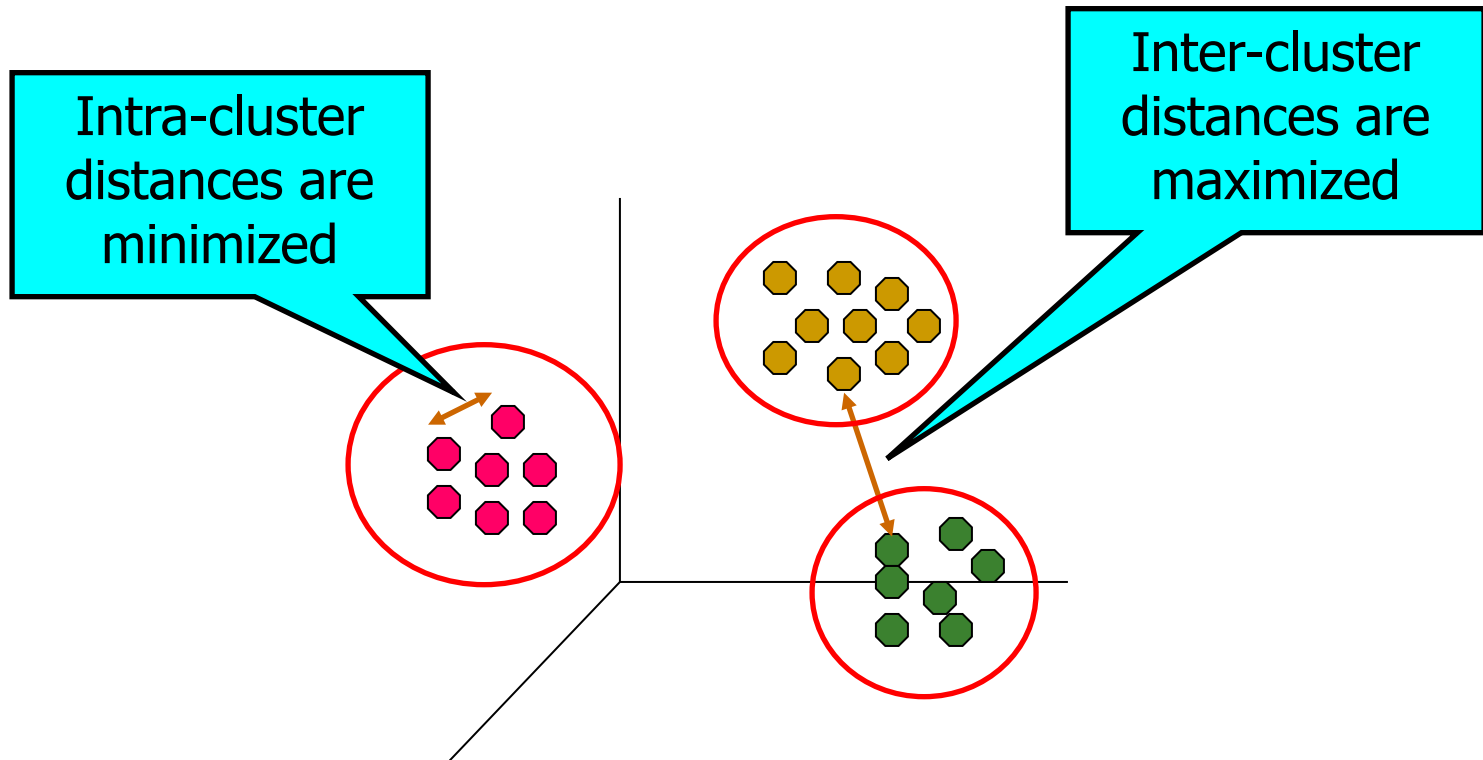


# Lecture 4


## Machine Learning: Unsupervised Learning

Phayung Meesad, Ph.D.  
King Mongkut's University of Technology  
North Bangkok (KMUTNB)  
Bangkok Thailand

- Unsupervised Learning is a machine learning algorithm that learns data without targets.
- Objects in the same cluster are more similar to each other than to those in other clusters.



# Clustering Process

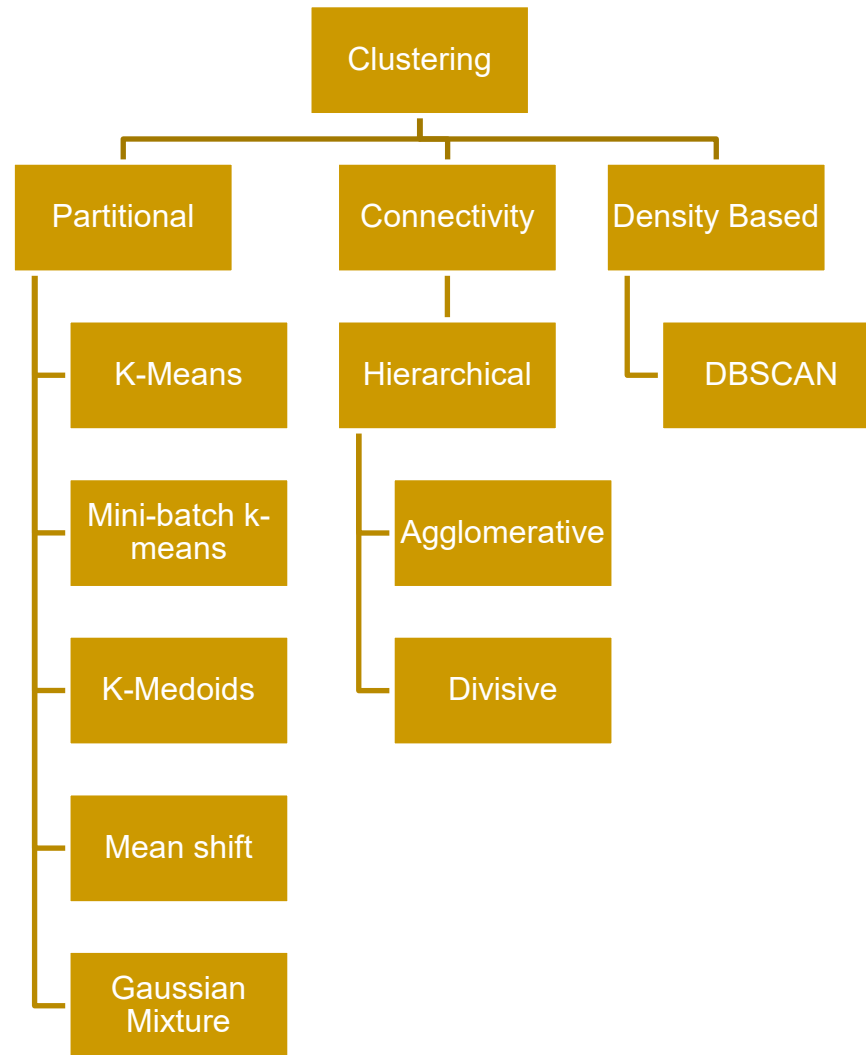
A horizontal flowchart illustrating the clustering process. It consists of three sequential steps, each represented by a yellow chevron pointing to the right, which is itself inside a white rounded rectangle with a yellow border. The steps are: 'Data Preprocessing', 'Clustering Algorithm', and 'Resulted Clusters'.

```
graph LR; A[Data Preprocessing] --> B[Clustering Algorithm]; B --> C[Resulted Clusters];
```

