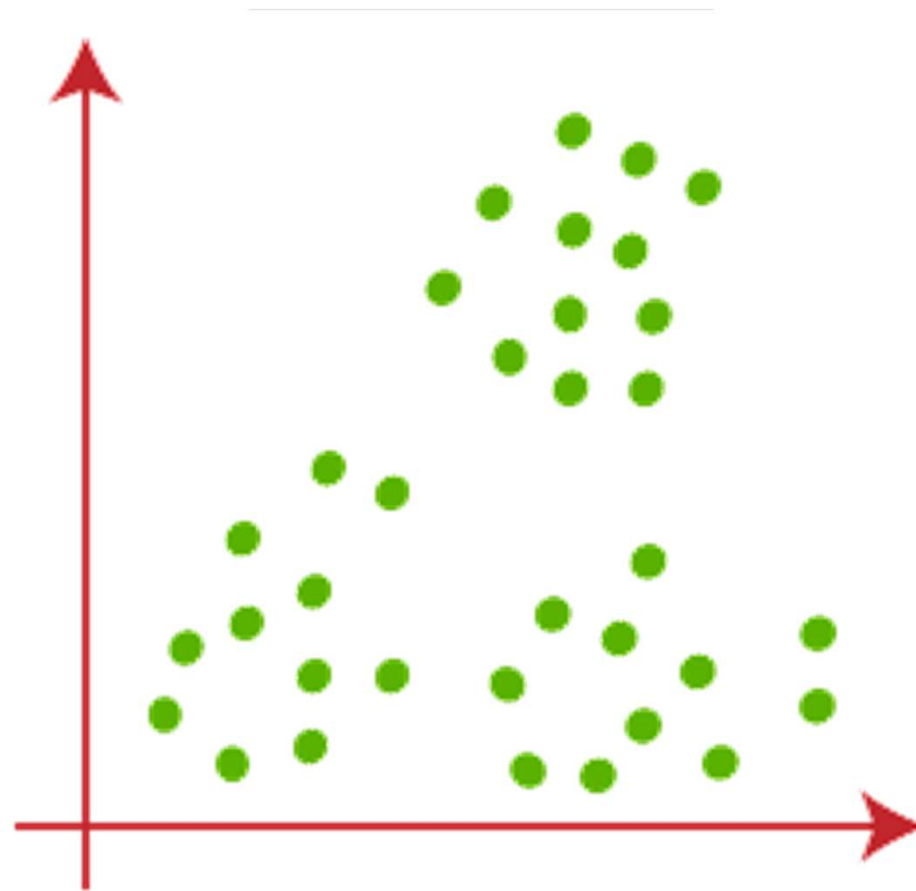


K-means clustering

Python

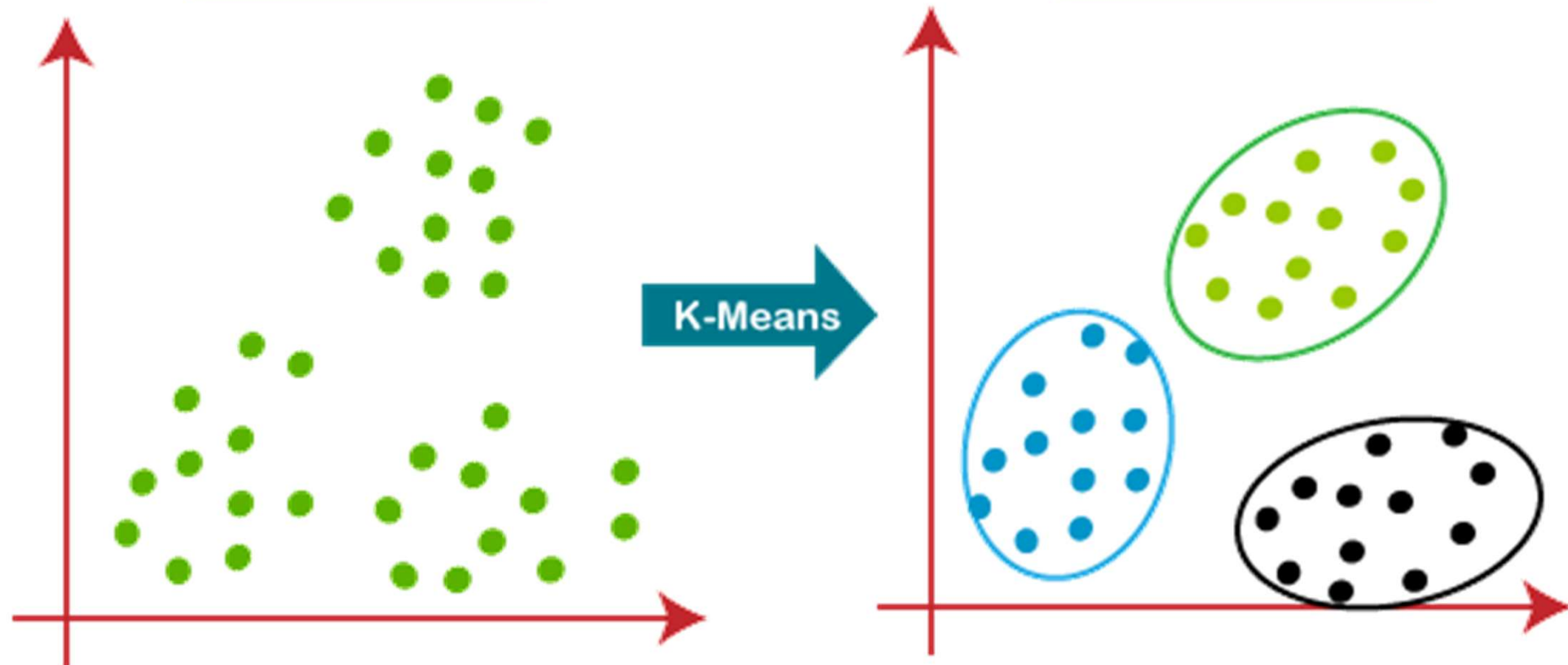
Wrocław, 2022



Before K-Means

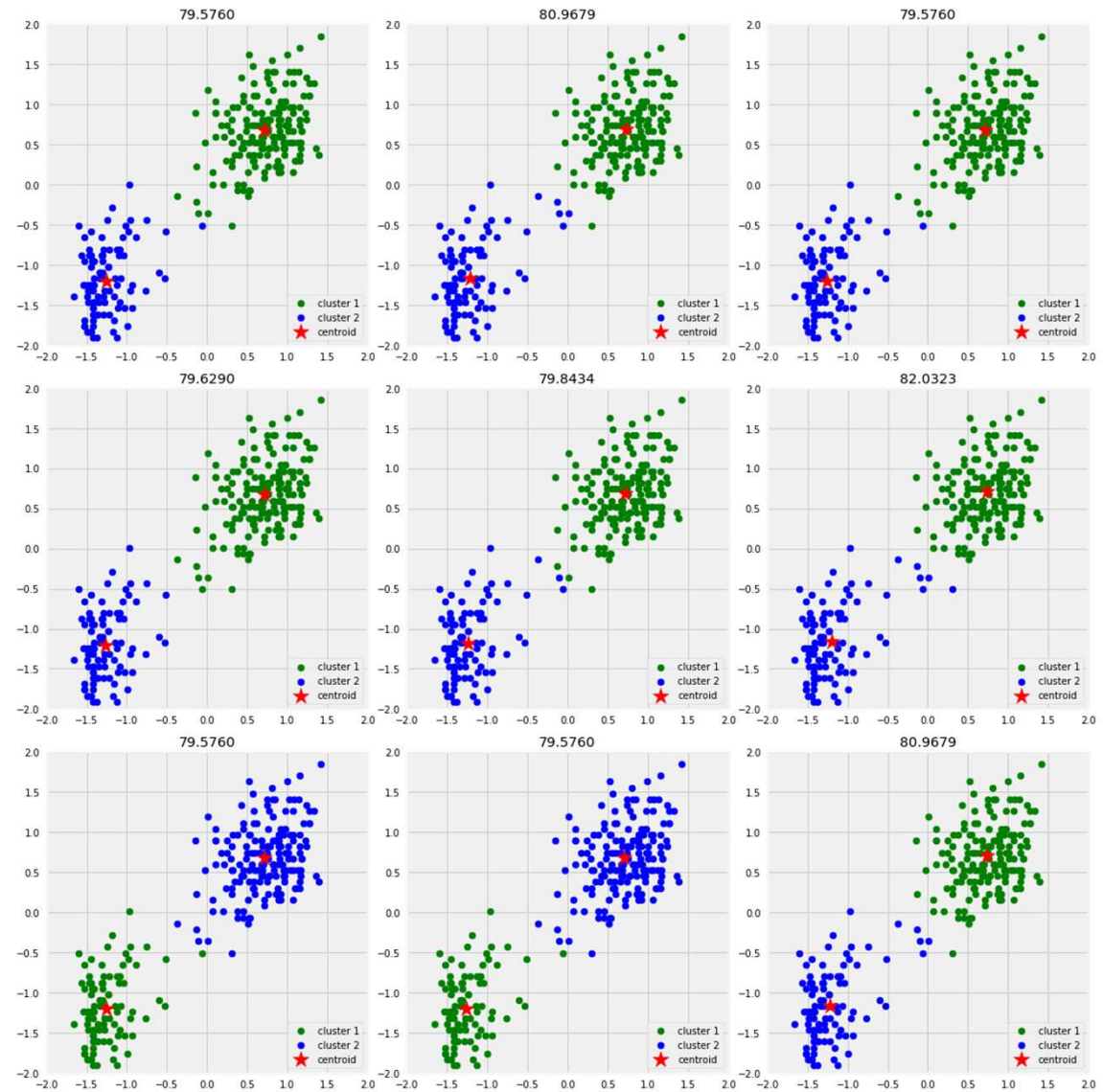
After K-Means

K-Means



Clustering

