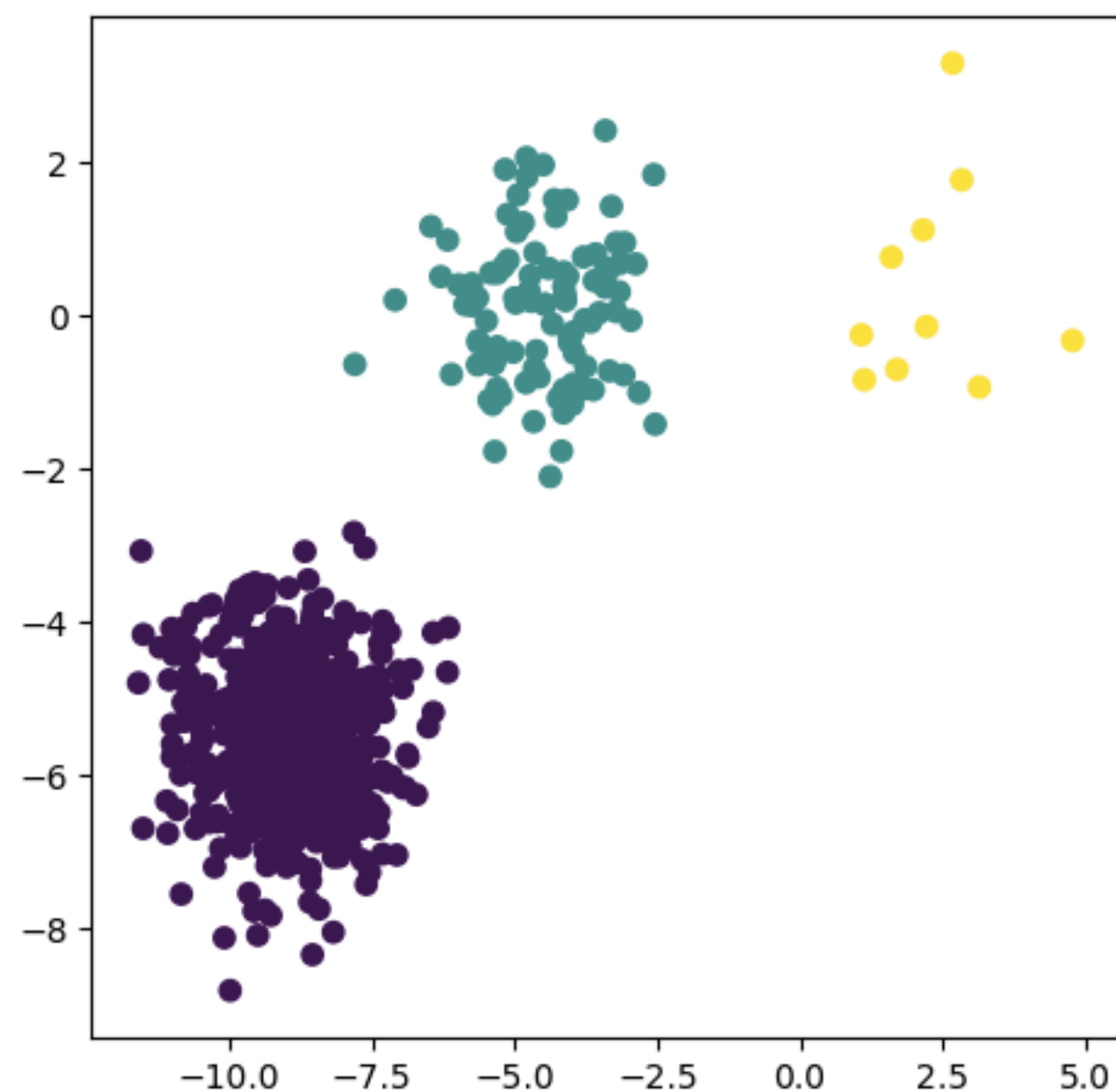


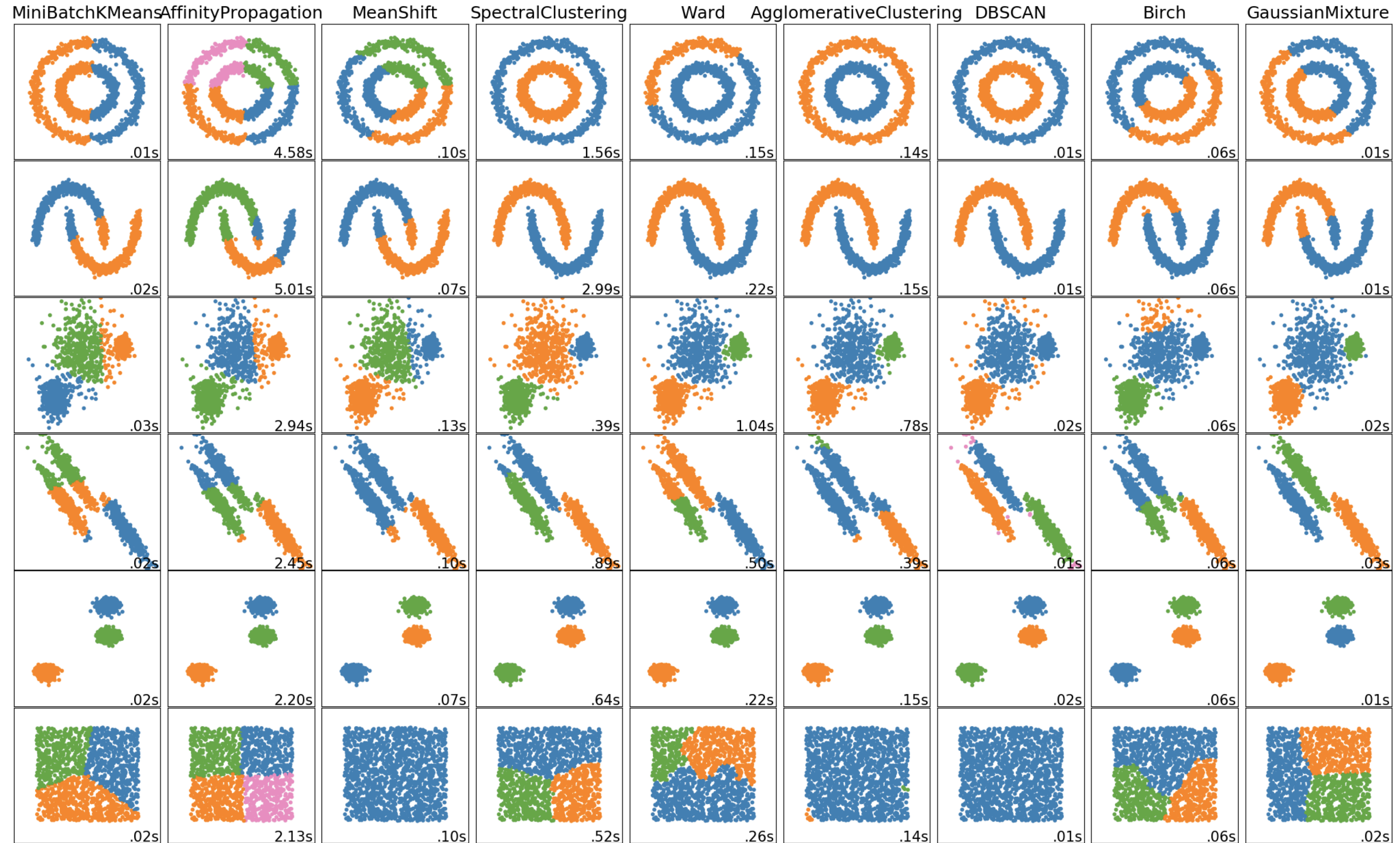
非監督式學習-分群

Unsupervised Learning-Clustering

分群 (Clustering)

- 非監督式：沒有 y ，只有 X (features)
- 透過資料彼此的距離分群

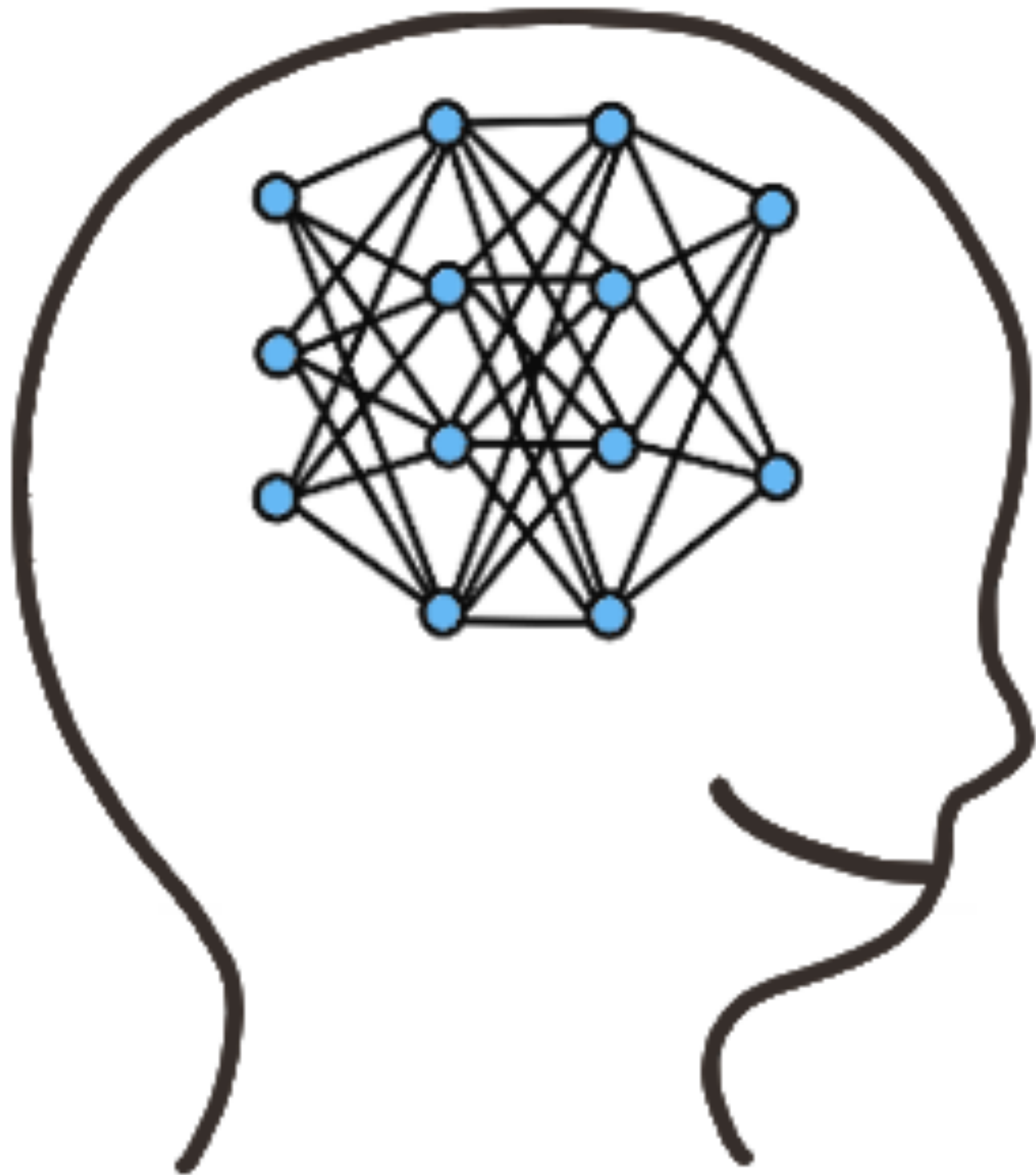






Clustering with sklearn

Method name	Parameters	Scalability	Usecase	Geometry (metric used)
K-Means	number of clusters	Very large <code>n_samples</code> , medium <code>n_clusters</code> with MiniBatch code	General-purpose, even cluster size, flat geometry, not too many clusters	Distances between points
Affinity propagation	damping, sample preference	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with <code>n_samples</code>	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium <code>n_samples</code> , small <code>n_clusters</code>	Few clusters, even cluster size, non- flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large <code>n_samples</code> and <code>n_clusters</code>	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large <code>n_samples</code> , medium <code>n_clusters</code>	Non-flat geometry, uneven cluster sizes	Distances between nearest points
Gaussian mixtures	many	Not scalable	Flat geometry, good for density estimation	Mahalanobis distances to centers
Birch	branching factor, threshold, optional global clusterer.	Large <code>n_clusters</code> and <code>n_samples</code>	Large dataset, outlier removal, data reduction.	Euclidean distance between points



