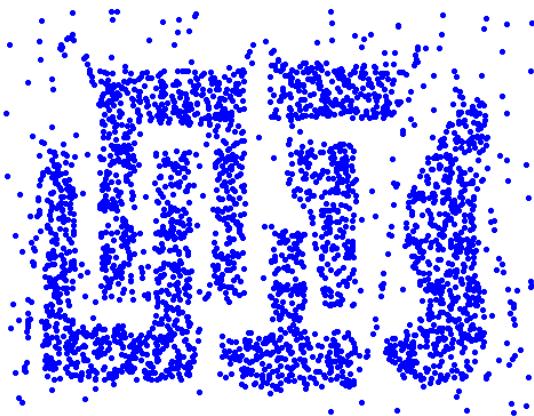
Original Points



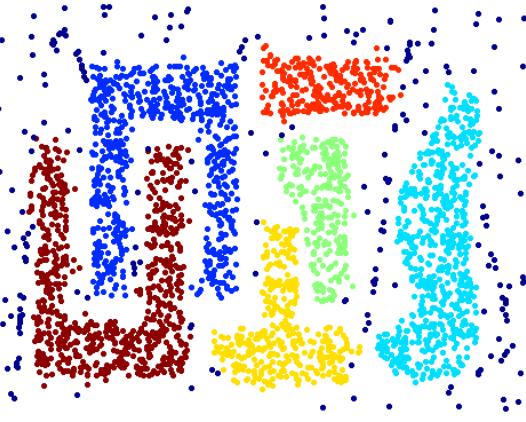
Point types: core,
border and outliers

$\varepsilon = 10$, MinPts = 4

When DBSCAN Works Well



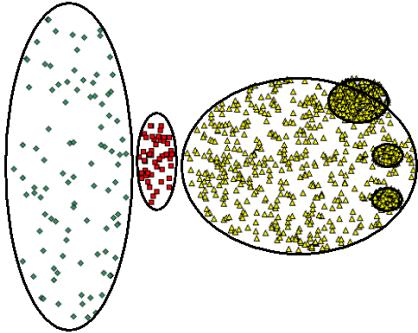
Original Points



Clusters

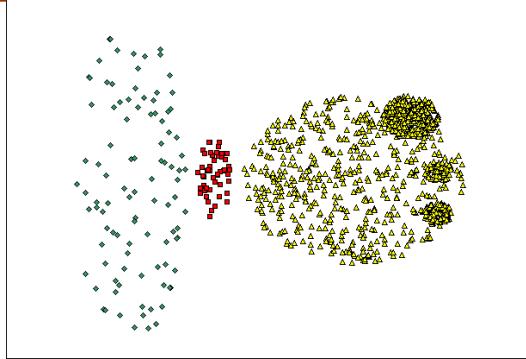
- Resistant to Noise
- Can handle clusters of different shapes and sizes

When DBSCAN Does NOT Work Well

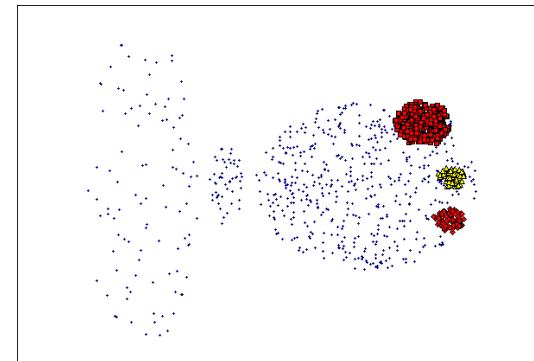


Original Points

- Cannot handle Varying densities
- sensitive to parameters



(MinPts=4, Eps=9.92).



(MinPts=4, Eps=9.75)

K-MEANS AND DBSCAN in Scikit-Learn

首先先引入本教學會用到的套件

```
In [1]: # Start from importing necessary packages.
import warnings
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.cm as cm

from IPython.display import display
from sklearn import metrics # for evaluations
from sklearn.datasets import make_blobs, make_circles # for generating experimental data
from sklearn.preprocessing import StandardScaler # for feature scaling
from sklearn.cluster import KMeans
from sklearn.cluster import DBSCAN

# make matplotlib plot inline (Only in Ipython).
warnings.filterwarnings('ignore')
%matplotlib inline
```

產生含有三個 clusters 的二維資料

In [2]:

```
# Generate data.  
# `random_state` is the seed used by random number generator for reproducibility (default=None).  
X, y = make_blobs(n_samples=5000,  
                  n_features=2,  
                  centers=3,  
                  random_state=170)  
  
# Print data like Ipython's cell output (Only in Ipython, otherwise use `print`).  
display(X)  
display(y)  
  
array([[-4.01009423, -1.01473496],  
       [ 1.00550526,  0.13163222],  
       [ 2.06563121, -0.24527689],  
       ...,  
       [-5.09493013,  1.47160372],  
       [-9.61459714, -4.91848716],  
       [-7.72675795, -5.86656563]])  
  
array([1, 2, 2, ..., 1, 0, 0])
```

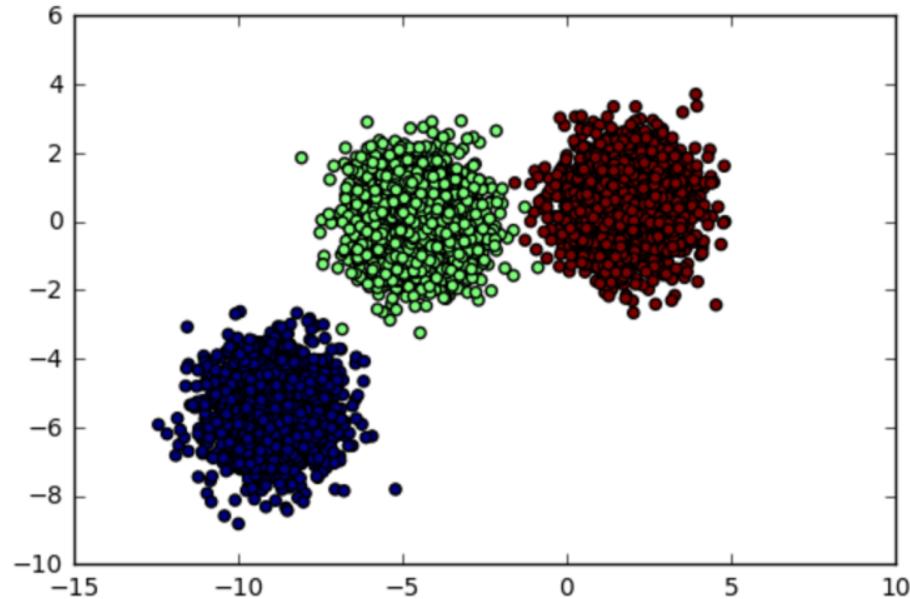
← 這是我們的訓練資料 X，共 5000 筆，每筆資料是二維，因此 X 是 (5000 x 2) 的矩陣。

← 這是我們每筆資料對應的 cluster label y (ground truth)，在此例有三個 clusters，因此 label = {0, 1, 2}，共 5000 筆，因此 y 是 5000 維的向量。

讓我們把資料視覺化一下，看看長什麼樣子。

```
In [3]: # Plot the data distribution (ground truth) using matplotlib `scatter(axis-x, axis-y, color)`.  
plt.scatter(X[:,0], X[:,1], c=y)
```

```
Out[3]: <matplotlib.collections.PathCollection at 0x111d4c310>
```



```
In [4]: # Perform K-means on our data (Train centroids)
kmeans = KMeans(n_clusters=3,
                 n_init=3,
                 init='random',
                 tol=1e-4,
                 random_state=170,
                 verbose=True).fit(X)
```

```
Initialization complete
Iteration 0, inertia 93044.046
Iteration 1, inertia 11941.892
Iteration 2, inertia 9733.271
Converged at iteration 2
Initialization complete
Iteration 0, inertia 49491.222
Iteration 1, inertia 43229.847
Iteration 2, inertia 42926.507
Iteration 3, inertia 42574.503
Iteration 4, inertia 41999.156
Iteration 5, inertia 39853.967
Iteration 6, inertia 28163.261
Iteration 7, inertia 11338.225
Iteration 8, inertia 9733.783
Converged at iteration 8
Initialization complete
Iteration 0, inertia 35757.303
Iteration 1, inertia 10003.824
Iteration 2, inertia 9733.154
Converged at iteration 2
```

