



TEXAS A&M UNIVERSITY  
Engineering

# **ECEN 758 Data Mining and Analysis: Lecture 11, Density-based Clustering**

---

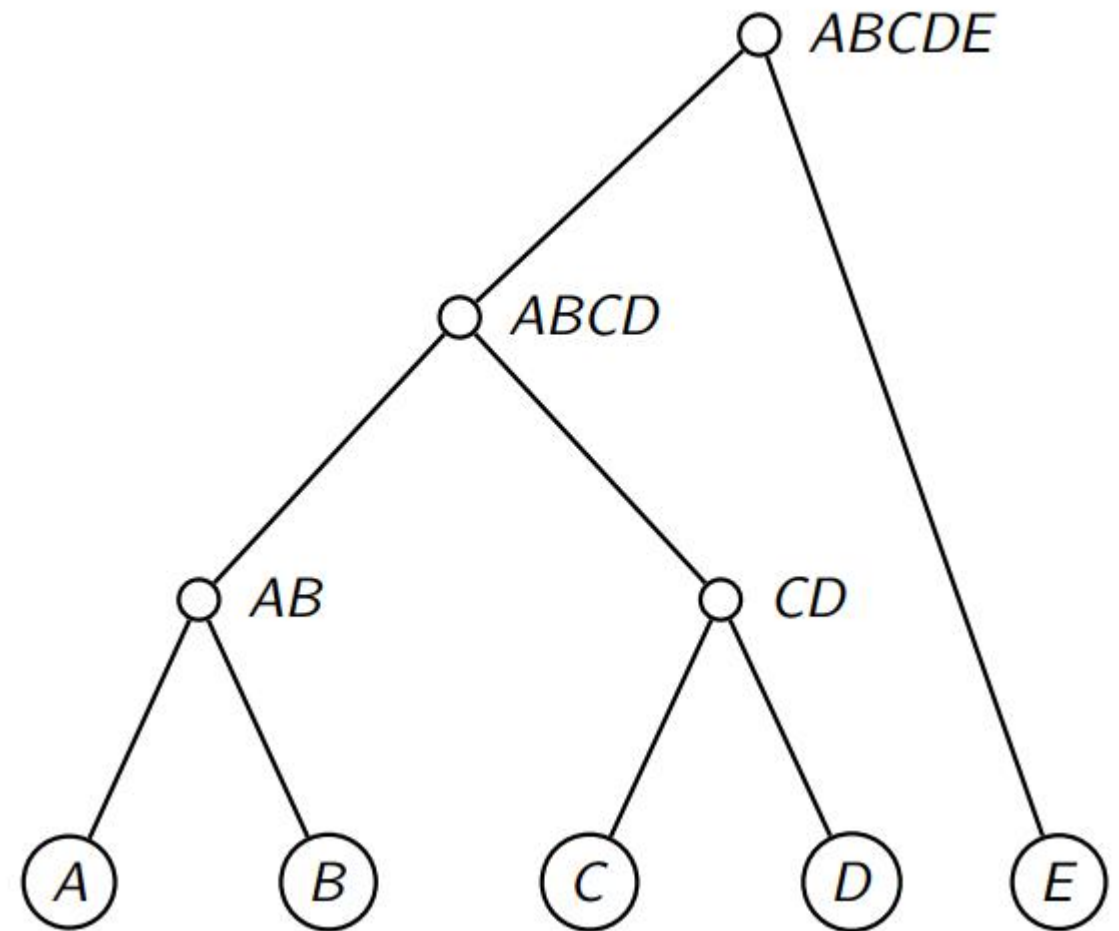
Joshua Peeples, Ph.D.

Assistant Professor

Department of Electrical and Computer Engineering

- Assignment #2 due this Friday (09/27)
  - Q1.4: Only need to show FP-tree and conditional pattern base for final solution

- Hierarchical Clustering



# Today



TEXAS A&M UNIVERSITY  
Engineering

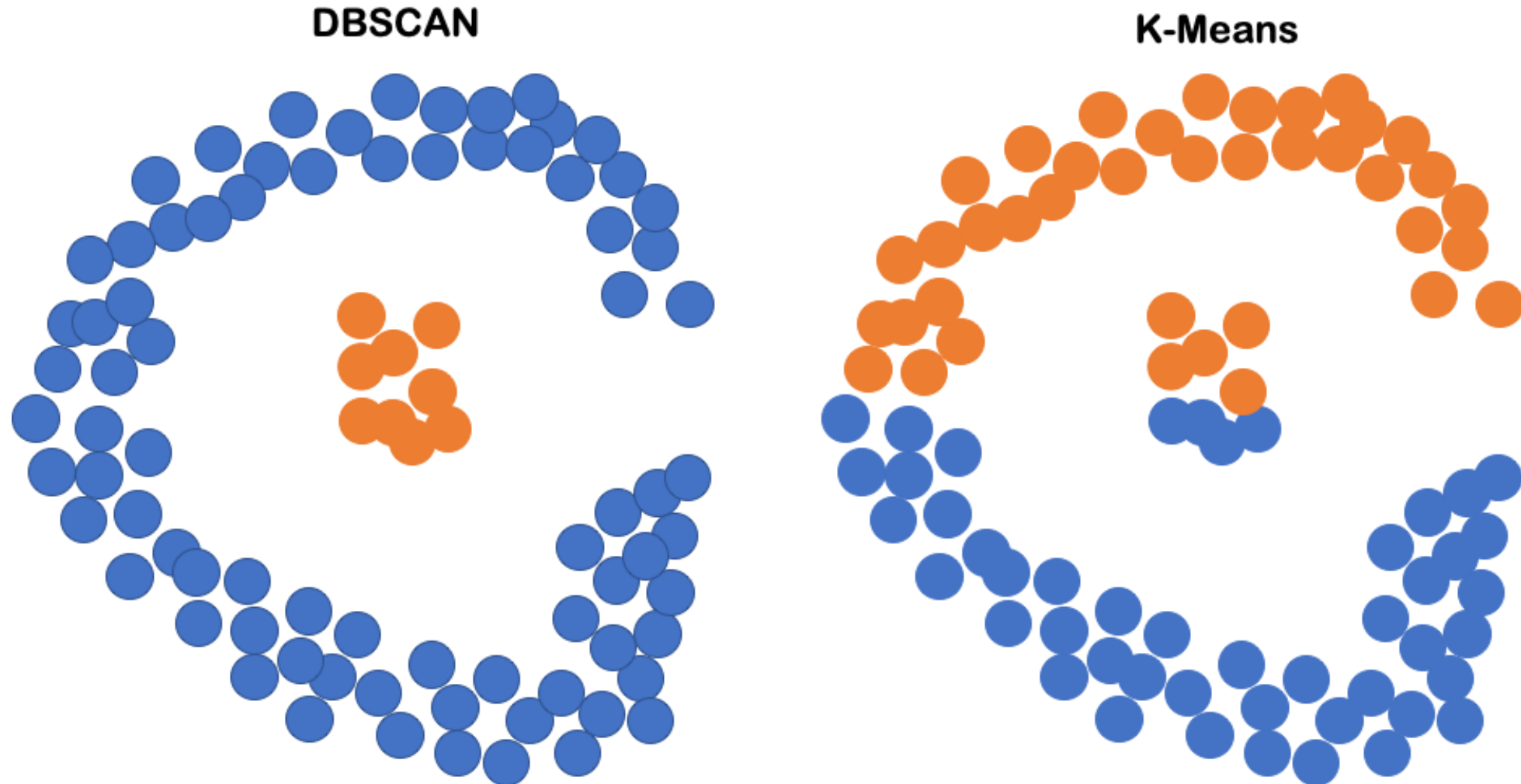
- Density-based Clustering
- Reading: ZM Chapter 15

- We will discuss several variants of clustering
  - Representative-based Clustering
  - Hierarchical Clustering
  - **Density-based Clustering**



# Density-based Clustering Overview

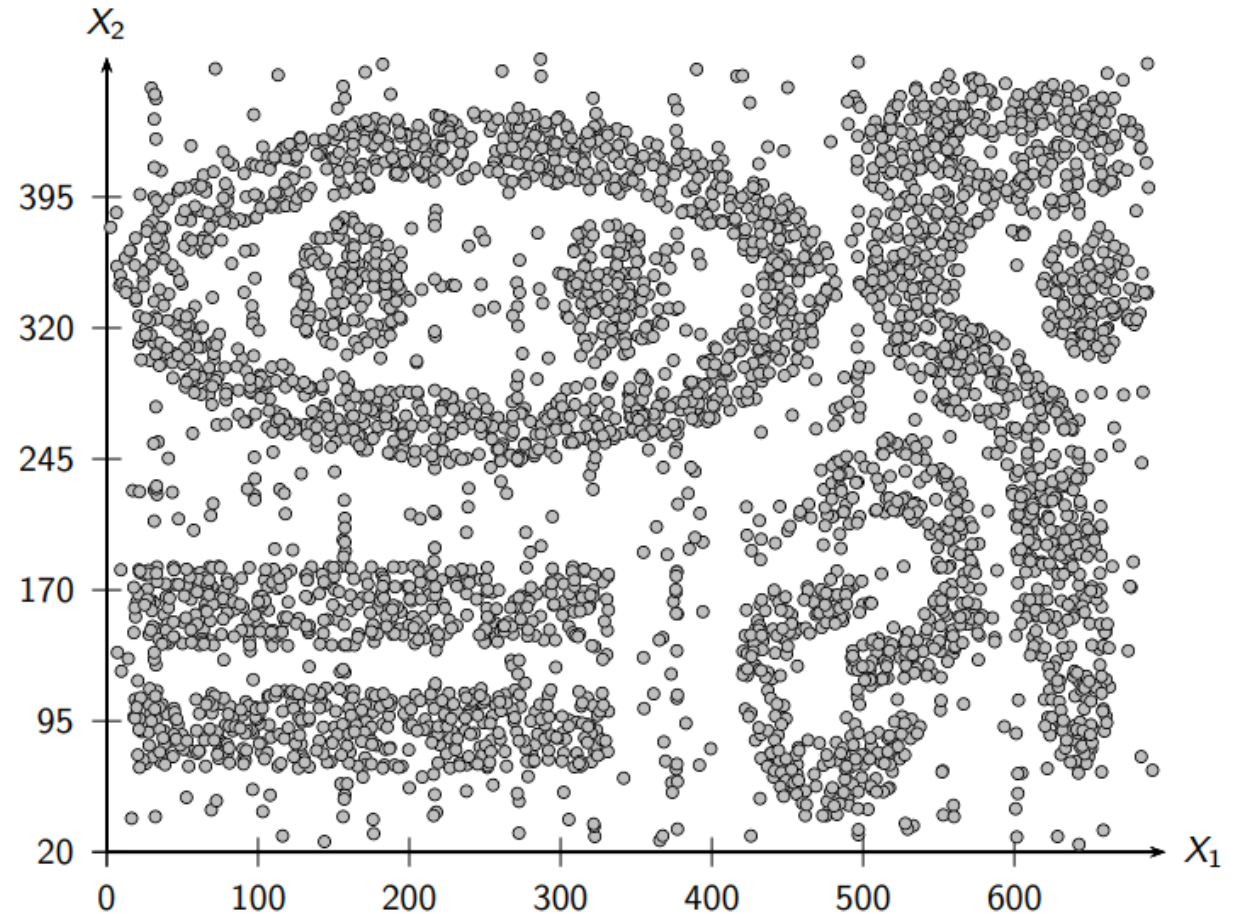
# Shortcomings of k-Means



# Density-based Clustering



- Able to find nonconvex cluster
- Distance-based method may have trouble with nonconvex clusters