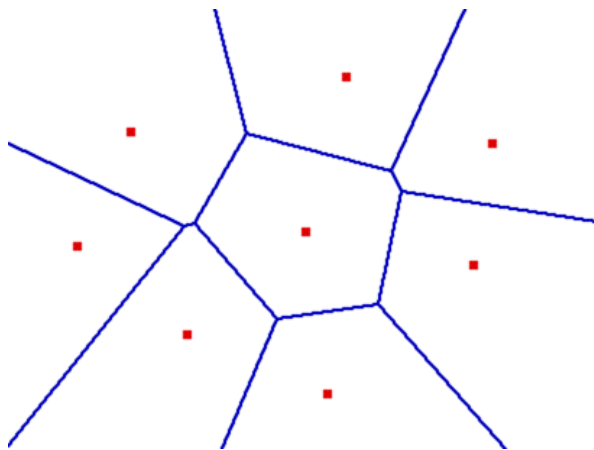
Data  
Preprocessing

Clustering  
Algorithm

Resulted  
Clusters



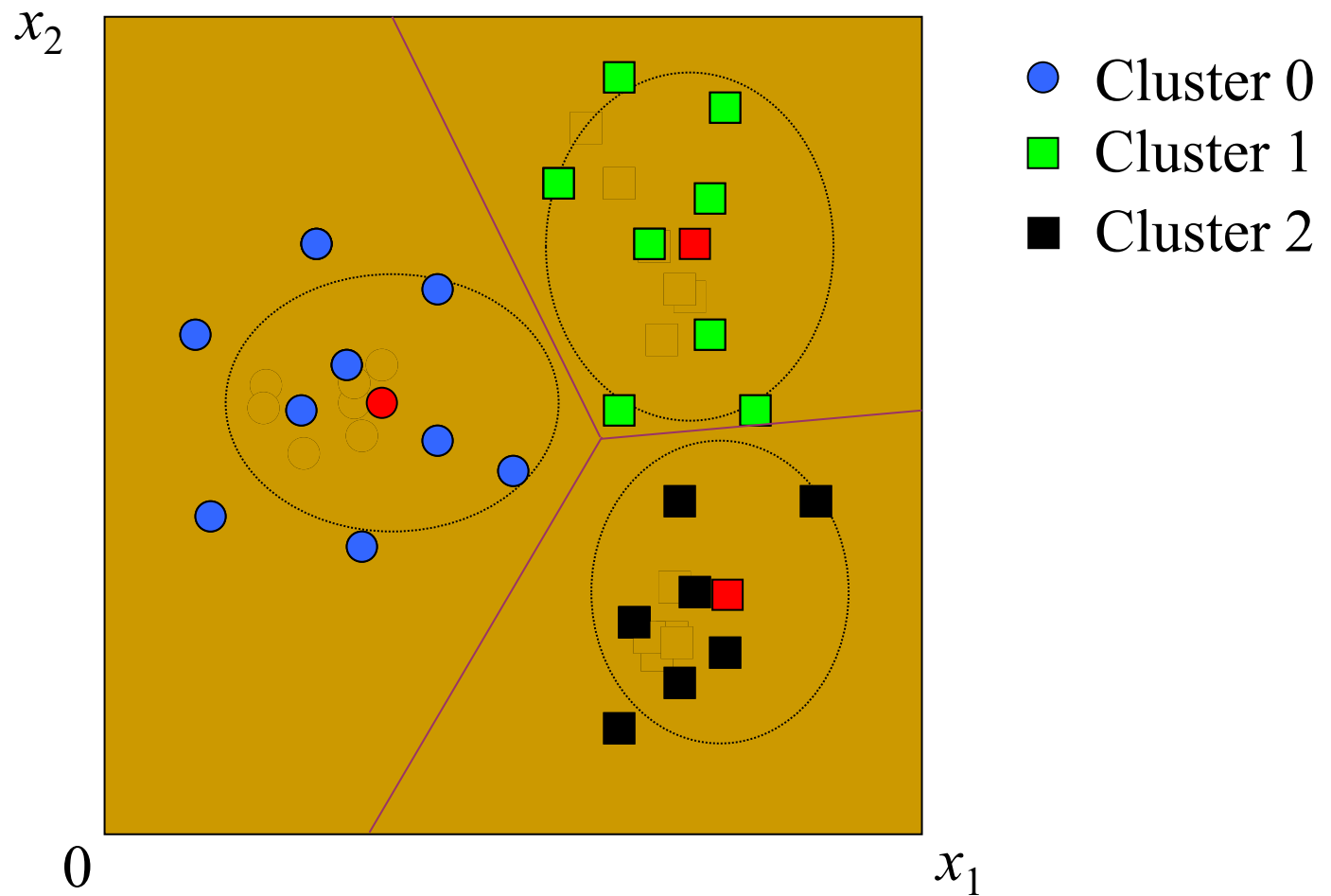
- Centroids are centers or means of the clusters. Each data point is assigned to the closest centroid.

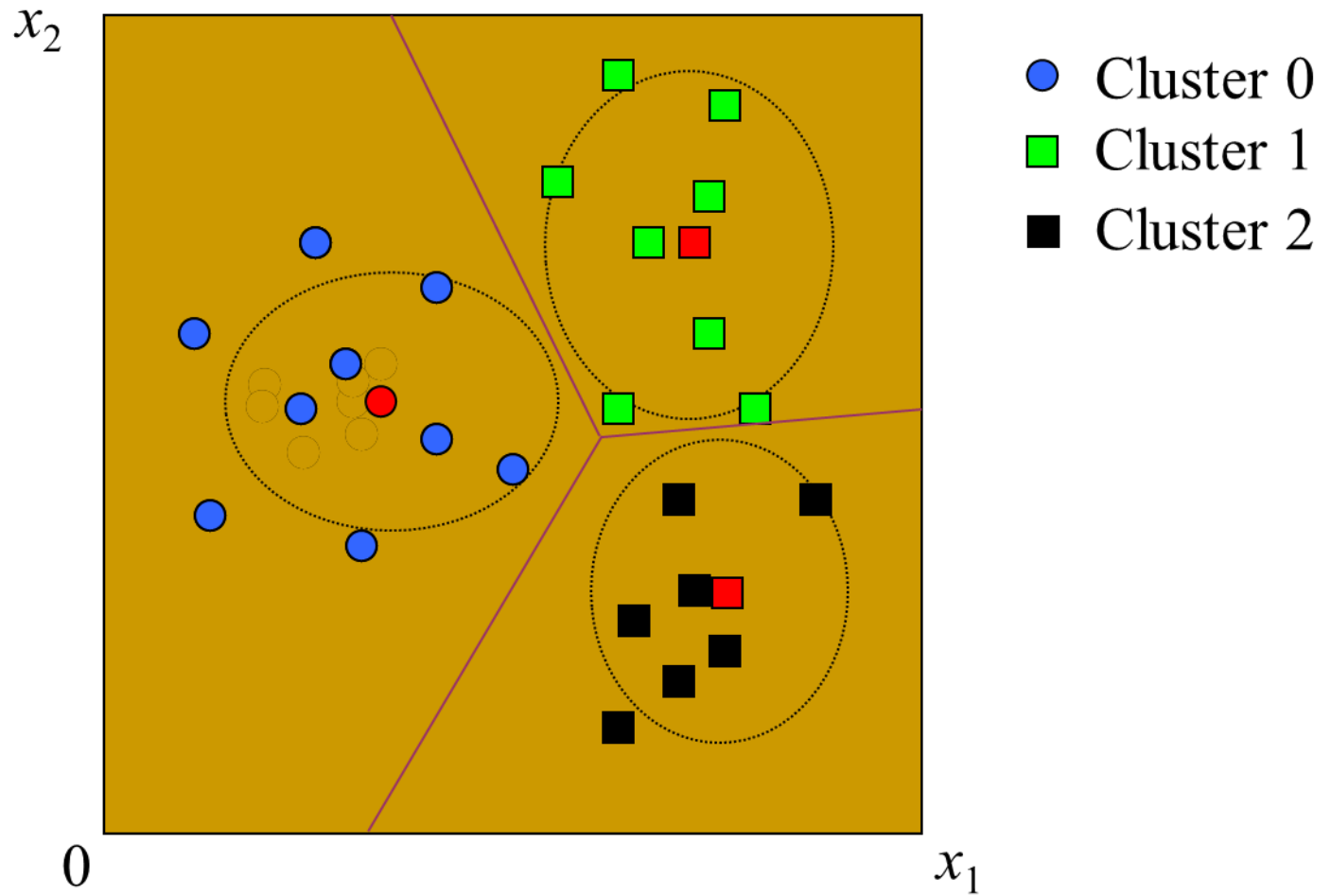


Partitional method

$$\text{similarity} = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

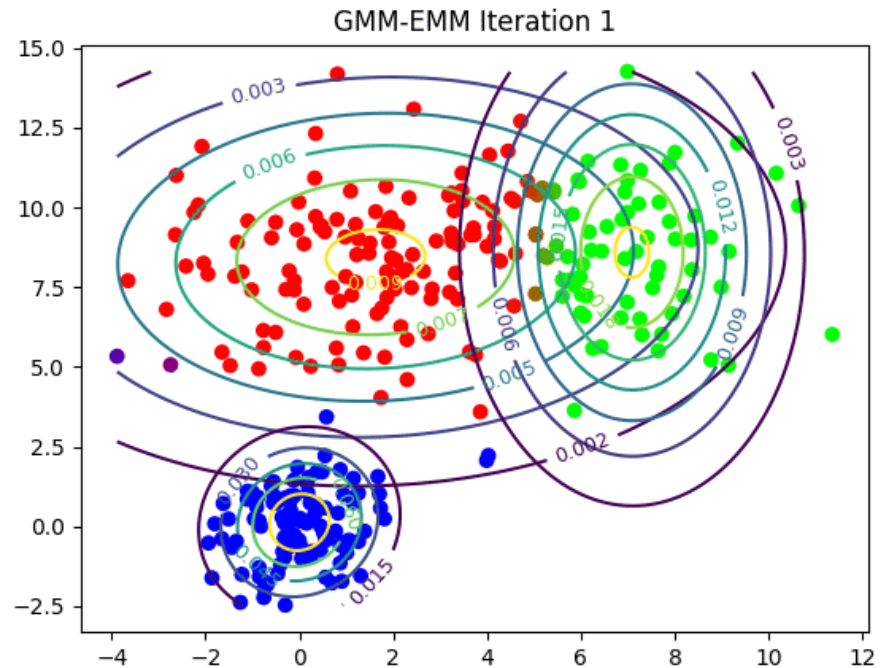
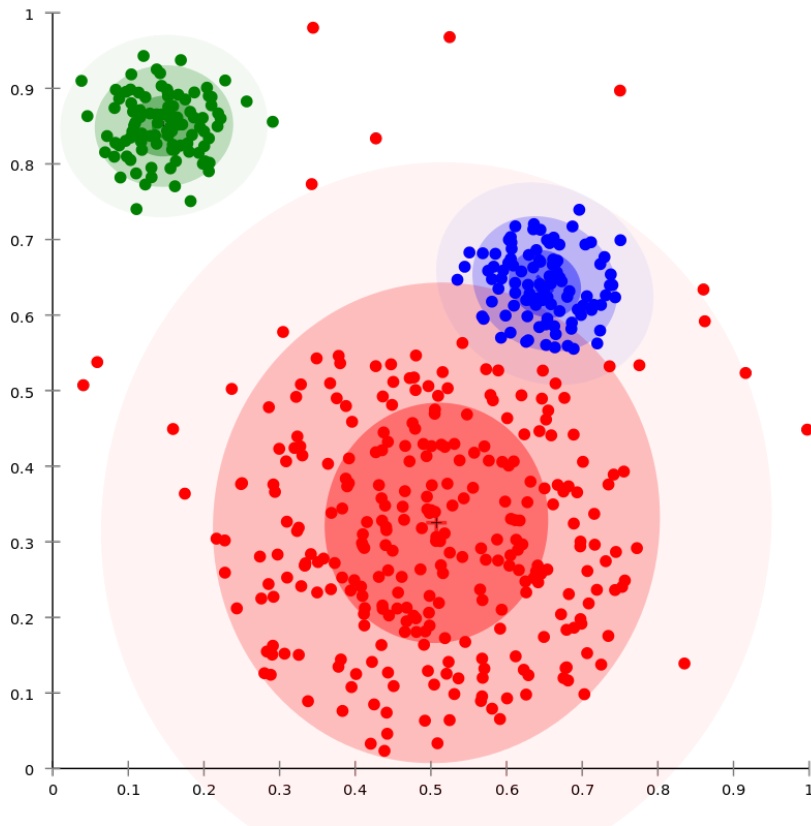
$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}.$$