K-means



Kmeans

Notes

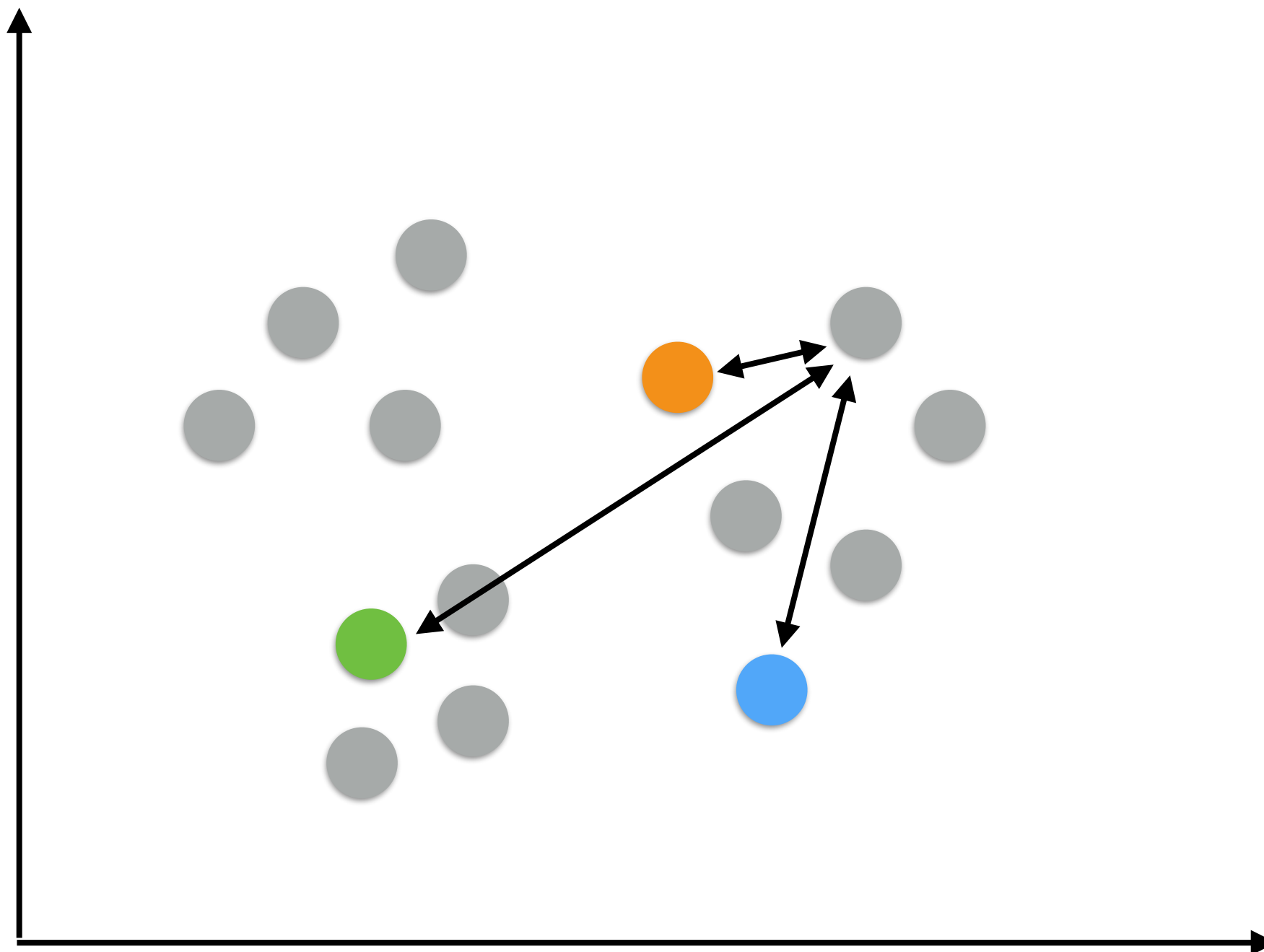
- 設 $K = 3$ ，起始隨機挑選3個點作為集群中心點(centroid)
 - ▶ 需做資料標準化，讓Kmeans計算距離時使用相同的尺規。
- repeat：將附近的點根據與這3個中心點的距離分配到這三群，並重新計算中心點，直到收斂為止
- 收斂：得到與各集群中心點距離和最小值
- Kmeans 使用歐幾里德距離(Euclidean Distance)

$$d(x, y) := \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 + \cdots + (x_n - y_n)^2} = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$