Centroid

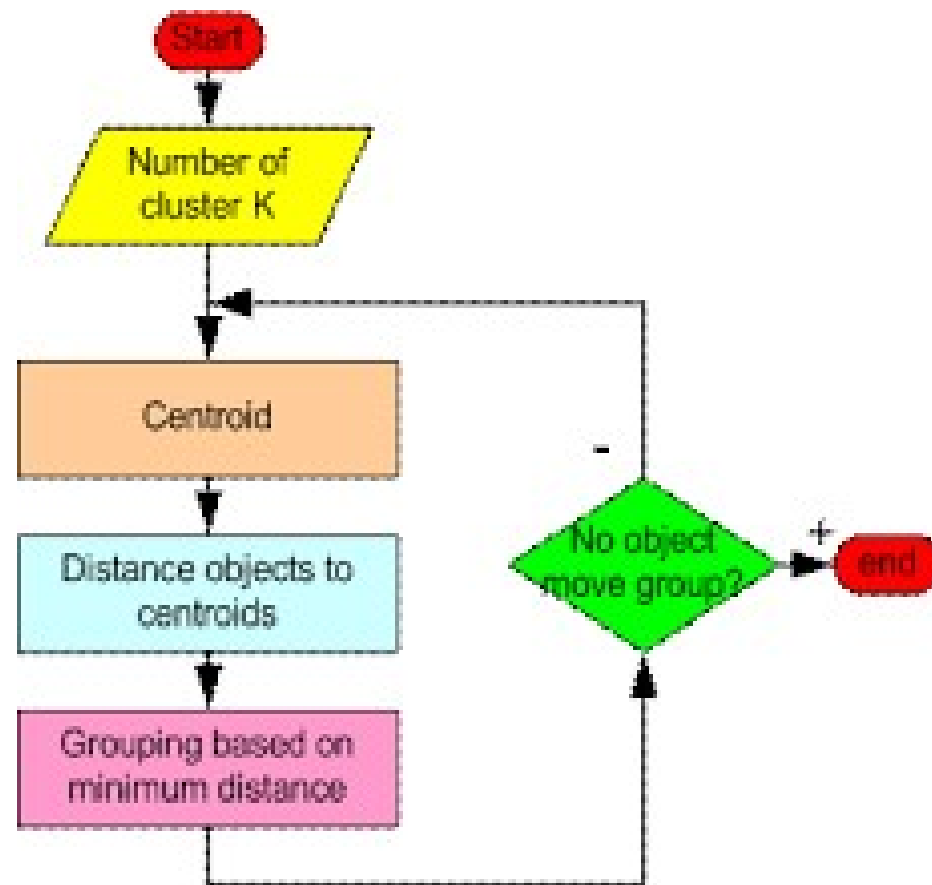


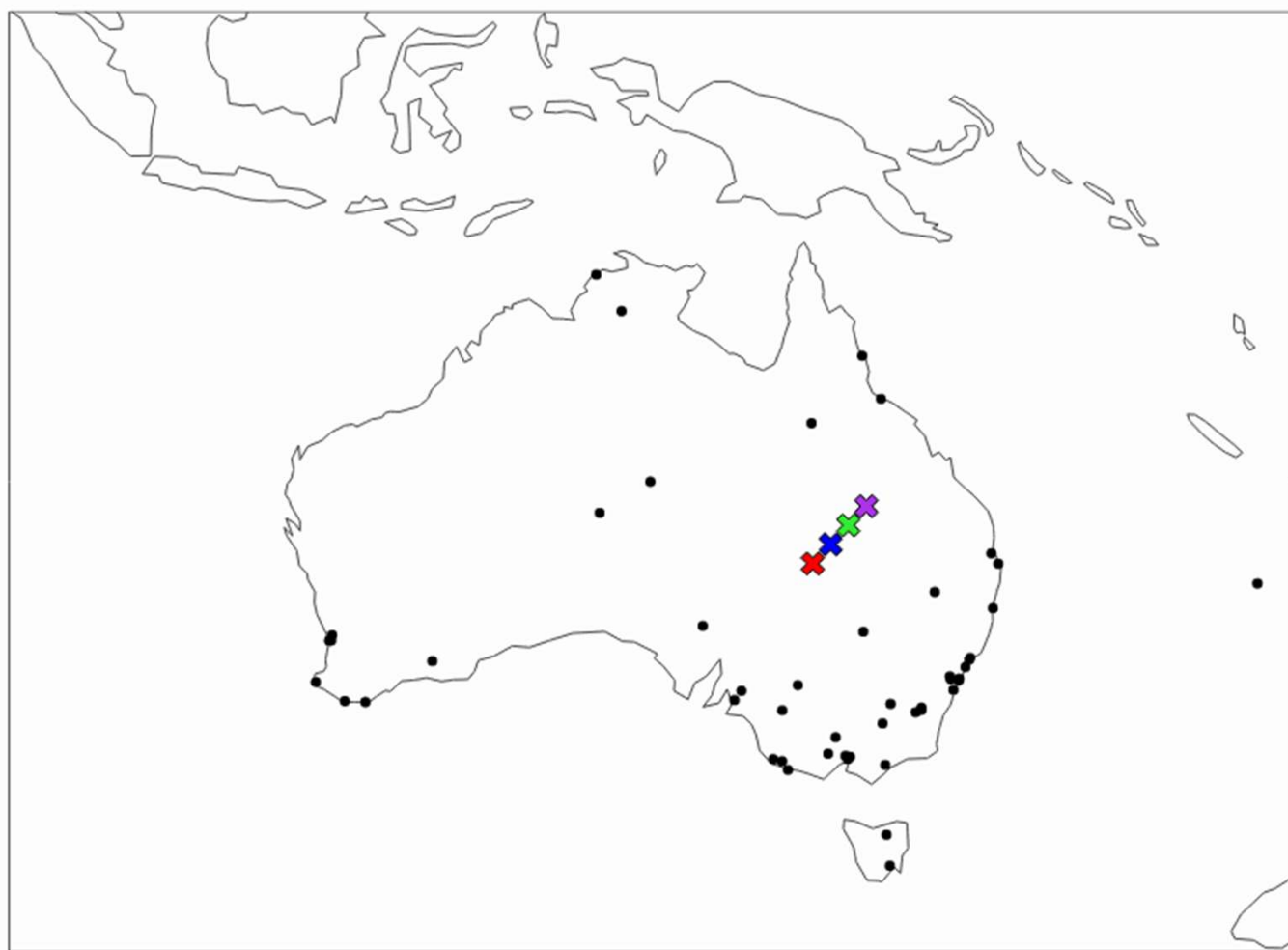
K-means clustering

Algorithm 1 *k*-means algorithm

- 1: Specify the number k of clusters to assign.
 - 2: Randomly initialize k centroids.
 - 3: **repeat**
 - 4: **expectation:** Assign each point to its closest centroid.
 - 5: **maximization:** Compute the new centroid (mean) of each cluster.
 - 6: **until** The centroid positions do not change.
-