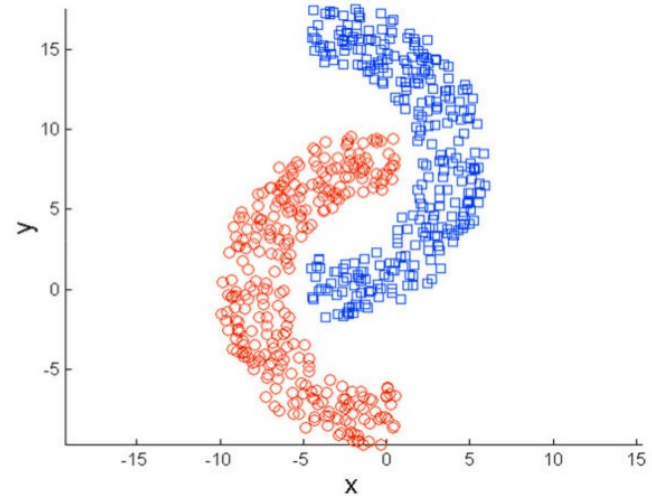
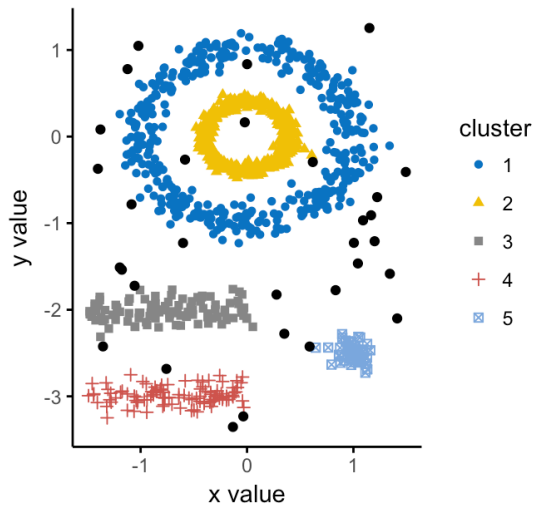
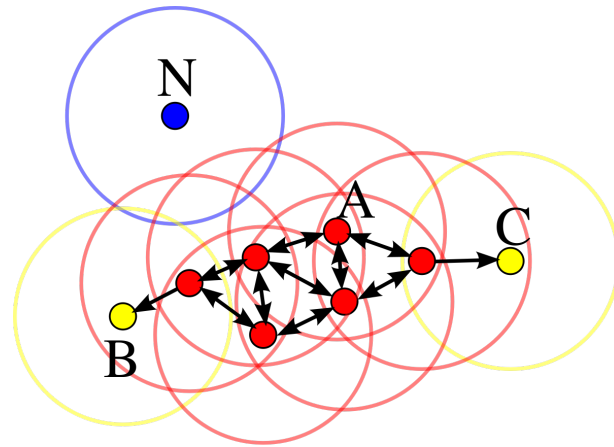




- Gaussian Mixture Models
- Expectation Maximization Algorithm



- The data points closer in data space exhibit more similarity to each other than the data points lying farther away.
- Two approaches:
  - ❑ 1) starting with classifying all data points into separate clusters & then aggregating them as the distance decreases;
  - ❑ 2) all data points are classified as a single cluster and then partitioned as the distance increases.



- K-Means
- Affinity propagation
- Mean-shift
- Spectral clustering
- Gaussian Mixture Models (GMM)
- Hierarchical clustering
- BIRCH
- DBSCAN

- This algorithm requires the number of clusters to be specified.
- It scales well to large number of samples and has been used across a large range of application areas in many different fields.
- The KMeans algorithm clusters data by trying to separate samples in  $n$  groups of equal variance, minimizing a criterion known as the inertia or within-cluster sum-of-squares.

- The K-means algorithm divides a set of  $N$  samples  $X$  into  $K$  disjoint clusters  $C$ , each described by the mean  $\mu_j$  of the samples in the cluster.
- The means are commonly called the cluster “centroids”.
- The K-means algorithm aims to choose centroids that minimise the inertia, or within-cluster sum of squared criterion:

- $$\sum_{i=0}^n \min_{\mu_i \in C} \left( \|x_j - \mu_i\|^2 \right)$$

- Partitioning clustering approach
  - Each cluster is associated with a centroid (center point)
  - Each point is assigned to the cluster with the closest centroid
  - Number of clusters,  $K$ , must be specified

---
- 1: Select  $K$  points as the initial centroids.
  - 2: **repeat**
  - 3:   Form  $K$  clusters by assigning all points to the closest centroid.
  - 4:   Recompute the centroid of each cluster.
  - 5: **until** The centroids don't change

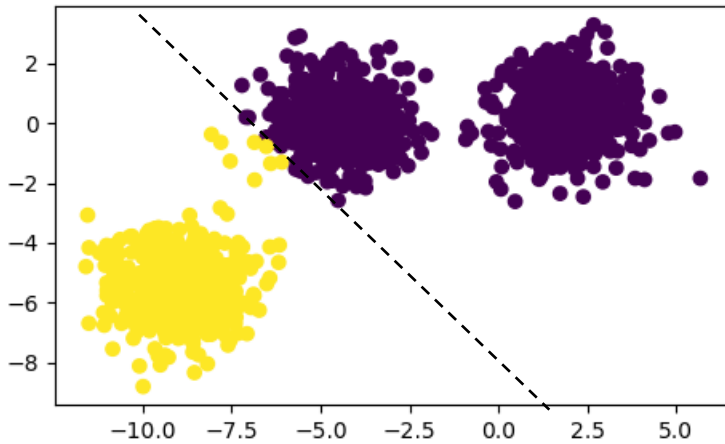
---

- Inertia, or the within-cluster sum of squares criterion, can be recognized as a measure of how internally coherent clusters are.
- Drawbacks:
  - ❑ Inertia makes the assumption that clusters are convex and isotropic. It responds poorly to elongated clusters, or manifolds with irregular shapes.
  - ❑ Inertia is not a normalized metric
  - ❑ Curse of dimensionality

