

K-means

Li, Jia

DSAA 5002

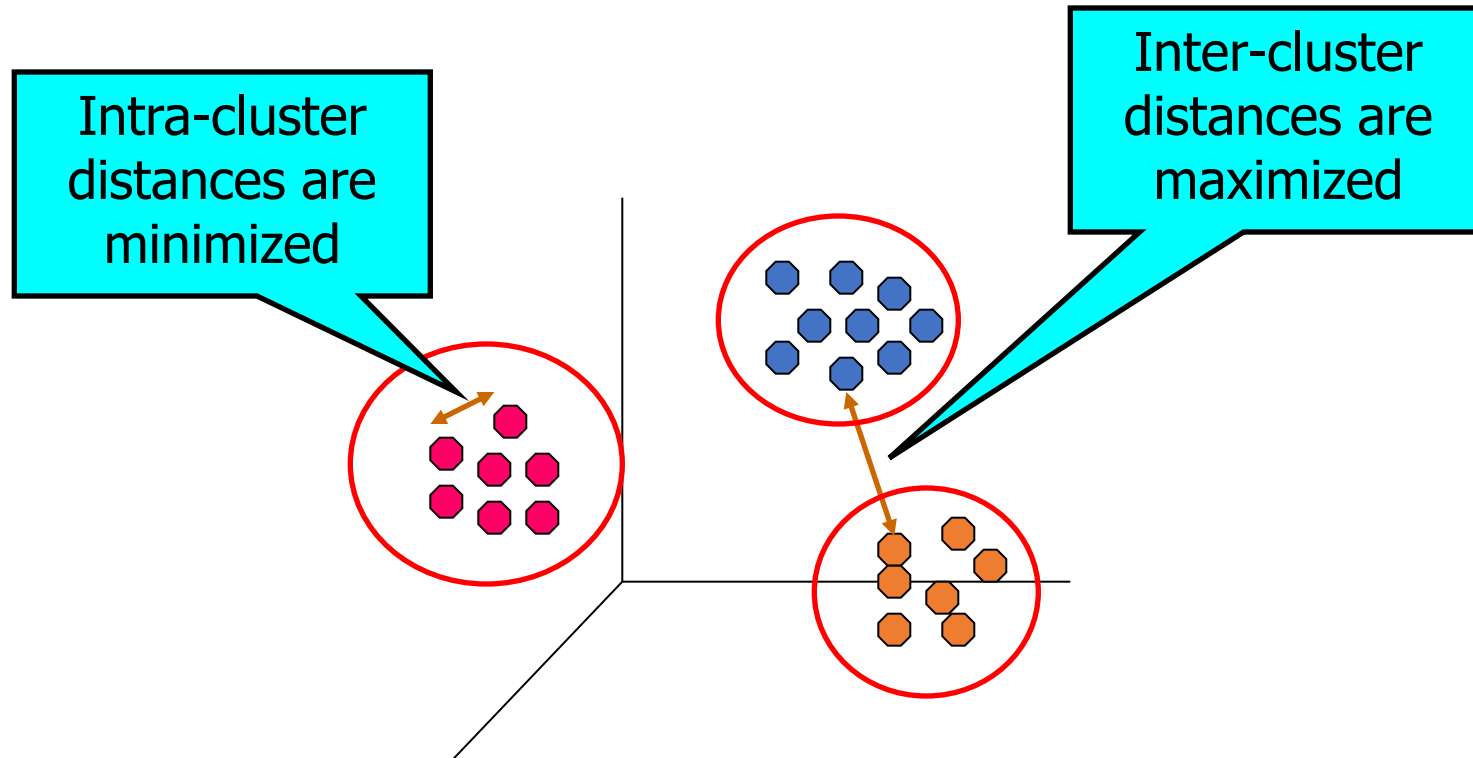
The Hong Kong of Science and Technology (Guangzhou)

2025 Fall

Sep 29

What is Cluster Analysis?

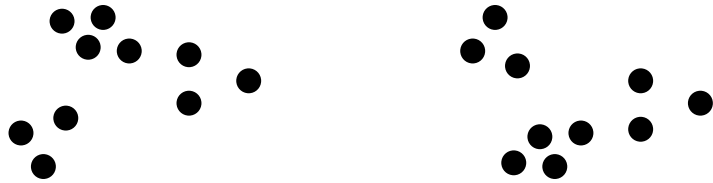
- Finding groups of objects such that the objects in a group will be similar (or related) to one another and different from (or unrelated to) the objects in other groups



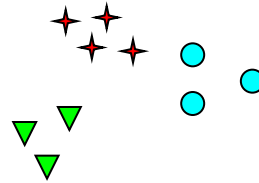
Different Learning Paradigm

- **Supervised learning:** learning with a teacher
 - You had training data which was (feature, label) pairs and the goal was to learn a mapping from features to labels
- **Unsupervised learning:** learning without a teacher
 - Only features and no labels
- Why is unsupervised learning useful?
 - Discover hidden structures in the data — clustering
 - Visualization — dimensionality reduction
 - lower dimensional features might help learning

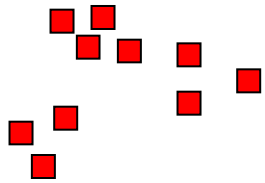
Notion of a Cluster can be Ambiguous



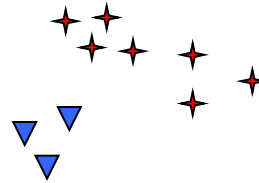
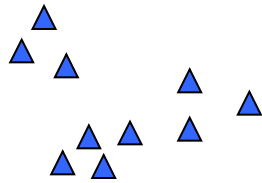
How many clusters?