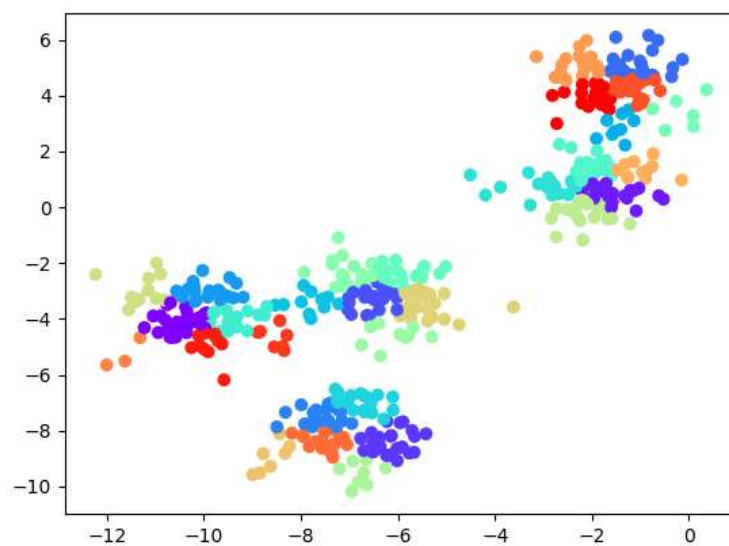
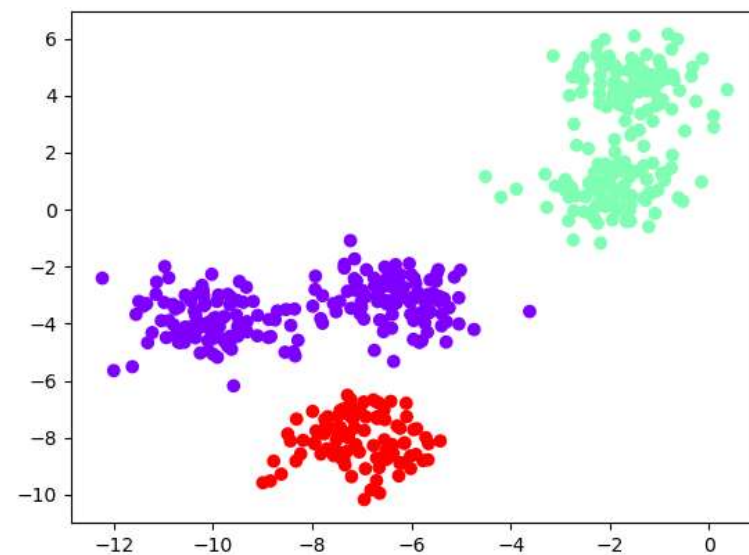
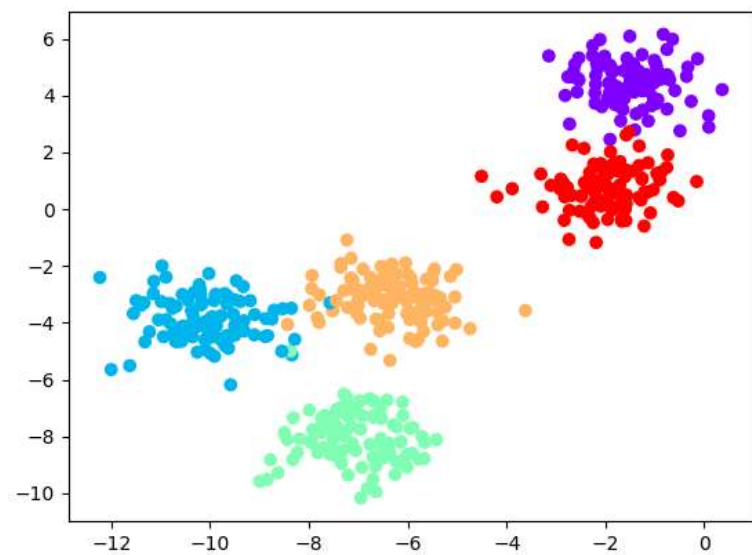
K-means clustering





How good is clustering?