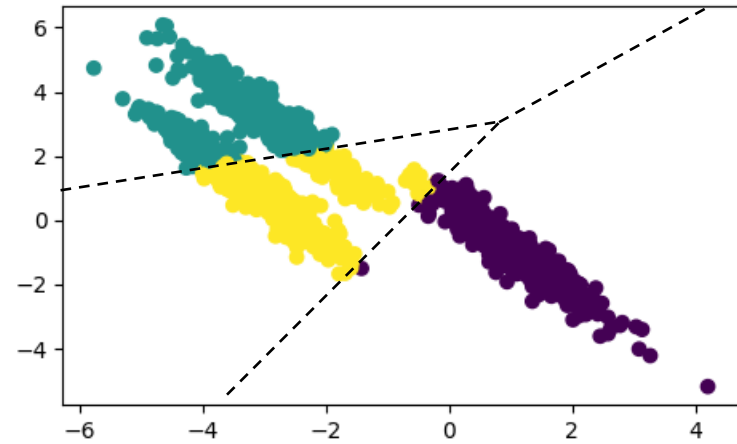


# Drawbacks of KMeans

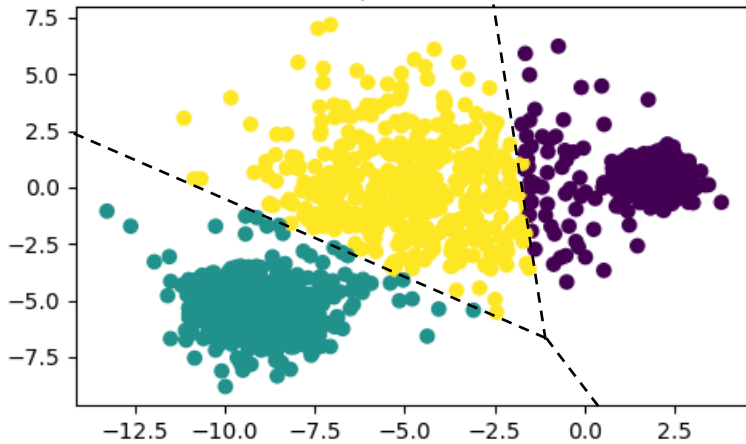
Incorrect Number of Blobs



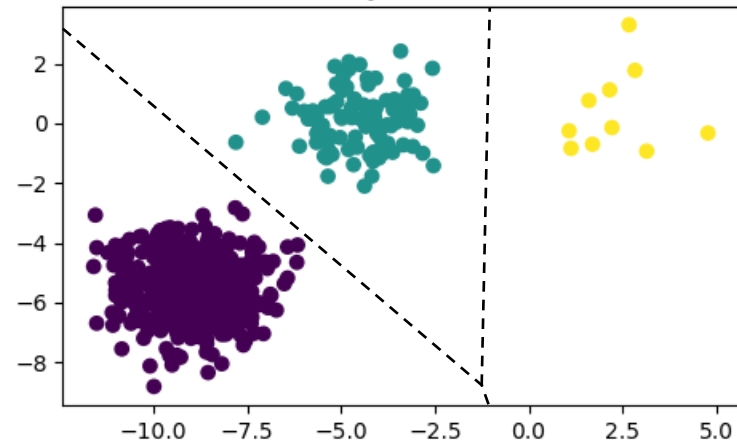
Anisotropically Distributed Blobs

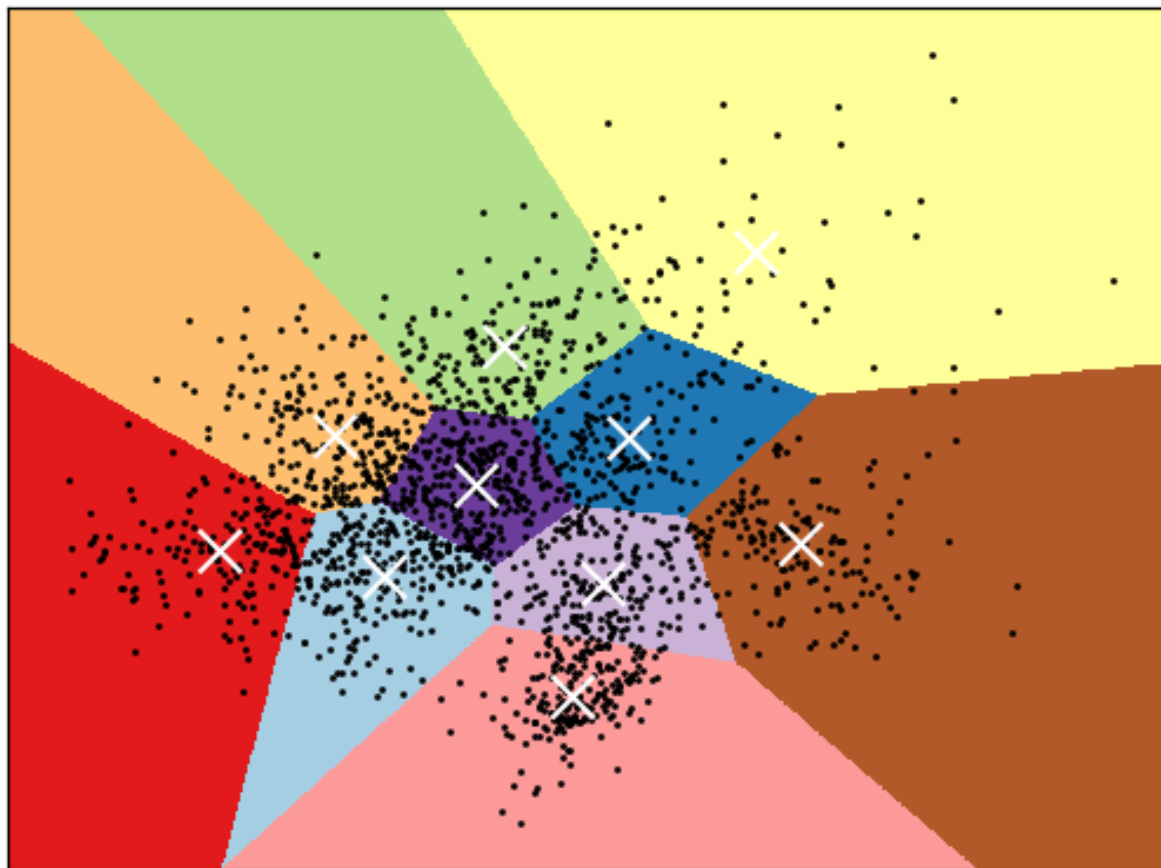


Unequal Variance



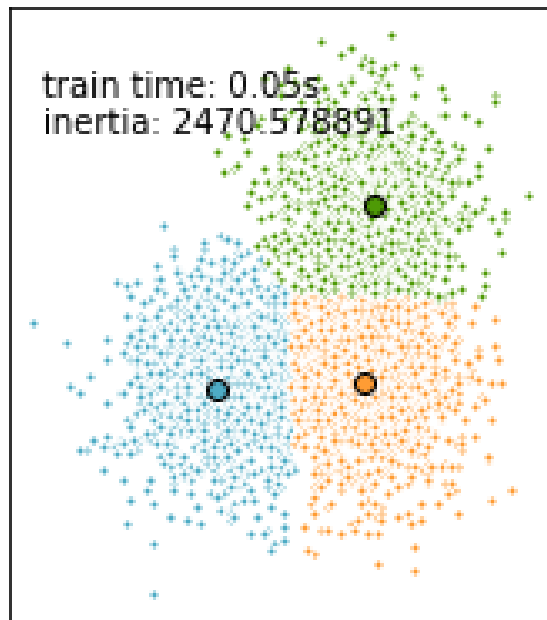
Unevenly Sized Blobs



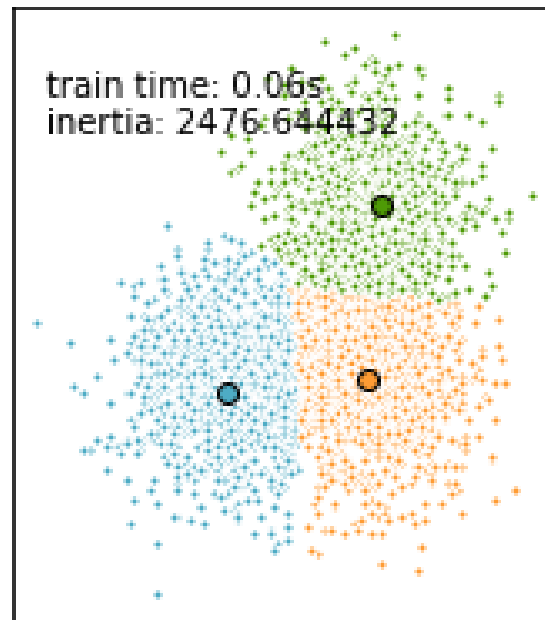


- Mini-batches are subsets of the input data, randomly sampled in each training iteration.
- The algorithm iterates between two steps.
- Step 1:  $b$  samples are drawn randomly from the dataset, to form a mini-batch. These are then assigned to the nearest centroid.
- Step 2: the centroids are updated. For each sample in the mini-batch, the assigned centroid is updated by taking the streaming average of the sample and all previous samples assigned to that centroid.

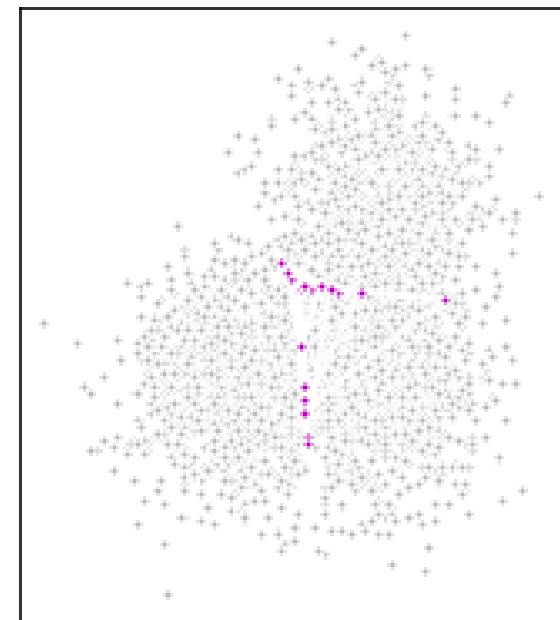
KMeans



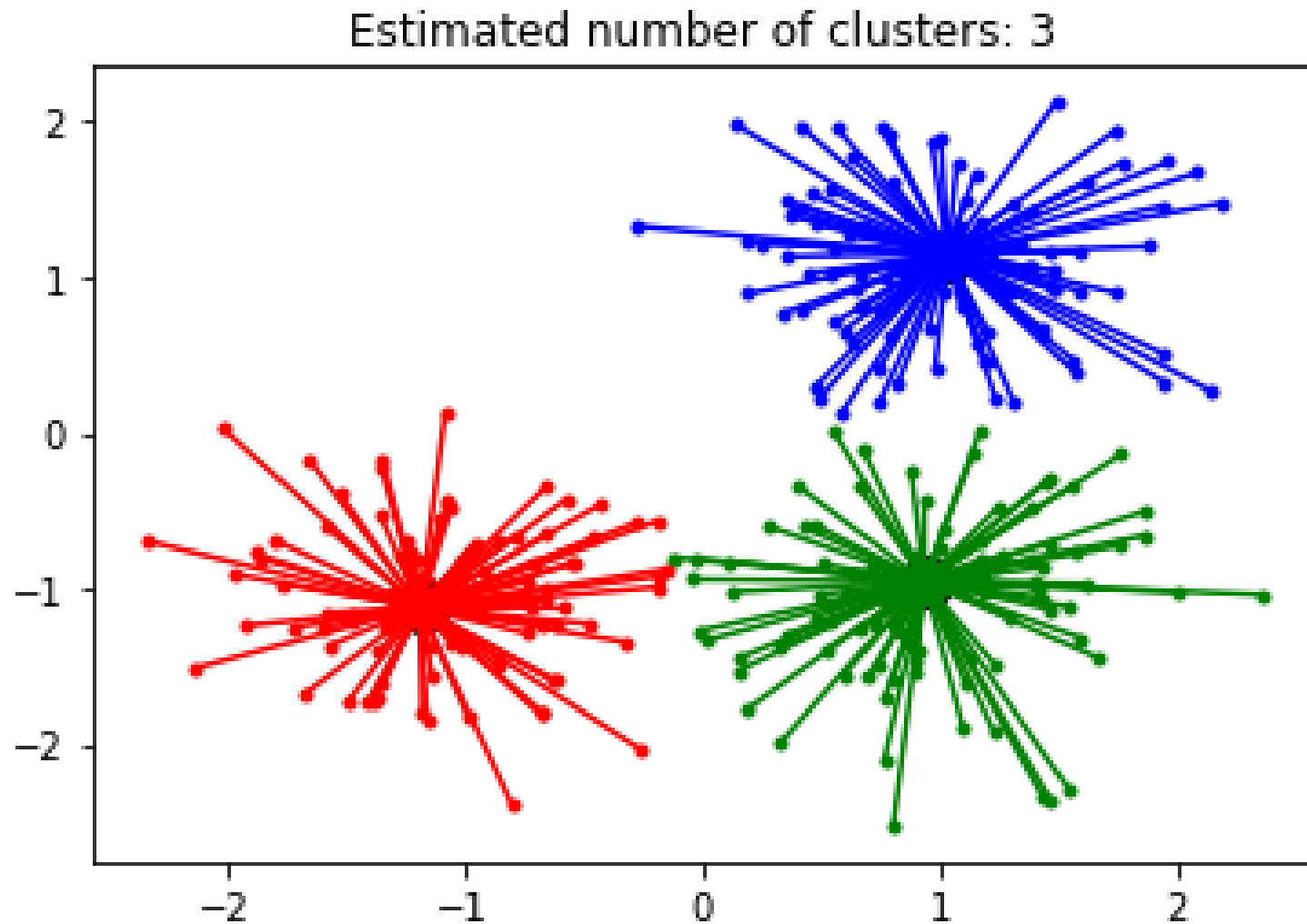
MiniBatchKMeans



Difference



- A dataset is described using a small number of exemplars, which are identified as those most representative of other samples.
- The messages sent between pairs represent the suitability for one sample to be the exemplar of the other, which is updated in response to the values from other pairs.
- Updating happens iteratively until convergence to the final clustering.