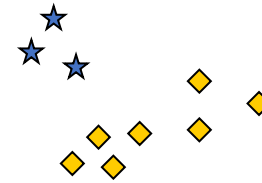
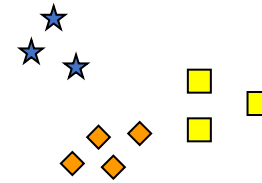
Six Clusters



Two Clusters



Four Clusters



K-means Clustering

- Partitional 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
- The basic algorithm is very simple

-
- 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
-

Lloyd's Algorithm

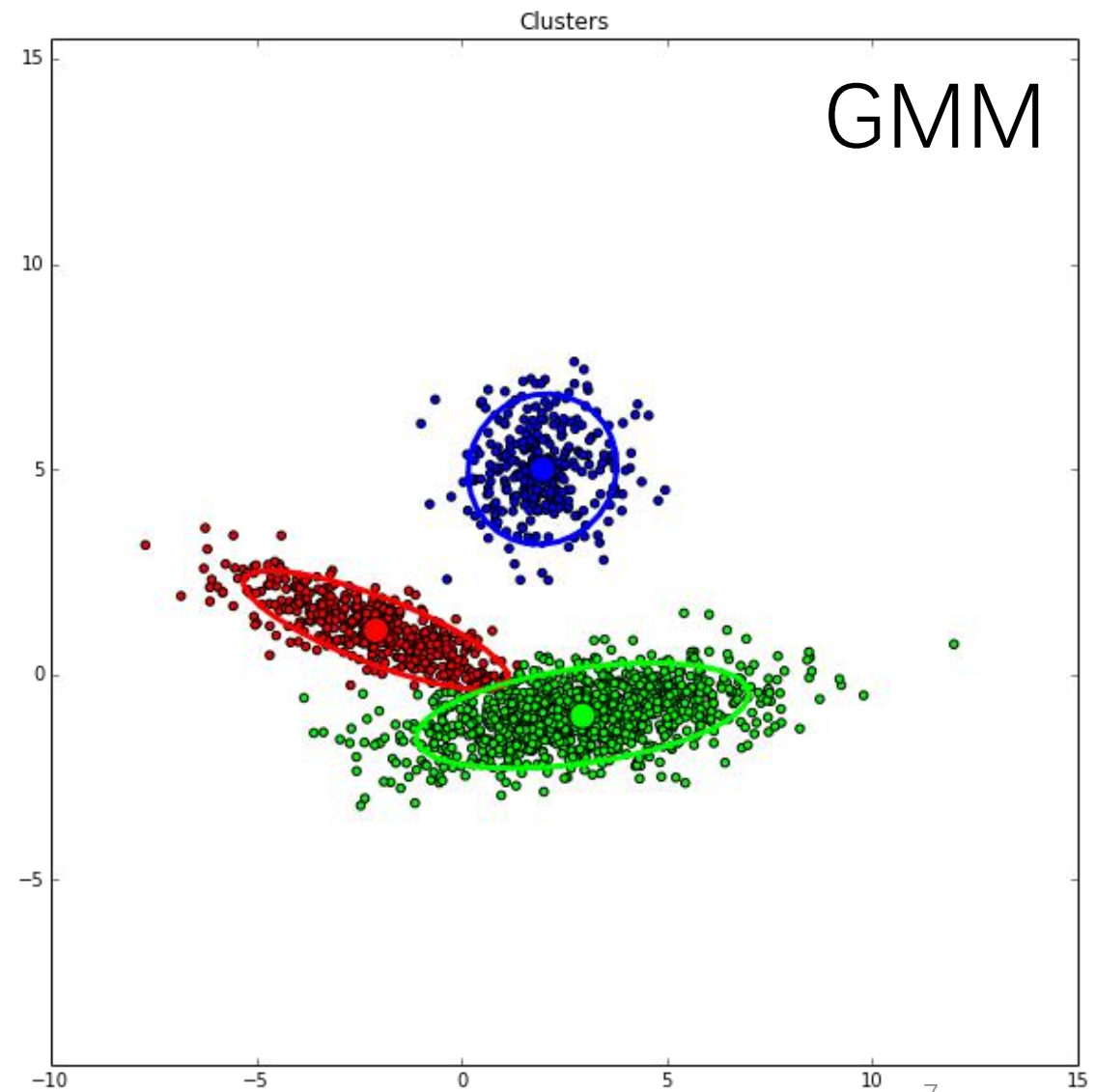
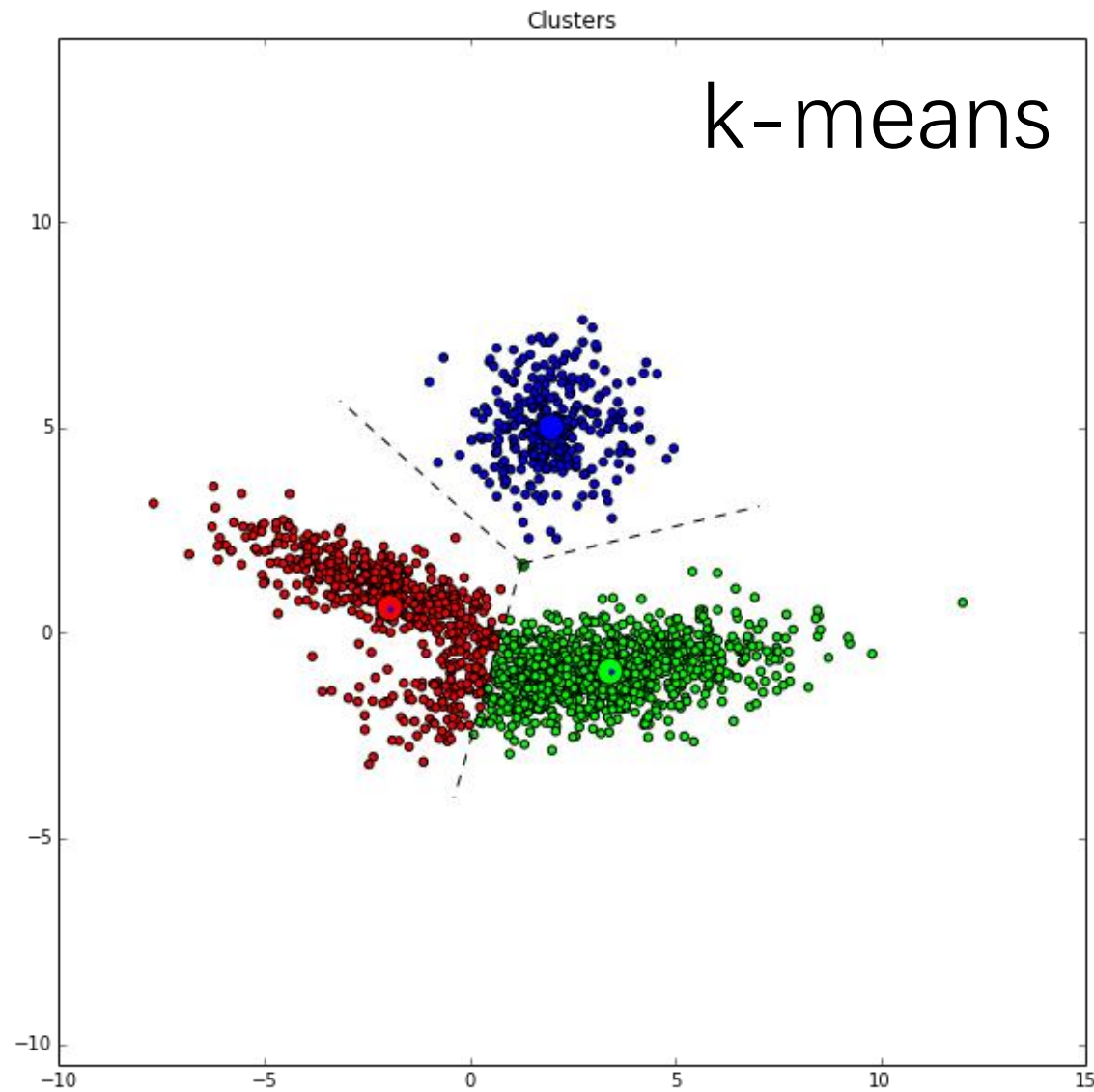
- Initialize k **centers** by picking k points **randomly** among all the points
- Repeat till convergence (or **max iterations**)
 - Assign each point to the nearest **center** (**assignment step**)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} ||\mathbf{x} - \mu_i||^2$$