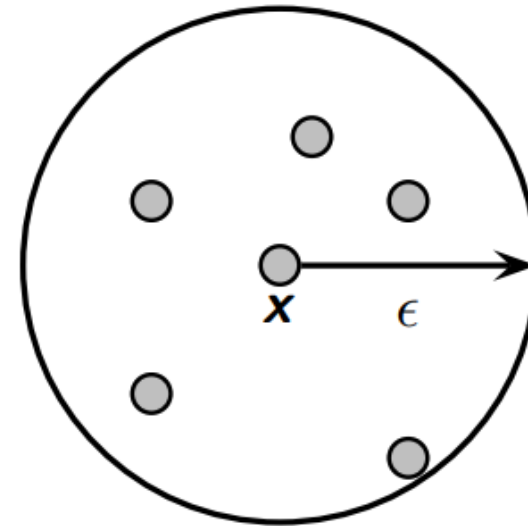
# Density-based spatial clustering of applications with noise (DBSCAN) Algorithm

# DBSCAN Approach



- Define a ball with radius  $\epsilon$  around a point  $\mathbf{x}$ 
  - $\epsilon$  neighborhood of  $\mathbf{x}$
- $\delta$  represents distance function

$$N_{\epsilon}(\mathbf{x}) = B_d(\mathbf{x}, \epsilon) = \{\mathbf{y} \mid \delta(\mathbf{x}, \mathbf{y}) \leq \epsilon\}$$

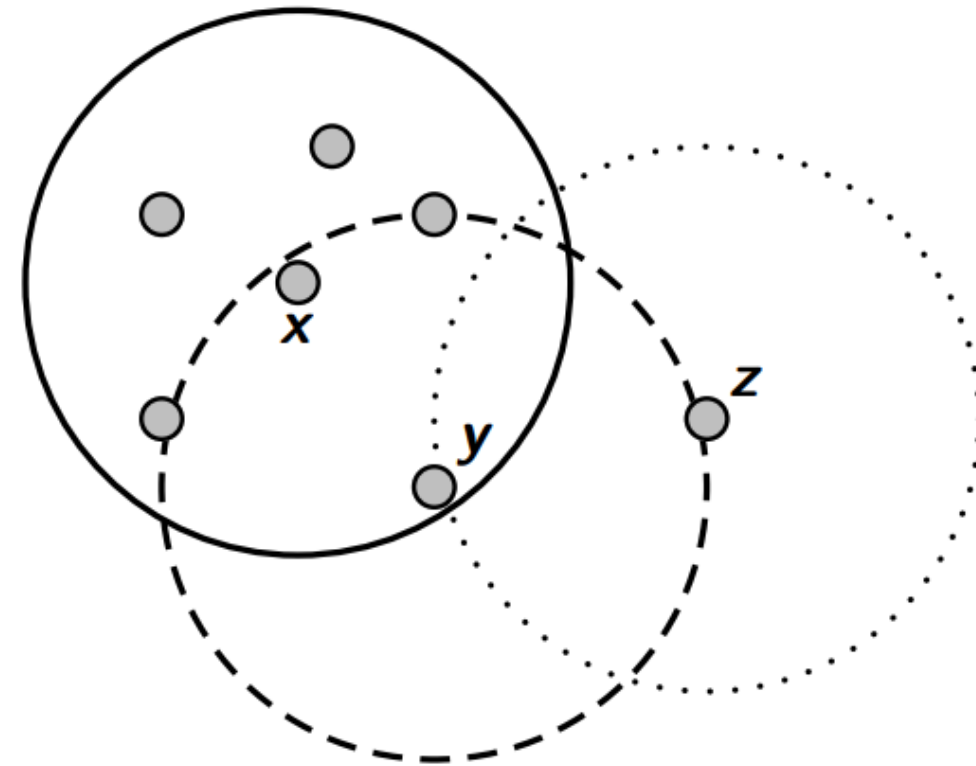


(a) Neighborhood of a Point

# DBSCAN Approach



- Core point:
  - There are at least *minpts* in its  $\epsilon$ -neighborhood
- Border point:
  - Does not meet the *minpts* threshold but is in  $\epsilon$ -neighborhood of core point
- Noise point
  - Neither core or border point



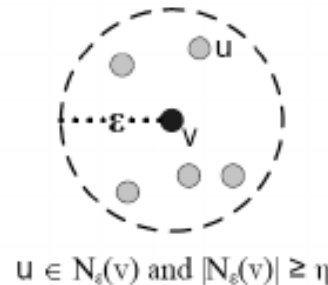
(b) Core, Border, and Noise Points

# DBSCAN Approach



- Directly density reachable:
  - If in a point's  $\epsilon$ -neighborhood and is a core point
- Density reachable:
  - If there are a set of core of points leading from **v** to **u**
- Density connected
  - If there exist a core point **m** such that both **v** and **u** are density reachable
- Density-based cluster defined for maximal set of density connected points

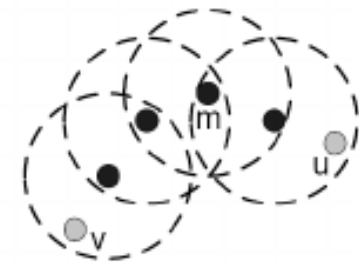
u is  
directly density  
reachable  
from v



u is  
density  
reachable  
from v



u is  
density  
connected  
from v

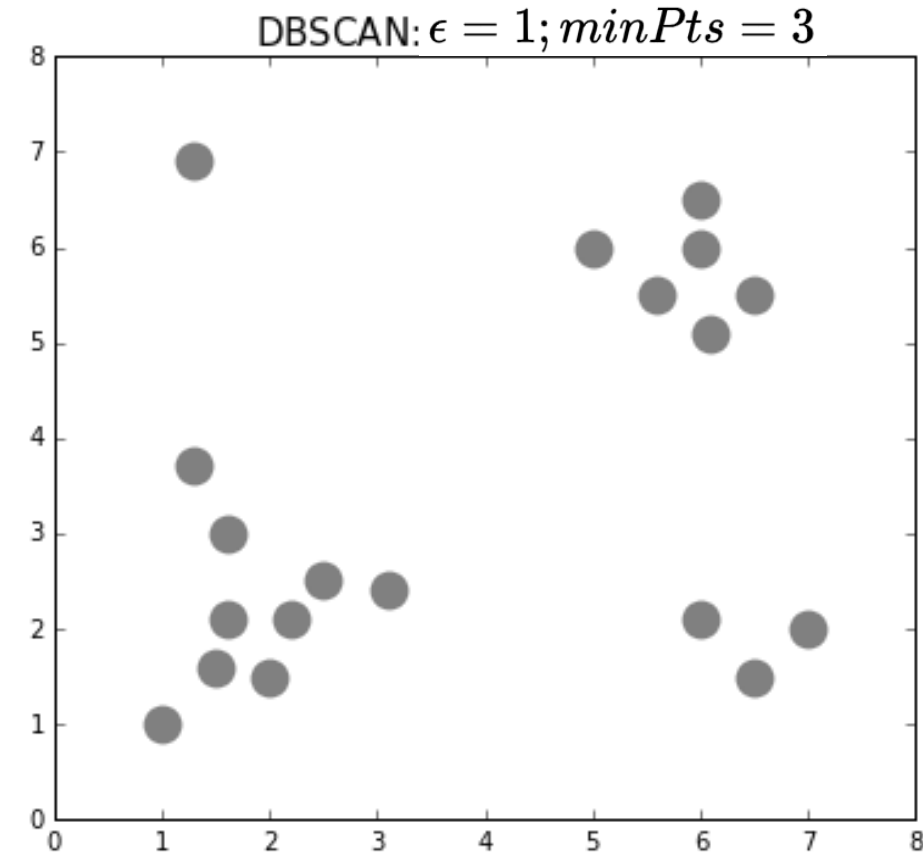


● core  
● border

# DBSCAN Algorithm



- Compute  $\epsilon$ -neighborhood for each point
- Check if this is a core point
- Recursively find all density connected point and assign to same cluster
- Border points may be reachable by core points in multiple clusters
  - Assign to one or to all overlapping clusters (“soft”)