Kmeans Algorithm

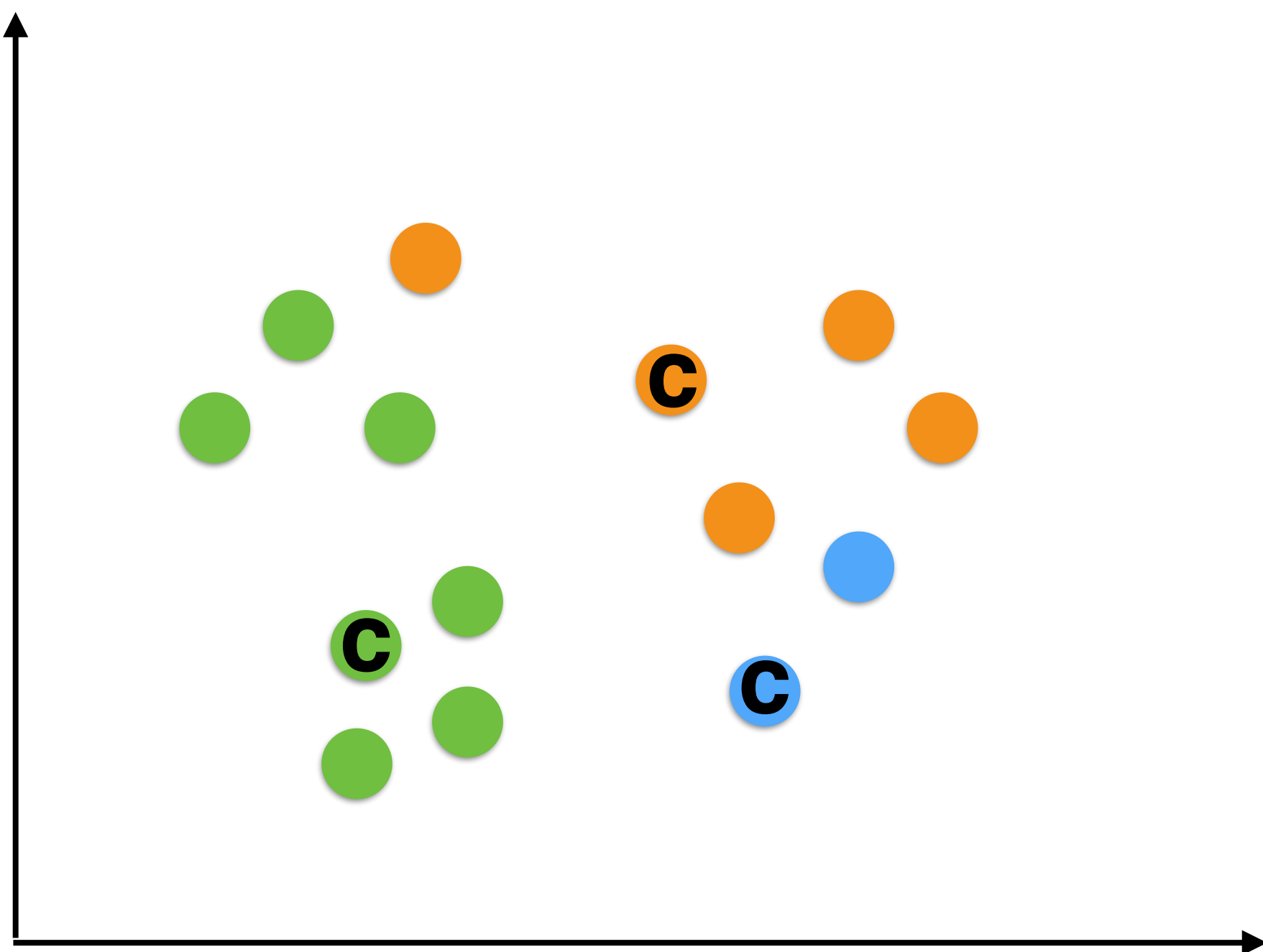
- choose k centroids





Kmeans Algorithm

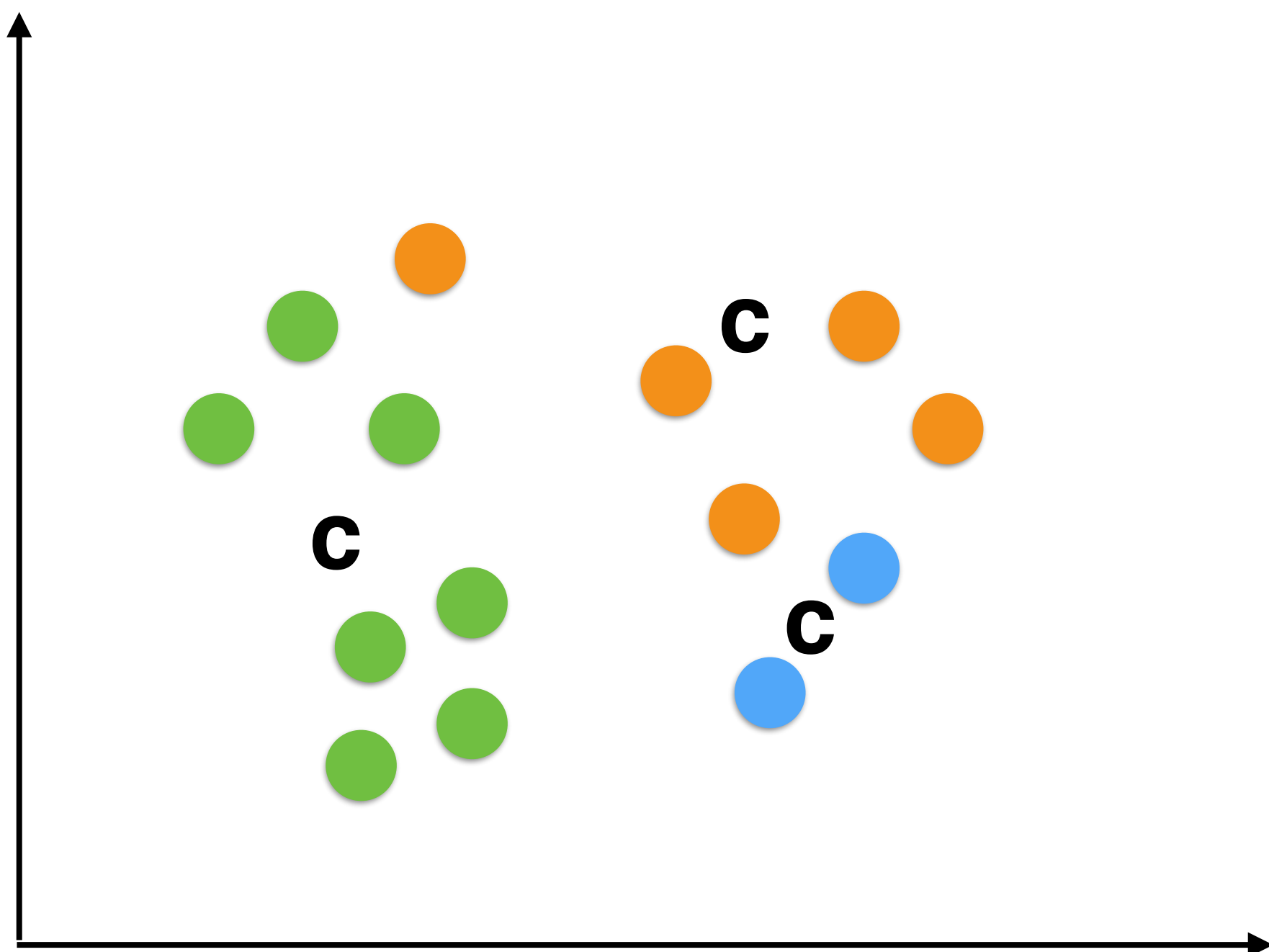
- 將資料點指定給距離最近的集群中心點





Kmeans Algorithm

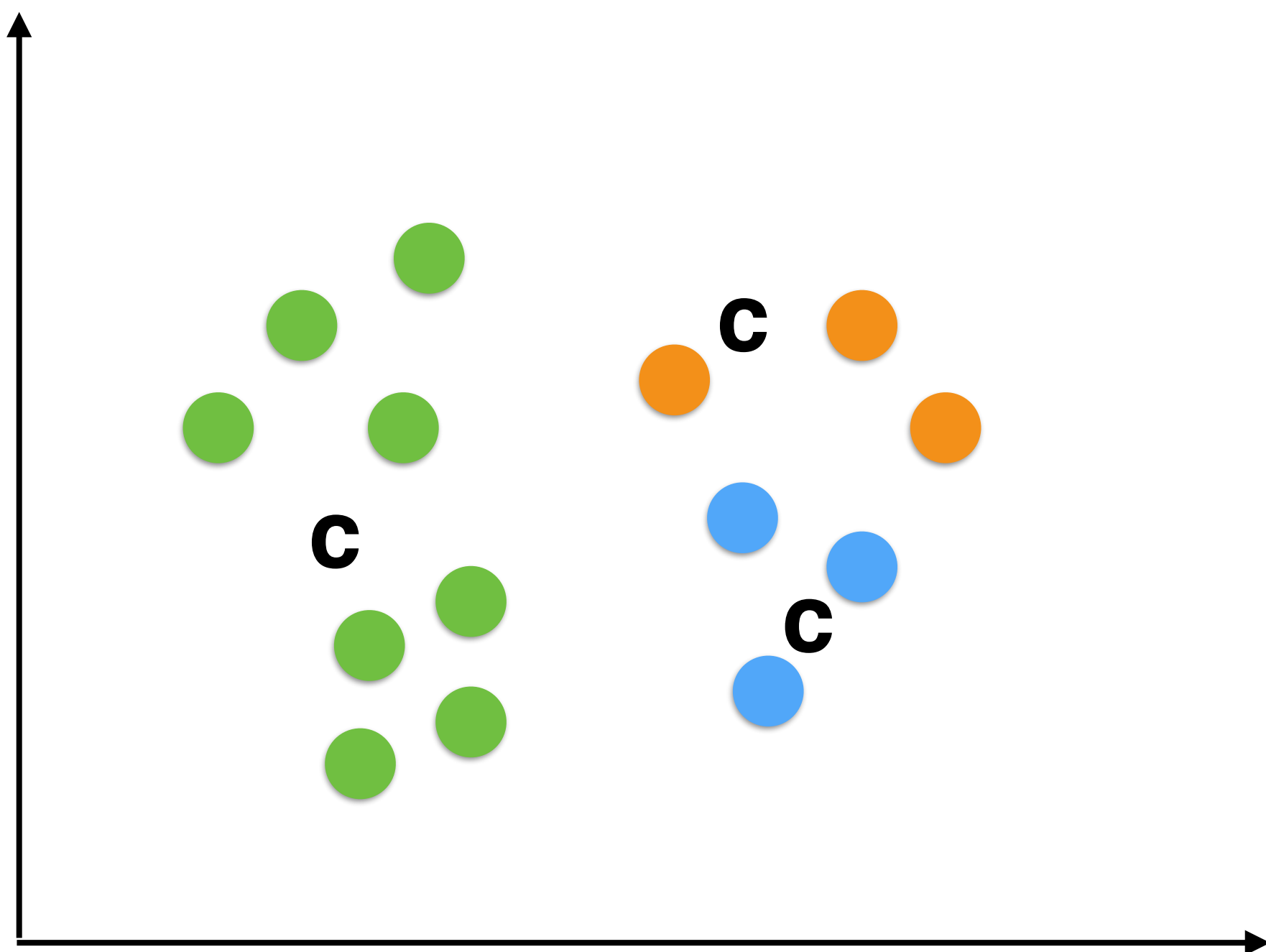
- 重新計算並設定新的集群中心點





Kmeans Algorithm

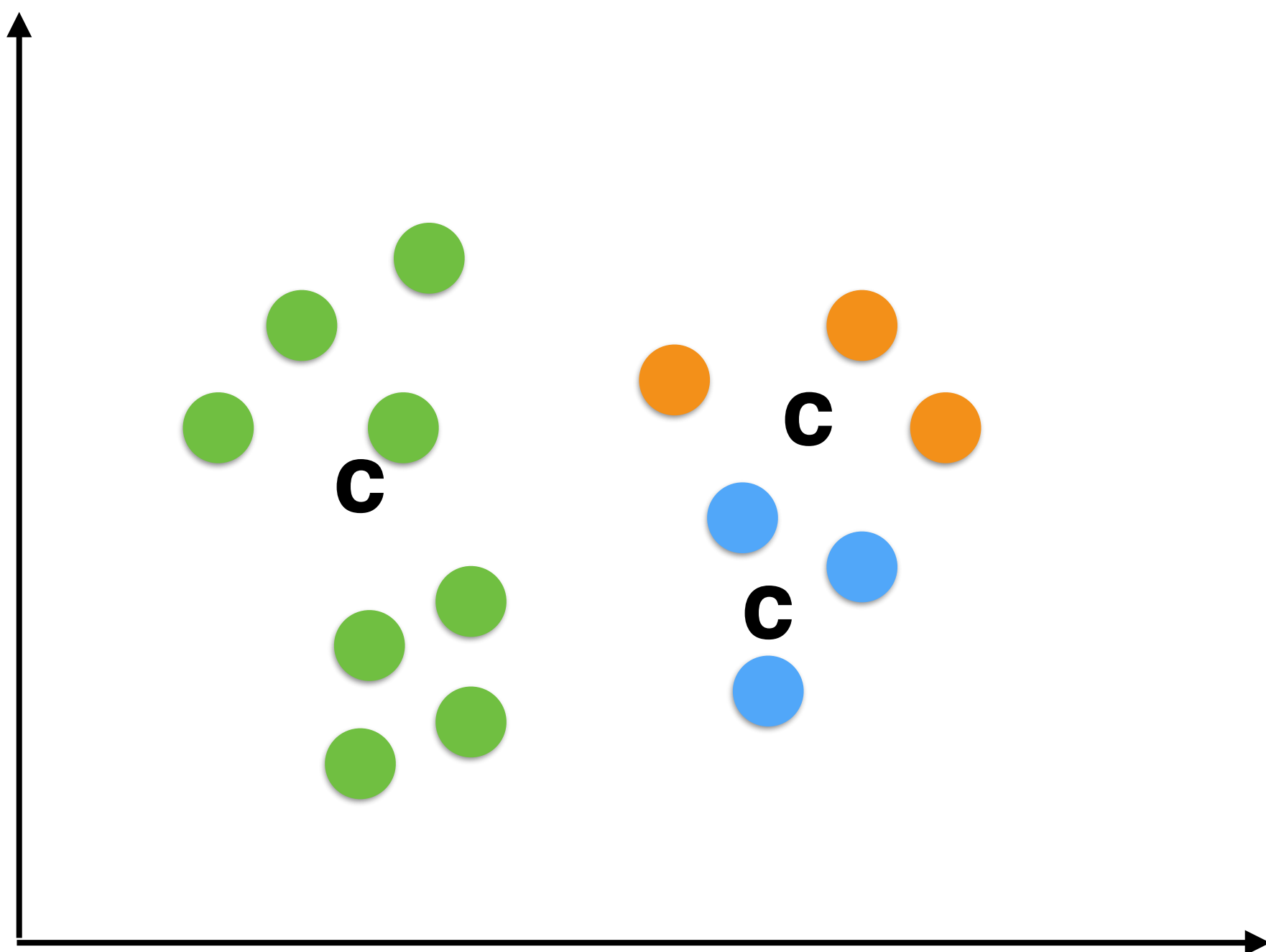
- 重新指定每個點到距離最近的新集群中心





Kmeans Algorithm

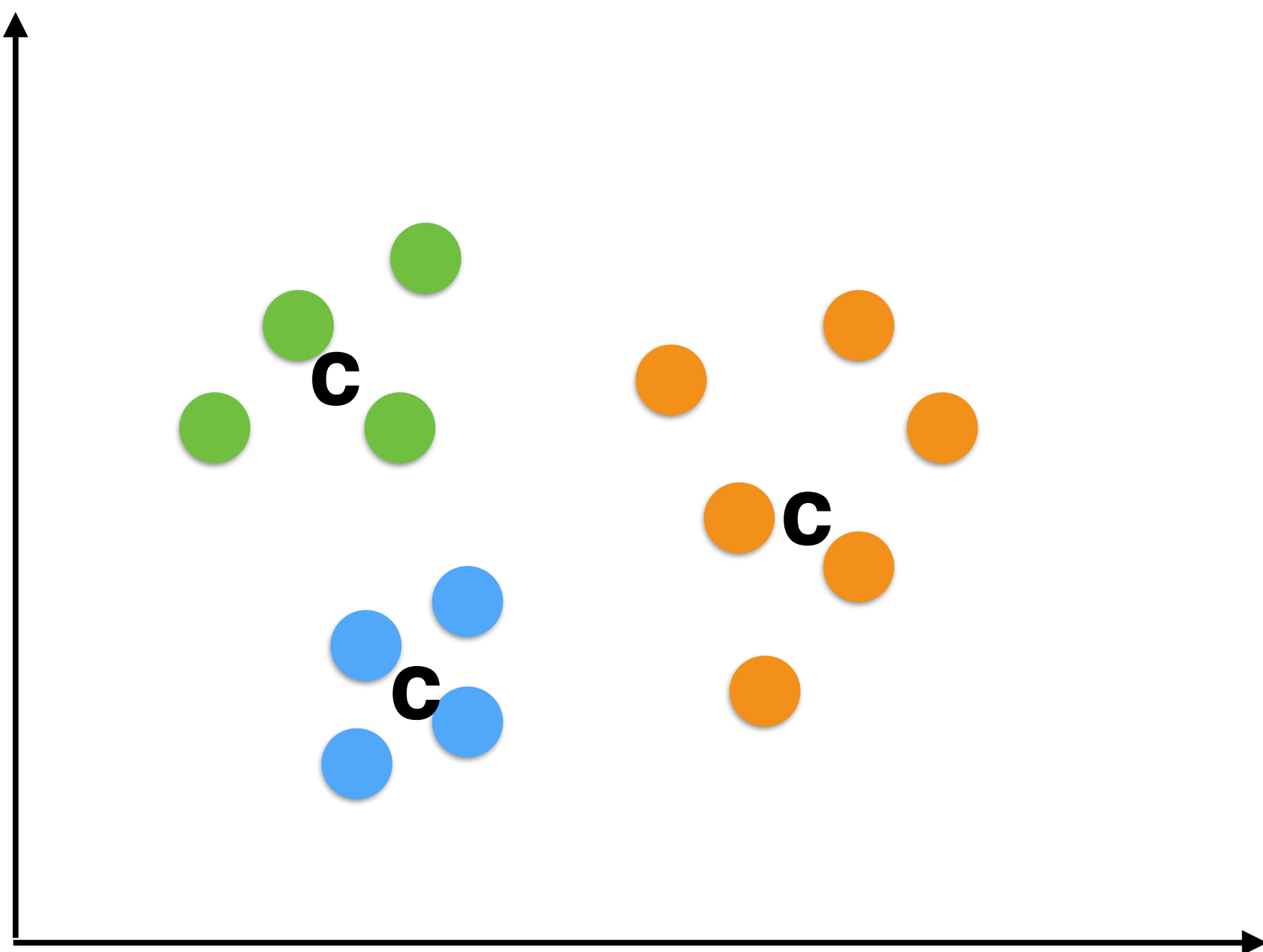
- 重新計算新的集群中心點





Kmeans Algorithm

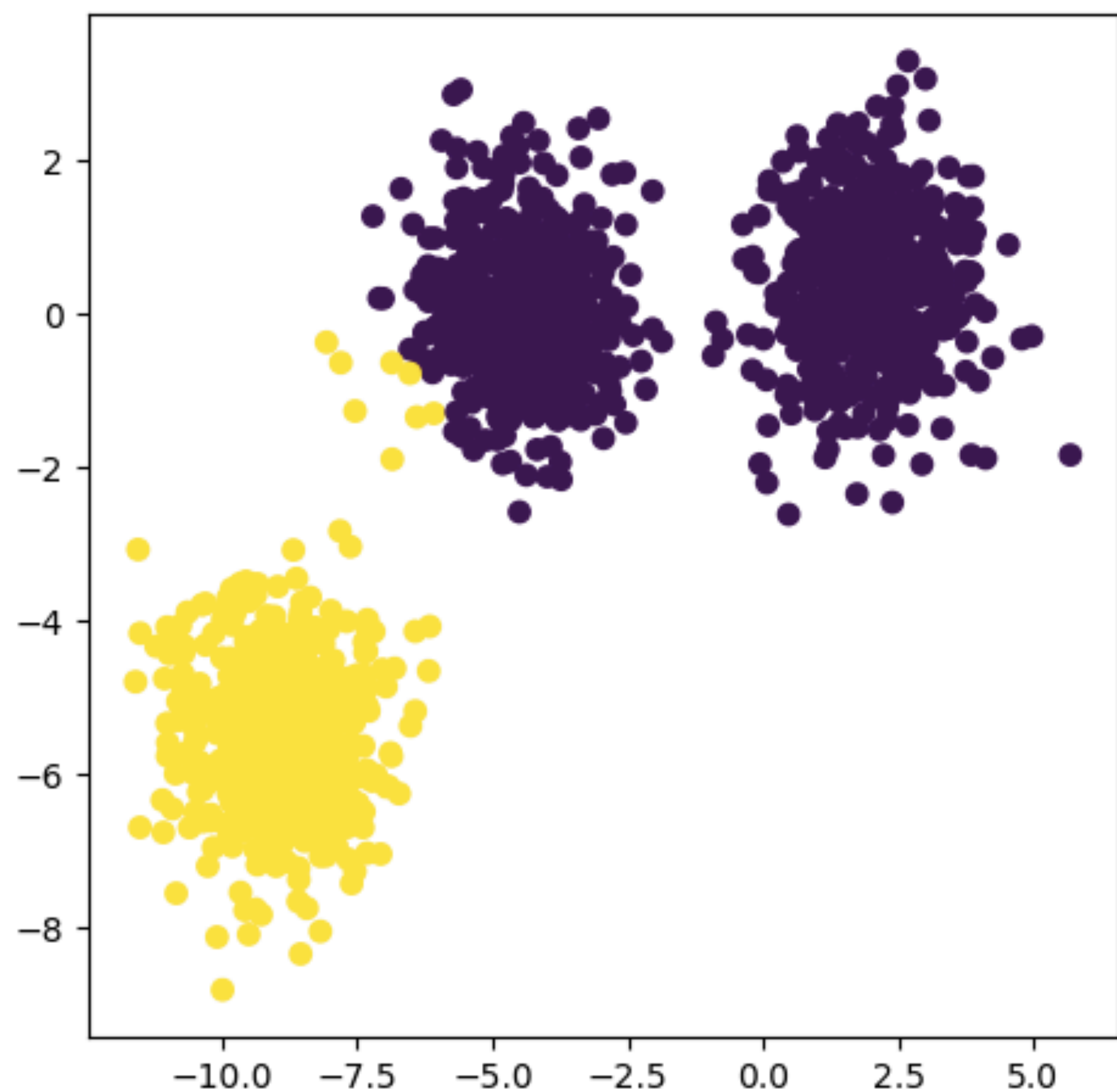
- 重複計算集群中心點、分群直達收斂(各點與集群中心距離和最小)





Problem 1 - k值

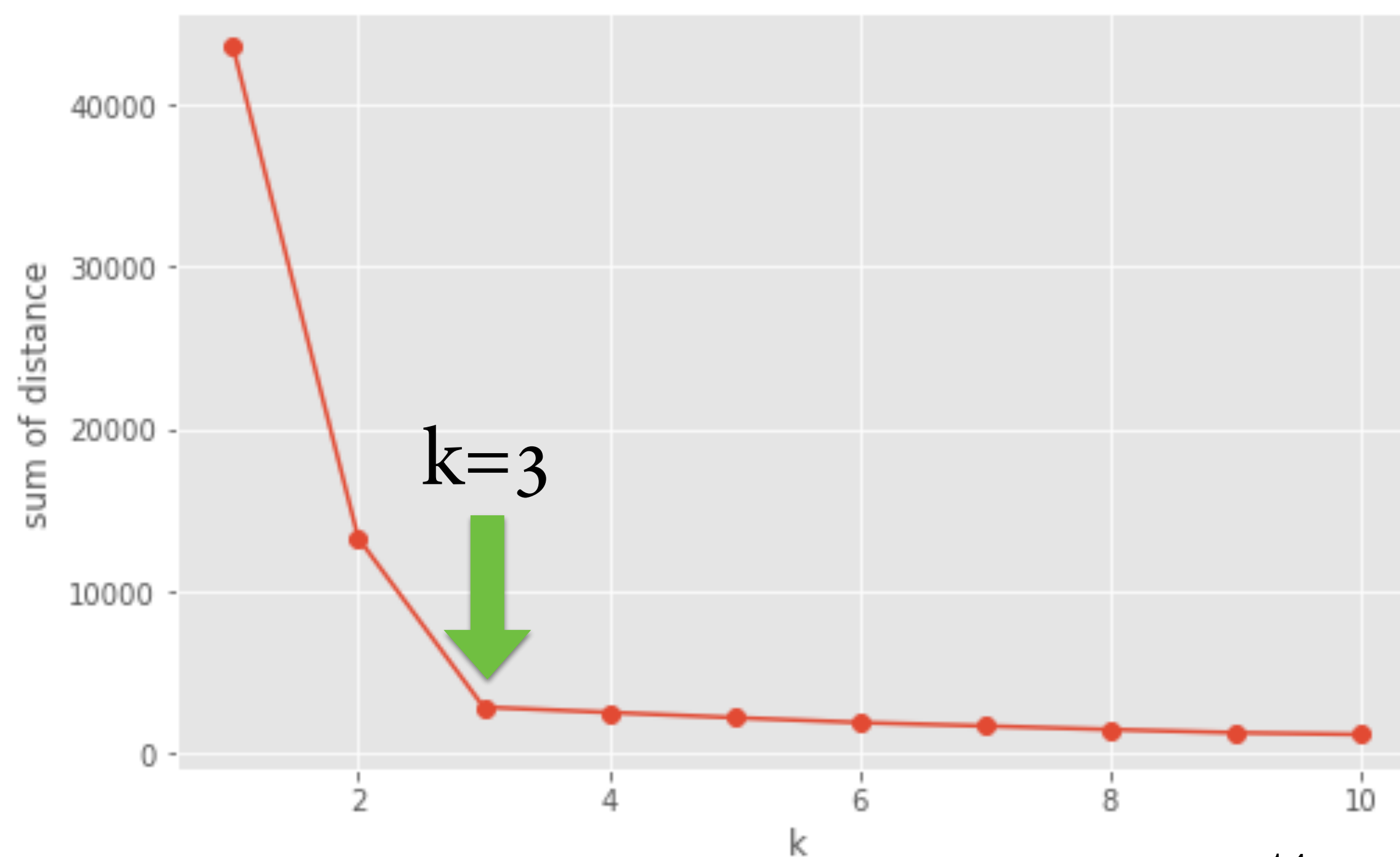
- 如何選擇好的k值?





