

Representative-Based Clustering

Data Mining and Text Mining



For each cluster there is a point (or a parameter vector) that summarizes it

- Given a dataset of N instances, and a desired number of clusters k , this class of algorithms generates a partition C of N in k clusters $\{C_1, C_2, \dots, C_k\}$
- For each cluster there is a point that summarizes the cluster
- The common choice being the mean of the points in the cluster

$$\mu_i = \frac{1}{n_i} \sum_{x_j \in C_i} x_j$$

where $n_i = |C_i|$ and μ_i is the centroid

Brute-force Approach

Generate all the possible clustering $C = \{C_1, C_2, \dots, C_k\}$

Then, select the best one

Unfortunately, there are $O(k^N/k!)$ possible partitions

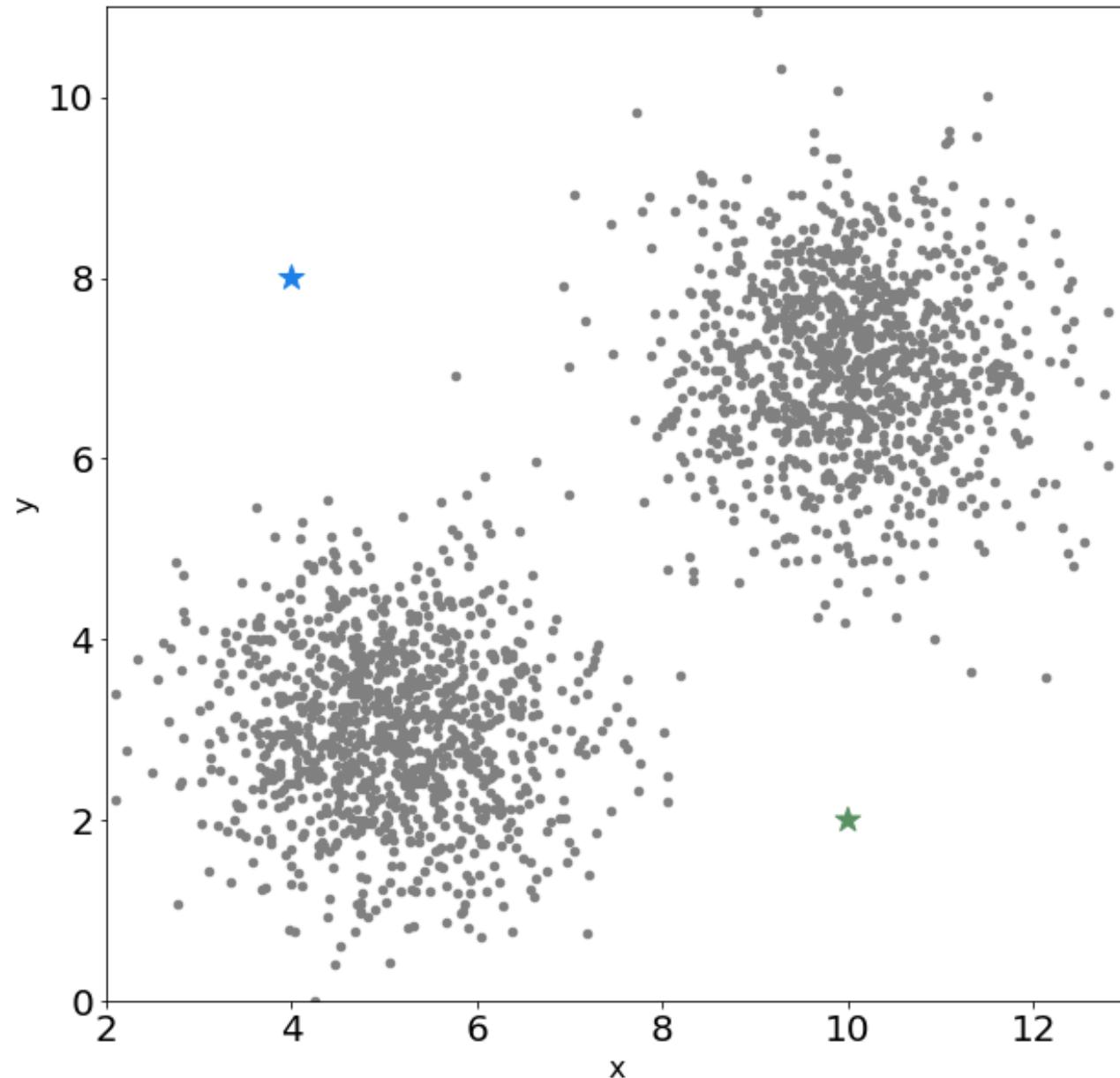
K-Means algorithm

- Greedy iterative approach to find a clustering that minimizes the SSE objective:

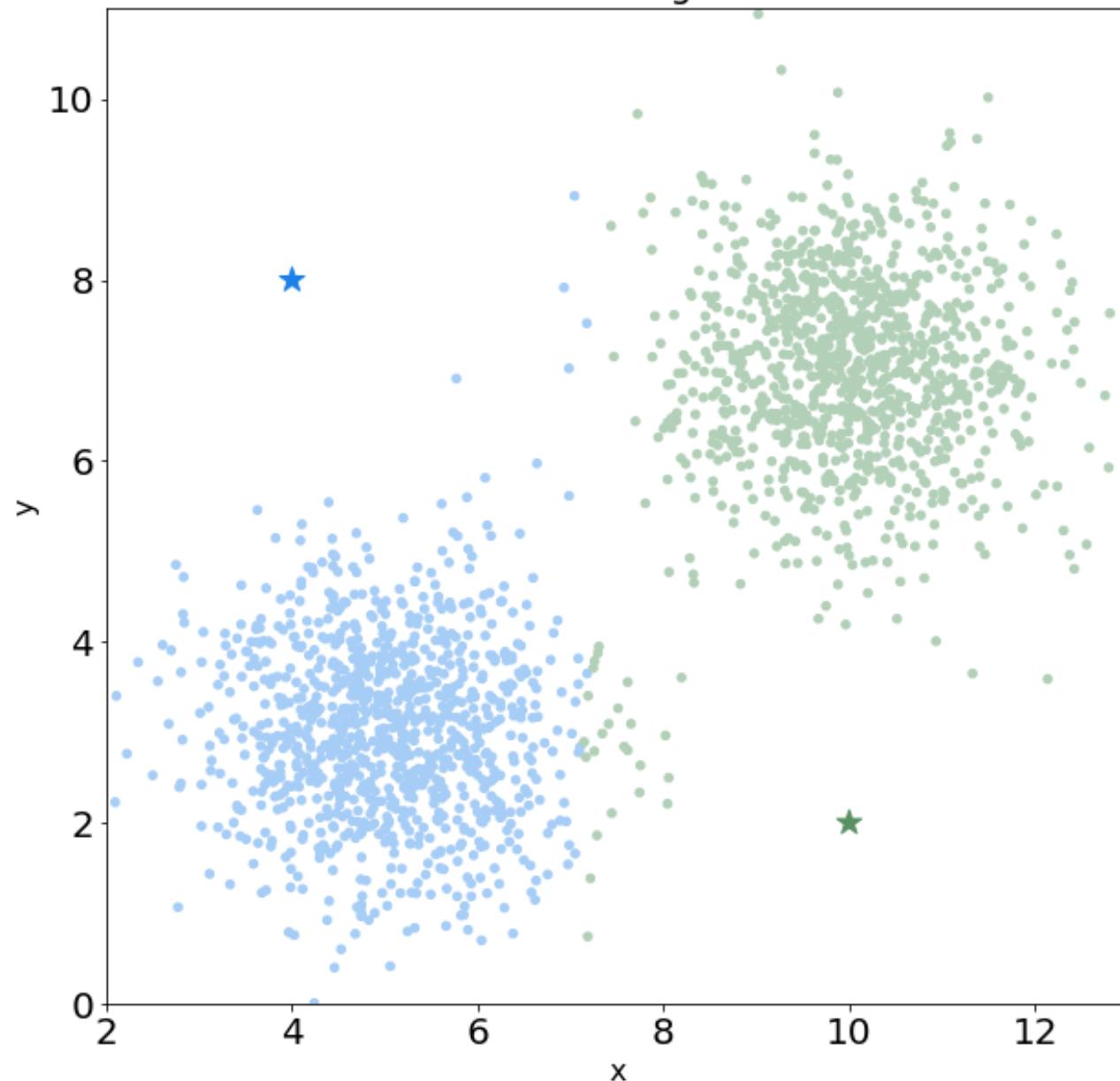
$$\text{SSE}(C) = \sum_{i=1}^k \sum_{x_j \in C_i} ||x_j - \mu_i||^2$$

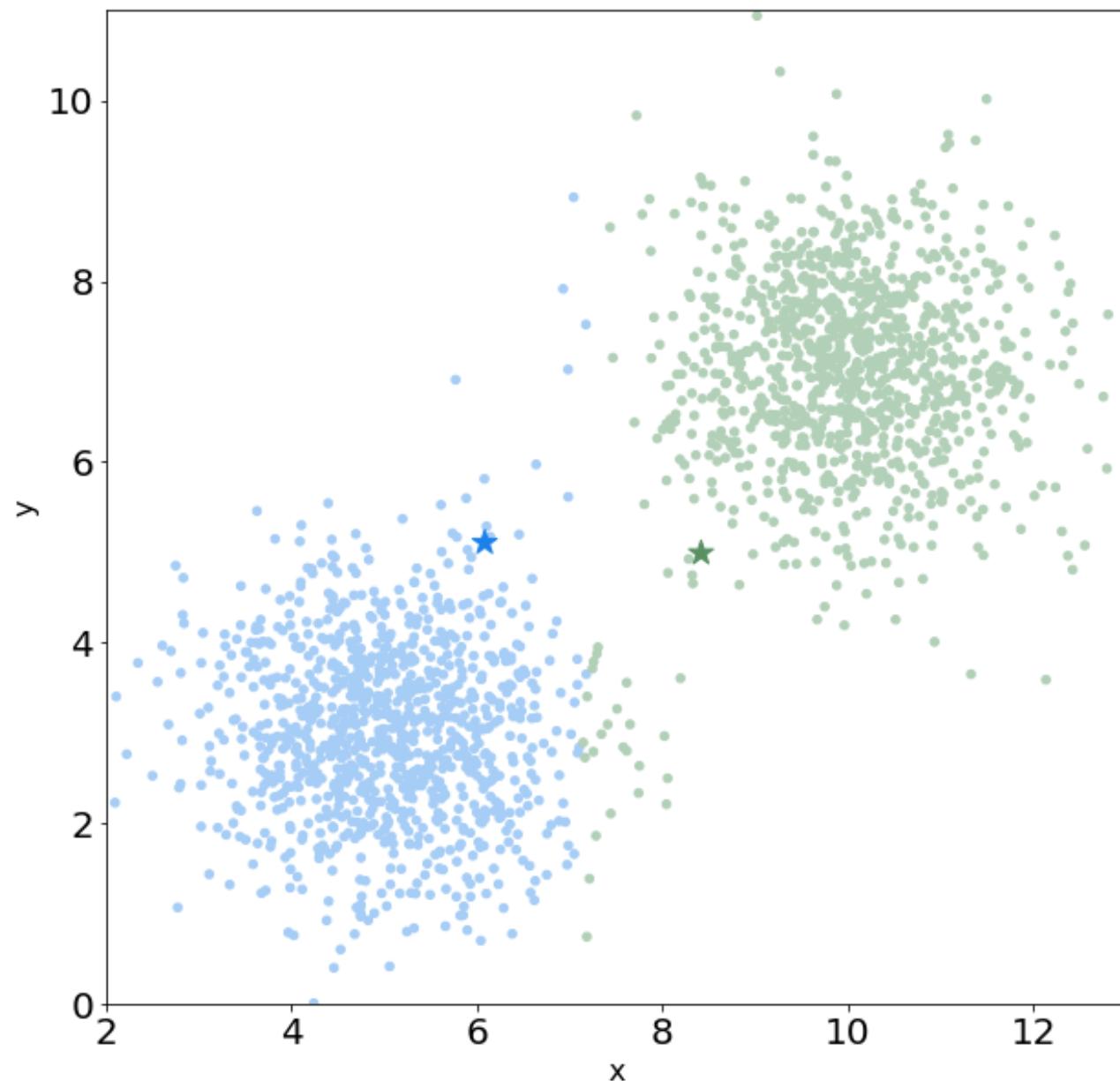
- The goal of the clustering process is thus to find

$$C^* = \arg \min_C \text{SSE}(C)$$

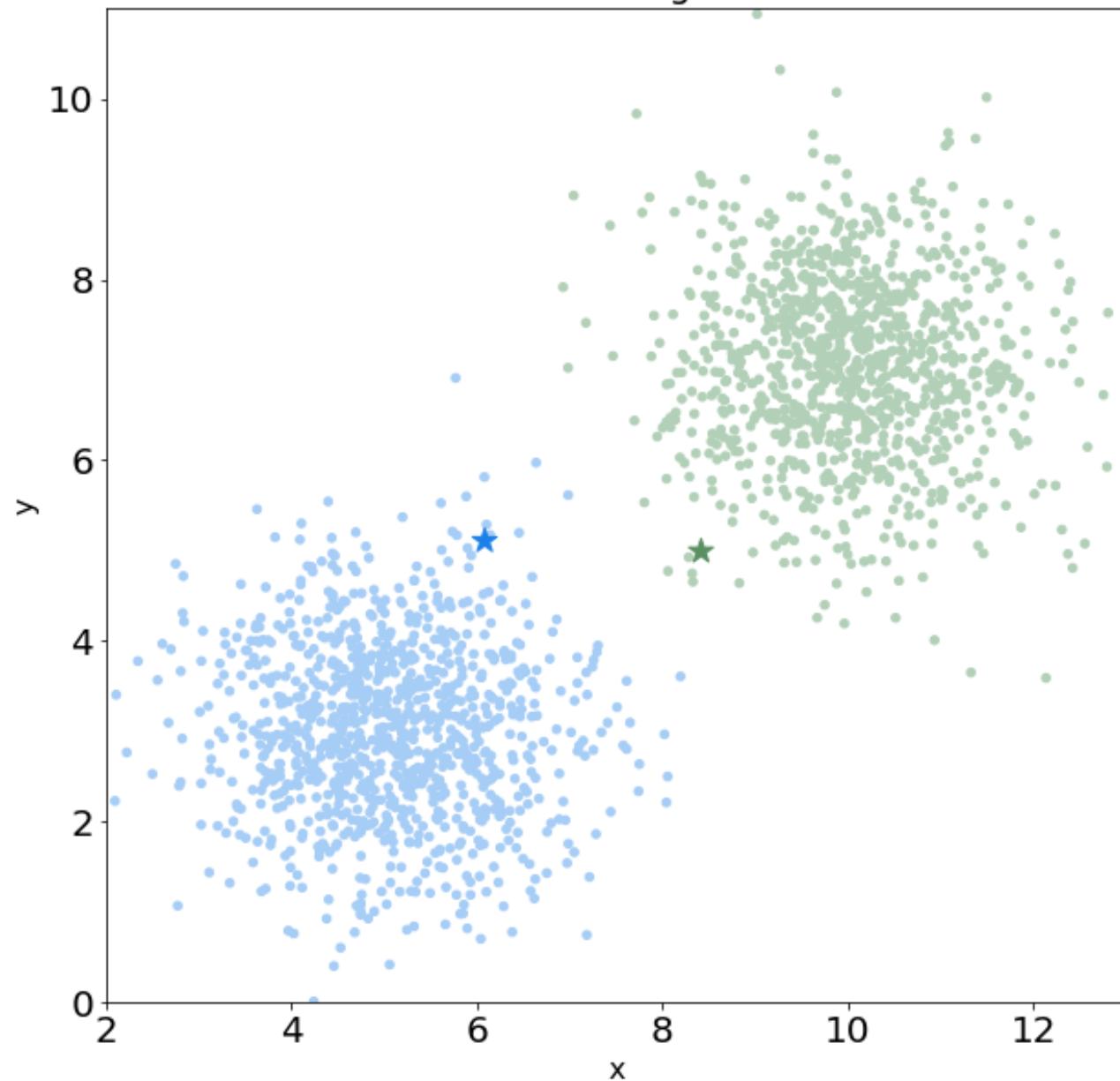


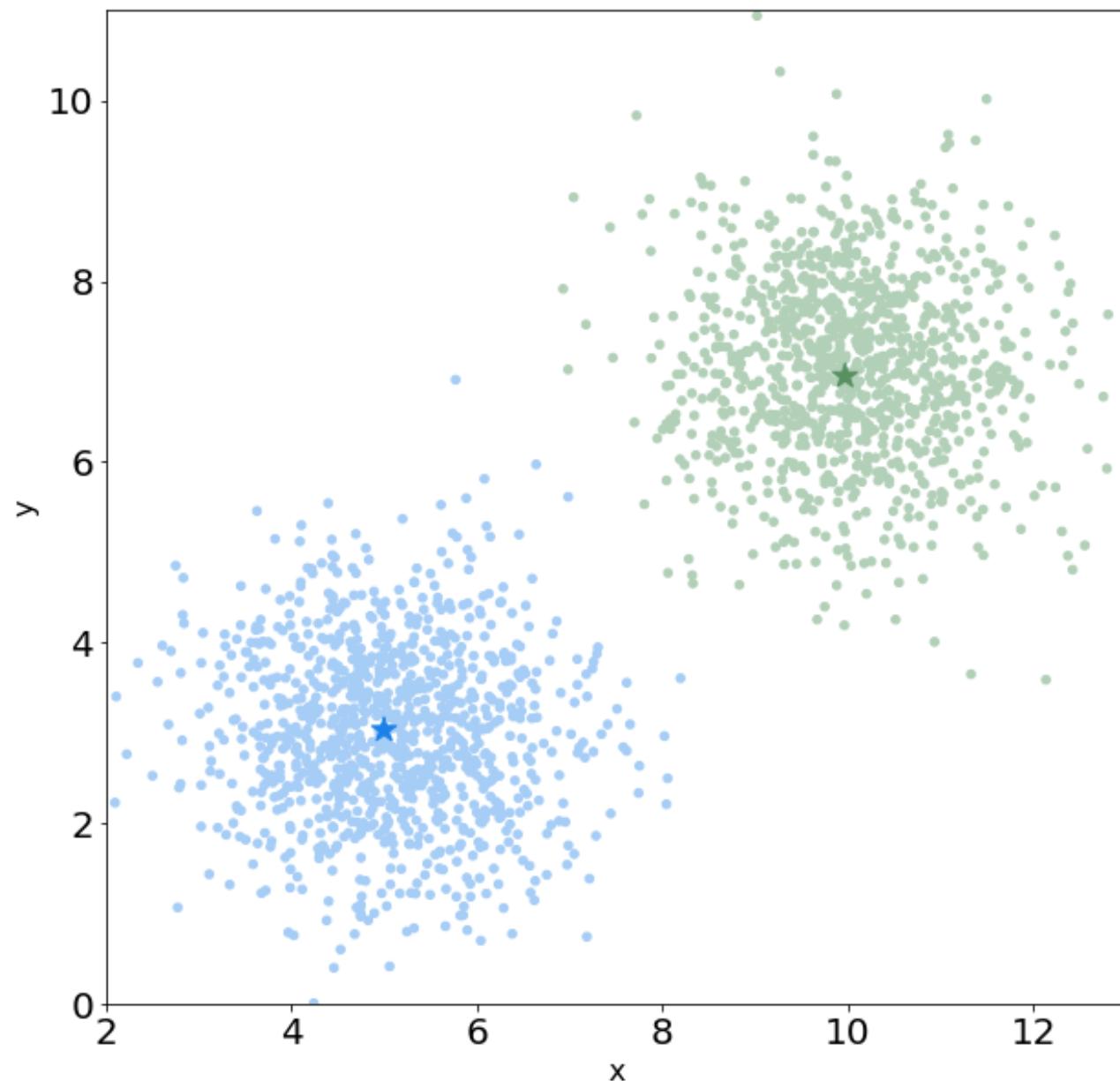
Iteration 0 - Assign Cluster



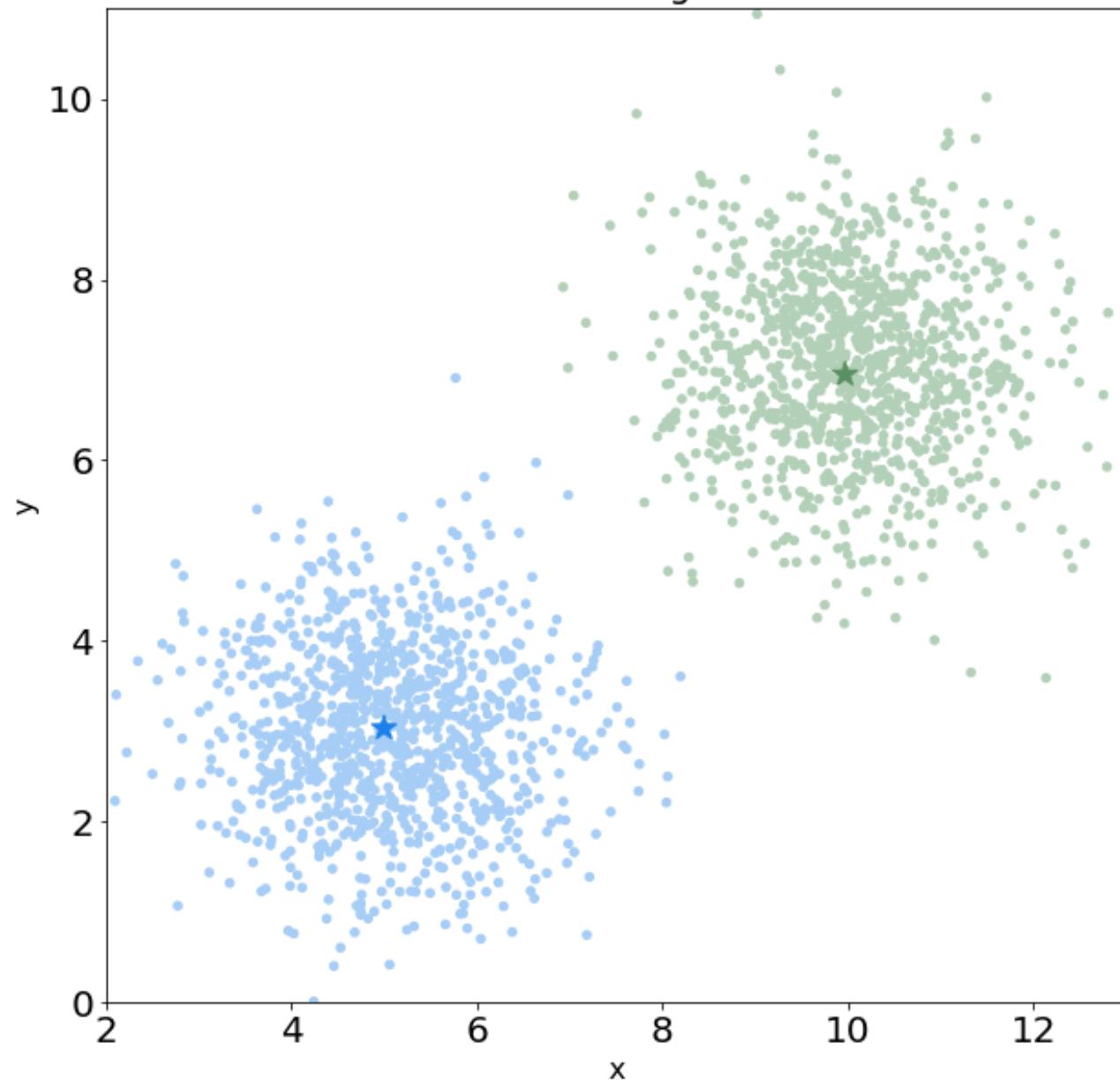


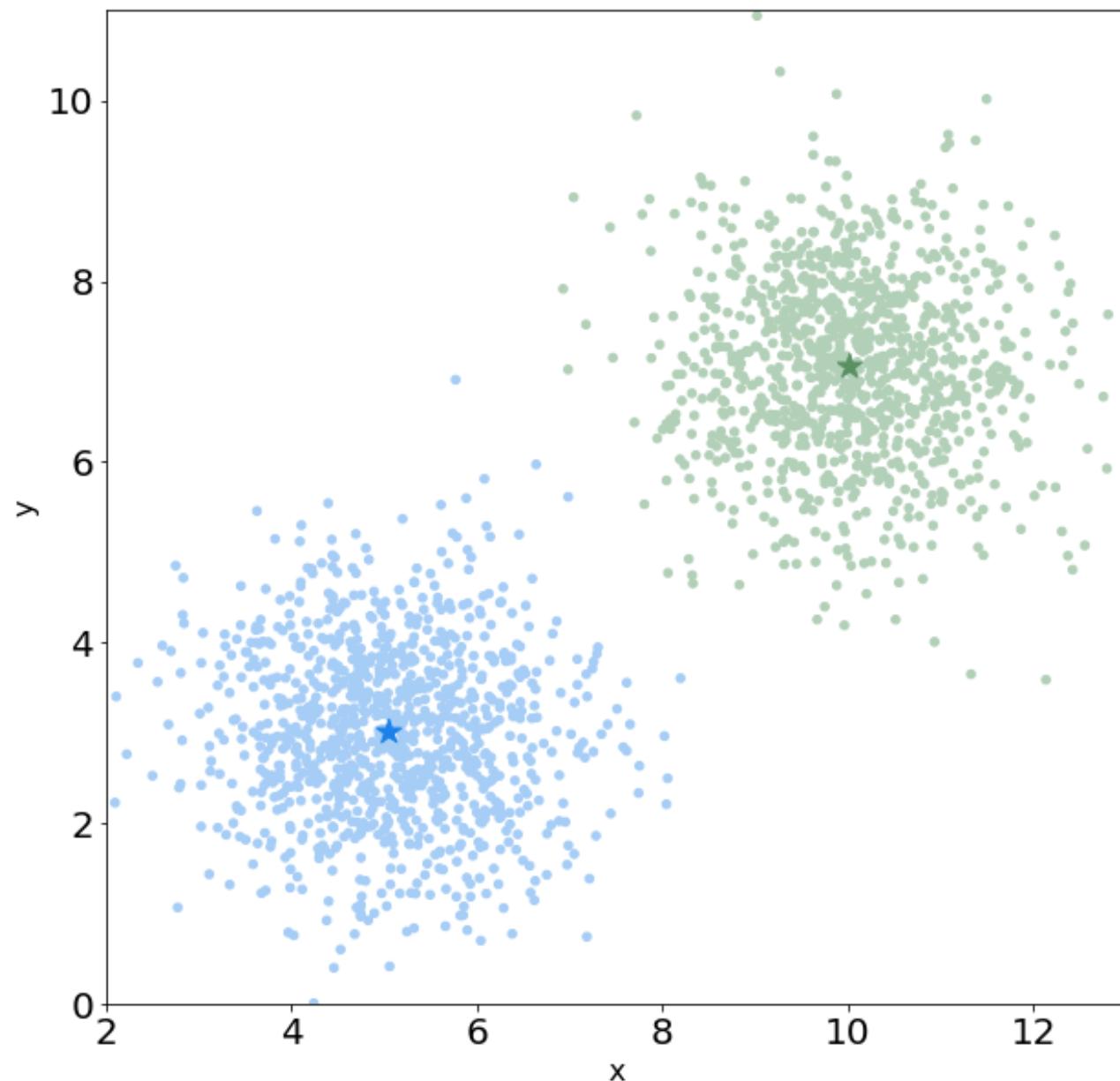
Iteration 1 - Assign Cluster





Iteration 2 - Assign Cluster





ALGORITHM 13.1. K-means Algorithm

K-MEANS (\mathbf{D}, k, ϵ):

```
1  $t = 0$ 
2 Randomly initialize  $k$  centroids:  $\mu_1^t, \mu_2^t, \dots, \mu_k^t \in \mathbb{R}^d$ 
3 repeat
4    $t \leftarrow t + 1$ 
5    $C_j \leftarrow \emptyset$  for all  $j = 1, \dots, k$ 
     // Cluster Assignment Step
6   foreach  $\mathbf{x}_j \in \mathbf{D}$  do
7      $j^* \leftarrow \operatorname{argmin}_i \left\{ \|\mathbf{x}_j - \mu_i^{t-1}\|^2 \right\}$  // Assign  $\mathbf{x}_j$  to closest centroid
8      $C_{j^*} \leftarrow C_{j^*} \cup \{\mathbf{x}_j\}$ 
     // Centroid Update Step
9   foreach  $i = 1$  to  $k$  do
10     $\mu_i^t \leftarrow \frac{1}{|C_i|} \sum_{\mathbf{x}_j \in C_i} \mathbf{x}_j$ 
11 until  $\sum_{i=1}^k \|\mu_i^t - \mu_i^{t-1}\|^2 \leq \epsilon$ 
```

- Most widely known representative-based algorithm
- Assumes an Euclidean space but sometimes it can be extended to the non-Euclidean case
- Employs a greedy iterative approaches that minimizes the SSE objective. Accordingly it can converge to a local optimal instead of a globally optimal clustering.

- Cluster assignment takes $O(nkd)$ time since, for each of the n points, it computes its distance to each of the k clusters, which takes d operations in d dimensions
- The centroid re-computation step takes $O(nd)$ time because it adds at total of n d -dimensional points
- Assuming that there are t iterations, the total time for K-means is given as $O(tnkd)$.
- In terms of the I/O cost it requires $O(t)$ full database scans, because we have to read the entire database in each iteration.

Centroid initialization

- **Solution 1**

- Pick points that are as far away from one another as possible.

- **Solution 2**

Pick the first point at random;

WHILE there are fewer than k points DO

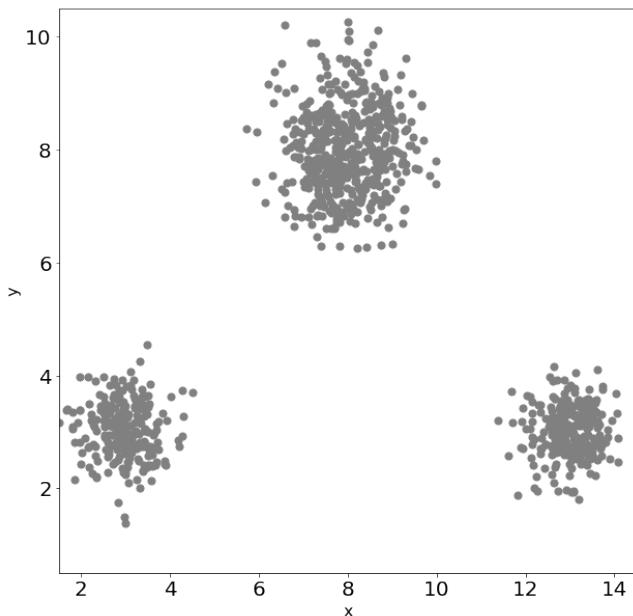
 Add the point whose minimum distance
 from the selected points is as large as
 possible;

END;

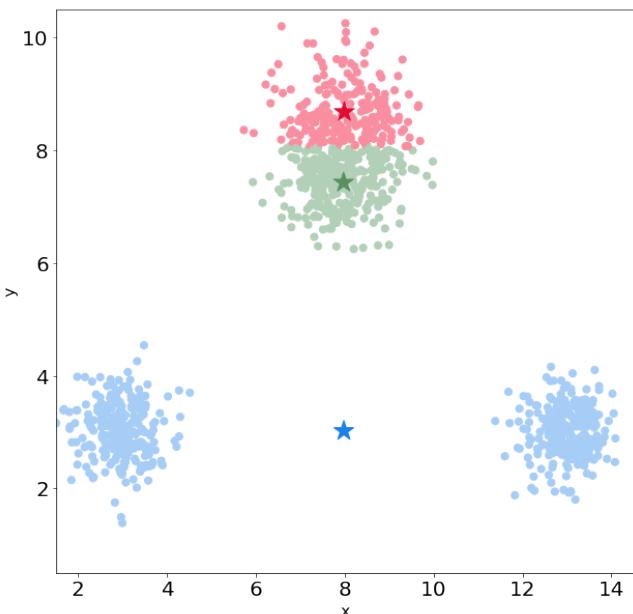
- **Solution 3**

- Cluster a sample of the data, perhaps hierarchically, so there are k clusters. Pick a point from each cluster, perhaps that point closest to the centroid of the cluster.

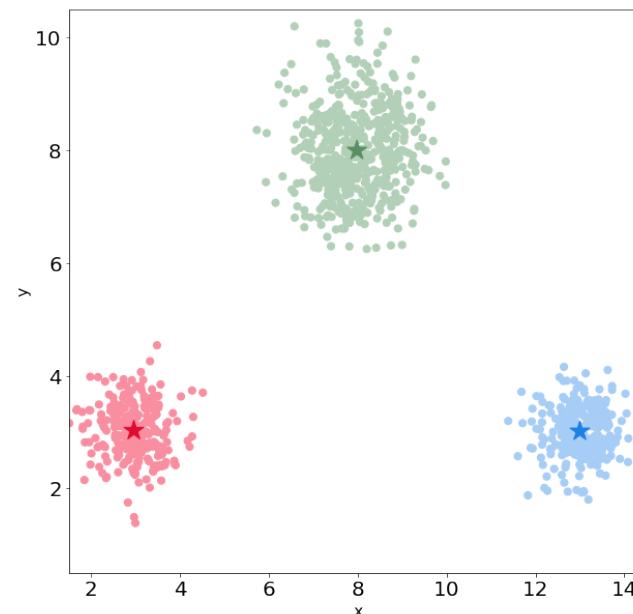
Initial centroids matters!