# DBSCAN Algorithm



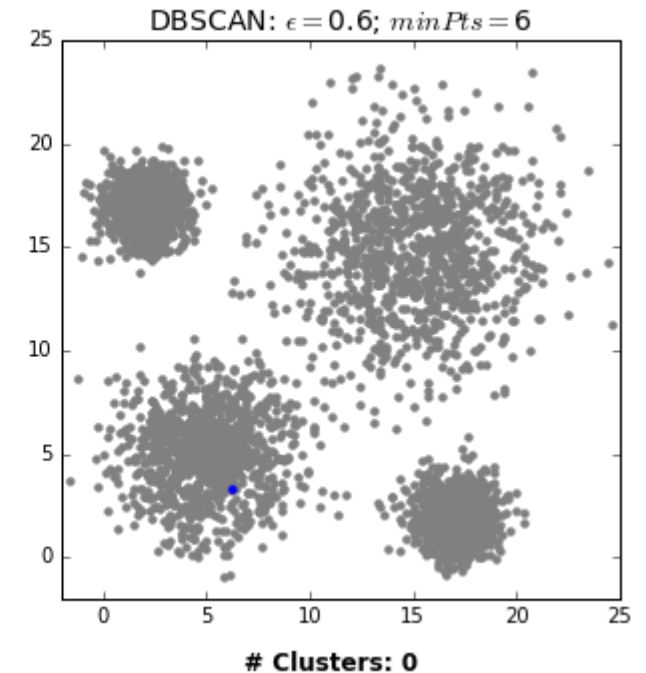
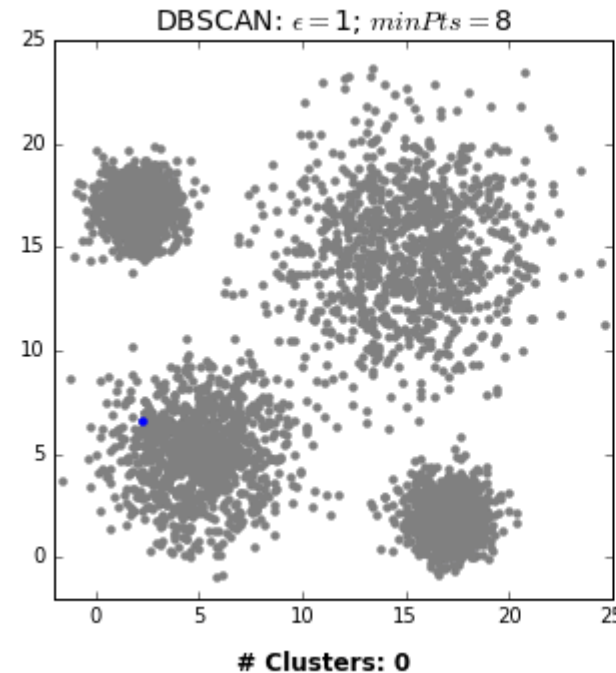
```
dbscan ( $D, \epsilon, minpts$ ):  
1  $Core \leftarrow \emptyset$   
2 foreach  $x_i \in D$  do // Find the core points  
3   Compute  $N_\epsilon(x_i)$   
4    $id(x_i) \leftarrow \emptyset$  // cluster id for  $x_i$   
5   if  $N_\epsilon(x_i) \geq minpts$  then  $Core \leftarrow Core \cup \{x_i\}$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

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

# DBSCAN Algorithm



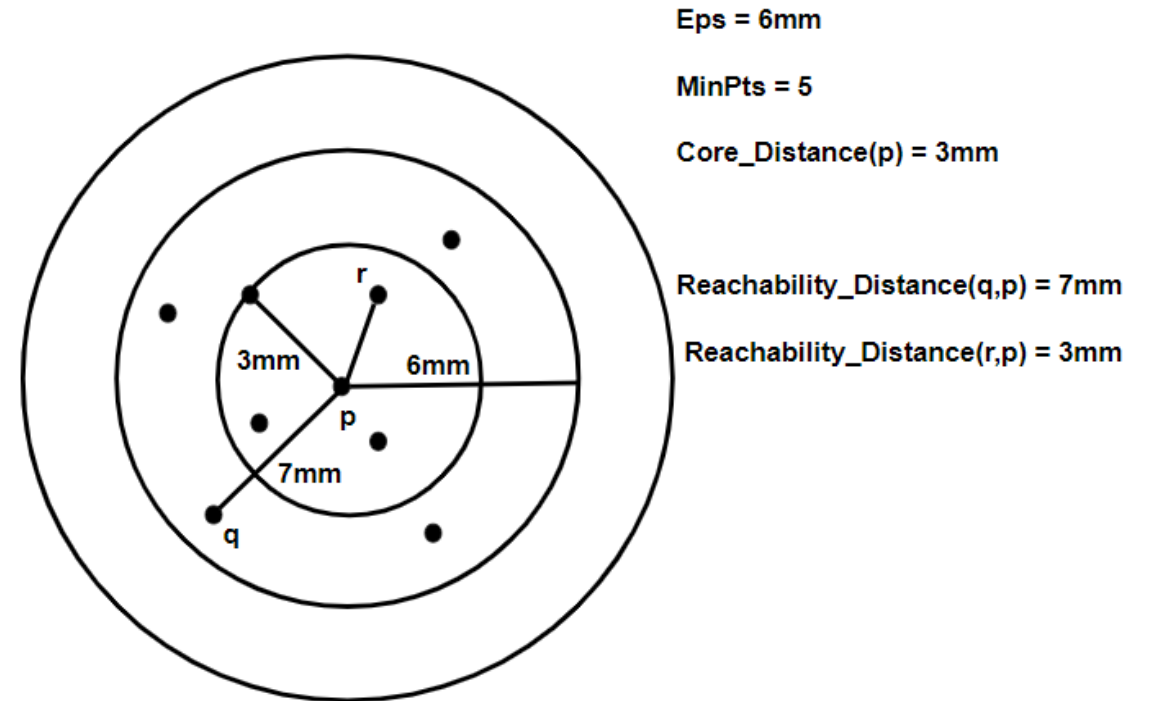
- Advantages
  - Find nonconvex clusters
  - Do not need to set the number of cluster
- Disadvantages
  - Very sensitive to selection of  $\epsilon$  and minPts
  - Difficult for data with clusters of varying density



- Ordering points to identify the clustering structure (OPTICS)
- Hierarchical DB (HDBSCSAN)



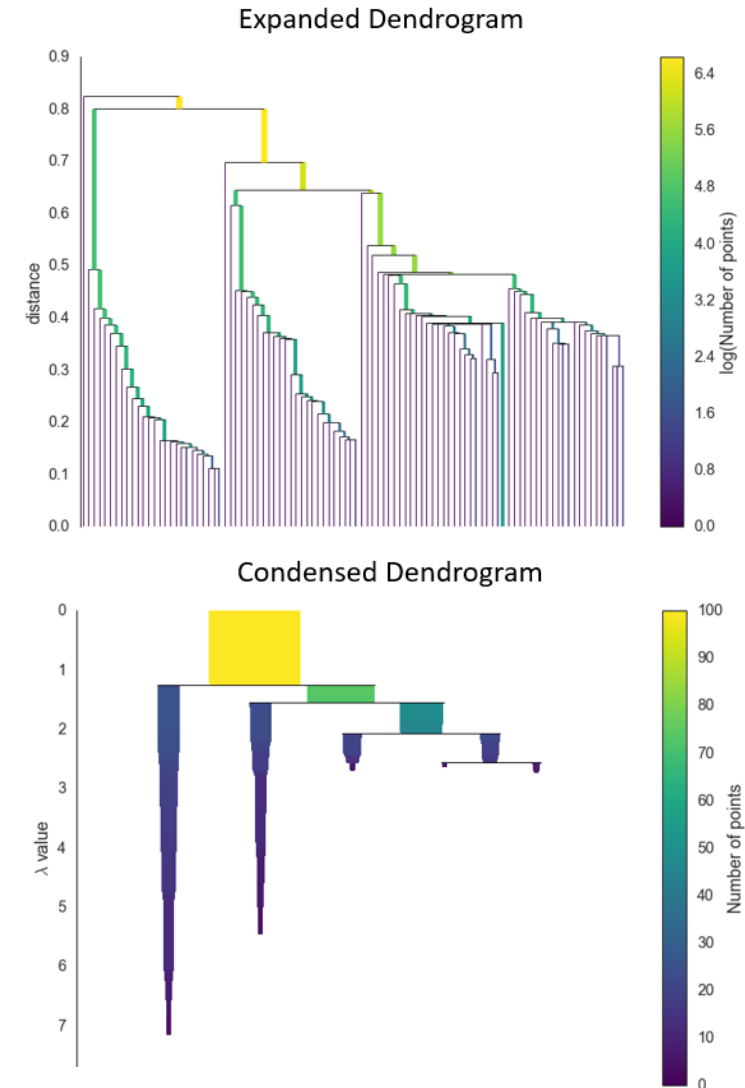
- Account for clustering of varying densities and shapes
- Creates an ordered list of points through reachability plot
- Reachability: measure of how easy it is to reach point



# HDBSCAN



- Transform space based on density-sparsity
- Build minimum spanning tree of distance weighted graph
- Construct hierarchy of connected components