How to choose number of clusters

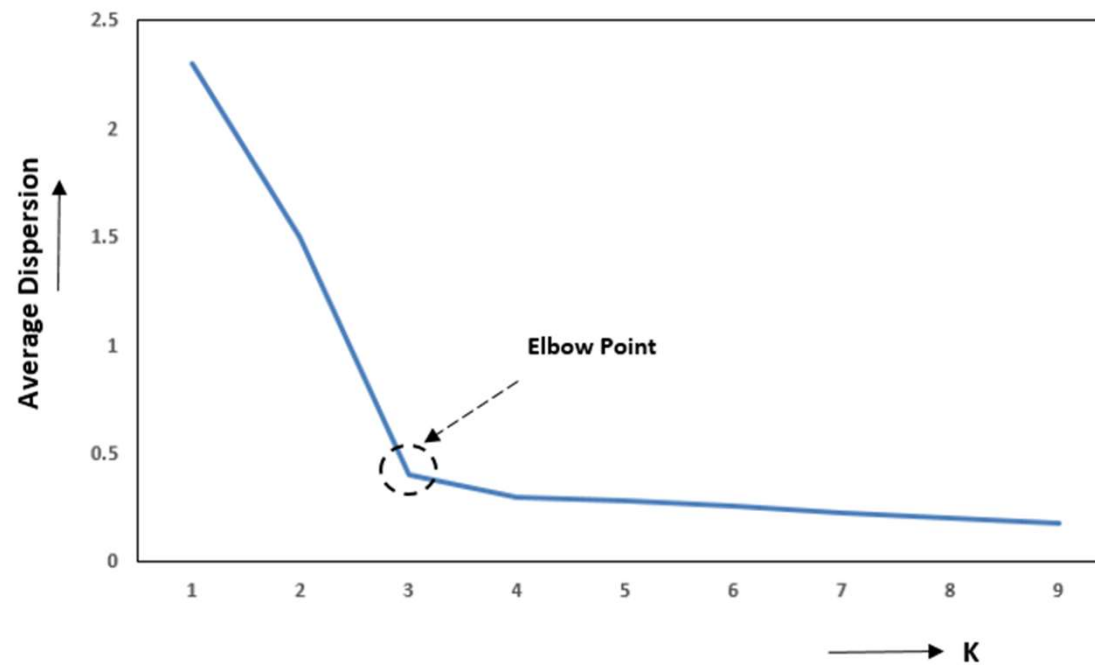


Elbow method

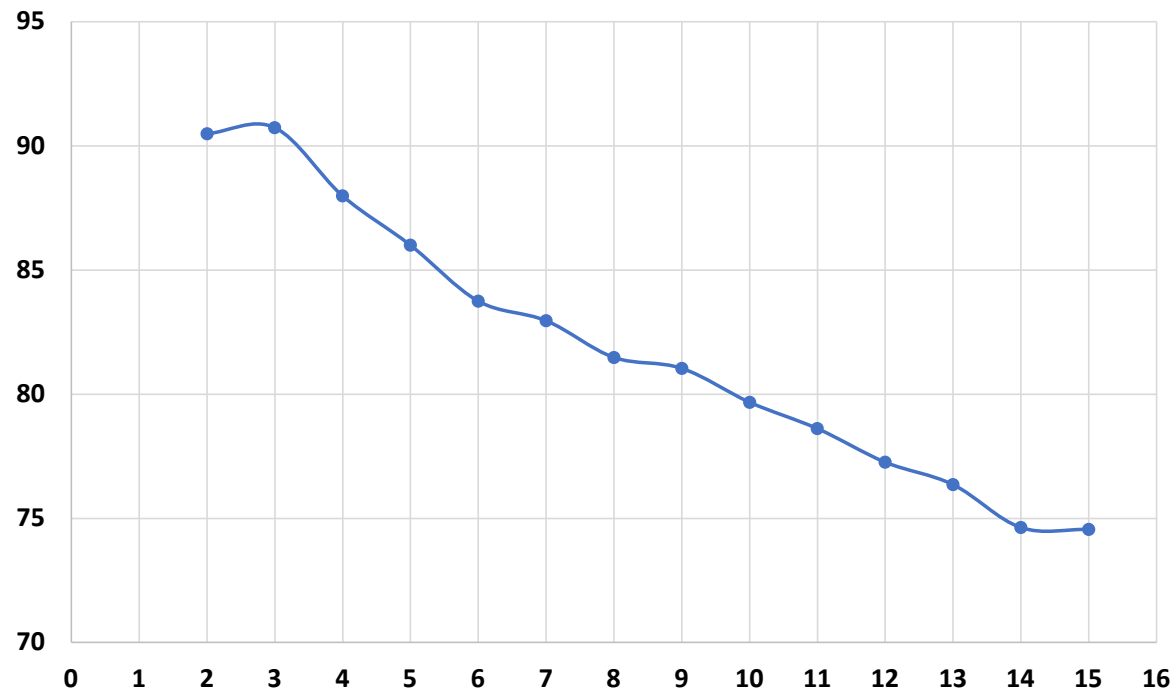
$$WCSS(k) = \sum_{j=1}^k \sum_{\mathbf{x}_i \in \text{cluster } j} \|\mathbf{x}_i - \bar{\mathbf{x}}_j\|^2,$$

where $\bar{\mathbf{x}}_j$ is the sample mean in cluster j

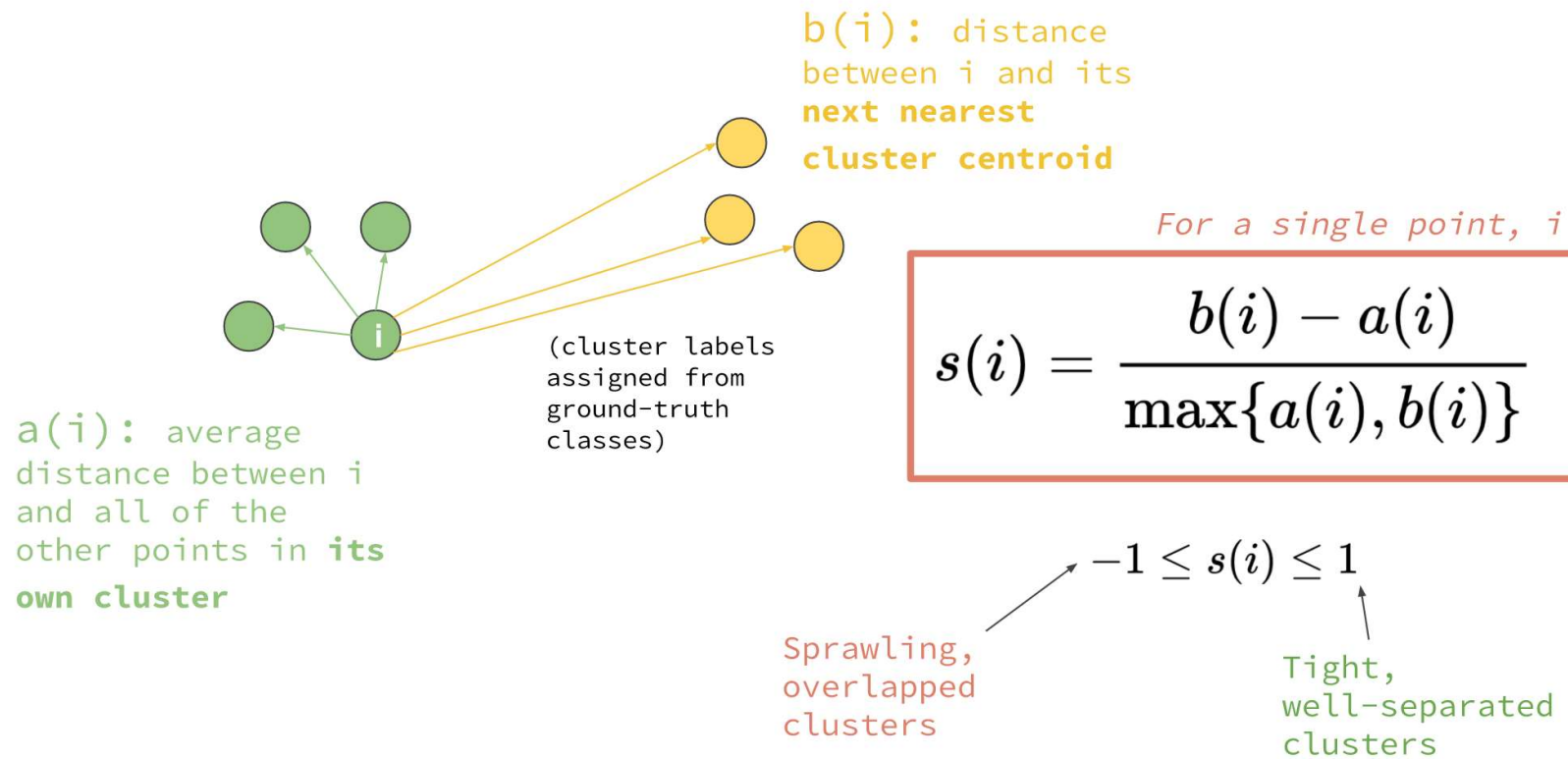
Elbow Method for selection of optimal “K” clusters