Original Points

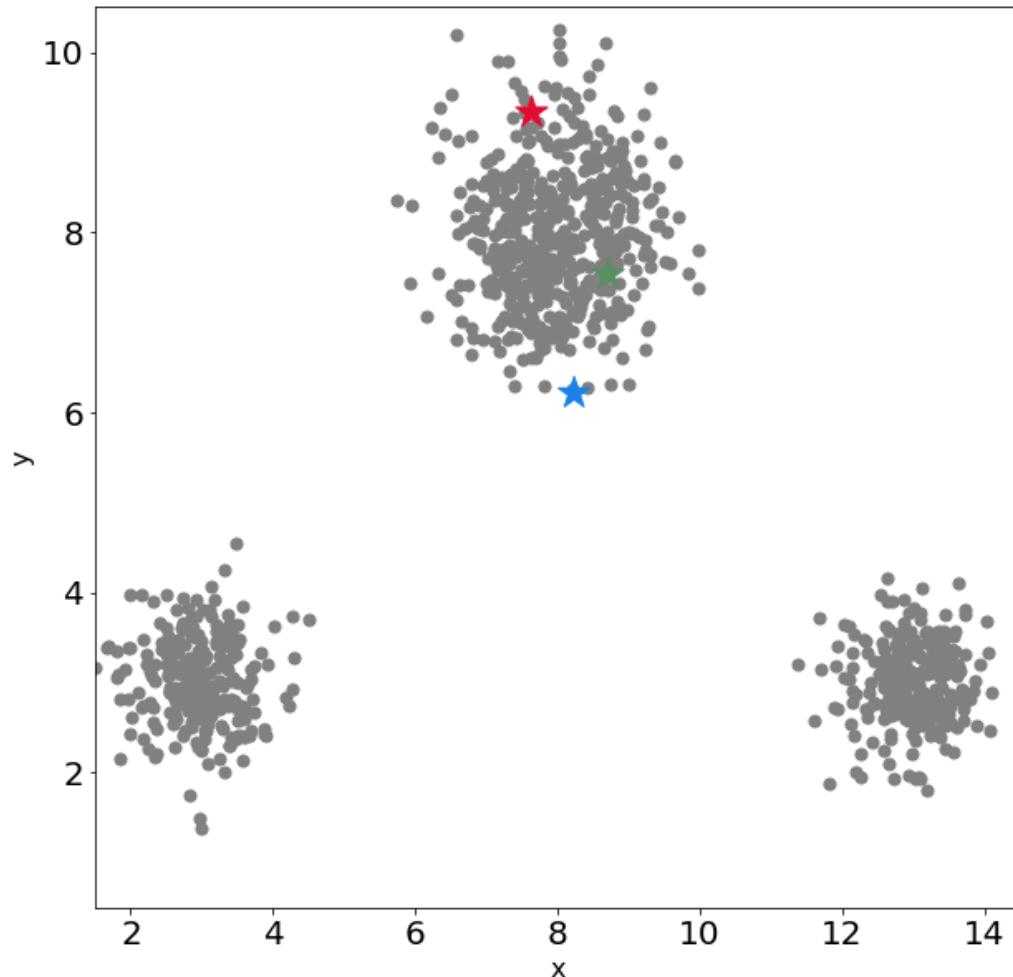


Sub-optimal Clustering

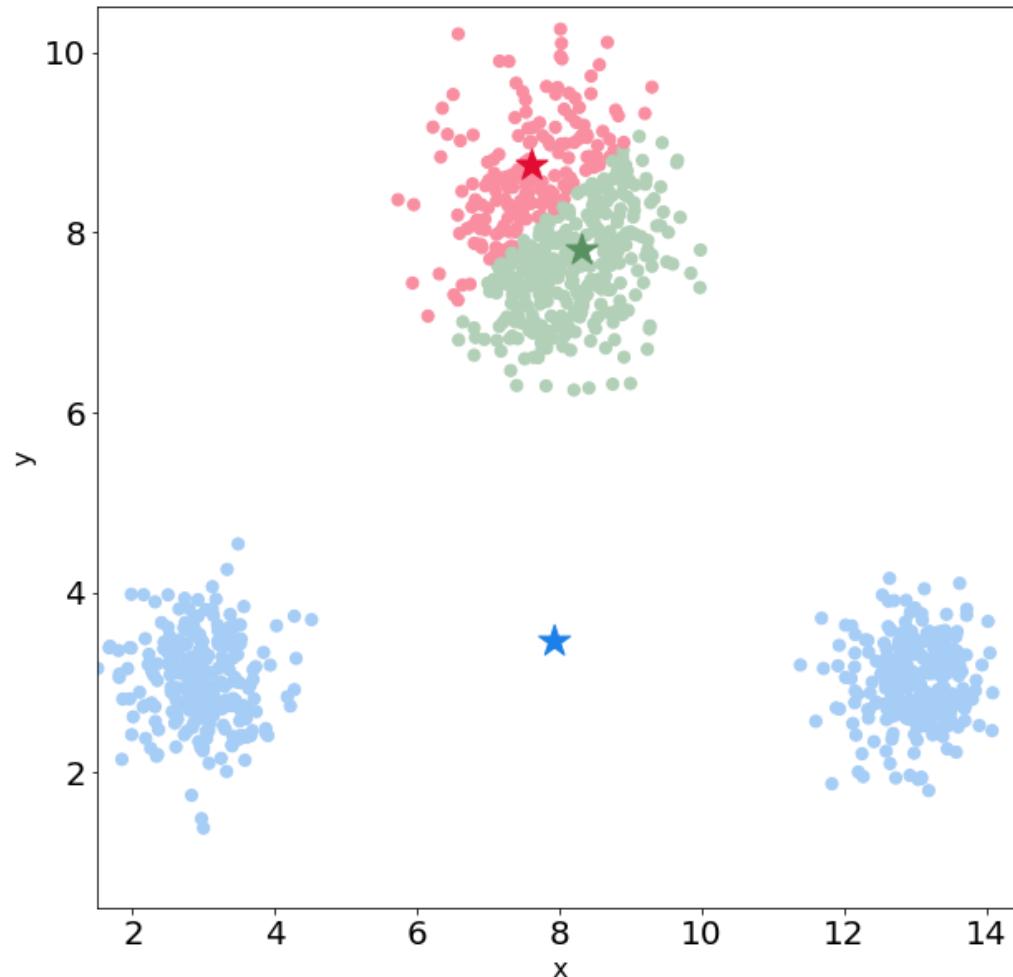


Optimal Clustering

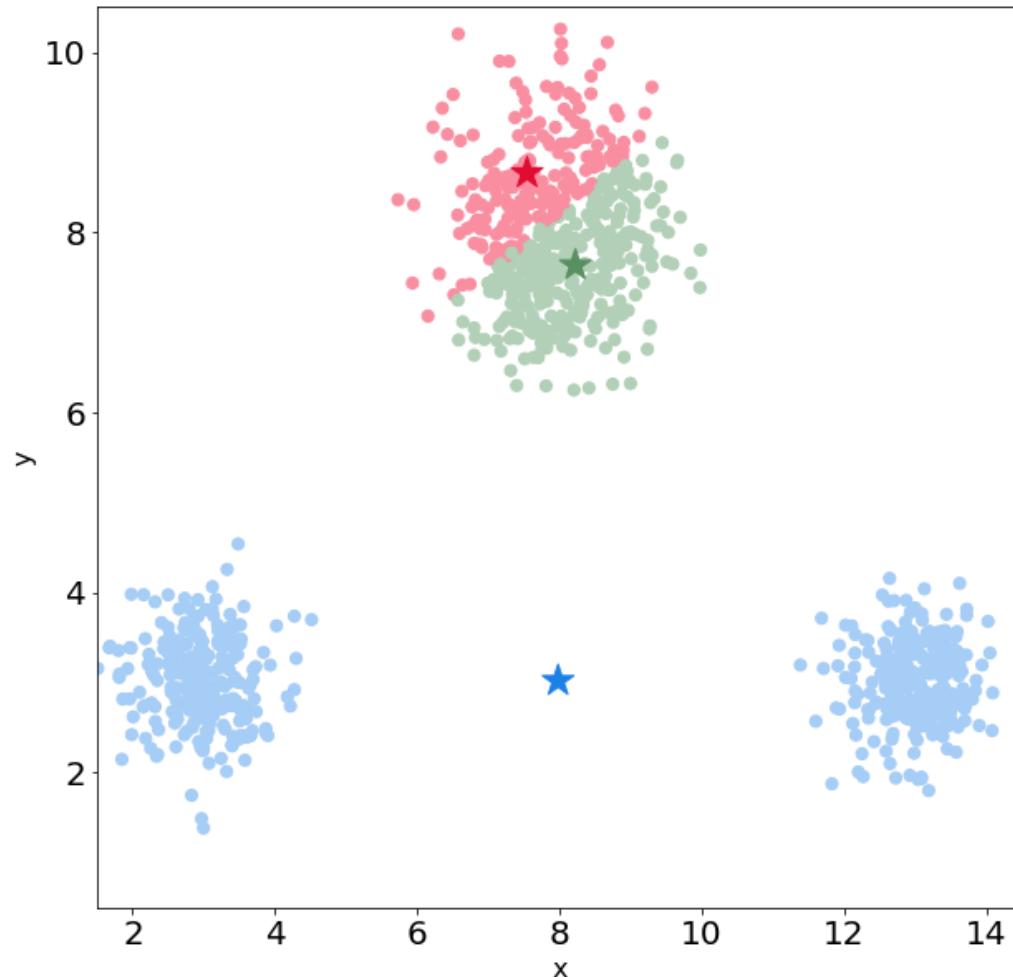
Run # 1



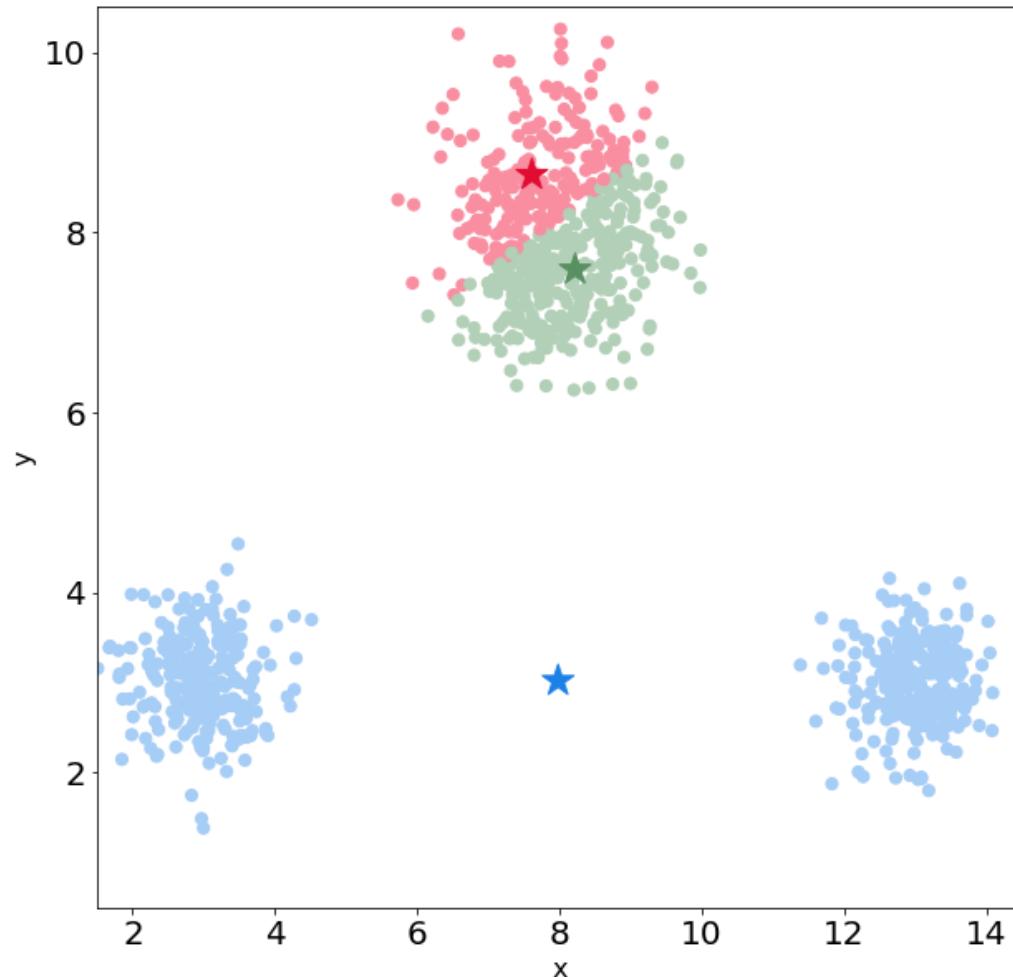
Run #1 – Start



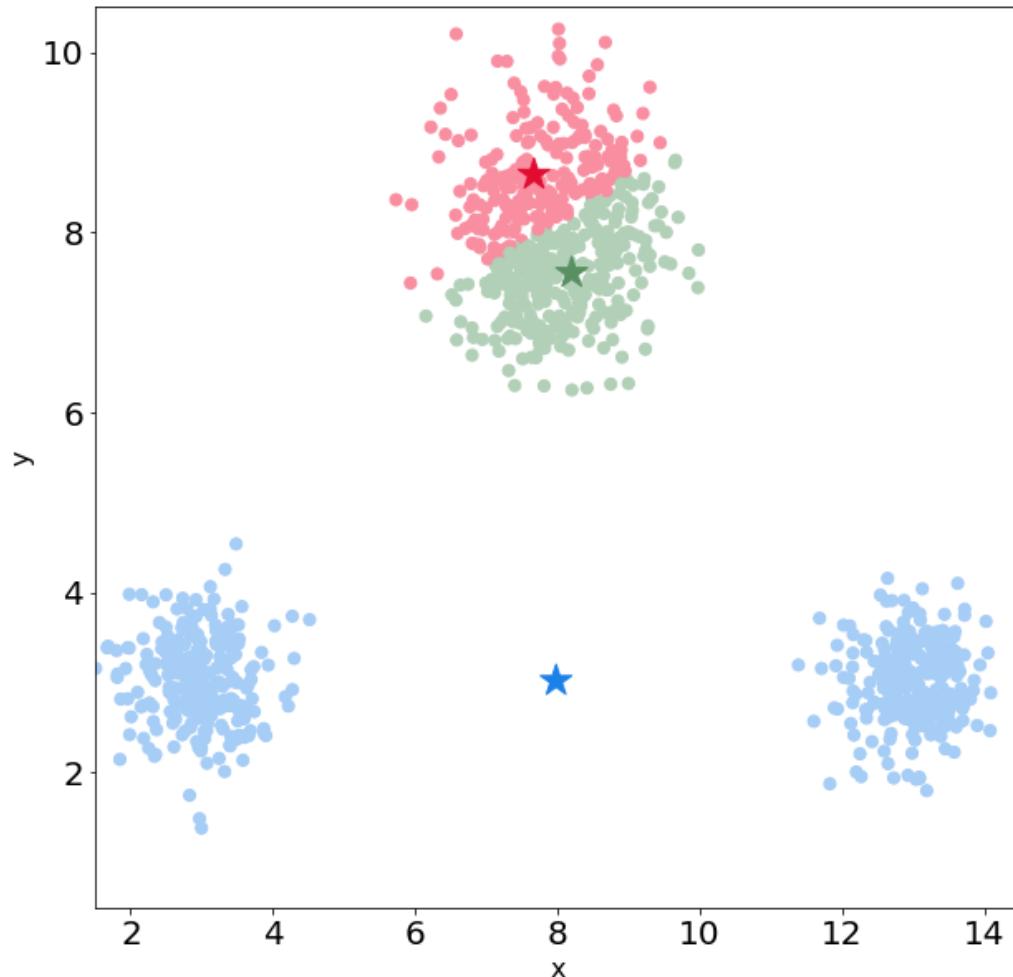
Run #1 – Step 1



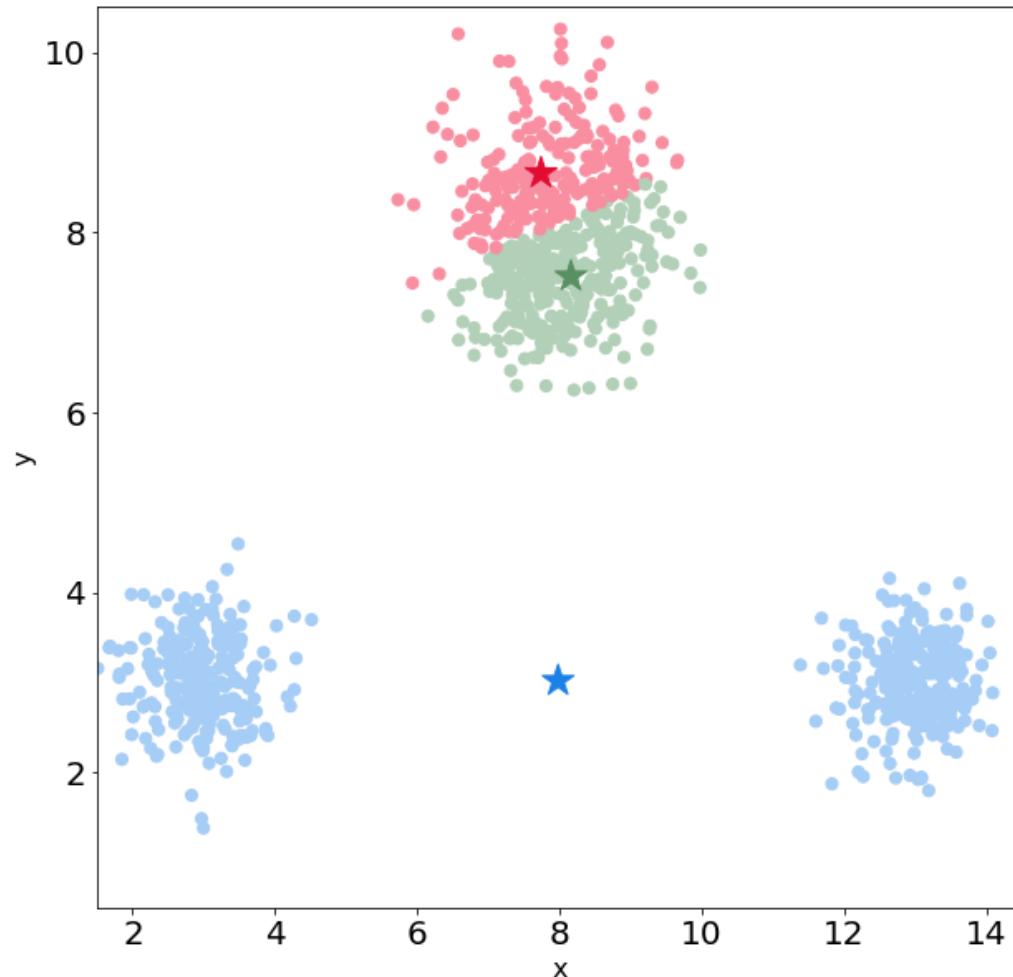
Run #1 – Step 2



Run #1 – Step 3

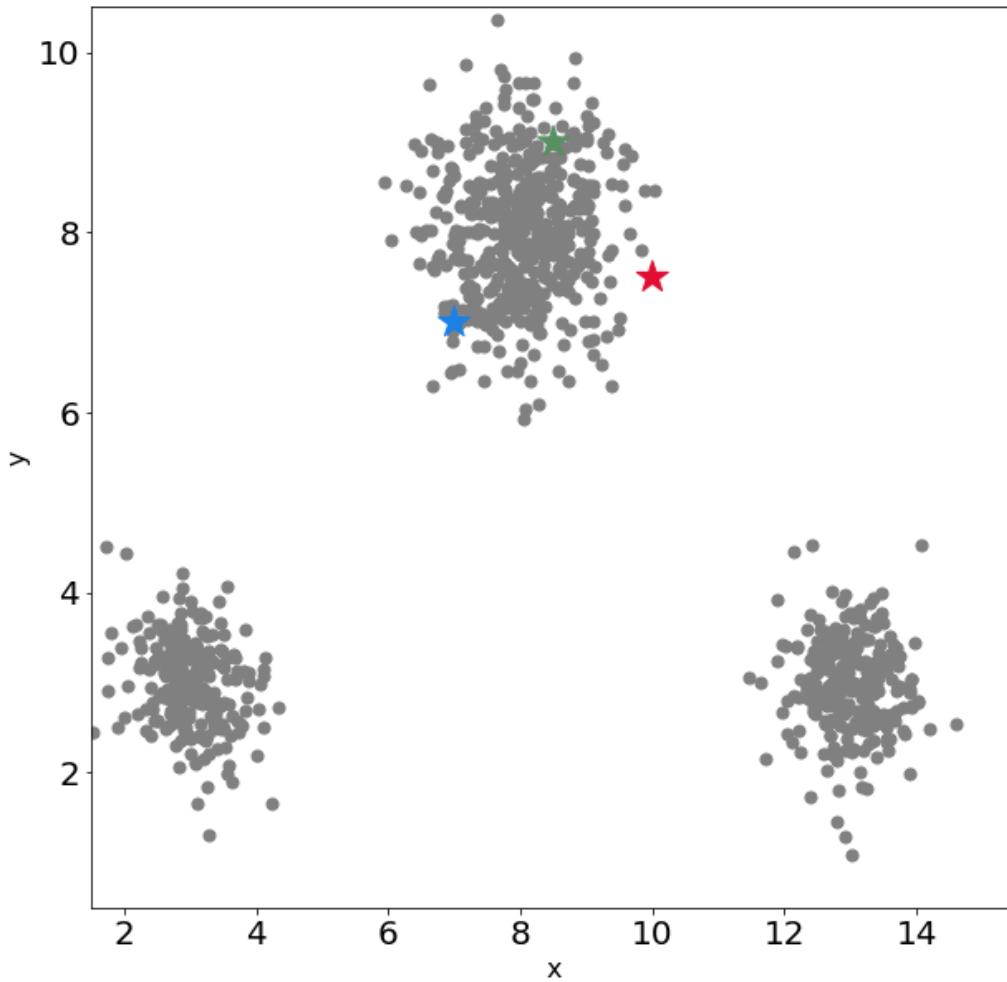


Run #1 – Step 4

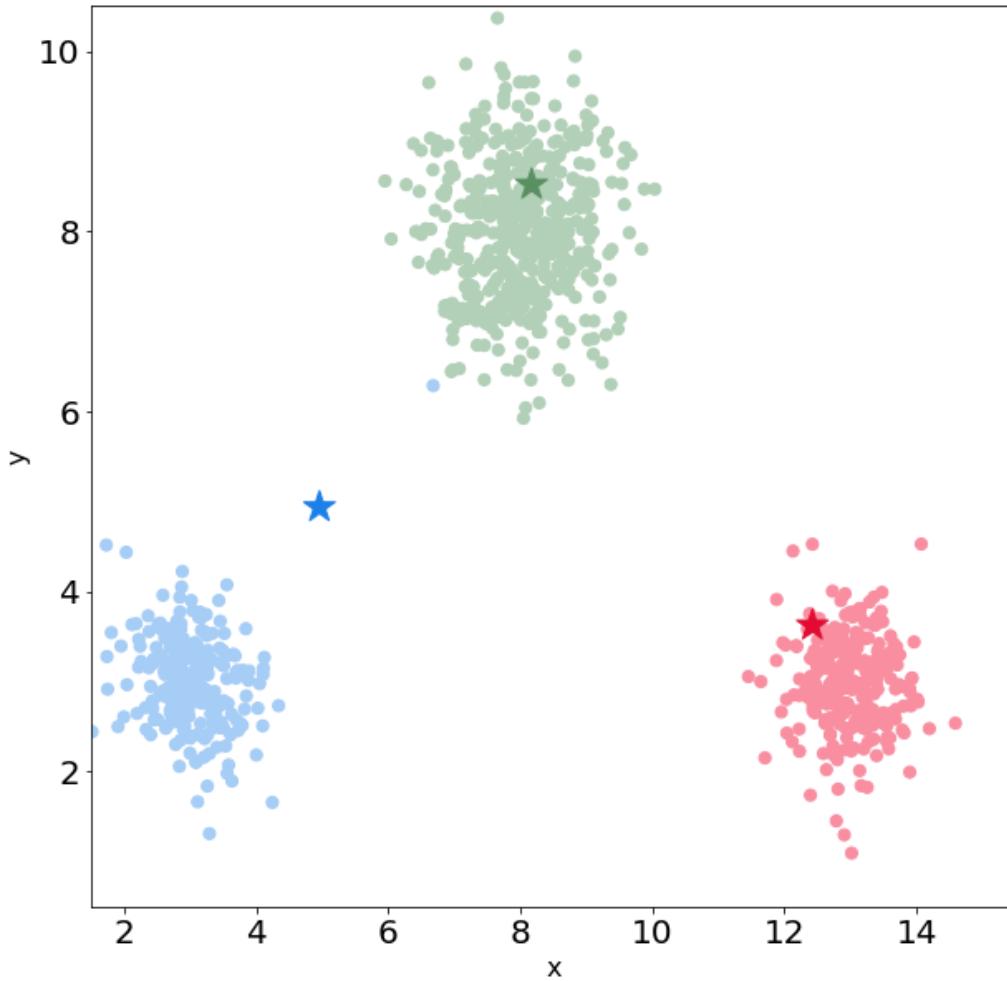


Run #1 – Step 5

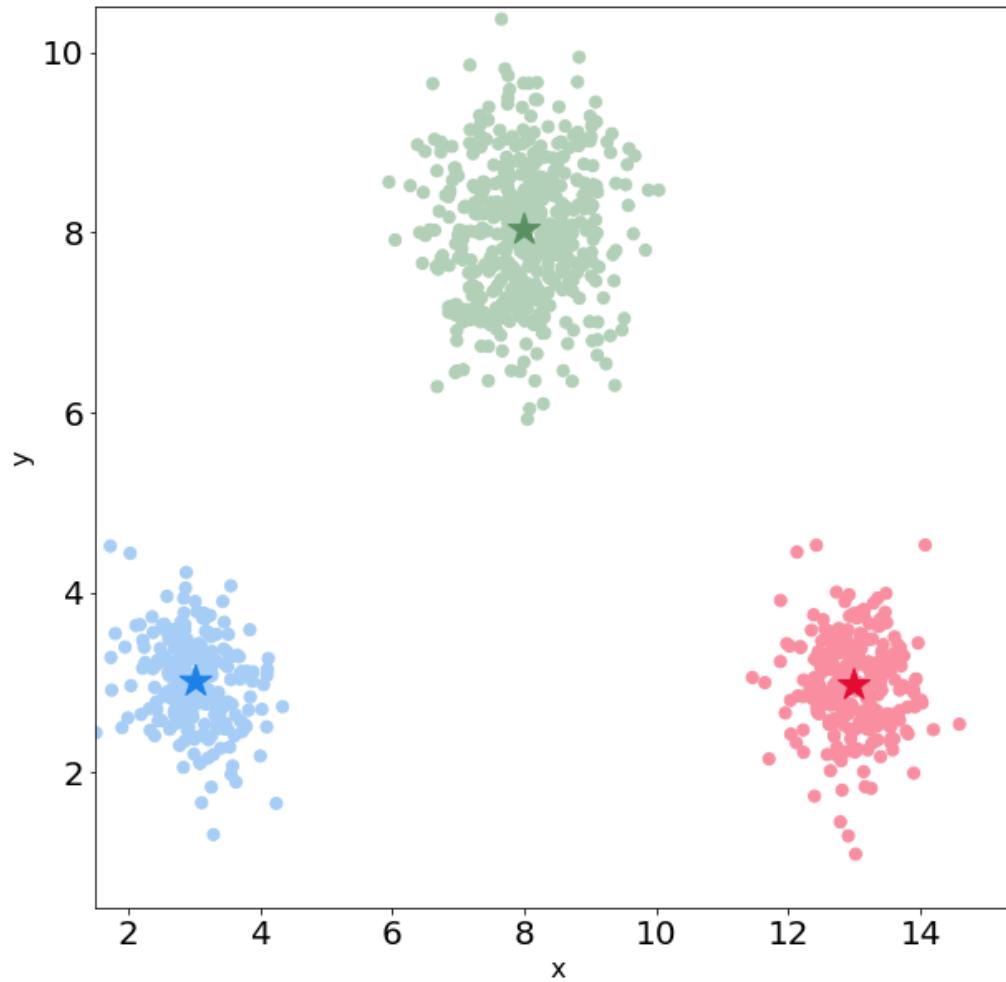
Run #2



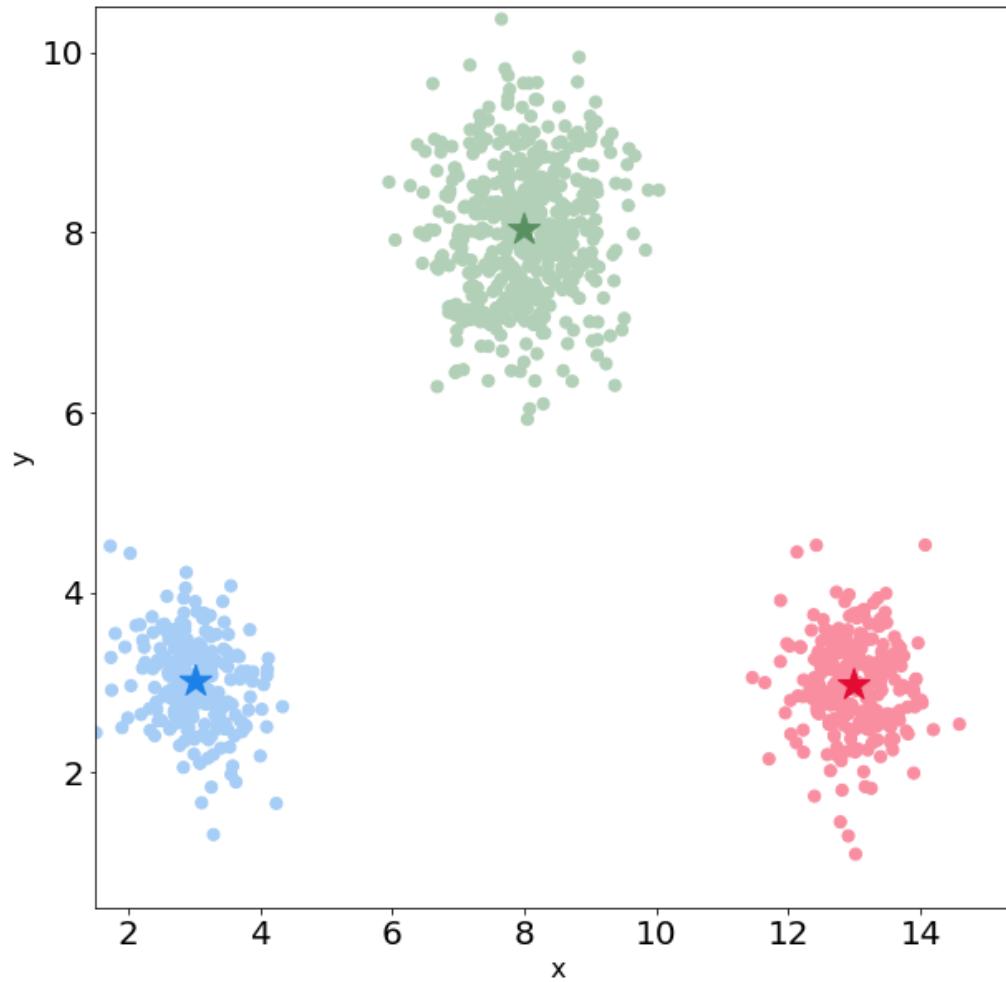
Run #2 – Start



Run #2 – Step 1



Run #2 – Step 2



Run #2 – Step 3

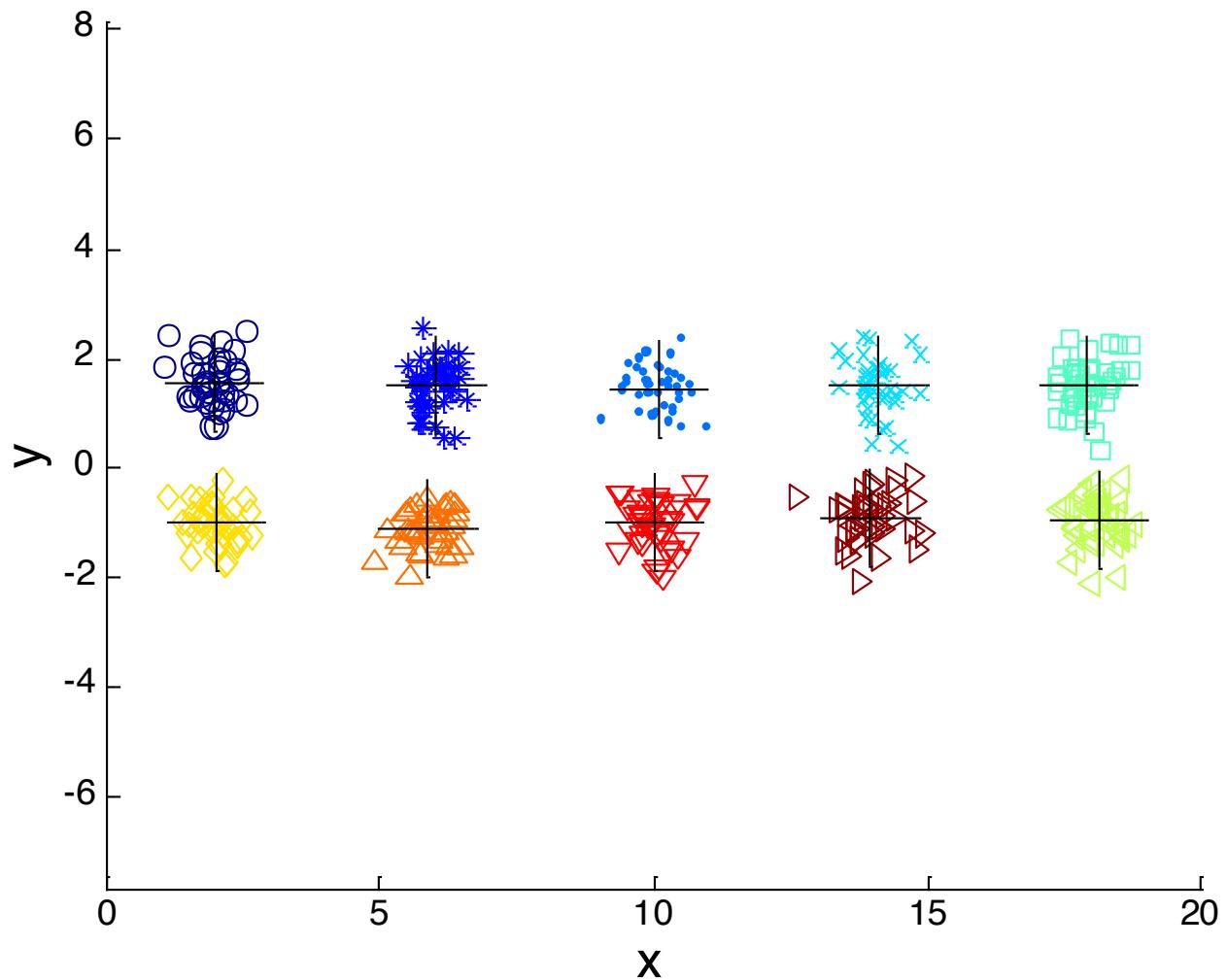
Why Selecting the Best Initial Centroids is Difficult?