- Condense hierarchy based on minimum cluster size
- Extract stable clusters from condensed tree





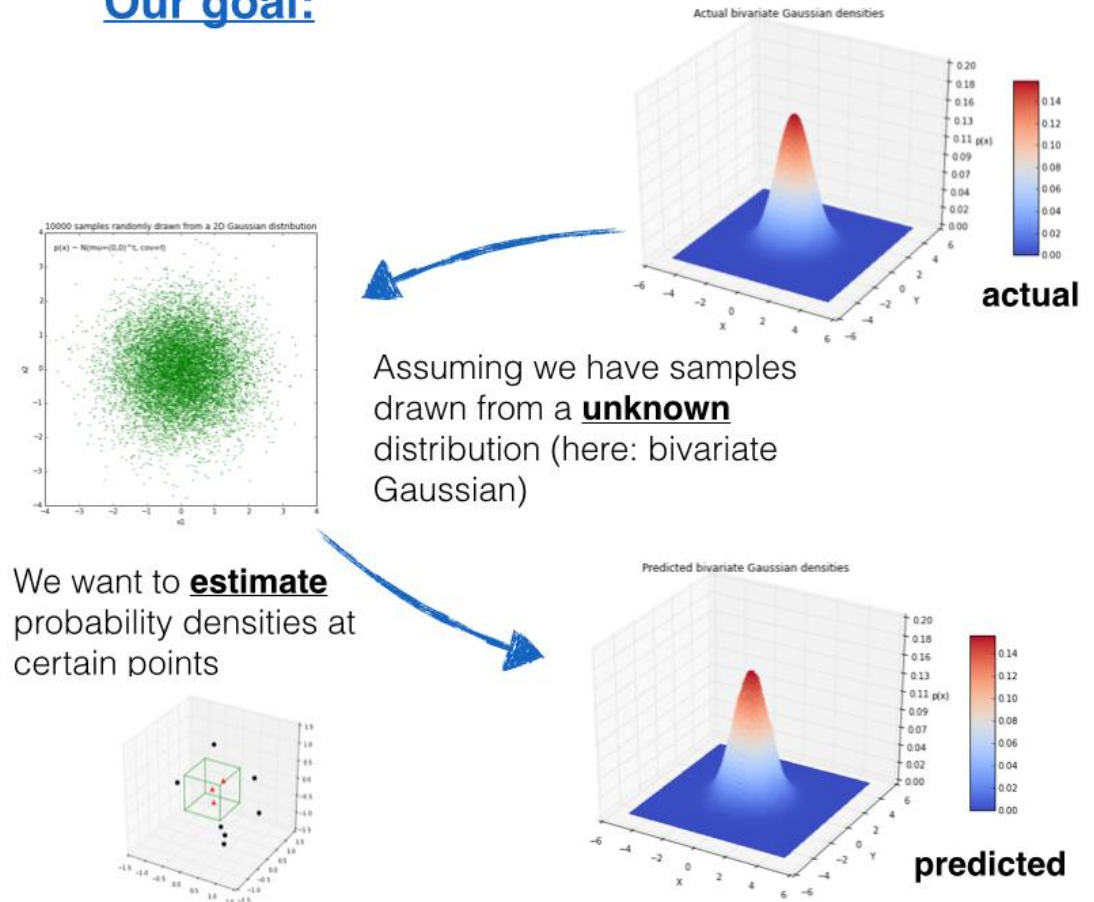
# Kernel Density Estimation

# Kernel Density Estimation (KDE)



- Close connection between density-based clustering and density estimation
- KDE seeks to determine unknown PDF by finding dense regions of points

Our goal:



- Can use CDF to estimate PDF
- $k$  is the number of points that lie in a window of width  $h$  centered on  $x$
- Density estimation is ratio of fraction of points in window and volume of window

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

$$\hat{f}(x) = \frac{\hat{F}(x + \frac{h}{2}) - \hat{F}(x - \frac{h}{2})}{h} = \frac{k/n}{h} = \frac{k}{nh}$$

- Need to define kernel function
- Need to be non-negative and integrates to 1 for all values
- Discrete kernel uses indicator function

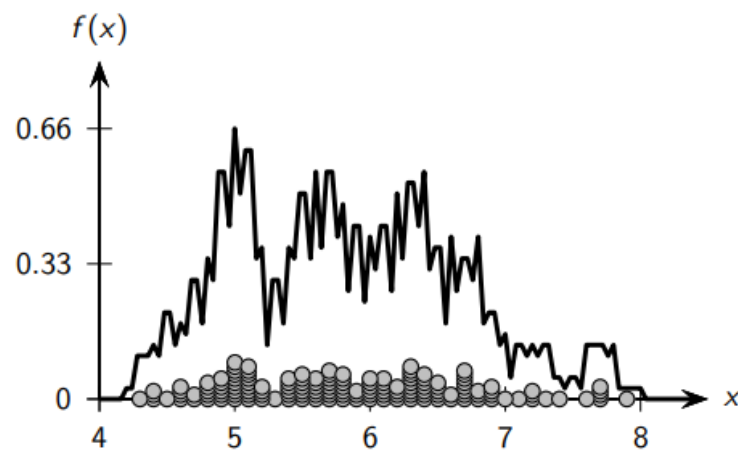
$$K(z) = \begin{cases} 1 & \text{If } |z| \leq \frac{1}{2} \\ 0 & \text{Otherwise} \end{cases}$$

Discrete kernel

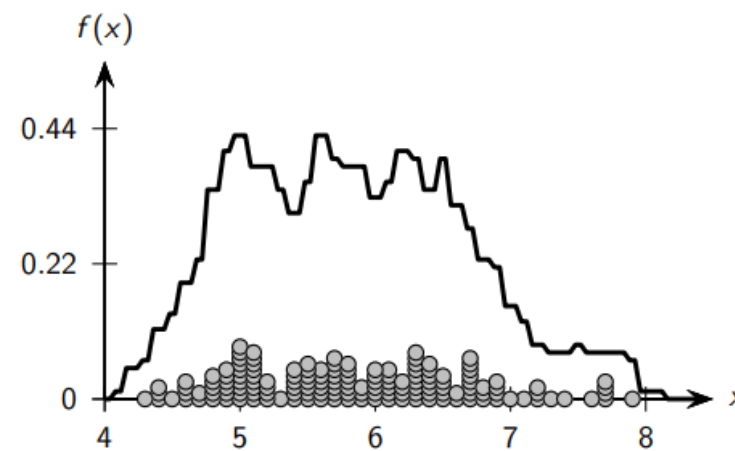
$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$

Density estimate

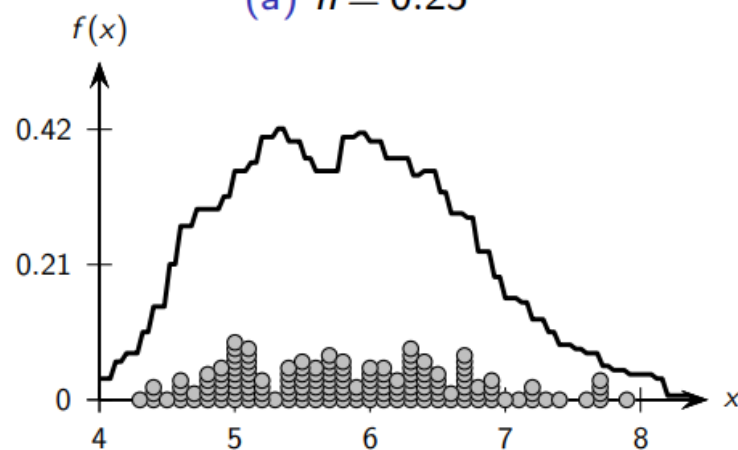
# Discrete KDE Example



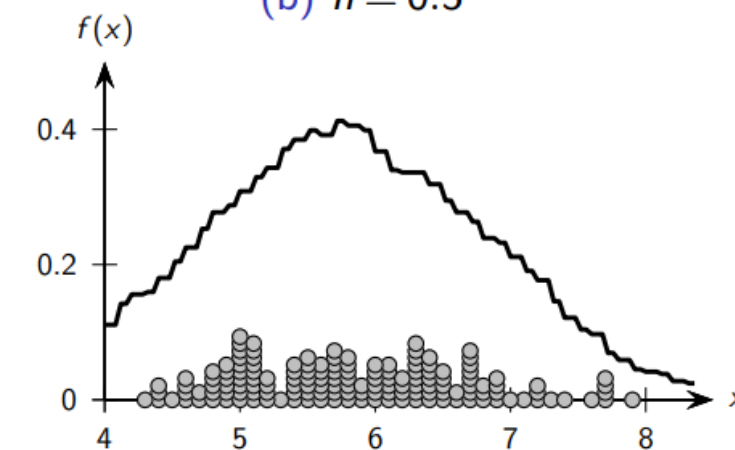
(a)  $h = 0.25$



(b)  $h = 0.5$