使用 scikit-learn 所提供的 KMeans 來跑我們的資料。

參數解釋

- n_clusters : 你想在這資料上分幾群，也就是「選擇 “k”」啦。
- n_init : 你想讓 scikit-learn 的 k-means 演算法「跑幾次」，這邊我們是設定成 3，意思就是它會跑 3 次的 k-means，並且選擇最低 inertia 的那一次 (這邊的 inertia 目標就是 minimize sum of squared distance, distance 越小當然就越好囉！)，當作 k-means 的訓練完後的結果，也就是留著最好的 centroids。
- init : 還記得 k-means 演算法一開始會隨機選擇 centroids 嗎？這邊就是設定「怎麼選擇 centroids 的演算法」，有 “random” 跟 “k-means++” (後面會解釋) 可以選擇。
- tol : k-means 收斂時距離門檻值的比例，比例越高會越早視為收斂結束。
- random_state : 紿一固定值可以讓我們一直跑出同樣的結果，通常 debug 才會這麼用，一般來說不需要。
- verbose : 設定 True 則會 print 出訓練 k-means 時的過程，反之則 False。

```
In [5]: # Retrieve predictions and cluster centers (centroids).  
display(kmeans.labels_)  
display(kmeans.cluster_centers_)
```

```
array([0, 2, 2, ..., 0, 1, 1], dtype=int32)
```

```
array([[-4.55676387,  0.04603707],  
      [-8.94710203, -5.51613184],  
      [ 1.89450492,  0.5009336 ]])
```

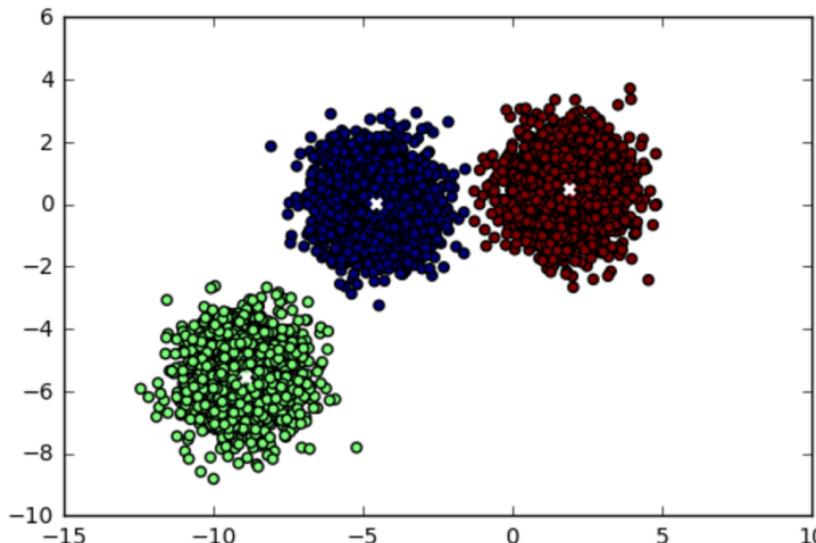
← 使用 “kmeans.labels_” 來取得我們預測的結果，同樣是 5000 維的向量。

← 使用 “kmeans.cluster_centers_” 來取得訓練出來的 centroids，有三個 clusters，資料為二維，因此是 (3x2) 的矩陣。

```
In [6]: # Plot the predictions.
```

```
plt.scatter(X[:,0], X[:,1], c=kmeans.labels_)  
plt.scatter(kmeans.cluster_centers_[:,0],  
            kmeans.cluster_centers_[:,1],  
            c='w', marker='x', linewidths=2)
```

```
Out[6]: <matplotlib.collections.PathCollection at 0x1125648d0>
```



← 把預測的結果視覺化出來，並且也順便把 centroids 紛畫上去。

注意：每次預測出來的 cluster 雖然跟 ground truth 畫的 cluster 不同，但這只是因為 k-means 隨機選擇 centroids 的關係，所以有可能每次跑出來的結果中，同樣的一群資料被分群到的 cluster label 會不同，但是重點是看整體的分群結果喔！顏色不同其實也沒關係。

```
In [7]: # We can make new predictions without re-run kmeans (simply find nearest centroids).
X_new = np.array([[10,10], [-10, -10], [-5, 10]])
y_pred = kmeans.predict(X_new)

""" The below code is equivalent to:
y_pred = KMeans(...).fit_predict(X), but this needs to fit kmeans again.
"""

display(y_pred)
```

```
array([2, 1, 0], dtype=int32) ← 訓練完 k-means 之後，使用 “kmeans.predict()” 來預測新資料 X_new。
```

```
In [8]: # We can get distances from data point to every centroid

""" The below code is equivalent to:
from sklearn.metrics.pairwise import euclidean_distances
euclidean_distances(X_new, kmeans.cluster_centers_)

"""

kmeans.transform(X_new)
```

```
Out[8]: array([[ 17.63464636,  24.48965134,  12.48724601],
               [ 11.42592142,   4.60582976,  15.86659553],
               [  9.96382639,  16.01030799,  11.73739582]]) ← 使用 “kmeans.transform()” 取得 X_new 中每個資料點到每個 centroid 的距離。如果你有 M 個資料點，N 個 clusters，則你會有 (MxN) 的矩陣，每個 row 對應至一個資料點到 N 個 clusters 的距離。
```

K-means++

K-means++ 演算法，更聰明初始化 centroid 的方式

A Smarter Way to Initialize Centroids: K-means++

Since *K-means* highly depends on the initialization of the centroids, the clustering results may be converged to a local minimum. We can address this by setting `init='kmeans++'` instead of `'random'`. *K-means++* initializes centroids in a smarter way to speed up convergence. The algorithm is as follows:

1. Randomly choose one centroid from the data points.
2. For each data point x_i , compute the distance $D(x_i, c_j)$ where c_j is nearest to x_i .
3. Randomly choose one new data point as a new centroid using *weighted probability distribution* proportional to $D(x_i, c_j)^2$.
4. Repeat steps 2 and 3 until k centroids have been chosen.
5. Now we have initialized centroids, run *K-means* algorithm.

簡單來說，kmeans++ 跟 kmeans 的差別是：

主要希望選到的 centroids 它們各自的距離越大越好，避免選到離彼此之間很近的 centroids。

In [9]:

```
# Perform K-means++ on our data.  
kmeans_plus_plus = KMeans(n_clusters=3,  
                           n_init=3,  
                           init='k-means++',  
                           tol=1e-4,  
                           random_state=170,  
                           verbose=True).fit(X)
```

```
Initialization complete  
Iteration 0, inertia 14504.164  
Iteration 1, inertia 9735.745  
Iteration 2, inertia 9733.168  
Converged at iteration 2  
Initialization complete  
Iteration 0, inertia 10751.002  
Iteration 1, inertia 9733.462  
Converged at iteration 1  
Initialization complete  
Iteration 0, inertia 12316.398  
Iteration 1, inertia 9733.520  
Converged at iteration 1
```

想使用 **k-means++**, 你只要設定 `init='k-means++'` 就可以囉！很簡單吧？而且有沒有發現 **k-means** 每次在訓練的時候，收斂速度變快了呢？這要歸功於 **k-means++** 一開始選了很好的 `centroids`, 導致後面在 `clustering` 的時候就會很快！

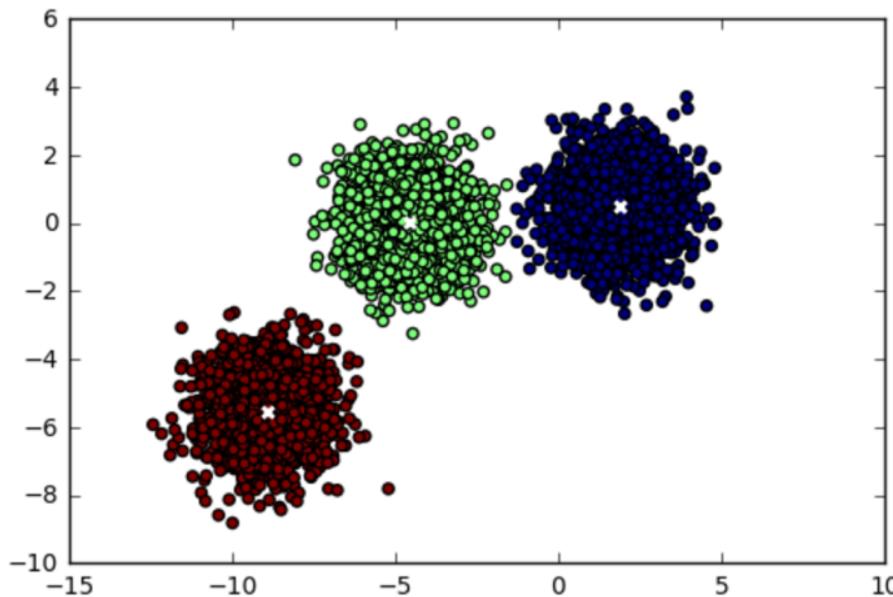
在 **scikit-learn** 中，你其實也不需要手動設定 `init='k-means++'`，因為它默認值就是如此囉！

You can see that *K-means++* converges much faster than *K-means*!

