

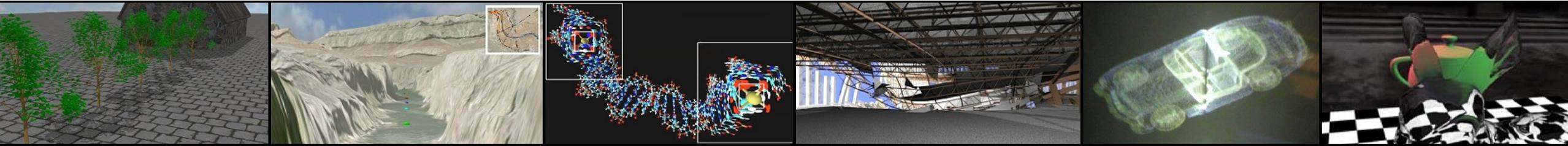
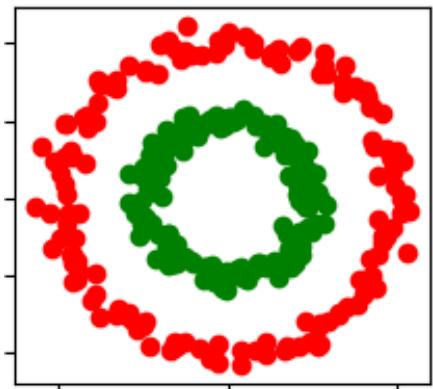
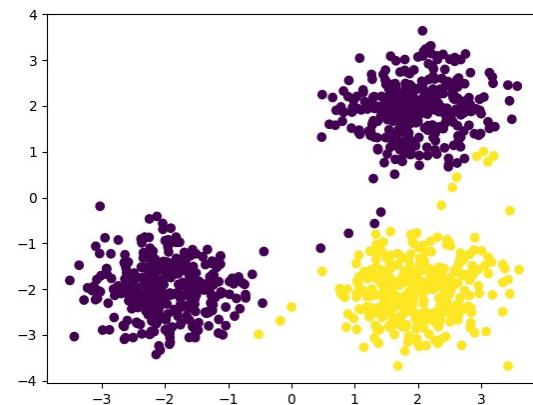
# COT 4521: INTRODUCTION TO COMPUTATIONAL GEOMETRY



## Clustering Algorithms

Paul Rosen  
Assistant Professor  
University of South Florida

Some slides from Mustafa Hajij

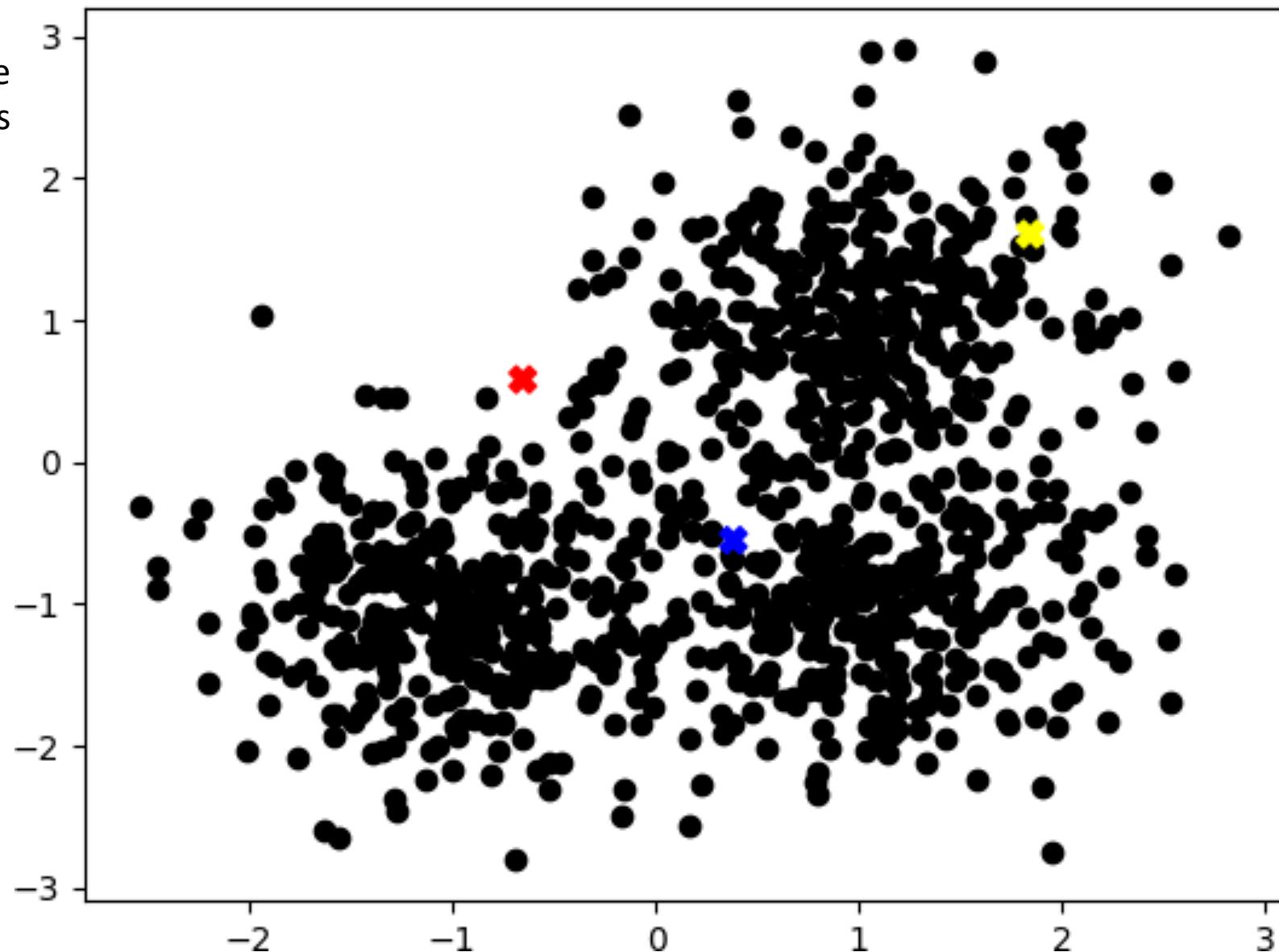


# CLUSTERING ALGORITHMS

- A CLUSTER IS A COLLECTION OF DATA OBJECTS
- A CLUSTERING ALGORITHM TRIES TO PUT SIMILAR OBJECTS TO ONE ANOTHER WITHIN THE SAME CLUSTER AND DISSIMILAR OBJECTS IN OTHER CLUSTERS
- CLUSTERING IS AN UNSUPERVISED CLASSIFICATION: THE DATA IS UNLABELED
- A CLUSTERING ALGORITHM TRIES TO UNDERSTAND WHAT KIND OF STRUCTURE IN THE DATA: WHAT SUB-POPULATION DOES THE DATA HAVE?

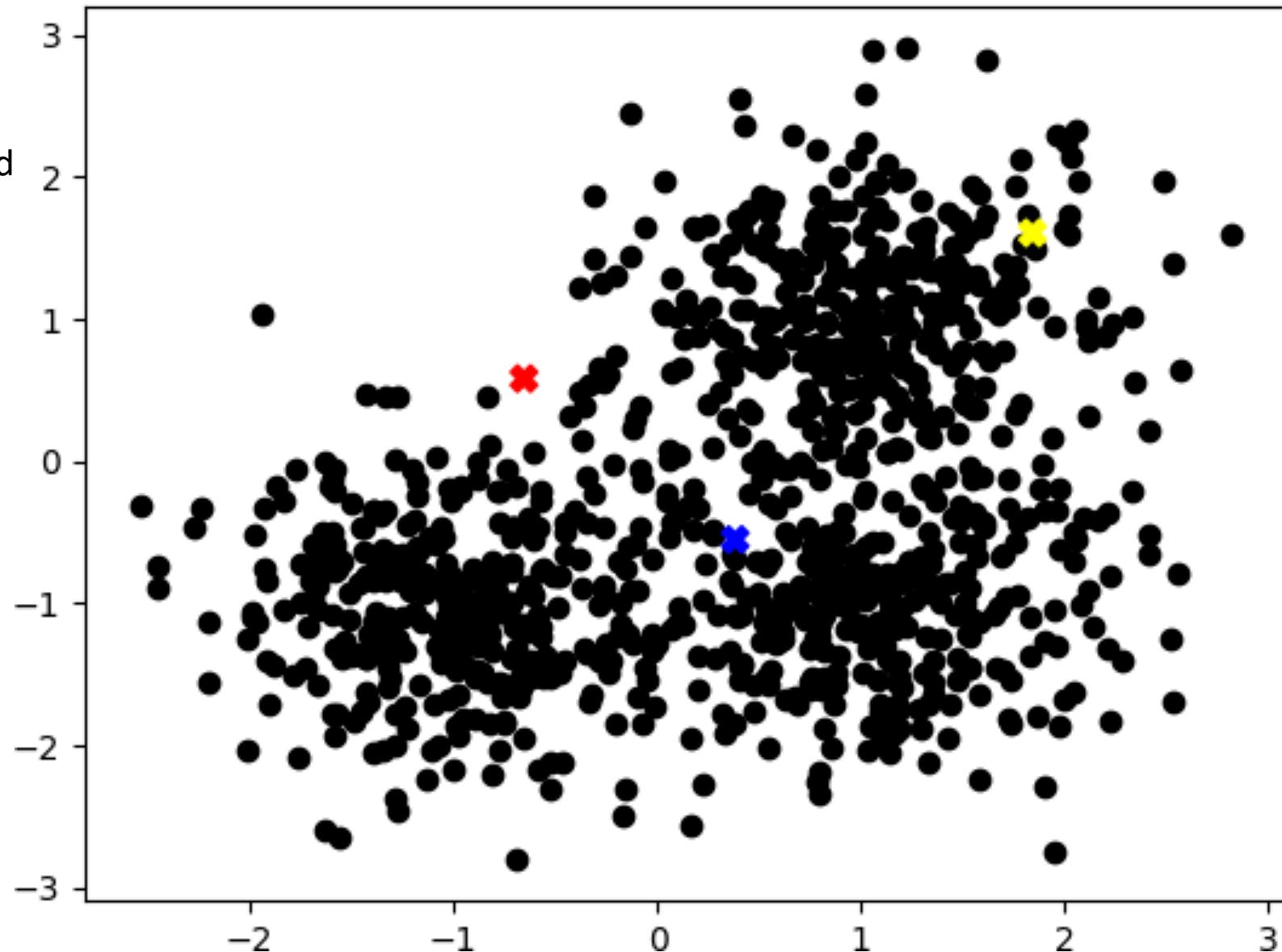
# K-MEANS ALGORITHM: EXAMPLE

Say that we have the following data points in the plan.



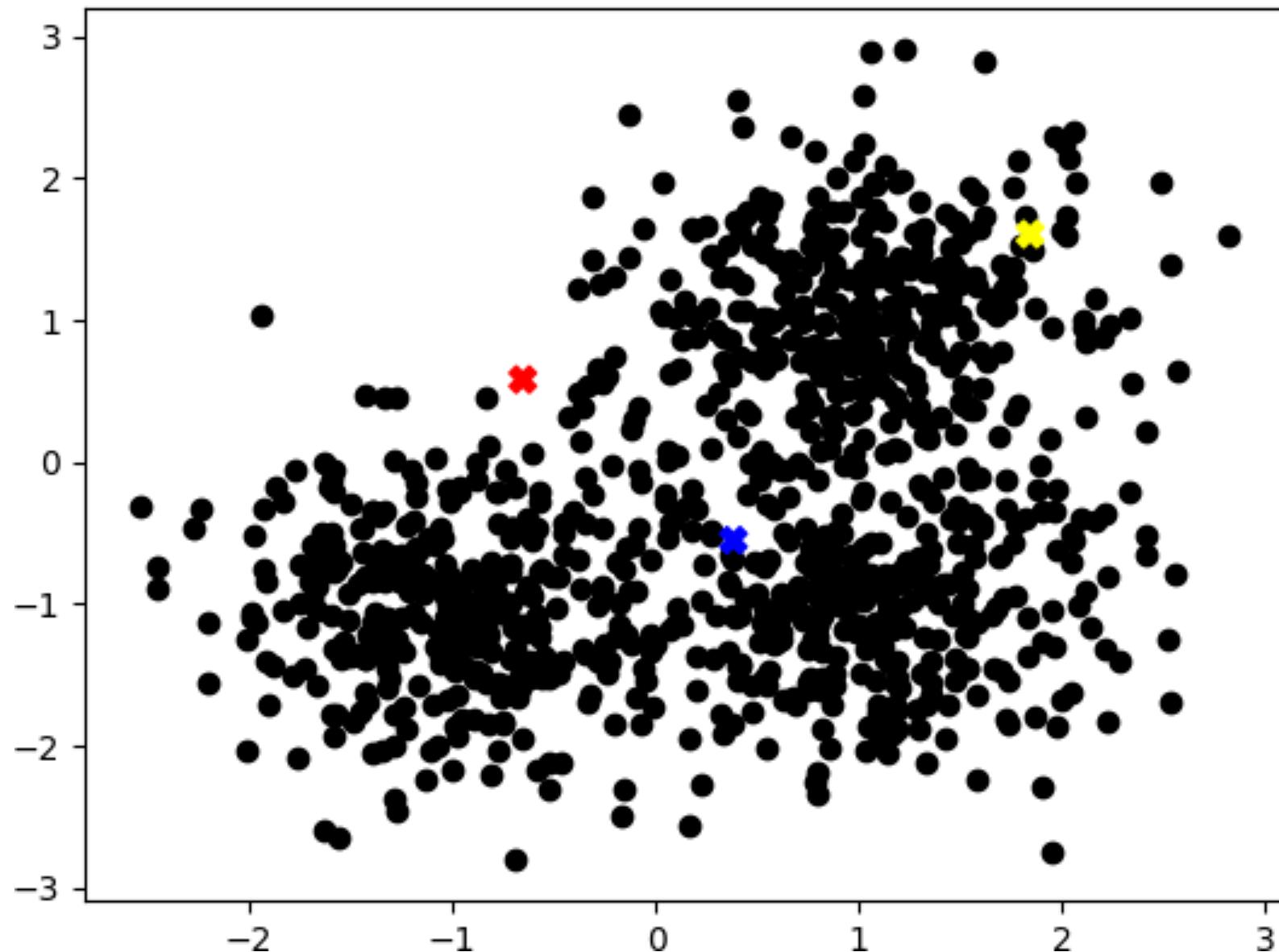
# K-MEANS ALGORITHM: EXAMPLE

K-means is a clustering algorithm that is best explained via an example



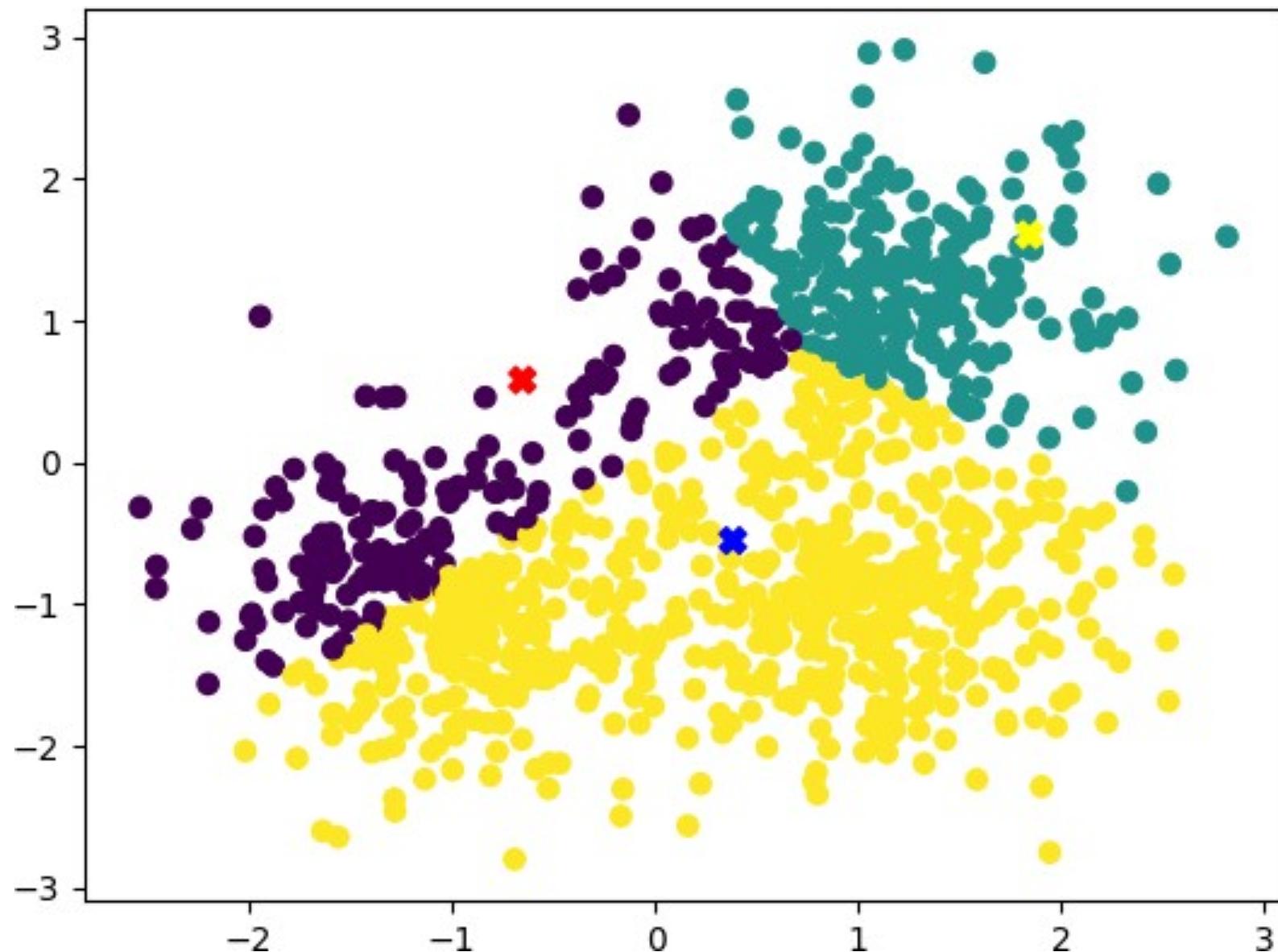
# K-MEANS ALGORITHM: EXAMPLE

Choose randomly 3  
centroids :  $c_1$ ,  $c_2$ ,  
 $c_3$  (the points  
appear in blue, red  
and yellow)