(c)  $h = 1.0$

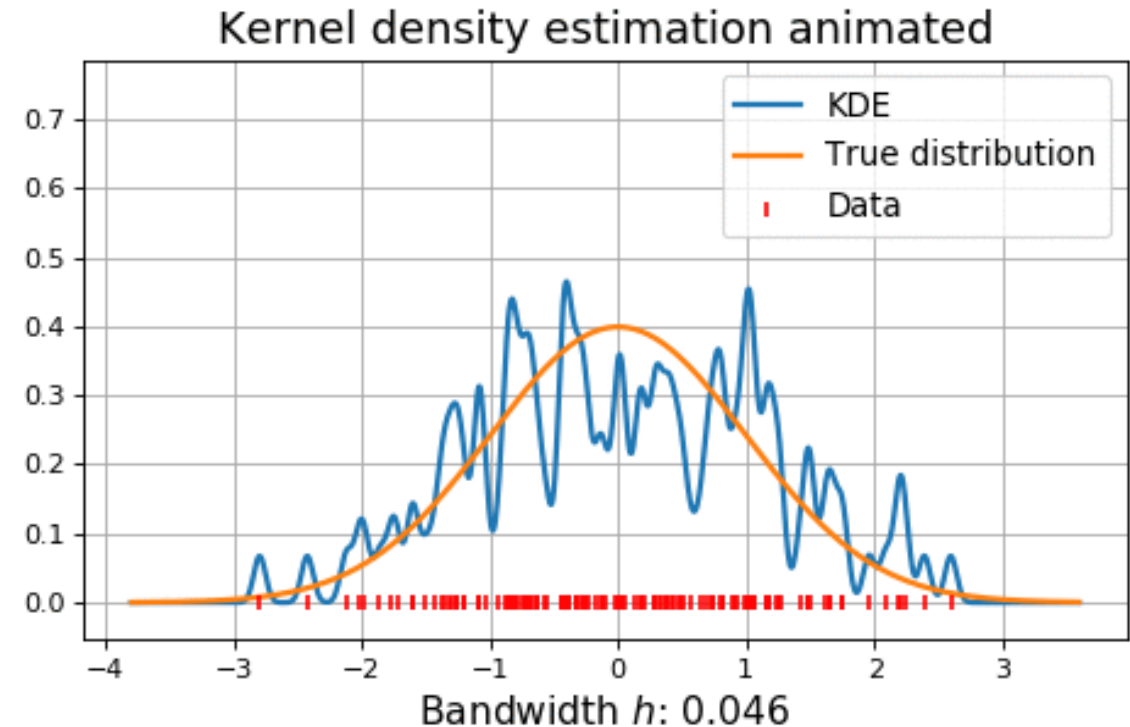


(d)  $h = 2.0$

# Gaussian Kernel Estimator



- Width parameter  $h$  controls spread or smoothness of estimate
- Discrete kernel function has abrupt changes
- Gaussian kernel provides smooth transition





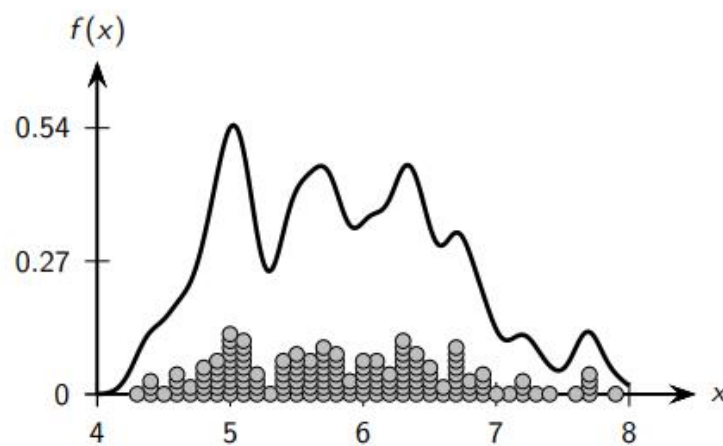
- Width parameter  $h$  controls spread or smoothness of estimate
- Discrete kernel function has abrupt changes
- Gaussian kernel provides smooth transition

$$K(z) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{z^2}{2} \right\}$$

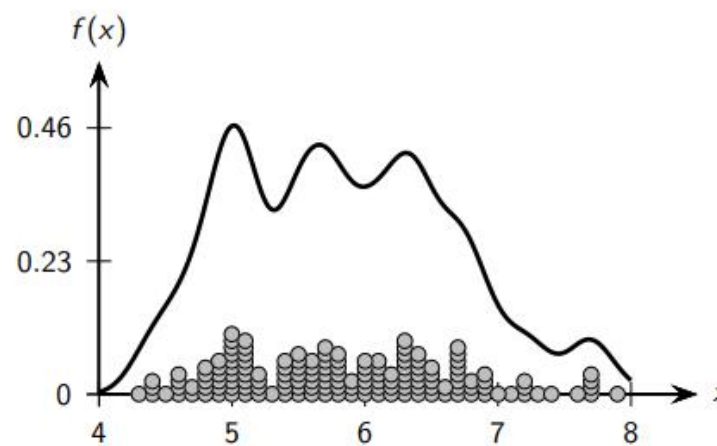
$$K \left( \frac{x - x_i}{h} \right) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{(x - x_i)^2}{2h^2} \right\}$$

Density estimate

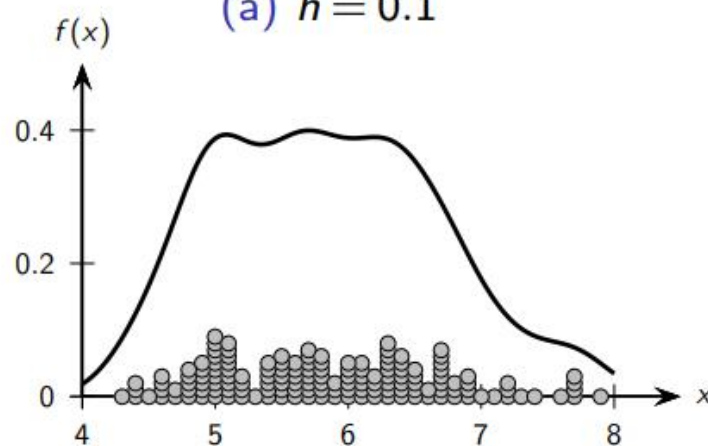
# Discrete KDE Example



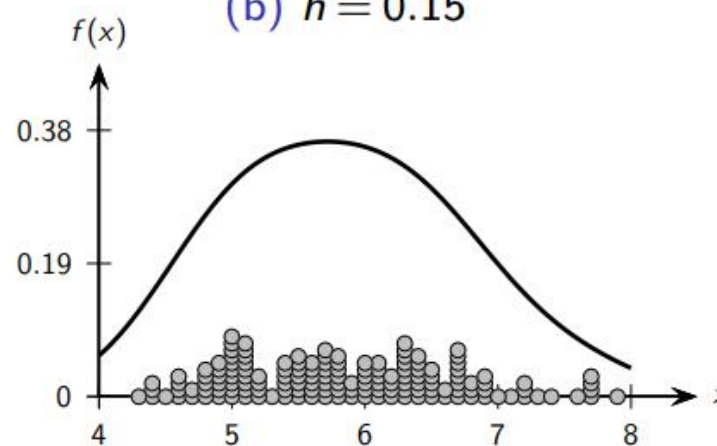
(a)  $h = 0.1$



(b)  $h = 0.15$



(c)  $h = 0.25$



(d)  $h = 0.5$

- Define d-dimension “window” as hypercube with edge length  $h$
- Density estimation is ratio of fraction of points in window and volume of window
- Kernel function still must integrate to 1

$$\text{vol}(H_d(h)) = h^d$$

Volume of hypercube

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)$$

Density estimate

$$\int K(\mathbf{z}) d\mathbf{z} = 1$$