Elbow method – possible problems

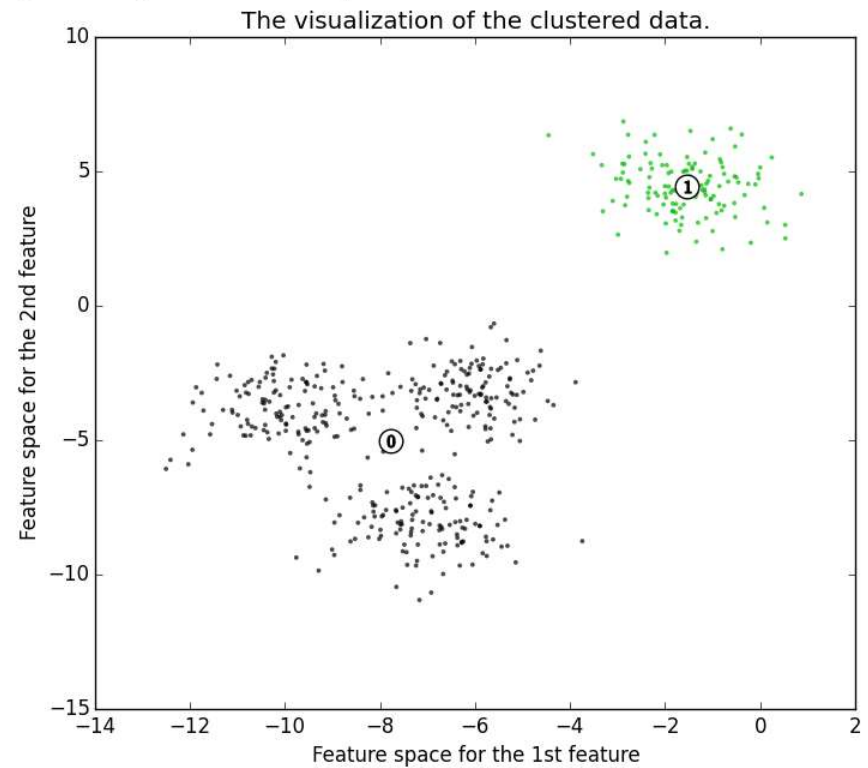
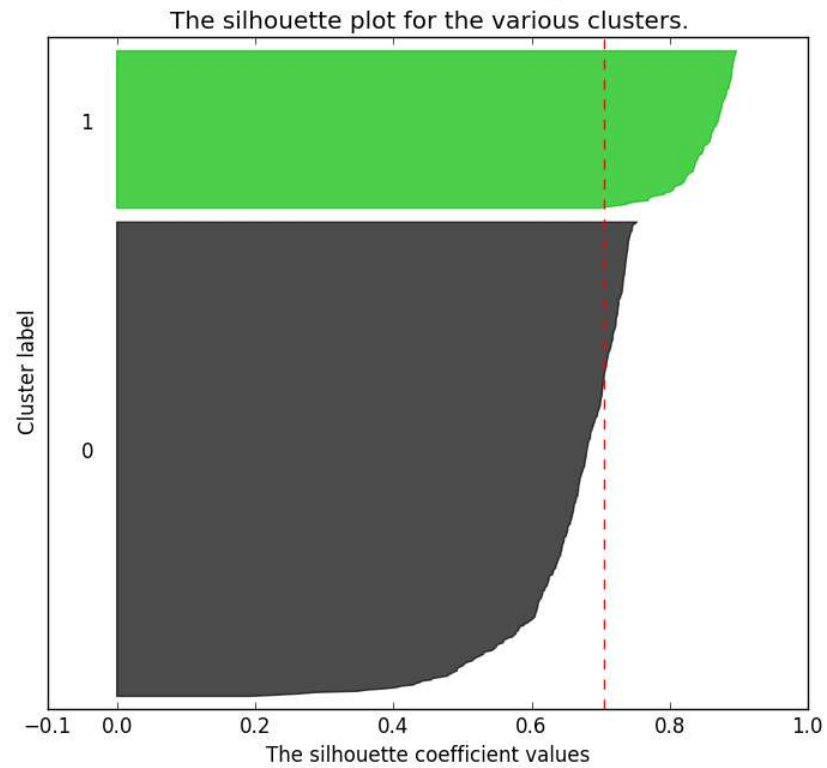