33

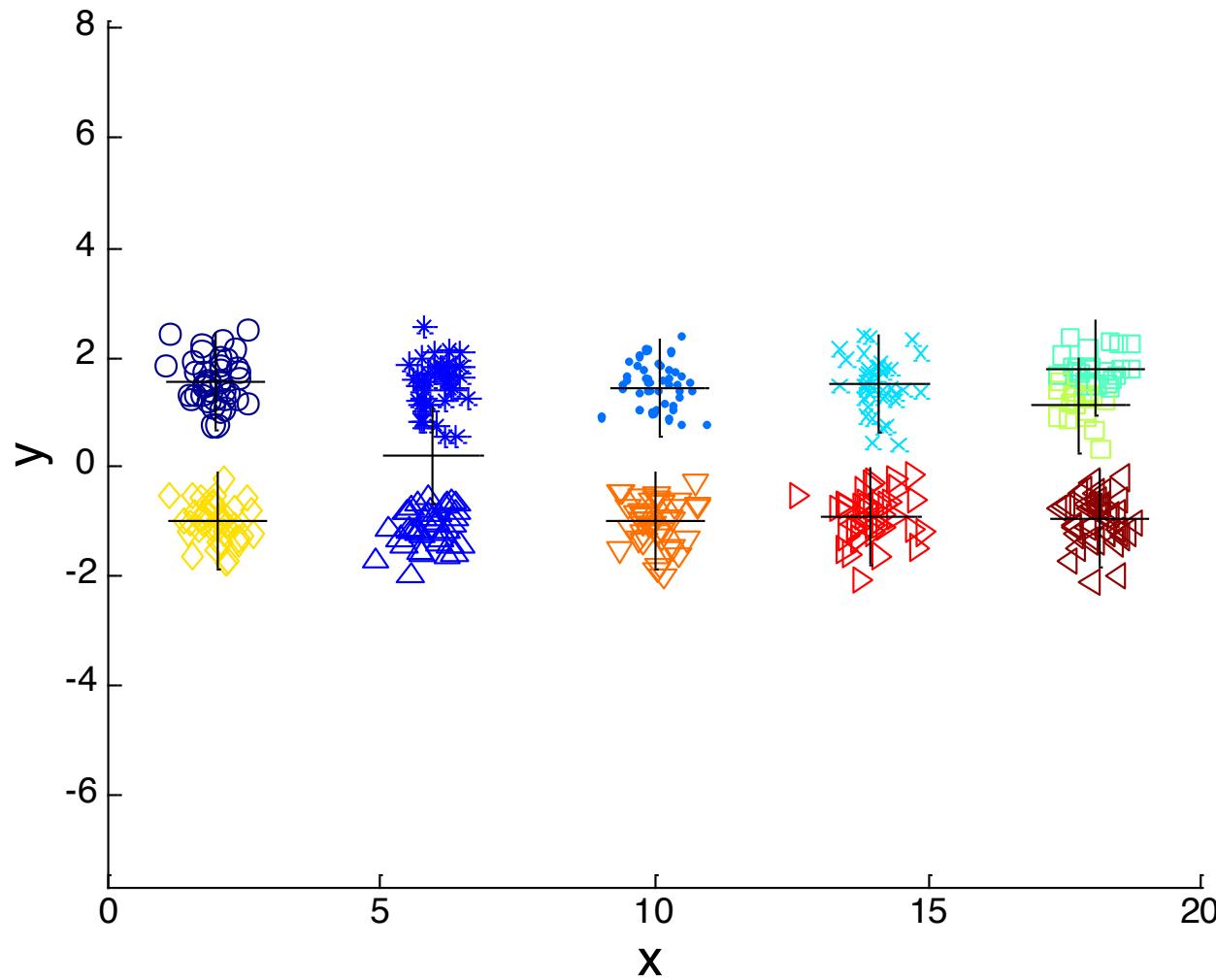
- If there are K ‘real’ clusters then the chance of selecting one centroid from each cluster is small.
- Chance is relatively small when K is large
- If clusters are the same size, n , then

$$P = \frac{\text{number of ways to select one centroid from each cluster}}{\text{number of ways to select } K \text{ centroids}} = \frac{K!n^K}{(Kn)^K} = \frac{K!}{K^K}$$

- For example, if $K = 10$, then probability = $10!/10^{10} = 0.00036$
- Sometimes the initial centroids will readjust themselves in ‘right’ way, and sometimes they don’t
- Consider an example of five pairs of clusters



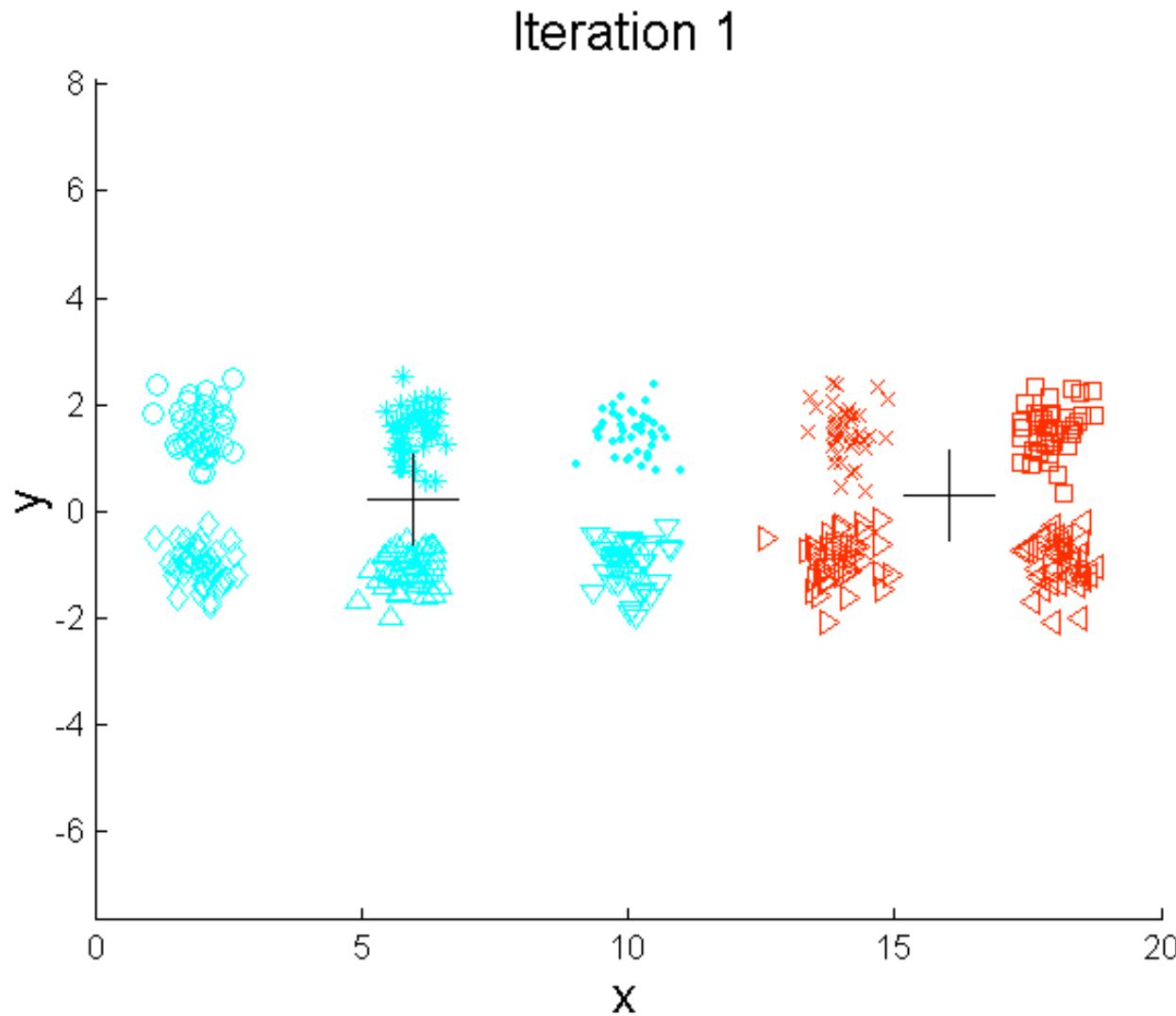
Starting with two initial centroids in one cluster of each pair of clusters

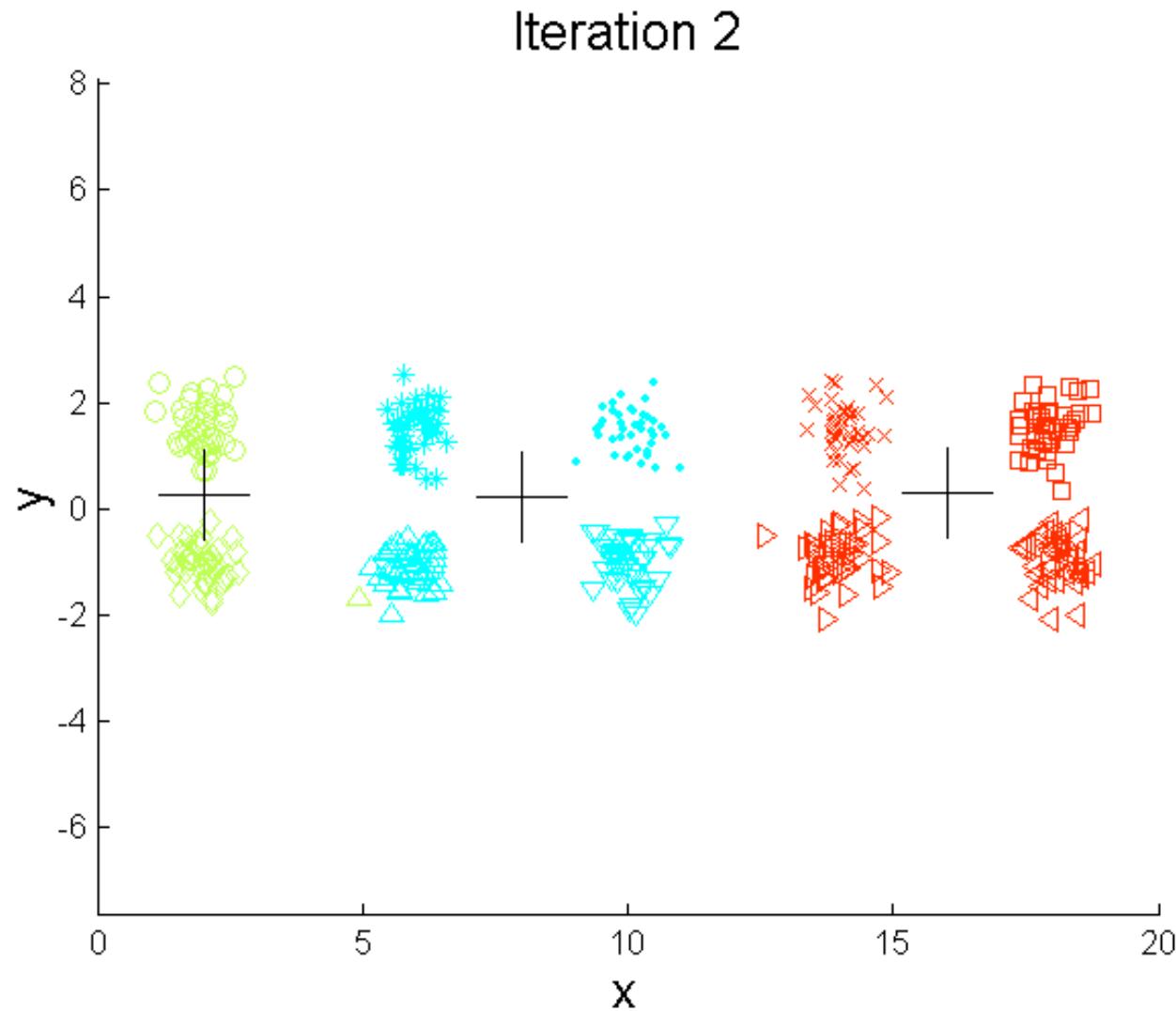


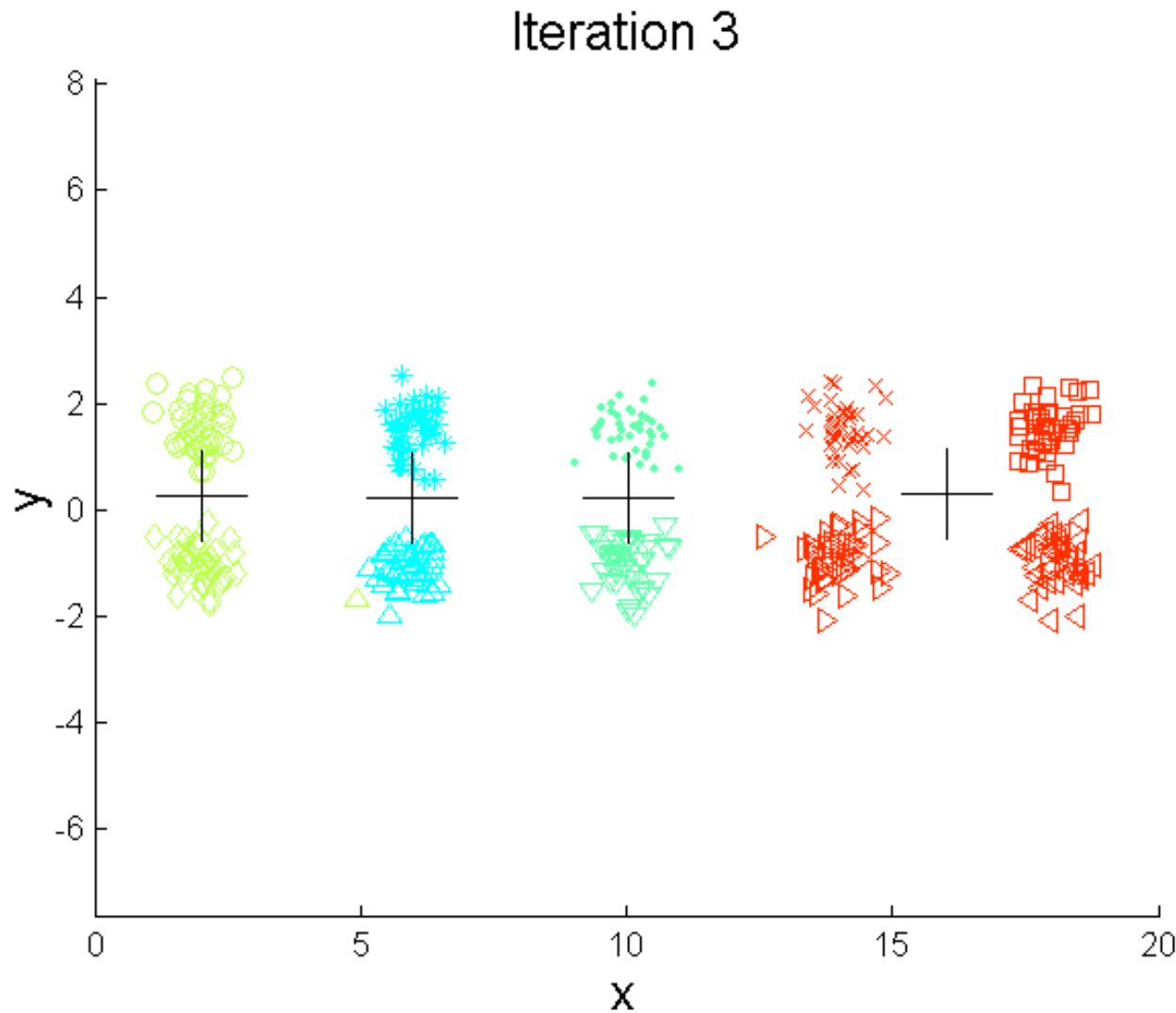
Starting with some pairs of clusters having three initial centroids, while other have only one.

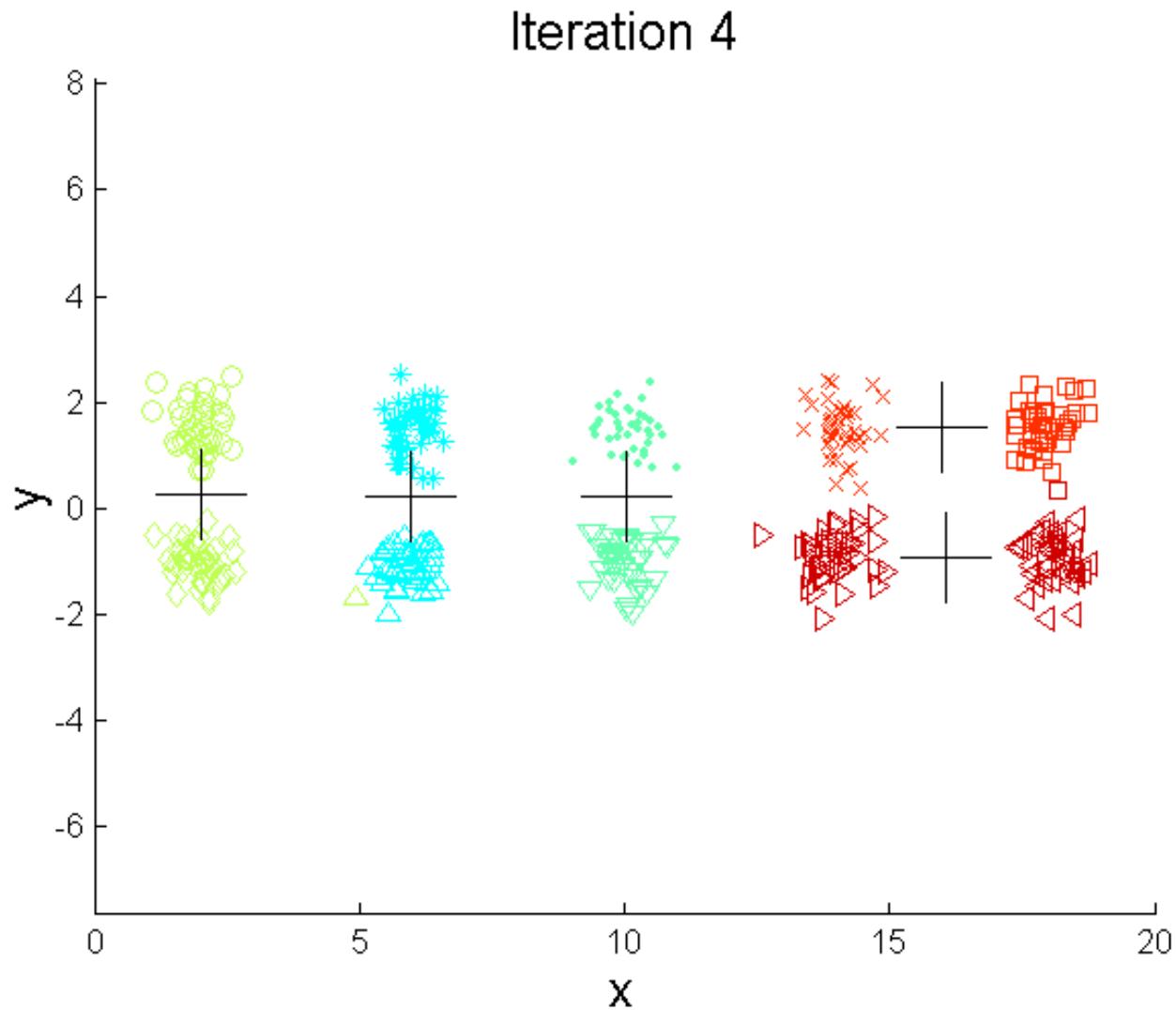
Algorithm 3 Bisecting K-means Algorithm.

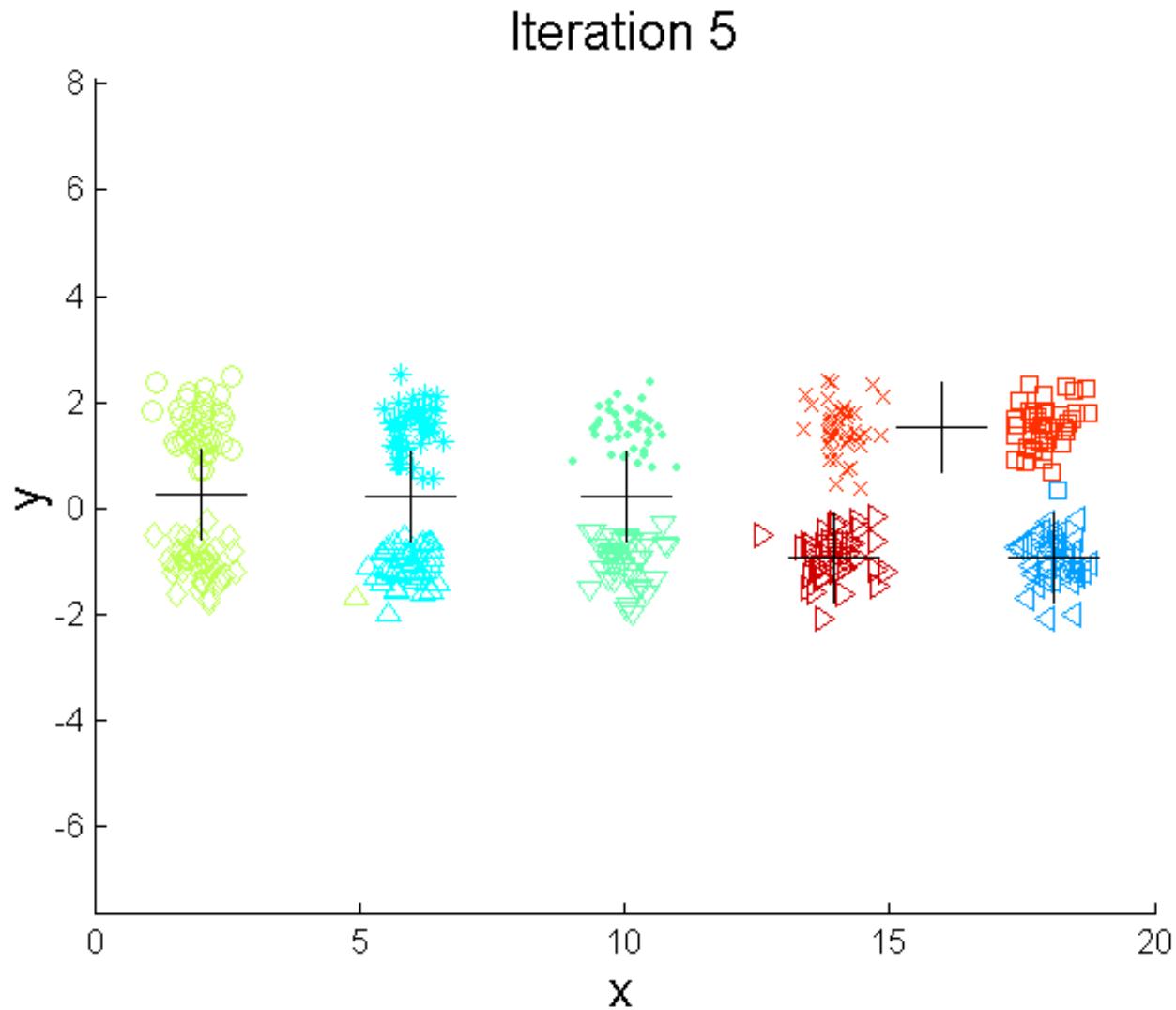
- 1: Initialize the list of clusters to contain the cluster containing all points.
 - 2: **repeat**
 - 3: Select a cluster from the list of clusters
 - 4: **for** $i = 1$ to *number_of_iterations* **do**
 - 5: Bisect the selected cluster using basic K-means
 - 6: **end for**
 - 7: Add the two clusters from the bisection with the lowest SSE to the list of clusters.
 - 8: **until** Until the list of clusters contains K clusters
-