K-means++ 視覺化的預測結果

```
In [10]: # Plot the predictions.  
plt.scatter(X[:,0], X[:,1], c=kmeans_plus_plus.labels_)  
plt.scatter(kmeans_plus_plus.cluster_centers_[:,0],  
            kmeans_plus_plus.cluster_centers_[:,1],  
            c='w', marker='x', linewidths=2)
```

```
Out[10]: <matplotlib.collections.PathCollection at 0x1127a4250>
```

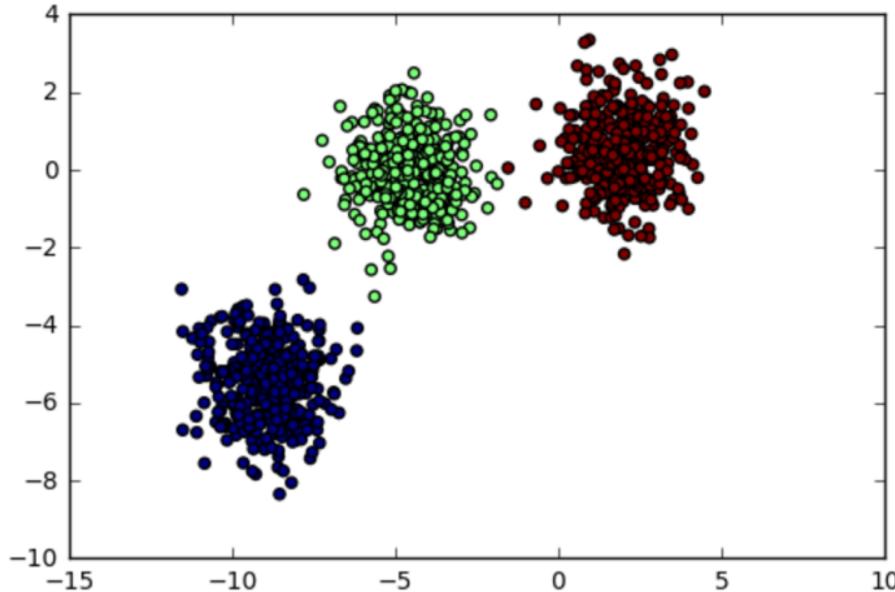


弱點一：需要選擇正確的“k”

Drawback 1: Need to choose a right number of clusters

```
In [11]: # Generate data.  
X, y = make_blobs(n_samples=1000,  
                  n_features=2,  
                  centers=3,  
                  random_state=170)  
  
# Plot the data distribution.  
plt.scatter(X[:,0], X[:,1], c=y)
```

```
Out[11]: <matplotlib.collections.PathCollection at 0x1129adc90>
```

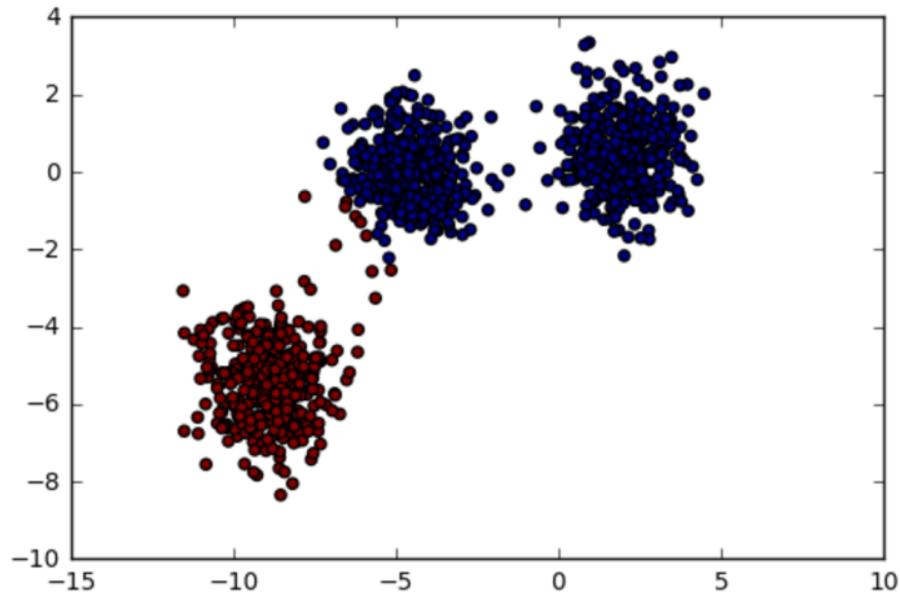


首先，我們一樣先產生三群資料，並且視覺化 ground truth 來看看。

如果我們選擇錯誤的“k”，會變得怎麼樣呢？

```
In [12]: # Run k-means on non-spherical data.  
y_pred = KMeans(n_clusters=2, random_state=170).fit_predict(X)  
  
# Plot the predictions.  
plt.scatter(X[:,0], X[:,1], c=y_pred)
```

```
Out[12]: <matplotlib.collections.PathCollection at 0x112b5bc90>
```



答案是：k 選錯，還是可以 work !

如果我們選擇 $k=2$ ，結果還是可以分群，你可以看到中間的 cluster 被歸類成藍色，是因為與右邊的 cluster 比較近的關係。

雖然能 work，但與 ground truth 明顯不符，因此未來如果要預測新資料，也許會表現得不好。

那 ... 我們要怎麼選擇一個好的 k 呑？

Solution: Measuring Cluster Quality to Determine the Number of Clusters

Supervised method

Homogeneity: Each cluster contains only members of a single class.

Completeness: All members of a given class are assigned to the same cluster.

如果你的訓練資料中已經含有 ground truth, 也就是知道哪一筆資料點是屬於哪一個 cluster, 那你可以使用 “supervised” 的方式來評估 clusters 的好壞！主要有兩種指標：

- Homogeneity : 每個 cluster 中的資料點最好都要是同一個 label。
- Completeness : 每個資料點所被預測出來的 cluster label 最好跟 ground truth 一樣（猜對）。

What if we don't have ground truth?

Unsupervised method

Silhouette Coefficient: Evaluate how well the **compactness** and the **separation** of the clusters are. (Note that the notation below is consistent with the above content.) Using *Silhouette Coefficient*, we can choose an optimal value for number of clusters.

$a(x_i)$ denotes the **mean intra-cluster distance**. Evaluate the compactness of the cluster to which x_i belongs. (The smaller the more compact)

$$a(x_i) = \frac{\sum_{x_k \in C_j, k \neq i} D(x_i, x_k)}{|C_j| - 1}$$

For the data point x_i , calculate its average distance to all the other data points in its cluster. (Minusing one in denominator part is to leave out the current data point x_i)

$b(x_i)$ denotes the **mean nearest-cluster distance**. Evaluate how x_i is separated from other clusters. (The larger the more separated)

$$b(x_i) = \min_{C_j: 1 \leq j \leq k, x_i \notin C_j} \left\{ \frac{\sum_{x_k \in C_j} D(x_i, x_k)}{|C_j|} \right\}$$

For the data point x_i and all the other clusters not containing x_i , calculate its average distance to all the other data points in the given clusters. Find the minimum distance value with respect to the given clusters.

Finally, *Silhouette Coefficient*: $s(x_i) = \frac{b(x_i) - a(x_i)}{\max\{a(x_i), b(x_i)\}}$, $-1 \leq s(x_i) \leq 1$. Want $a(x_i) < b(x_i)$ and $a(x_i) \rightarrow 0$ so as to $s(x_i) \rightarrow 1$.

如果很不幸地，你的訓練資料中沒有 ground truth，沒關係！我們還有“unsupervised”的方式可以評估 clusters 的好壞！

Unsupervised method

Silhouette coefficient: a value for interpretation and validation of consistency within clusters of data, which provides a succinct graphical representation of how well each object lies within its cluster.

The silhouette value is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation). The silhouette ranges from -1 to $+1$, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighboring clusters. If most objects have a high value, then the clustering configuration is appropriate. If many points have a low or negative value, then the clustering configuration may have too many or too few clusters.

Silhouette coefficient

For each datum \mathbf{i} , let $a(\mathbf{i})$ be the average distance between \mathbf{i} and all other data within the same cluster.