- Estimate the **mean** of each group (**update step**)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} ||\mathbf{x} - \mu_i||^2$$

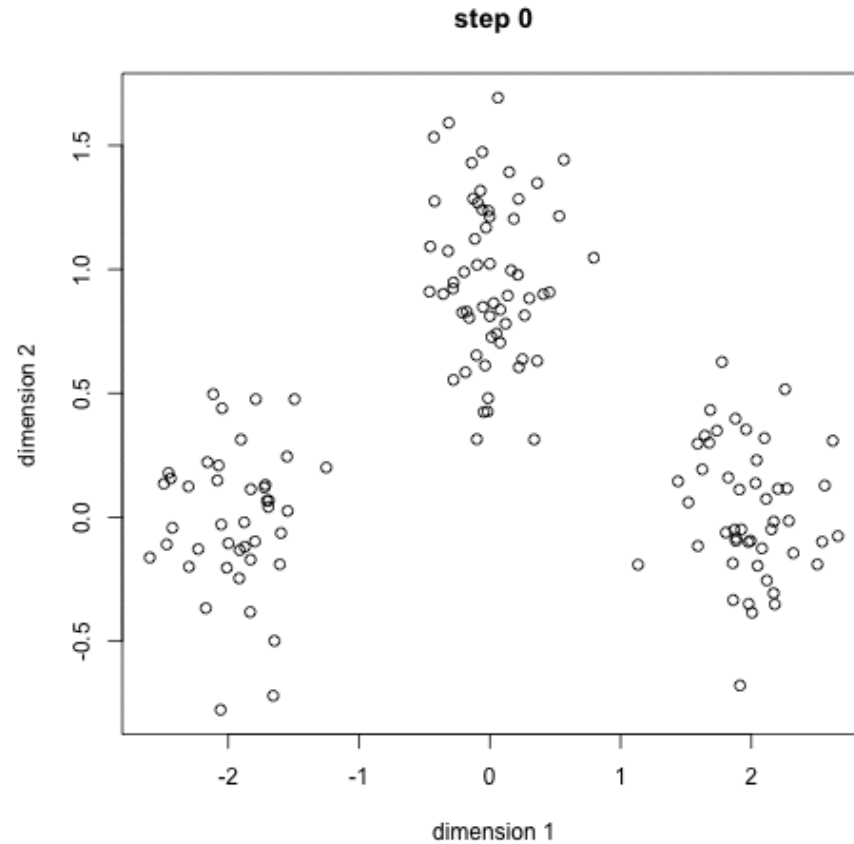
K-means vs GMM



Details

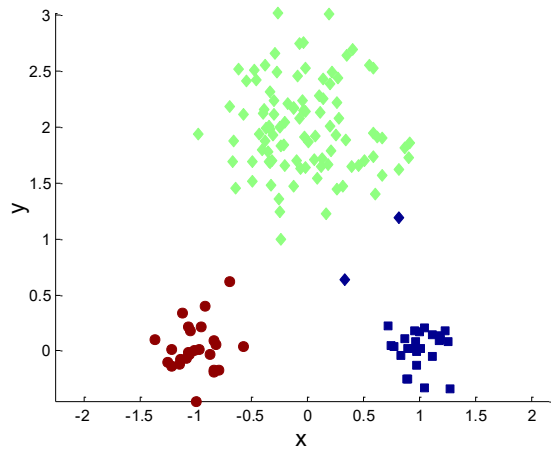
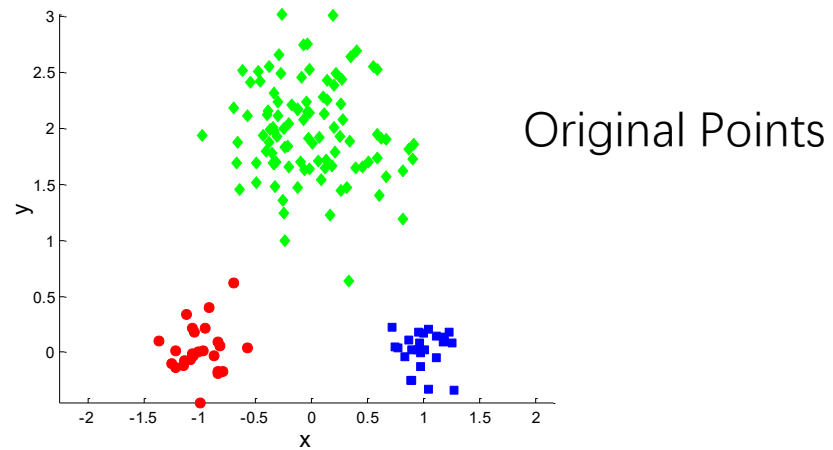
- Initial centroids are often chosen randomly.
 - Clusters produced vary from one run to another.
- The centroid is (typically) the mean of the points in the cluster.
- 'Closeness' is measured by Euclidean distance, cosine similarity, correlation, etc.
- K-means will converge for common similarity measures mentioned above.
- Most of the convergence happens in the first few iterations.
 - Often the stopping condition is changed to 'Until relatively few points change clusters'
- Complexity is $O(n * K * I * d)$
 - n = number of points, K = number of clusters,
 I = number of iterations, d = number of attributes