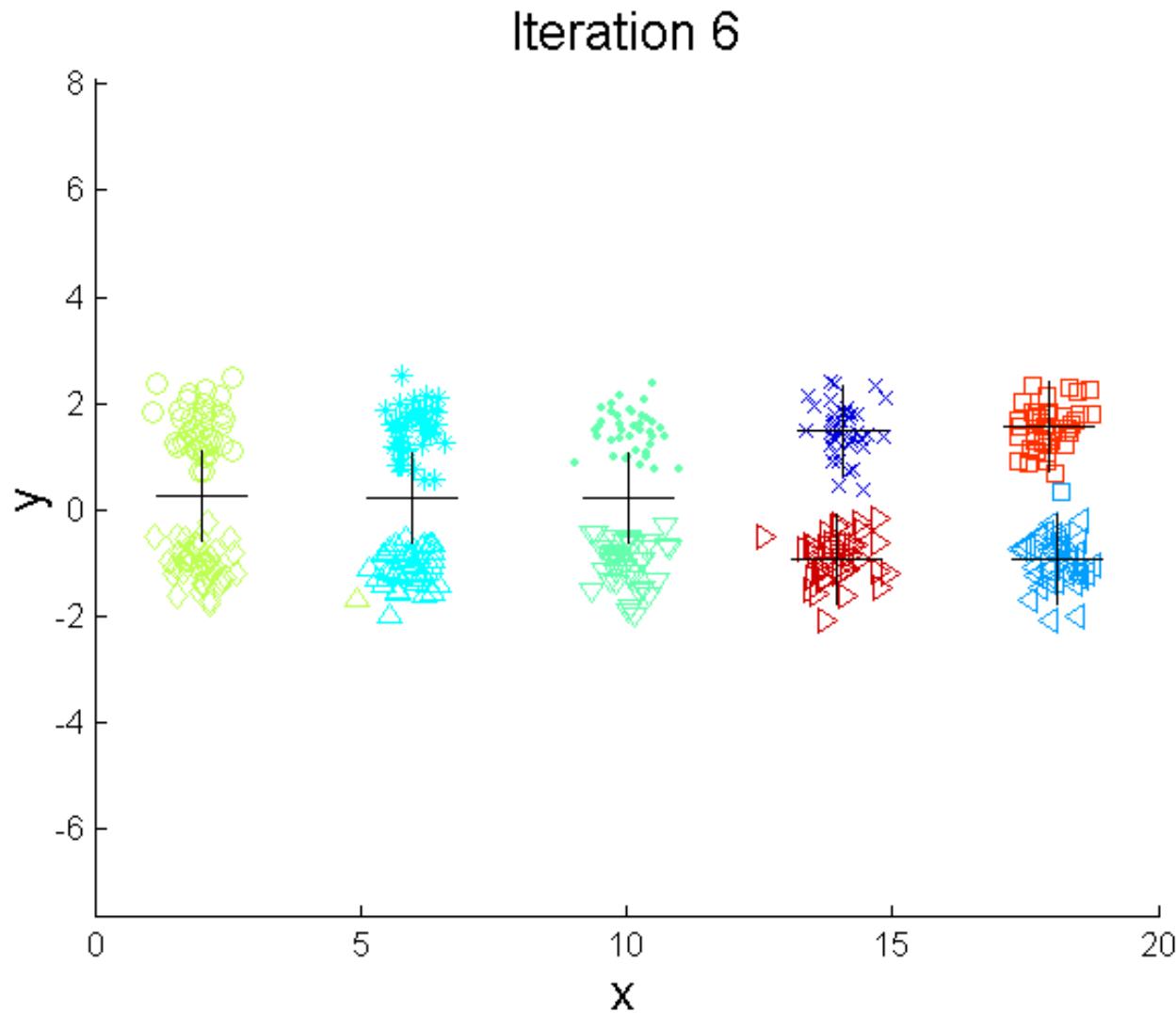


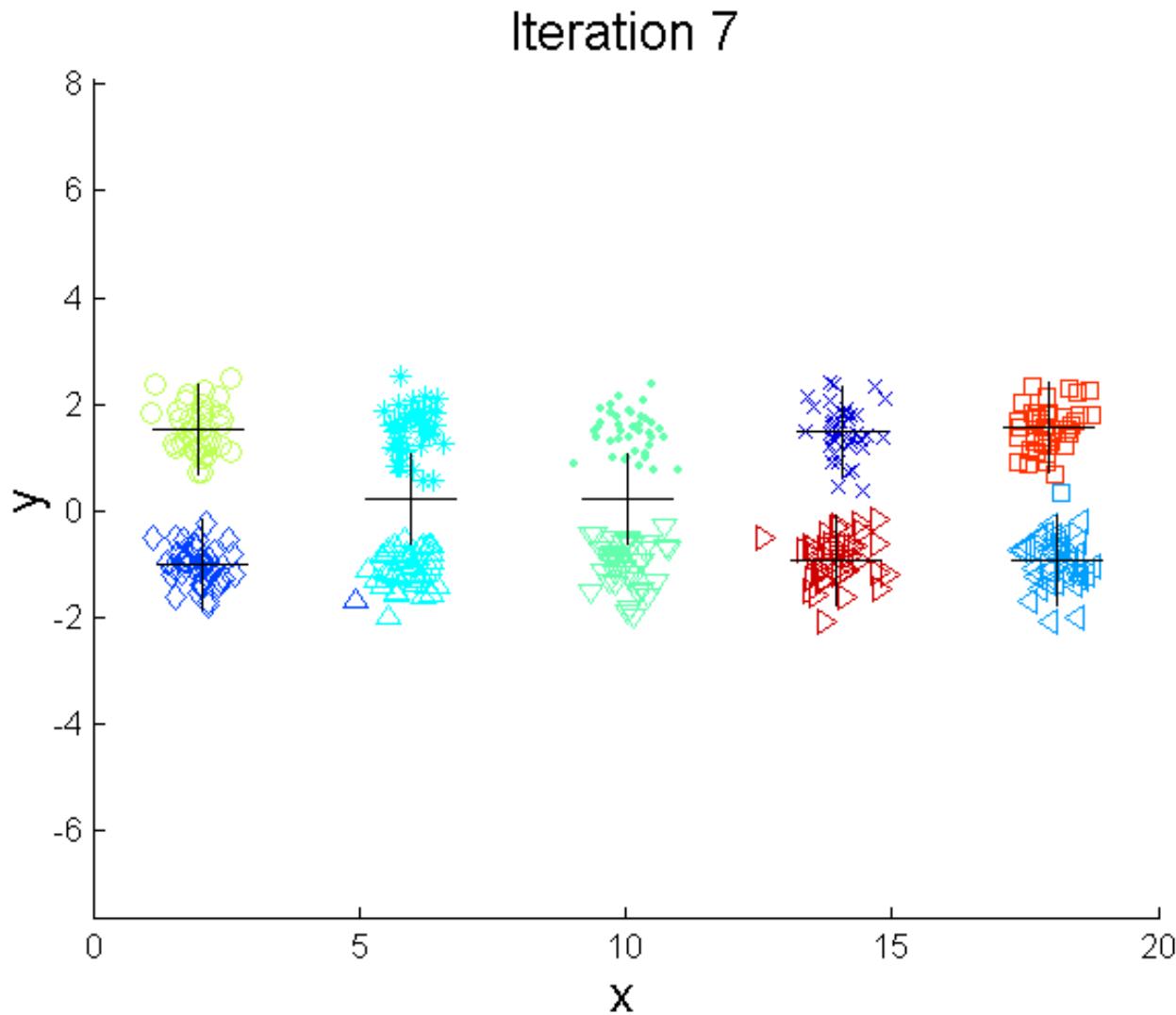


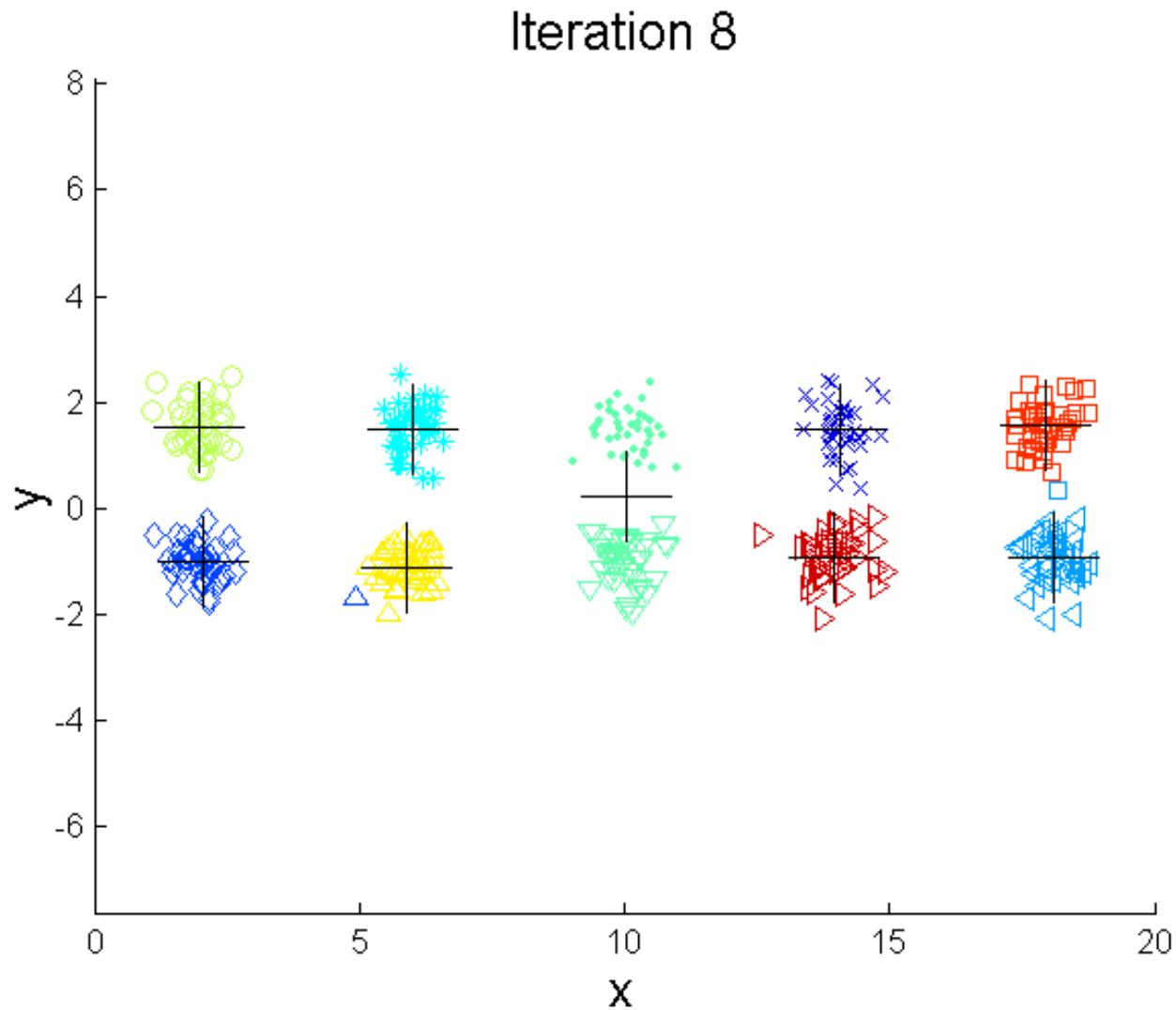


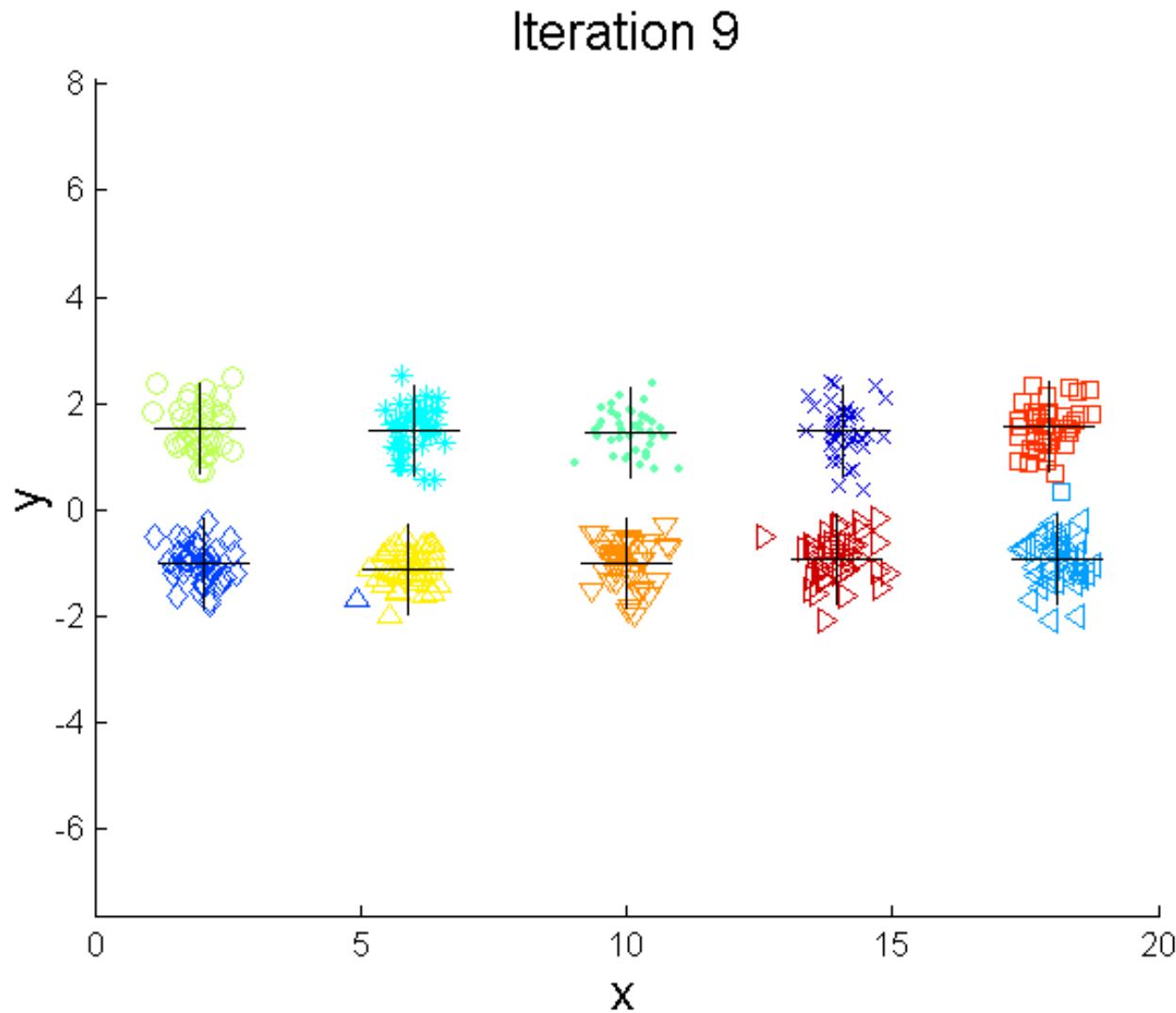


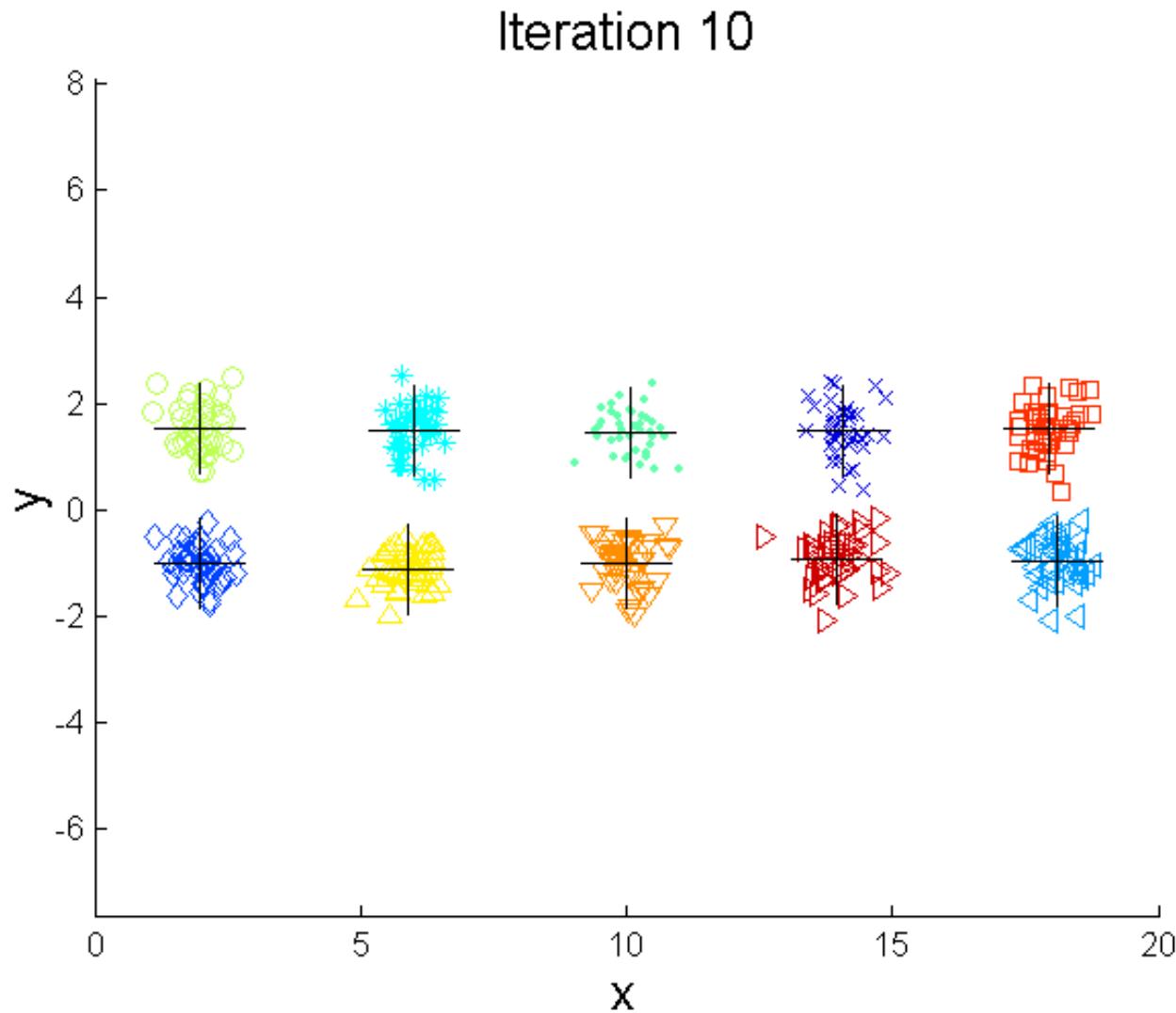












- Multiple runs, helps, but probability is not on your side
- Sample and use another clustering method (hierarchical?) to determine initial centroids
- Select more than k initial centroids and then select among these initial centroids
- Postprocessing
- Bisecting K-means, not as susceptible to initialization issues

Extensions

- In the basic K-means algorithm, centroids are updated after all points are assigned to a centroid
- An alternative is to update the centroids after each assignment (incremental approach)
 - Each assignment updates zero or two centroids
 - More expensive
 - Introduces an order dependency
 - Never get an empty cluster
 - Can use “weights” to change the impact

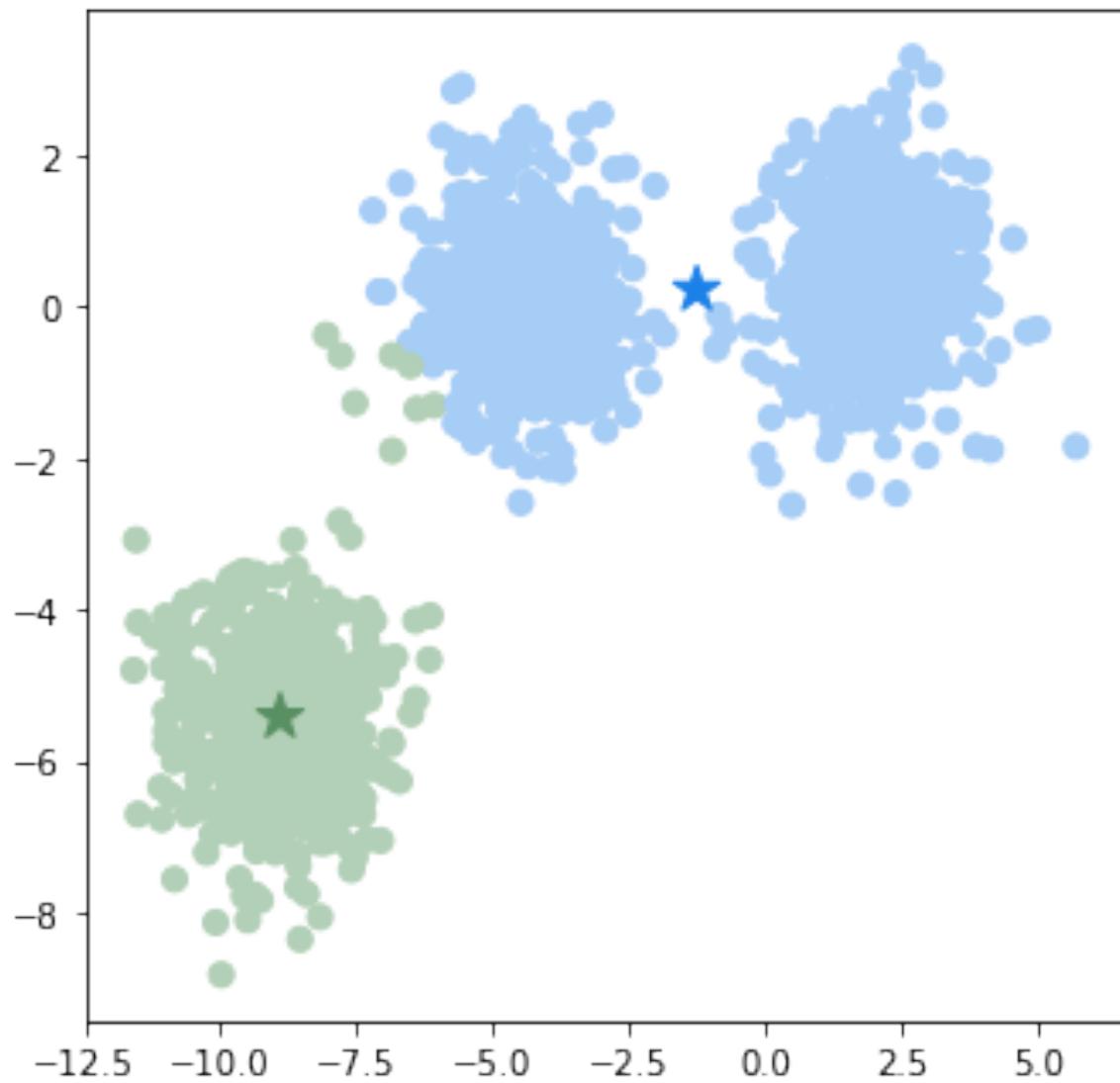
- Pre-processing
 - Normalize the data
 - Eliminate outliers
- Post-processing
 - Eliminate small clusters that may represent outliers
 - Split ‘loose’ clusters, i.e., clusters with relatively high SSE
 - Merge clusters that are ‘close’ and that have relatively low SSE
 - These steps can be used during the clustering process

- A few variants of the k-means which differ in
 - Selection of the initial k means
 - Dissimilarity calculations
 - Strategies to calculate cluster means
- Handling categorical data: k-modes
 - Replacing means of clusters with modes
 - Using new dissimilarity measures to deal with categorical objects
 - Using a frequency-based method to update modes of clusters
 - A mixture of categorical and numerical data: k-prototype method

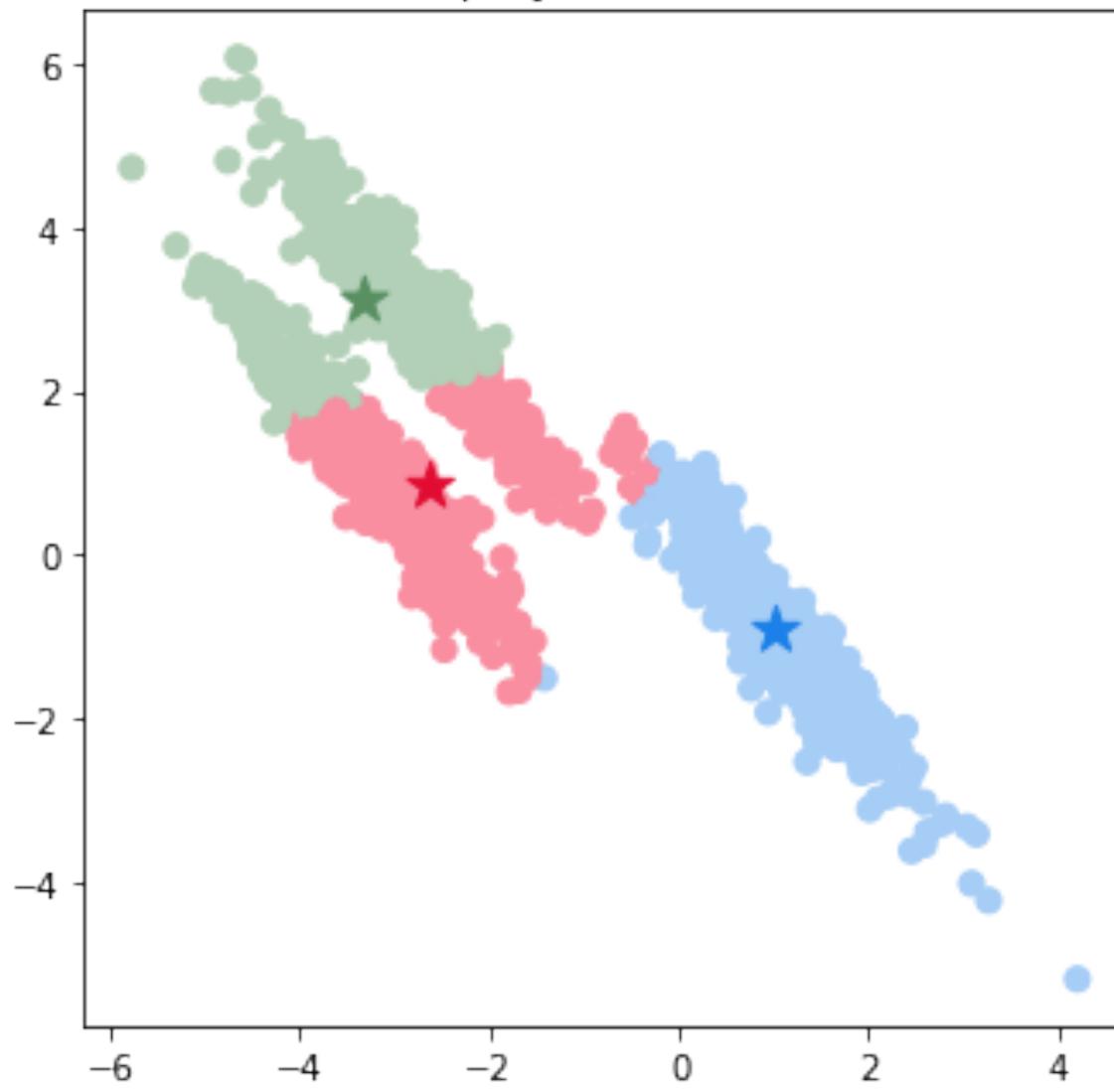
Limitations

- K-means has problems when clusters are of differing
 - Sizes
 - Densities
 - Non-globular shapes
- K-means has also problems when the data contains outliers.

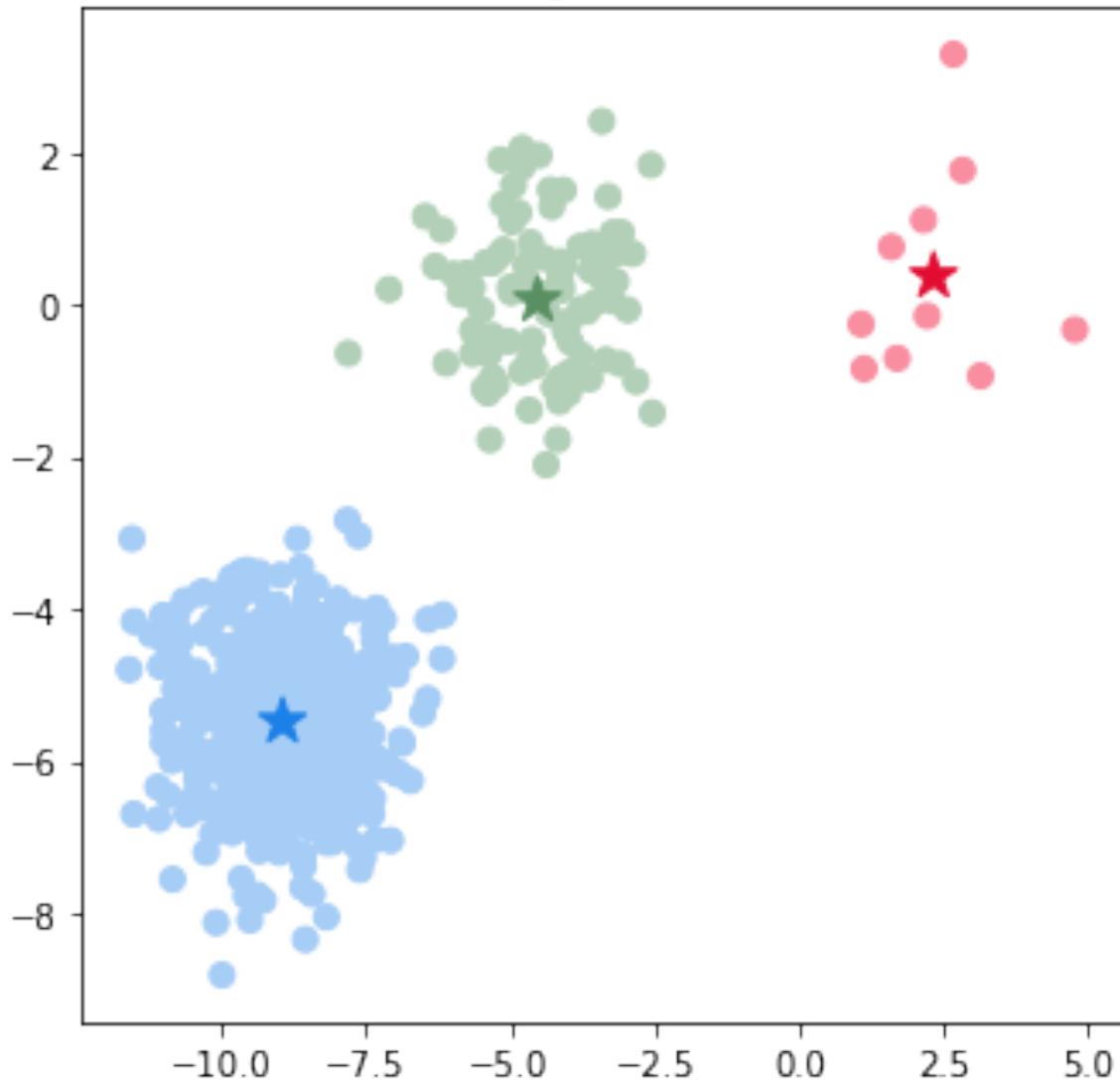
Incorrect Number of Blobs

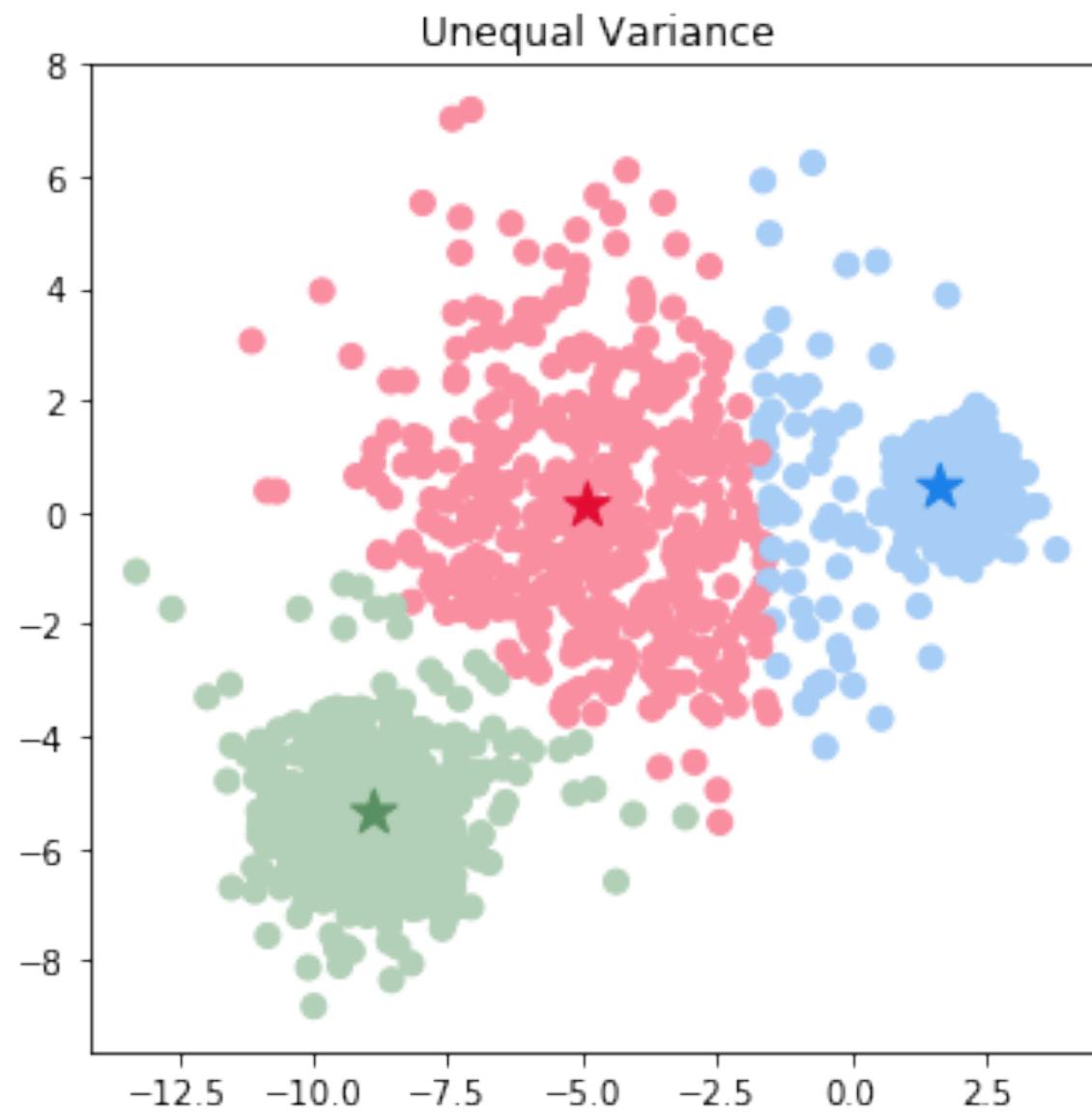


Anisotropicly Distributed Blobs

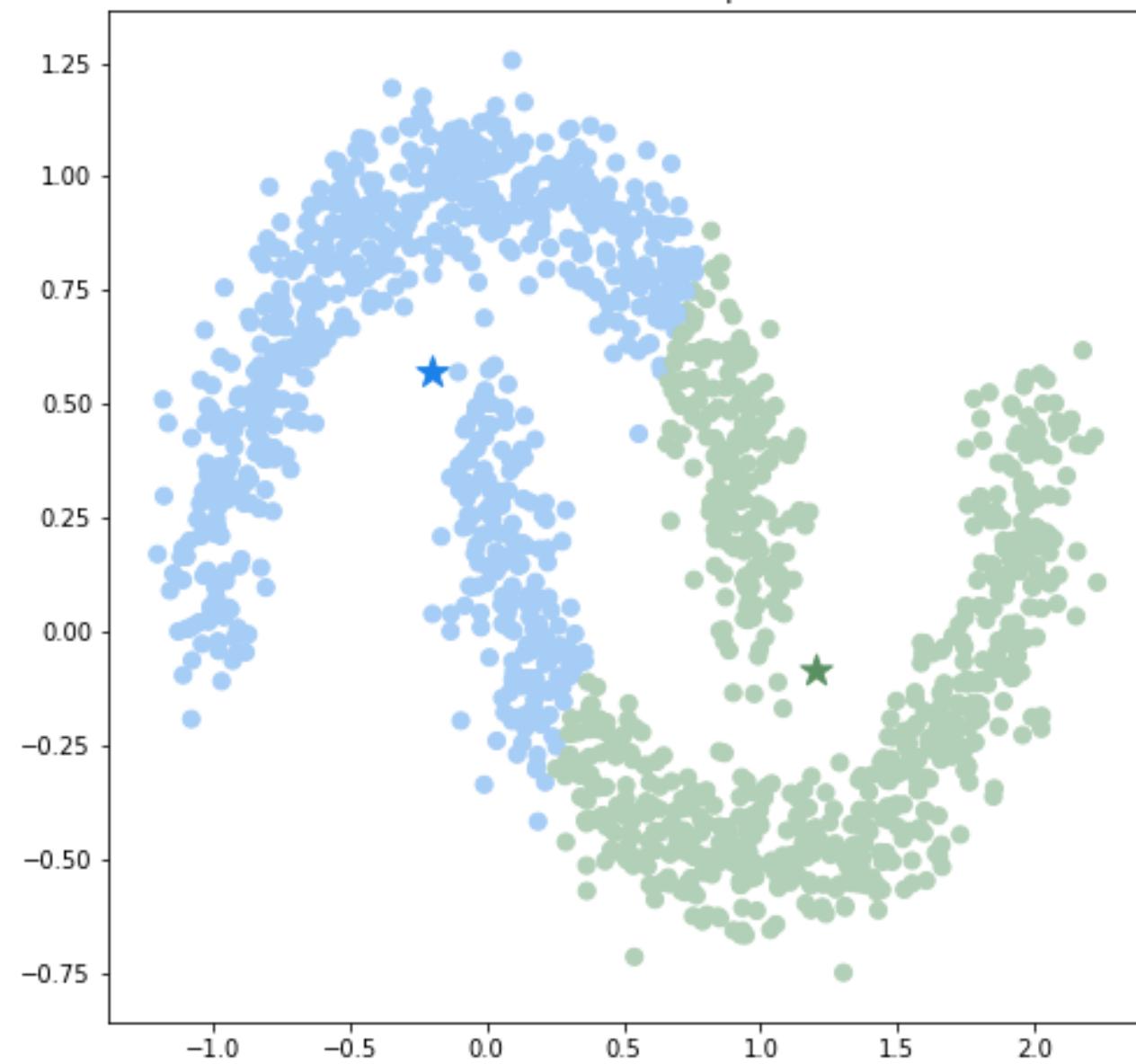


Unevenly Sized Blobs

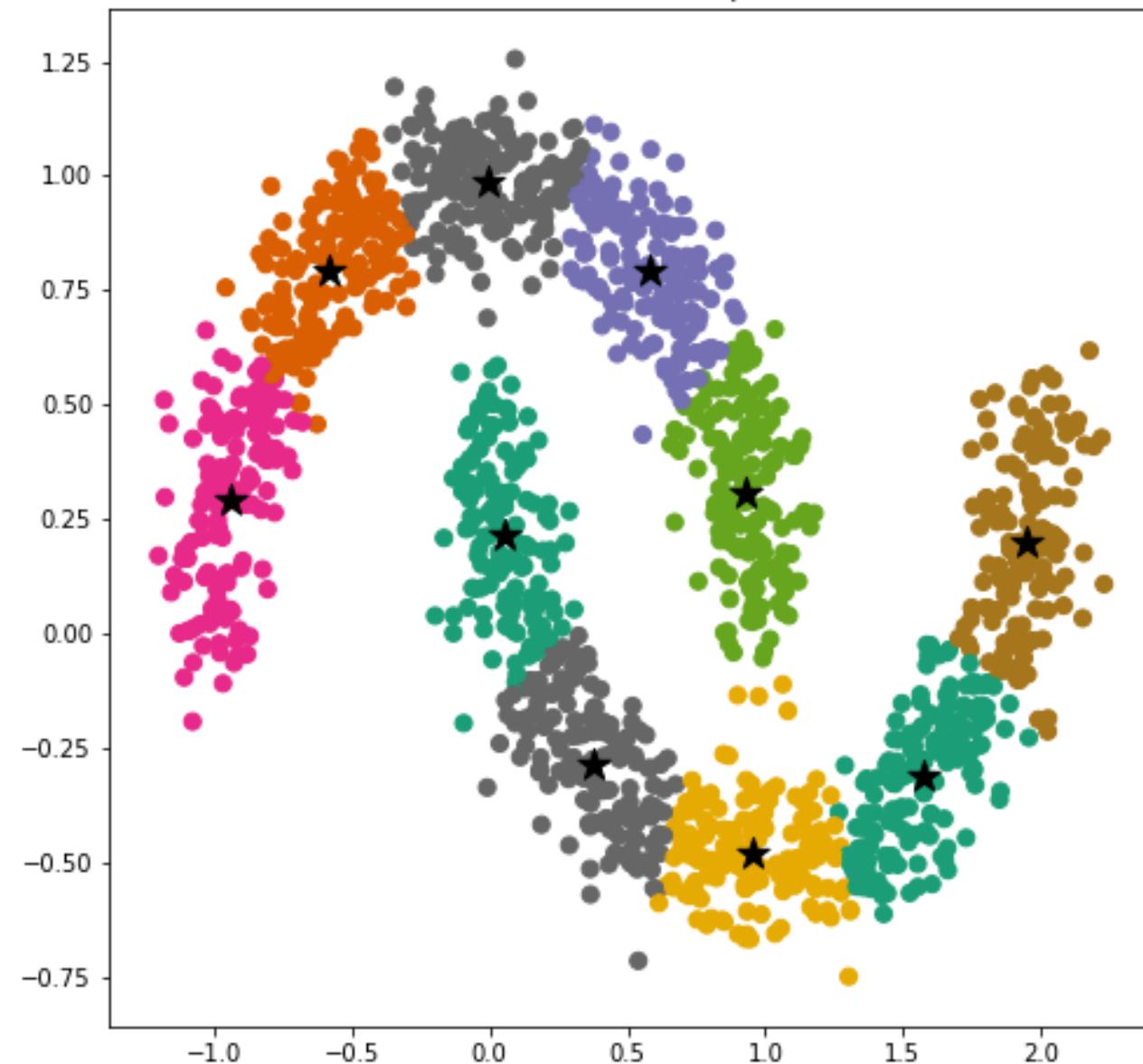




Non Globular Shapes



Non Globular Shapes

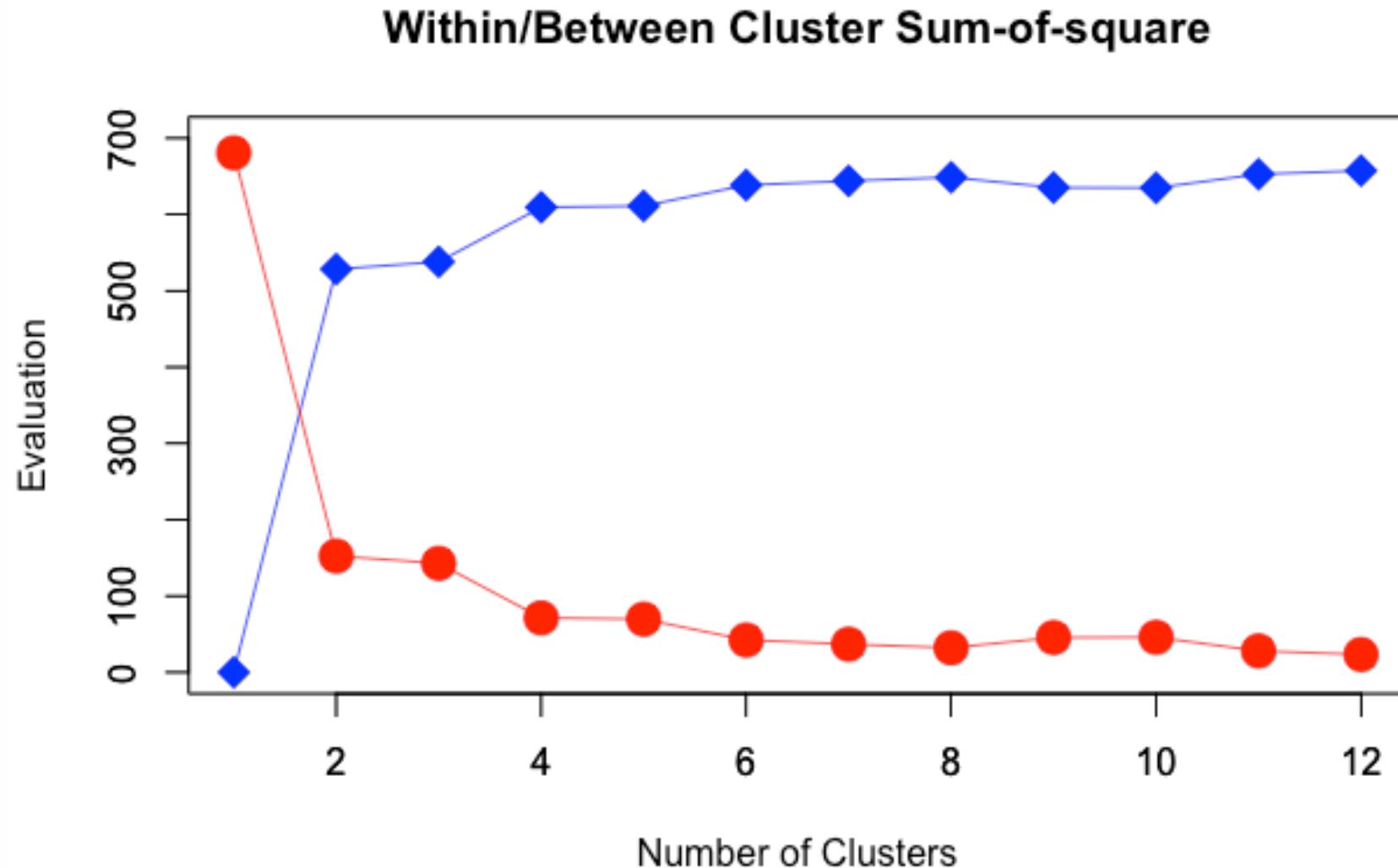


We can use higher values of k and then preprocess the resulting clusters and merge them when we deem it useful.