Silhouette



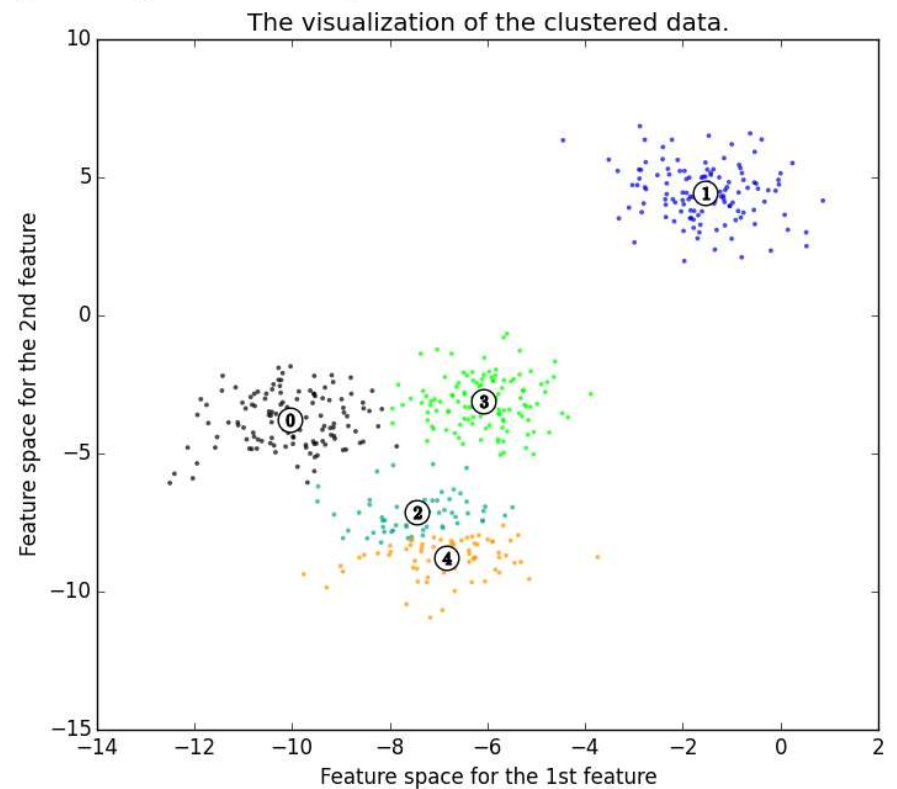
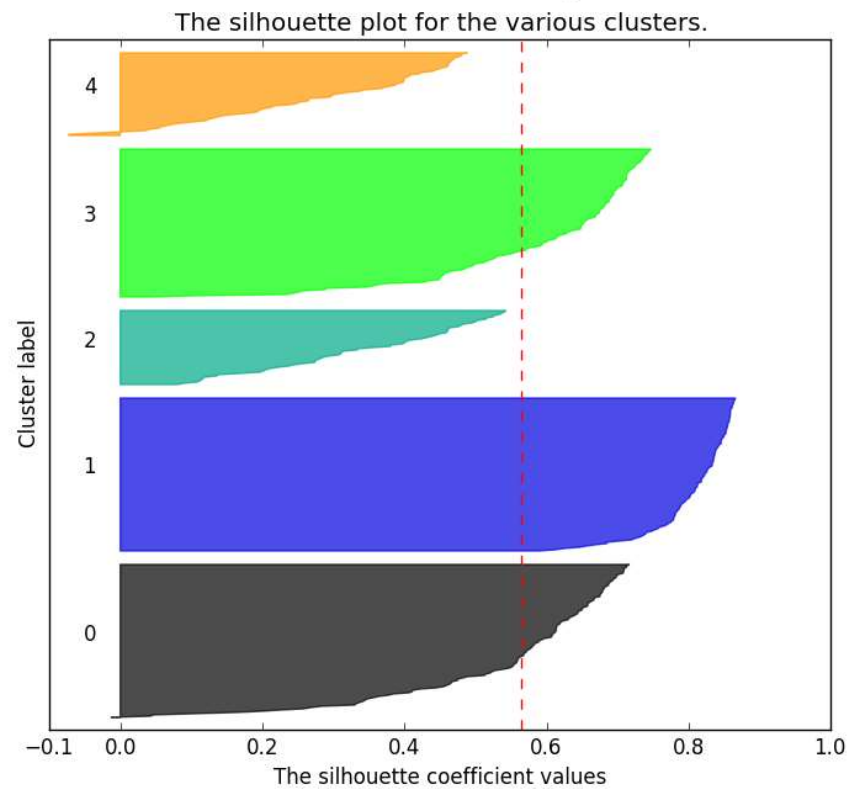
Silhouette