K-means in Action

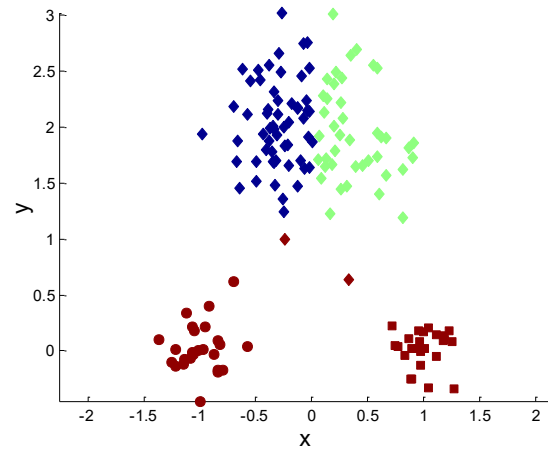


<http://simplystatistics.org/2014/02/18/k-means-clustering-in-a-gif/>

Two different K-means Clusterings

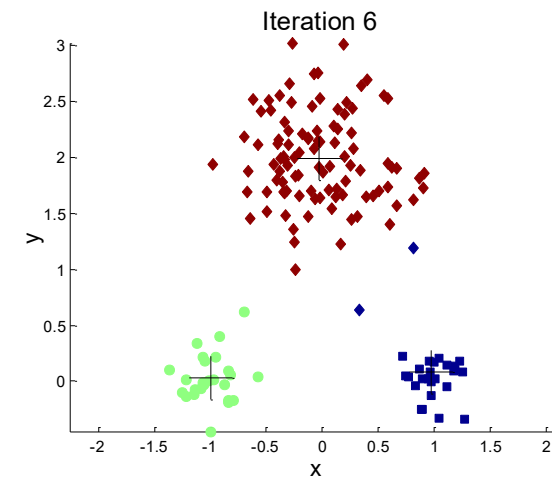
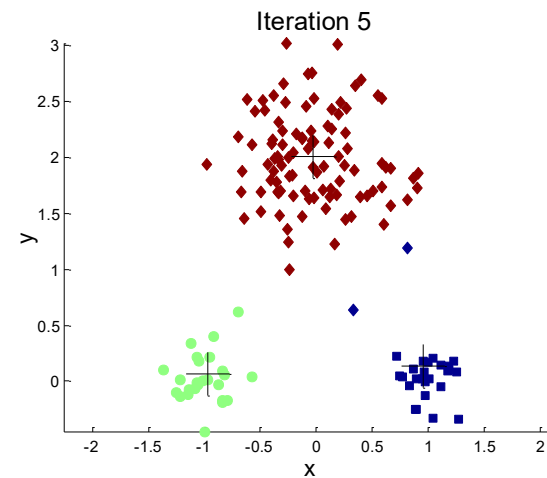
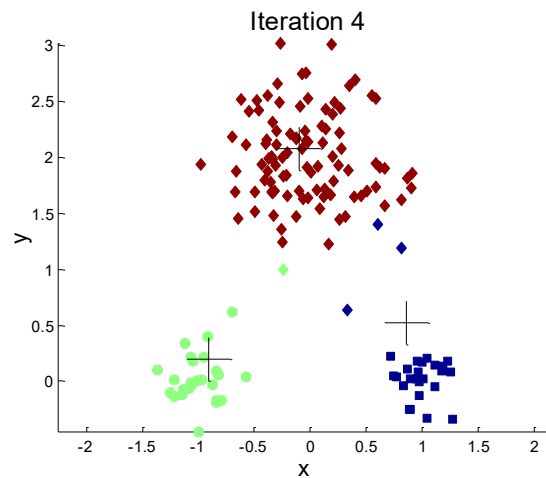
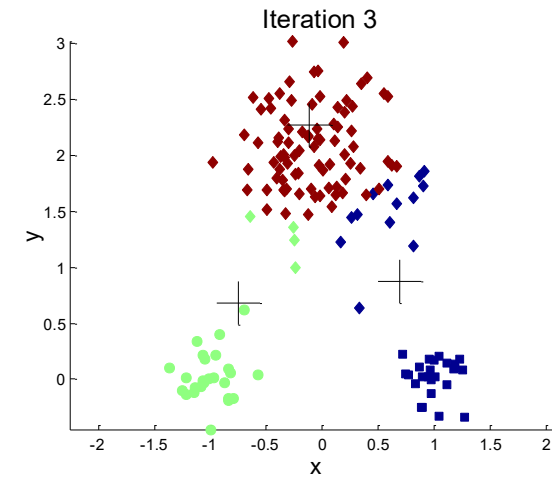
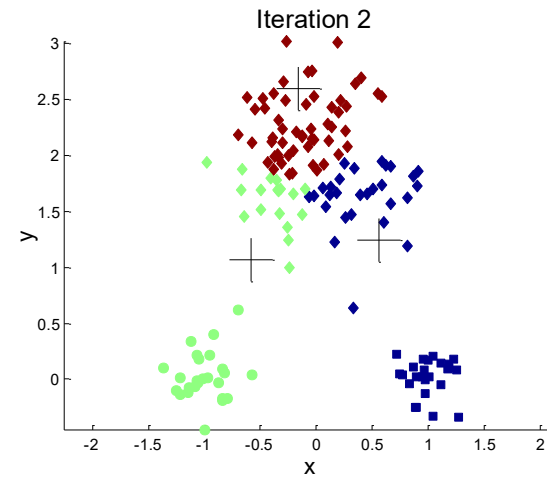
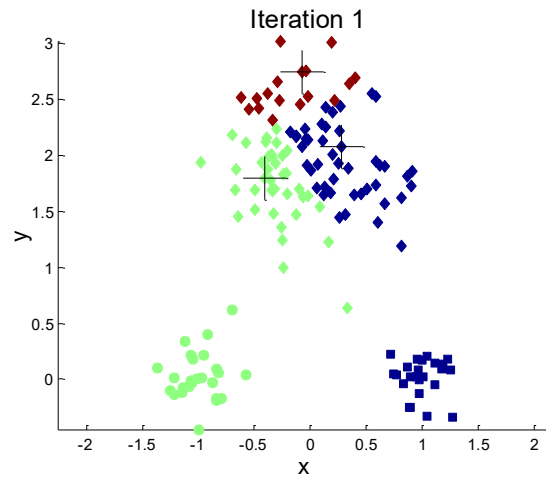