Let $b(\mathbf{i})$ be the smallest average distance of \mathbf{i} to all points in any other cluster, of which \mathbf{i} is not a member.

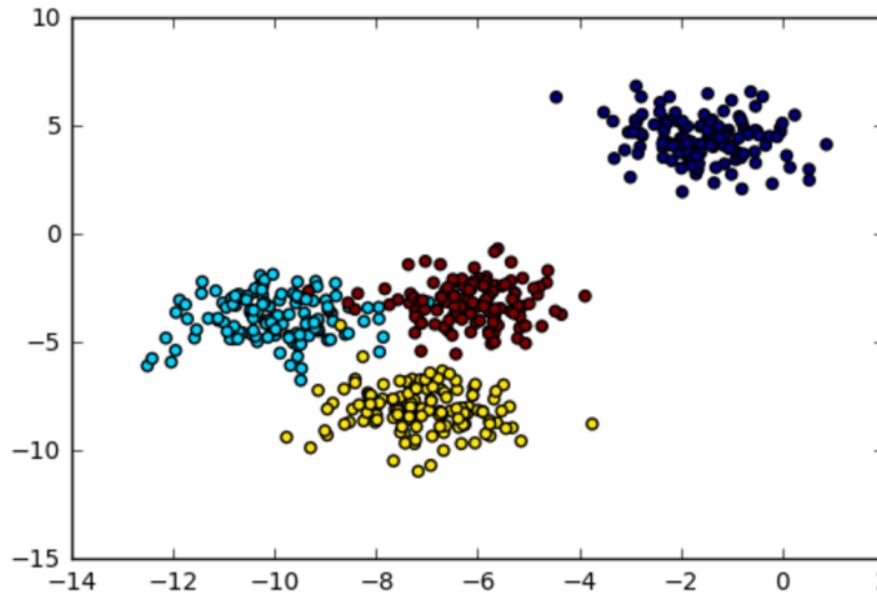
We now define a silhouette:

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

```
In [13]: # Generate data.
# This particular setting has one distinct cluster and 3 clusters placed close together.
X, y = make_blobs(n_samples=500,
                  n_features=2,
                  centers=4,
                  cluster_std=1,
                  center_box=(-10.0, 10.0),
                  shuffle=True,
                  random_state=1)

# Plot the data distribution.
plt.scatter(X[:,0], X[:,1], c=y)
```

Out[13]: <matplotlib.collections.PathCollection at 0x112cccd950>



首先，我們先產生四群 **clusters** 的資料，並且視覺化出來，注意到一點是我們故意產生出三群 **clusters** 是稍微有重疊到的，方便我們用上述所提到的 **cluster quality measurements** 來跑看看如果選擇不同的“ k ”，系統跑出來的 **cluster quality** 指數是好還是壞。

```
In [14]: # List of number of clusters
range_n_clusters = [2, 3, 4, 5, 6]

# For each number of clusters, perform Silhouette analysis and visualize the results.
for n_clusters in range_n_clusters:

    # Perform k-means.
    kmeans = KMeans(n_clusters=n_clusters, random_state=10)
    y_pred = kmeans.fit_predict(X)

    # Compute the cluster homogeneity and completeness.
    homogeneity = metrics.homogeneity_score(y, y_pred)
    completeness = metrics.completeness_score(y, y_pred)

    # Compute the Silhouette Coefficient for each sample.
    s = metrics.silhouette_samples(X, y_pred)

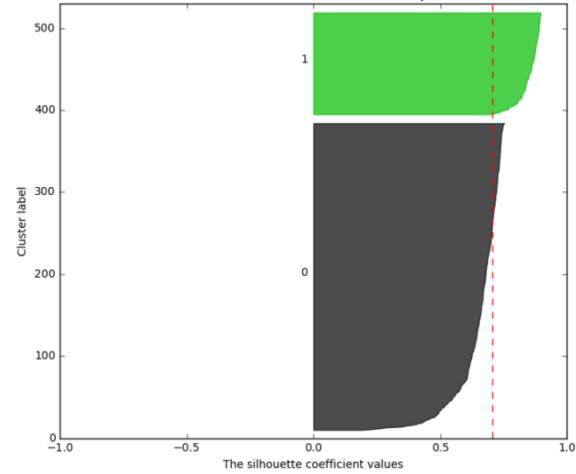
    # Compute the mean Silhouette Coefficient of all data points.
    s_mean = metrics.silhouette_score(X, y_pred)
```

上述程式碼解釋：

- 首先我們先設定要跑的 “k” 的範圍，此例我們嘗試跑 $k=2, 3, 4, 5, 6$ 各個不同的 case，並且使用 for 迴圈，分別計算不同 “k” 之分群結果的 homogeneity, completeness 跟 silhouette coefficient。
- 使用 “homogeneity_score(y, y_pred)” 與 “completeness_score(y, y_pred)” 即可計算對應的兩種指標
- 注意到 scikit-learn 在計算 silhouette coefficient 時有提供兩種選擇：
 - 使用 “silhouette_samples(X, y_pred)” 來計算每個資料點的 silhouette coefficient。
 - 使用 “silhouette_score(X, y_pred)” 來計算每個資料點的 silhouette coefficient 加總之後的平均值。

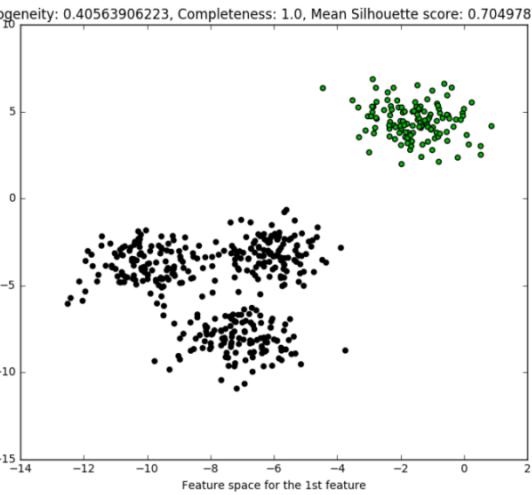
Silhouette analysis for K-Means clustering with n_clusters: 2

Silhouette Coefficient for each sample



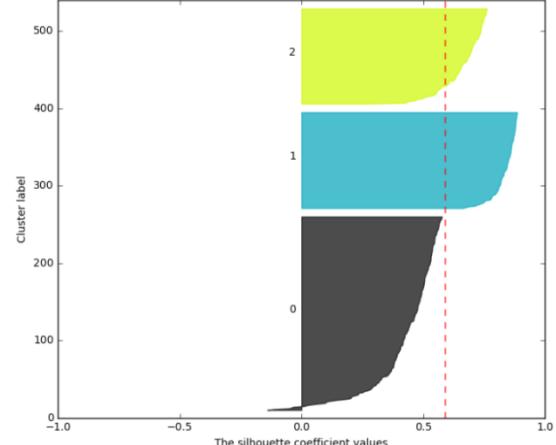
Homogeneity: 0.40563906223, Completeness: 1.0, Mean Silhouette score: 0.704978749608

Feature space for the 2nd feature



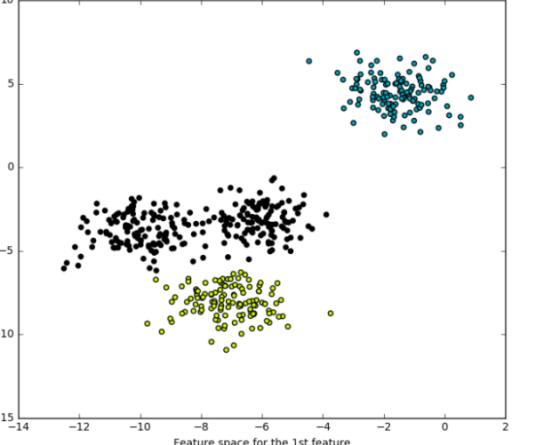
Silhouette analysis for K-Means clustering with n_clusters: 3

Silhouette Coefficient for each sample



Homogeneity: 0.725794887698, Completeness: 0.969029755722, Mean Silhouette score: 0.588200401213