$$K(\mathbf{z}) = \begin{cases} 1 & \text{If } |z_j| \leq \frac{1}{2}, \text{ for all dimensions } j = 1, \dots, d \\ 0 & \text{Otherwise} \end{cases}$$

Discrete kernel

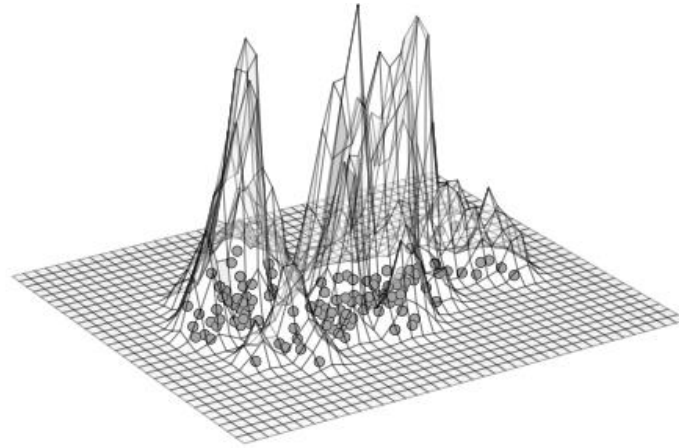
$$K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) = \frac{1}{(2\pi)^{d/2}} \exp\left\{-\frac{(\mathbf{x} - \mathbf{x}_i)^T(\mathbf{x} - \mathbf{x}_i)}{2h^2}\right\}$$

Gaussian kernel

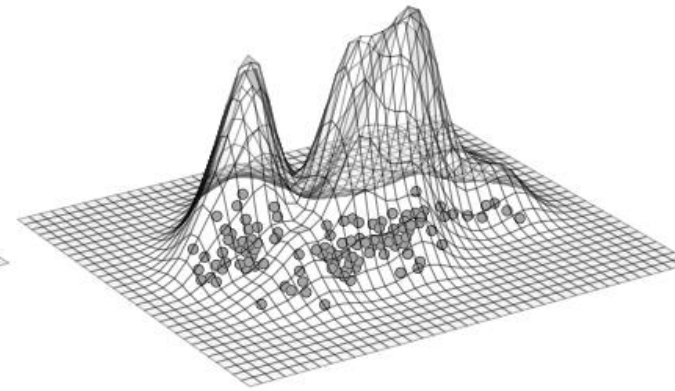
# Gaussian KDE (Iris 2D)



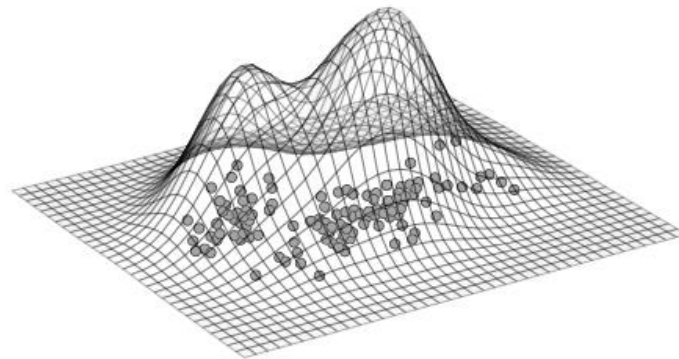
TEXAS A&M UNIVERSITY  
Engineering



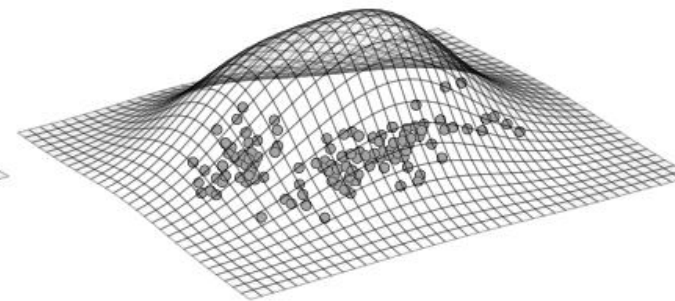
(a)  $h = 0.1$



(b)  $h = 0.2$



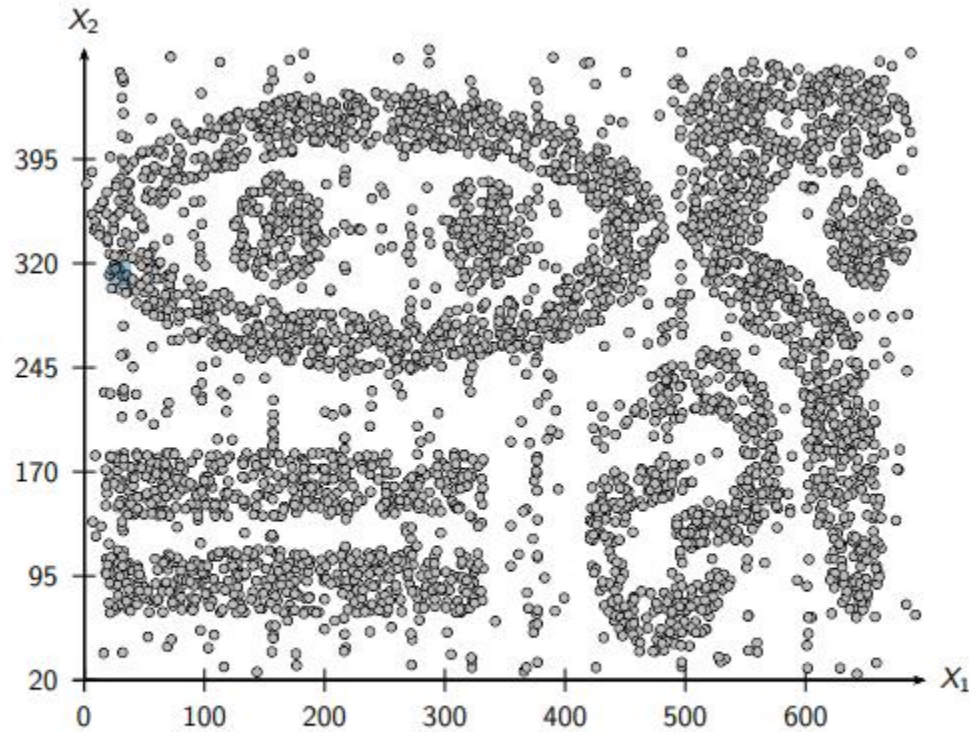
(c)  $h = 0.35$



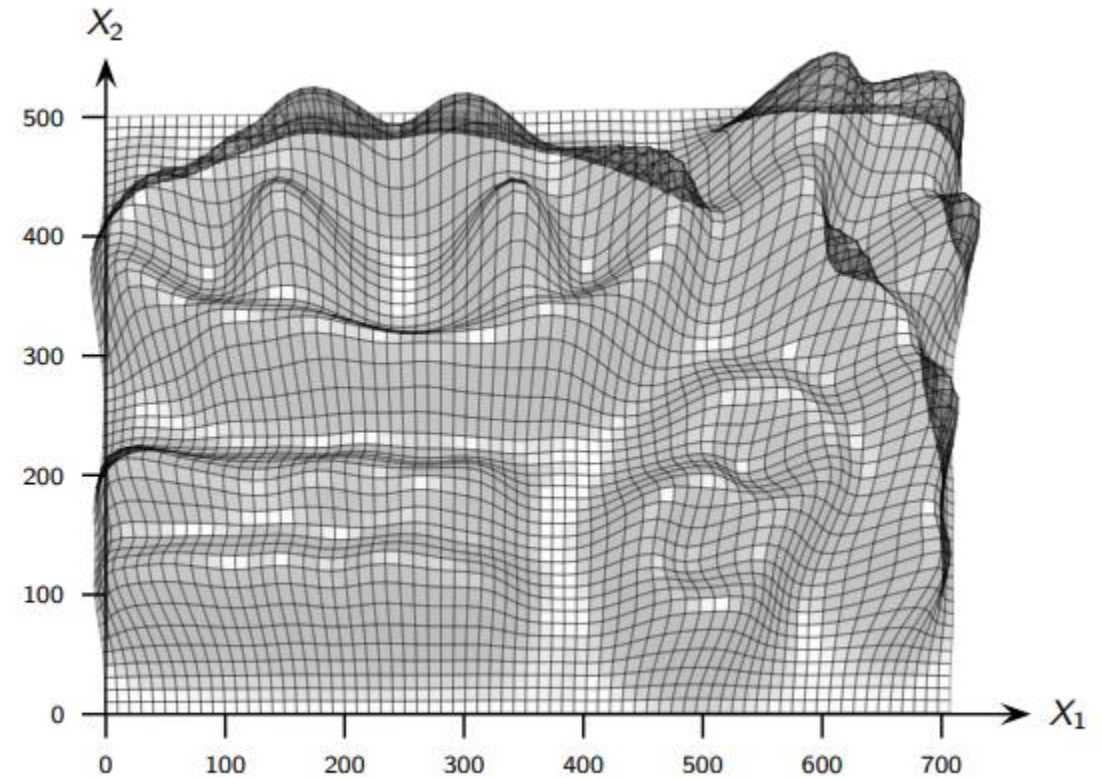
(d)  $h = 0.6$



# Gaussian KDE



(a) Original Points



(b) Gaussian Density Estimation

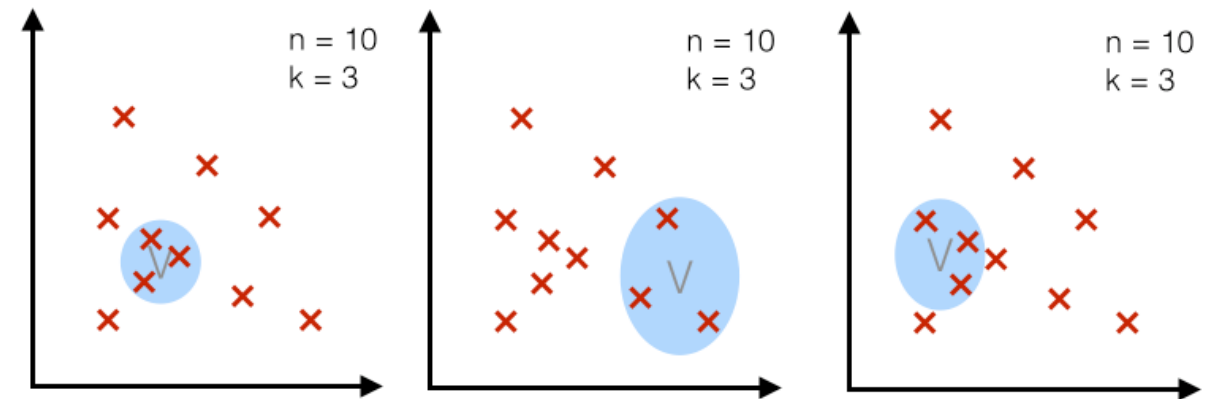
# Nearest Neighbor Density Estimation



TEXAS A&M UNIVERSITY  
Engineering

- Previously, fixed volume by fixing  $h$
- Alternative: fix number of points and adapt volume
- $k$  is number of neighbors

$$\hat{f}(\mathbf{x}) = \frac{k}{n \text{vol}(S_d(h_{\mathbf{x}}))}$$

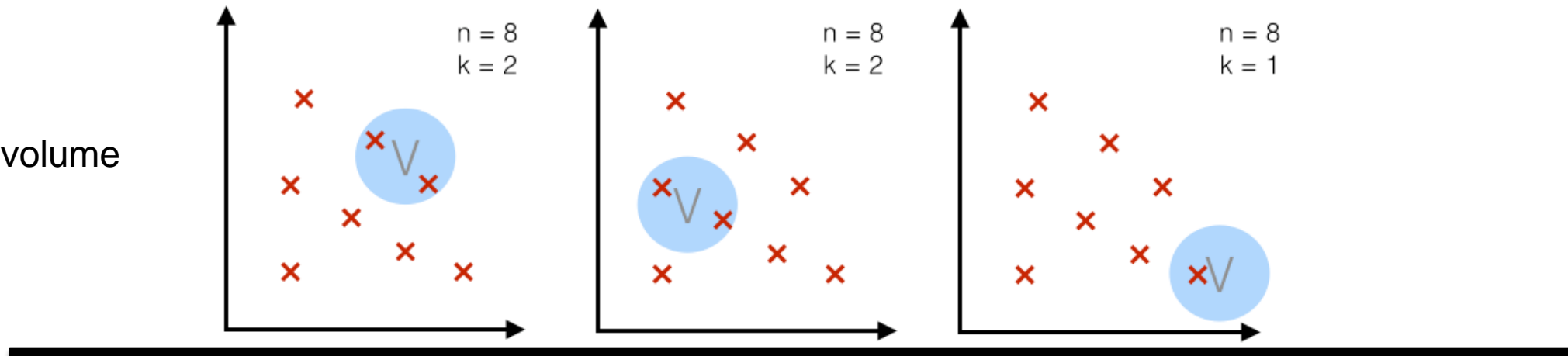


# Nearest Neighbor Density Estimation

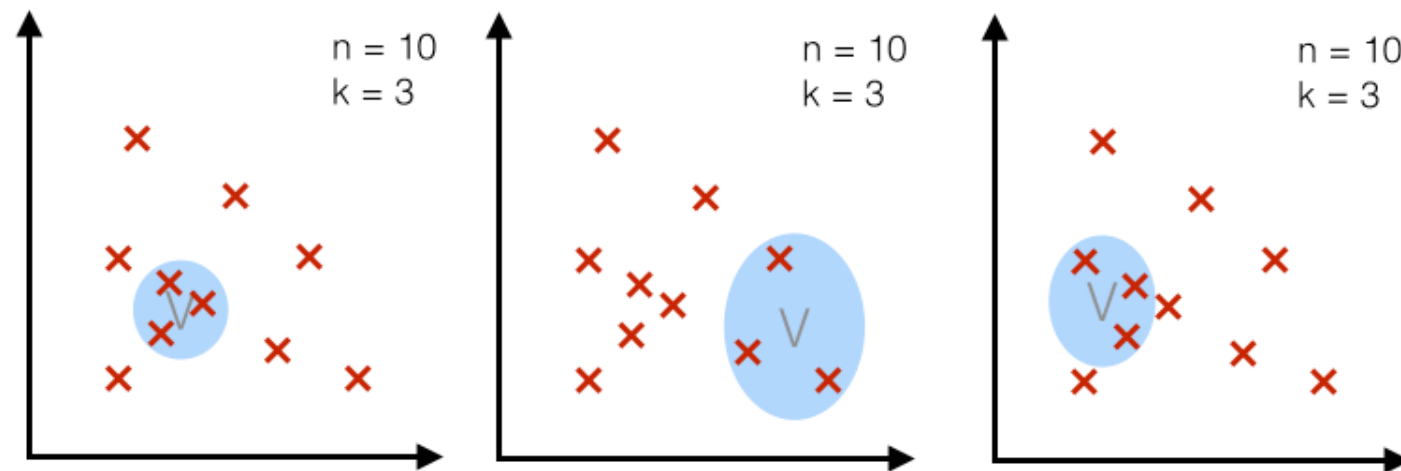


TEXAS A&M UNIVERSITY  
Engineering