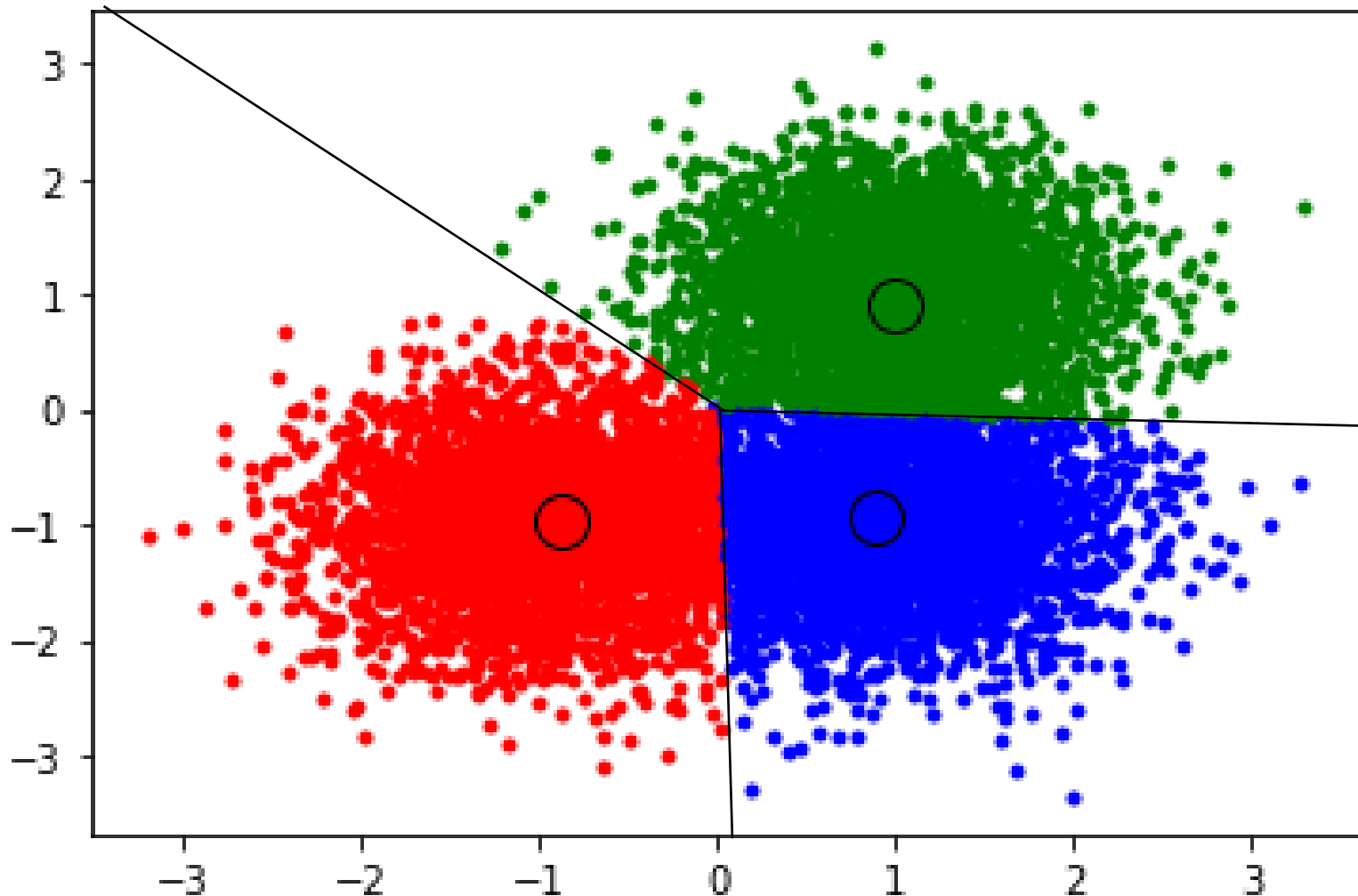


- MeanShift clustering is a centroid based algorithm, updating candidates for centroids to be the mean of the points within a given region.
- Then filtered in a post-processing stage to eliminate near-duplicates to form the final set of centroids.

- $x_i^{\{t+1\}} = x_i^t + m(x_i^t)$

- $$m(x_i) = \frac{\sum_{x_j \in N(x_i)} K(x_j - x_i) x_j}{\sum_{x_j \in N(x_i)} K(x_j - x_i)}$$

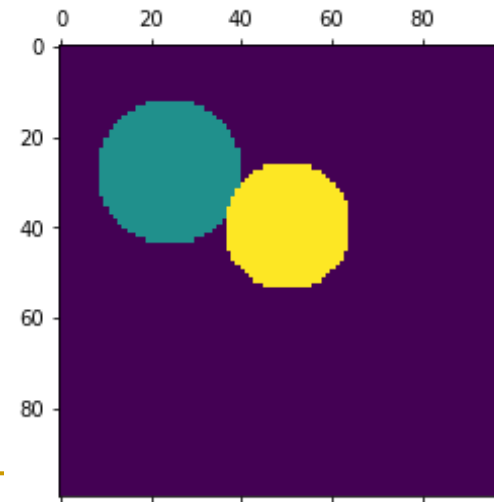
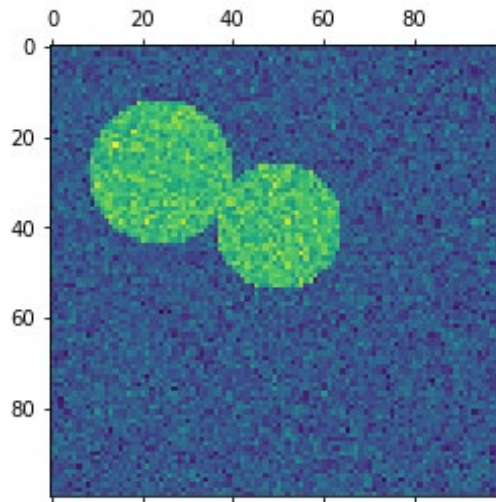
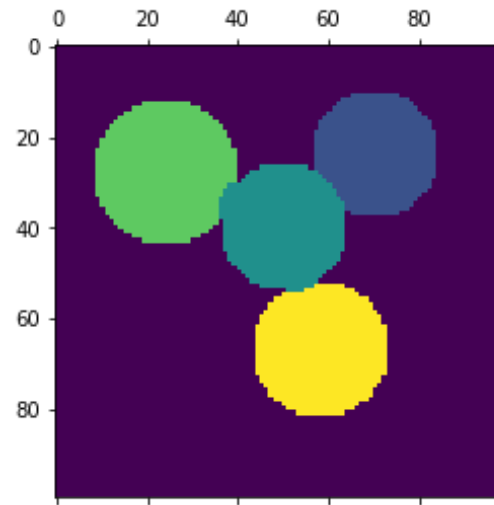
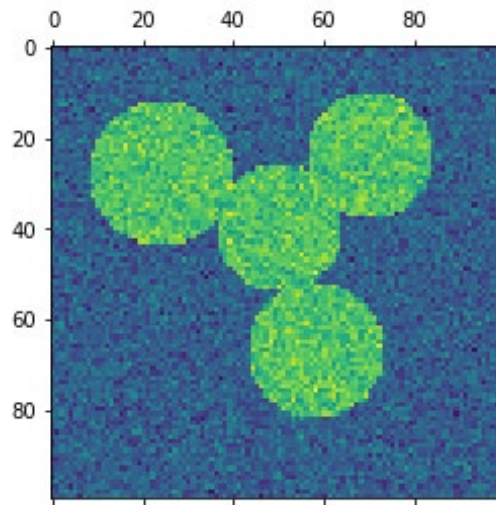
Estimated number of clusters: 3





- SpectralClustering does a low-dimension embedding of the affinity matrix between samples, followed by a KMeans in the low dimensional space.
- SpectralClustering requires the number of clusters to be specified.
- It works well for a small number of clusters.
- This criteria is especially interesting when working on images: graph vertices are pixels, and edges of the similarity graph are a function of the gradient of the image.

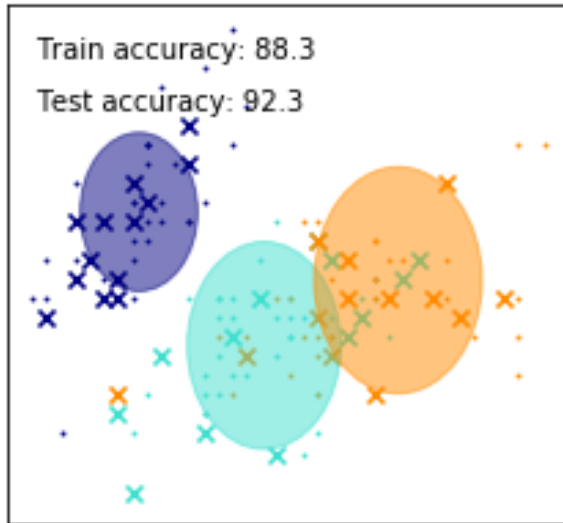
# Spectral clustering



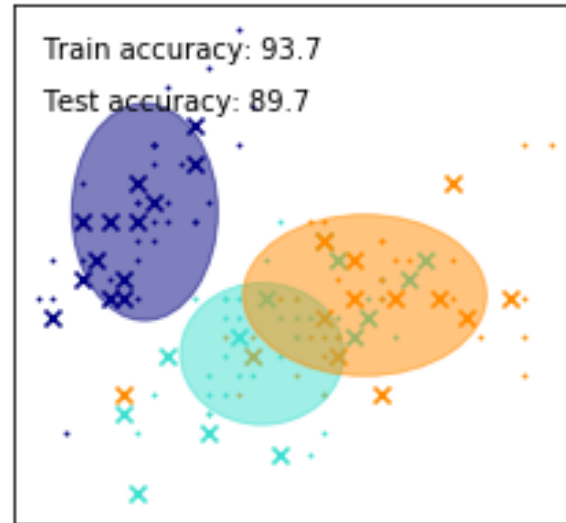
- The GaussianMixture object implements the expectation-maximization (EM) algorithm for fitting mixture-of-Gaussian models.
- It can also draw confidence ellipsoids for multivariate models, and compute the Bayesian Information Criterion to assess the number of clusters in the data.