Tuning k

- 若已經確定要分幾群(e.g. L、M、S)，可指定k值
- 計算不同k值之各集群總距離和





Problem 2 - 起始點與Kmeans++

- 隨機起始點若挑不好，可能導致不同的結果、花更多的時間才達到收斂
- 修正方法：Kmeans ++
 - 以距離較遠優先挑選初始中心點
 1. 隨機挑選一個中心點
 2. 計算所有資料對這些中心點的距離
 3. 將距離較遠的點加權(提高被挑中的機率)
 4. 重複以上步驟直到挑選完k個中心點再開始

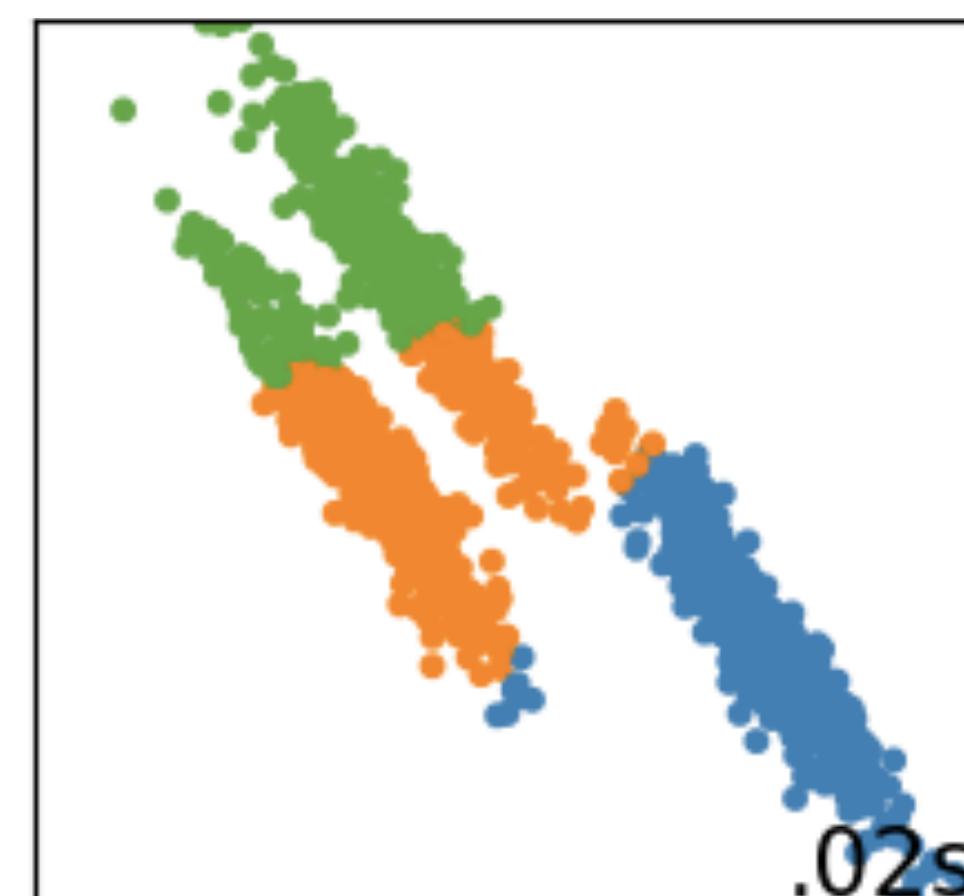
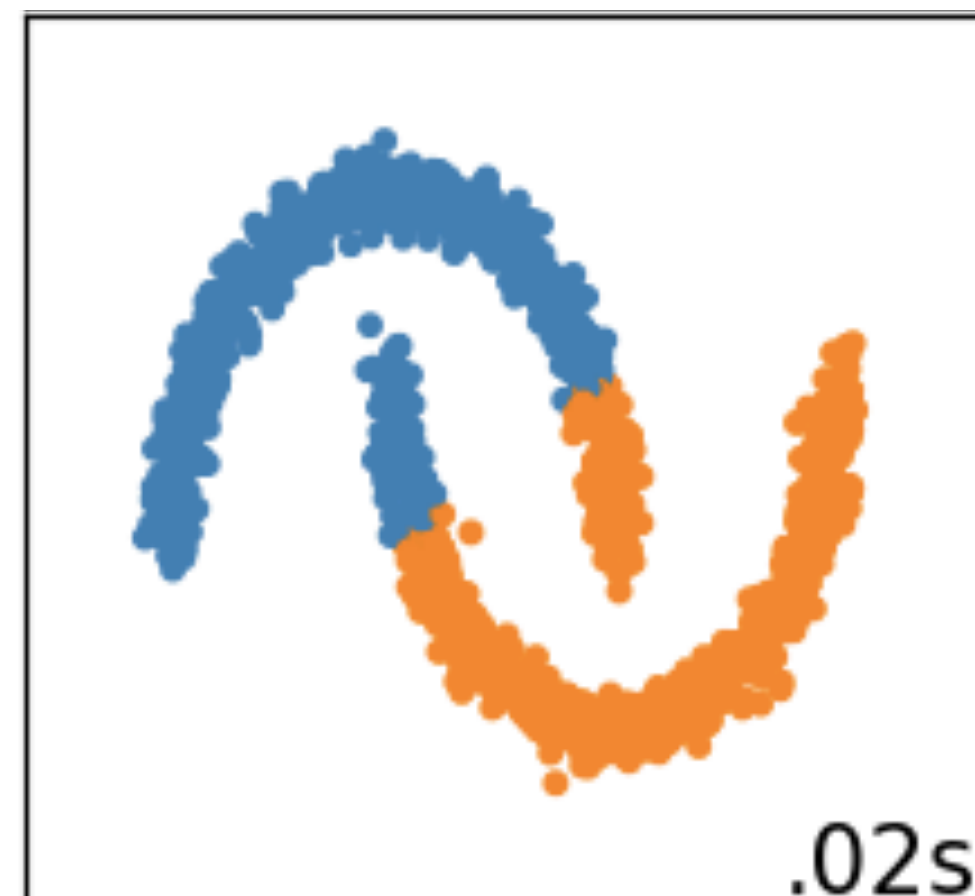
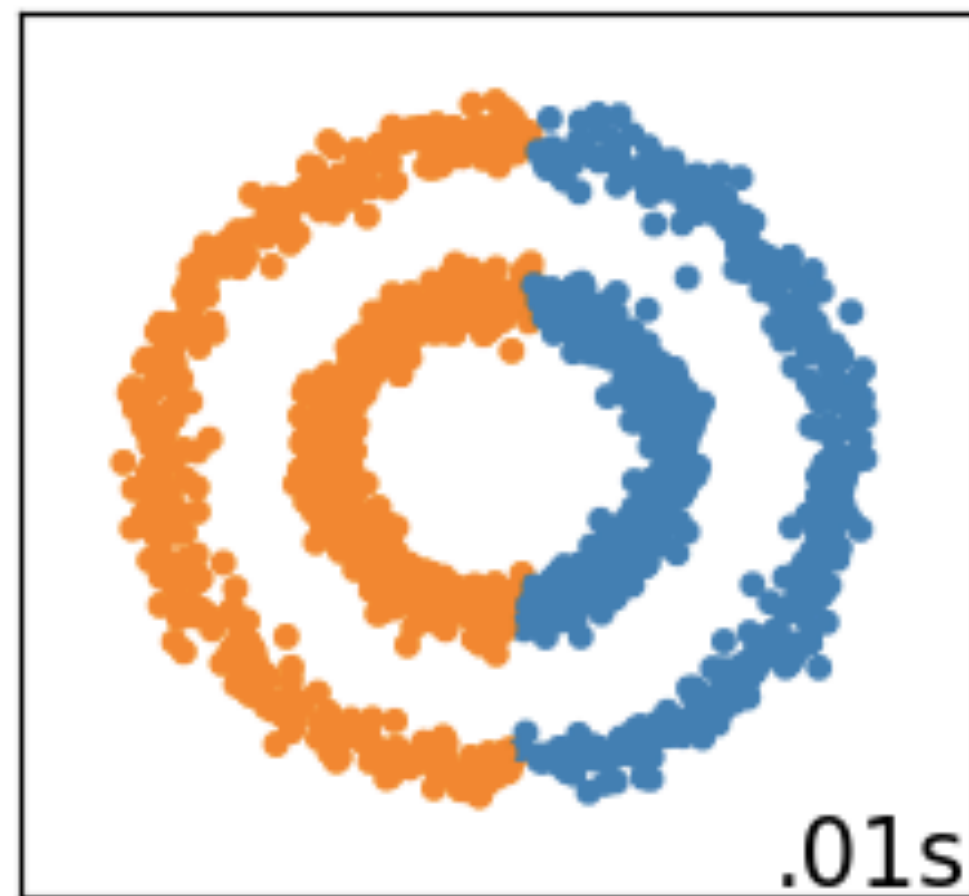


Kmeans with sklearn

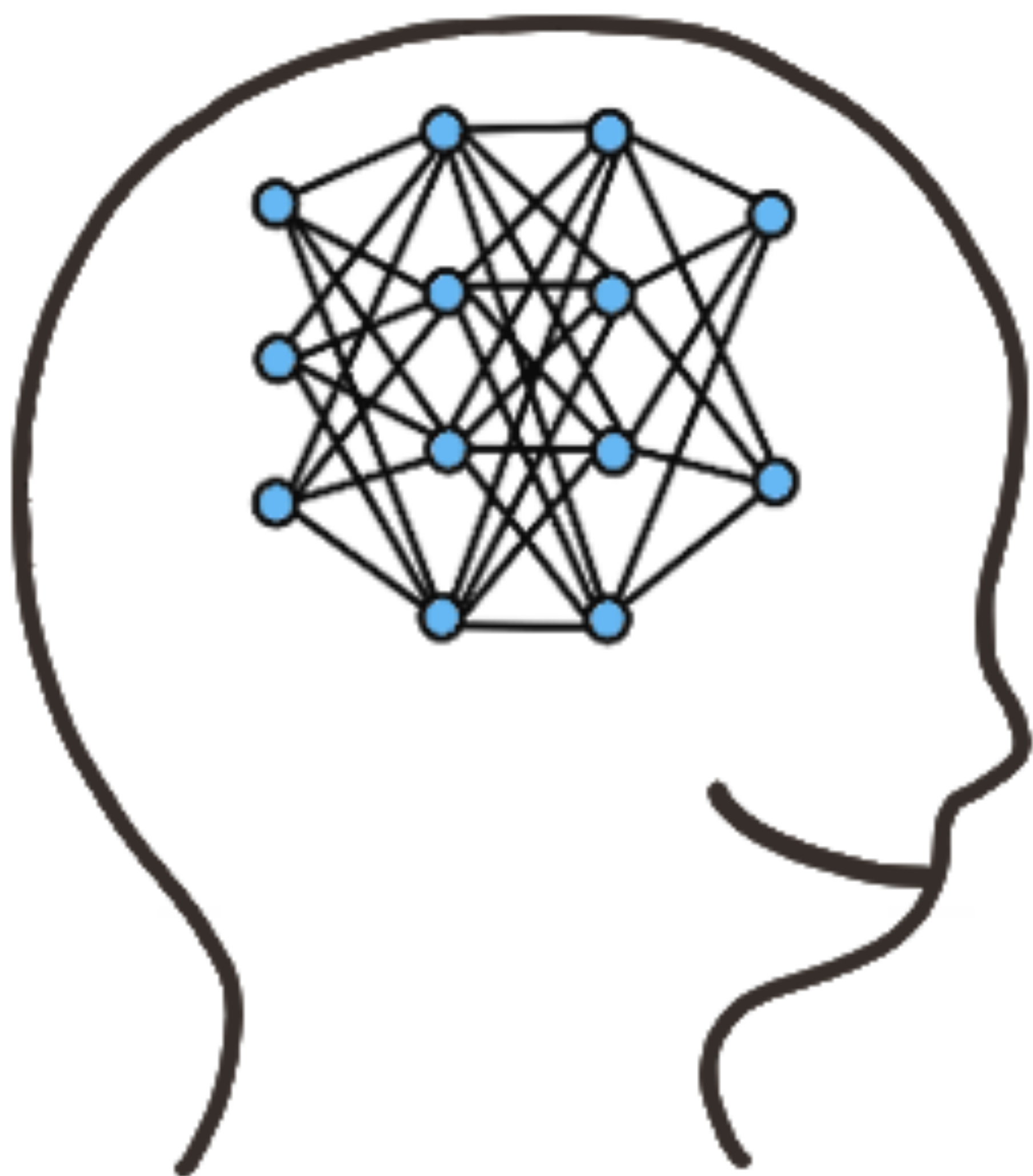
- `class sklearn.cluster.KMeans(n_clusters=8, init='k-means++', n_init=10, max_iter=300, tol=0.0001, precompute_distances='auto', verbose=0, random_state=None, copy_x=True, n_jobs=1, algorithm='auto')`
 - `n_clusters`: k 值(分群數)
 - `init`: 'k-means++', 'random'(原始隨機挑選初始中心點方法)
 - `max_iter`: 最多迭代次數(可避免出現無法收斂的情況)
- `attributes`:
 - `cluster_centers_`: 各集群中心點
 - `labels_`: 各資料點屬於的集群標注
 - `inertia_`: 資料與各集群距離總和