# K-MEANS ALGORITHM: EXAMPLE

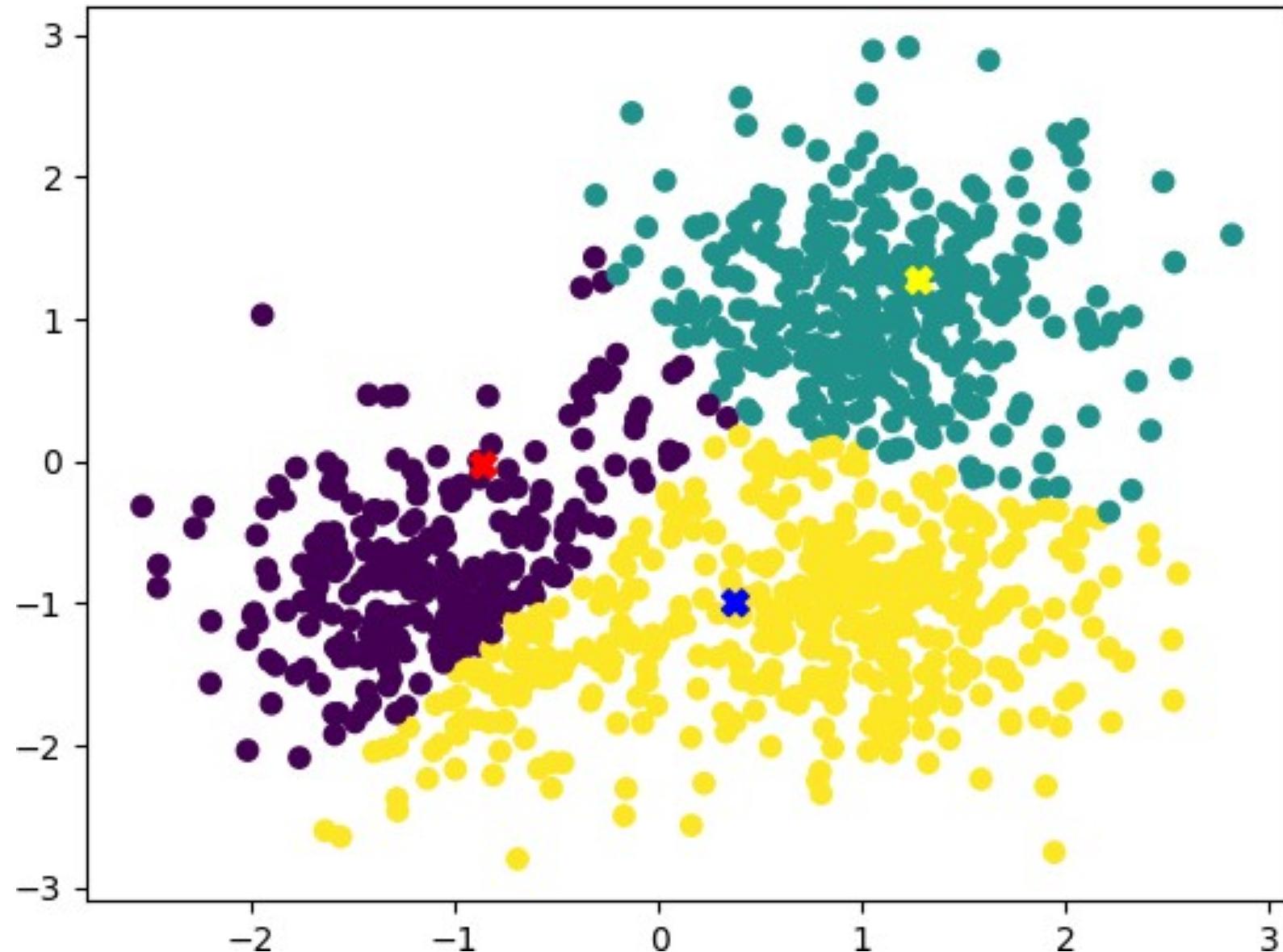
assign each point  $x$   
in the set to the  
closest centroid



# K-MEANS ALGORITHM: EXAMPLE

Update the  
centroids  $c_1, c_2,$   
 $c_3$  as follows :

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$



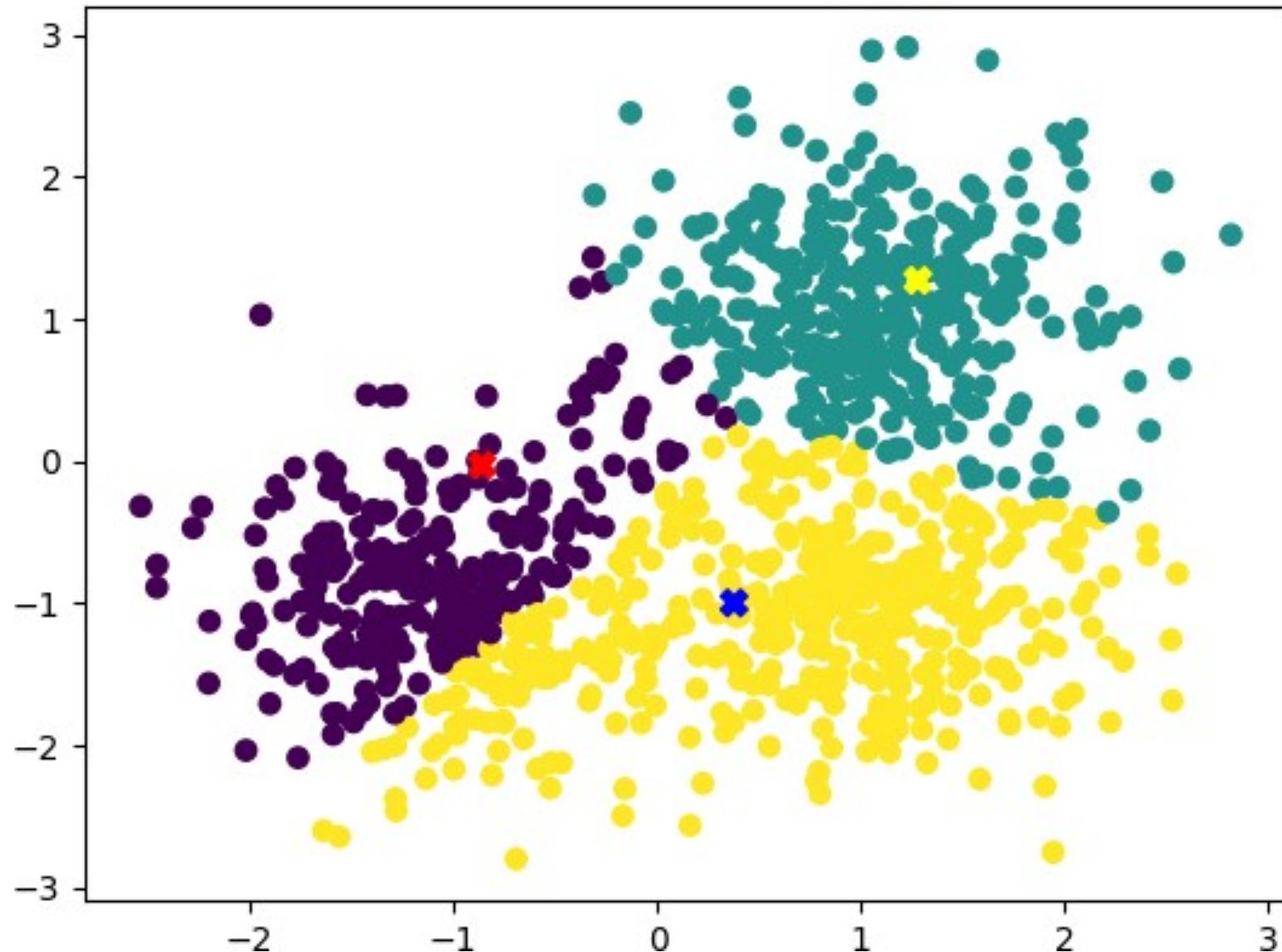
# K-MEANS ALGORITHM: EXAMPLE

Update the centroids  $c_1, c_2, c_3$  as follows :

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

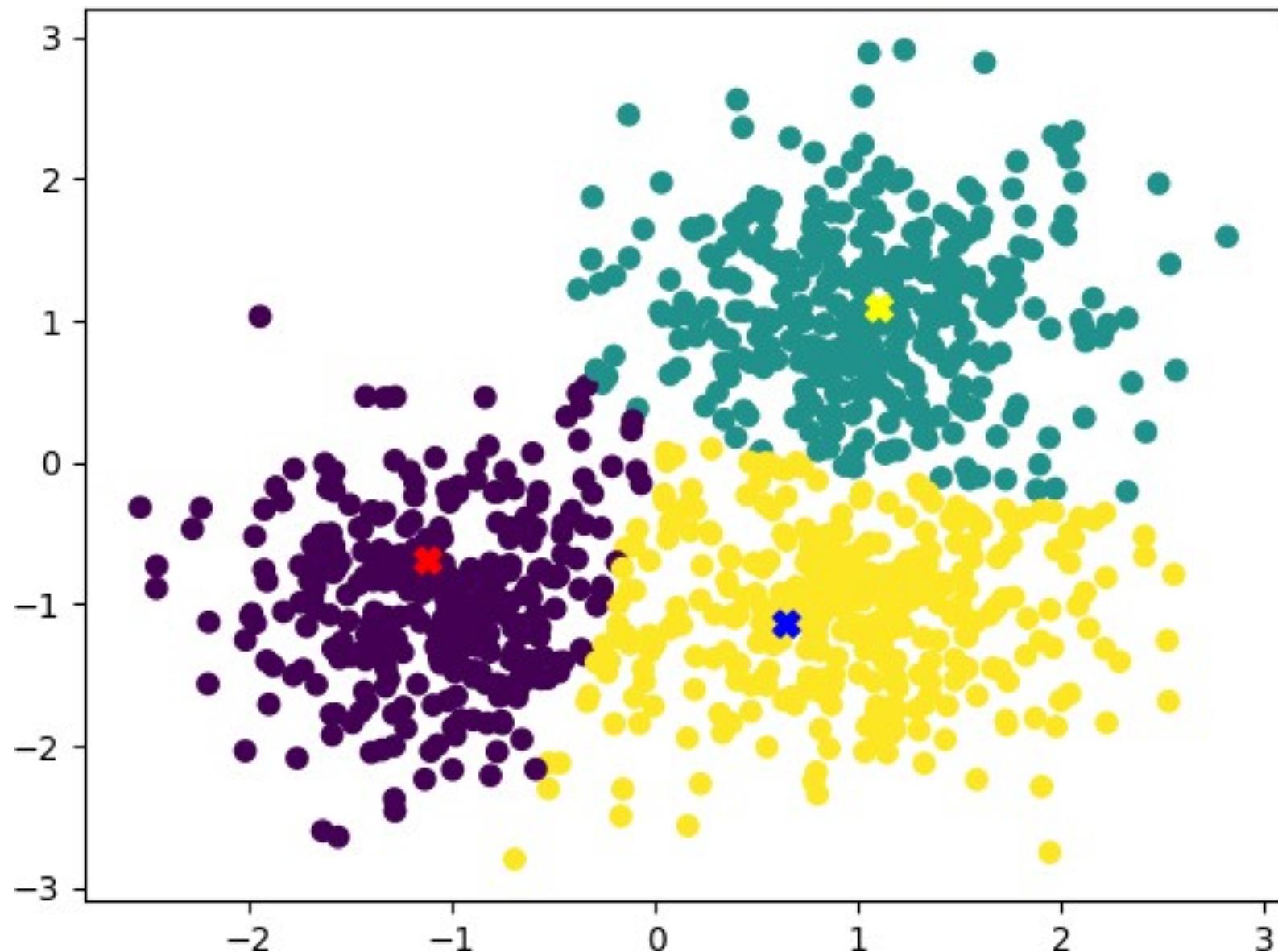
In other words, the new center is the average of the cluster members