Optimal Clustering



Sub-optimal Clustering

Importance of Choosing Initial Centroids



Evaluating K-means Clusters

- Most common measure is Sum of Squared Error (SSE)
 - For each point, the error is the distance to the nearest cluster
 - To get SSE, we square these errors and sum them.

$$SSE = \sum_{i=1}^K \sum_{x \in C_i} dist^2(m_i, x)$$

- x is a data point in cluster C_i and m_i is the representative point for cluster C_i
 - can show that m_i corresponds to the center (mean) of the cluster
- Given two clusters, we can choose the one with the smallest error
- One easy way to reduce SSE is to increase K , the number of clusters
 - A good clustering with smaller K can have a lower SSE than a poor clustering with higher K

Solutions to Initial Centroids Problem

- Multiple runs
 - Helps, but probability is not on your side
- Sample and use hierarchical clustering to determine initial centroids
- Select more than k initial centroids and then select among these initial centroids
 - Select most widely separated
- Postprocessing
- Bisecting K-means
 - Not as susceptible to initialization issues

Bisecting K-means

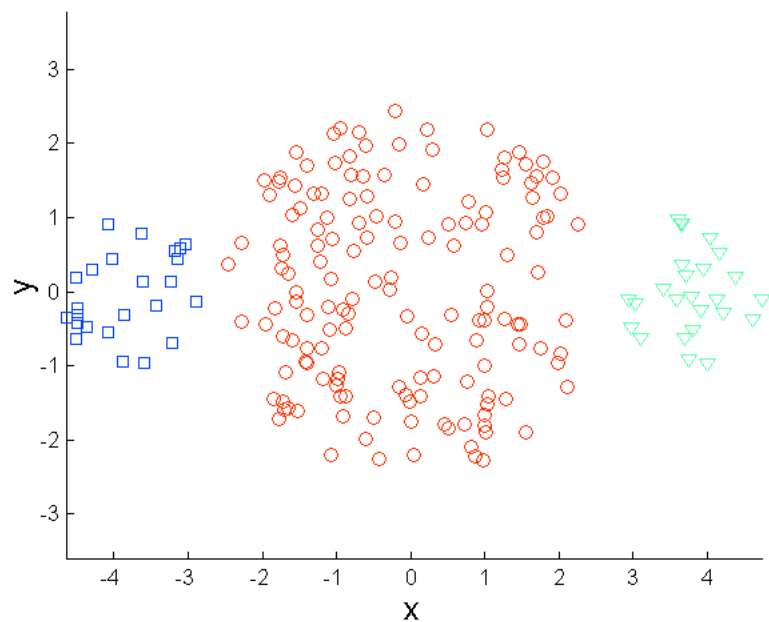
- Bisecting K-means algorithm
 - Variant of K-means that can produce a partitional or a hierarchical clustering

```
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
```