- **Strength**
 - Relatively efficient
 - 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, then what about categorical data?
 - Need to specify k, the number of clusters, in advance
 - Unable to handle noisy data and outliers
 - Not suitable to discover clusters with non-convex shapes
- **Advantages**
 - Simple, understandable
 - Items automatically assigned to clusters
- **Disadvantages**
 - Must pick number of clusters before hand
 - All items forced into a cluster
 - Too sensitive to outliers

- A few variants of the k-means which differ in
 - Selection of the initial k means
 - Dissimilarity calculations
 - Strategies to calculate cluster means
- Handling categorical data: k-modes
 - Replacing means of clusters with modes
 - Using new dissimilarity measures to deal with categorical objects
 - Using a frequency-based method to update modes of clusters
 - A mixture of categorical and numerical data: k-prototype method

How do we choose k ?



Mean Shift Clustering

Mean shift is an iterative, nonparametric,
and versatile algorithm

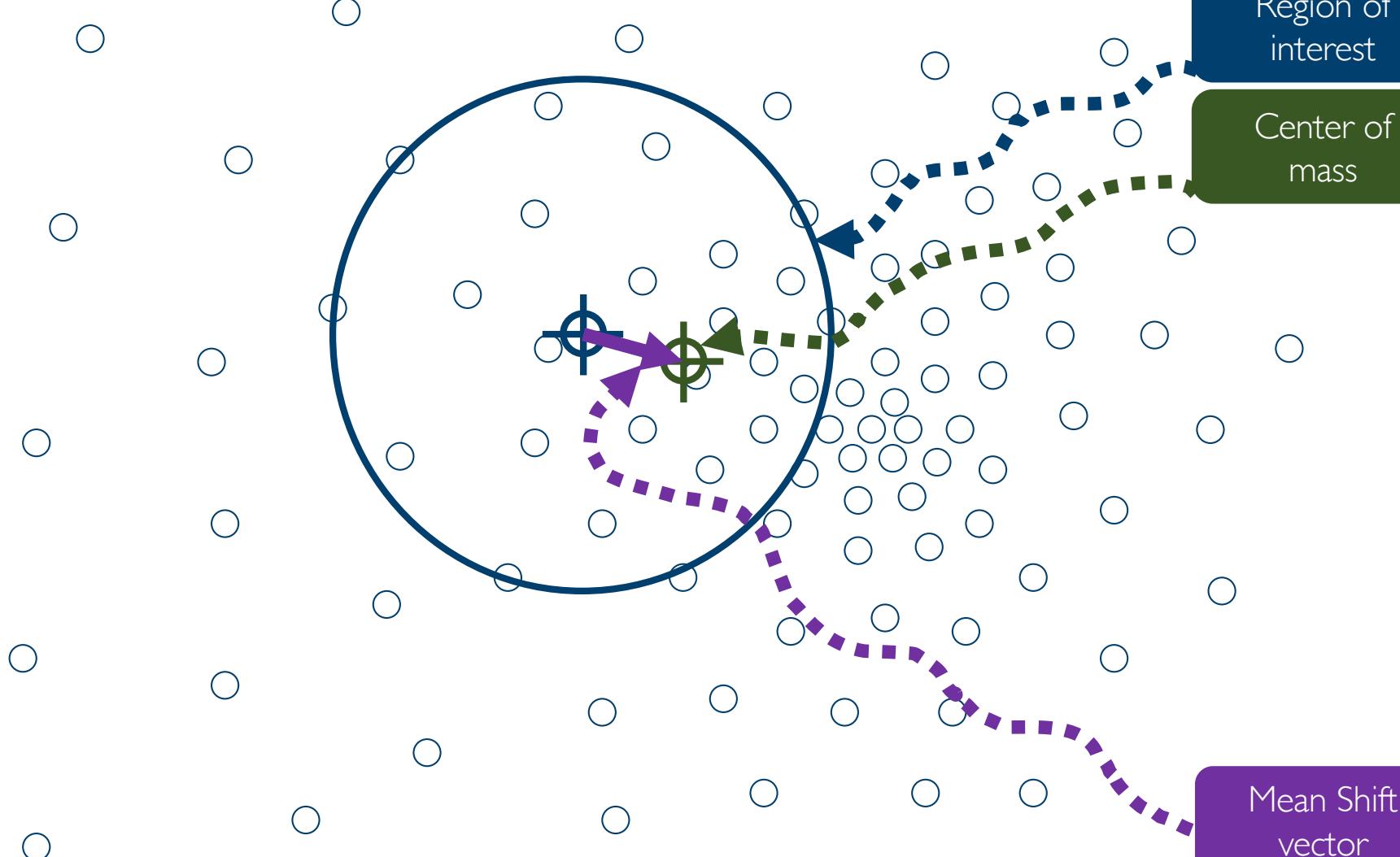
It searches for the mode (i.e., the point of
highest density) of a data distribution

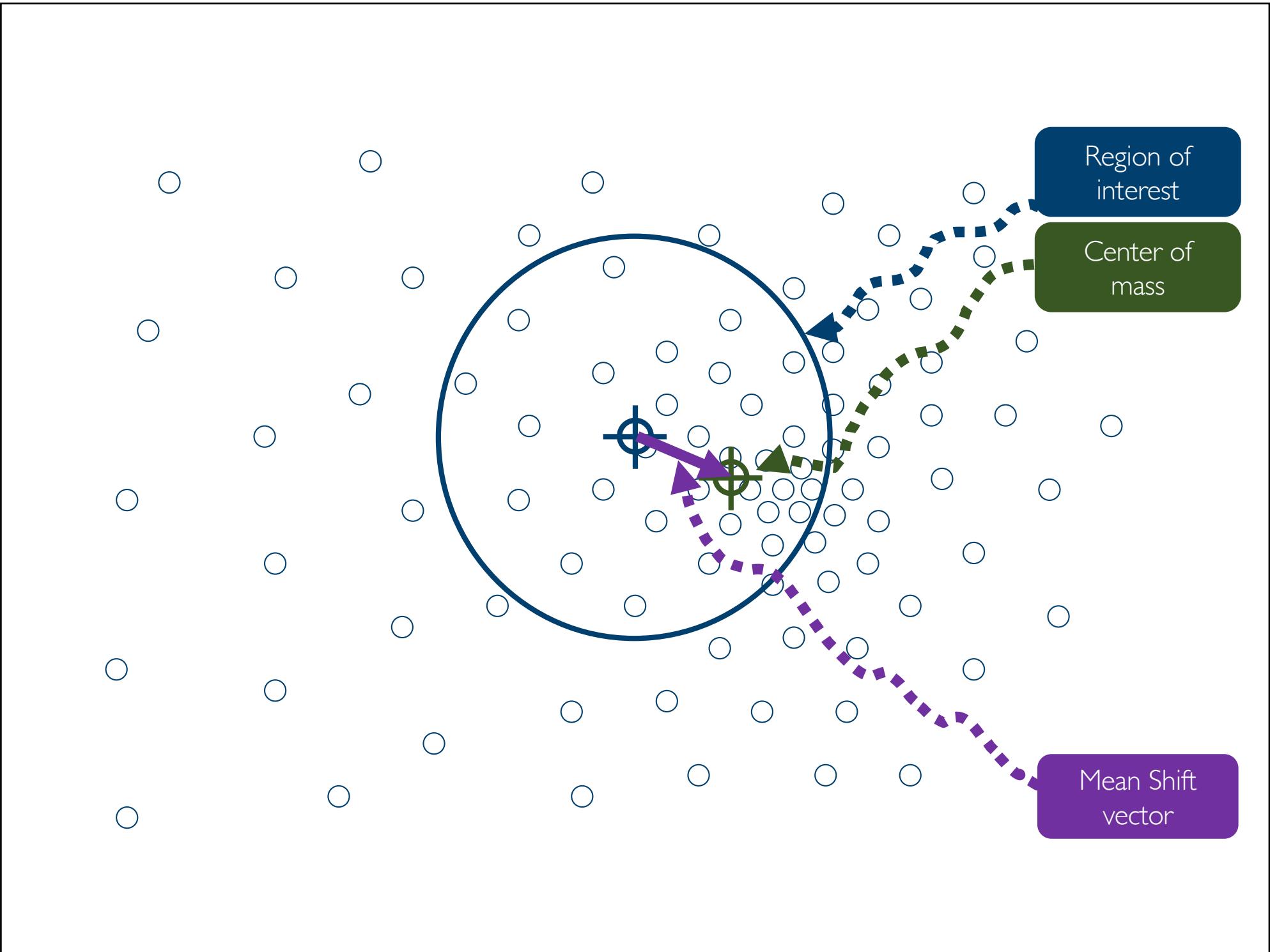
Underlying intuition

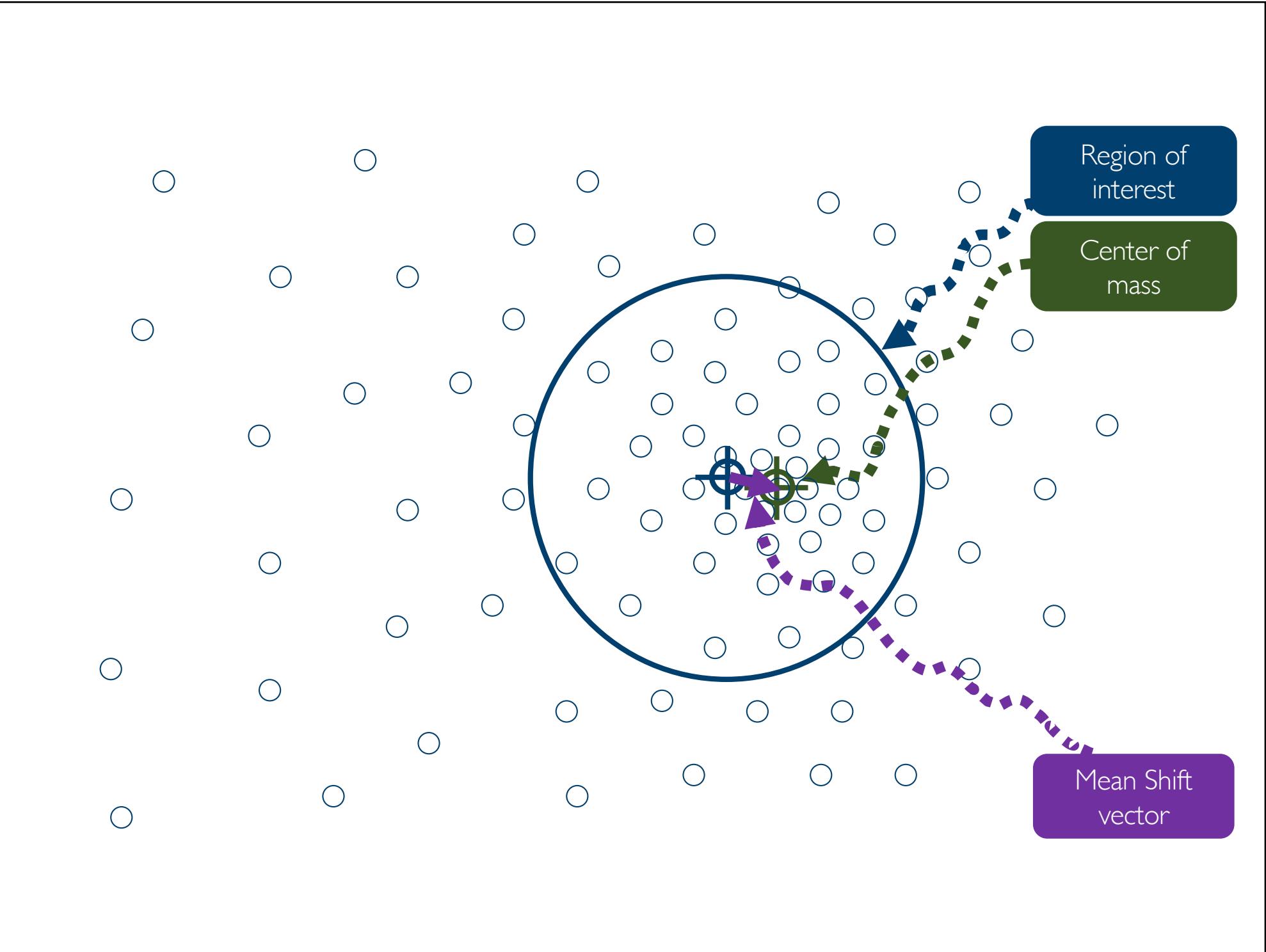


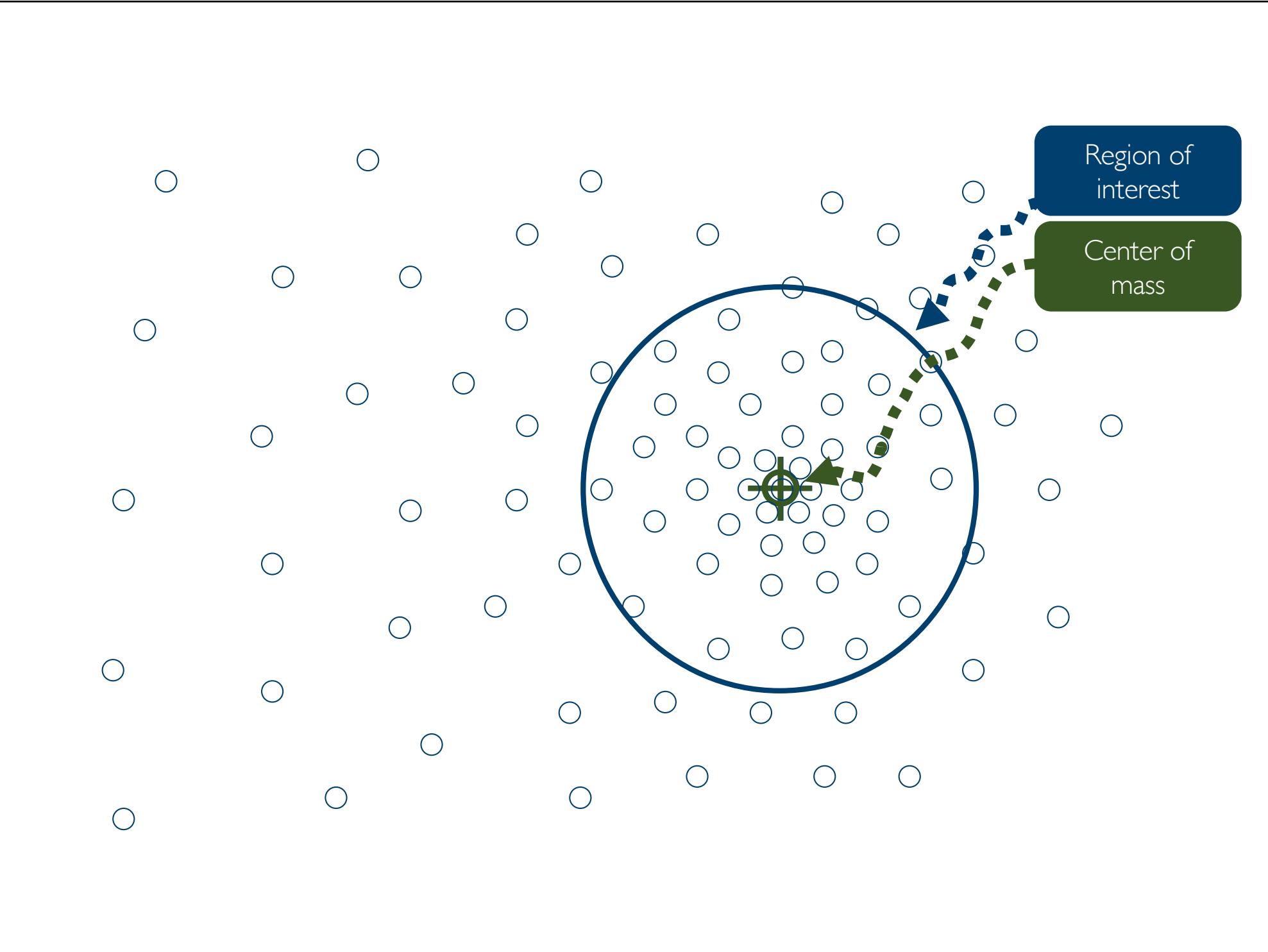
The objective is to find the densest region

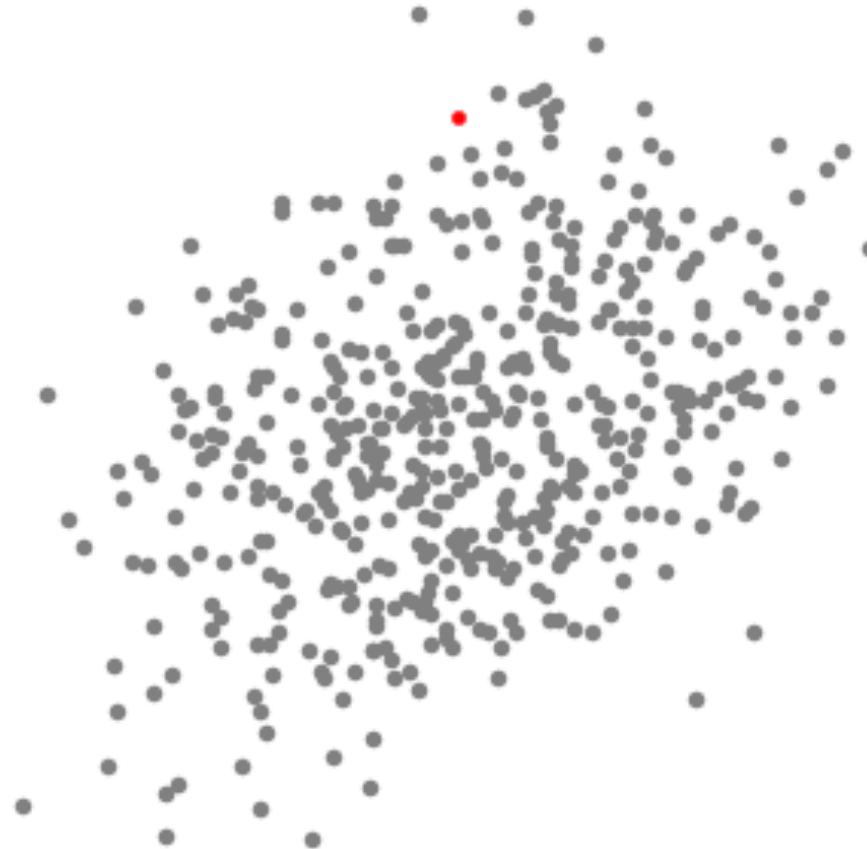
A scatter plot showing a collection of data points represented by blue circles with black outlines. The points are distributed across a white background, with a higher density of points forming a central cluster. The points are scattered more sparsely towards the edges of the plot area.









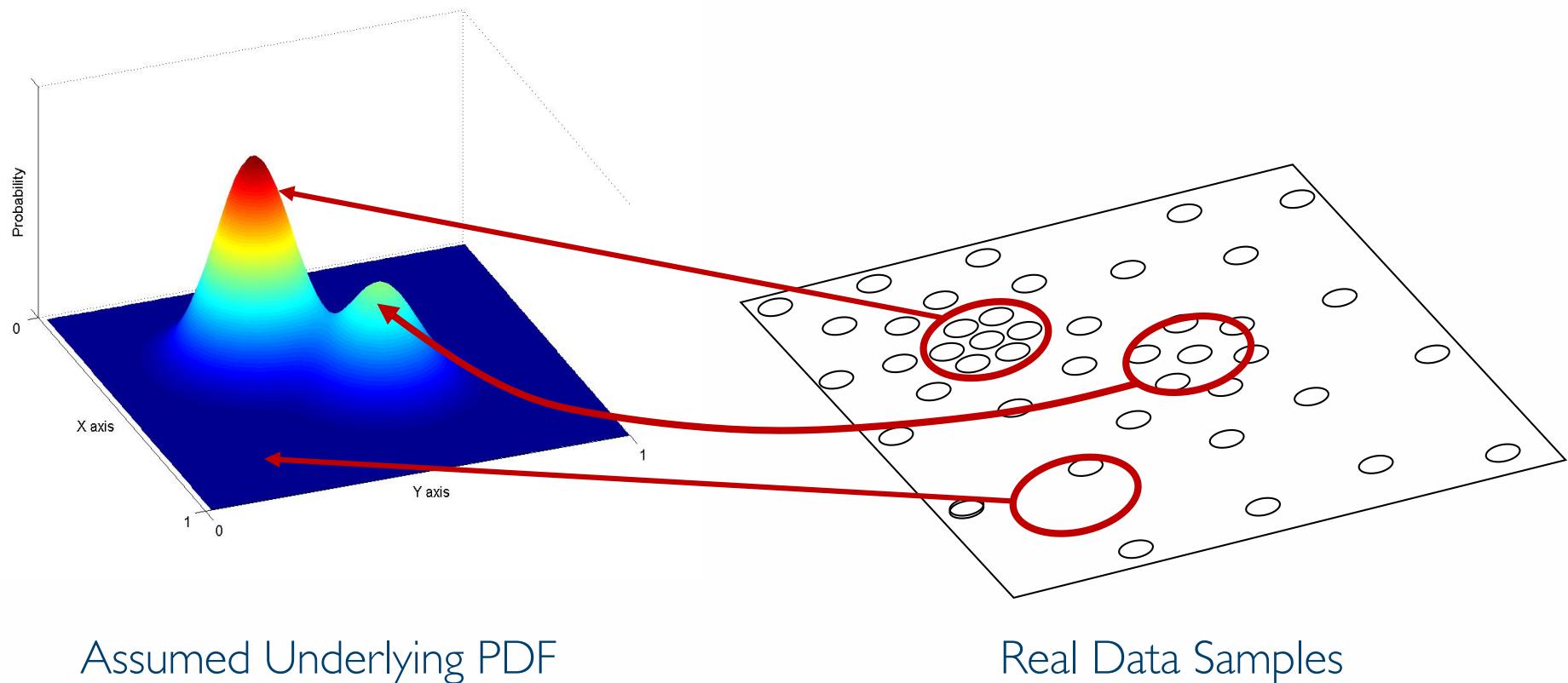


Clustering with Scikit with GIFs

<https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/>

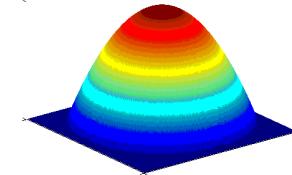
1. Choose a search window size (the bandwidth)
2. Choose the initial location of the search window
3. Compute the mean location in the search window
4. Center the search window at the location computed in Step 3
5. Repeat Steps 3 and 4 until convergence

- We assume that the data points are sampled from an underlying probability density function (PDF)



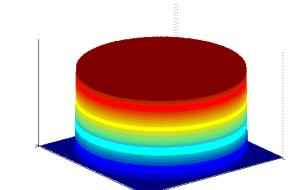
- Epanechnikov Kernel

$$K_E(x) = \begin{cases} c(1 - ||x||^2) & ||x|| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$



- Uniform Kernel

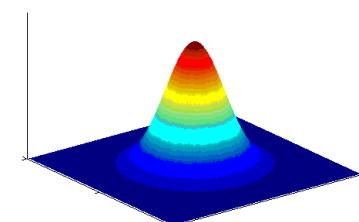
$$K_U(x) = \begin{cases} c & ||x|| \leq 1 \\ 0 & \text{otherwise} \end{cases}$$



- Normal/Gaussian Kernel

$$K_N(x) = c \cdot e^{-0.5||x||^2}$$

$$K_G\left(\frac{x - x_i}{h}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(x - x_i)^2}{2h^2}}$$



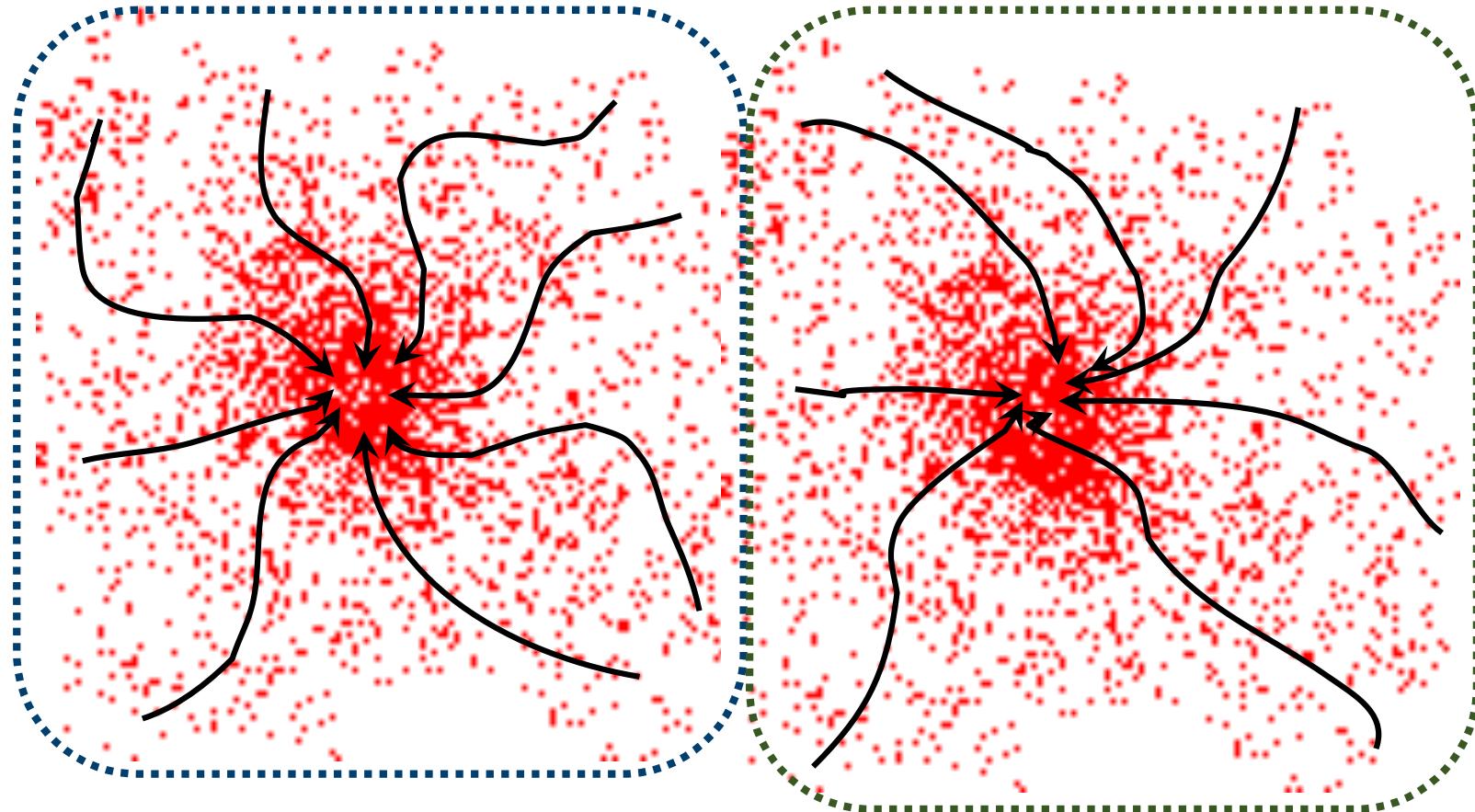
Mean Shift Segmentation

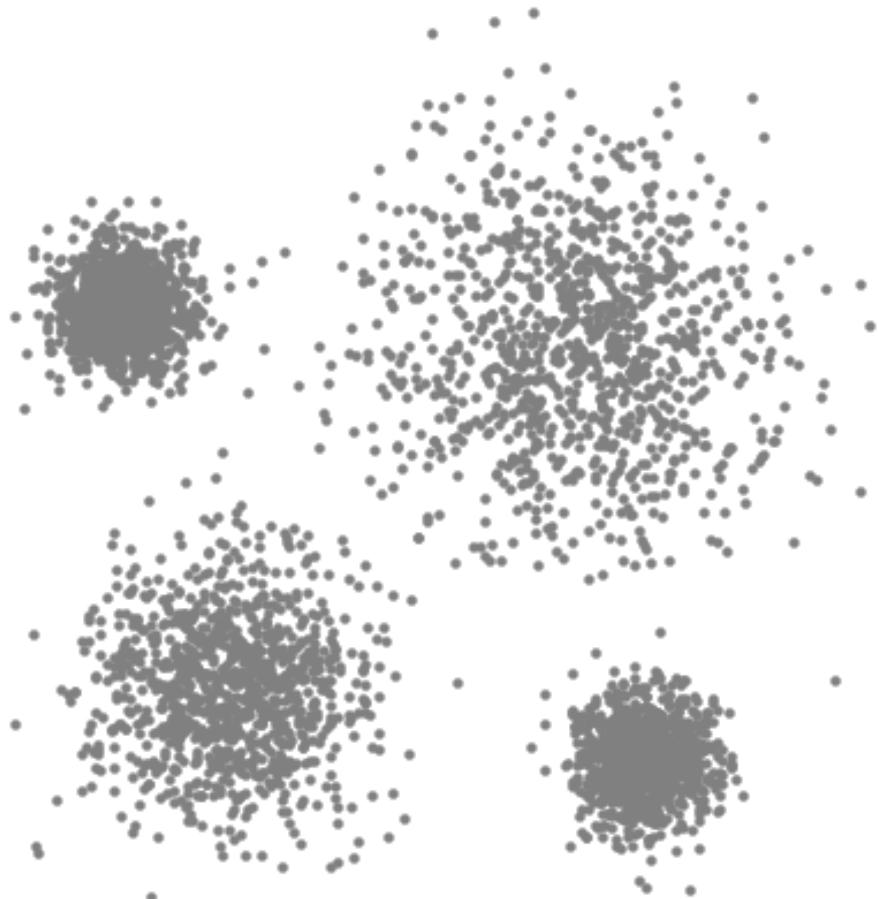


D. Comaniciu and P. Meer, Mean Shift: A Robust Approach
toward Feature Space Analysis, PAMI 2002.

Mean Shift Clustering

Trajectories that lead to the same mode
(attraction basin) should be in the same cluster





Clustering with Scikit with GIFs

<https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/>

1. Choose kernel and bandwidth
2. For each point
 - a. Center a window on that point
 - b. Compute the mean of the data in the search window
 - c. Center the search window at the new mean location
 - d. Repeat (b,c) until convergence
3. Assign points that lead to nearby modes to the same cluster

- **Number of clusters and shape**
 - k-means assumes that the number of clusters is already known and the clusters are shaped spherically (or elliptically)
 - Mean shift does not assume anything about number of clusters. Instead, the number of modes give the number of clusters
 - Also, since mean shift is based on density estimation, it can handle arbitrarily shaped clusters.
- **Initialization and outliers**
 - k-means is sensitive to initialization
 - Mean shift is robust to initializations. Typically, mean shift is run for each point or sometimes points are selected uniformly from the feature space
 - Similarly, k-means is sensitive to outliers while Mean Shift is less sensitive.
- **Time complexity**
 - k-means is fast and has a time complexity $O(knT)$ where k is the number of clusters, n is the number of points and T is the number of iterations.
 - Classic mean shift is computationally expensive with a time complexity $O(Tn^2)$.