assign each point  $x$  in the set to the closest centroid



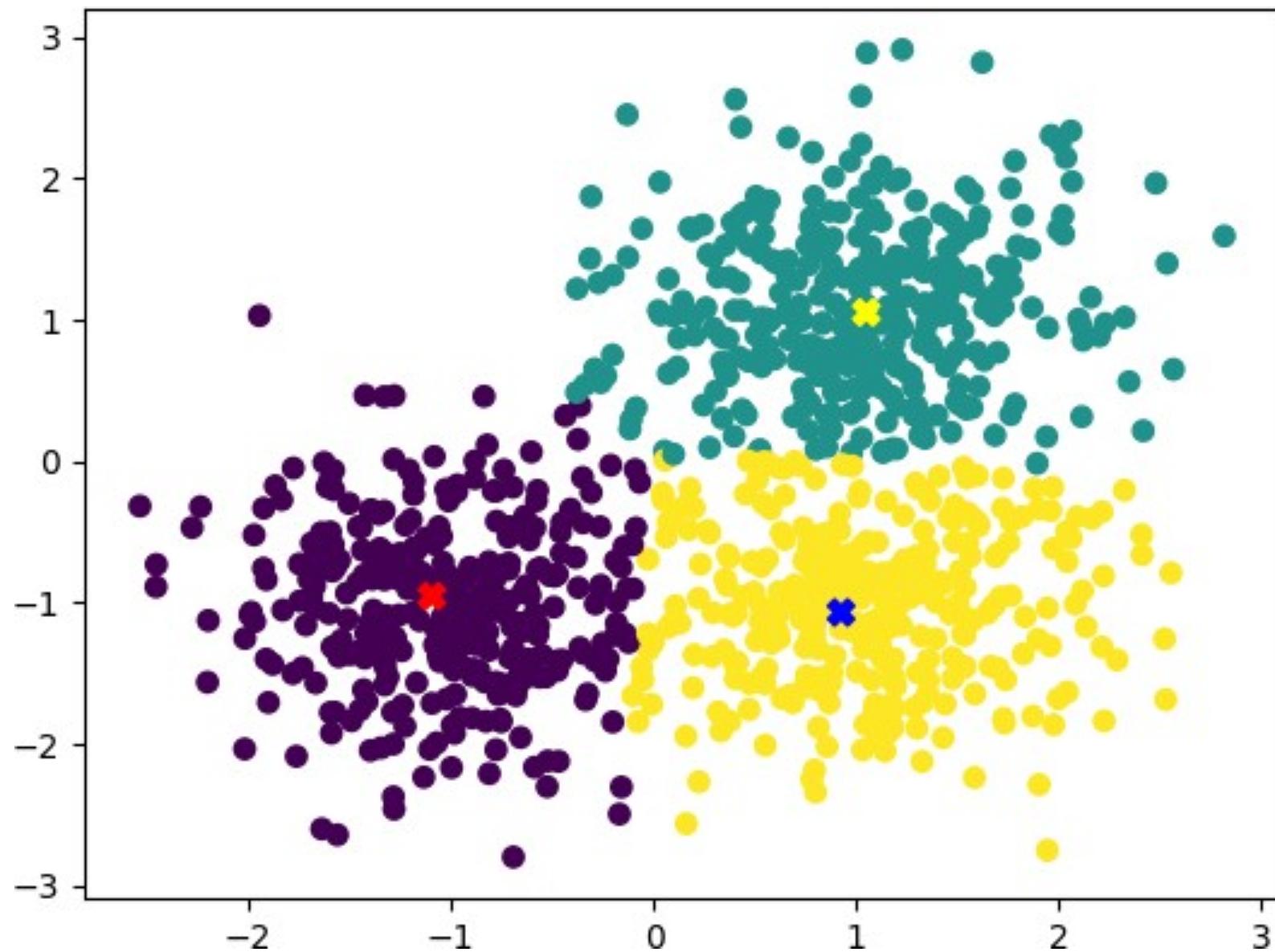
# K-MEANS ALGORITHM: EXAMPLE

Repeat until convergence



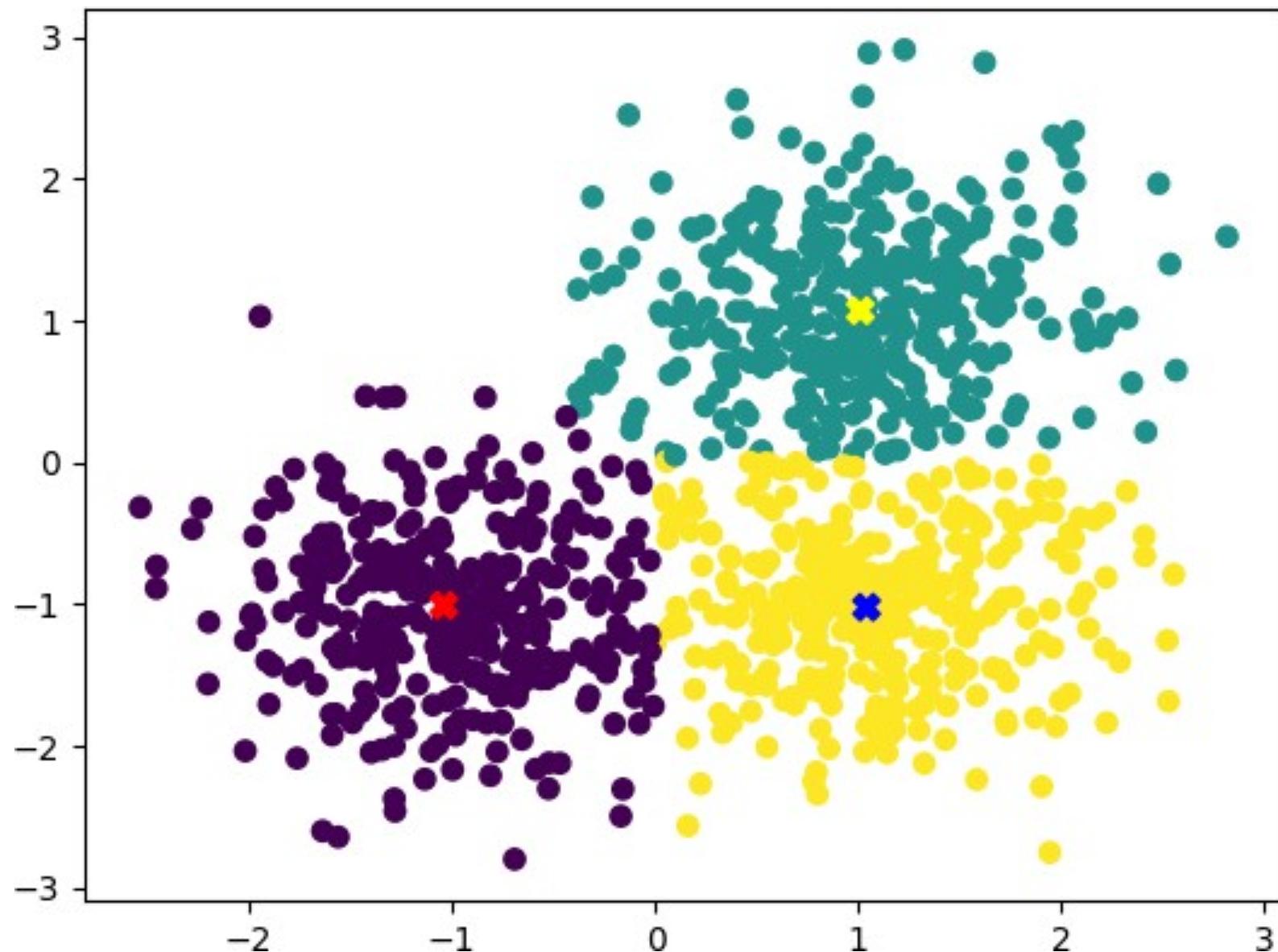
# K-MEANS ALGORITHM: EXAMPLE

Repeat until convergence



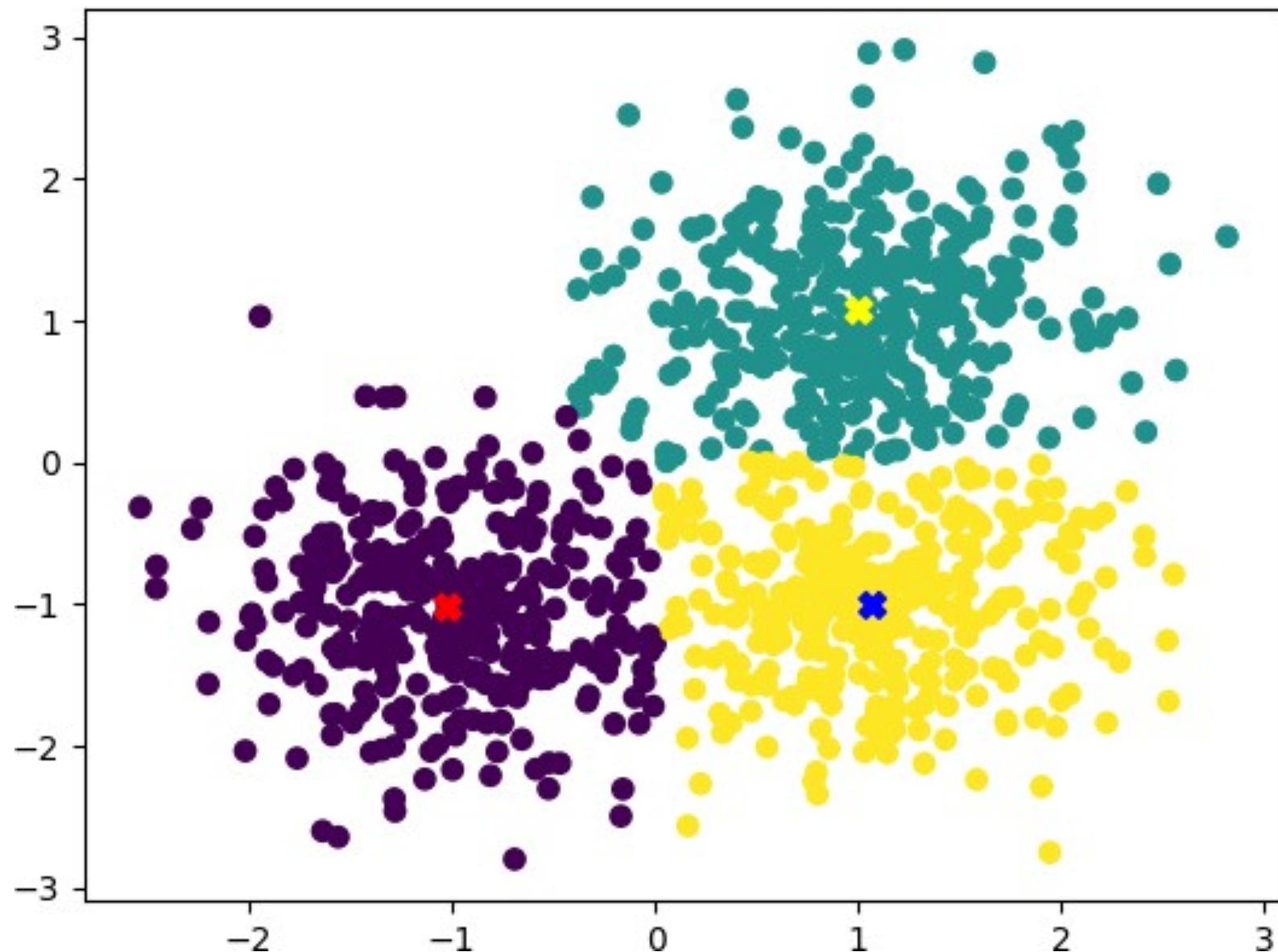
# K-MEANS ALGORITHM: EXAMPLE

Repeat until convergence



# K-MEANS ALGORITHM: EXAMPLE

Repeat until convergence



# NOTATIONS

- THE  $n^{th}$  EUCLIDIAN SPACE WILL BE DENOTED BY  $R^n$ . A POINT IN  $R^n$  WILL BE DENOTED BY  $x$ . IN THIS LECTURE THE TERM DATA, OR THE TRAINING SET,  $X$  WILL MEAN A FINITE SET OF POINTS IN  $R^n$ .
- IN OTHER WORDS,  $X = \{x_1, \dots, x_m\}$  WHERE  $x_j$  is a point in  $R^n$
- THE EUCLIDIAN DISTANCE BETWEEN TWO POINTS  $x$  AND  $y$  IN  $R^n$  WILL BE DENOTED BY  $d(x, y)$ .

# K-MEANS ALGORITHM

- THE K-MEANS ALGORITHM TAKES TWO INPUTS:
  - A parameter  $K$ , which is the number of clusters one wants to find in the data.
  - The training set  $X$  of the points. Here  $X = \{x_1, \dots, x_m\}$  where  $x_j$  is a point in  $R^n$ .
- THE ALGORITHM RETURNS A LABEL PARTITIONING  $X$  INTO  $K$ -CLUSTERS.

# K-MEANS ALGORITHM

- CHOOSE RANDOMLY K CENTROIDS :  $c_1, c_2, \dots, c_K$  IN  $R^n$
- REPEAT UNTIL CONVERGENCE
  - We assign each point  $x$  in the set to the closest centroid.  Cluster assignment step
  - Update the centroids  $c_1, c_2, \dots, c_K$  as follows:  Centroid move step
$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$
    - Here  $S_i$  is the cluster associated with the centroid  $c_i$
- Convergence:
  - None of the cluster assignments change
  - Centroids do not change

# K-MEANS ALGORITHM

- CHOOSE RANDOMLY K CENTROIDS :  $c_1, c_2, \dots, c_K$  IN  $R^n$
- REPEAT UNTIL CONVERGENCE
  - We assign each point  $x$  in the set to the closest centroid.
  - Update the centroids  $c_1, c_2, \dots, c_K$  as follows:

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

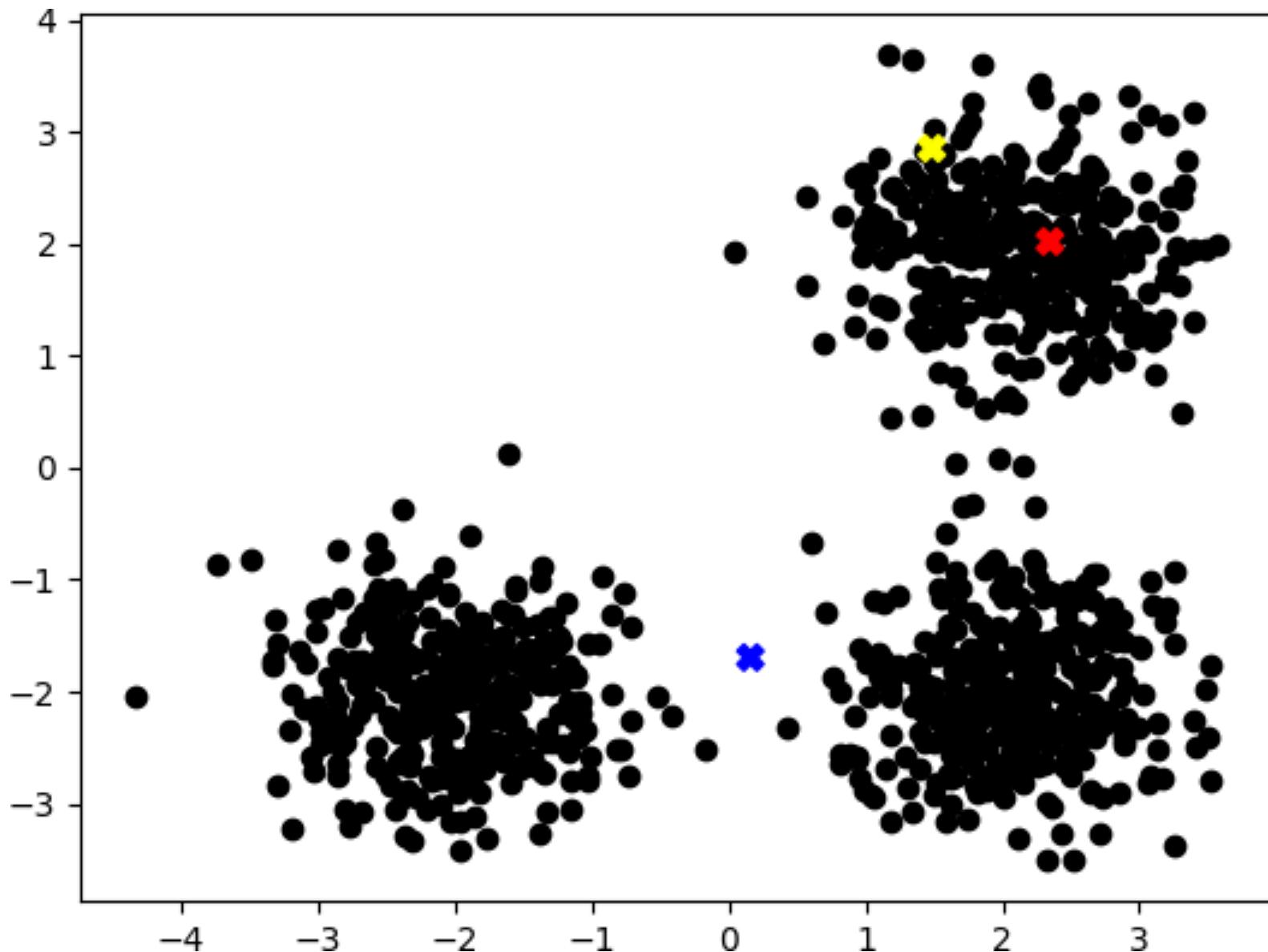
- Here  $S_i$  is the cluster associated with the centroid  $c_i$
- Convergence:
  - None of the cluster assignments change
  - Centroids do not change