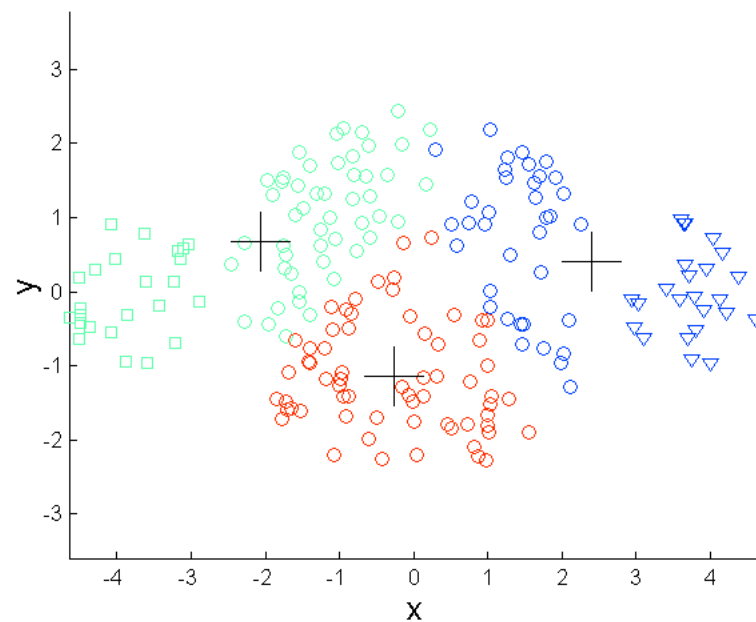
Limitations of K-means

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

Limitations of K-means: Differing Sizes

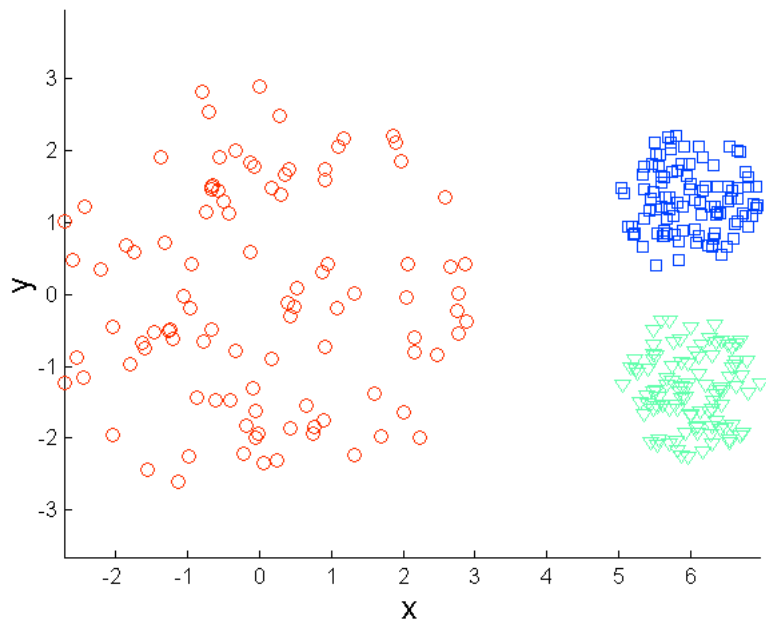


Original Points

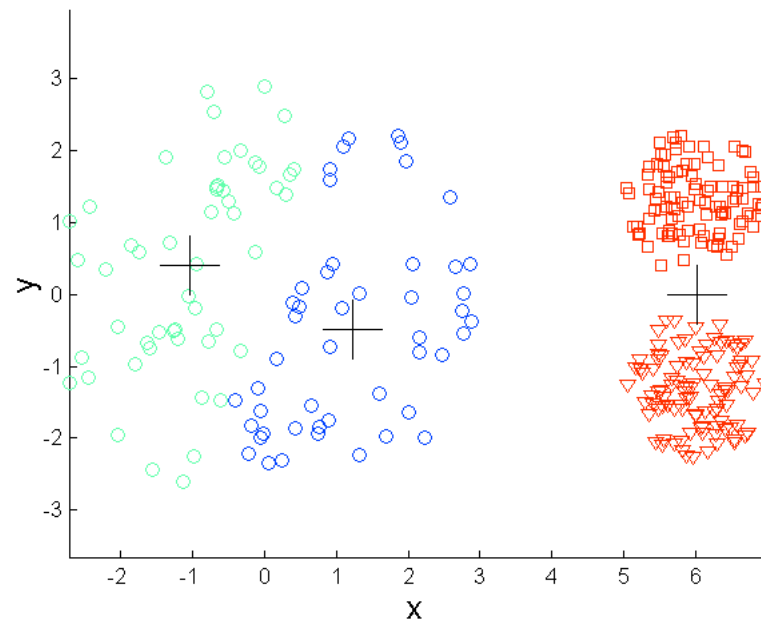


K-means (3 Clusters)

Limitations of K-means: Differing Density

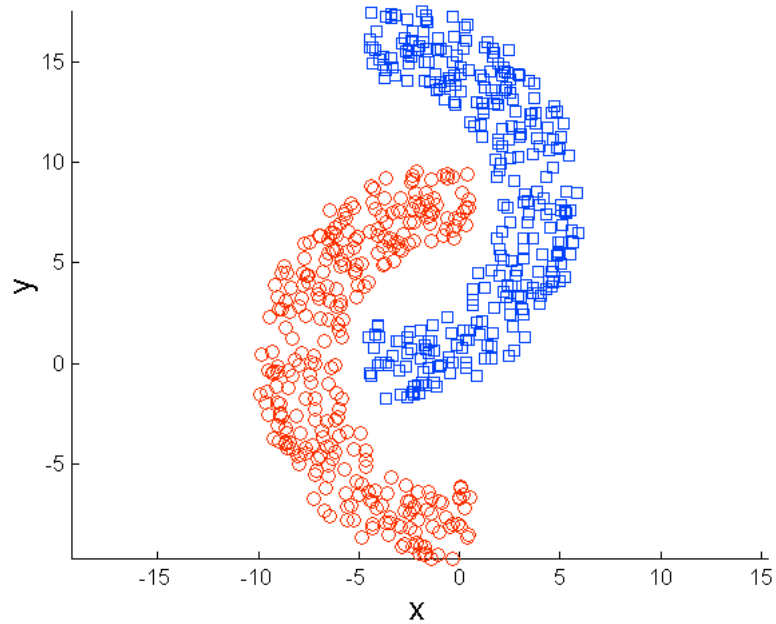


Original Points

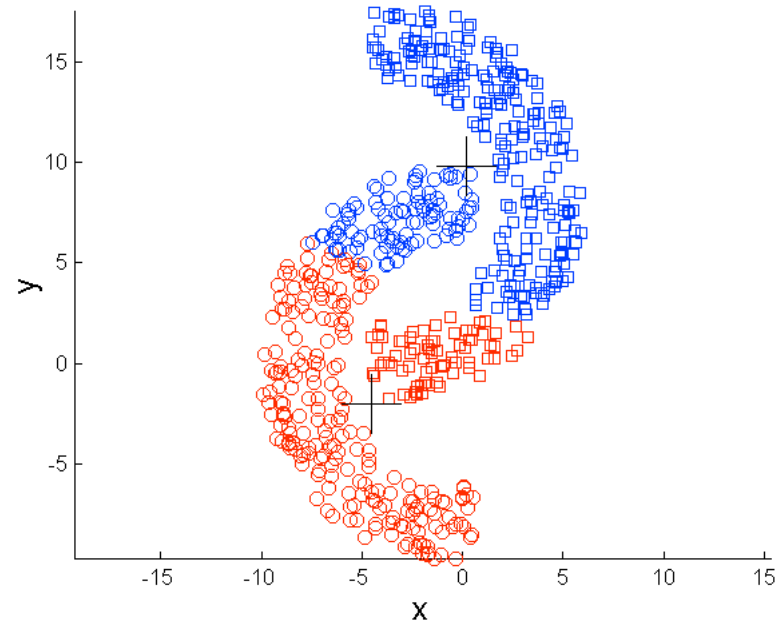


K-means (3 Clusters)