- **Strength**

- Does not assume any prior shape on data clusters
- Can handle arbitrary feature spaces
- Only the window size h to choose
- The window size h (or bandwidth) has a physical meaning

- **Weaknesses**

- The window size h (or bandwidth) selection is not trivial
- Inappropriate window size can cause modes to be merged, or generate additional “shallow” modes
- Adaptive window size can help
- Not suitable for high-dimensional features

- MeanShift clustering runs at $O(Tn^2)$, k-Means at $O(knT)$ where T is the number of iterations and n is the number of data points
- The functions to estimate the bandwidth scale particularly badly (see sklearn estimate_bandwidth documentation)

Expectation Maximization

- k-means assigns each point to only one cluster (hard assignment)
- The approach can be extended to consider soft assignment of points to clusters, so that each point has a probability of belonging to each cluster
- We assume that each cluster C_i is characterized by a multivariate normal distribution and thus identified by
 - The mean vector μ_i
 - The covariance matrix Σ_i
- A clustering is identified by a vector of parameter θ defined as

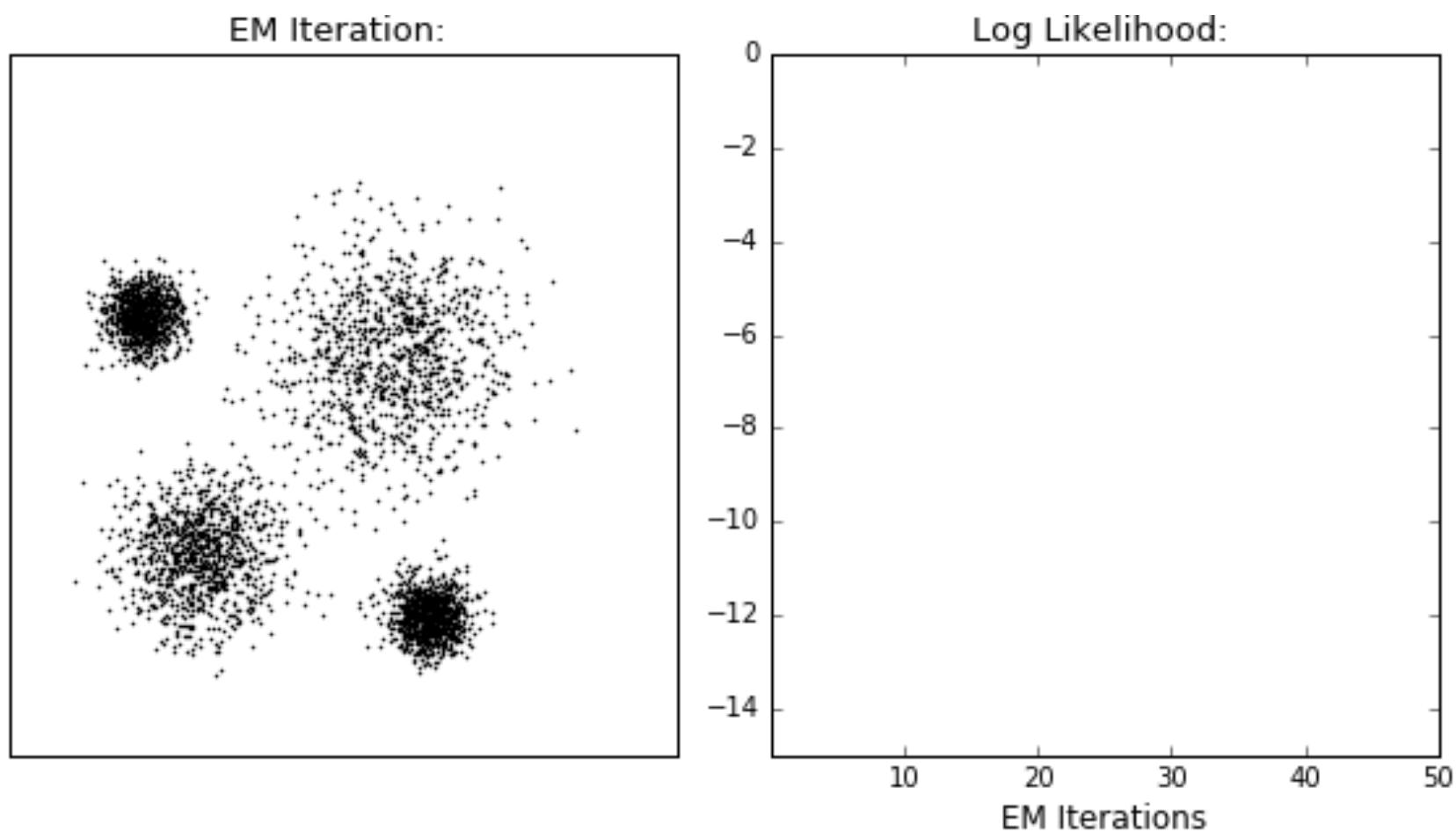
$$\theta = \{\mu_i \ \Sigma_i \ P(C_i)\}$$

where $P(C_i)$ are the prior probability of all the clusters C_i which sum up to one

- The goal of maximum likelihood estimation (MLE) is to choose the parameters θ that maximize the likelihood, that is

$$\theta^* = \arg \max_{\theta} P(D|\theta)$$

- General idea
 - Starts with an initial estimate of the parameter vector
 - Iteratively rescores the patterns against the mixture density produced by the parameter vector
 - The rescored patterns are used to update the parameter updates
 - Patterns belonging to the same cluster, if they are placed by their scores in a particular component



Clustering with Scikit with GIFs
<https://dashee87.github.io/data%20science/general/Clustering-with-Scikit-with-GIFs/>

Example in One Dimension

- Consider a dataset consisting of a single attribute $X = \{x_1, \dots, x_n\}$
- Clusters will be represented using univariate normal,

$$f_i(x) = f(x|\mu_i, \sigma_i^2) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp \left\{ -\frac{(x - \mu_i)^2}{2\sigma_i^2} \right\}$$

- Each cluster will be represented by the parameters $\{\mu_i, \sigma_i^2, P(C_i)\}$
- We initialize the parameters as follows, the mean μ_i is selected uniformly at random, σ_i^2 is initialized to 1, the cluster $P(C_i)$ probability with $1/k$ (where k is the number of clusters).

- For each cluster, $1 \dots k$, we can use the current estimate of the parameters $\{\mu_i, \sigma_i^2, P(C_i)\}$ to compute the posterior probabilities,

$$P(C_i|x_j) = \frac{f(x_j|\mu_i, \sigma_i^2) \cdot P(C_i)}{\sum_{a=1}^k f(x_j|\mu_a, \sigma_a^2) \cdot P(C_a)}$$

- We denote with
 - $w_{i,j}$ the values $P(C_i|x_j)$
 - w_i the weight vector for cluster C_i over all the n points, that is, $(w_{i,1} \dots w_{i,n})^T$

- Given all the posterior probability values w_{ij} , the maximization step computes the maximum likelihood estimates of the cluster parameters $\{\mu_i \ \sigma_i^2 \ P(C_i)\}$
- The means are updated as,

$$\mu_i = \frac{\sum_{j=1}^n w_{ij} \cdot x_j}{\sum_{j=1}^n w_{ij}}$$

which can be rewritten as

$$\mu_i = \frac{\mathbf{w}_i^T \mathbf{X}}{\mathbf{w}_i^T \mathbf{1}}$$

- Similarly, for σ_i^2

$$\sigma_i^2 = \frac{\sum_{j=1}^n w_{ij}(x_j - \mu_i)^2}{\sum_{j=1}^n w_{ij}} \quad \text{or} \quad \sigma_i^2 = \frac{\mathbf{w}_i^T Z_i^s}{\mathbf{w}_i^T \mathbf{1}}$$

where $Z_i = (x_1 - \mu_i, \dots, x_n - \mu_i)^\top$

- Finally,

$$P(C_i) = \frac{\sum_{j=1}^n w_{ij}}{\sum_{a=1}^k \sum_{j=1}^n w_{aj}} = \frac{\sum_{j=1}^n w_{ij}}{\sum_{j=1}^n 1} = \frac{\sum_{j=1}^n w_{ij}}{n}$$

- The expectation and maximization steps are repeated until convergence (e.g. until the means change a little between steps)

Example with one variable
(Example 13.4 from the textbook)

- Data Points

$$x_1 = 1.0 \quad x_2 = 1.3 \quad x_3 = 2.2 \quad x_4 = 2.6 \quad x_5 = 2.8$$

$$x_6 = 5.0 \quad x_7 = 7.3 \quad x_8 = 7.4 \quad x_9 = 7.5 \quad x_{10} = 7.7 \quad x_{11} = 7.9$$

- Initialization of the parameters

$$\mu_1 = 6.63$$

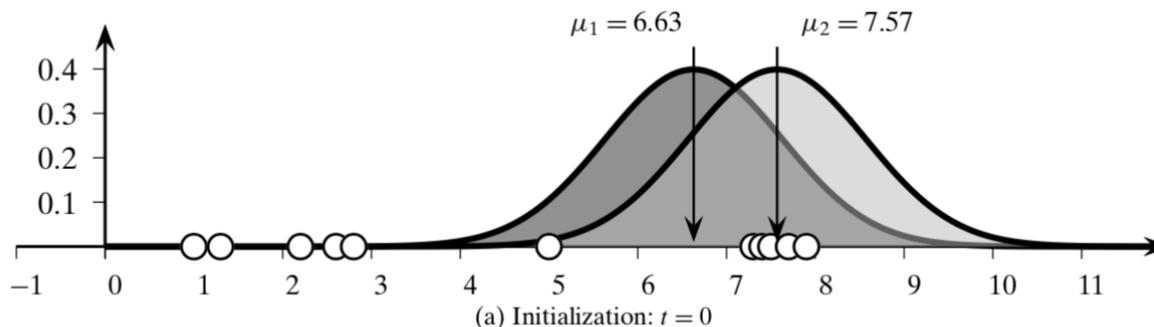
$$\sigma_1^2 = 1$$

$$P(C_2) = 0.5$$

$$\mu_2 = 7.57$$

$$\sigma_2^2 = 1$$

$$P(C_1) = 0.5$$



- After one expectation-maximization iteration the parameters are

$$\mu_1 = 3.72$$

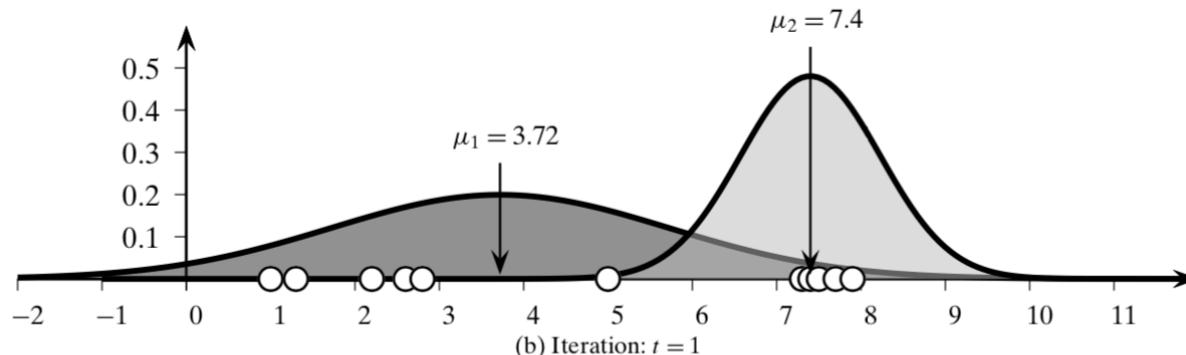
$$\sigma_1^2 = 6.13$$

$$P(C_1) = 0.71$$

$$\mu_2 = 7.4$$

$$\sigma_2^2 = 0.69$$

$$P(C_2) = 0.29$$



- The procedure converges after five iterations to

$$\mu_1 = 2.48$$

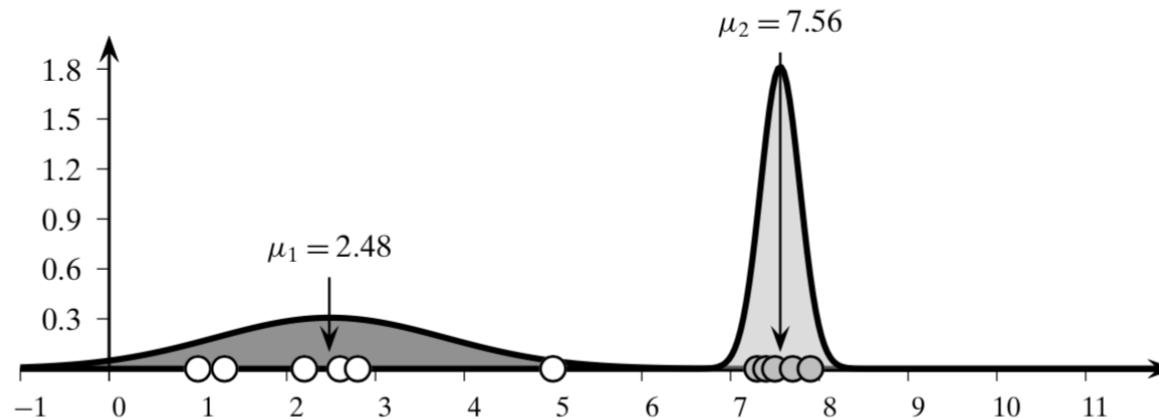
$$\sigma_1^2 = 1.69$$

$$P(C_1) = 0.55$$

$$\mu_2 = 7.56$$

$$\sigma_2^2 = 0.05$$

$$P(C_2) = 0.45$$



- EM Clustering returns probabilities, it does not assign points to clusters.
- In this example, we assigned elements to the cluster with the highest posterior probability

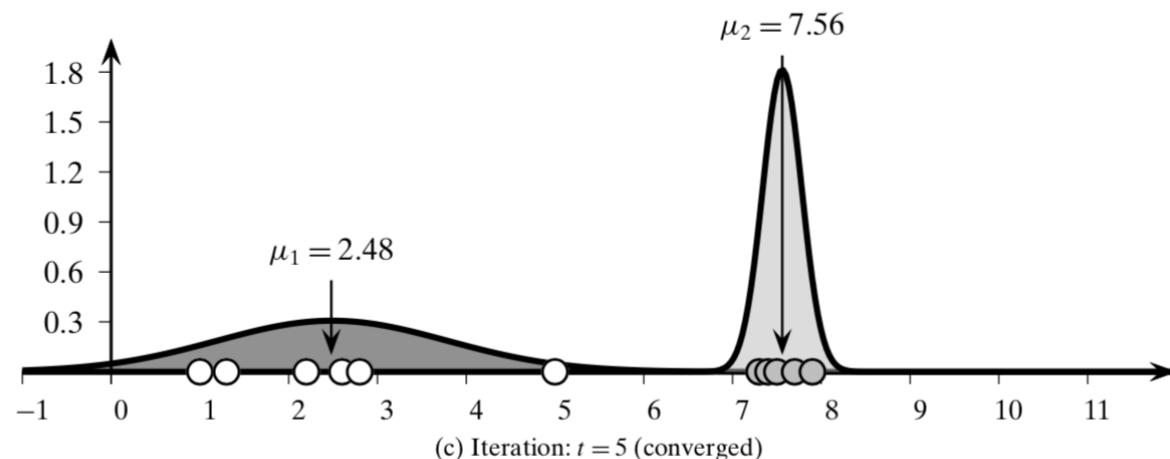
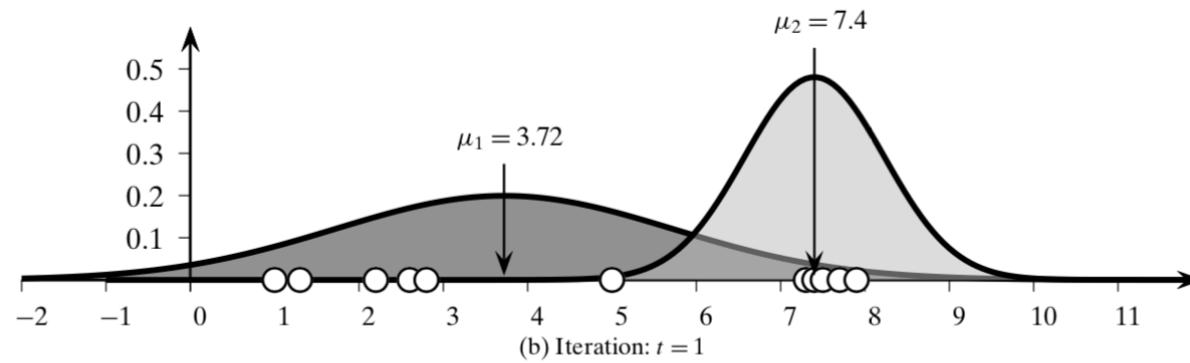
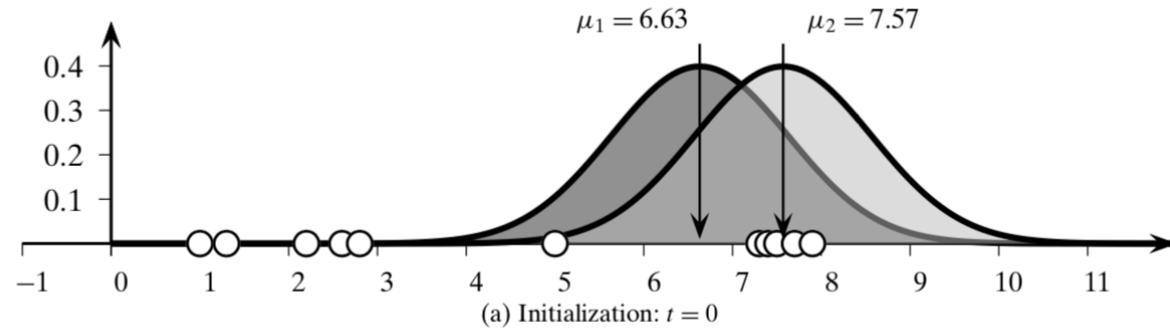


Figure 13.4. EM in one dimension.

Example from Chapter 13 of the textbook

Example with two variables

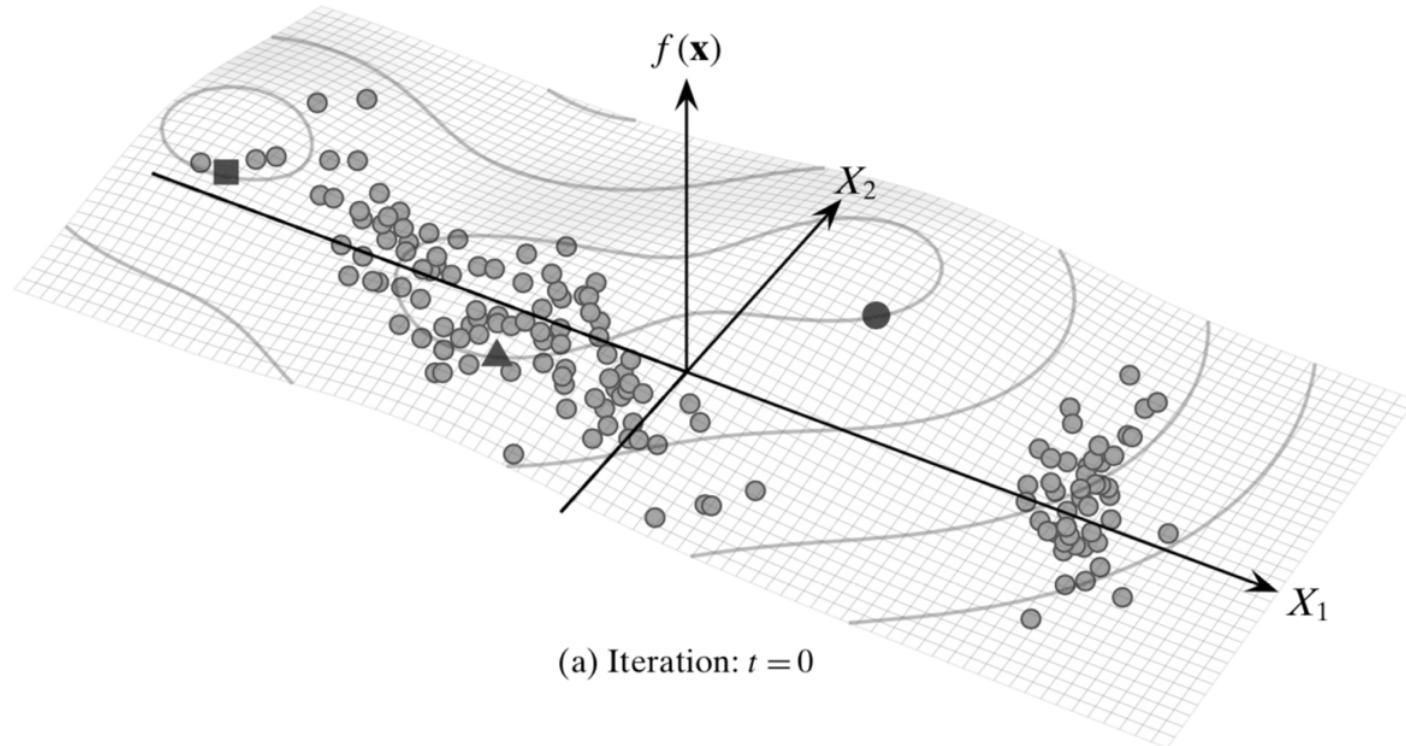


Figure 13.5 from Chapter 13 of the textbook

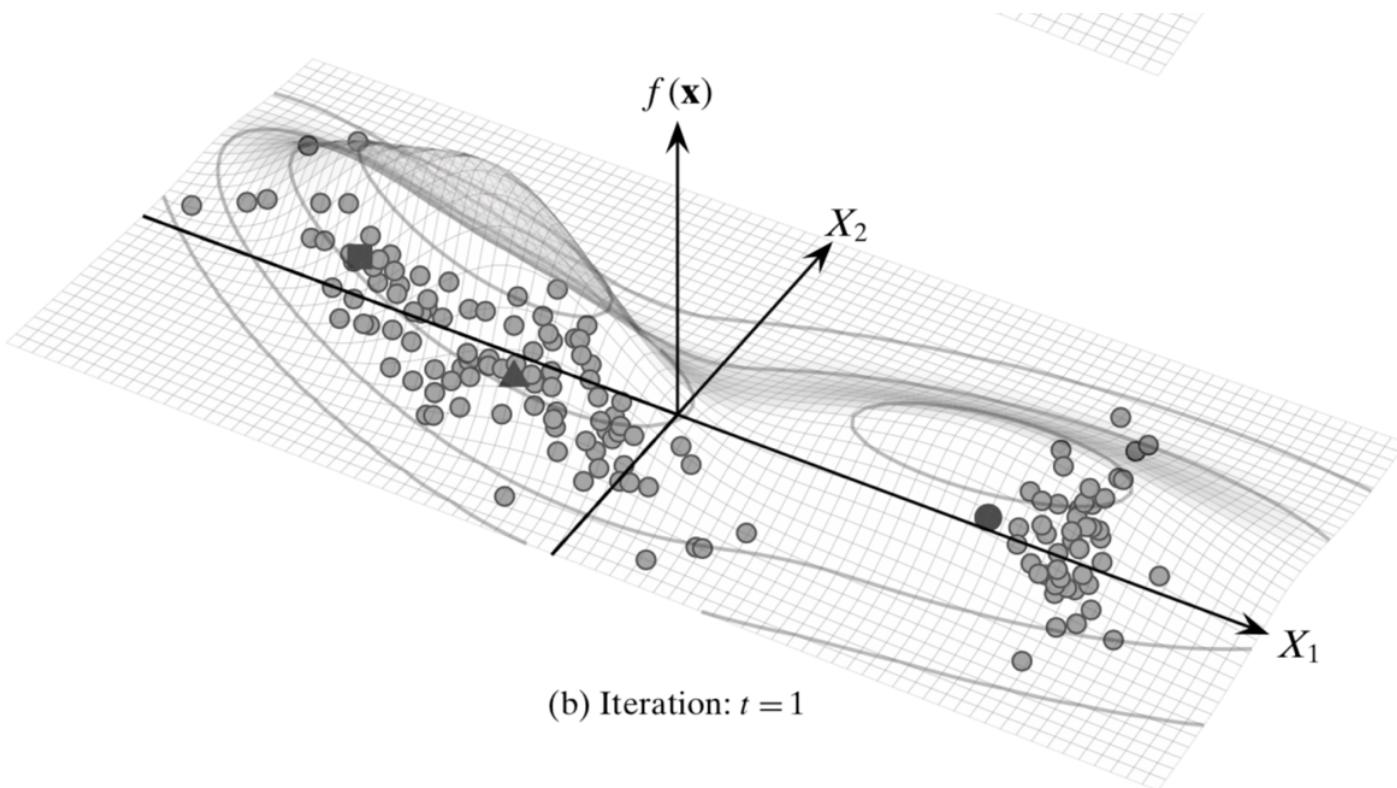


Figure 13.5 from Chapter 13 of the textbook

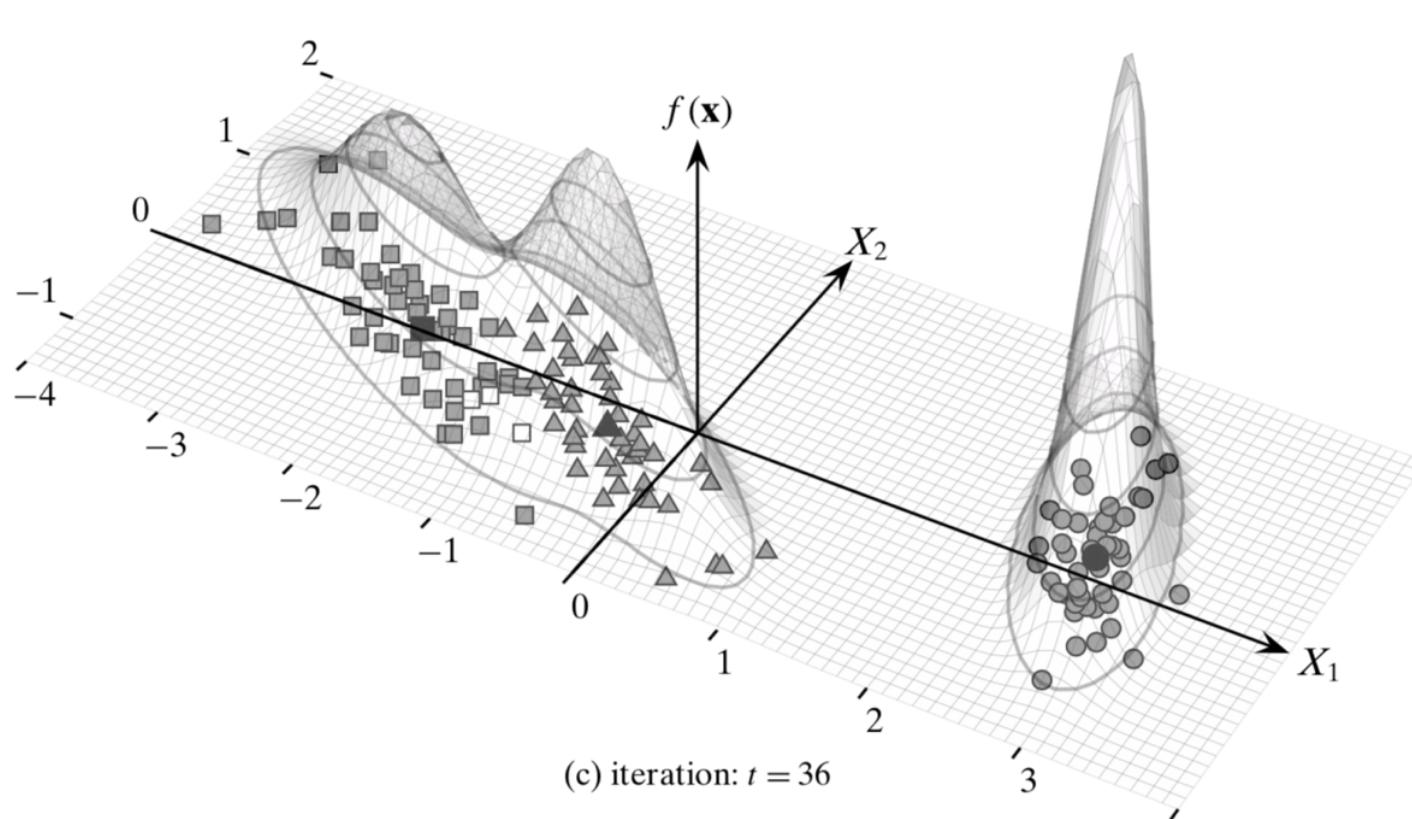


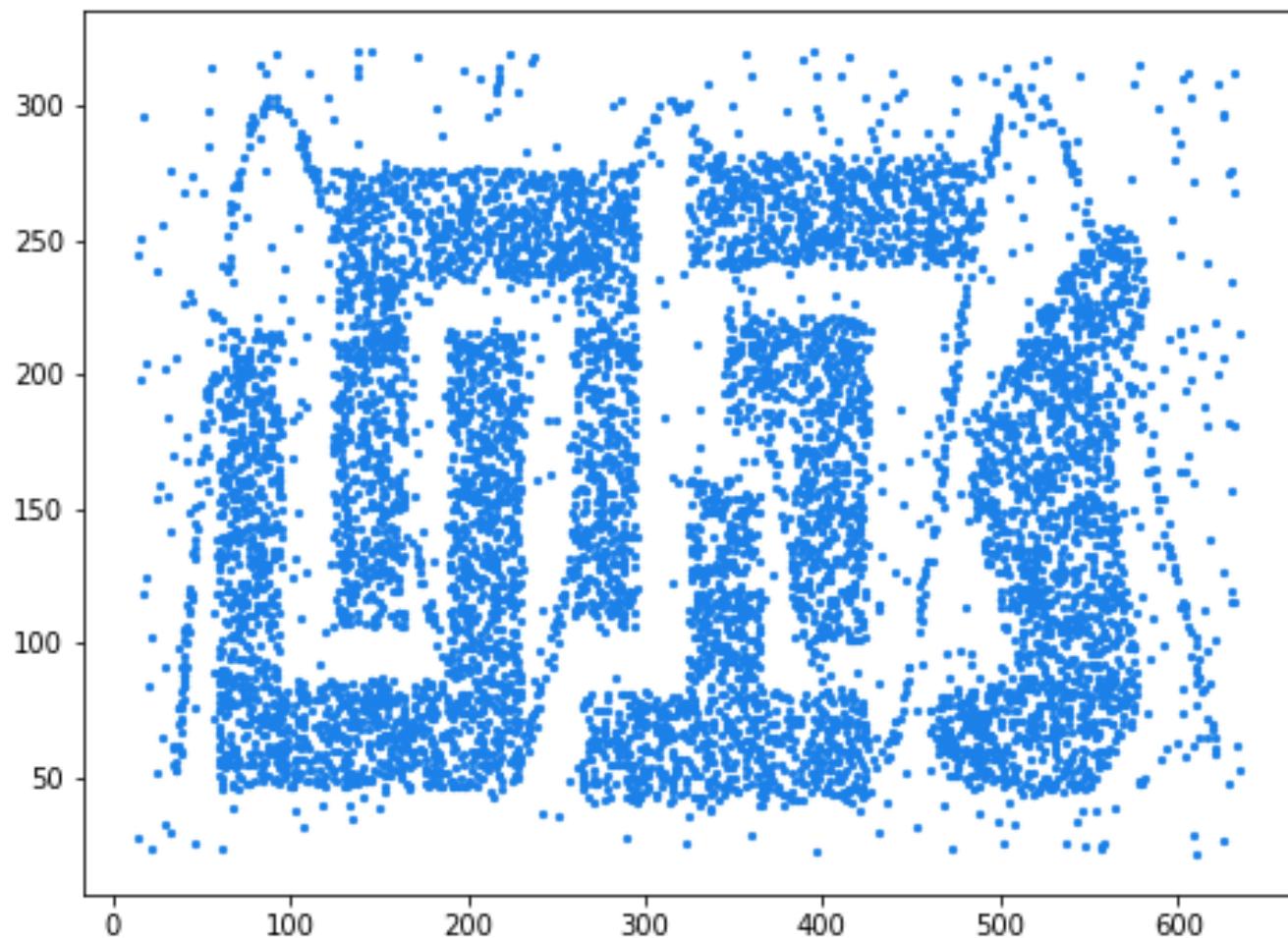
Figure 13.5 from Chapter 13 of the textbook