This is where *mean* comes from

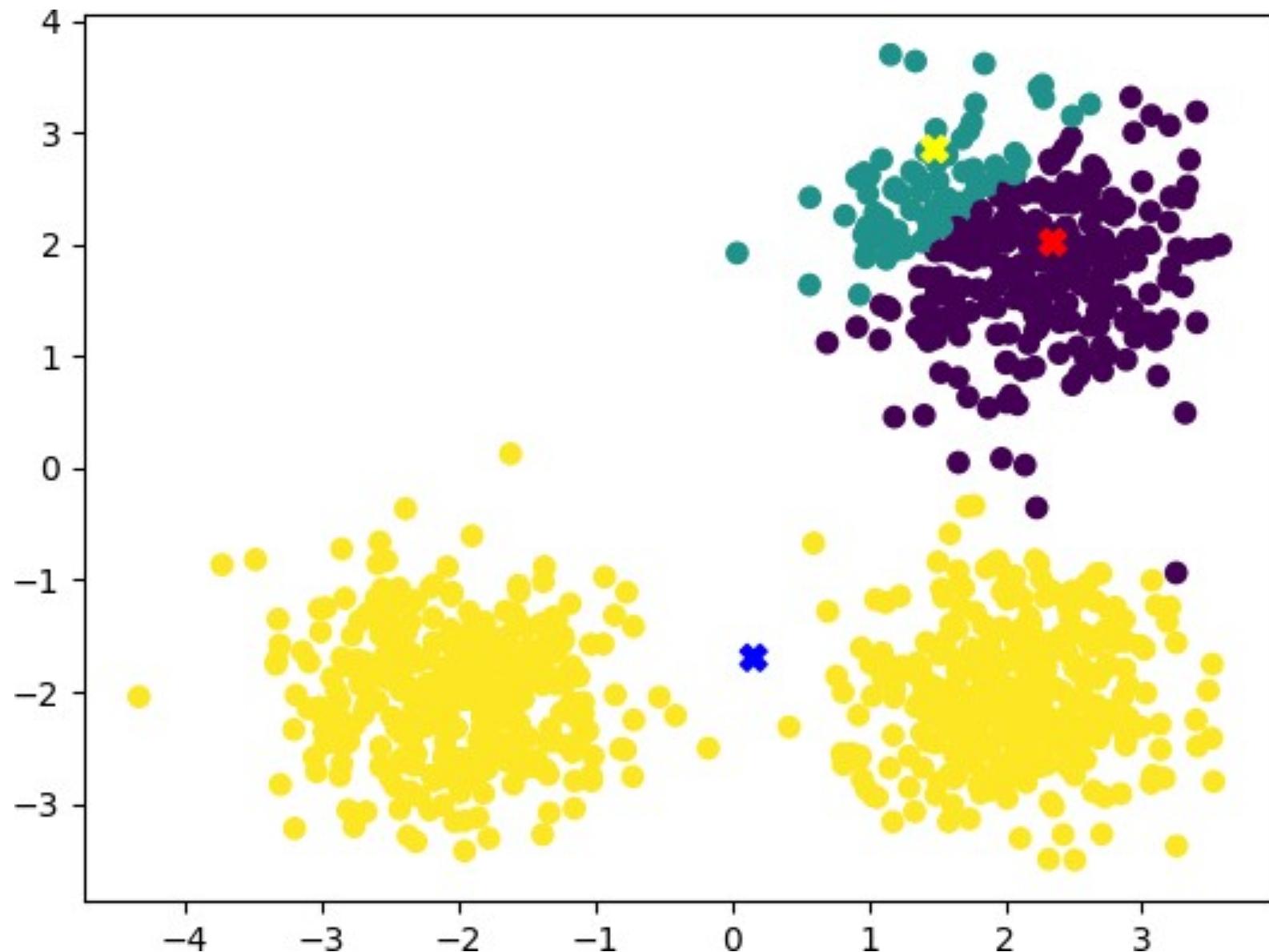


# K-MEANS ALGORITHM: EXAMPLE

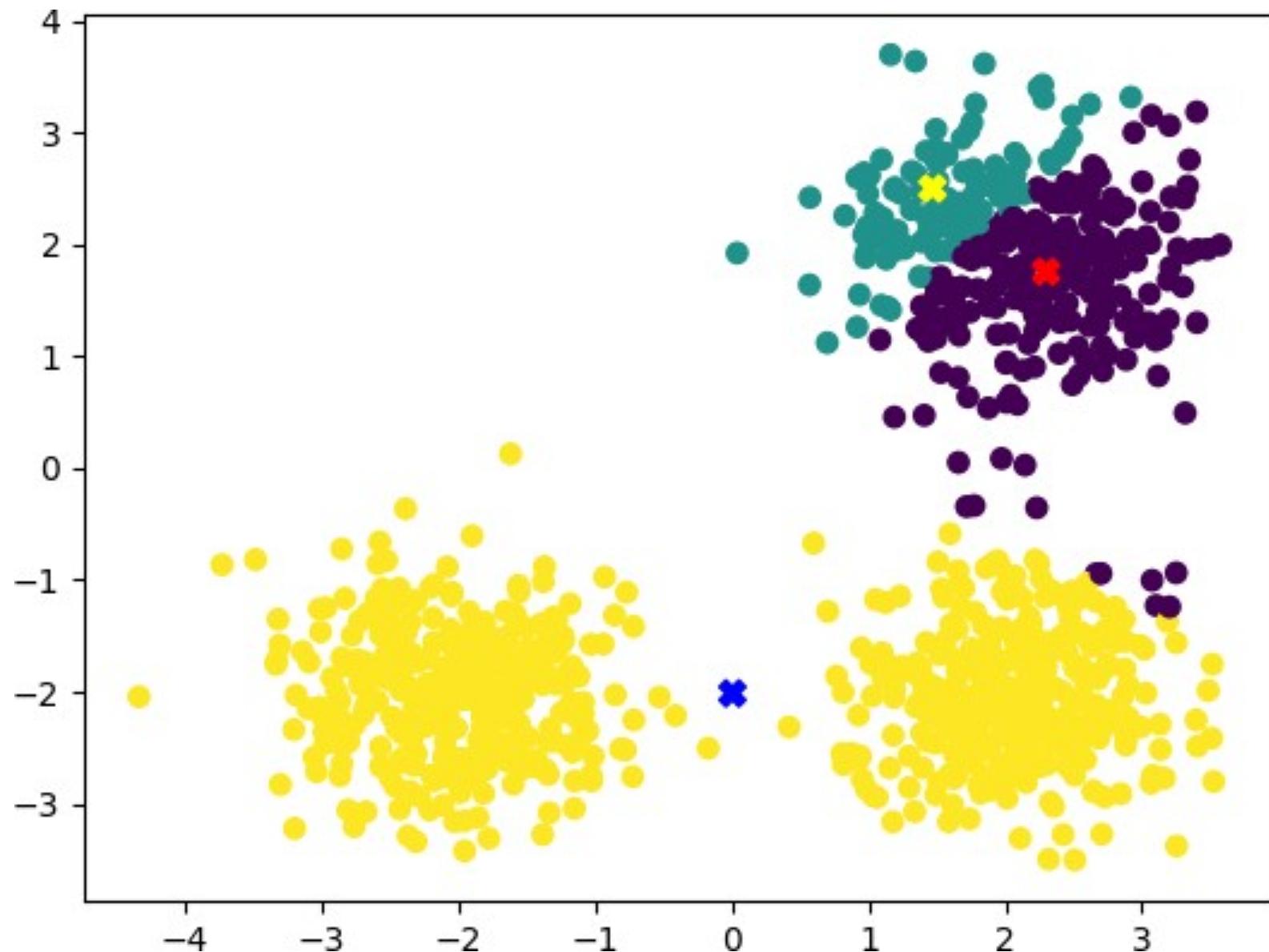
Lets consider  
another example



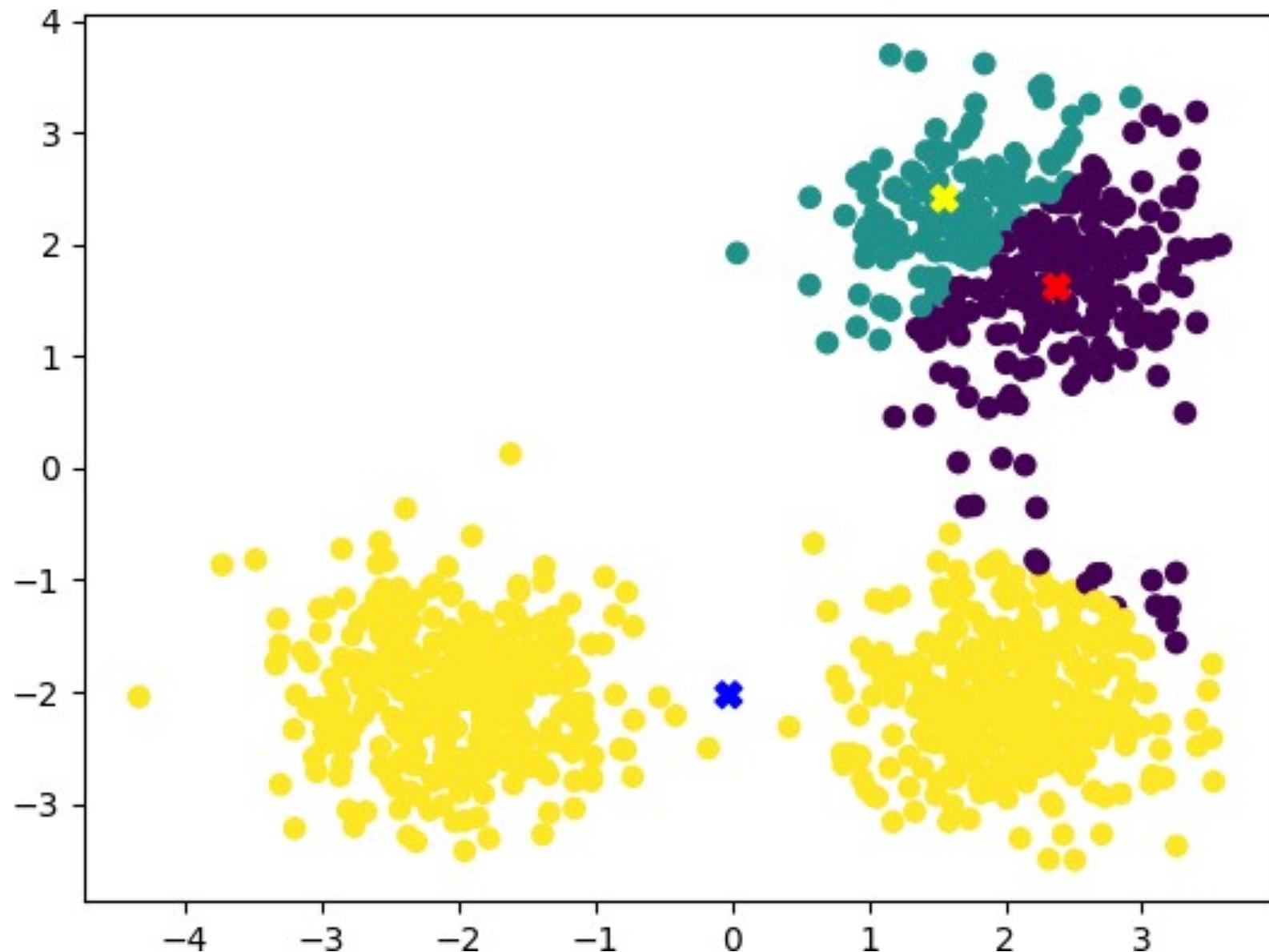
# K-MEANS ALGORITHM: EXAMPLE



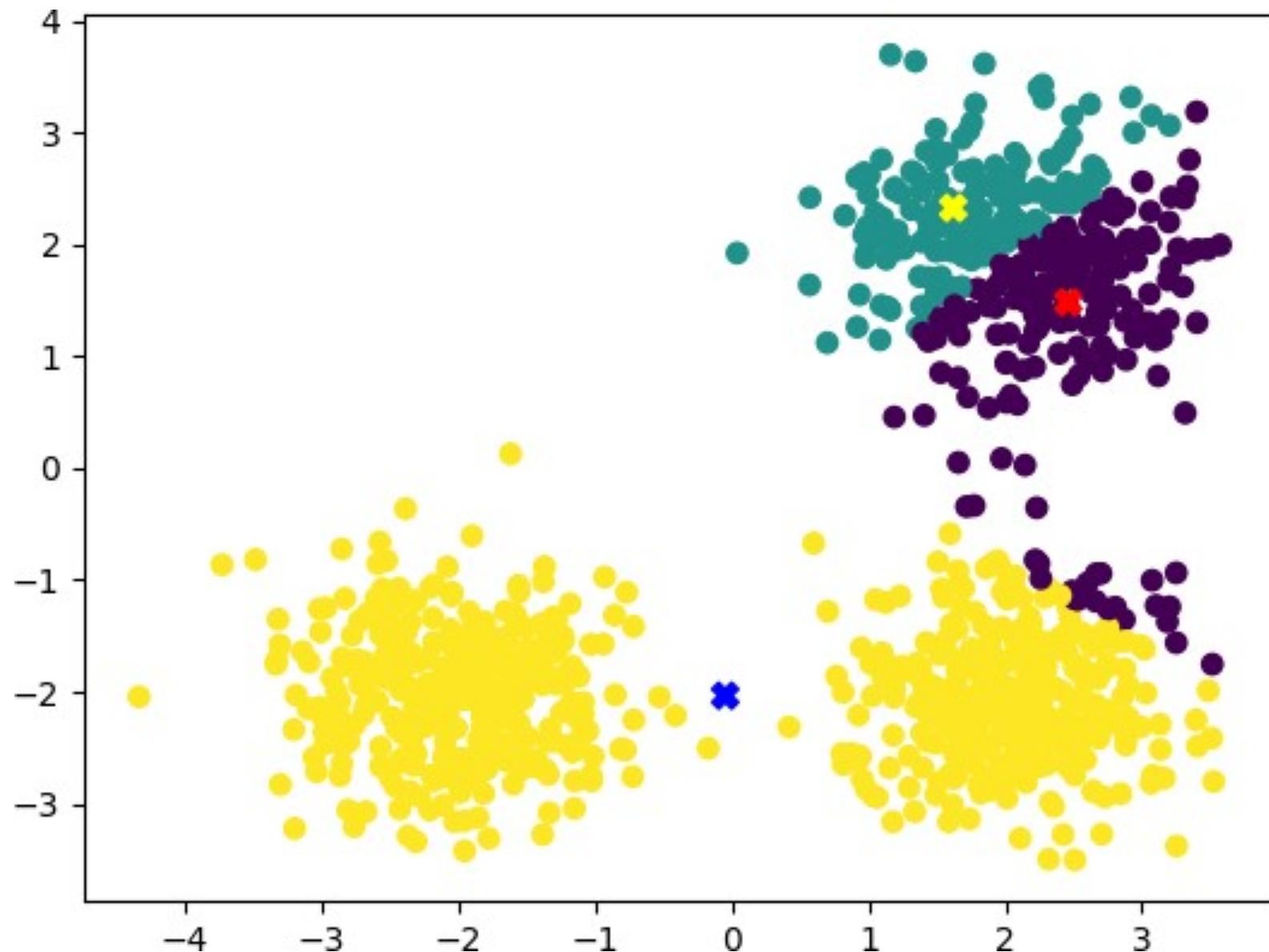
# K-MEANS ALGORITHM: EXAMPLE



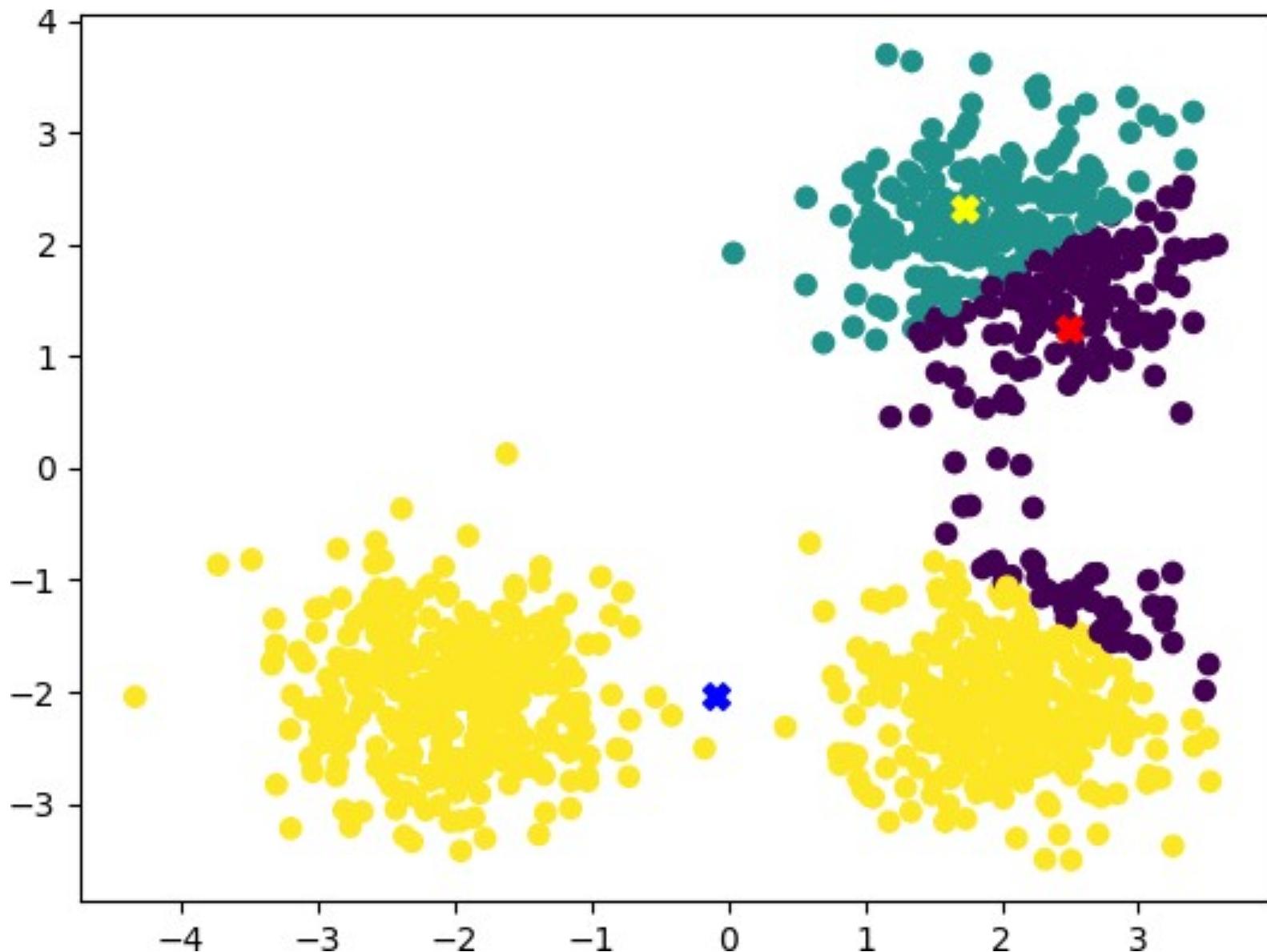
# K-MEANS ALGORITHM: EXAMPLE



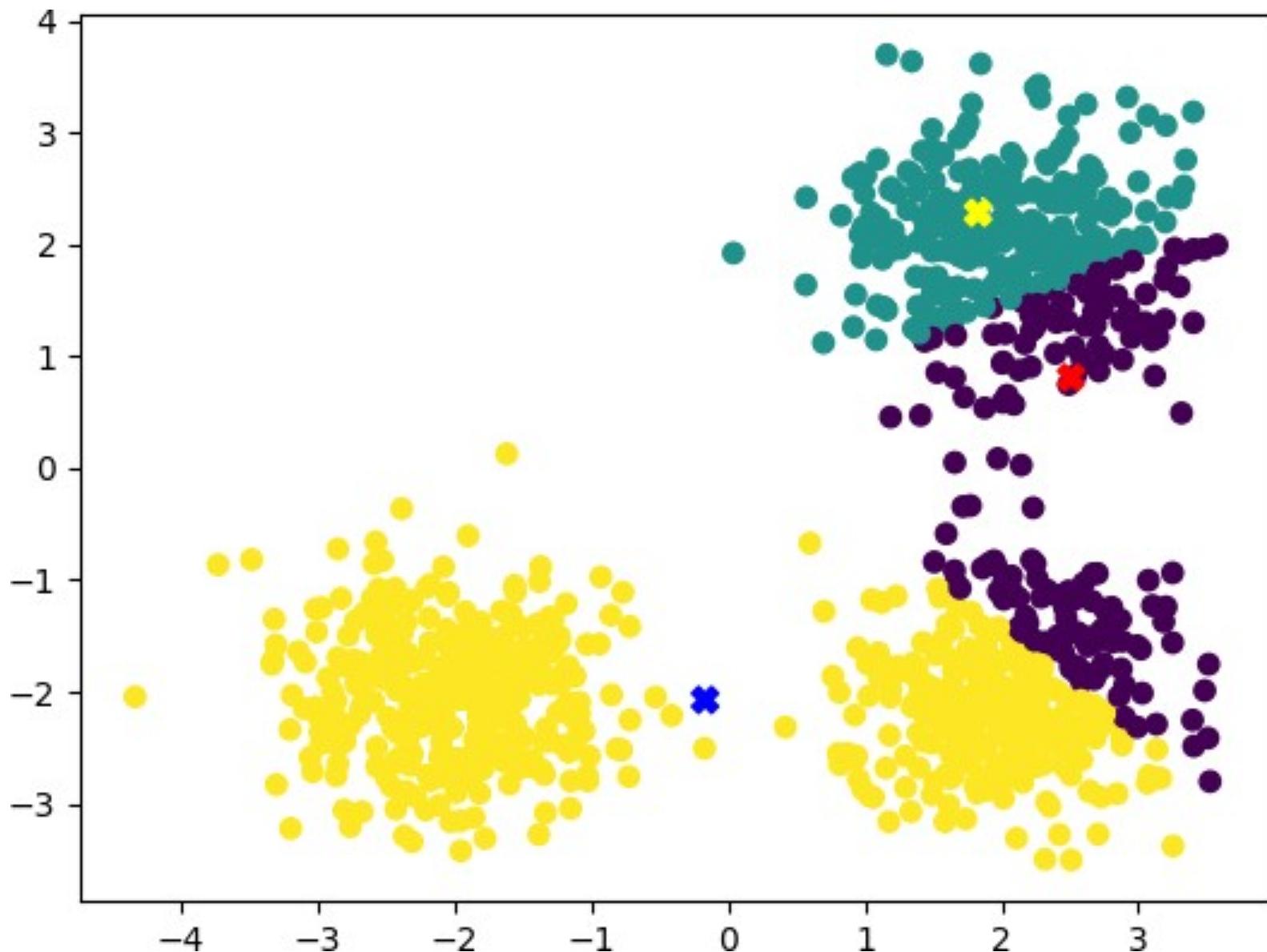
# K-MEANS ALGORITHM: EXAMPLE



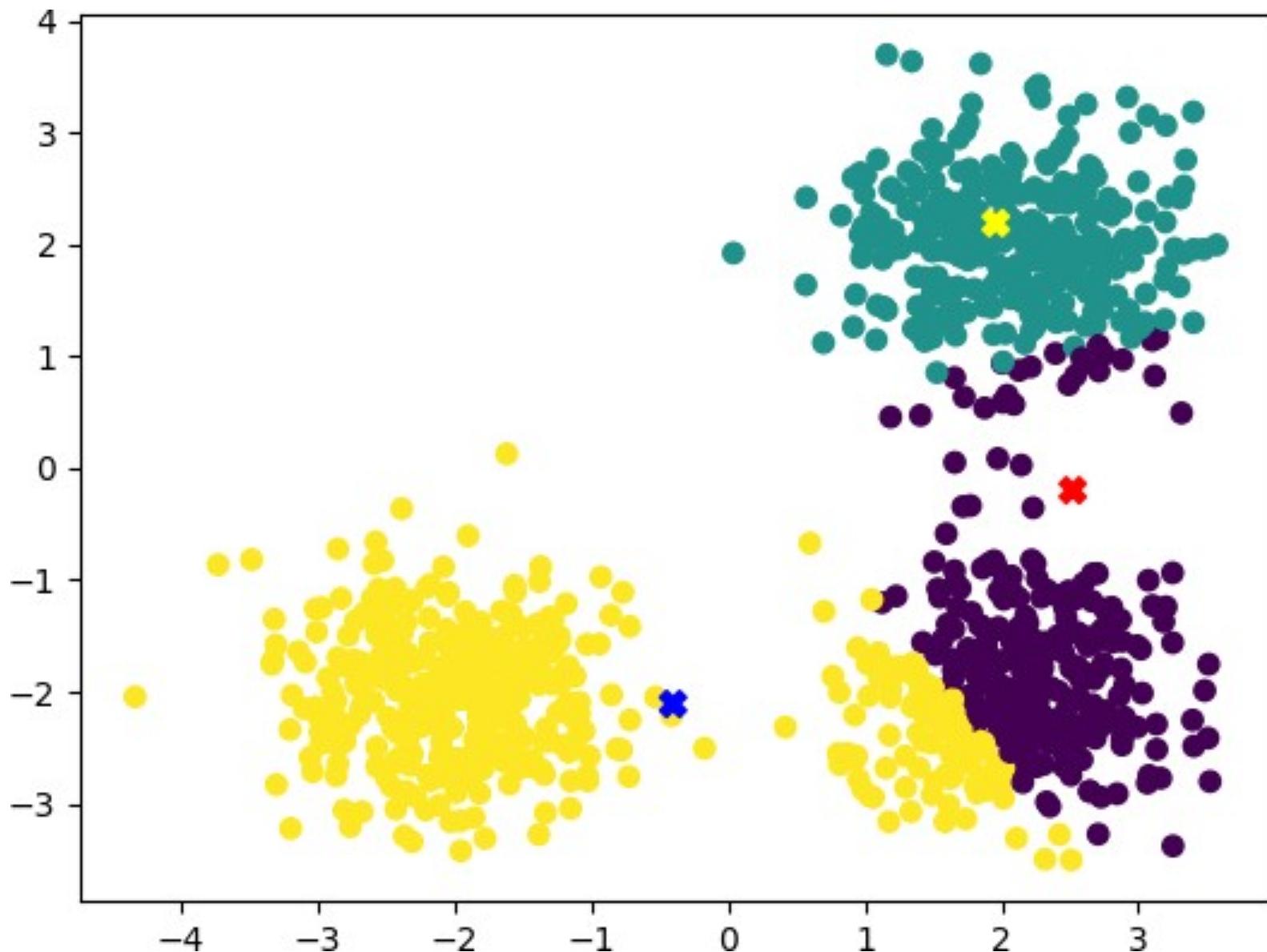
# K-MEANS ALGORITHM: EXAMPLE



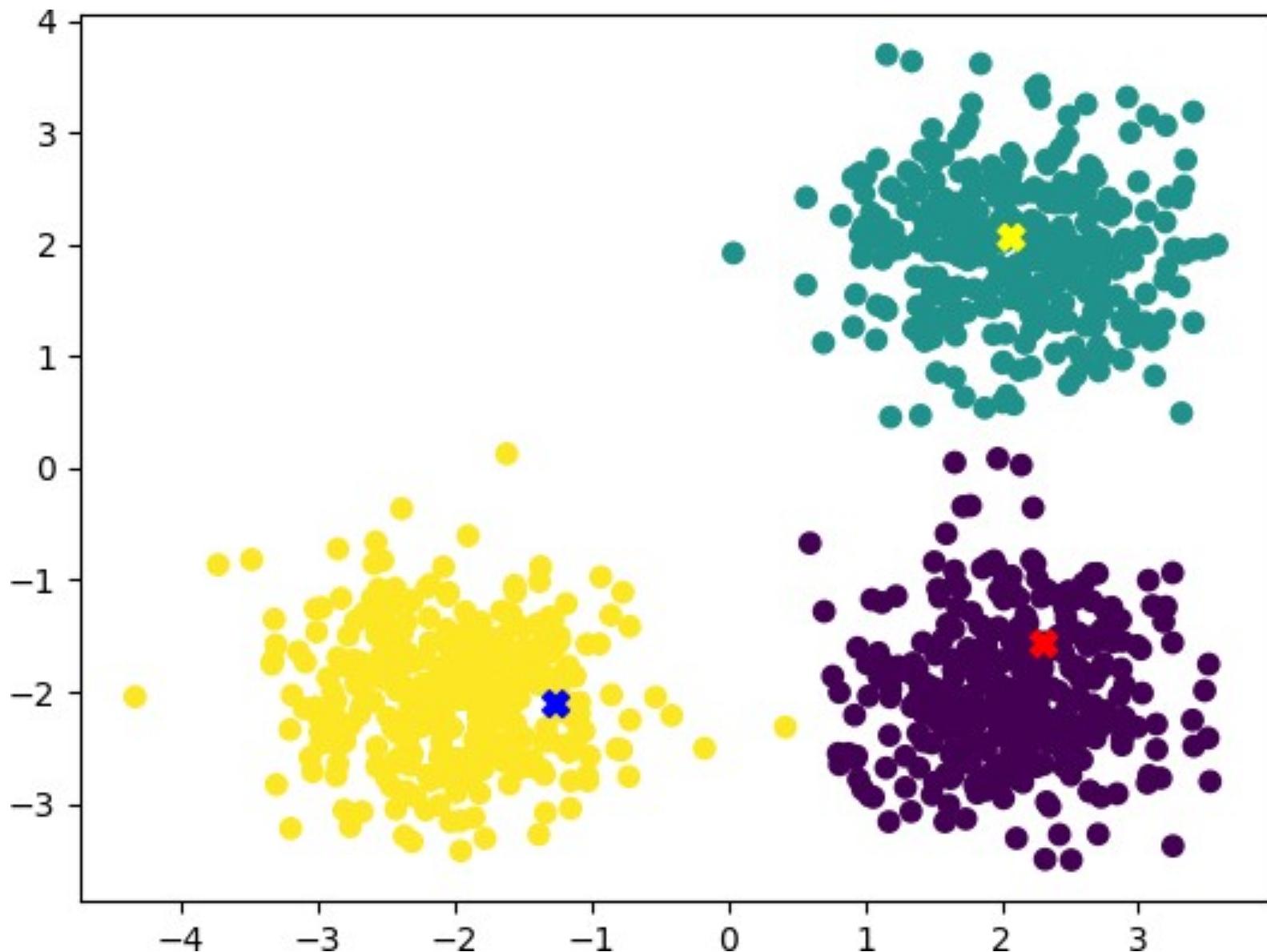
# K-MEANS ALGORITHM: EXAMPLE



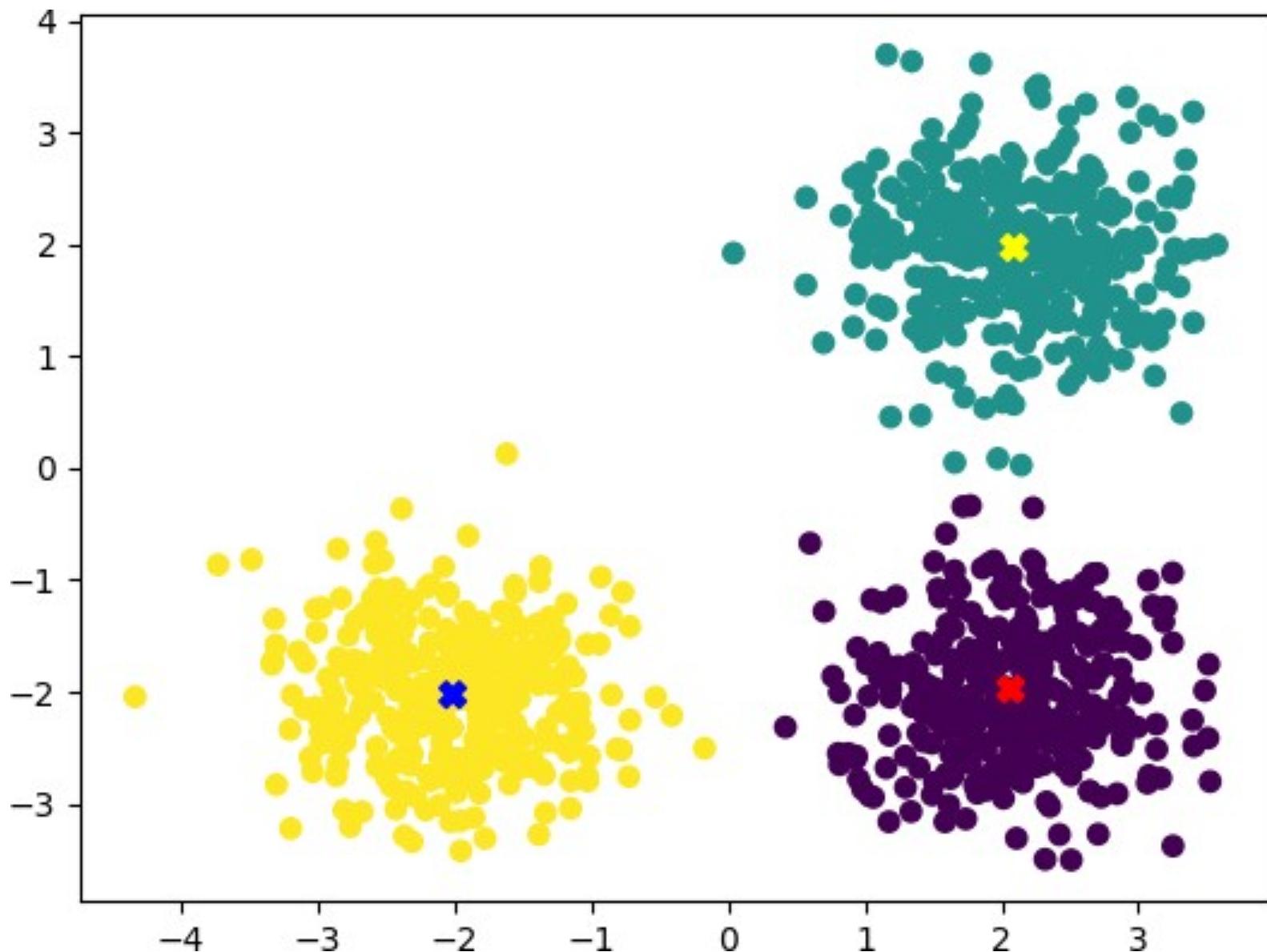
# K-MEANS ALGORITHM: EXAMPLE



# K-MEANS ALGORITHM: EXAMPLE



# K-MEANS ALGORITHM: EXAMPLE



## K-MEANS ALGORITHM: COMPLEXITY

- CHOOSE RANDOMLY K CENTROIDS :  $c_1, c_2, \dots, c_K$  IN  $R^n$
- REPEAT UNTIL CONVERGENCE
  - We assign each point  $x$  in the set to the closest centroid.
  - Update the centroids  $c_1, c_2, \dots, c_K$  as follows:

$$c_i = \frac{1}{|S_i|} \sum_{x_j \in S_i} x_j$$

- Here  $S_i$  is the cluster associated with the centroid  $c_i$

## K-MEANS ALGORITHM: COMPLEXITY

- COMPLEXITY IS  $O(m * K * I * n)$ 
  - $m$  = number of points in the data set
  - $K$  = number of clusters
  - $I$  = number of iterations
  - $n$  = number of attributes=number of features= the dimension of the space  $R^n$
- How CAN WE DO BETTER?

## K-MEANS ALGORITHM: CONVERGENCE

- WHAT EXACTLY IS THE OPTIMIZATION FUNCTION OF THIS ALGORITHM ?
  - Cost function of K-means:  $\sum_i \sum_{x \in S_i} d(c_i, x)$