其他問題

- 無法正確分群非圓狀的資料散佈，還可能無法成功收斂



- 容易受極端值影響

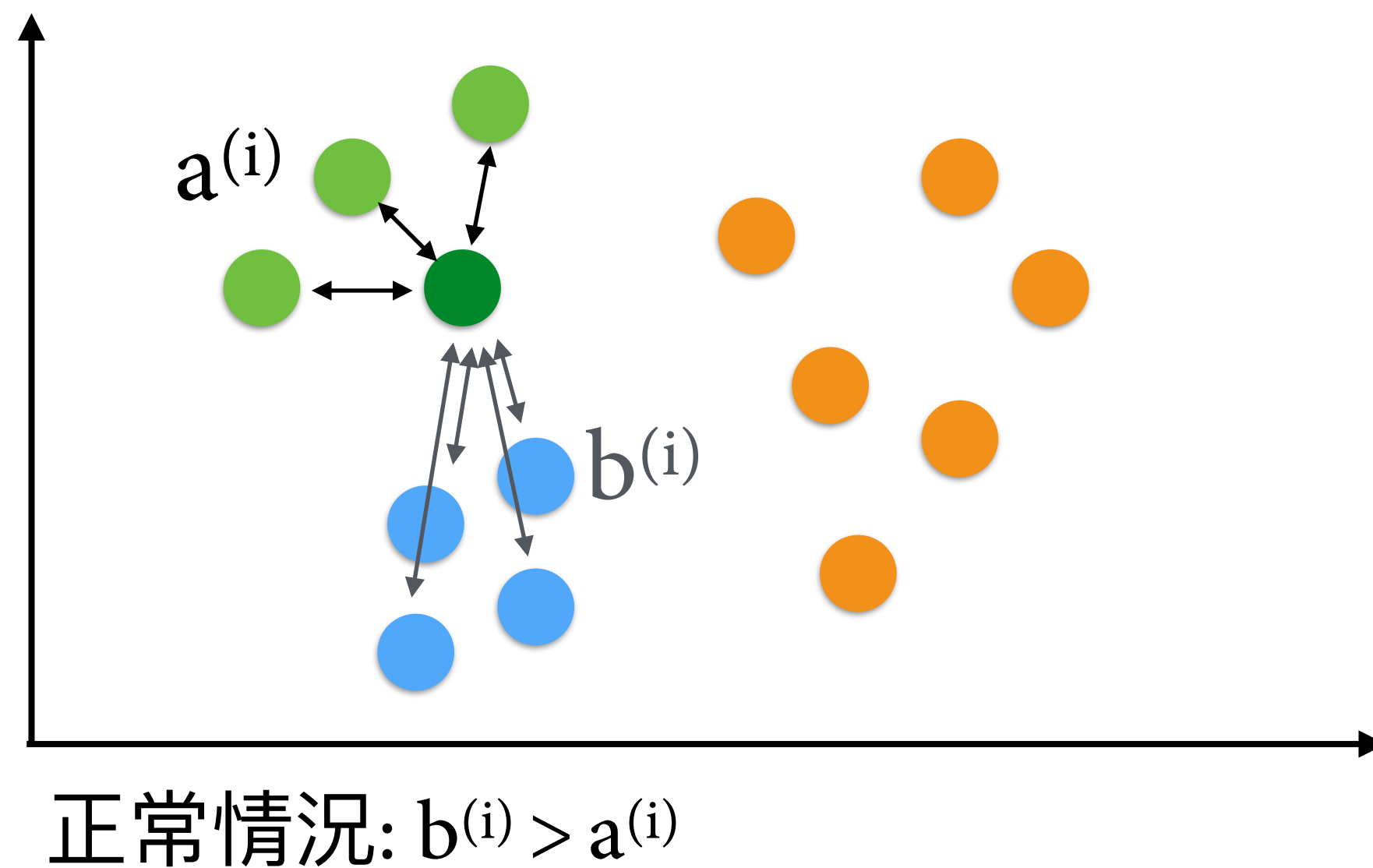


輪廓圖分析

Silhouette Analysis

輪廓係數 (silhouette coefficient)

- 集群內聚(對內) $a^{(i)}$: 該資料與同集群中其他樣本的平均距離
- 集群分離(對外) $b^{(i)}$: 該資料與其他最近集群中所有樣本的平均距離



輪廓係數 (silhouette coefficient)

- 輪廓係數

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

$$s(i) = \begin{cases} 1 - a(i)/b(i), & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ b(i)/a(i) - 1, & \text{if } a(i) > b(i) \end{cases}$$

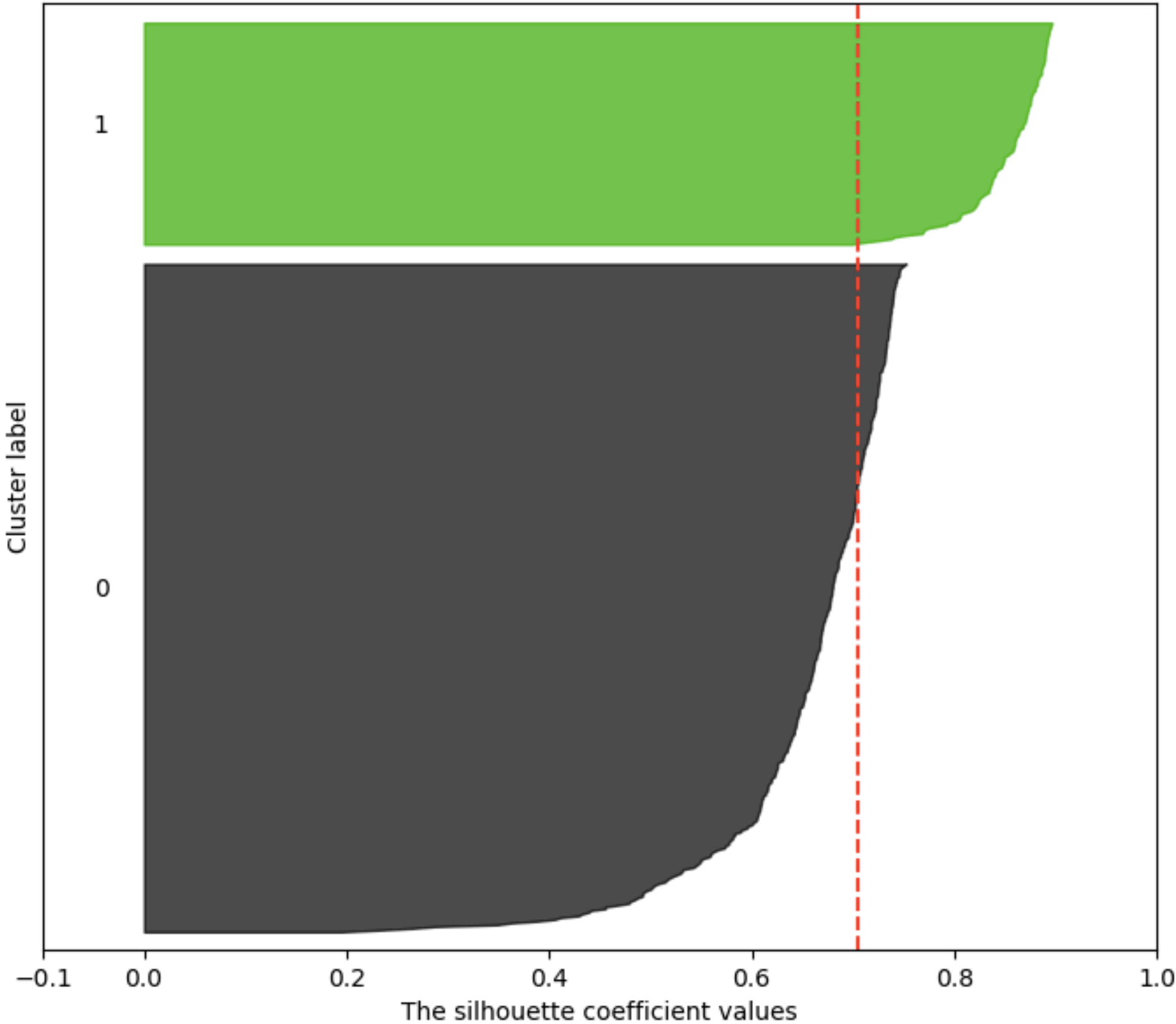
- $s^{(i)}$ 範圍：-1 ~ +1
- $s^{(i)}=0$: 內聚程度=分離程度
- $s^{(i)}=1$: 理想值(分離程度>>內聚程度)



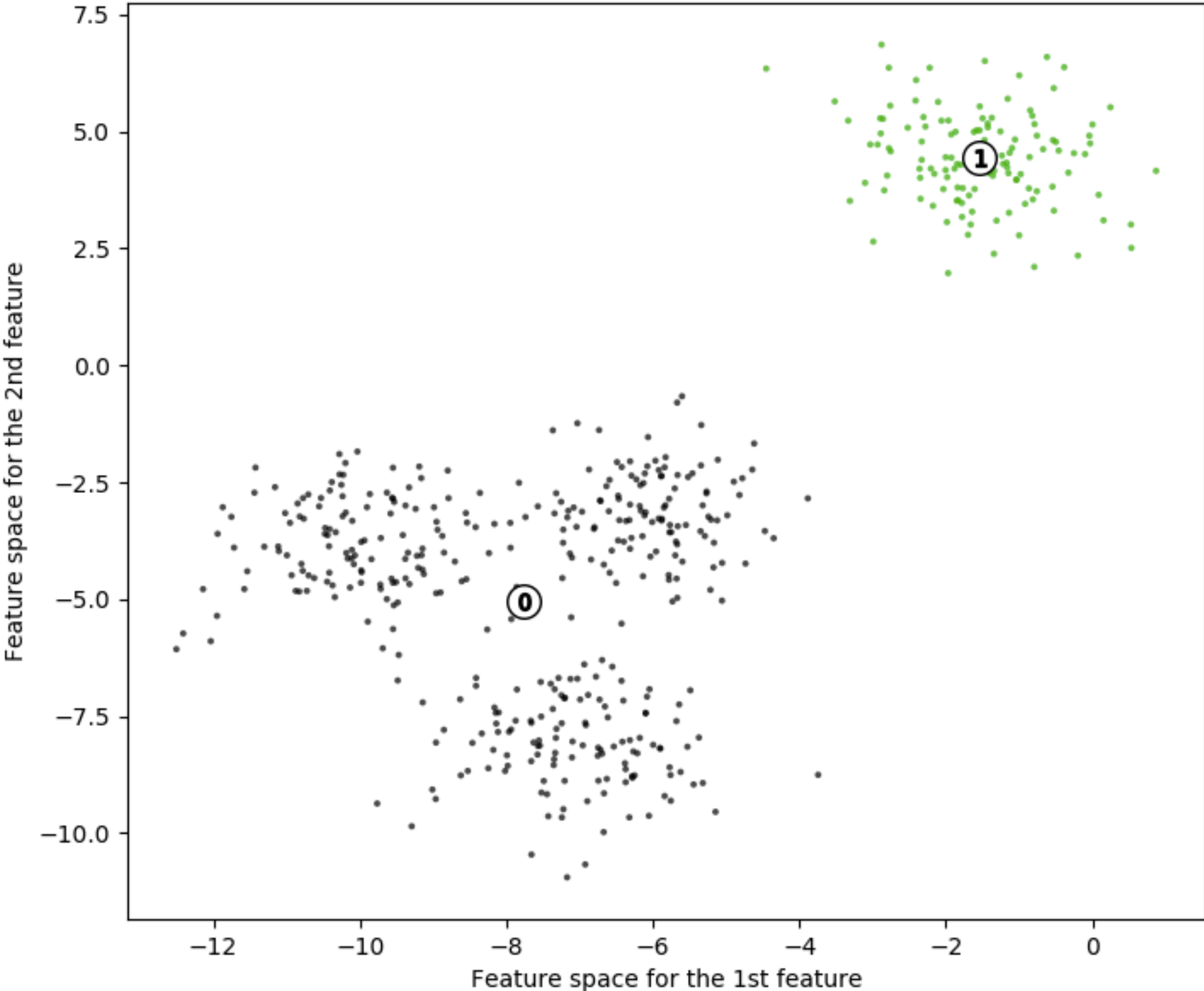
輪廓圖

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 2$

The silhouette plot for the various clusters.