Feature space for the 2nd feature



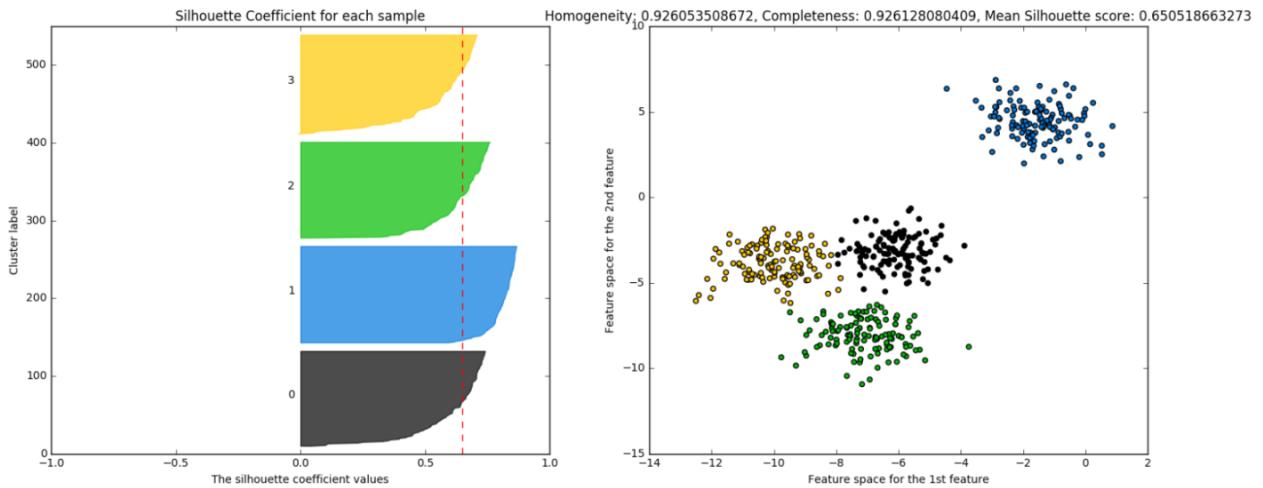
由於視覺化的程式碼稍微複雜，有興趣請查看 [Ipython](#)。

我們先來講講怎麼看這張圖吧！

右邊很簡單，就是分群後的結果，主要是左邊：

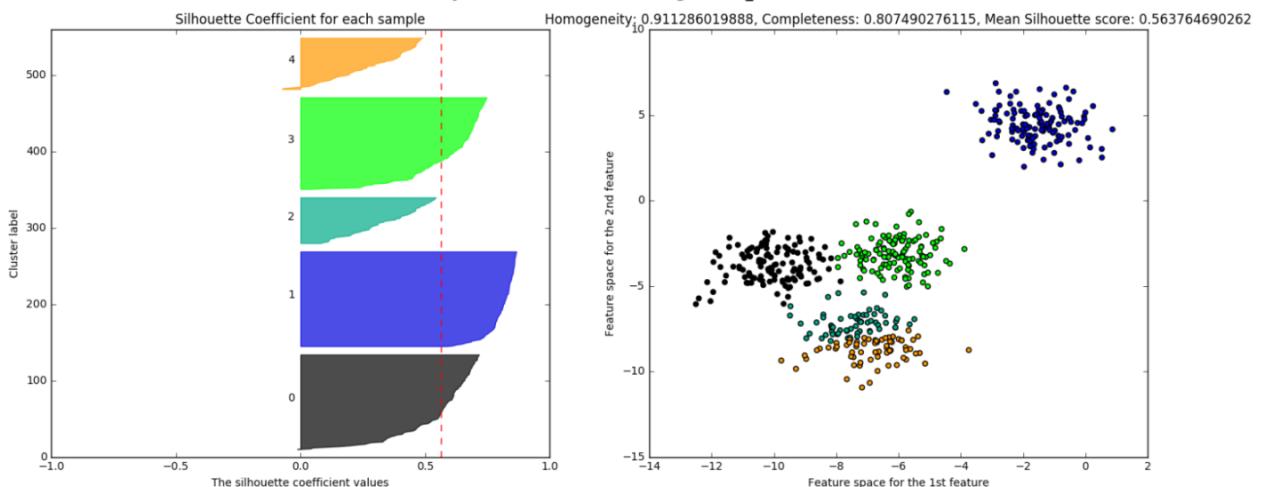
- 橫軸代表 silhouette coefficient。
- 縱軸代表 data points。
- 相同 cluster 的資料點畫成相同顏色，且該 cluster 內的資料點之 silhouette coefficient 皆排序過。
- 注意：紅色的垂直虛線是“silhouette_score()”算出來的全部 silhouette coefficient 加總之平均。

Silhouette analysis for K-Means clustering with n_clusters: 4



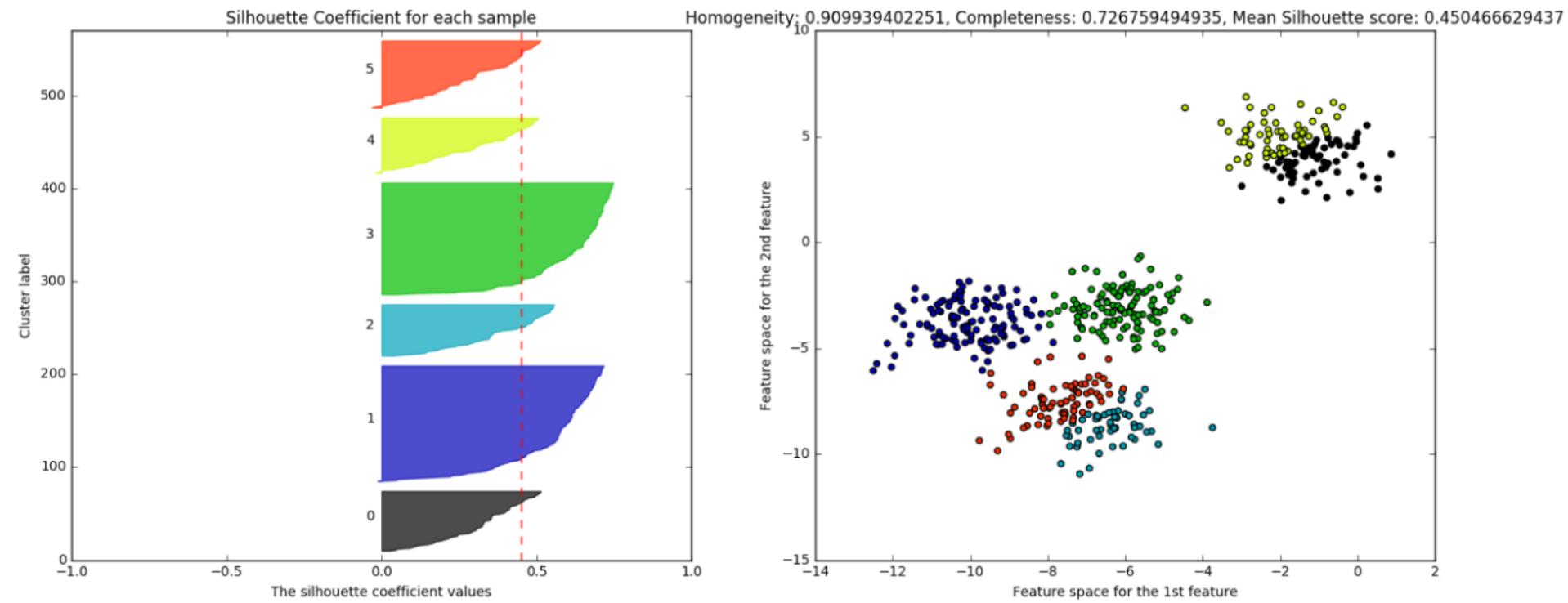
← k=4 時，每個 cluster 內的資料點之 silhouette coefficient 普遍高於平均值（超出紅色垂直線），因此以 k=4 來講，分群效果還算不錯。

Silhouette analysis for K-Means clustering with n_clusters: 5



← k=5 時，可以注意到有兩群（土黃色以及蒂芬尼藍）的平均 silhouette coefficient 很低，因為兩群在右圖中重疊很多（影響到 silhouette coefficient 中計算 mean nearest cluster distance $b(x)$ 的指標），因此 k=5 不是一個好選擇。

Silhouette analysis for K-Means clustering with n_clusters: 6



The silhouette plot shows that the `n_clusters` value of 3, 5 and 6 are a bad pick for the given data due to the presence of clusters with above average silhouette scores and also due to wide fluctuations in the size of the silhouette plots. Silhouette analysis is more ambivalent in deciding between 2 and 4.