  - *THIS IS THE TOTAL SQUARED DISTANCE FROM THE CENTROID TO THE POINTS OF THE CLUSTER ASSOCIATED TO THE CENTROID*

## K-MEANS ALGORITHM: CONVERGENCE

- SINCE WE START WITH RANDOM CENTERS EVERY TIME WE RUN THIS ALGORITHM:
  - Is it guaranteed to give the same clustering configuration?
  - Is the algorithm guaranteed to converge ?

## K-MEANS ALGORITHM: CONVERGENCE

- THE ALGORITHM CONVERGES (IN THE SENSE THAT EACH ITERATION MINIMIZES THE COST FUNCTION ABOVE) BUT IT CONVERGES TO A LOCAL MIN.
- WHICH MEANS THAT:
  - (1) the solution might not be the optimal solution and
  - (2) one might get different results for different initial starts.



# K-MEANS ALGORITHM: CONVERGENCE

- WHAT EXACTLY IS THE OPTIMIZATION FUNCTION OF THIS ALGORITHM ?
  - Cost function of K-means:  $\sum_i \sum_{x \in S_i} d(c_i, x)$
  - *THIS IS THE TOTAL SQUARED DISTANCE FROM THE CENTROID TO THE POINTS OF THE CLUSTER ASSOCIATED TO THE CENTROID*

Choose randomly k centroids :  $c_1, c_2, \dots, c_K$  in  $R^n$   
Repeat until convergence

Assign each point  $x$  to the closest centroid

Update the centroids  $c_1, c_2, \dots, c_K$

Minimize the cost  
function with respect  
to the clusters

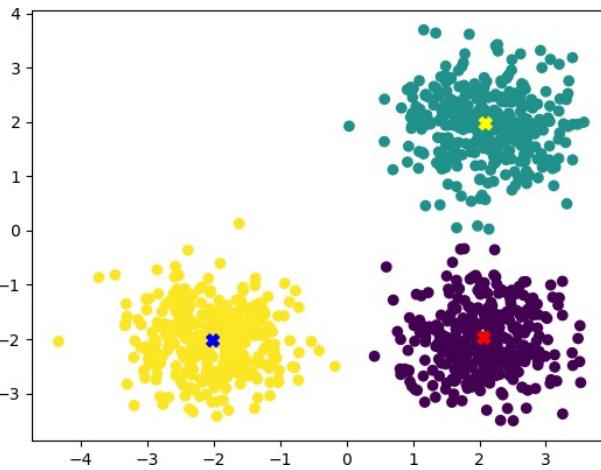
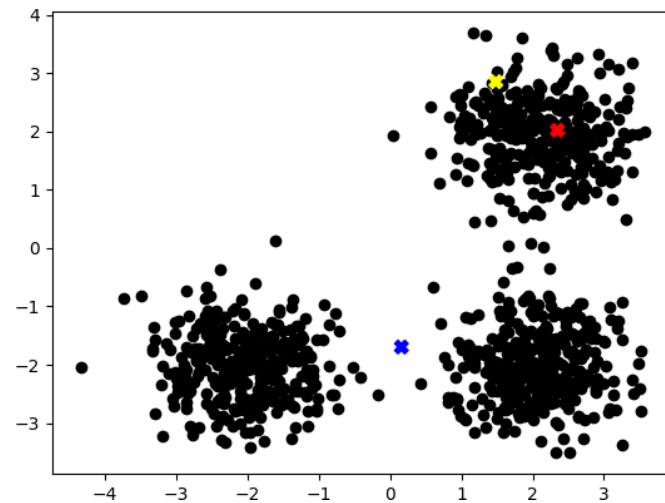
Minimize the cost  
function with respect  
to the centroids



# K-MEANS ALGORITHM: PROBLEMS

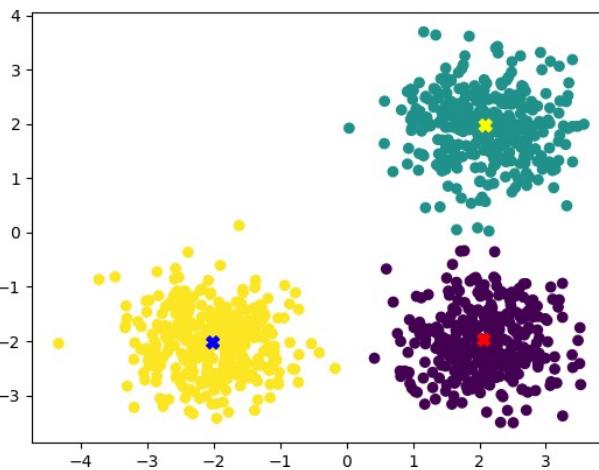
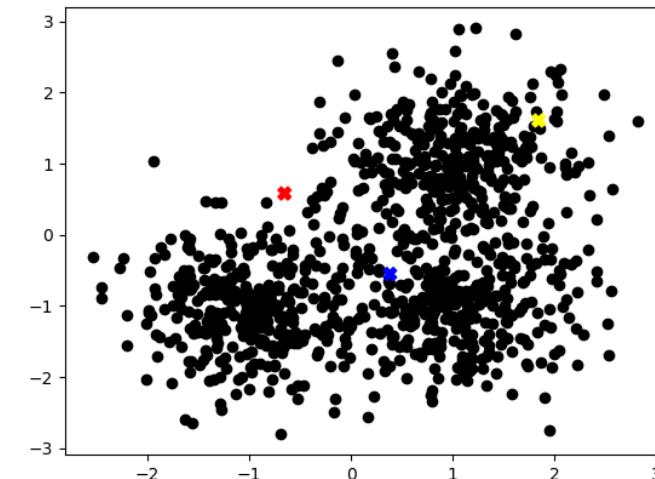
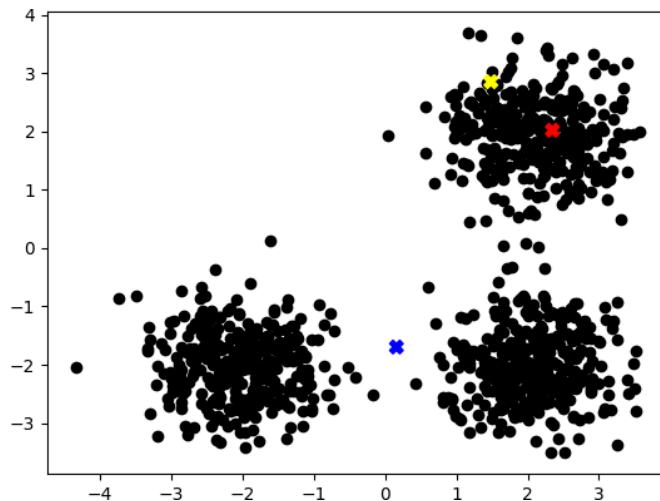
- OUTLIERS
- CLUSTERS WITH DIFFERENT DENSITIES
- NON-CONVEX SHAPES
- CLUSTERS WITH DIFFERENT SIZES
- MAY CONVERGE TO LOCAL OPTIMUM