Fixed volume



Fixed  $k$





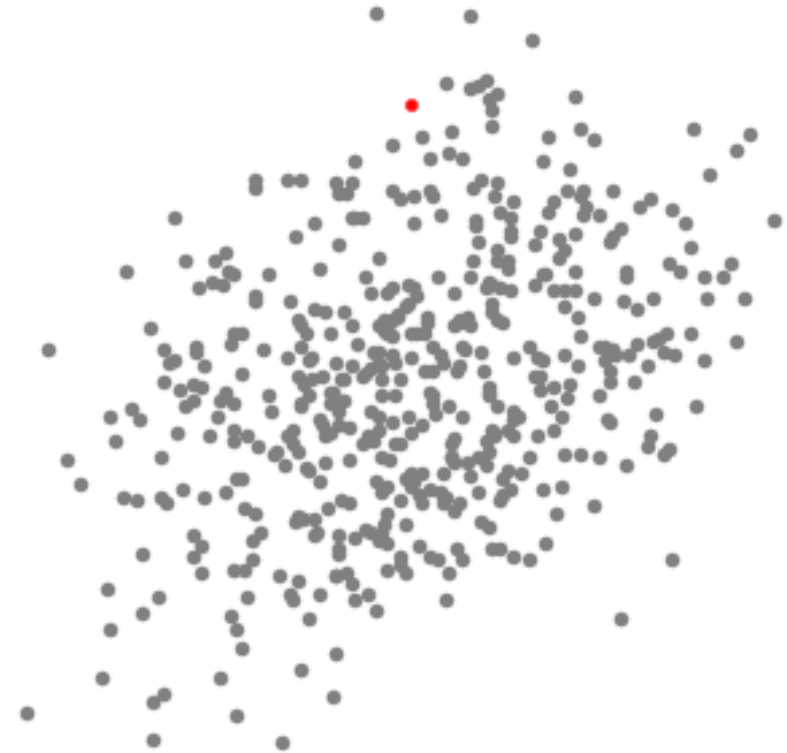


# Mean Shift Algorithm

# Mean Shift Algorithm



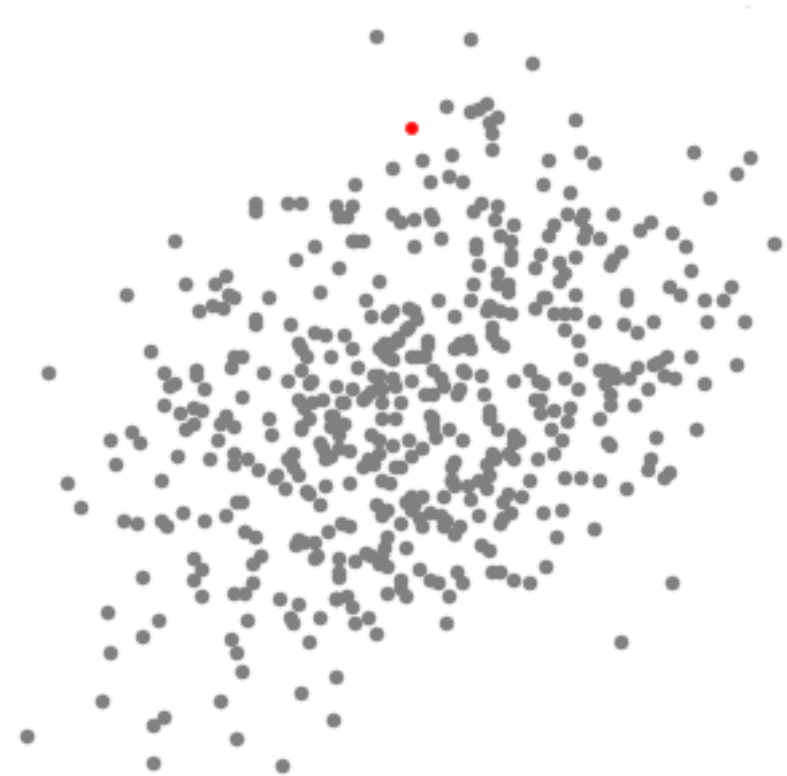
- Useful for clusters with arbitrary shapes and not well-separated data
- Idea: shift each data toward the mode of the distribution of points within a given radius



# Mean Shift Algorithm Steps



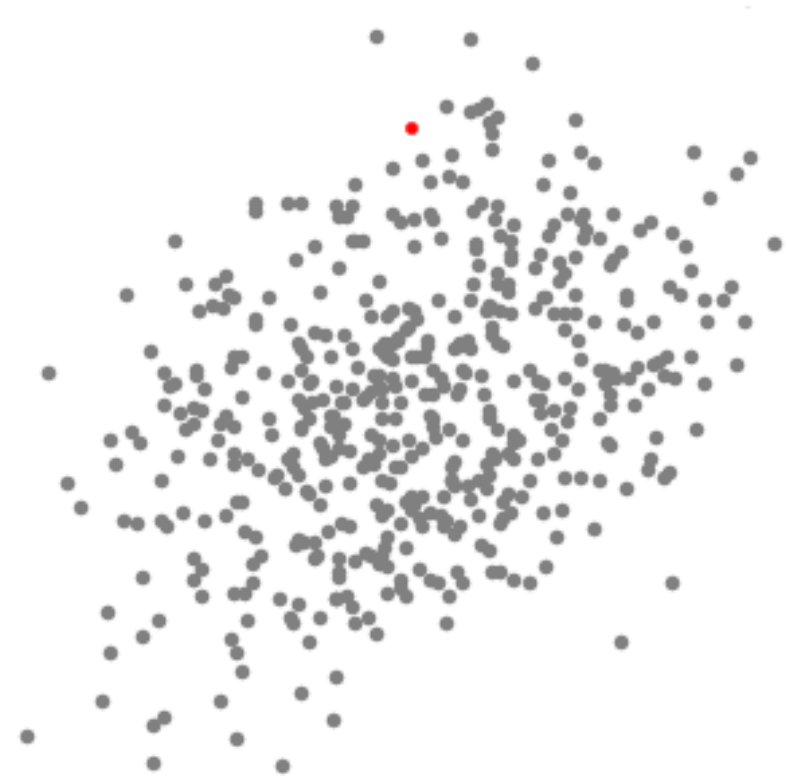
- Initialize data points as cluster centers
- Repeat until convergence
  - Compute the mean of all data points within a certain radius (kernel)
  - Shift the data to the mean
  - Identify the cluster centroids as points that have not moved
  - Return cluster assignments



# Mean Shift Algorithm Steps



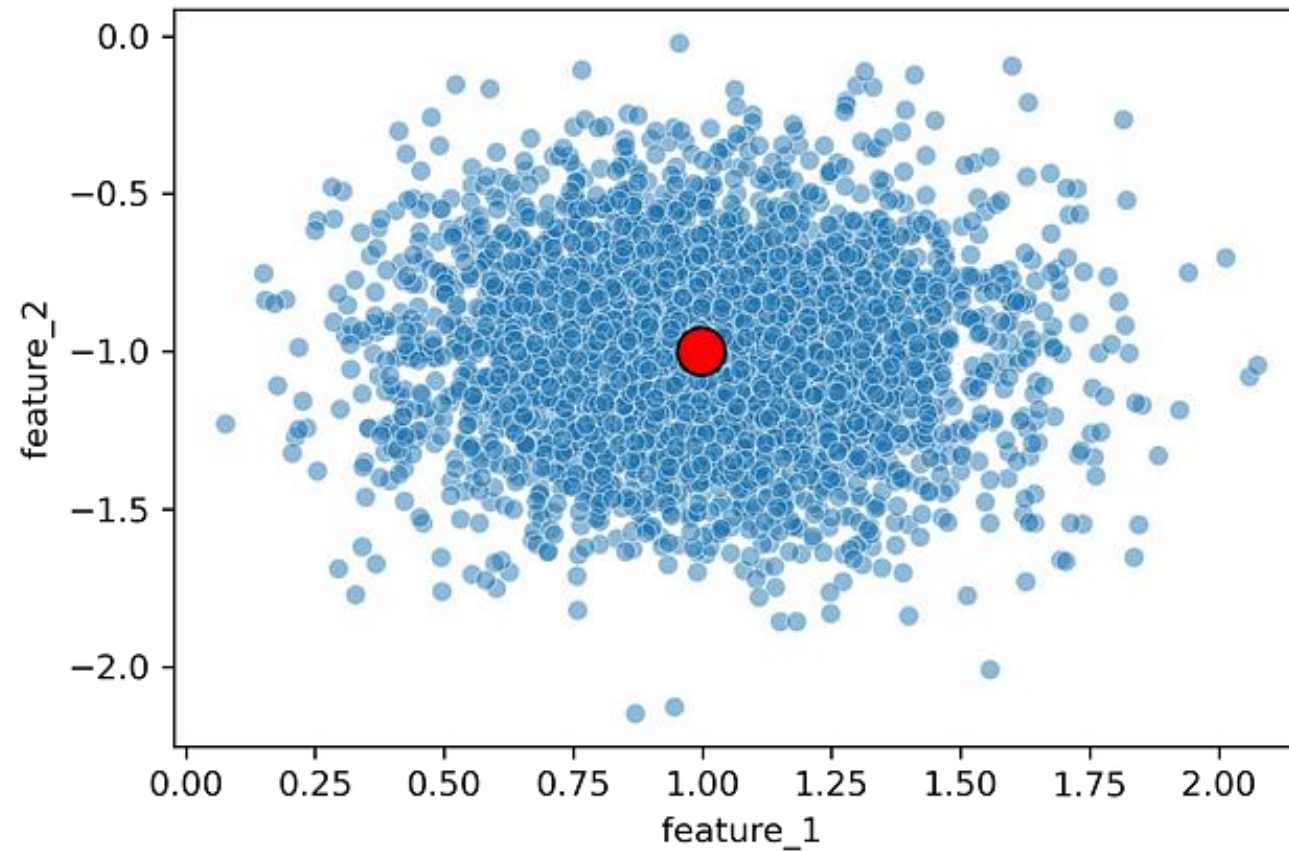
- Initialize data points as cluster centers
- Repeat until convergence
  - **Compute the mean of all data points within a certain radius (kernel)**
  - Shift the data to the mean
  - Identify the cluster centroids as points that have not moved
  - Return cluster assignments



# Mean Shift Algorithm: Kernel Function



TEXAS A&M UNIVERSITY  
Engineering

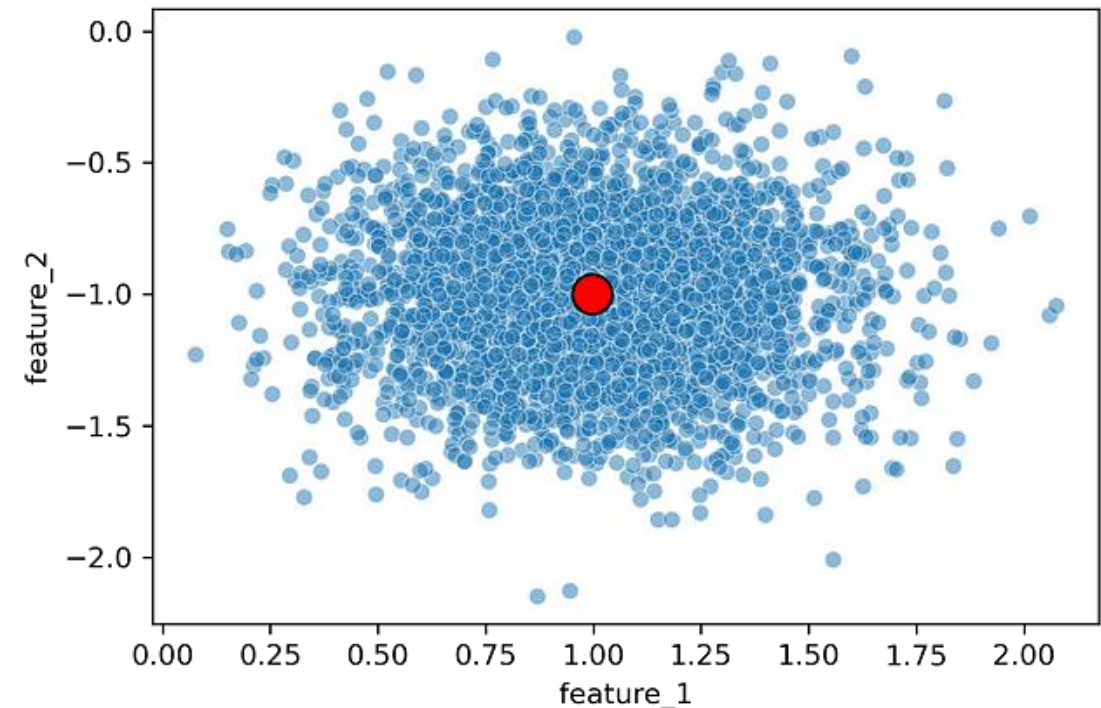


# Mean Shift Algorithm: Kernel Function



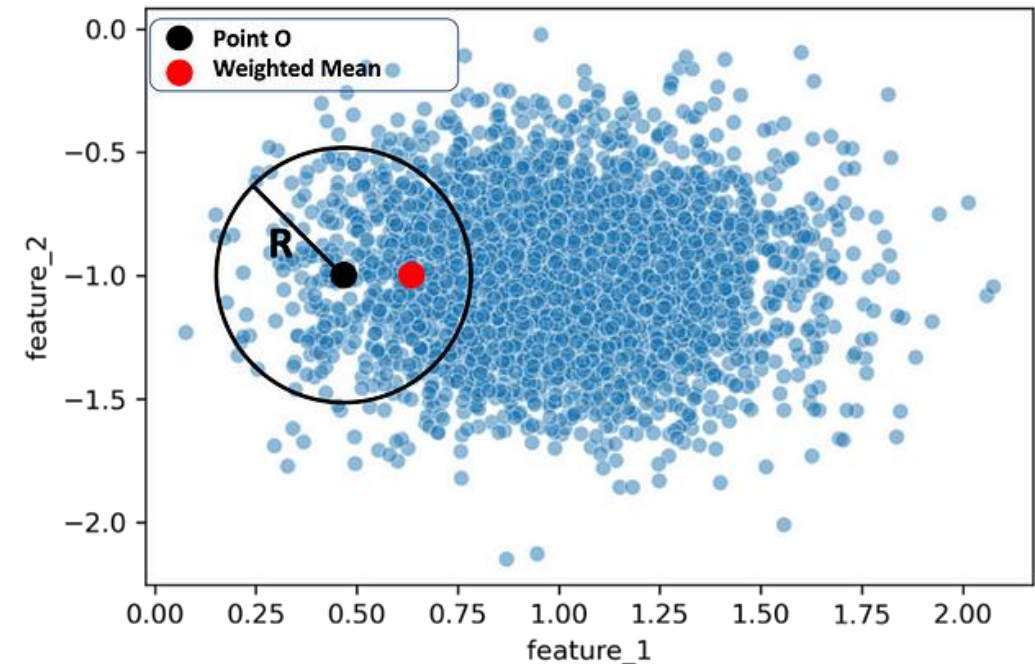
- Easily compute mean of data
- Will be more meaningful if weighted mean