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

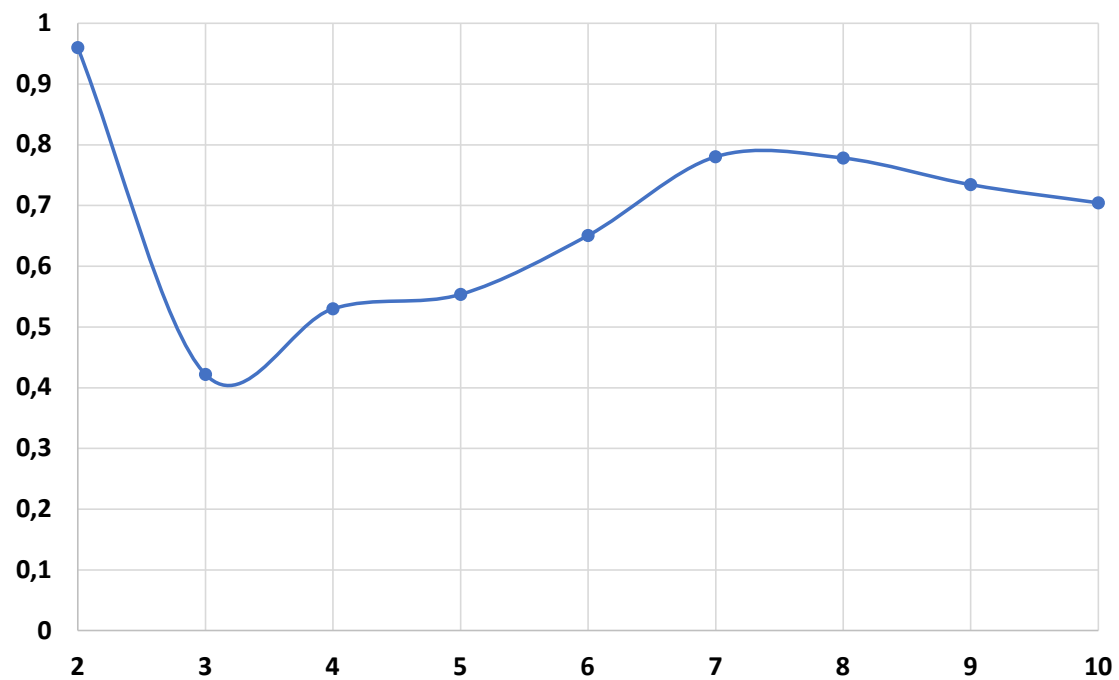


Silhouette

Silhouette analysis for KMeans clustering on sample data with $n_clusters = 5$



Silhoutte



Other methods...

Współczynnik Calińskiego-Harabasza

$$\frac{SS_B}{SS_W} \cdot \frac{N - k}{k - 1}$$

K – liczba klastrow,

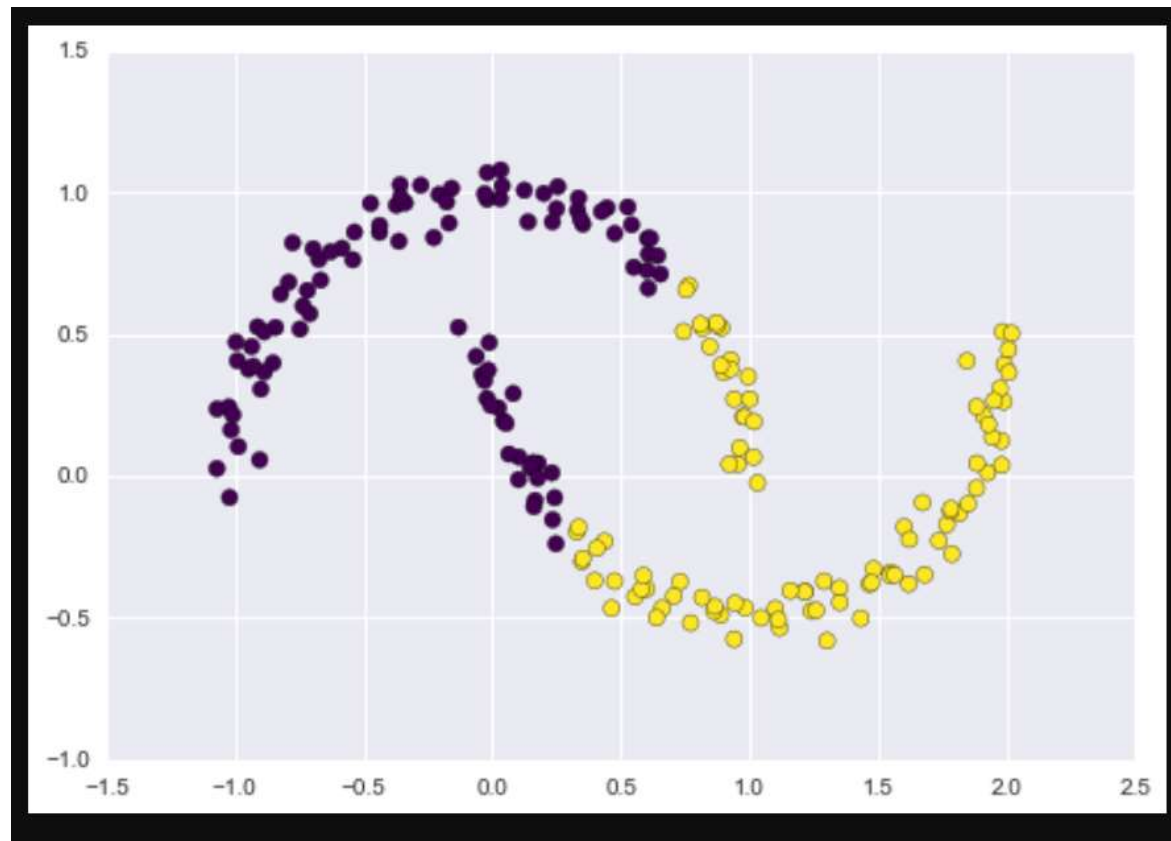
N – liczba obserwacji,

SS_W – wariancja wewnątrz-klastrowa

SS_B – wariancja między-klastrowa

- Im wyższa wartość tym lepiej

K-means clustering - problems



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