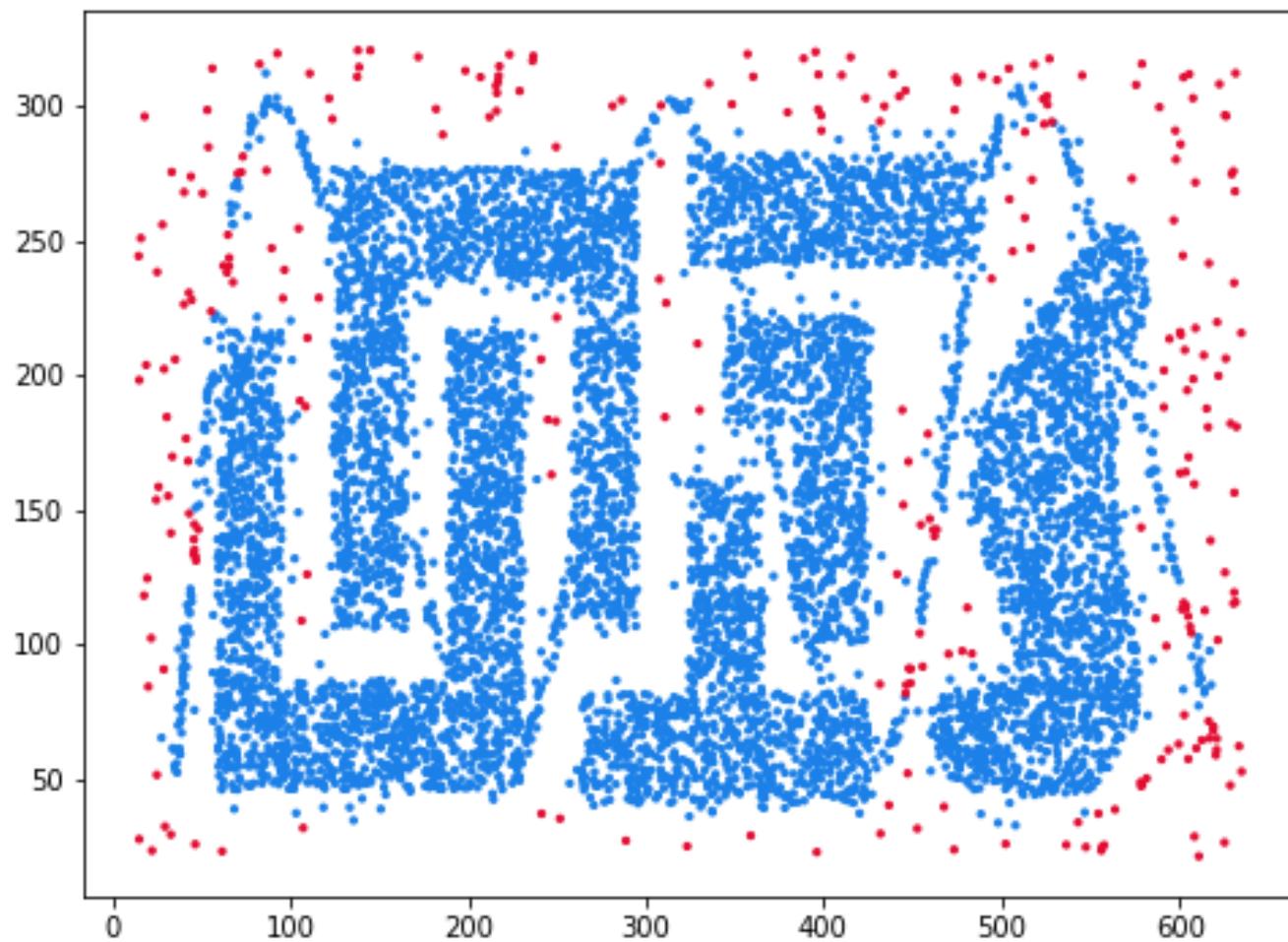
Beyond convex clusters

Representative-based clustering methods are suitable for finding ellipsoid-shaped clusters, or at best convex clusters

For nonconvex clusters, these methods have trouble finding the true clusters

Density-based methods can mine such nonconvex clusters





- 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

- The neighborhood within a radius ϵ of a given object is called the ϵ -neighborhood of the object

$$N_\epsilon(x) = \{y | \delta(x, y) \leq \epsilon\}$$

- **Core Point**
 - If the ϵ -neighborhood of an object contains at least **minpts** objects, then the object is a core object
- **Border Point**
 - If its ϵ -neighborhood does not contain at least **minpts**, but it is inside the neighborhood of a core point
- **Noise Point**
 - If its ϵ -neighborhood does not contain at least **minpts** and it does not belong to the neighborhood of a core point

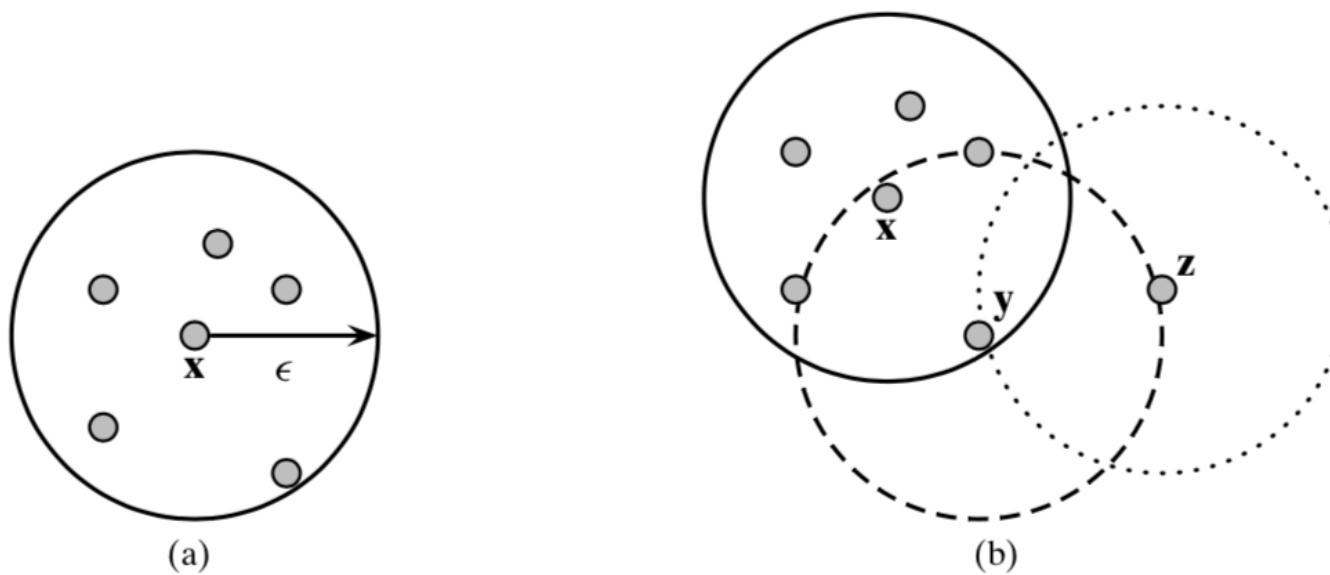
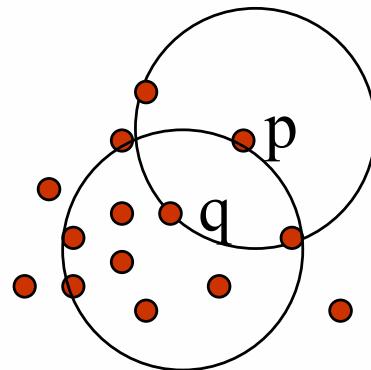


Figure 15.2. (a) Neighborhood of a point. (b) Core, border, and noise points.

Core, border and noise points when minpts is 6
“Data Mining and Analysis” by M.J. Zaki & W. Meira

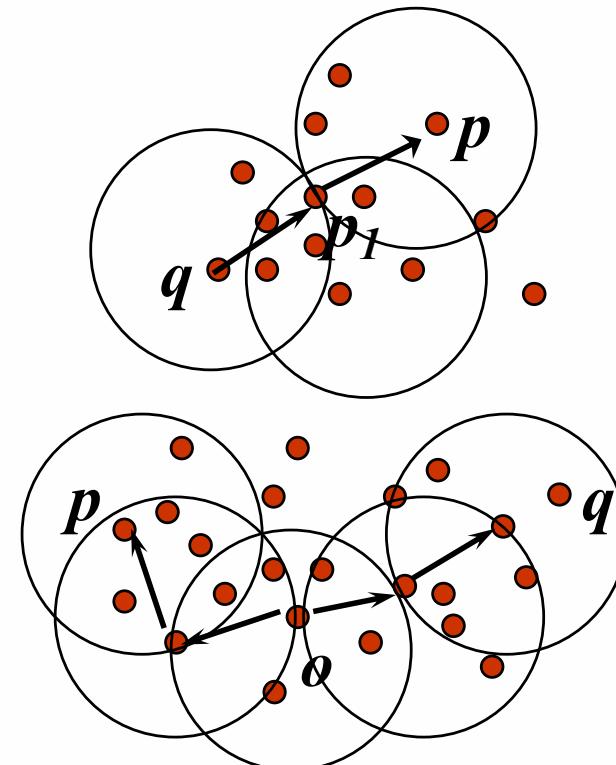
- Directly density reachable
 - An object x is directly density-reachable from object y if x is within the ε -neighborhood of y and y is a core object
- Density Reachable
 - An object x is density-reachable from object y if there is a chain of objects x_1, \dots, x_n where $x_1=x$ and $x_n=y$ such that x_{i+1} is directly density reachable from x_i
- Density Connected
 - An object p is density-connected to q with respect to ε and MinPts if there is an object o such that both p and q are density reachable from o
- Density-Based Cluster
 - A density-based cluster is defined as a maximal set of density connected points.

Directly density-reachable



MinPts = 5

Density-reachable



Density-connected

- Density corresponds to have at least minpts points within a specified radius ϵ
- A border point has fewer than minpt within ϵ , 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

ALGORITHM 15.1. Density-based Clustering Algorithm

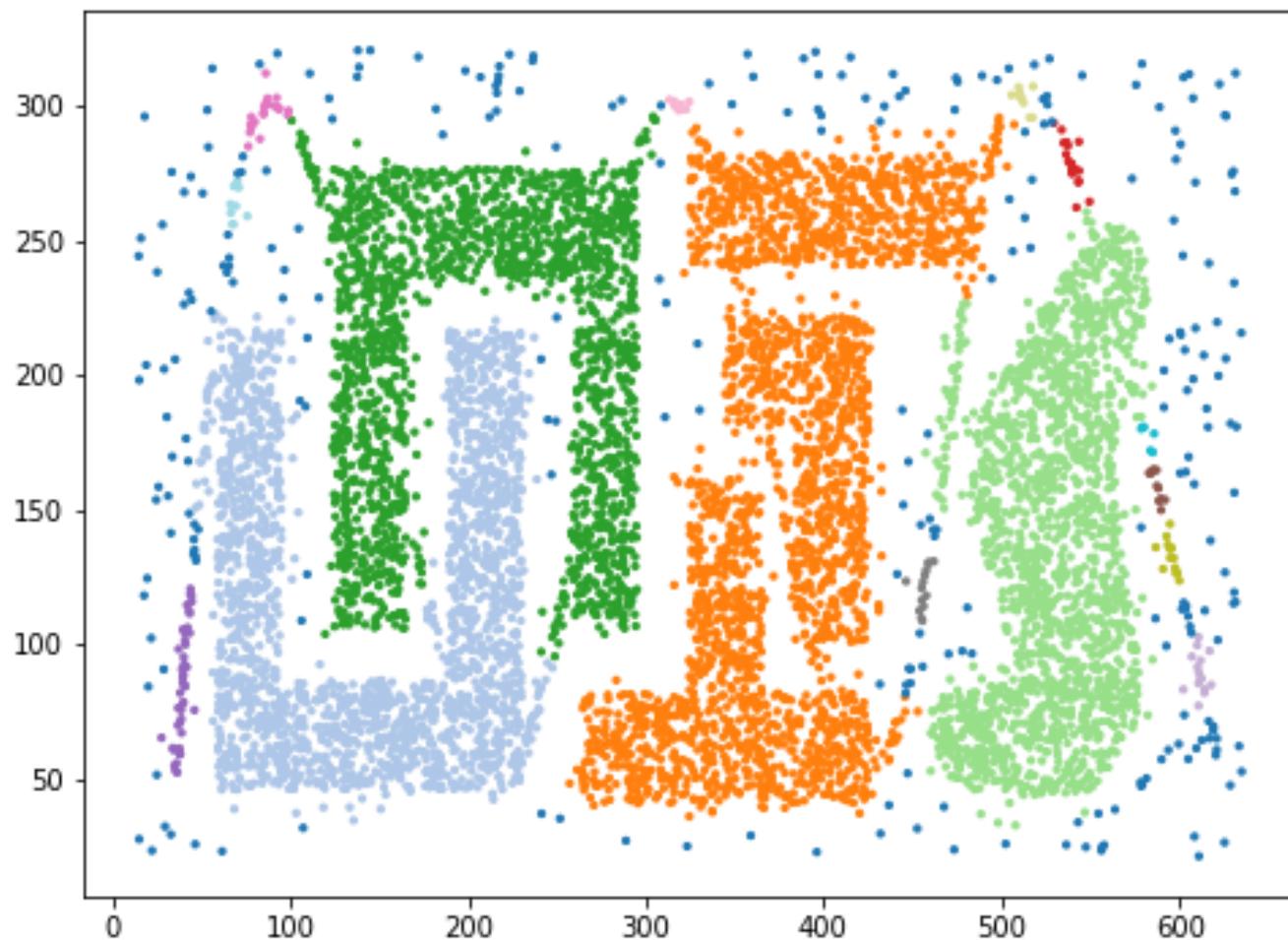
DBSCAN ($\mathbf{D}, \epsilon, minpts$):

```
1  $Core \leftarrow \emptyset$ 
2 foreach  $\mathbf{x}_i \in \mathbf{D}$  do // Find the core points
3   Compute  $N_\epsilon(\mathbf{x}_i)$ 
4    $id(\mathbf{x}_i) \leftarrow \emptyset$  // cluster id for  $\mathbf{x}_i$ 
5   if  $N_\epsilon(\mathbf{x}_i) \geq minpts$  then  $Core \leftarrow Core \cup \{\mathbf{x}_i\}$ 
6    $k \leftarrow 0$  // cluster id
7   foreach  $\mathbf{x}_i \in Core$ , such that  $id(\mathbf{x}_i) = \emptyset$  do
8      $k \leftarrow k + 1$ 
9      $id(\mathbf{x}_i) \leftarrow k$  // assign  $\mathbf{x}_i$  to cluster id  $k$ 
10    DENSITYCONNECTED ( $\mathbf{x}_i, k$ )
11   $\mathcal{C} \leftarrow \{C_i\}_{i=1}^k$ , where  $C_i \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = i\}$ 
12   $Noise \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = \emptyset\}$ 
13   $Border \leftarrow \mathbf{D} \setminus \{Core \cup Noise\}$ 
14 return  $\mathcal{C}, Core, Border, Noise$ 
```

DENSITYCONNECTED (\mathbf{x}, k):

```
15 foreach  $\mathbf{y} \in N_\epsilon(\mathbf{x})$  do
16    $id(\mathbf{y}) \leftarrow k$  // assign  $\mathbf{y}$  to cluster id  $k$ 
17   if  $\mathbf{y} \in Core$  then DENSITYCONNECTED ( $\mathbf{y}, k$ )
```

- DBSCAN needs to compute the ε -neighborhood for each point
- If the dimensionality is not too high this can be done efficiently using a spatial index structure in $O(n \log n)$
- If the dimensionality is high, it takes $O(n^2)$ to compute the neighborhood for each point.
- Once the ε -neighborhood has been computed the algorithm needs only a single pass over all the points to find the density connected clusters.
- The overall complexity of DBSCAN is between $O(n \log n)$ and in the worst-case.



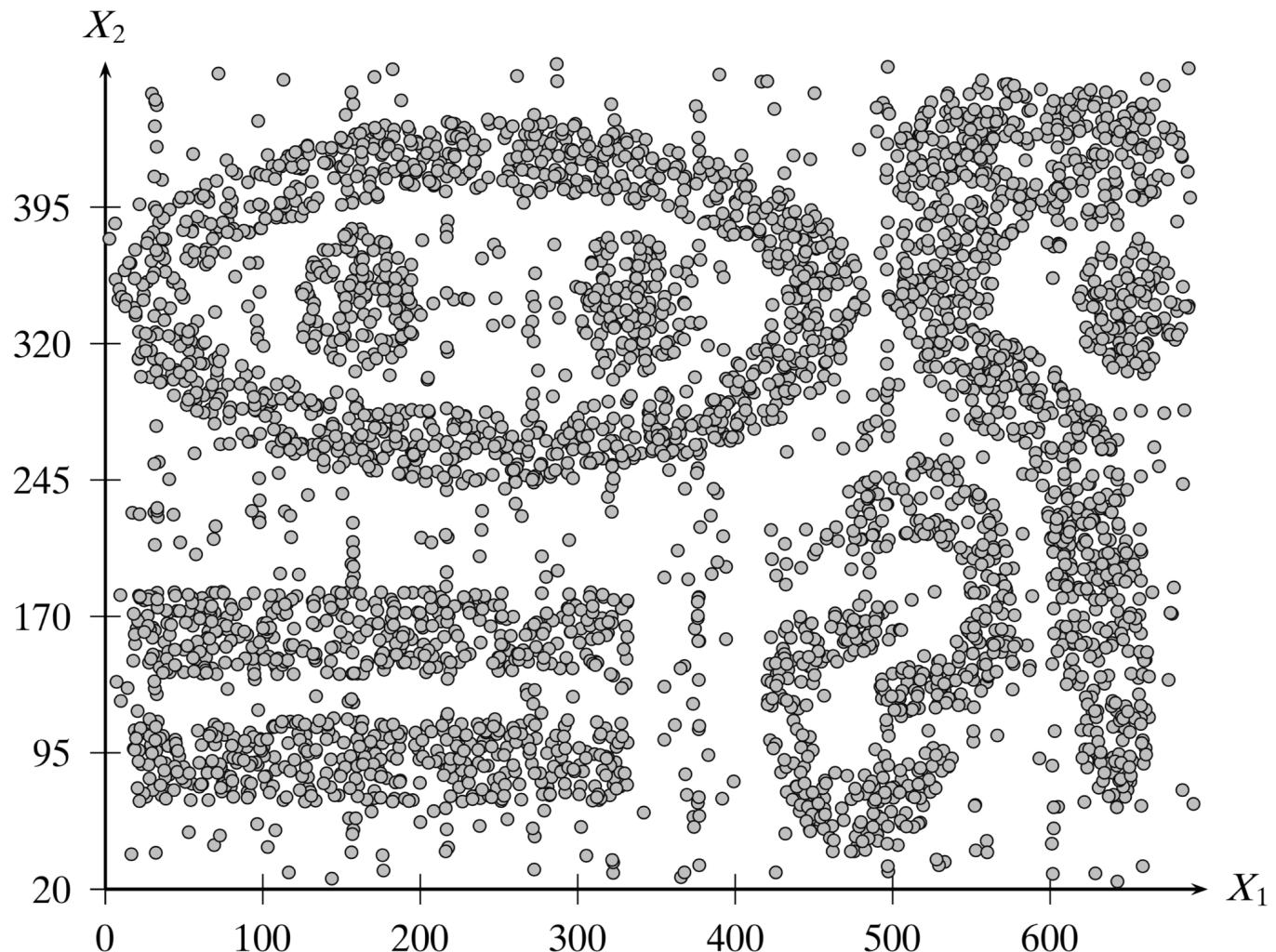


Figure 15.1 from Chapter 15 of the textbook “Data Mining and Analysis” by M.J. Zaki & W. Meira

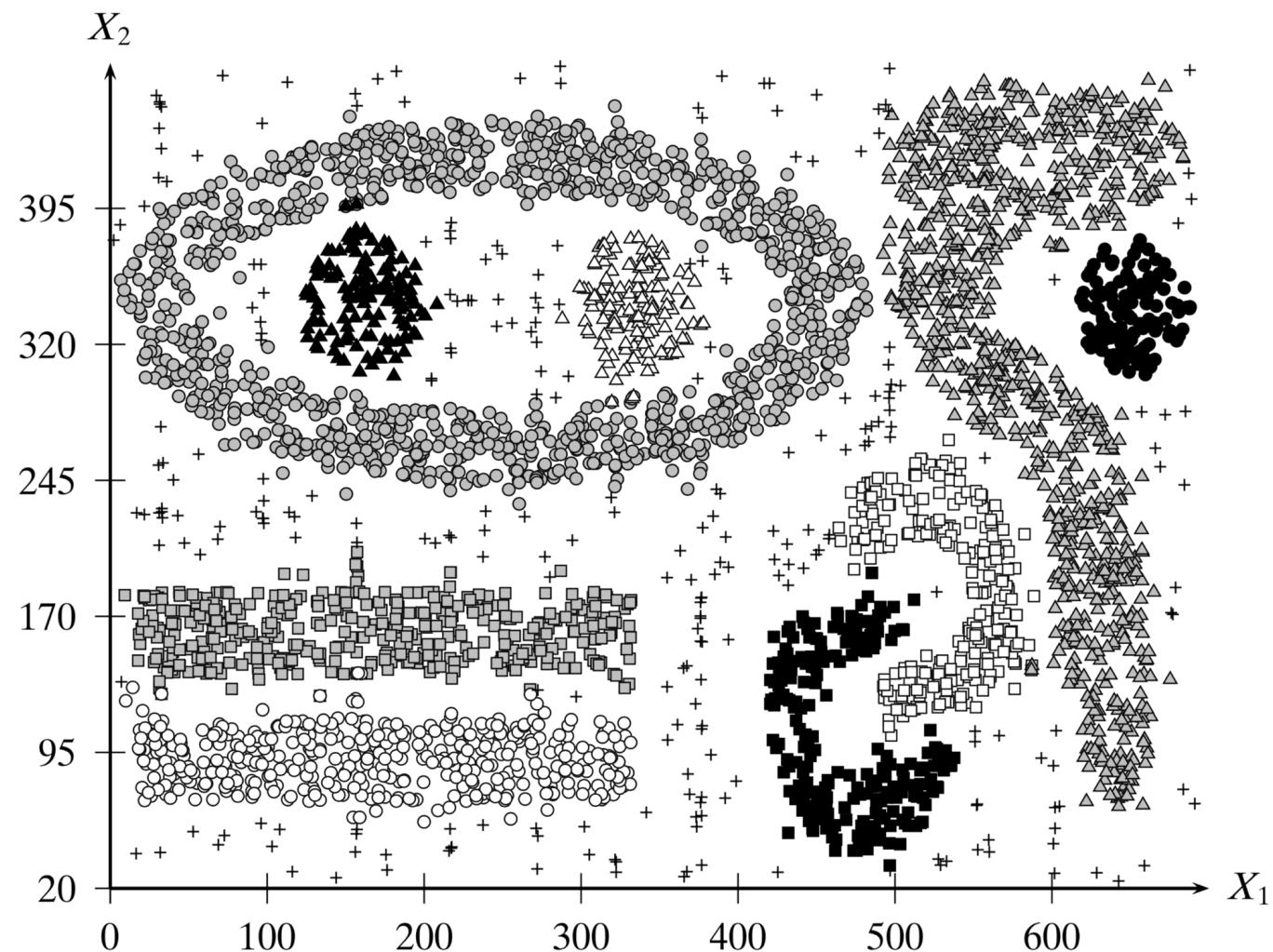


Figure 15.3 from Chapter 15 of the textbook “Data Mining and Analysis” by M.J. Zaki & W. Meira

Let's play with DBSCAN

<https://www.naftaliharris.com/blog/visualizing-dbscan-clustering/>

DBScan fails when applied
to data of varying density

HDBSCAN

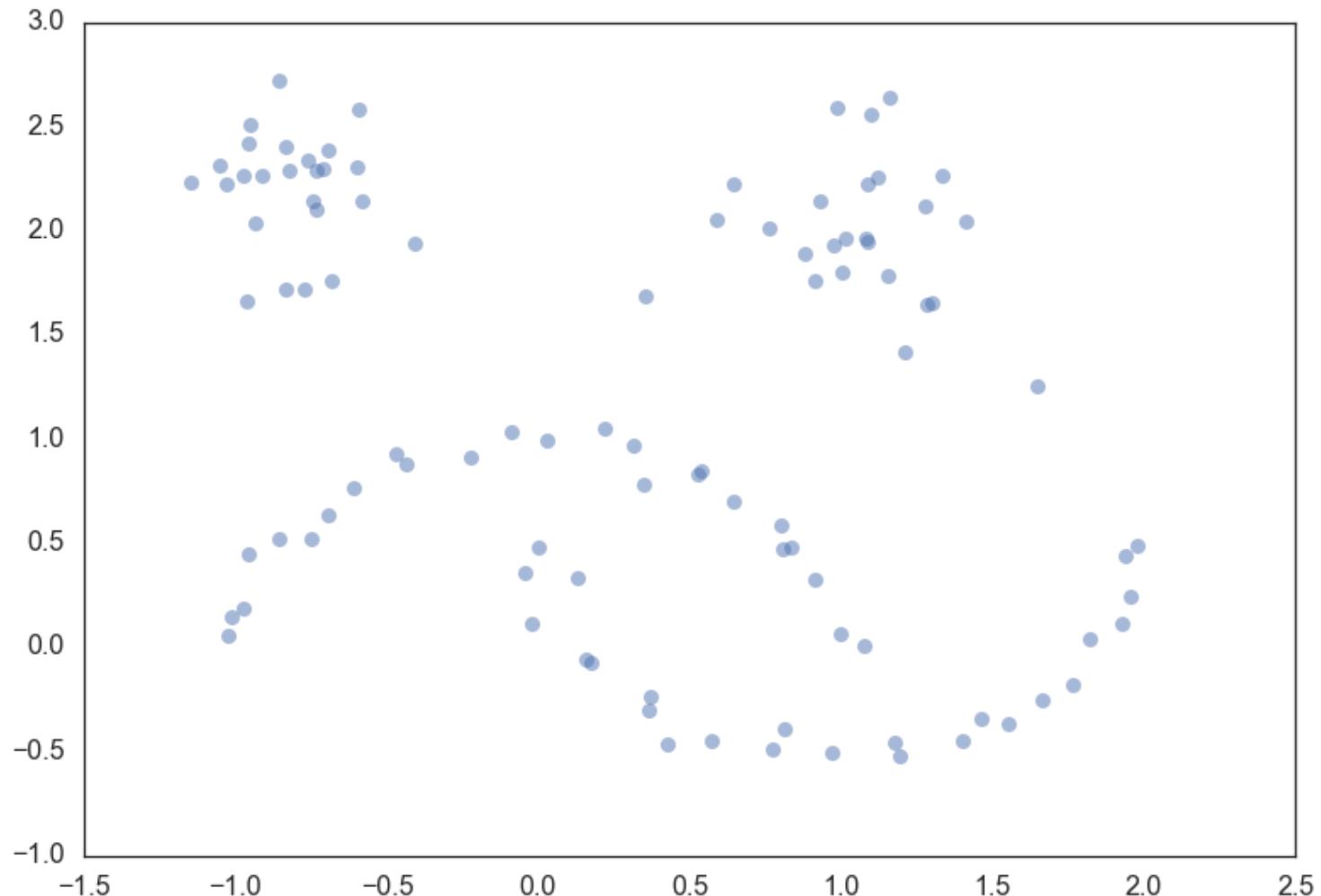
- Extends DBSCAN by converting it into a hierarchical clustering algorithm, and then using a technique to extract a flat clustering based in the stability of clusters
- It defines the “core distance” of a point as the minimum distance of a point to its k -th nearest neighbor or $\text{core}_k(x)$
- Thus, points in high density areas will have a low core distance, while point in low density areas will have a high core distance
- This provides a local inexpensive estimate of density

- This new distance metric is defined as,

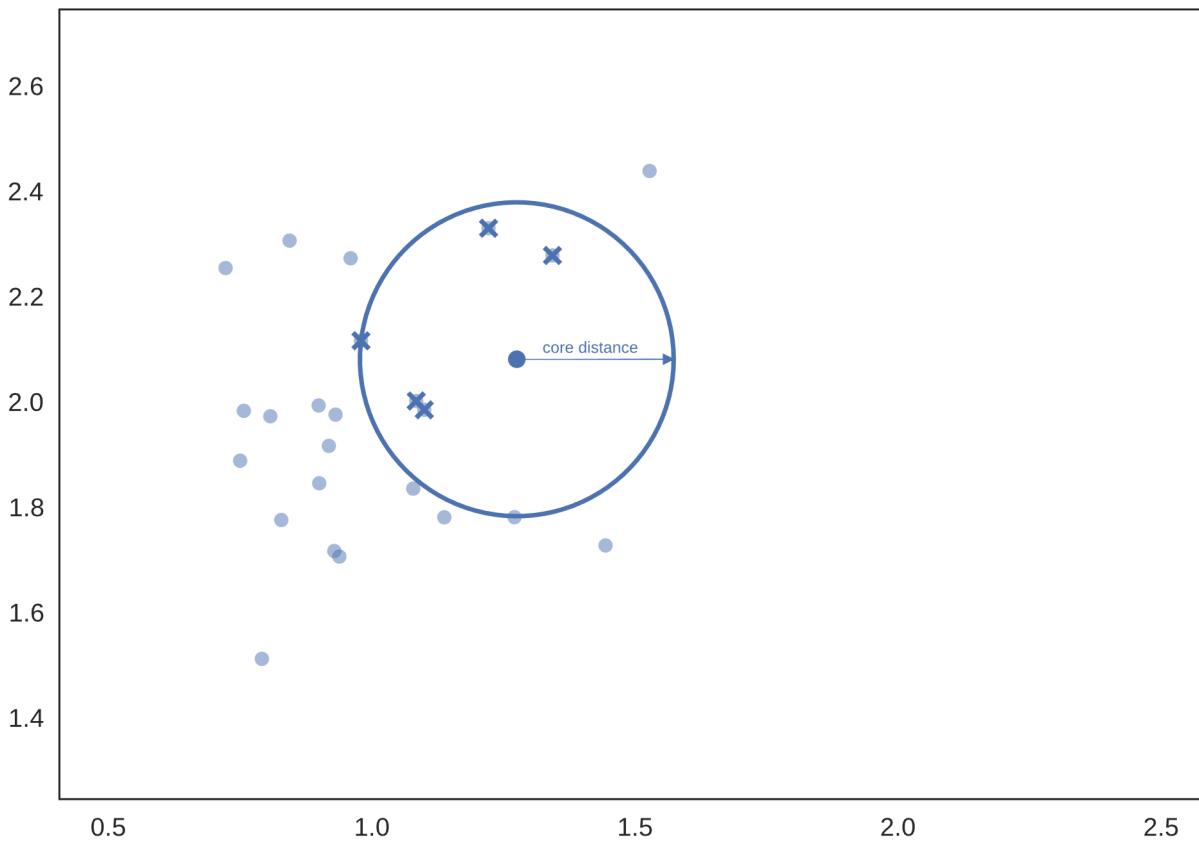
$$d_{mreach-k}(a, b) = \max\{core_k(a), core_k(b), d(a, b)\}$$

using the core distance and the original distance metrics d

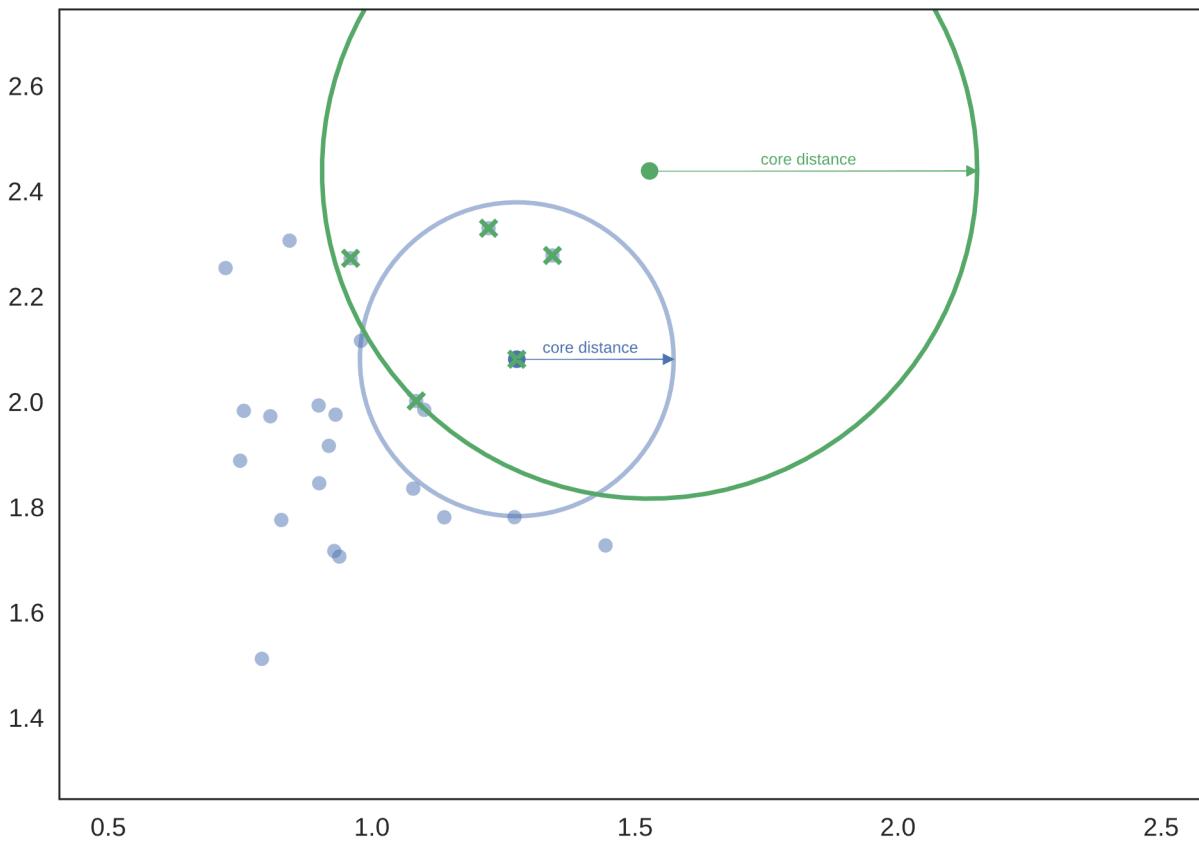
- This distance keeps dense points with lower core distance at the same distance while pushing away sparser points (at least at their core distance from any other point).
- Note that as with DBSCAN high values of k focus the algorithm toward very dense region