總體來說， $k=3, 5, 6$ 分群效果不佳，而 $k=2, 4$ 還不錯，至於要選 2 還是 4，就看應用在什麼問題上了。

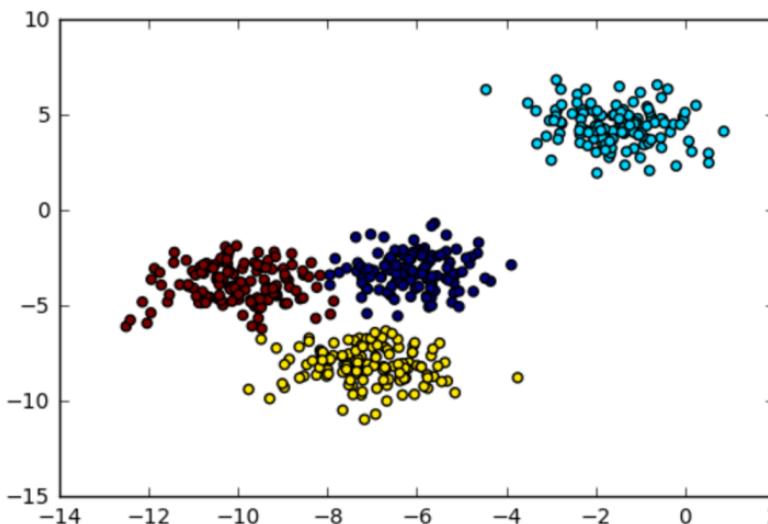
弱點二：不能處理 Noise data
與 Outliers

```
In [15]: # Generate data.
# This particular setting has one distinct cluster and 3 clusters placed close together.
# (Same as the above example)
X, y = make_blobs(n_samples=500,
                   n_features=2,
                   centers=4,
                   cluster_std=1,
                   center_box=(-10.0, 10.0),
                   shuffle=True,
                   random_state=1)

# Perform k-means with n_clusters=4
kmeans = KMeans(n_clusters=4, random_state=10)
y_pred = kmeans.fit_predict(X)

# Plot the prediction
plt.scatter(X[:,0], X[:,1], c=y_pred)
```

Out[15]: <matplotlib.collections.PathCollection at 0x1066c8e50>



一樣使用上面同樣的資料，並且使用 $k=4$ 去分群，可以發現其實還是有些點離它們各自的 cluster 很遠，但是 k-means 依然將他們分群至 cluster 內。

那 ... 我們要怎麼把這些離得比較遠的點，也就是所謂的 outliers 紿偵測出來呢？

Solution: Use Distance Threshold to Detect Noise data and Outliers

However, we can detect the noises/outliers conditioning on whether the distance between the data point x_i and the centroid c_j of x_i 's corresponding cluster is larger than the average distance in the cluster. That is to say:

$$x_i = \begin{cases} \text{Outlier}, & \text{if } D(x_i, c_j) > \frac{1}{|Cluster_j|} \sum_{k=0, k \neq i}^{|Cluster_j|} D(x_k, c_j) \text{ where } c_j \in Cluster_j \\ \text{Non-outlier}, & \text{otherwise} \end{cases}$$

Let's begin to find out the outliers of each cluster.

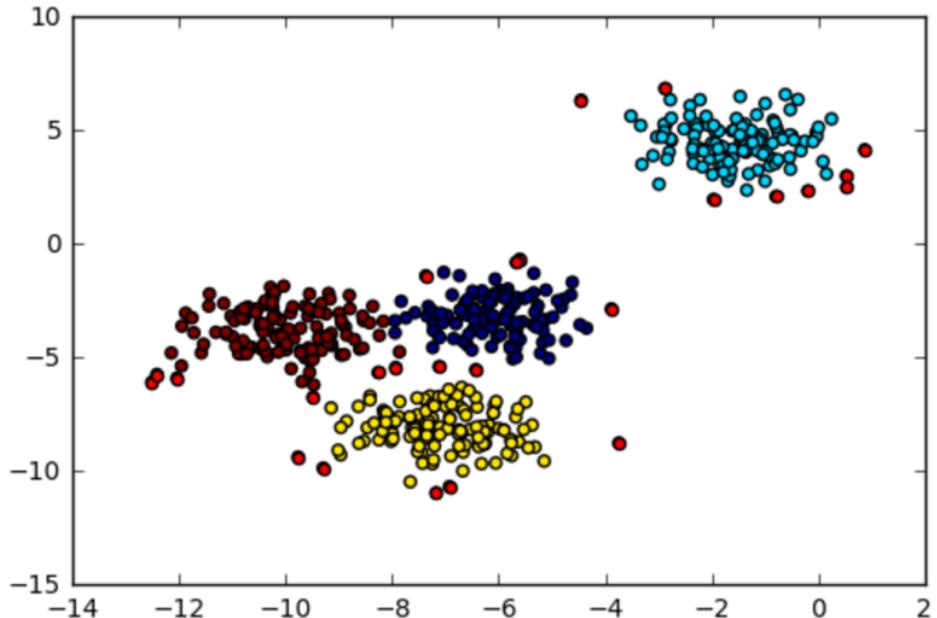
你可以藉由計算每個 cluster 中的 average distance 當作判斷一個資料點是否是 outlier 的 distance threshold (有時候 distance threshold 還會再乘上一個 ratio 來控制 threshold 大小)。

```
In [16]: # Ratio for our distance threshold, controlling how many outliers we want to detect.  
distance_threshold_ratio = 2.0  
  
# Plot the prediction same as the above.  
plt.scatter(X[:,0], X[:,1], c=y_pred)  
  
# For each ith cluster, i=0~3 (we have 4 clusters in this example).  
for i in [0, 1, 2, 3]:  
  
    # Retrieve the indexes of data points belong to the ith cluster.  
    # Note: `np.where()` wraps indexes in a tuple, thus we retrieve indexes using `tuple[0]`  
    indexes_of_X_in_ith_cluster = np.where(y_pred == i)[0]  
  
    # Retrieve the data points by the indexes  
    X_in_ith_cluster = X[indexes_of_X_in_ith_cluster]  
  
    # Retrieve the centroid.  
    centroid = kmeans.cluster_centers_[i]  
  
    # Compute distances between data points and the centroid.  
    # Same as: np.sqrt(np.sum(np.square(X_in_ith_cluster - centroid), axis=1))  
    # Note: distances.shape = (X_in_ith_cluster.shape[0], 1). A 2-D matrix.  
    distances = metrics.pairwise.euclidean_distances(X_in_ith_cluster, centroid)  
  
    # Compute the mean distance for ith cluster as our distance threshold.  
    distance_threshold = np.mean(distances)  
  
    # Retrieve the indexes of outliers in ith cluster  
    # Note: distances.flatten() flattens 2-D matrix to vector, in order to compare with scalar `distance_threshold`.  
    indexes_of_outlier = np.where(distances.flatten() > distance_threshold * distance_threshold_ratio)[0]  
  
    # Retrieve outliers in ith cluster by the indexes  
    outliers = X_in_ith_cluster[indexes_of_outlier]  
  
    # Plot the outliers in ith cluster.  
    plt.scatter(outliers[:,0], outliers[:,1], c='r')
```

← 1. 取得 cluster 內所有資料點，以及該 cluster 的 centroid。

← 2. 計算 cluster 內所有點到 centroid 的距離，並且平均起來，得出 distance threshold。

← 3. 找出資料點到 centroid 的距離如果大於 (distance threshold * ratio)，則視為 outliers。



← 藉由上述方法，我們成功將 **outliers** (紅色) 紿偵測出來了！不過其實值得提醒的一點是：當你的 **cluster** 之中佔有太多的 **outliers**，可想而知，**k-means** 分群出來得結果還是會不理想喔！因為在計算 **distance threshold** 時，還是會將這些 **outliers** 列入計算，因此大大影響了 **distance threshold**。

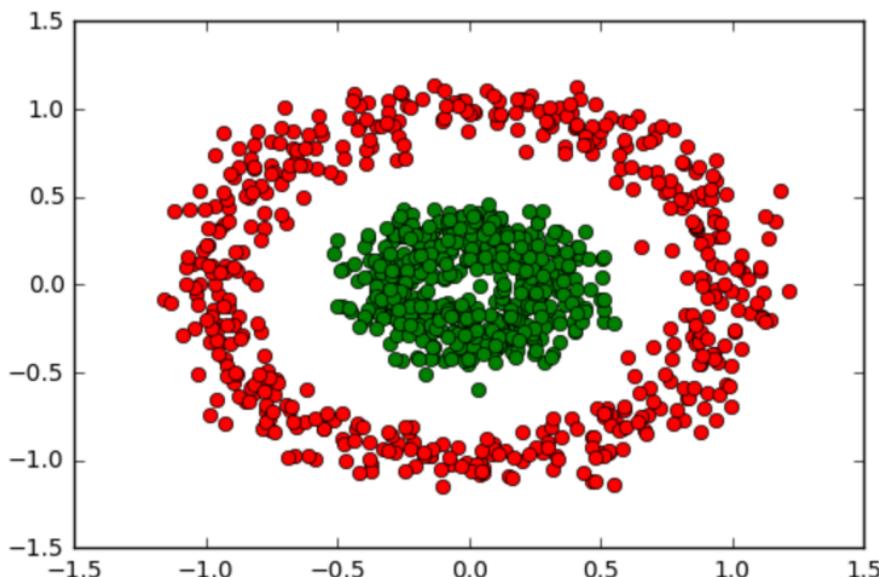
As we've mentioned about measuring cluster quality analysis, you can run different settings of `distance_threshold_ratio` to find out the best cluster quality.