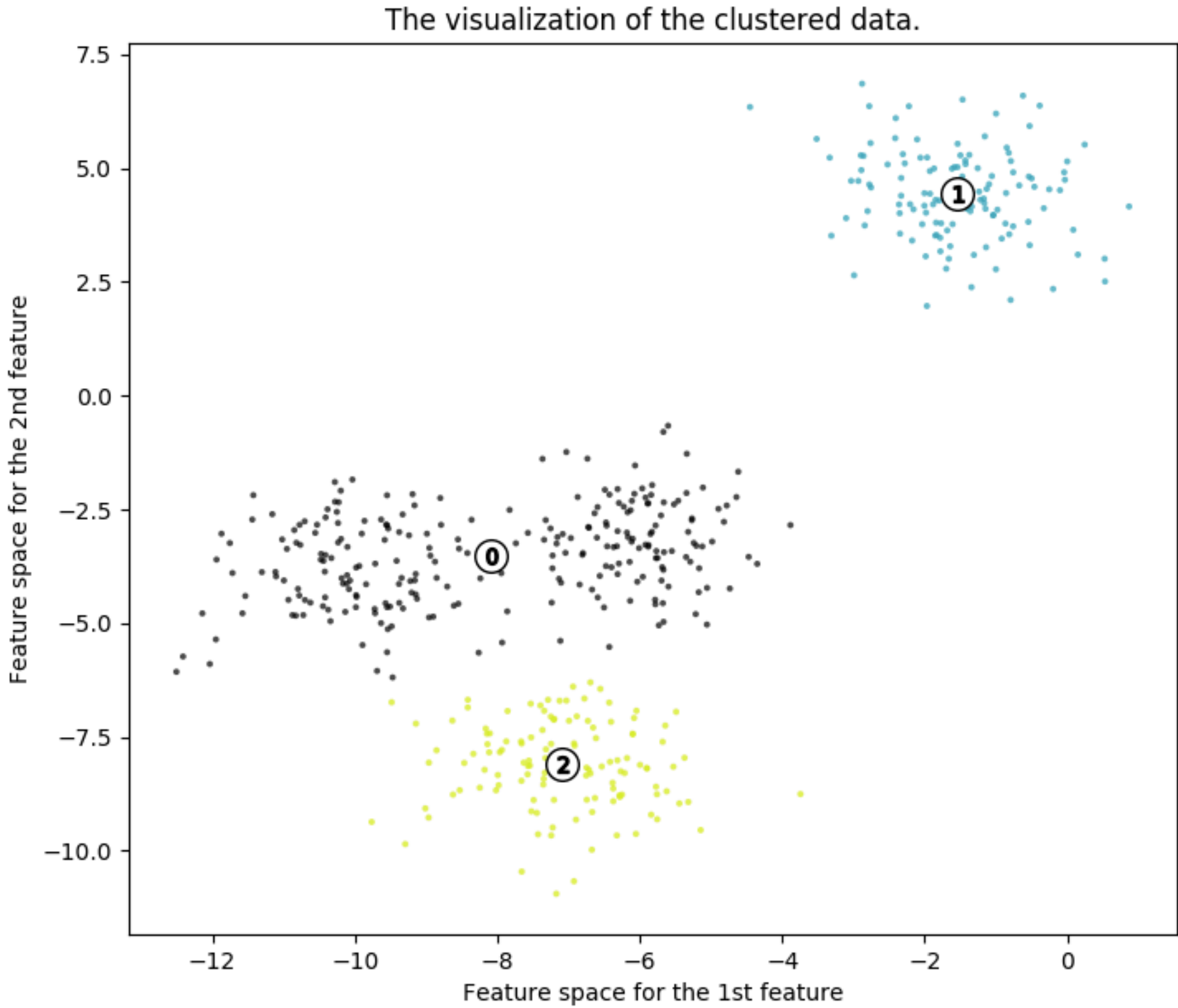
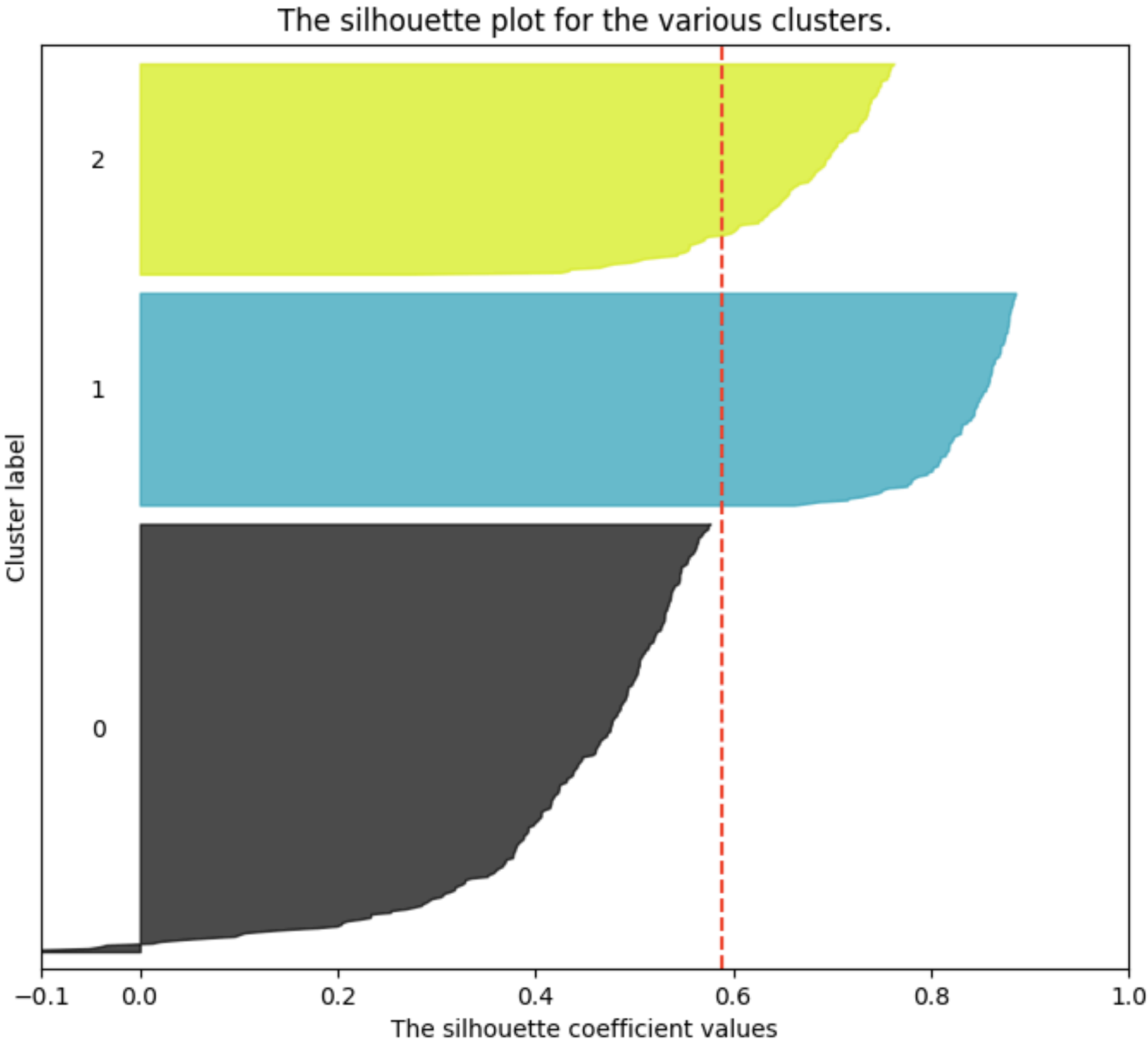
The visualization of the clustered data.





輪廓圖

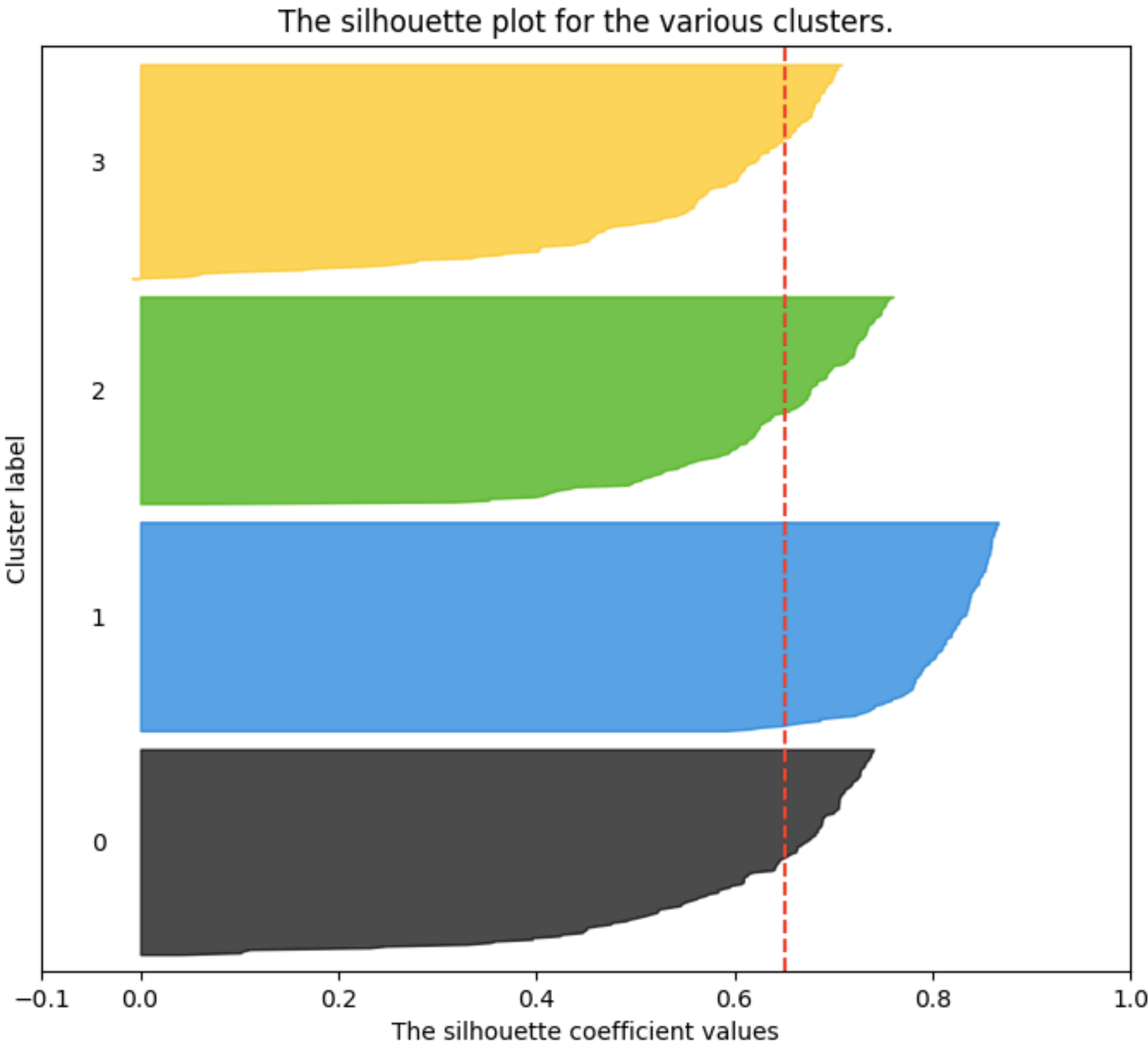
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 3$





輪廓圖

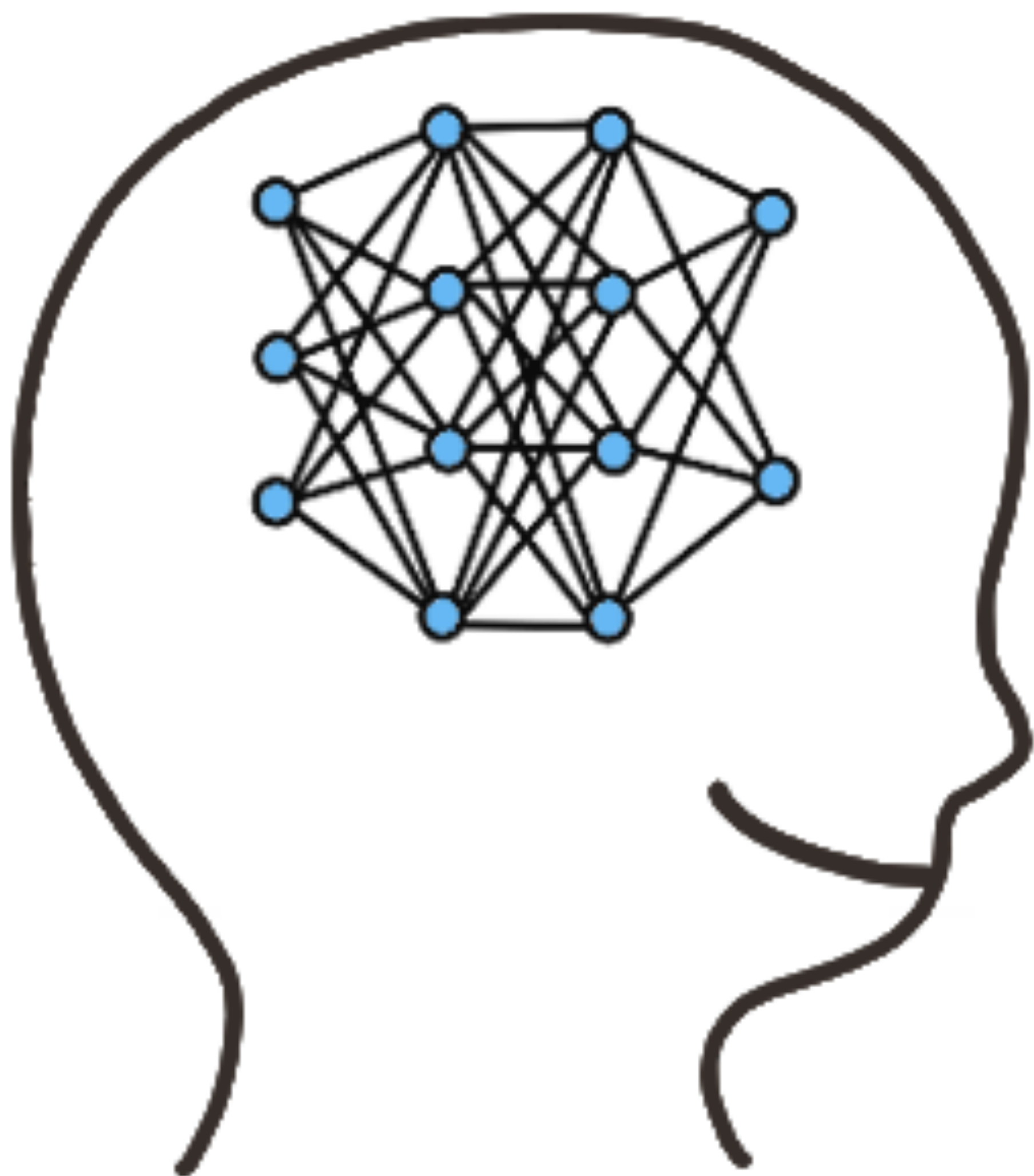
Silhouette analysis for KMeans clustering on sample data with $n_clusters = 4$





Silhouette coefficient with sklearn

- `sklearn.metrics.silhouette_samples(X, labels, metric='euclidean', **kwargs)`
 - 計算每筆樣本的輪廓係數(Silhouette Coefficient)
- `sklearn.metrics.silhouette_score(X, labels, metric='euclidean', sample_size=None, random_state=None, **kwargs)`
 - 計算所有樣本的輪廓係數平均(Mean of Silhouette Coefficient)



階層式分群

Hierarchical Clustering



階層式分群

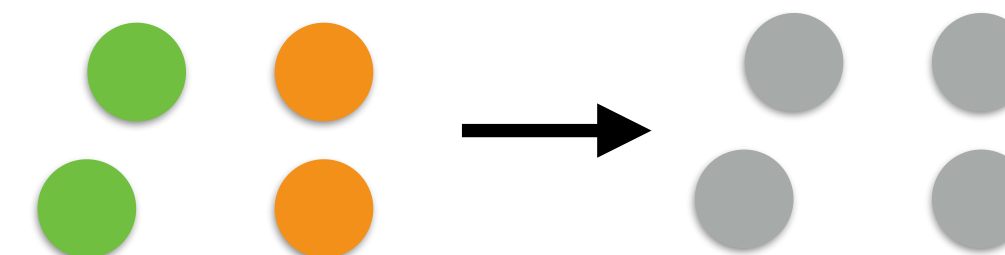
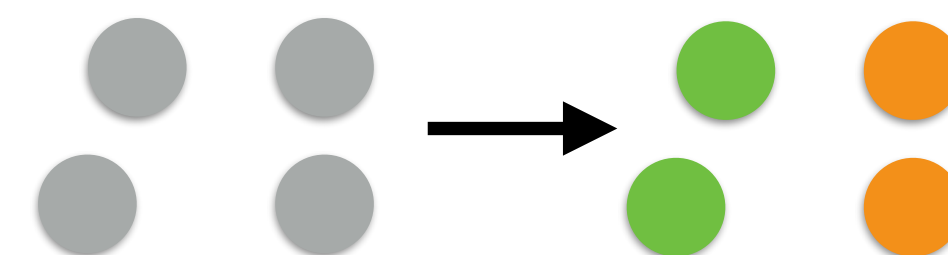
- 兩種方式：