Limitations of K-means: Non-global Shapes

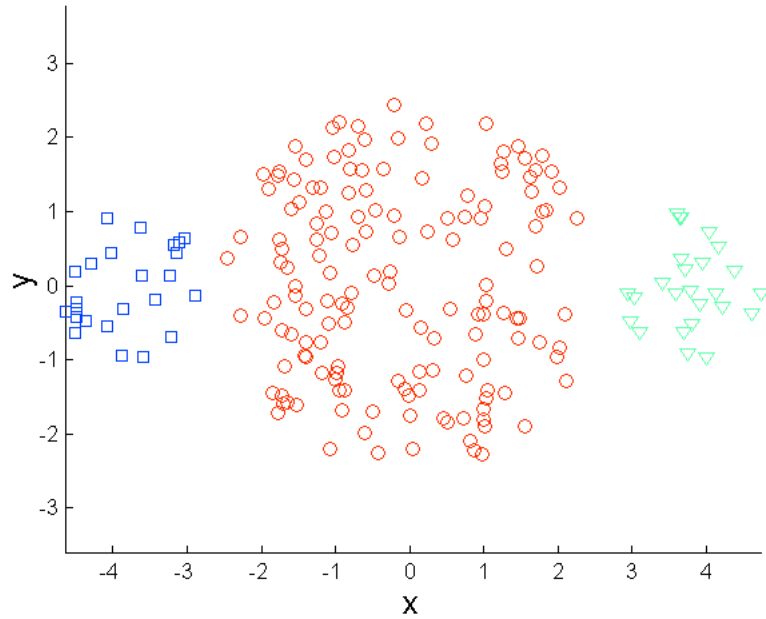


Original Points

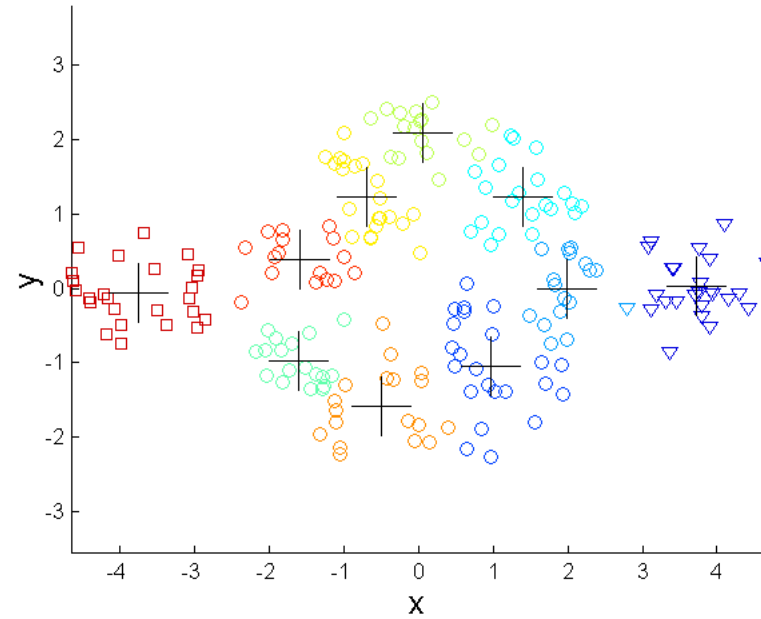


K-means (2 Clusters)

Overcoming K-means Limitations



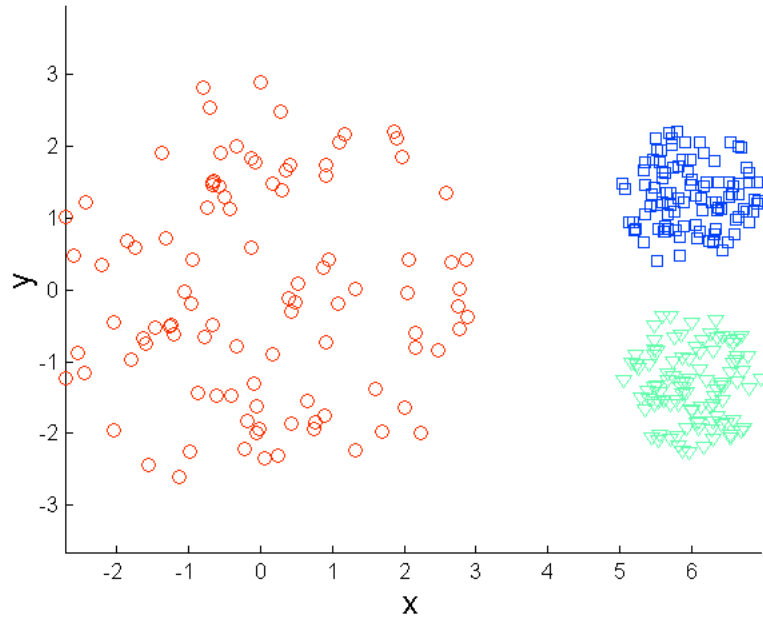
Original Points



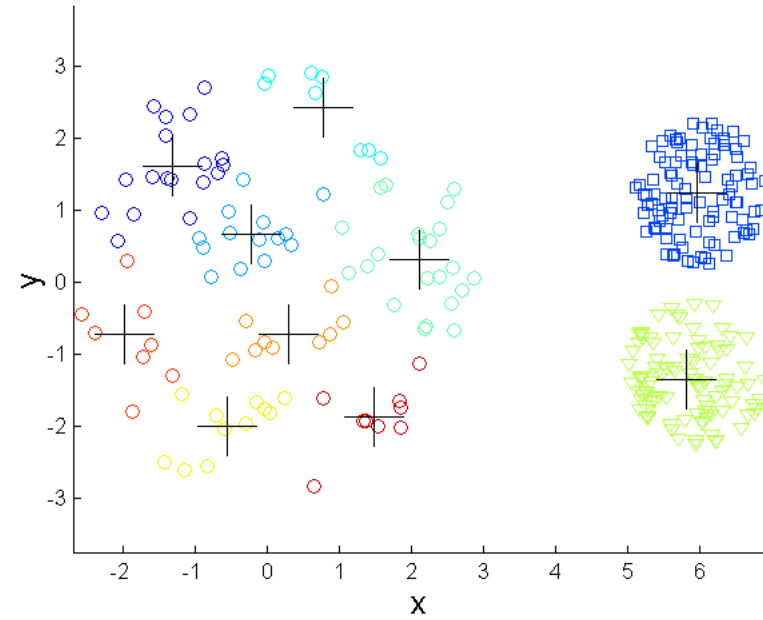
K-means Clusters

One solution is to use many clusters.
Find parts of clusters, but need to put together.