# Gaussian Mixture Model

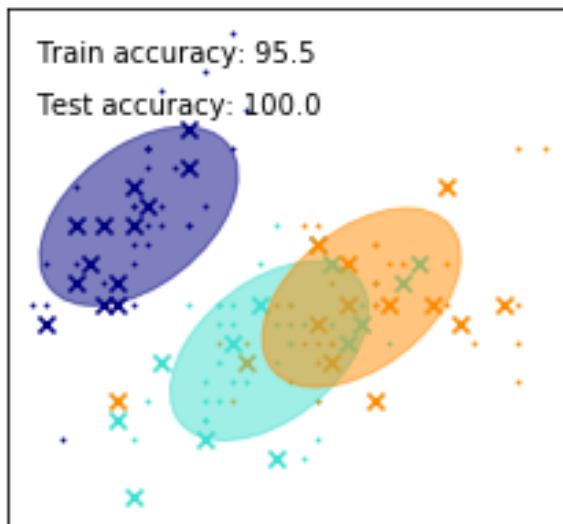
spherical



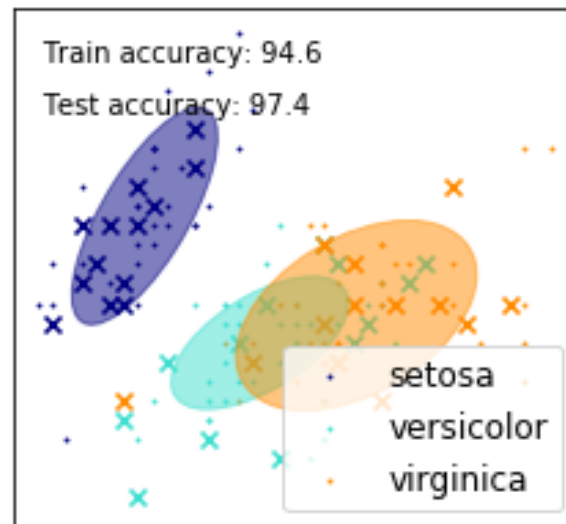
diag



tied



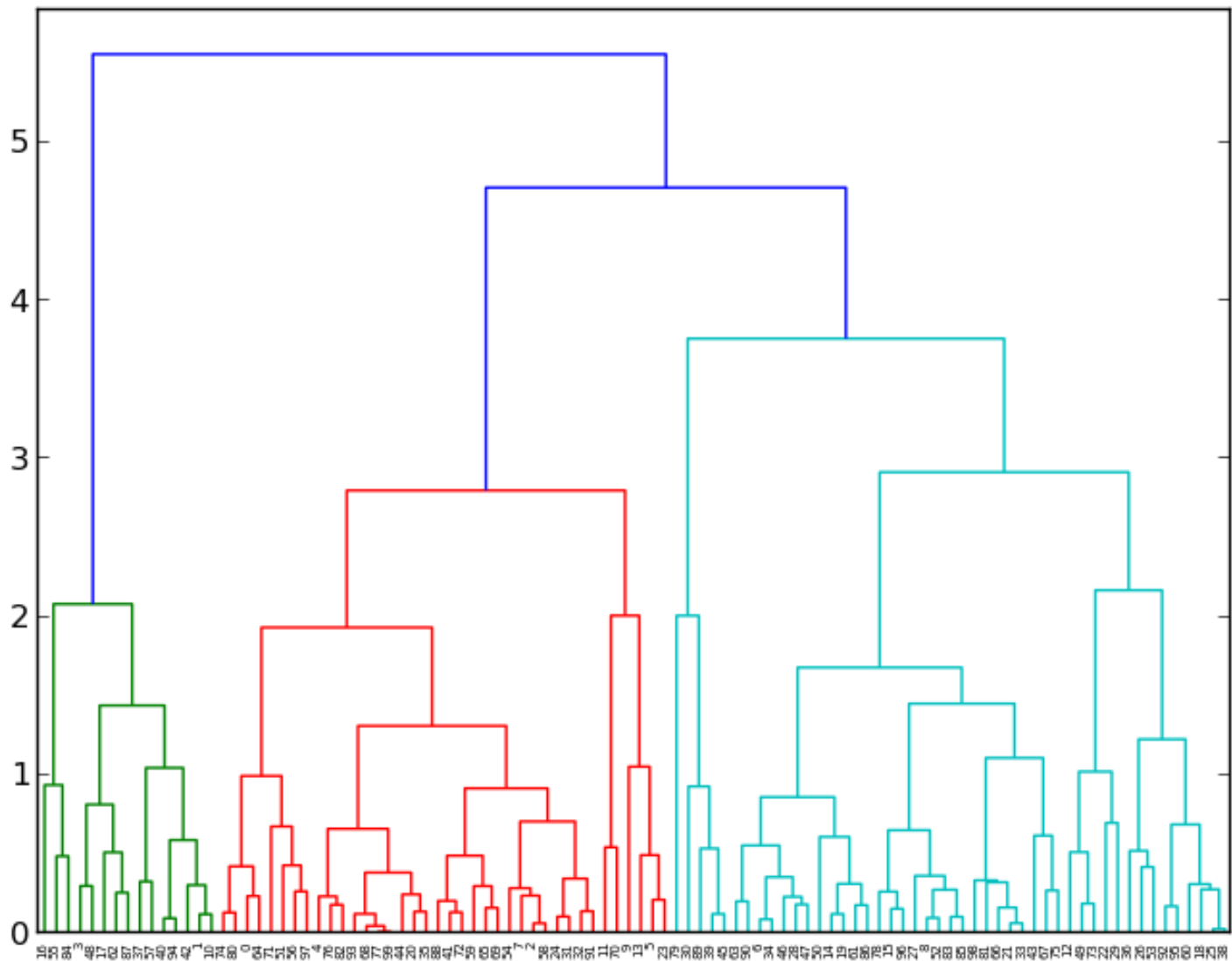
full



- Hierarchical clustering builds nested clusters by merging or splitting them successively.
- This hierarchy of clusters is represented as a tree (or dendrogram).
- The root of the tree is the unique cluster that gathers all the samples, the leaves being the clusters with only one sample.
- The Agglomerative Clustering object performs a hierarchical clustering using a bottom up approach: each observation starts in its own cluster, and clusters are successively merged together.

- Agglomerative Clustering can also scale to large number of samples when it is used jointly with a connectivity matrix, but is computationally expensive when no connectivity constraints are added between samples.
- It considers at each step all the possible merges.

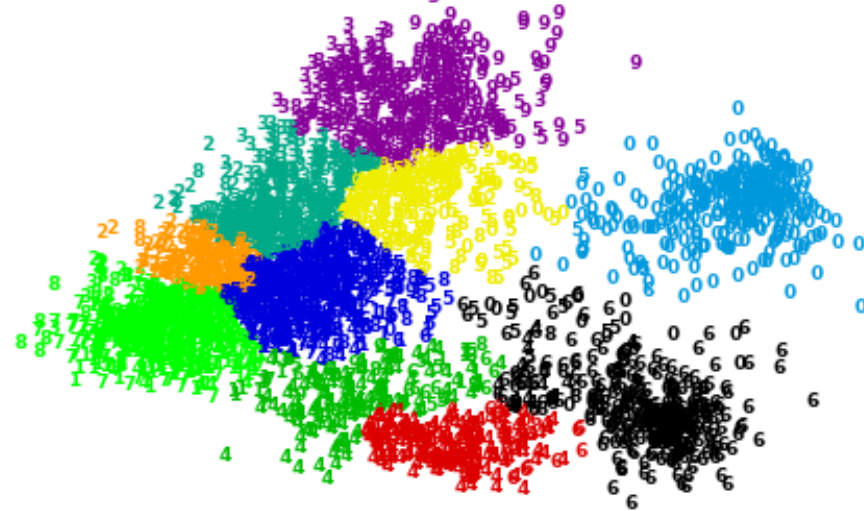
- The linkage criteria determines the metric used for the merge strategy.
- Minimum or single linkage
- Maximum or complete linkage minimizes the maximum distance between observations of pairs of clusters.
- Average linkage minimizes the average of the distances between all observations of pairs of clusters.
- Ward hierarchical clustering minimizes the sum of squared differences within all clusters.
- It is a variance-minimizing approach and in this sense is similar to the k-means objective function but tackled with an agglomerative hierarchical approach.