- Top-down/群由大變小(分裂)

- 分離分層集群(Divisive Hierarchical Clustering)

- Bottom-up/群由小變大(合併)

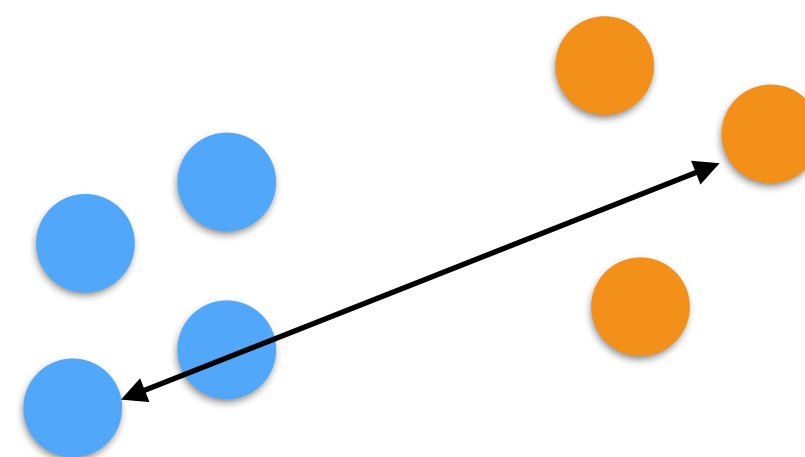
- 凝聚分層集群(Agglomerative Hierarchical Clustering)





凝聚分層

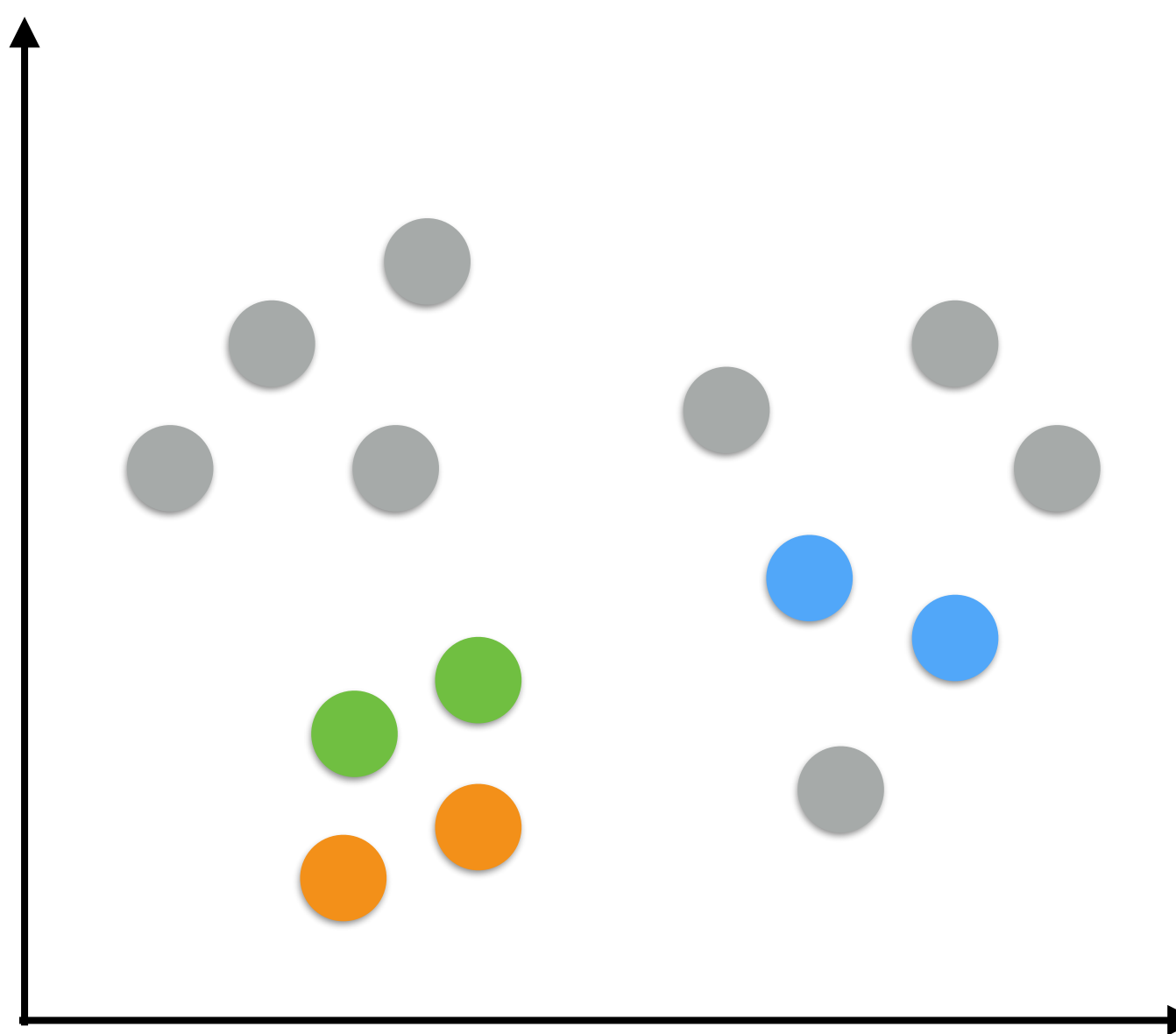
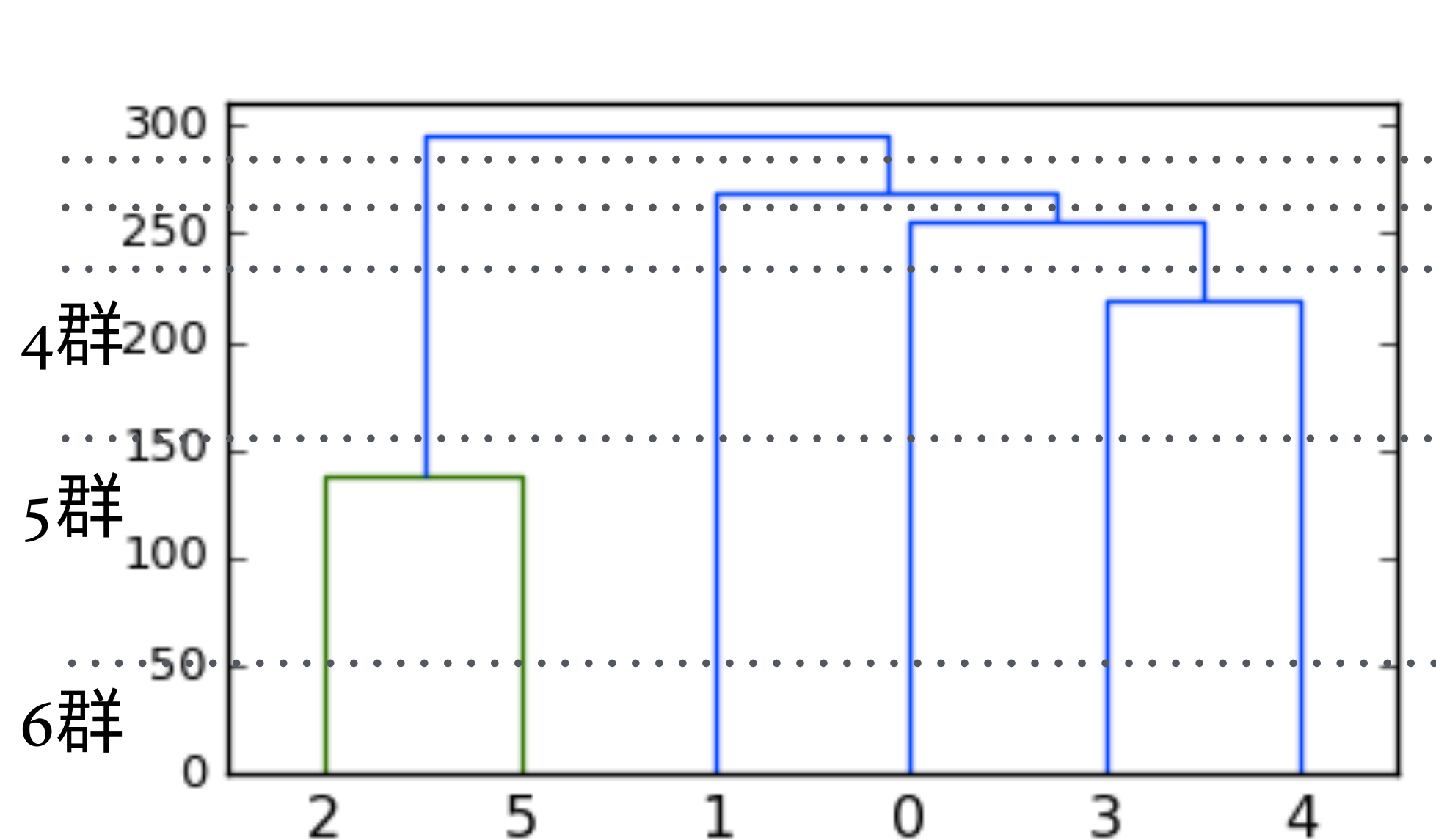
- 每個資料點都當成獨自一群
- 從彼此距離最近的群開始合併，直到產生一定量的集群
- 計算距離方式
 - Ward: 集群聚合後的變異(variance)/距離平方和最小
 - Average: 集群中各點與各點資料點的平均距離
 - Complete: 集群中兩點最遠的距離





凝聚分層

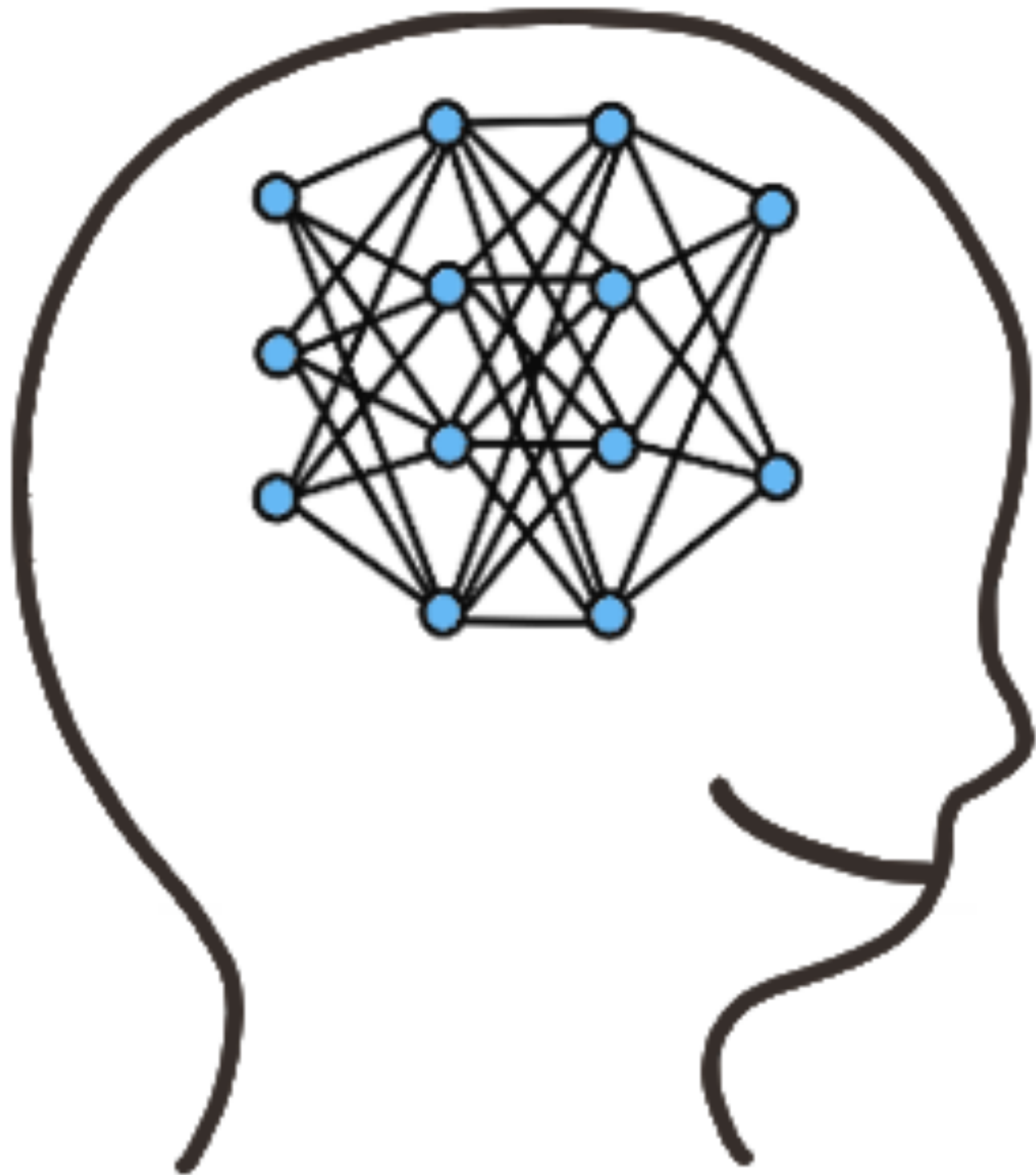
- 從每個點都是獨立集群開始依距離——合併集群





Agglomerative Clustering with sklearn

- `class sklearn.cluster.AgglomerativeClustering(n_clusters=2, affinity='euclidean', memory=None, connectivity=None, compute_full_tree='auto', linkage='ward', pooling_func=<function mean>)`
 - `n_clusters`: 集群數
 - `affinity`: 距離計算方法 euclidean(歐幾里德距離)
 - `linkage`: 'ward', 'average', 'complete'
- `scipy.cluster.hierarchy.dendrogram(Z, p=30, truncate_mode=None, color_threshold=None, get_leaves=True, orientation='top', labels=None, count_sort=False, distance_sort=False, show_leaf_counts=True, no_plot=False, no_labels=False, leaf_font_size=None, leaf_rotation=None, leaf_label_func=None, show_contracted=False, link_color_func=None, ax=None, above_threshold_color='b')`
 - 畫階層樹狀圖