Overcoming K-means Limitations

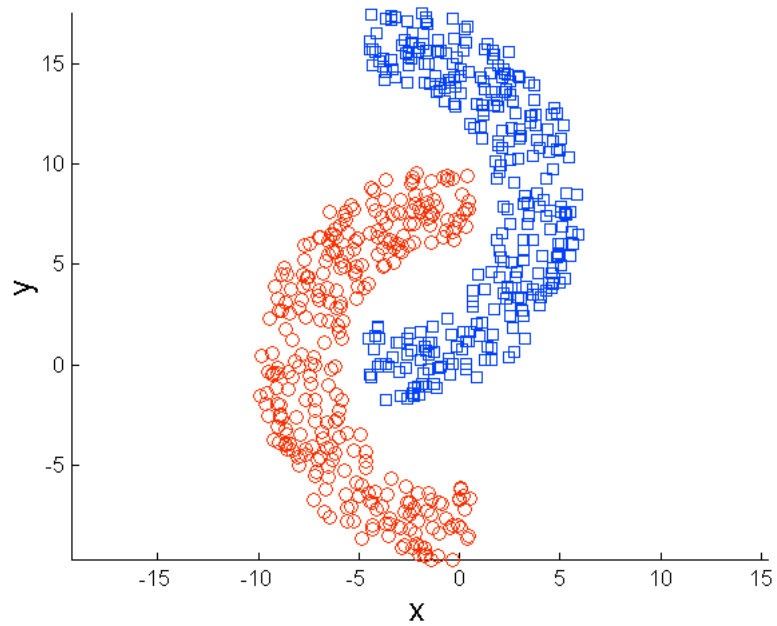


Original Points

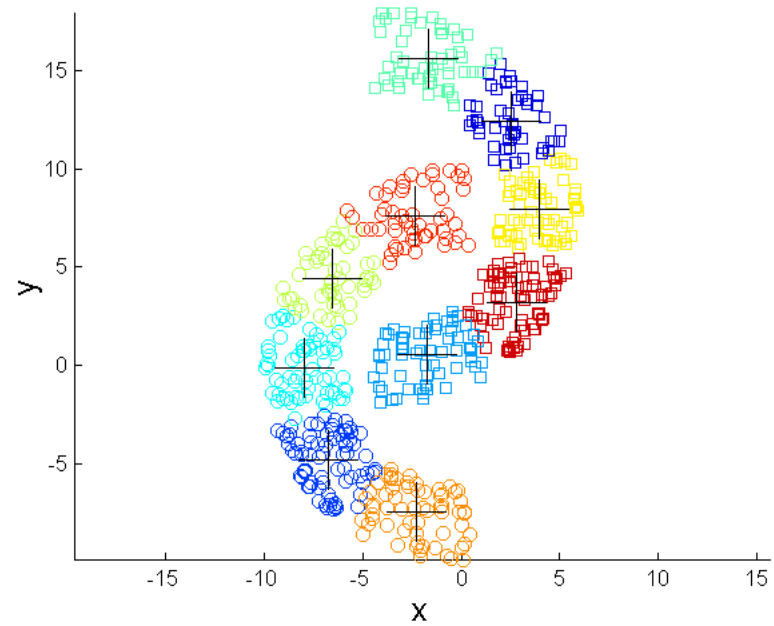


K-means Clusters

Overcoming K-means Limitations



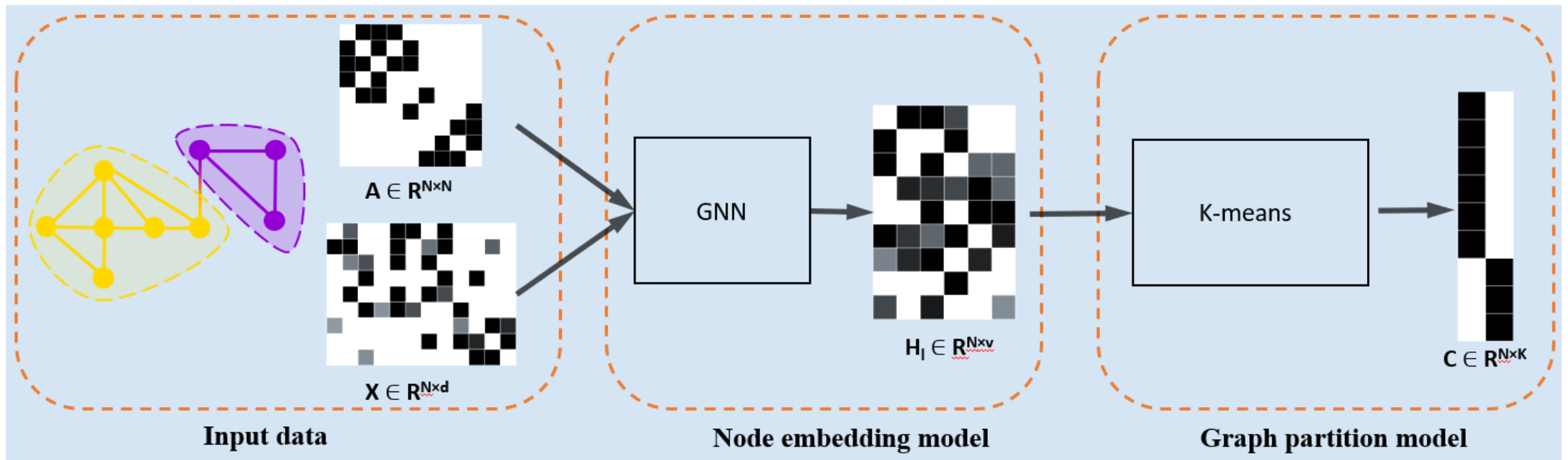
Original Points



K-means Clusters

K-means for Graph Clustering

- Node embedding model
 - VGAE (Kipf et al. 2017)
 - Node2vec (Aditya et al. 2016)
- Graph partition model could use other clustering algorithms.



Slides Credit

- [1] Tan et al. K-means in Introduction to Data Mining.
- [2] Subhransu Maji. Clustering in CMPSCI689