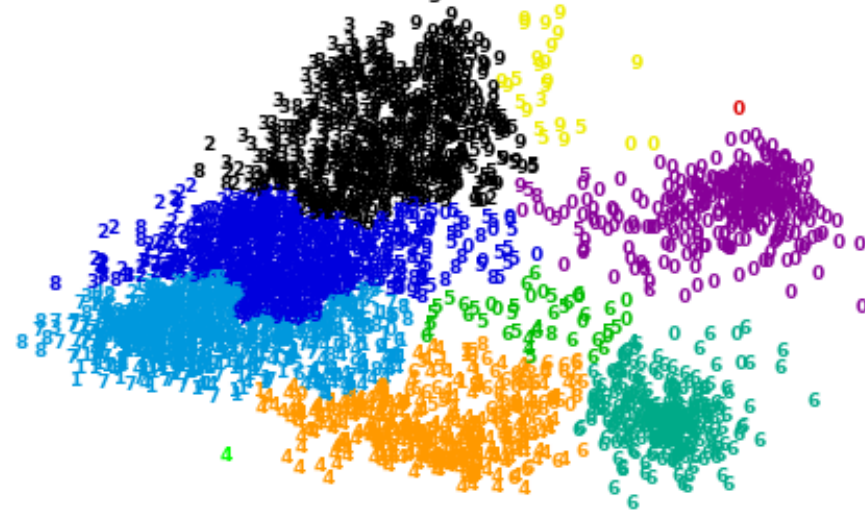


# Hierarchical Clustering Comparison

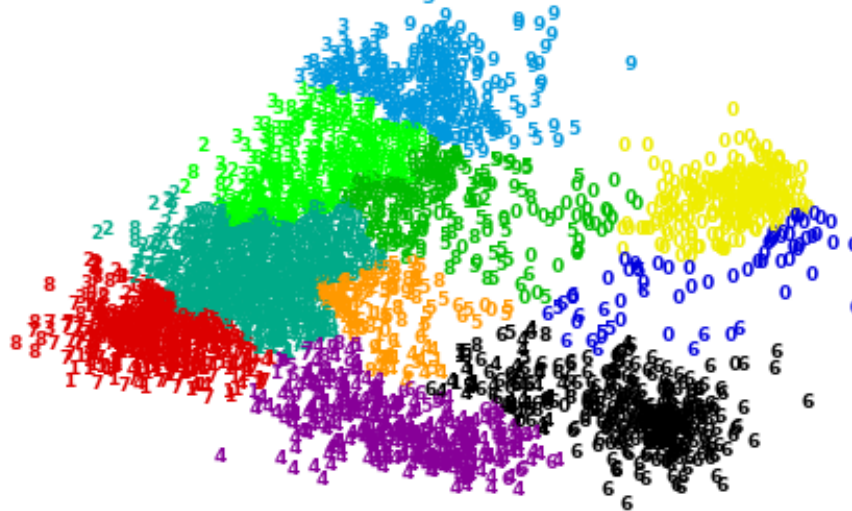
ward linkage



average linkage



complete linkage

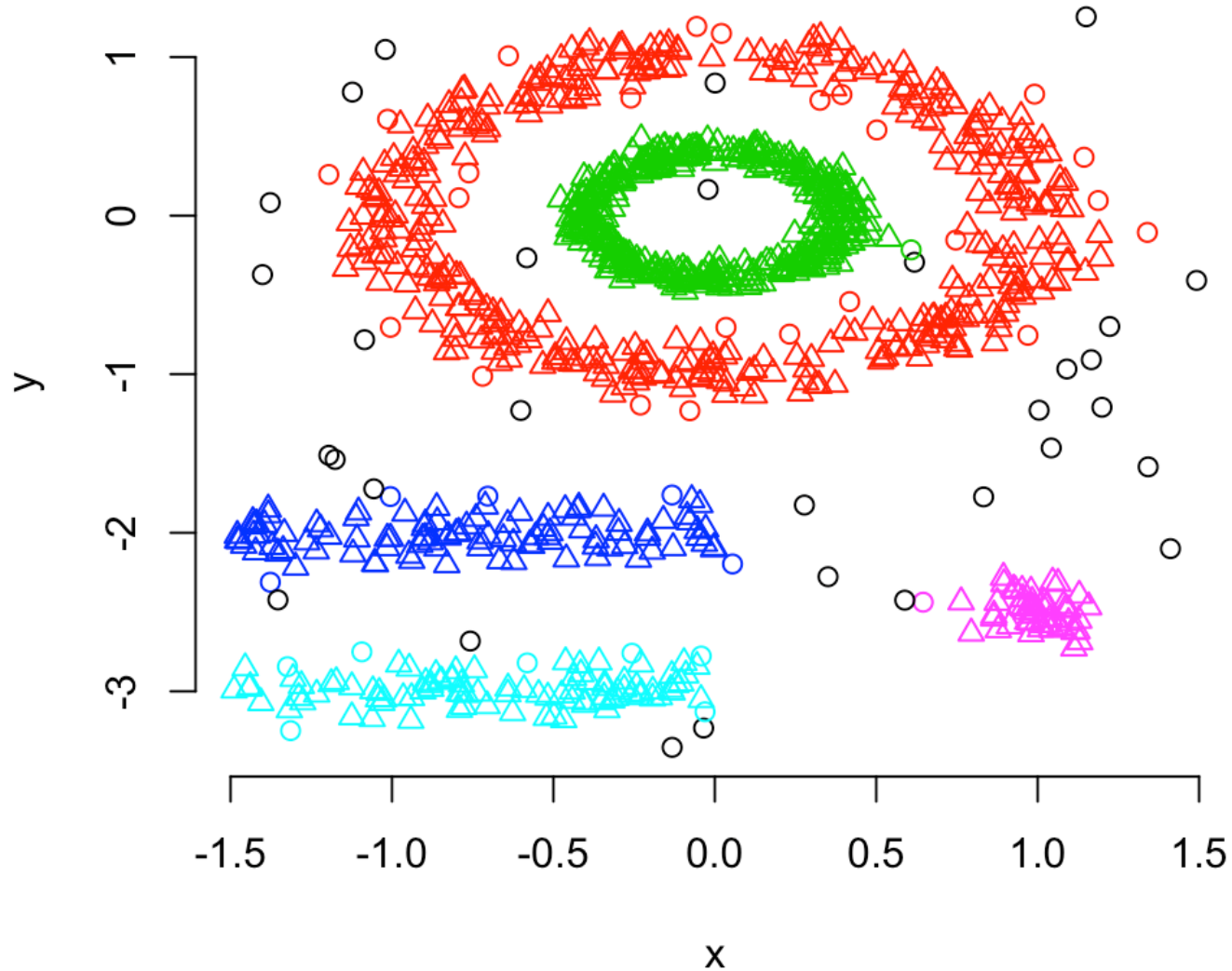


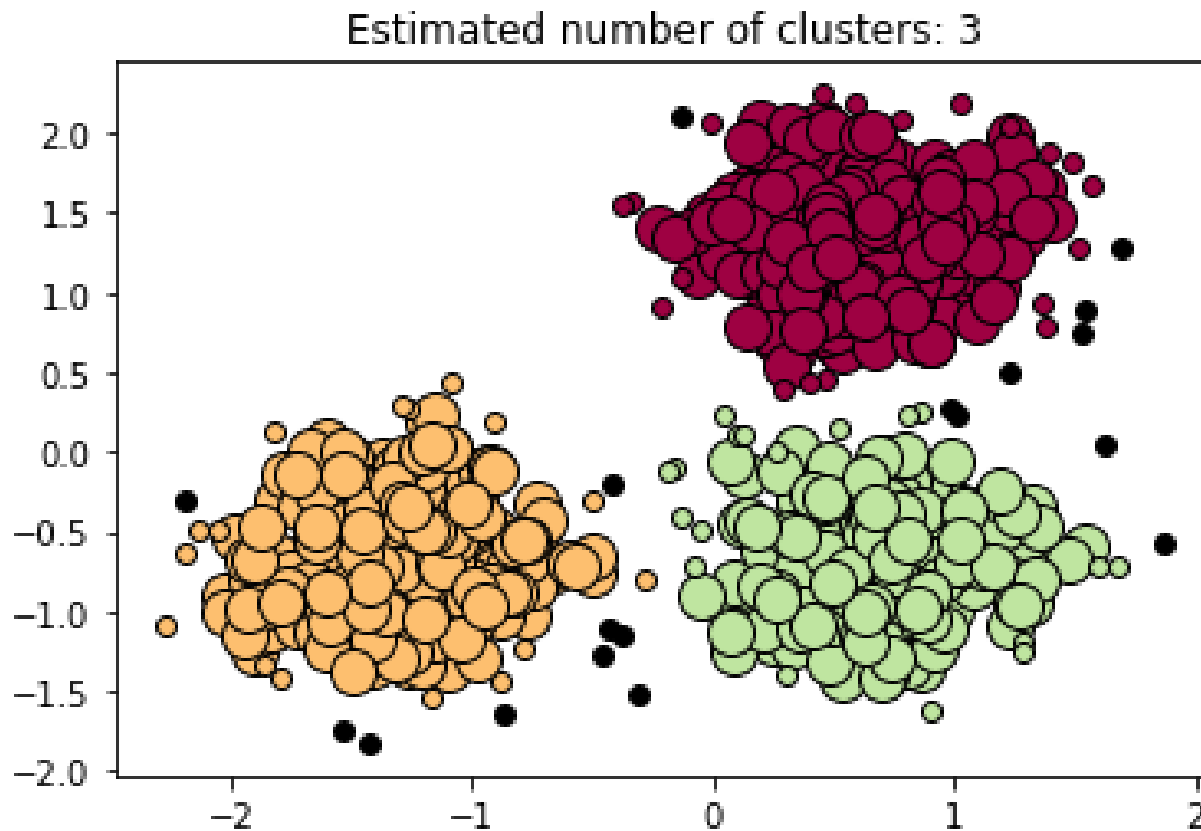
- The BIRCH builds a tree called the Characteristic Feature Tree (CFT) for the given data.
- The data is essentially lossy compressed to a set of Characteristic Feature nodes (CF Nodes).
- The CF Nodes have a number of subclusters called Characteristic Feature subclusters (CF Subclusters) and these CF Subclusters located in the non-terminal CF Nodes can have CF Nodes as children.
- The CF Subclusters hold the necessary information for clustering which prevents the need to hold the entire input data in memory. This information includes:

- Number of samples in a subcluster.
- Linear Sum - an N-dimensional vector holding the sum of all samples
- Squared Sum - Sum of the squared L2 norm of all samples.
- Centroids - To avoid recalculation linear sum / n\_samples.
- Squared norm of the centroids.
- The BIRCH algorithm has two parameters: the threshold and the branching factor.
- The branching factor limits the number of subclusters in a node and the threshold limits the distance between the entering sample and the existing subclusters.

- The DBSCAN algorithm views clusters as areas of high density separated by areas of low density.
- Clusters found by DBSCAN can be any shape.
- The central component to the DBSCAN is the concept of core samples, which are samples that are in areas of high density.
- A cluster is a set of core samples, each close to each other (measured by some distance measure) and a set of non-core samples that are close to a core sample (but are not themselves core samples).
- There are two parameters to the algorithm, `min_samples` and `eps`, which define formally what we mean when we say dense.
- Higher `min_samples` or lower `eps` indicate higher density necessary to form a cluster.

- A core sample as being a sample in the dataset such that there exist  $\text{min\_samples}$  other samples within a distance of  $\text{eps}$ , which are defined as neighbors of the core sample.
- The core sample is in a dense area of the vector space.
- A cluster is a set of core samples that can be built by recursively taking a core sample, finding all of their neighbors that are core samples, and so on.
- A cluster also has a set of non-core samples, which are samples that are neighbors of a core sample in the cluster but are not themselves core samples called border.
- Any sample that is not a core sample or border is an outlier.