# K-MEANS ALGORITHM: REMARKS

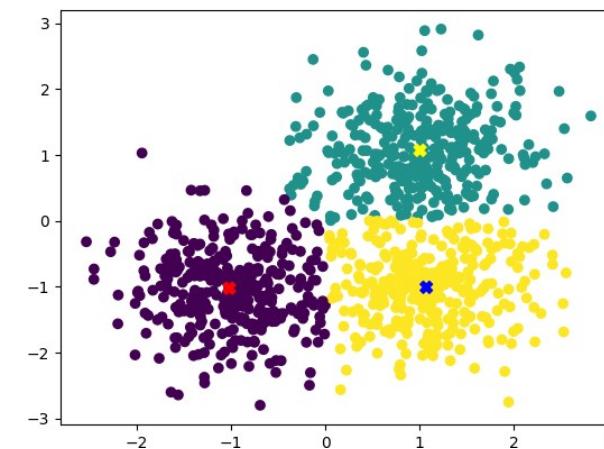


K-means is clearly applicable to data sets where the clusters are very well-separated

# K-MEANS ALGORITHM: REMARKS



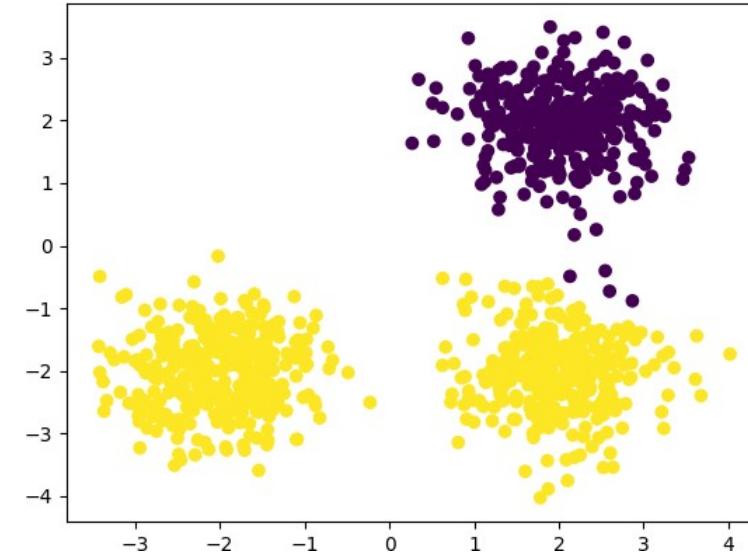
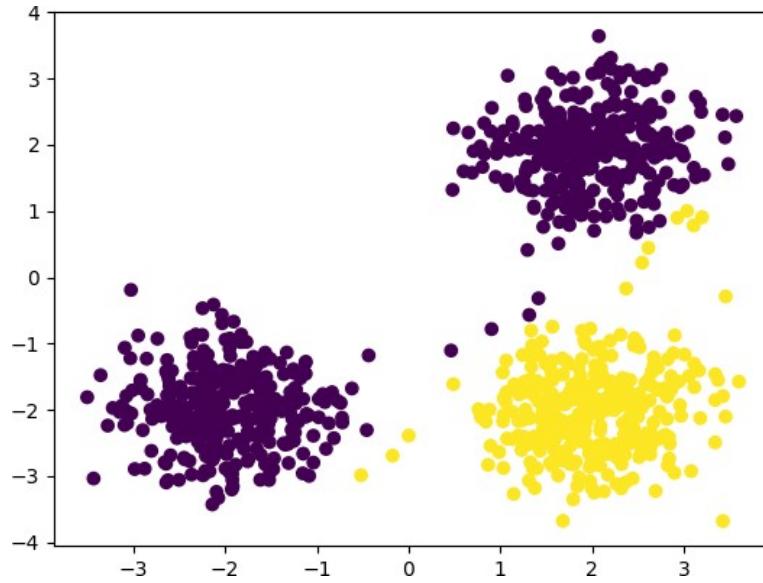
K-means is clearly applicable to data sets where the clusters are very well-separated



K-Means is often applied to data that have no clear clustering structure



# K-MEANS ALGORITHM: REMARKS



- SELECTING K IS CHALLENGING
- ONE SOLUTION FOR THIS IS TO RUN K-MEANS SEVERAL TIMES AND PICK THE ATTEMPT THAT MINIMIZES THE COST FUNCTION

## K-MEANS ALGORITHM: REMARKS

- NEARBY POINTS MAY NOT END IN THE SAME CLUSTER



## K-MEANS ALGORITHM: REMARKS

- NEARBY POINTS MAY NOT END IN THE SAME CLUSTER



- IT IS POSSIBLE THAT K-MEANS GETS STUCK IN THIS LOCAL OPTIMUM.

## K-MEANS ALGORITHM: REMARKS

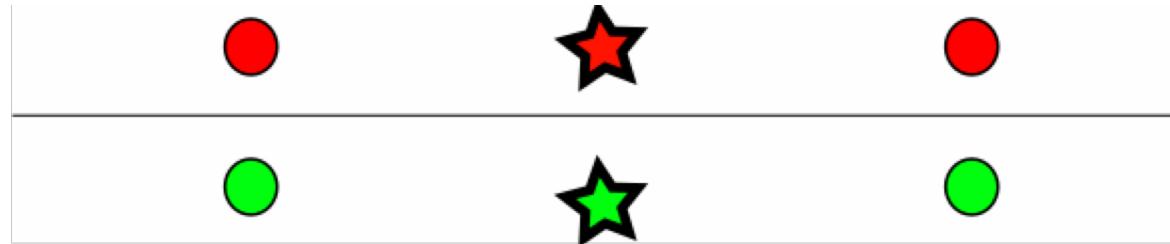
- NEARBY POINTS MAY NOT END IN THE SAME CLUSTER



- IT IS POSSIBLE THAT K-MEANS GETS STUCK IN THIS LOCAL OPTIMUM.
- How WOULD YOU SOLVE THIS PROBLEM ?

## K-MEANS ALGORITHM: REMARKS

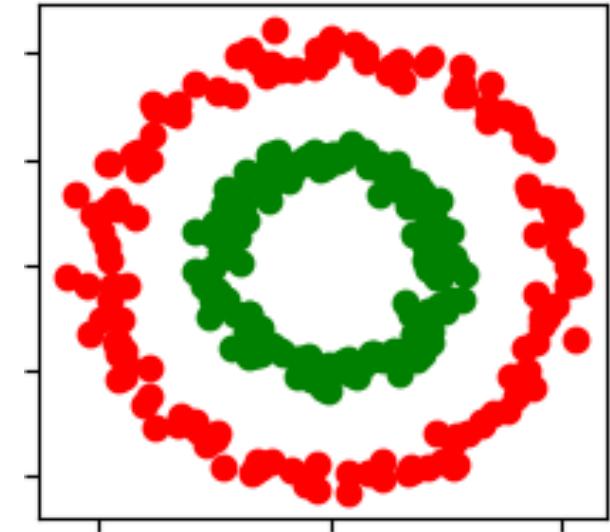
- NEARBY POINTS MAY NOT END IN THE SAME CLUSTER



- IT IS POSSIBLE THAT K-MEANS GETS STUCK IN THIS LOCAL OPTIMUM.
- How WOULD YOU SOLVE THIS PROBLEM ?
  - Try different initialization and choose the one that produces the best result (wrt the cost function).

## K-MEANS ALGORITHM: REMARKS

- USING EUCLIDIAN DISTANCE K-MEANS WILL NOT GIVE THE NATURAL CLUSTERS FOR THIS SET.



## K-MEANS ALGORITHM: REMARKS

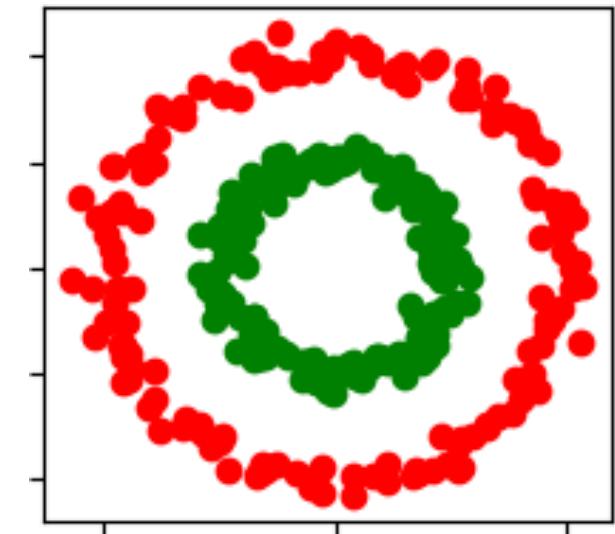
- USING EUCLIDIAN DISTANCE K-MEANS WILL NOT GIVE THE NATURAL CLUSTERS FOR THIS SET.

- ONE CAN TRY TO CHANGE THE FEATURES:

$$(x, y) \rightarrow \left( \sqrt{x^2 + y^2}, \tan^{-1} \frac{y}{x} \right)$$

(CONVERT THE FEATURE VECTOR TO POLAR COORDINATES)

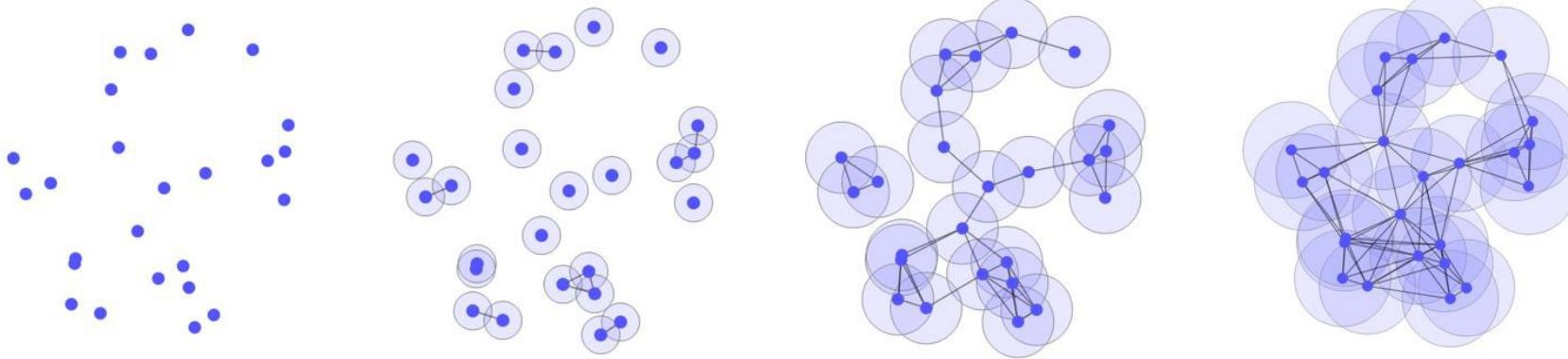
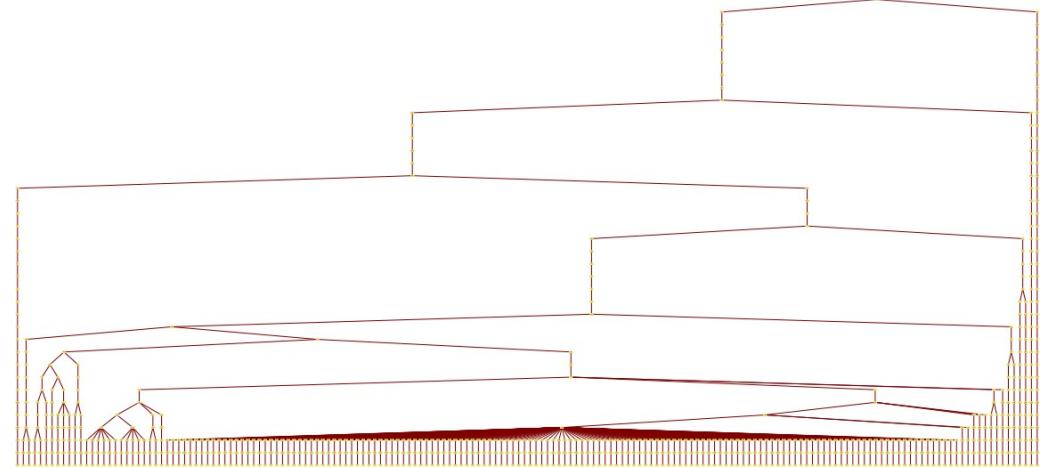
- OR CHANGE THE DISTANCE FUNCTION USED IN THE K-MEANS ALGORITHM.



## K-MEANS VARIATIONS

- ALMOST EVERY ASPECT OF K-MEANS HAS BEEN ALTERED AND CHANGED TO PERFORM OTHER CLUSTERING TASKS.
  - Distance function: Any function that satisfy distance axioms can be used instead of the Euclidian distance.
  - Cost function
  - Initialization heuristics
  - Efficiency
  - Centroid definition: [K-Medians, K-mediods](#)

# HIERARCHICAL CLUSTERING



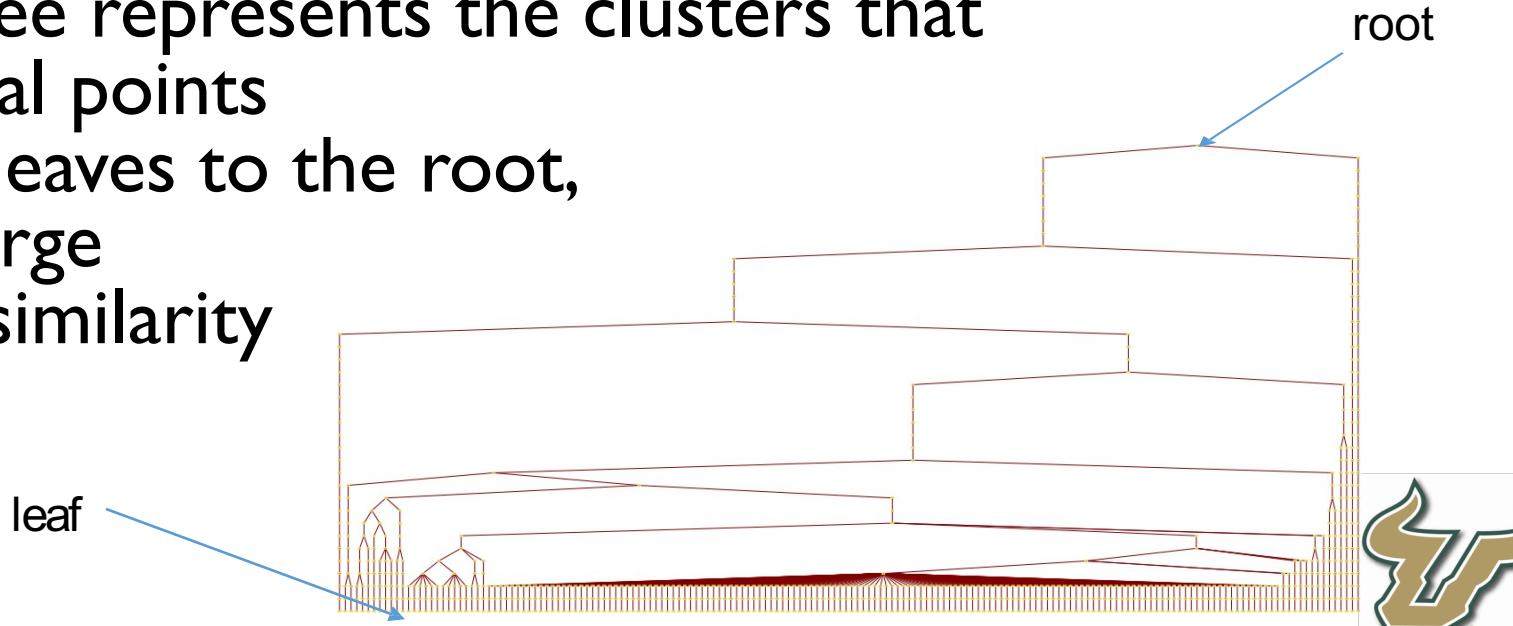
# HIERARCHICAL CLUSTERING

- HIERARCHICAL CLUSTERING IS A FAMILY OF CLUSTERING ALGORITHMS THAT BUILD A TREE CLUSTERS.
- IT IS USUALLY DONE BY MERGING OR SPLITTING THE CLUSTERS SUCCESSIVELY.
  - Agglomerative: bottom up (Merge)
  - Divisive: top down (Split)



# HIERARCHICAL CLUSTERING

- HIERARCHICAL CLUSTERING IS USUALLY REPRESENTED BY A TREE CALLED THE DENDROGRAM THAT REPRESENTS THE CLUSTERING AT ALL LEVELS.
- EACH NODE IN THE TREE REPRESENTS A CLUSTER
  - The root of the tree represents the cluster that contains all points
  - The leaves of the tree represents the clusters that contain the individual points
  - As we go from the leaves to the root, clusters start to merge according to some similarity criterion.



## GENERAL STEPS

- **GENERAL STEPS IN A STANDARD HIERARCHICAL AGGLOMERATIVE CLUSTERING ALGORITHM:**
  1. Compute the distance matrix, or the dissimilarity, between all points in the input data.
    - Choice of metric will impact the results largely.
  2. Initialize every point in the dataset to be its own cluster
  3. Compute the distance between all clusters
  4. Combine the closest two clusters: the two clusters  $c_i$  and  $c_j$  with  $\min_{i,j} d(c_i, c_j)$
  5. Remove the clusters  $c_i$  and  $c_j$  and add the cluster  $c_i + c_j$ .
  6. Go back to 3 and repeat until we have a single cluster.