DBSCAN

Density-based Spatial Clustering
of Applications with Noise



DBSCAN

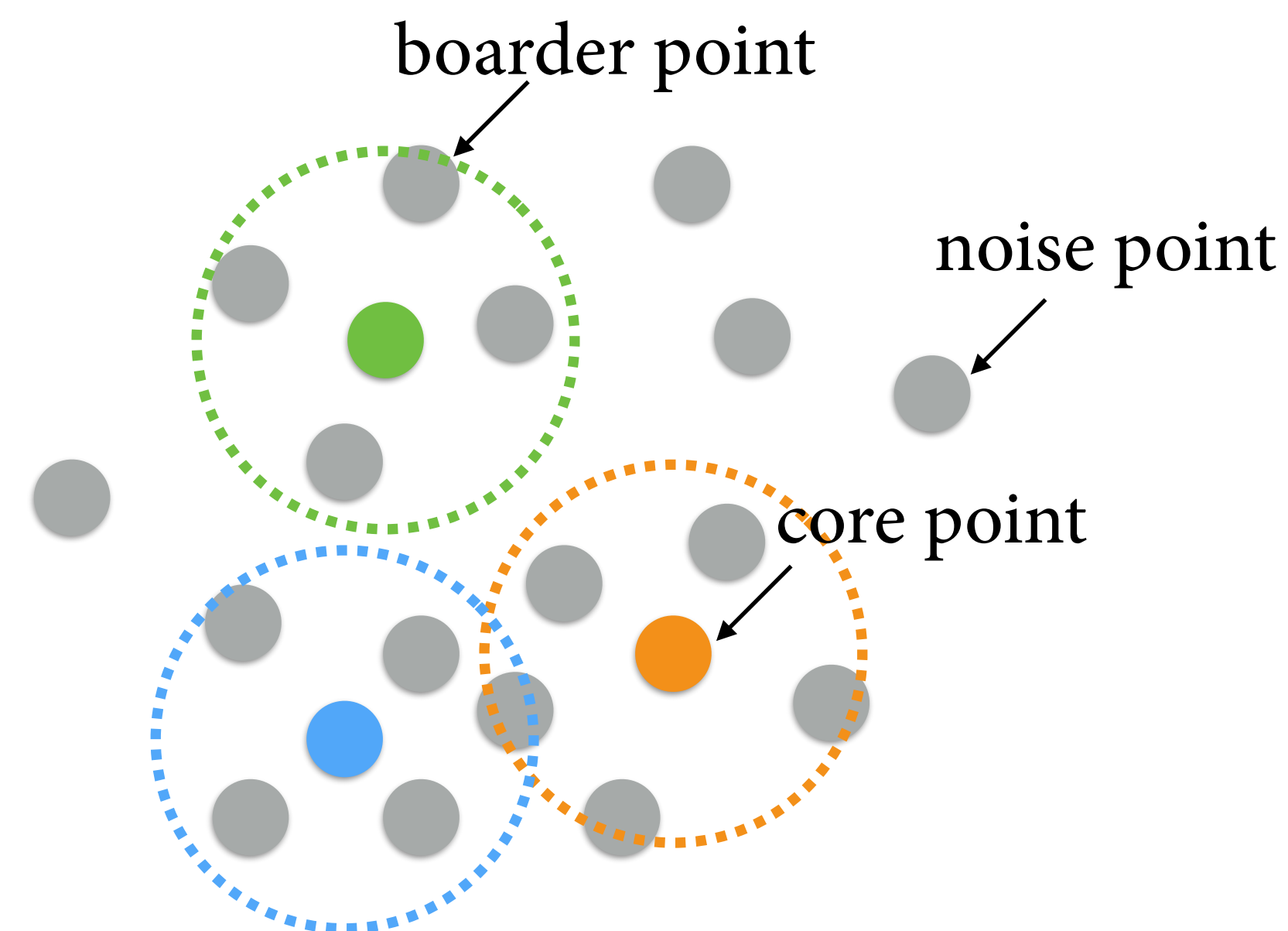
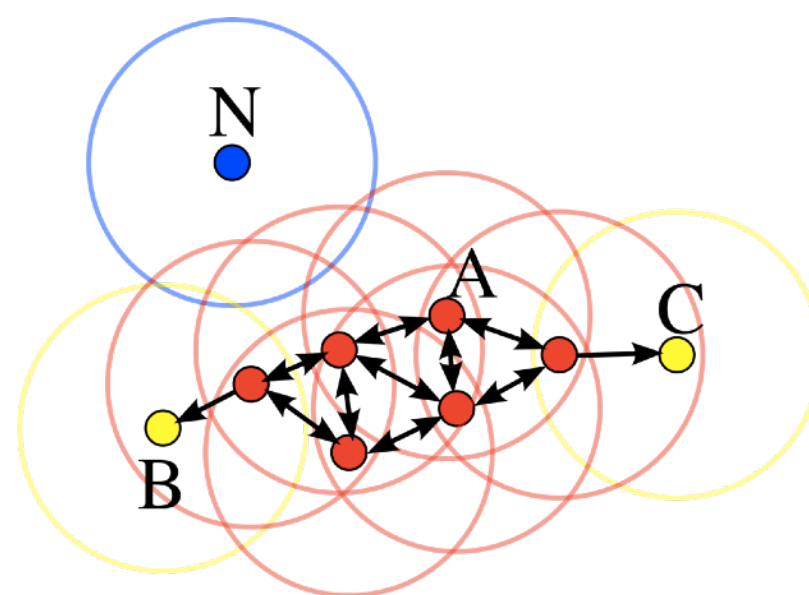
- 根據密度(指定半徑內)形成集群

- 核心點(core point): 半徑內包含超過密度域值(minPts)的點
- 邊緣點(border point): 核心點半徑內包含的點

- 雜訊點(noise point)

- 集群方法

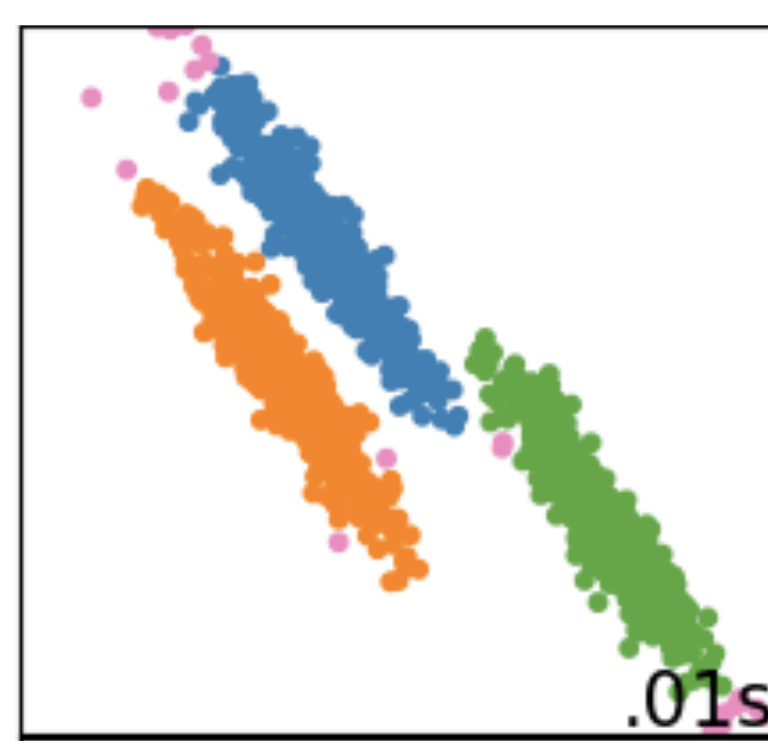
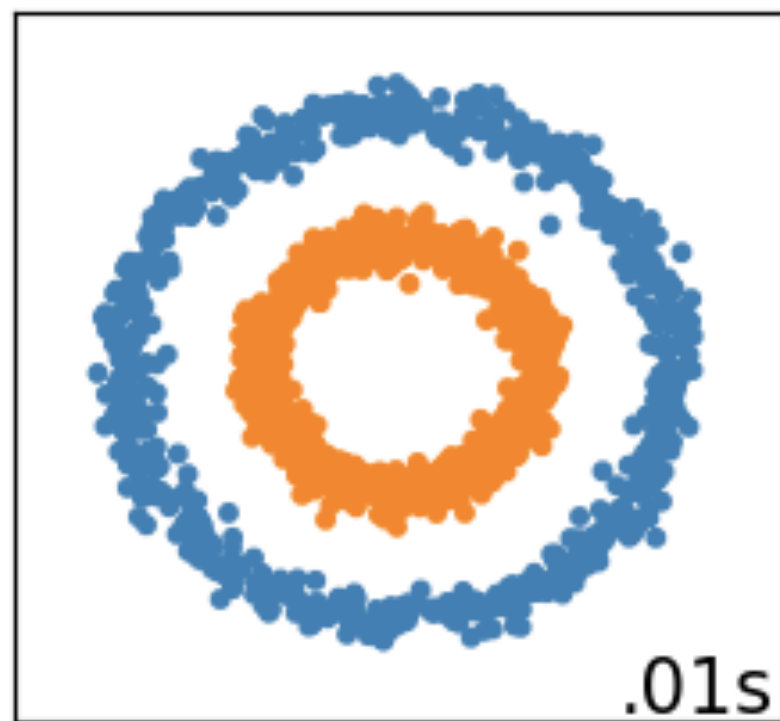
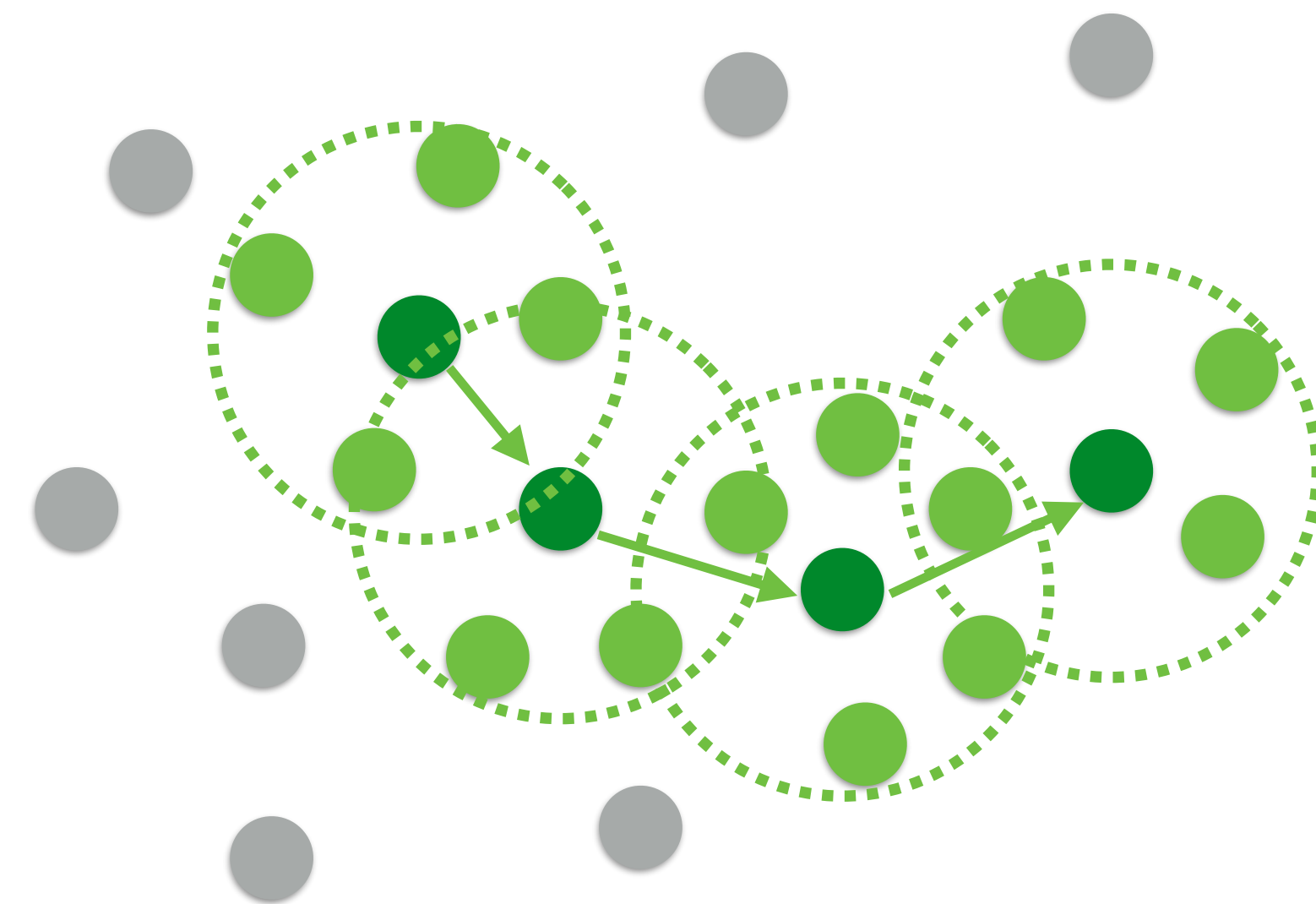
- 每個核心點或兩核心點距離小於半徑，形成一集群
- 指定邊緣點到核心點集群





特點

- 可判別雜訊點(noise)，使分群不受雜訊影響
- 不用預先設定分群數
- 可順利將特殊形狀分群(如: 半月形、甜甜圈形)





DBSCAN with sklearn

- `class sklearn.cluster.DBSCAN(eps=0.5, min_samples=5, metric='euclidean', metric_params=None, algorithm='auto', leaf_size=30, p=None, n_jobs=1)`
 - eps: 半徑
 - min_samples: minPts 半徑內至少要有多少點才能被判定為core point
- attributes:
 - labels_: 分群標籤(若為noise則回傳-1)