- Adjusted Mutual Information
- Rand index
- Calinski and Harabaz score
- Davies-Bouldin score
- Completeness metric of a cluster labeling given a ground truth.
- Contingency matrix describing the relationship between labels.
- Fowlkes Mallows Score

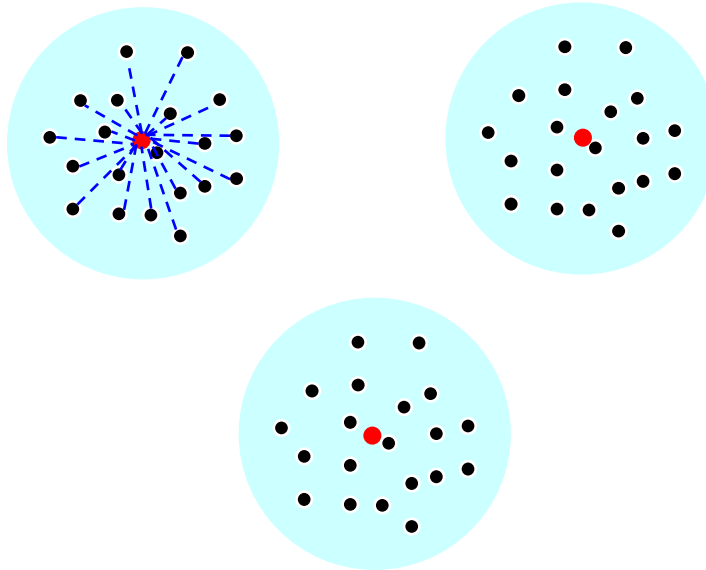


- Homogeneity
- Completeness
- V-Measure scores
- Mutual Information between two clusterings
- Normalized Mutual Information between two clusterings
- Silhouette Score
- Silhouette Samples

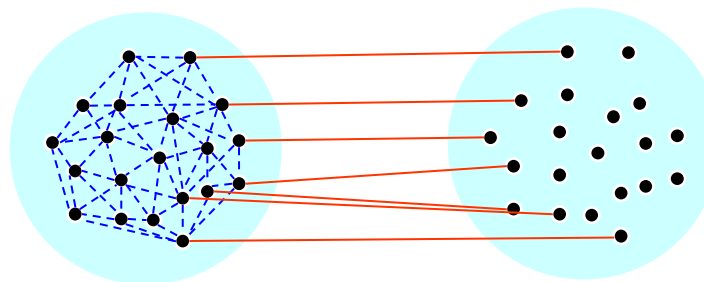
- Cohesion: Measures how closely related are objects in a cluster
- Separation: Measure how distinct or well-separated a cluster is from other clusters
- Sum Squared Error (SSE)
  - Cohesion is measured by the within cluster sum of squares (SSE) 
$$WSS = \sum_i \sum_{x \in C_i} (x - m_i)^2$$
  - Separation is measured by the between cluster sum of squares 
$$BSS = \sum_i |C_i| (m - m_i)^2$$

# Sum Squared Error (SSE)

$$SSE = \sum_{j=1}^C \sum_{i=1}^{N_j} \left( dist(\mathbf{x}_i, \mathbf{m}_j) \right)^2$$



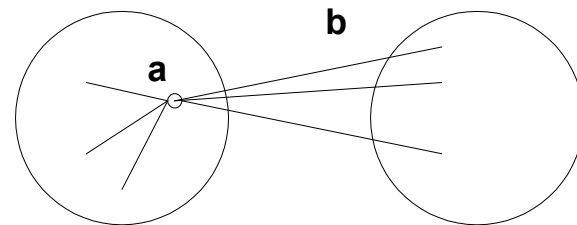
- A proximity graph based approach
  - Cluster cohesion is the sum of the weight of all links within a cluster.
  - Cluster separation is the sum of the weights between nodes in the cluster and nodes outside the cluster.



- Silhouette combine ideas of both cohesion and separation, but for individual points, as well as clusters and clusterings
- For an individual point,  $i$ 
  - Calculate  $a$  = average distance of  $i$  to the points in its cluster
  - Calculate  $b$  = min (average distance of  $i$  to points in another cluster)
  - The silhouette coefficient for a point is then given by

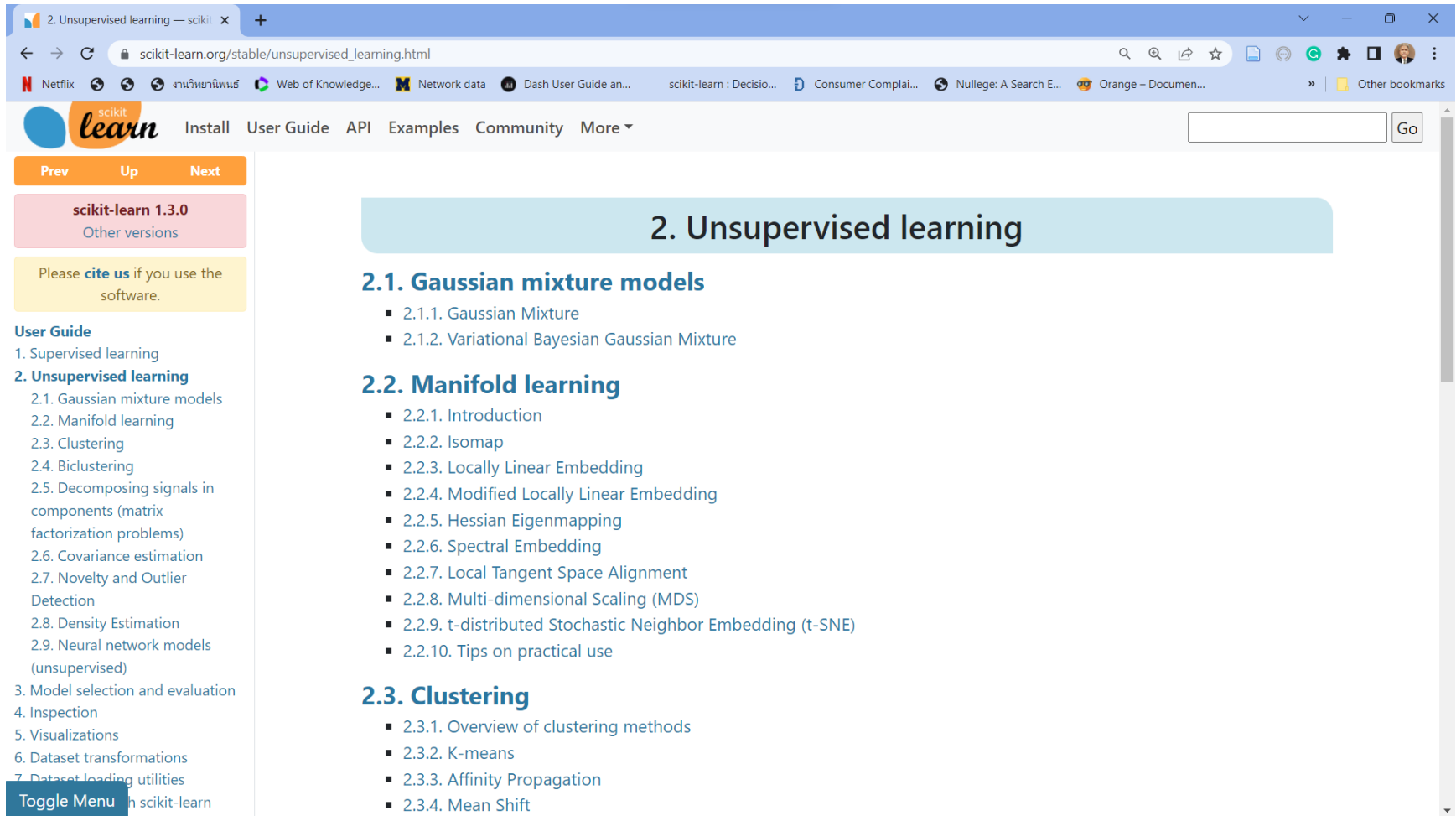
$$s = 1 - a/b \quad \text{if } a < b, \quad (\text{or } s = b/a - 1 \quad \text{if } a \geq b, \text{ not the usual case})$$

- Typically between 0 and 1.
- The closer to 1 the better.



- Can calculate the Average Silhouette width for a cluster or a clustering

- Large amount of data are without labels
- Clustering or Unsupervised learning, machine learning data without targets or labels
- Popular clustering algorithms are based on partitional techniques, connectivity or hierarchical techniques, and density based techniques.
- Each cluster technique has different drawbacks and advantages.
- To select which technique for a particular jobs, we must compare based on the clustering performances, such as SSE, Silhouette, etc.



The screenshot shows the scikit-learn website's "2. Unsupervised learning" page. The browser address bar shows the URL "scikit-learn.org/stable/unsupervised\_learning.html". The page features a navigation bar with links to "Install", "User Guide", "API", "Examples", "Community", and "More". A sidebar on the left contains a "User Guide" section with a list of topics, including "2. Unsupervised learning" which is currently selected. The main content area is titled "2. Unsupervised learning" and lists three sub-topics: "2.1. Gaussian mixture models", "2.2. Manifold learning", and "2.3. Clustering". Each sub-topic has a list of specific methods or techniques.

**2. Unsupervised learning**

- 2.1. Gaussian mixture models**
  - 2.1.1. Gaussian Mixture
  - 2.1.2. Variational Bayesian Gaussian Mixture
- 2.2. Manifold learning**
  - 2.2.1. Introduction
  - 2.2.2. Isomap
  - 2.2.3. Locally Linear Embedding
  - 2.2.4. Modified Locally Linear Embedding
  - 2.2.5. Hessian Eigenmapping
  - 2.2.6. Spectral Embedding
  - 2.2.7. Local Tangent Space Alignment
  - 2.2.8. Multi-dimensional Scaling (MDS)
  - 2.2.9. t-distributed Stochastic Neighbor Embedding (t-SNE)
  - 2.2.10. Tips on practical use
- 2.3. Clustering**
  - 2.3.1. Overview of clustering methods
  - 2.3.2. K-means
  - 2.3.3. Affinity Propagation
  - 2.3.4. Mean Shift