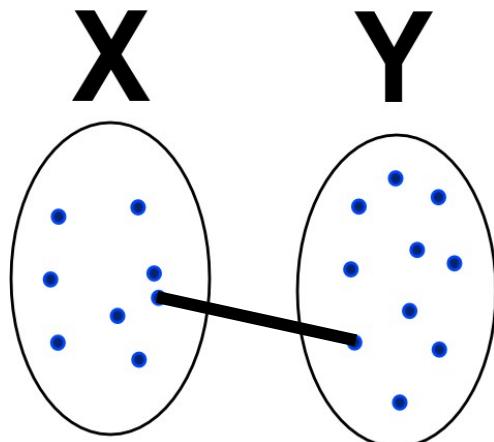
## GENERAL STEPS

- QUESTION : GIVEN TWO CLUSTERS X AND Y. HOW DO WE MEASURE THE DISTANCE BETWEEN THEM?
  - This is important because step 3 we need to measure the distance between clusters rather than points.



# DISTANCE BETWEEN CLUSTERS

- GIVEN TWO CLUSTERS  $X$  AND  $Y$ .
- HOW DO WE MEASURE THE DISTANCE BETWEEN THEM ?



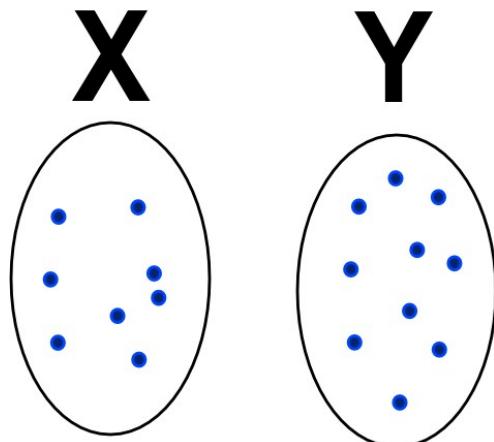
$$D(X, Y) = \min_{x \in X, y \in Y} d(x, y)$$

- ONE WAY IS TO MEASURE THE MINIMAL DISTANCE BETWEEN ALL POINTS OF  $X$  AND  $Y$ .
- THIS DISTANCE INDUCE SINGLE LINKAGE CLUSTERING.



# DISTANCE BETWEEN CLUSTERS

- GIVEN TWO CLUSTERS X AND Y.
- HOW DO WE MEASURE THE DISTANCE BETWEEN THEM ?



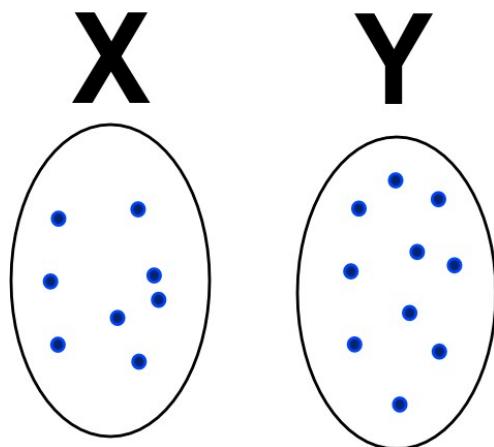
$$D(X, Y) = \frac{1}{|X||Y|} \sum_{x \in X, y \in Y} d(x, y)$$

- WE COULD ALSO CONSIDER THE MEAN DISTANCE BETWEEN THE POINTS OF THE CLUSTERS



# DISTANCE BETWEEN CLUSTERS

- GIVEN TWO CLUSTERS X AND Y.
- HOW DO WE MEASURE THE DISTANCE BETWEEN THEM ?

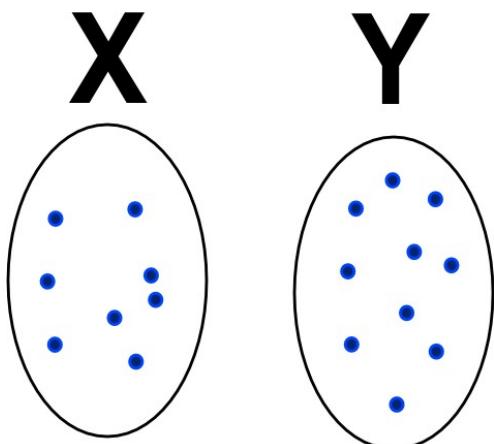


- THERE ARE OTHER MEASURES AS WELL :THE MINIMAL ENERGY CRITERION,  
THE DISTANCE BETWEEN THE CENTROIDS OF THE CLUSTERS



# DISTANCE BETWEEN CLUSTERS

- GIVEN TWO CLUSTERS X AND Y. HOW DO WE MEASURE THE DISTANCE BETWEEN THEM ?



- FOR EFFICIENT CALCULATIONS WE USUALLY REQUIRE THE FOLLOWING CONDITION:
  - KNOWING THE DISTANCE  $d(A, C)$ ,  $d(B, C)$
  - IMPLIES WE CAN CALCULATE IN CONSTANT TIME THE DISTANCE  $d(A + B, C)$
- THERE ARE OTHER MEASURES AS WELL : [THE MINIMAL ENERGY CRITERION](#), THE DISTANCE BETWEEN THE CENTROIDS OF THE CLUSTERS



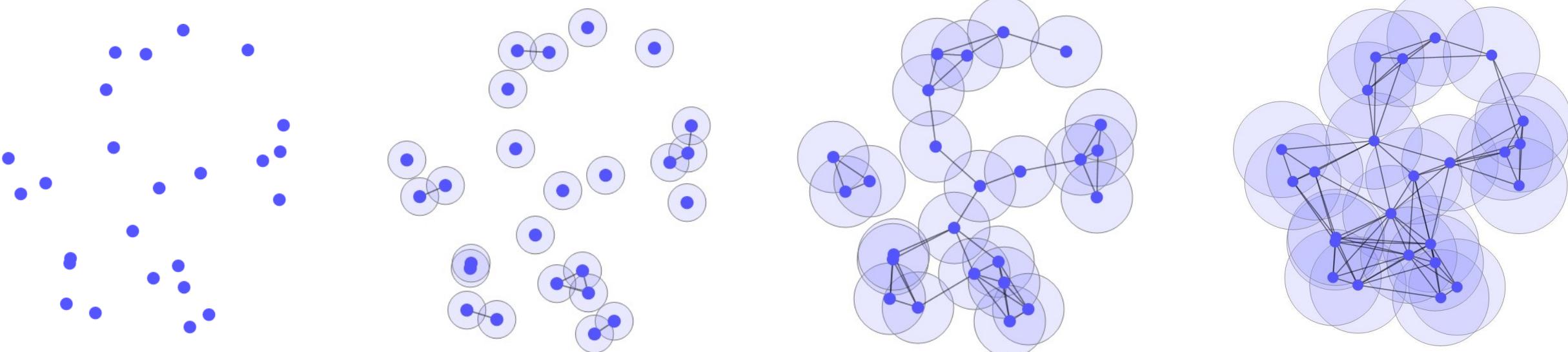
# SINGLE-LINKAGE CLUSTERING ALGORITHM

- THE INPUT OF THE ALGORITHM IS A DISTANCE MATRIX D, CONTAINS ALL DISTANCES  $d(i, j)$  BETWEEN THE POINTS. AT THE BEGINNING EACH POINT IS ITS OWN CLUSTER.



# SINGLE LINKAGE HIERARCHICAL CLUSTERING AND THE $\varepsilon$ -NEIGHBORHOOD GRAPH

- SUPPOSE THAT WE ARE GIVEN A SET OF POINTS  $X = \{p_1, p_2, \dots, p_n\}$  IN  $R^d$  WITH A DISTANCE FUNCTION  $d$  DEFINED ONE THEM.
- CONSIDER THE CONNECTED COMPONENTS OF THE  $\varepsilon$ -NEIGHBORHOOD GRAPH AS WE CONTINUOUSLY INCREASE  $\varepsilon$  FROM ZERO TO INFINITY.



Every point is a connected component

When  $\varepsilon$  is a little larger we start some clusters starts to get form

When  $\varepsilon$  is even larger we have few clusters. As the clusters get larger and larger

At some point all points become a part to of a single cluster

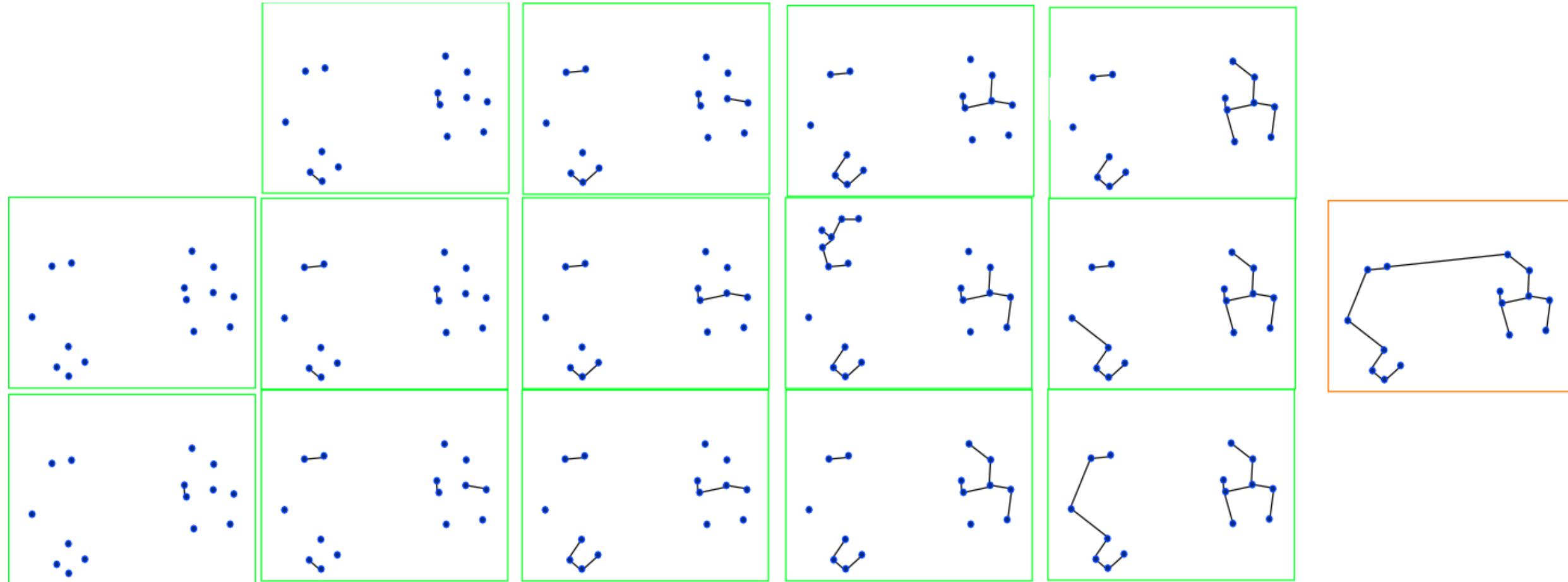


## KRUSKAL'S ALGORITHM

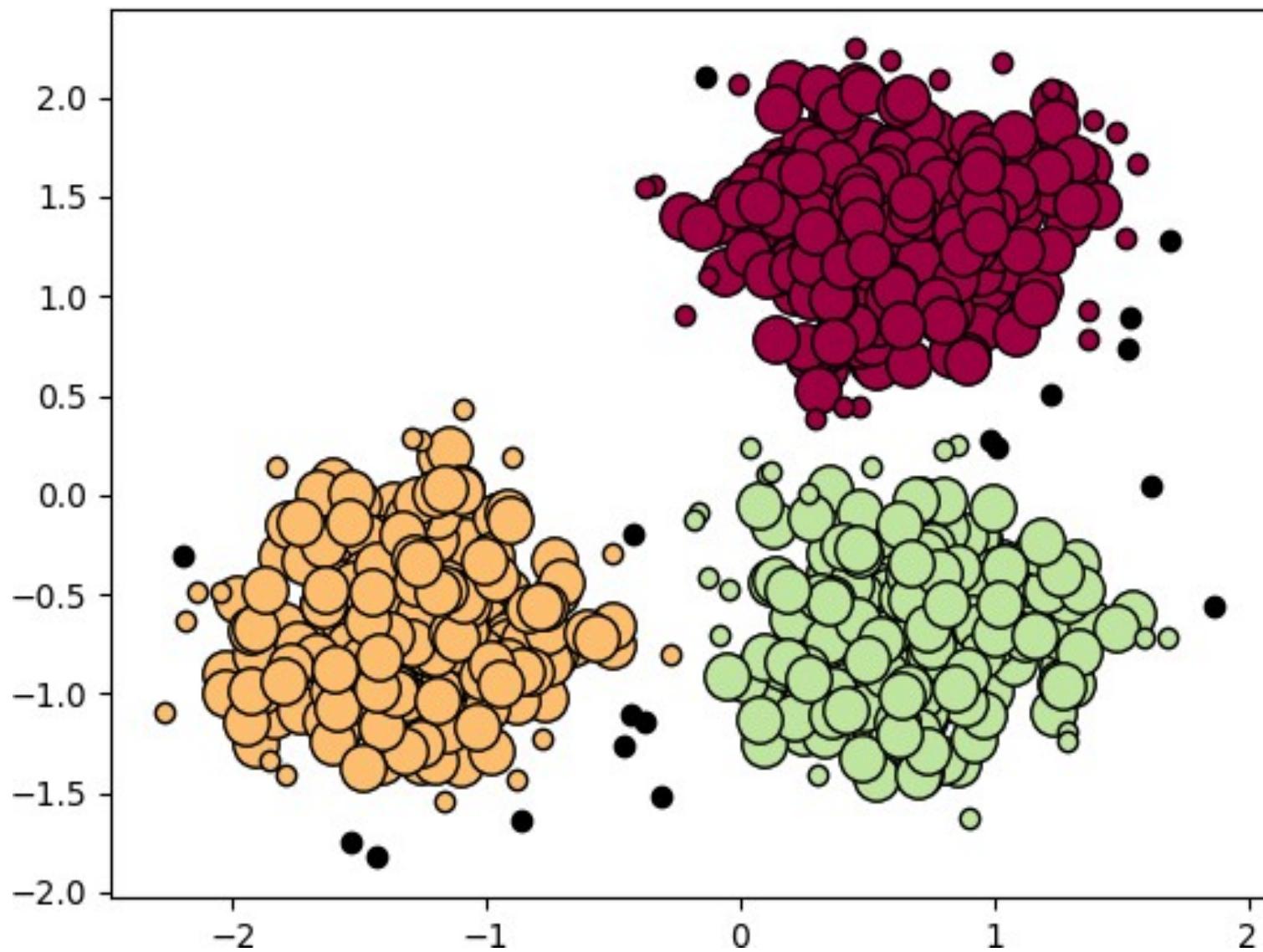
- LET  $G = (V, E, w)$  BE A CONNECTED WEIGHTED GRAPH.
- INFORMALLY, THE ALGORITHM CAN BE GIVEN BY THE FOLLOWING STEPS:
  - I. SET  $VT$  to be  $V$ , Set  $ET = \{\}$ . LET  $S = E$
  2. WHILE  $S$  IS NOT EMPTY AND  $T$  IS NOT A SPANNING TREE
    - I. Select an edge  $e$  from  $S$  with the minimum weight and delete  $e$  from  $S$ .
    2. If  $e$  connects two separate trees of  $T$  then add  $e$  to  $ET$



# SINGLE LINKAGE HIERARCHICAL CLUSTERING AND THE AND KRUSKAL'S ALGORITHM

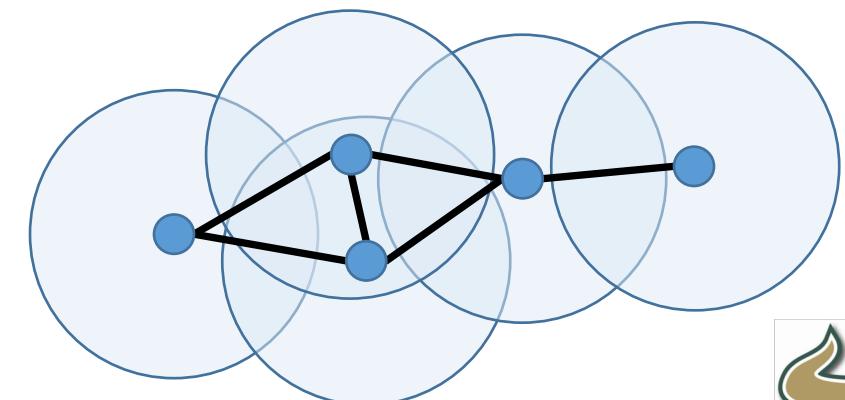


DBSCAN  
(DENSITY-BASED SPATIAL CLUSTERING OF APPLICATIONS WITH NOISE)



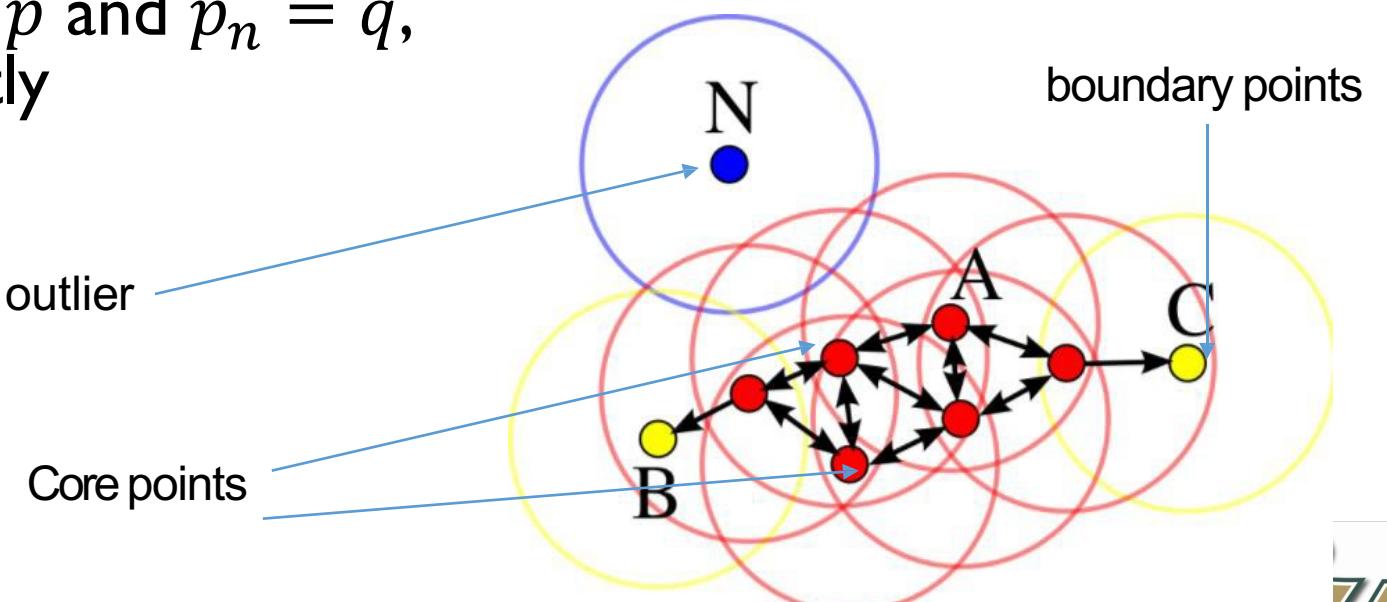
## $\varepsilon$ NEIGHBOR GRAPH

- DBSCAN USE AN  $\varepsilon$  NEIGHBOR GRAPH TO IDENTIFY CLUSTERS.
- FOR THE  $\varepsilon$  NEIGHBOR GRAPH:
  - Points from data are vertices
  - Edges are formed if  $d(x_i, x_j) \leq \varepsilon$  for all  $i \neq j$



# DBSCAN

- IN DBSCAN POINTS ARE CLASSIFIED AS FOLLOWS :
  - **core points**:A point p is said to be a core point if at least **min\_samples** points are within a distance  $\varepsilon$
  - **directly reachable**:A point q is directly reachable from p if the point q is within distance  $\varepsilon$  from point p, and p is a core point.
  - **Density-reachable points**:A point q is reachable from p if there is a path  $p_1, \dots, p_n$  with  $p_1 = p$  and  $p_n = q$ , where each  $p_{i+1}$  is directly reachable from  $p_i$  (all the points on the path must be core points, with the possible exception of q).