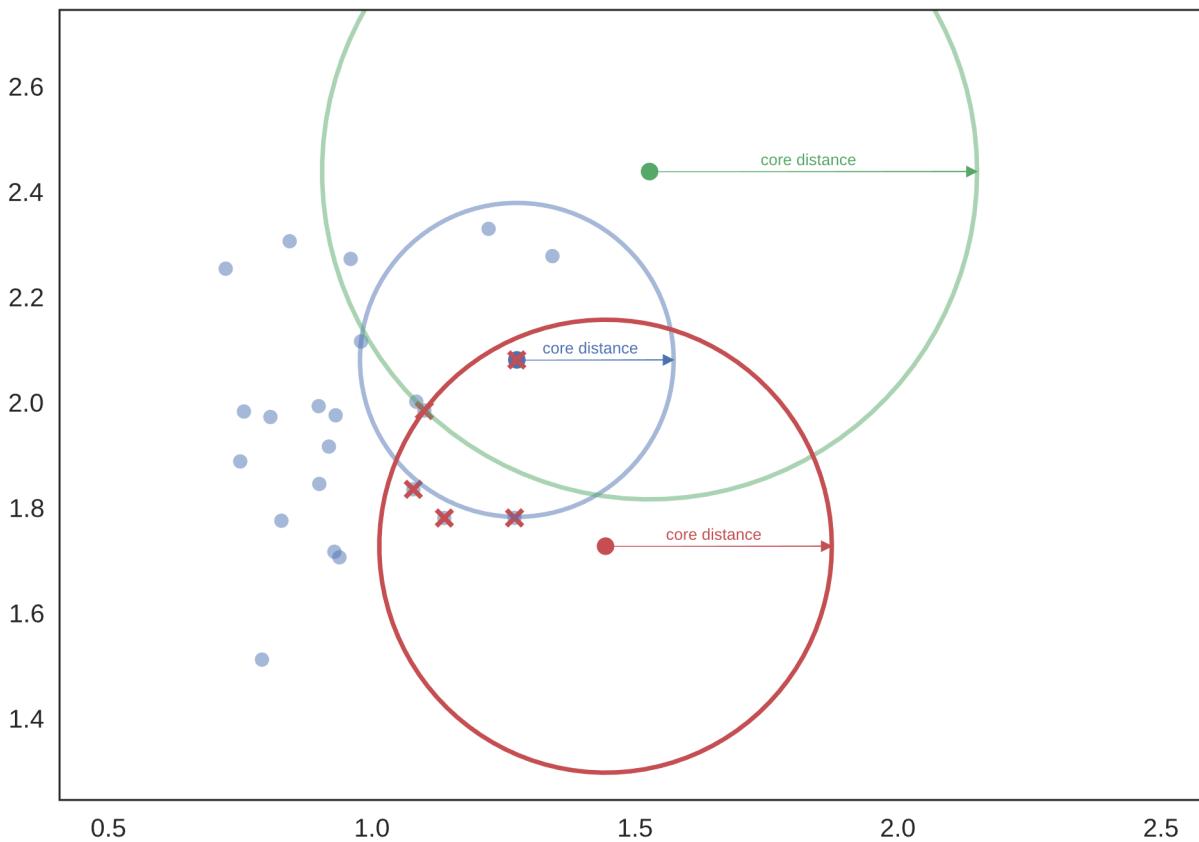
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html



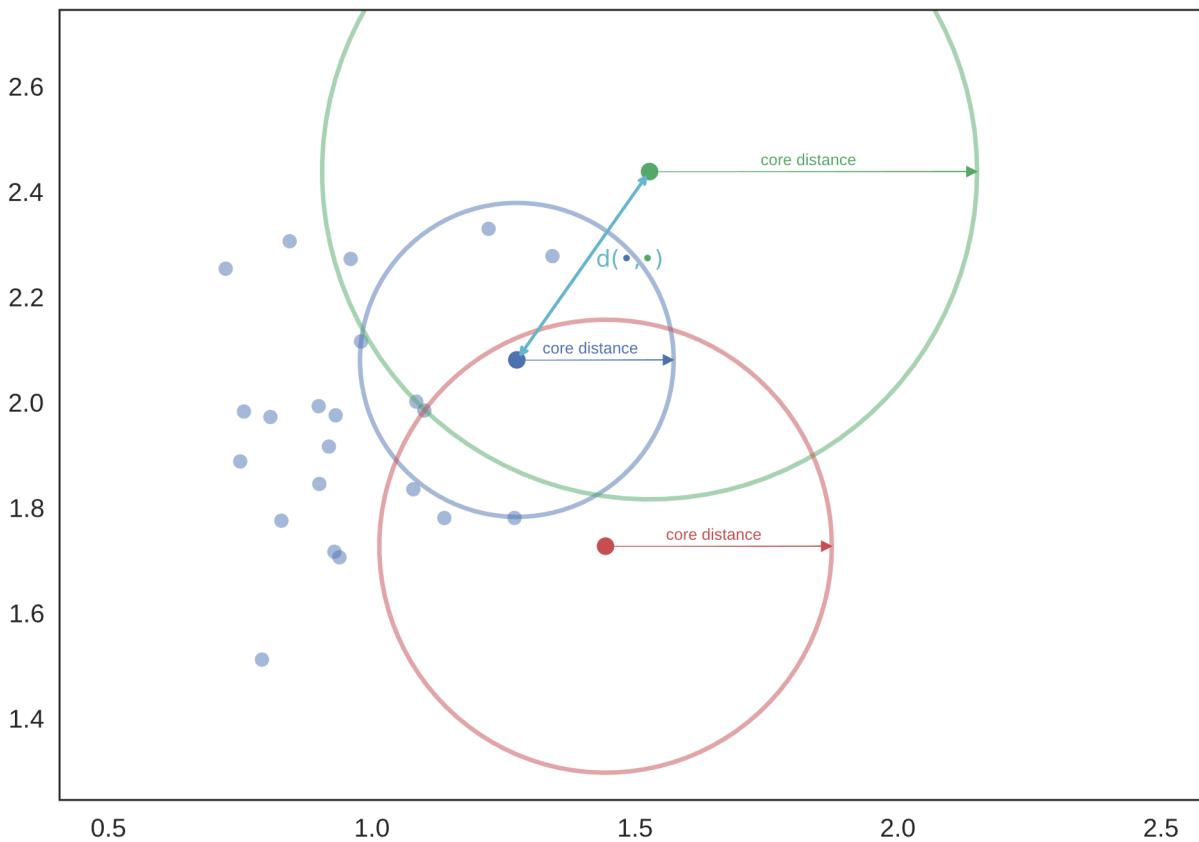
Let's consider a k equal to 5 and plot the core distance for one point.
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html



This point is in a lower density area and thus it has a higher core distance.
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

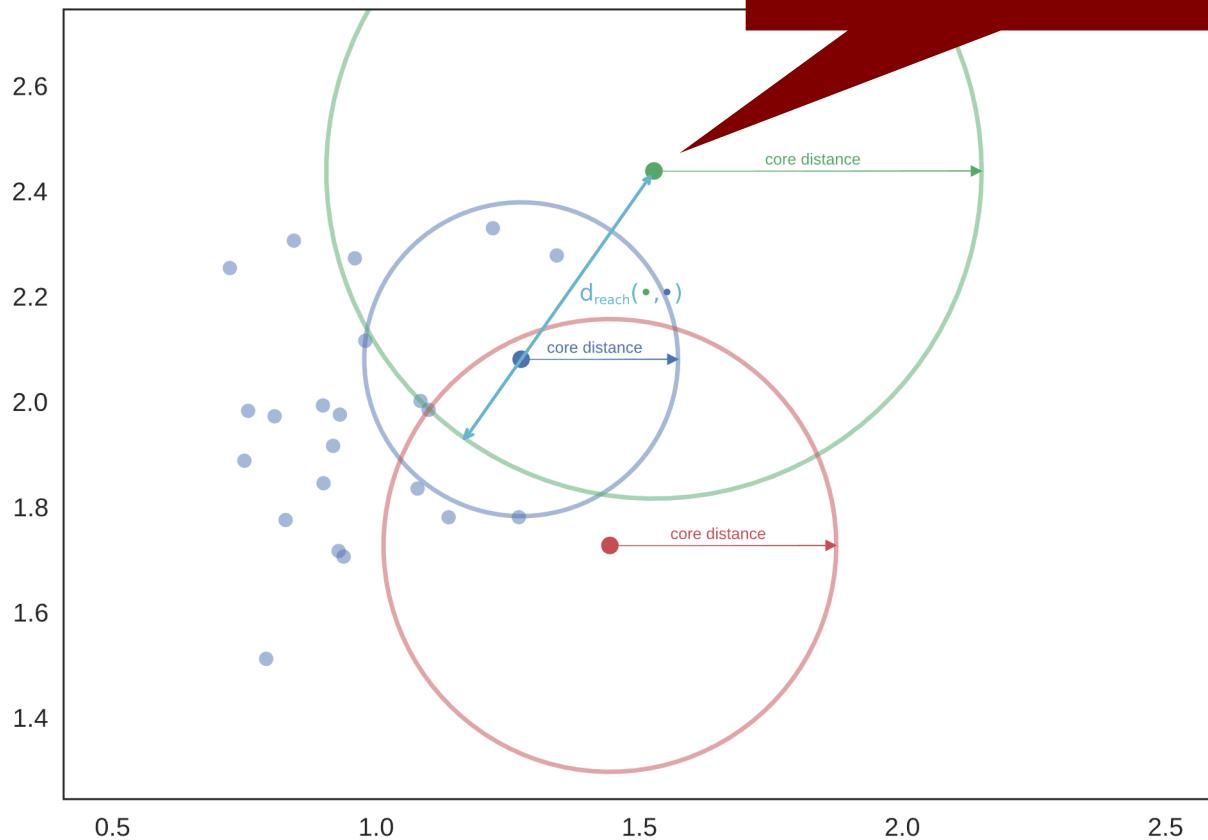


And let's add another point with its core distance
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

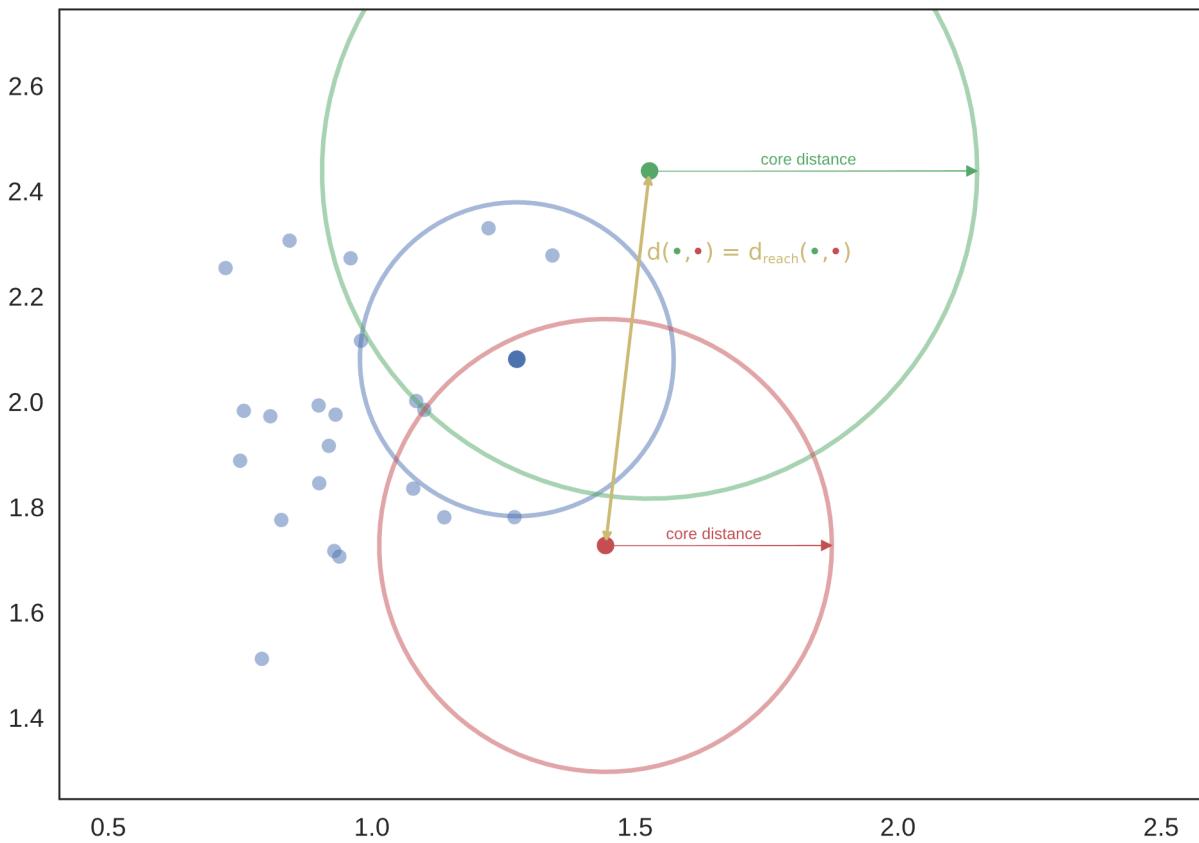


To compute the mutual reachability distance between blue and green points we draw $d(.,.)$ which crosses the blue circle but not the green one. So the green core distance is the largest value.

The two points are moved away by the new distance



To compute the mutual reachability distance between blue and green points we draw $d(.,.)$ which crosses the blue circle but not the green one. So the green core distance is the largest value.

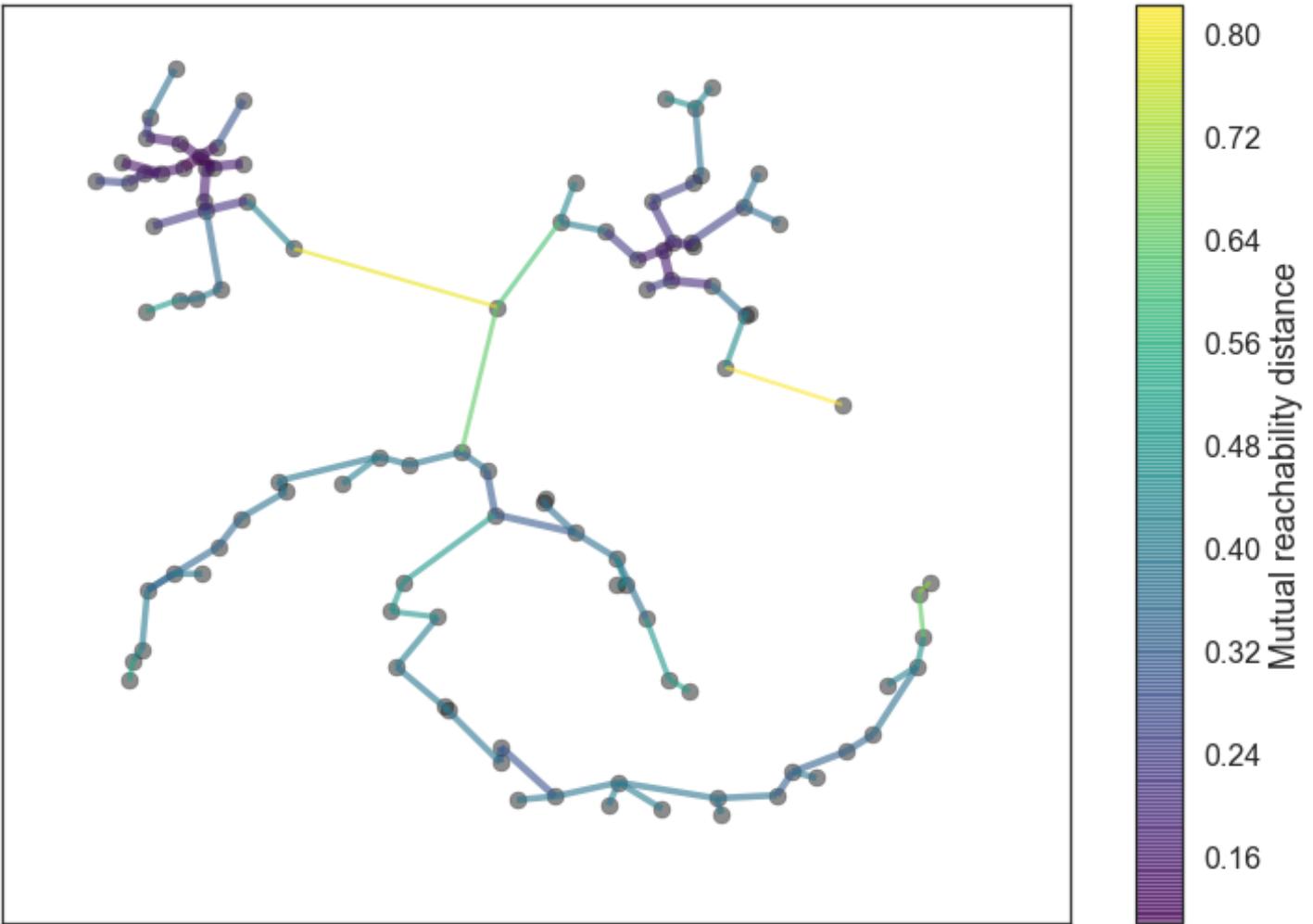


However, the mutual reachability distance between red and green is the original distance $d(\cdot, \cdot)$ which crosses both the red and green circle.

Next, Build the Minimum Spanning Tree using Prim's Algorithm

130

- Consider the data as a weighted graph
- Data points are vertices
- An edge between any two points with weight equal to the mutual reachability distance of those points
- Build the spanning tree one edge at a time by adding the lowest weight edge that connects the current tree to a vertex not yet in the tree

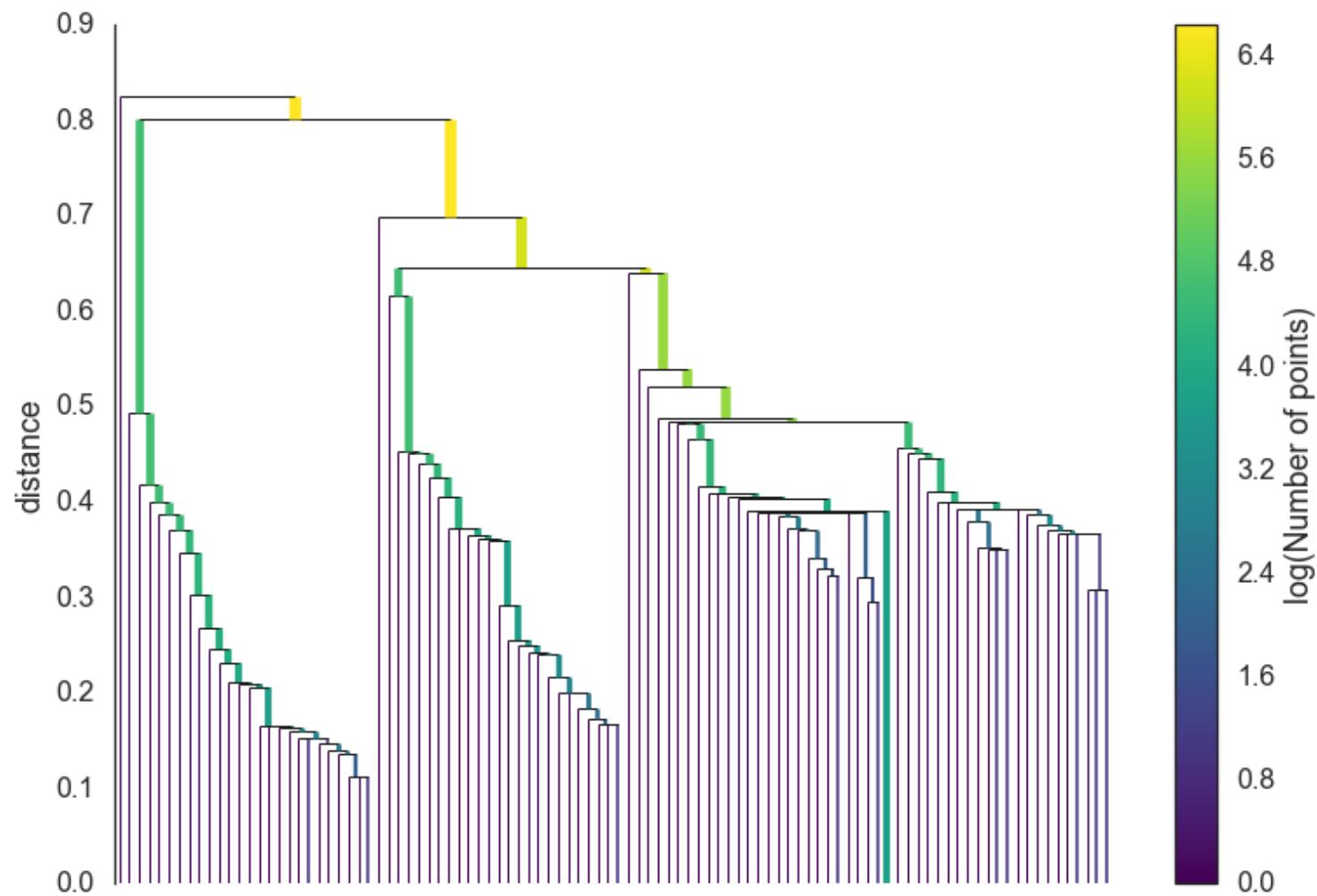


https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

HDBSCAN converts the minimum spanning tree into a hierarchy of connected components

It sorts the edges of the tree by distance
(in increasing order)

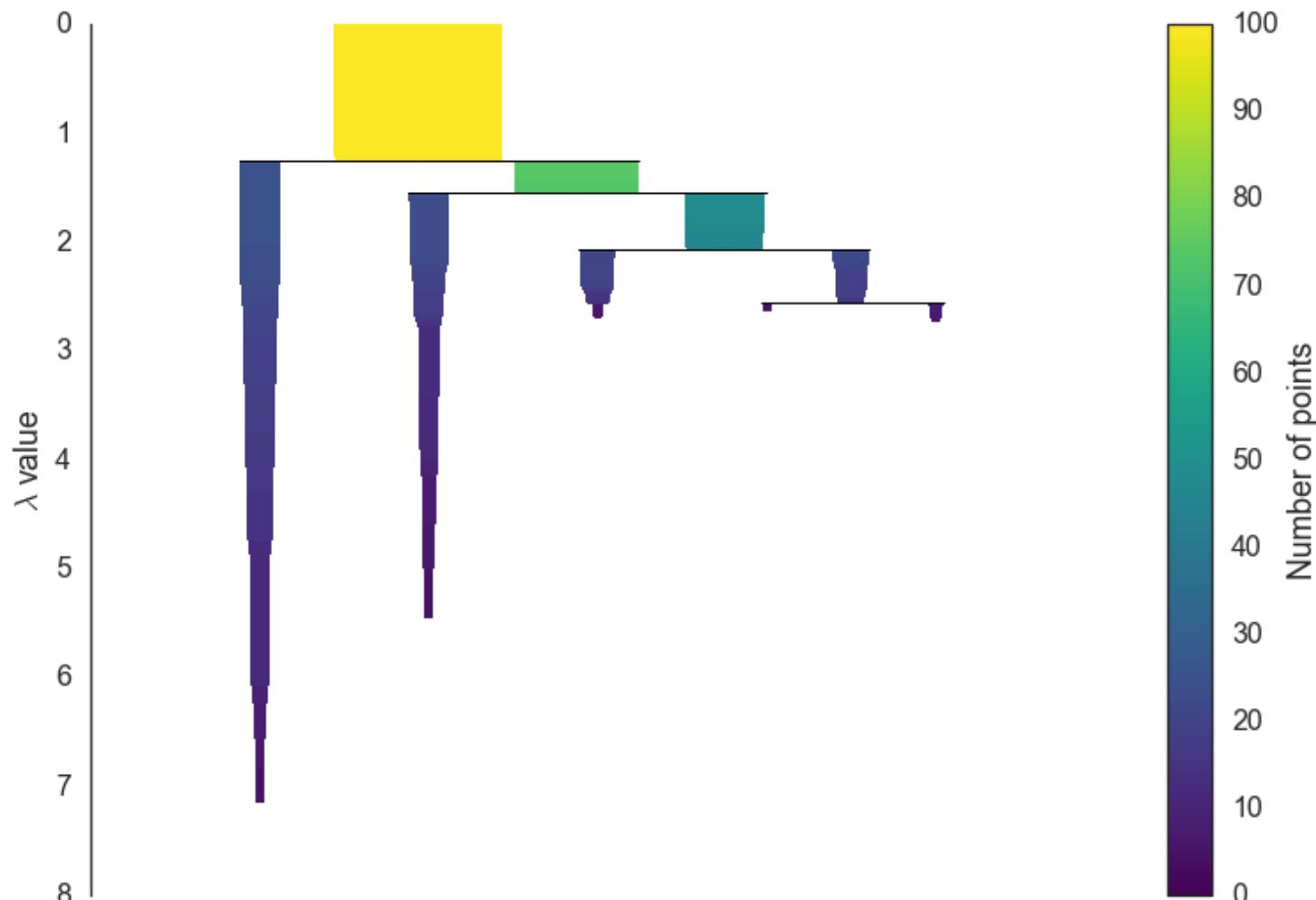
Then iterate creating a new merged cluster for each edge.



https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

We have a hierarchy now
we want a set of flat clusters

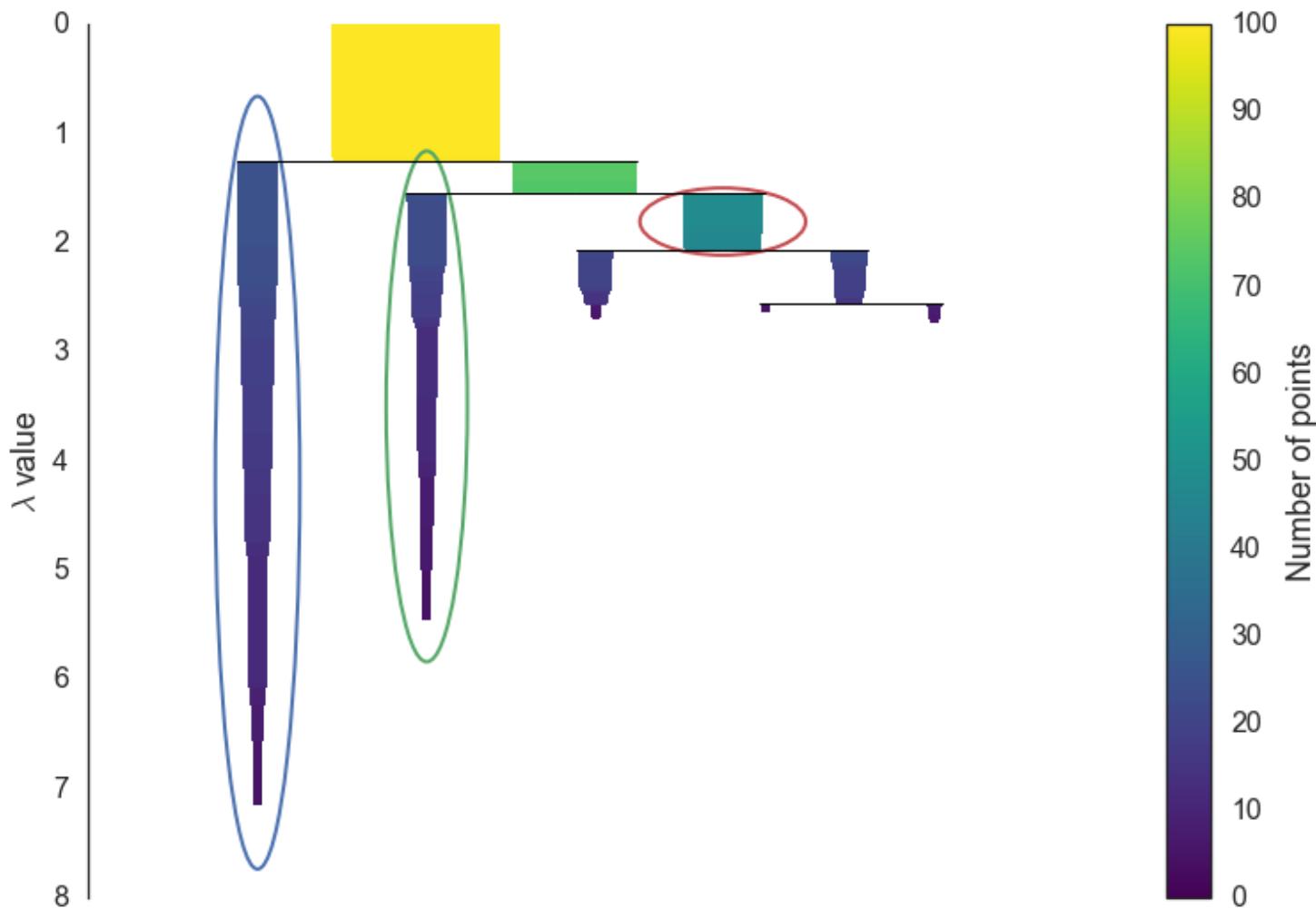
- We introduce the notion of “minimum cluster size” threshold
- We navigate the hierarchy and from the top and at each split we check the size of the merged clusters
- If a clusters has fewer points than the threshold, then the smaller clusters is eliminated (the points have “fallen out of the cluster”)
- Otherwise both clusters are maintained.
- Note that, which a point p ‘fell out of the cluster’ and at what distance value that happened is stored as $\lambda_p = 1/\text{distance}$
- At the end, a much smaller tree with a small number of nodes remains with node has data about how the size of the cluster at that node decreases over varying distance



Condensed tree ($\lambda = 1/\text{distance}$)

https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

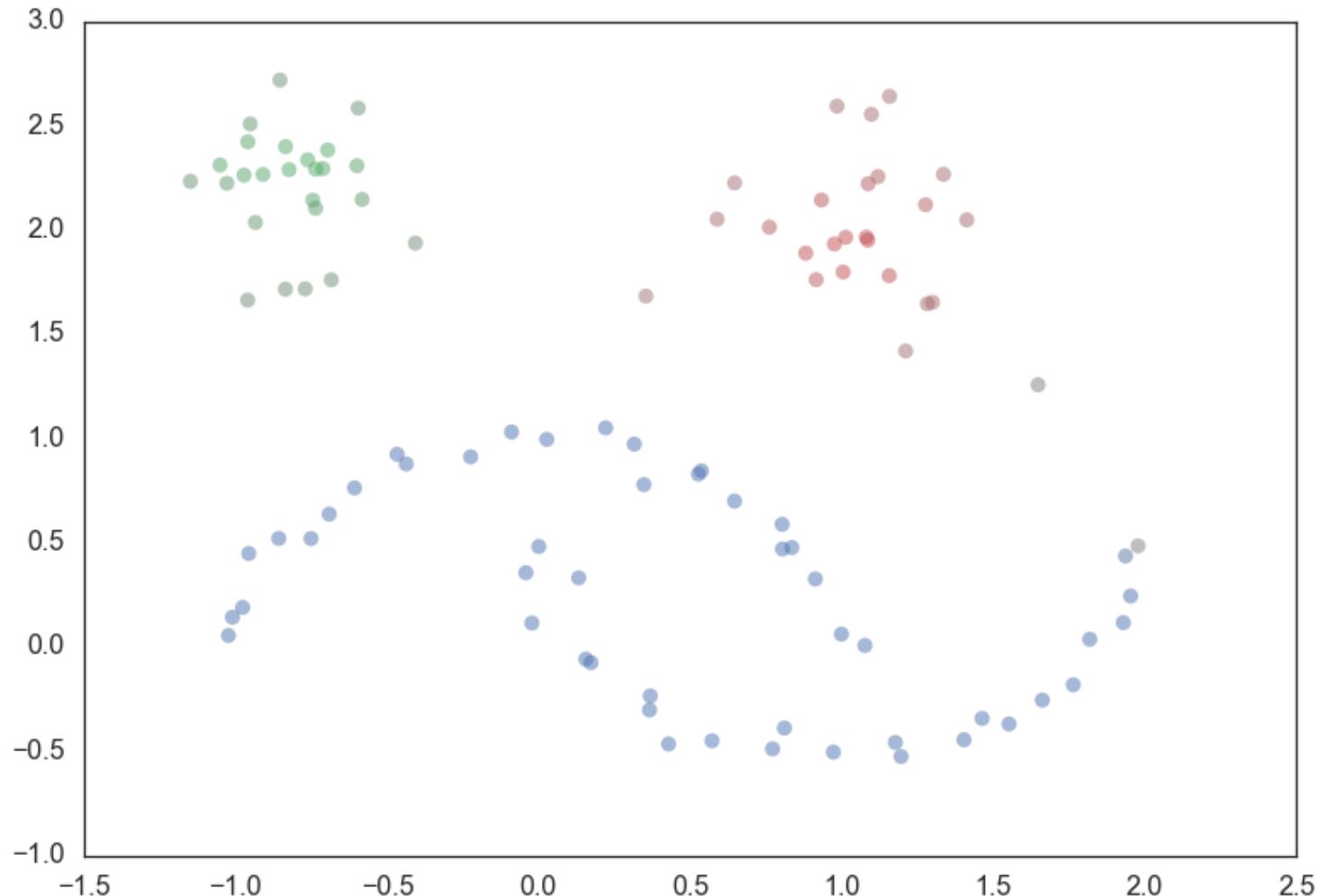
- Given the simplified dendrogram we want to select clusters that persist and have a longer lifetime since short lived clusters probably represent artifacts
- We may say that we want to choose those clusters that have the greatest area of ink in the plot. Also, we cannot select any a descendant of a cluster we chose.



By applying the two principles (larger total ink area and no descendants) we extract three clusters
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

Finally, we can assign to each point a probability to belong to a cluster based on its λ_p

HDBSCAN complexity is comparable and lower than DBSCAN



Finally, we can assign to each point a probability to belong to a cluster based on its λ_p
https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

- Birch and Chameleon are other hierarchical clustering algorithms
- There are also probabilistic hierarchical clustering algorithms
- OPTICS, and DENCLU are interesting density-based algorithms
- STING and CLIQUE are grid-based methods, where CLIQUE is also a subspace clustering algorithm
- ...



K-Means shouldn't be your first choice



SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>



K-Means probably
shouldn't be your
second choice either



SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>



Think hard about what
“cluster” means for
your application



SciPy 2016

<https://www.youtube.com/watch?v=AgPQ76Rli6A>

- “Data Mining and Analysis” by Zaki & Meira
 - Chapter 13
 - Chapter 15
- <http://www.dataminingbook.info>
- https://hdbscan.readthedocs.io/en/latest/how_hdbscan_works.html