弱點三：不能處理非球型分佈的資料集

In [17]:

```
# Generate non-spherical data.  
X, y = make_circles(n_samples=1000, factor=0.3, noise=0.1)  
  
# Plot the data distribution. (Here's another way to plot scatter graph)  
plt.plot(X[y == 0, 0], X[y == 0, 1], 'ro')  
plt.plot(X[y == 1, 0], X[y == 1, 1], 'go')
```

Out[17]: [`<matplotlib.lines.Line2D at 0x10f8e1fd0>`]



k-means 演算法根據其定義：此演算法目的是將資料點分成離它們最接近的 **mean** 值的 **cluster** 之中。這個假設建立於你的資料分布必須得是球型分佈，因此如果我們的資料長得像左圖這種同心圓資料分佈，就會輕易的使 k-means 分群的效果不佳。

讓我們來看看當 $k=2$ 時，k-means 會將資料分成什麼樣子。

After performing *K-means* on non-spherical data, the following result shows that it fails to cluster non-spherical data, since *K-means* has an assumption that the data distribution is spherical.

In [18]:

```
# Run k-means on non-spherical data.
y_pred = KMeans(n_clusters=2, random_state=170).fit_predict(X)

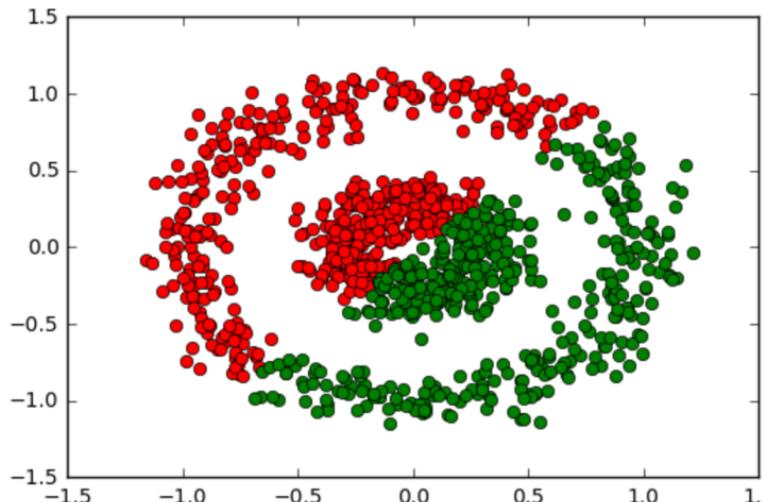
# Plot the predictions.
plt.plot(X[y_pred == 0, 0], X[y_pred == 0, 1], 'ro')
plt.plot(X[y_pred == 1, 0], X[y_pred == 1, 1], 'go')

# Print the evaluations
print('Homogeneity: {}'.format(metrics.homogeneity_score(y, y_pred)))
print('Completeness: {}'.format(metrics.completeness_score(y, y_pred)))
print('Mean Silhouette score: {}'.format(metrics.silhouette_score(X, y_pred)))
```

Homogeneity: 0.000184968858484

Completeness: 0.000185182645182

Mean Silhouette score: 0.295848480827



使用 **k-means** 分群後的效果並不佳，因為外圈與內圈的 **mean** 是差不多的，而 **k-means** 是基於 **euclidean distance** 來做分群，因此當 **k=2** 時，此同心圓有一半會跟其中一個 **centroid** 比較近，另一半會跟另一個 **centroid** 比較近，因此就被硬生生的砍成一半了。

Solution: Using Feature Transformation or Extraction Techiques Makes Data Clusterable

If you know that your clusters will always be concentric circles, you can simply convert your cartesian (x-y) coordinates to polar coordinates, and use only the radius for clustering - as you know that the angle theta doesn't matter.

Or more generally: use a suitable kernel for k-means clustering, e.g. use *Kernel PCA* to find a projection of the data that makes data linearly separable, or use another clustering algorithm, such as *DBSCAN*.

In [19]:

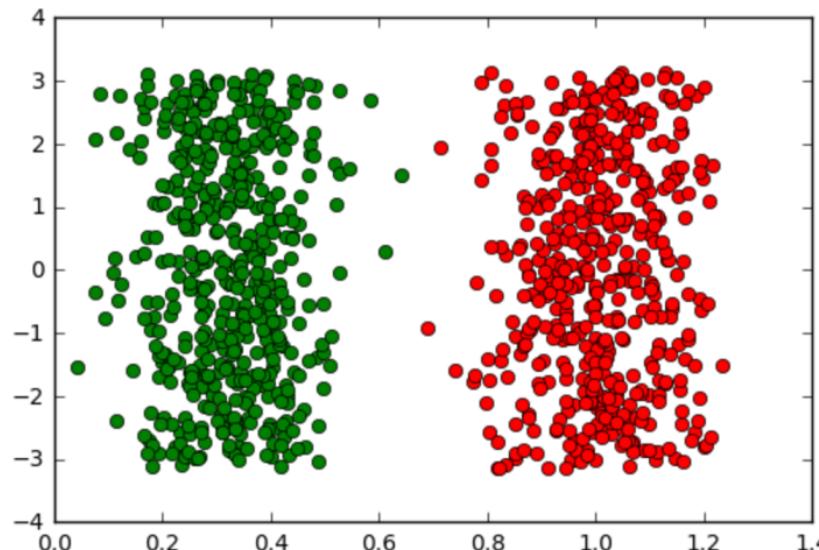
```
1 def cart2pol(x, y):
2     radius = np.sqrt(x**2 + y**2) ← 將 Carteisan 座標轉換成 polar 座標的 function。
3     theta = np.arctan2(y, x)
4     return radius, theta
5
6 X_transformed = np.zeros_like(X)
7 # Convert cartesian (x-y) to polar coordinates. ← 我們只看第一個 feature—radius, 第二個 feature
8 X_transformed[:,0], _ = cart2pol(X[:,0], X[:,1])就全部填 0, 最後直接丟到 k-means 做分群。
9
10 # Only use `radius` feature to cluster.
11 y_pred = KMeans(n_clusters=2).fit_predict(X_transformed)
12
13 plt.plot(X[y_pred == 0, 0], X[y_pred == 0, 1], 'ro')
14 plt.plot(X[y_pred == 1, 0], X[y_pred == 1, 1], 'go')
```

```
In [38]: def cart2pol(x, y):
    radius = np.sqrt(x**2 + y**2)
    theta = np.arctan2(y, x)
    return radius, theta

X_transformed = np.zeros_like(X)
X_transformed[:,0], X_transformed[:,1] = cart2pol(X[:,0], X[:,1])

plt.plot(X_transformed[y == 0, 0], X_transformed[y == 0, 1], 'ro')
plt.plot(X_transformed[y == 1, 0], X_transformed[y == 1, 1], 'go')
```

```
Out[38]: [<matplotlib.lines.Line2D at 0x1163fe890>]
```



We just successfully make data linearly separable by converting features (x-y) to (radius-theta) !

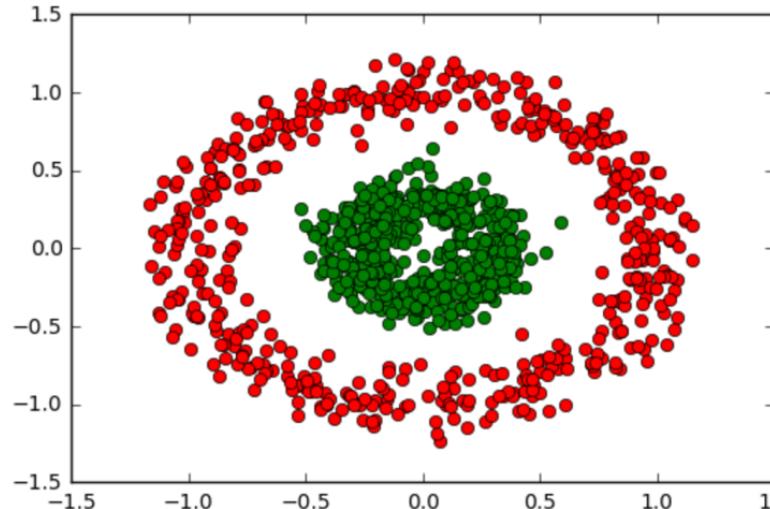
```
In [39]: def cart2pol(x, y):
    radius = np.sqrt(x**2 + y**2)
    theta = np.arctan2(y, x)
    return radius, theta

X_transformed = np.zeros_like(X)
# Convert cartesian (x-y) to polar coordinates.
X_transformed[:,0], _ = cart2pol(X[:,0], X[:,1])

# Only use `radius` feature to cluster.
y_pred = KMeans(n_clusters=2).fit_predict(X_transformed)

plt.plot(X[y_pred == 0, 0], X[y_pred == 0, 1], 'ro')
plt.plot(X[y_pred == 1, 0], X[y_pred == 1, 1], 'go')
```

```
Out[39]: [<matplotlib.lines.Line2D at 0x1181a2490>]
```



Now the data is successfully clustered!