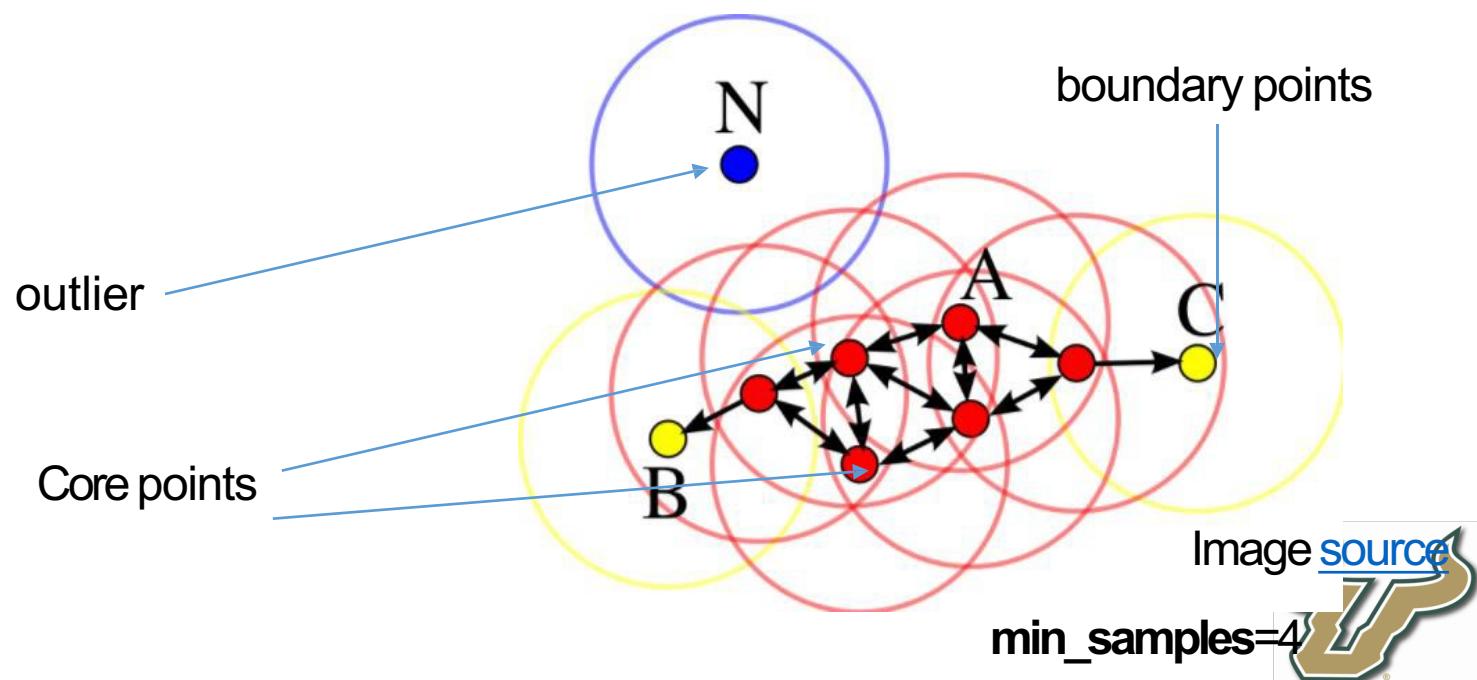


MIN\_SAMPLES=4



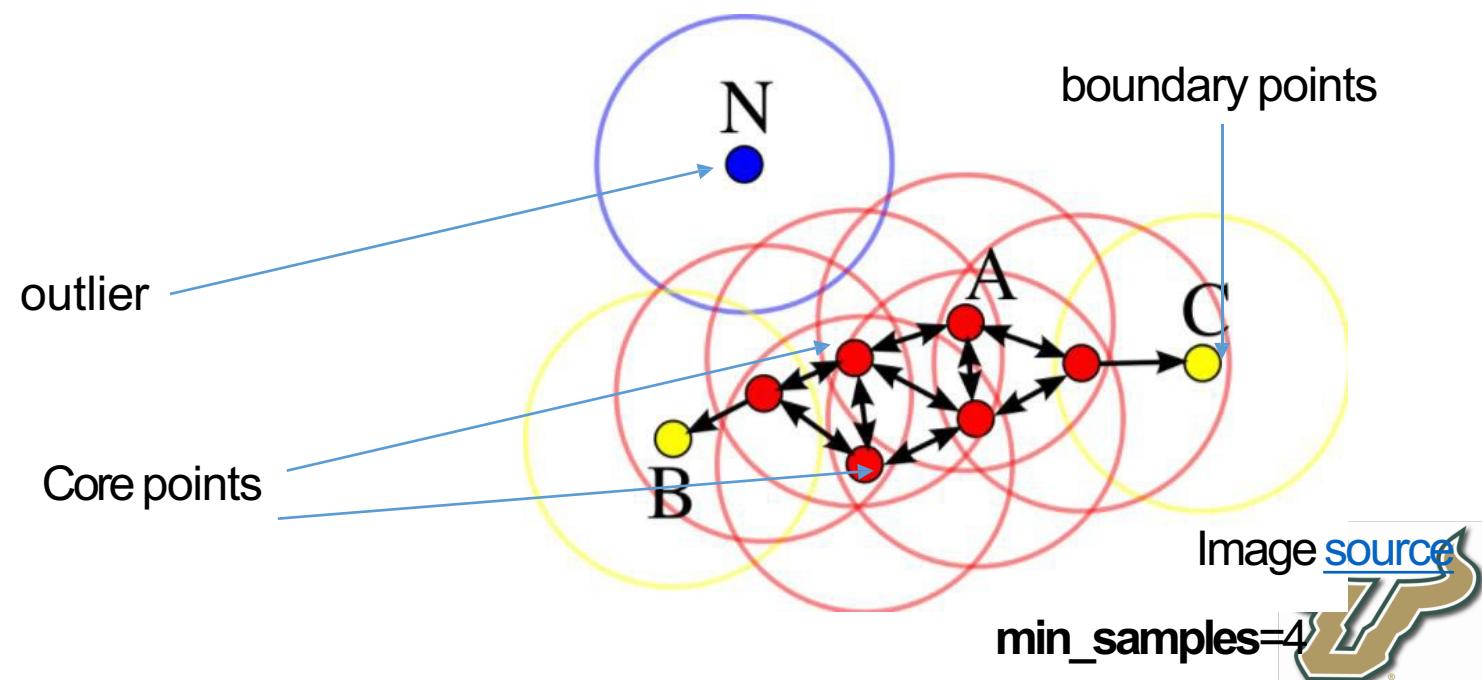
# DBSCAN

- IN DBSCAN POINTS ARE CLASSIFIED AS FOLLOWS :
  - **Outliers**:All points not reachable from any other point are called outliers.
  - **Boundary points**: Boundary points cannot reach other points but other points core points can reach them



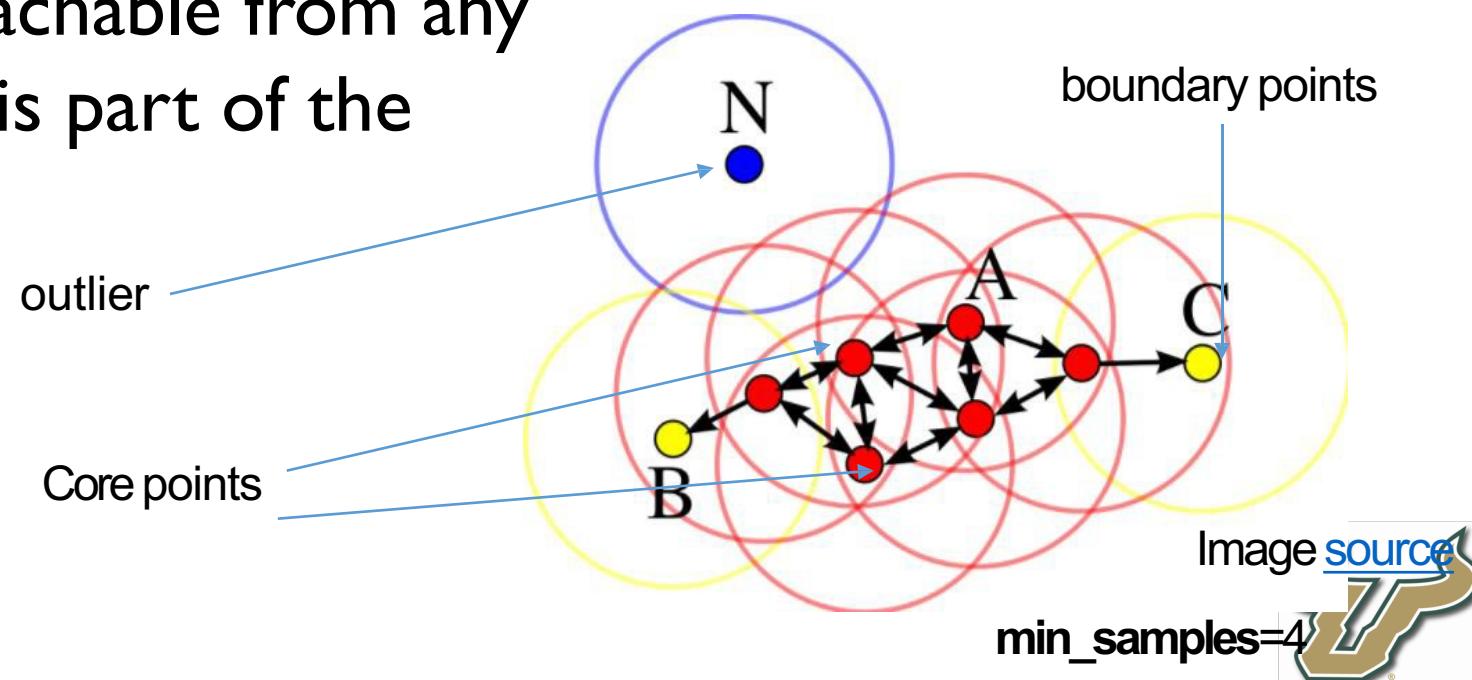
## DBSCAN-REACHABILITY

- IF P IS REACHABLE FROM Q. DOES THAT MEAN Q IS REACHABLE FROM P? EXPLAIN. CAN YOU GIVE AN EXAMPLE FROM THE POINTS BELOW?



# DBSCAN-DENSITY-CONNECTEDNESS

- TWO POINTS P AND Q ARE DENSITY-CONNECTED IF THERE IS A POINT O SUCH THAT BOTH P AND Q ARE REACHABLE FROM O.
- A CLUSTER THEN SATISFIES TWO PROPERTIES:
  - All points within the cluster are mutually density-connected.
  - If a point is density-reachable from any point of the cluster, it is part of the cluster as well.

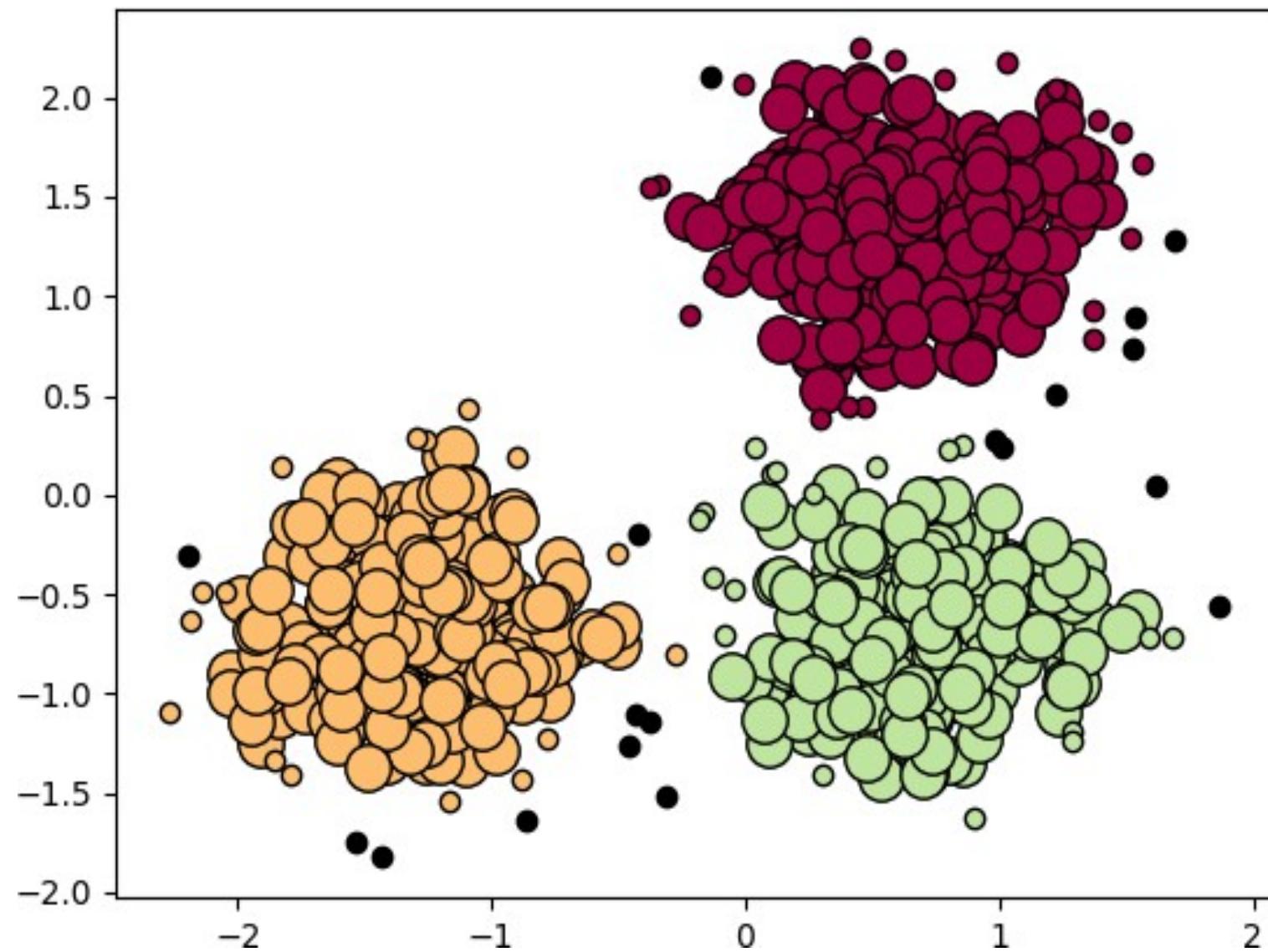


# DBSCAN-ALGORITHM

- INPUT A DATA  $X$ , A POSITIVE REAL NUMBER  $\epsilon$  AND A POSITIVE INTEGER **MIN\_SAMPLES**
  1. Find the  $\epsilon$  neighbor graph.
  2. Identify the core points with more than **MIN\_SAMPLES** neighbors.
  3. Find the connected components of core points on the neighbor graph, ignoring all non-core points.
  4. Assign each non-core point to a nearby cluster if the cluster is an  $\epsilon$  neighbor, otherwise assign it to noise.



# DBSCAN-SKLEARN



# COMPARISON BETWEEN CLUSTERING ALGORITHMS

Method name	Parameters	Scalability	Use case	Geometry (metric used)
K-Means	number of clusters	Very large n_samples , medium n_clusters 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 n_samples	Many clusters, uneven cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Mean-shift	bandwidth	Not scalable with n_samples	Many clusters, uneven cluster size, non-flat geometry	Distances between points
Spectral clustering	number of clusters	Medium n_samples , small n_clusters	Few clusters, even cluster size, non-flat geometry	Graph distance (e.g. nearest-neighbor graph)
Ward hierarchical clustering	number of clusters	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints	Distances between points
Agglomerative clustering	number of clusters, linkage type, distance	Large n_samples and n_clusters	Many clusters, possibly connectivity constraints, non Euclidean distances	Any pairwise distance
DBSCAN	neighborhood size	Very large n_samples , medium n_clusters	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 n_clusters and n_samples	Large dataset, outlier removal, data reduction.	Euclidean distance between points



[SOURCE](#)