DBSCAN:

Density-Based Spatial

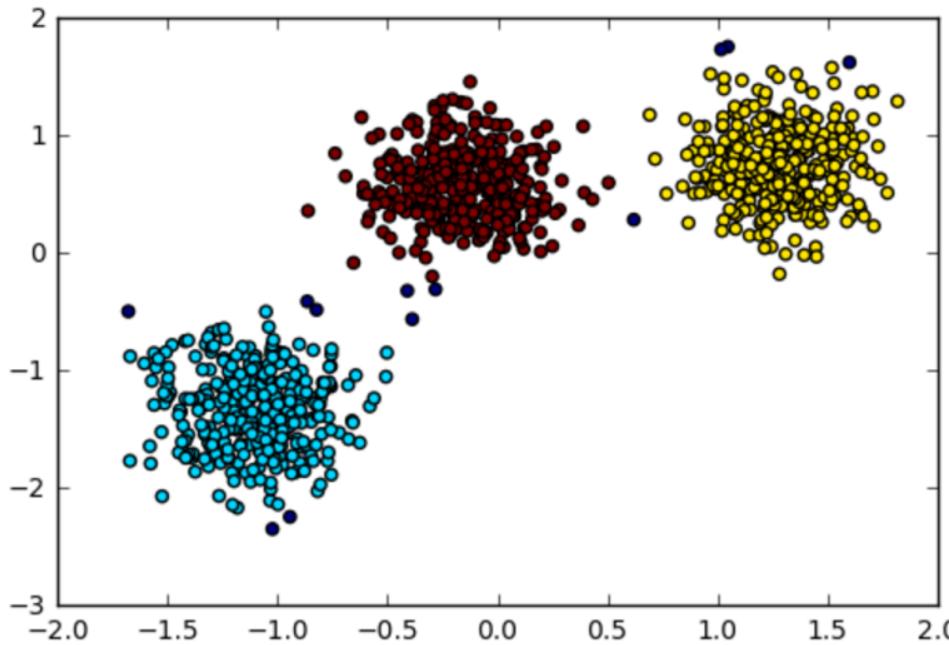
Clustering Algorithm with

Noise

Let's begin to perform DBSCAN on spherical data

```
In [20]: # Generate data with 3 centers.  
X, y = make_blobs(n_samples=1000,  
                  n_features=2,  
                  centers=3,  
                  random_state=170)           ← 1. 產生 (1000x2) 的訓練資料 X 及對應 ground truth y.  
  
# Standardize features to zero mean and unit variance.  
X = StandardScaler().fit_transform(X)          ← 2. 將資料 standardize 有助於避免某些維度因為  
                                                scaling 很大而影響 classification 或 clustering 的結果。  
  
# Perform DBSCAN on the data  
y_pred = DBSCAN(eps=0.3, min_samples=30).fit_predict(X)    ← 3. 使用 DBSCAN 來分群，其參數 "min_samples"  
                                                等同前面所介紹的 "MinPts".  
  
# Plot the predictions  
plt.scatter(X[:,0], X[:,1], c=y_pred)  
  
# Print the evaluations  
print('Number of clusters: {}'.format(len(set(y_pred[np.where(y_pred != -1)]))))  
print('Homogeneity: {}'.format(metrics.homogeneity_score(y, y_pred)))  
print('Completeness: {}'.format(metrics.completeness_score(y, y_pred)))  
print('Mean Silhouette score: {}'.format(metrics.silhouette_score(X, y_pred)))           ← 4. 印出 cluster quality 指數
```

Number of clusters: 3
Homogeneity: 0.982039796605
Completeness: 0.93781096194
Mean Silhouette score: 0.679854880541



The **black** data points denote the **outliers** in the above result.

Note that we don't need to specify the number of clusters with *DBSCAN* algorithm. Besides, *DBSCAN* is good at finding out the outliers without requiring some hacks like we did above in *K-means* section.

讓我們來試試 DBSCAN 跑在非球狀資料分布上

Now, let's try *DBSCAN* on non-spherical data.

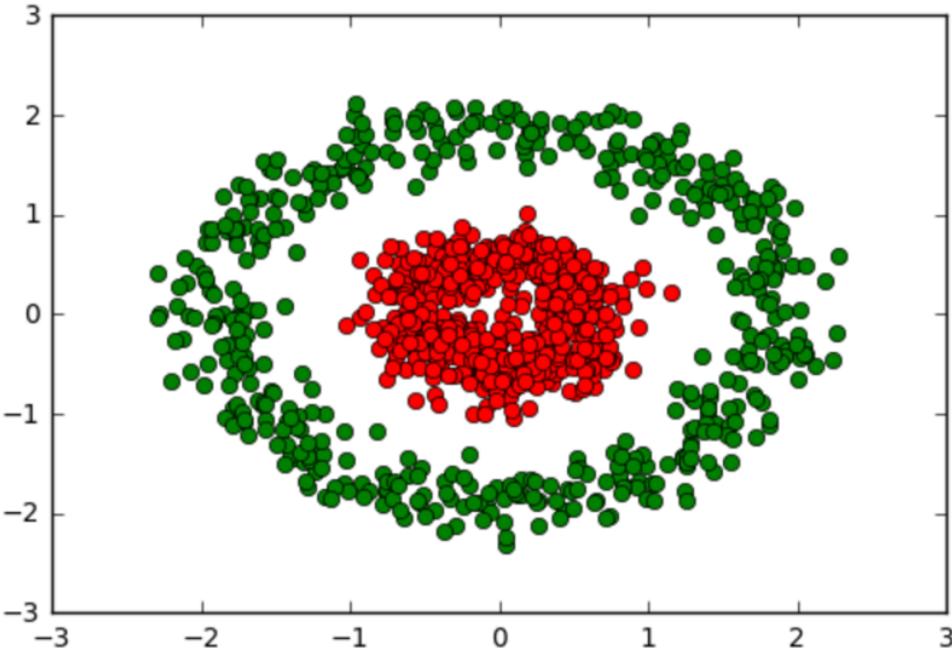
```
In [21]: # Generate non-spherical data.  
X, y = make_circles(n_samples=1000, factor=0.3, noise=0.1)  
  
# Standardize features to zero mean and unit variance.  
X = StandardScaler().fit_transform(X)  
  
# Perform DBSCAN on the data  
y_pred = DBSCAN(eps=0.3, min_samples=10).fit_predict(X)  
  
# Plot the data distribution. (Here's another way to plot scatter graph)  
plt.plot(X[y_pred == 0], X[y_pred == 0, 1], 'ro')  
plt.plot(X[y_pred == 1], X[y_pred == 1, 1], 'go')  
  
# Print the evaluations  
print('Number of clusters: {}'.format(len(set(y_pred[np.where(y_pred != -1)]))))  
print('Homogeneity: {}'.format(metrics.homogeneity_score(y, y_pred)))  
print('Completeness: {}'.format(metrics.completeness_score(y, y_pred)))  
print('Mean Silhouette score: {}'.format(metrics.silhouette_score(X, y_pred)))
```

Number of clusters: 2

Homogeneity: 0.996390359526

Completeness: 0.957691789224

Mean Silhouette score: 0.174843819318



Comparing to *K-means*, we can directly apply *DBSCAN* on this form of data distribution due to the density-based clustering criterion.

Note: It's worth mention that the *Silhouette score* is generally higher for **convex** clusters than other concepts of clusters, such as density based clusters.

**Thanks for Listening.
Any Questions ?**

Reference

- [scikit-learn](#)
- [How to understand the drawbacks of K-means](#)
- [The Sihouette Coefficient](#)
- [Cluster Analysis in Data Mining tutorial video series](#)
- [DBSCAN](#)