$$M_W = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$



# Mean Shift Algorithm: Kernel Function

- Can use “flat” or Gaussian kernel for weights



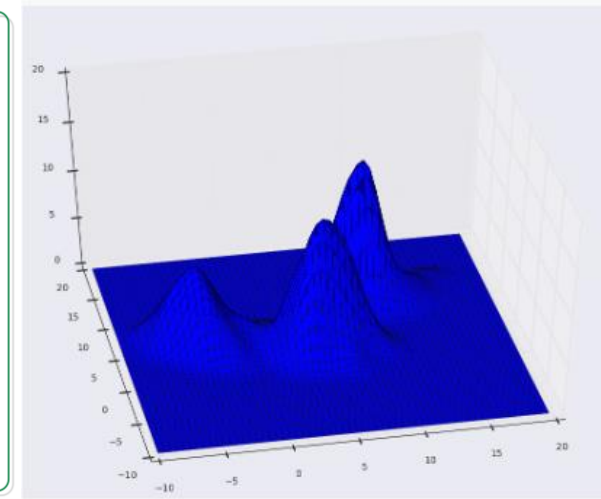
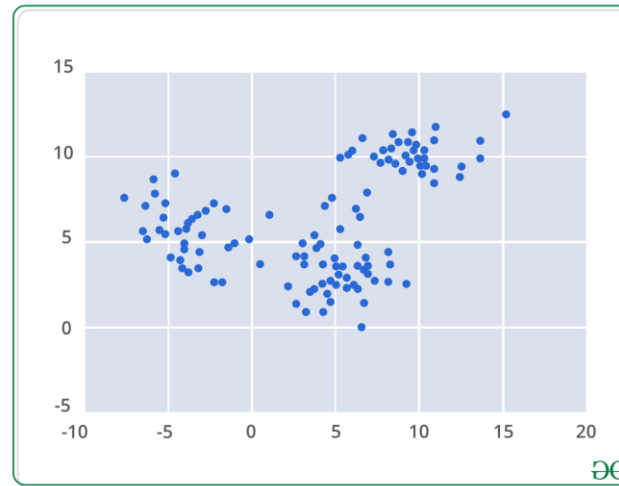
$$w(d) = \begin{cases} 1, & \text{if } d \leq R \\ 0, & \text{if } d > R \end{cases}$$

# Mean Shift Algorithm: Kernel Function



TEXAS A&M UNIVERSITY  
Engineering

- Can use “flat” or Gaussian kernel for weights



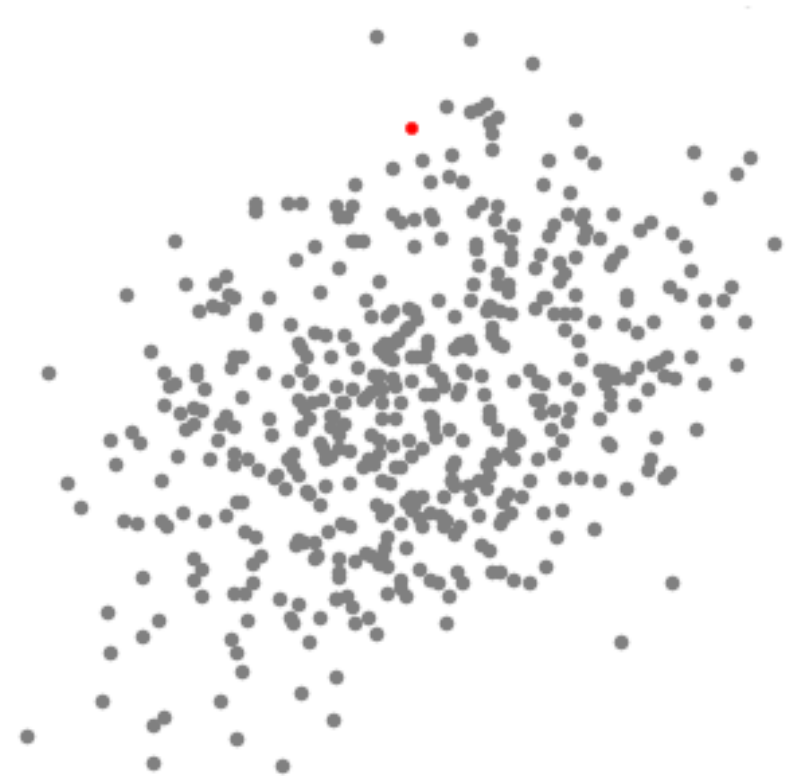
$$w(d) = e^{-\frac{d}{2\sigma^2}}$$



# Mean Shift Algorithm Steps



- Initialize data points as cluster centers
- Repeat until convergence
  - Compute the mean of all data points within a certain radius (kernel)
  - **Shift the data to the mean**
  - Identify the cluster centroids as points that have not moved
  - Return cluster assignments



# Mean Shift Algorithm

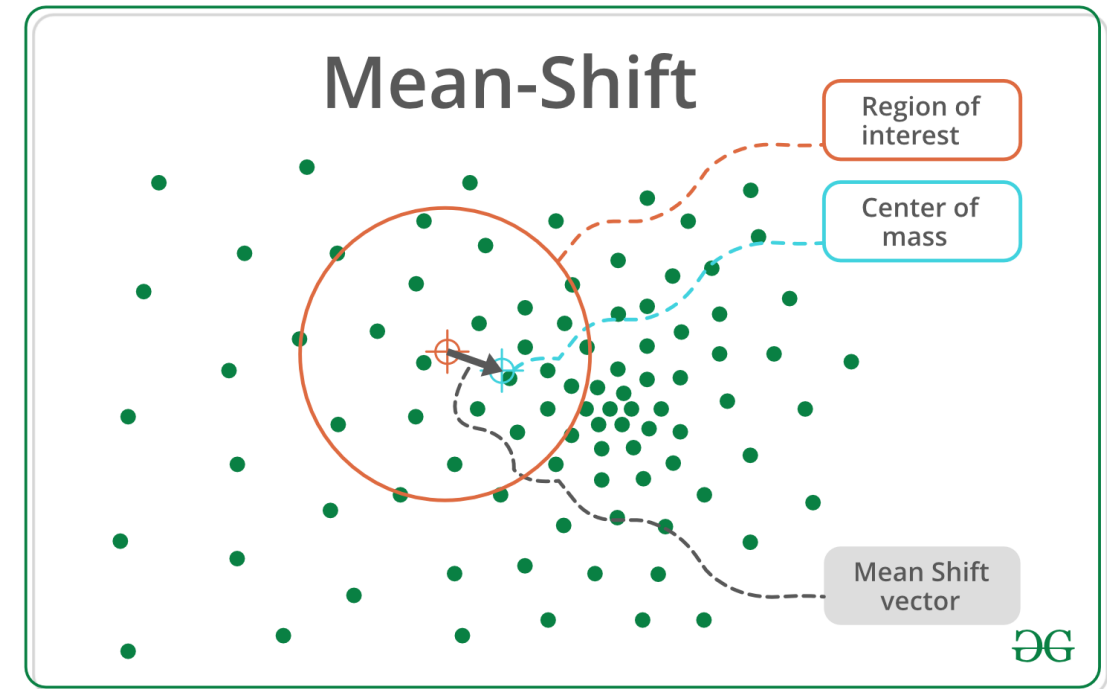
- “Shift” local mean based on neighborhood defined by kernel

$$x^{t+1} = x^t + m(x^t)$$

Update centroid location

$$m(x) = \frac{1}{|N(x)|} \sum_{x_j \in N(x)} x_j - x = \frac{\sum_{x_j \in N(x)} K(x_j - x) x_j}{\sum_{x_j \in N(x)} K(x_j - x)} - x$$

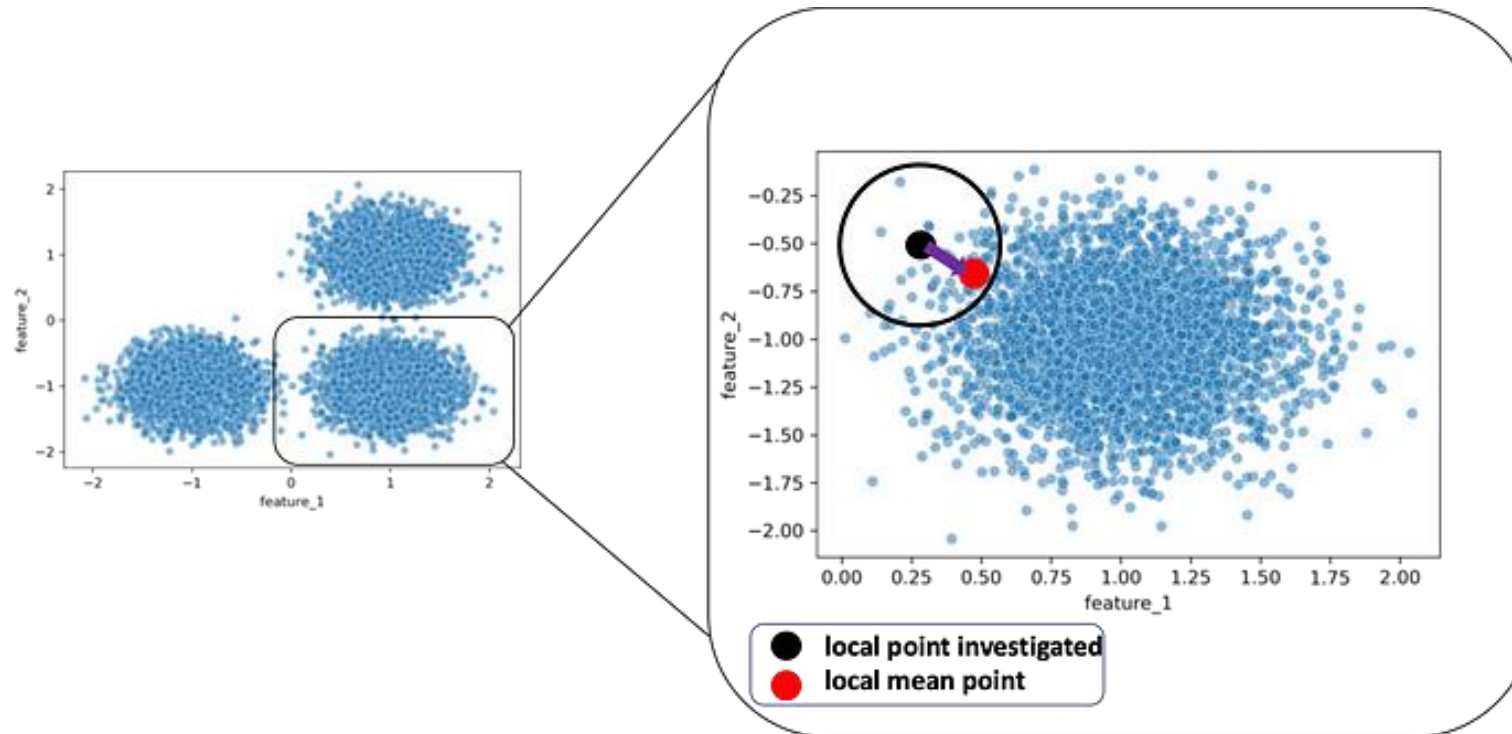
Mean-shift vector



# Mean Shift Algorithm: Shifting



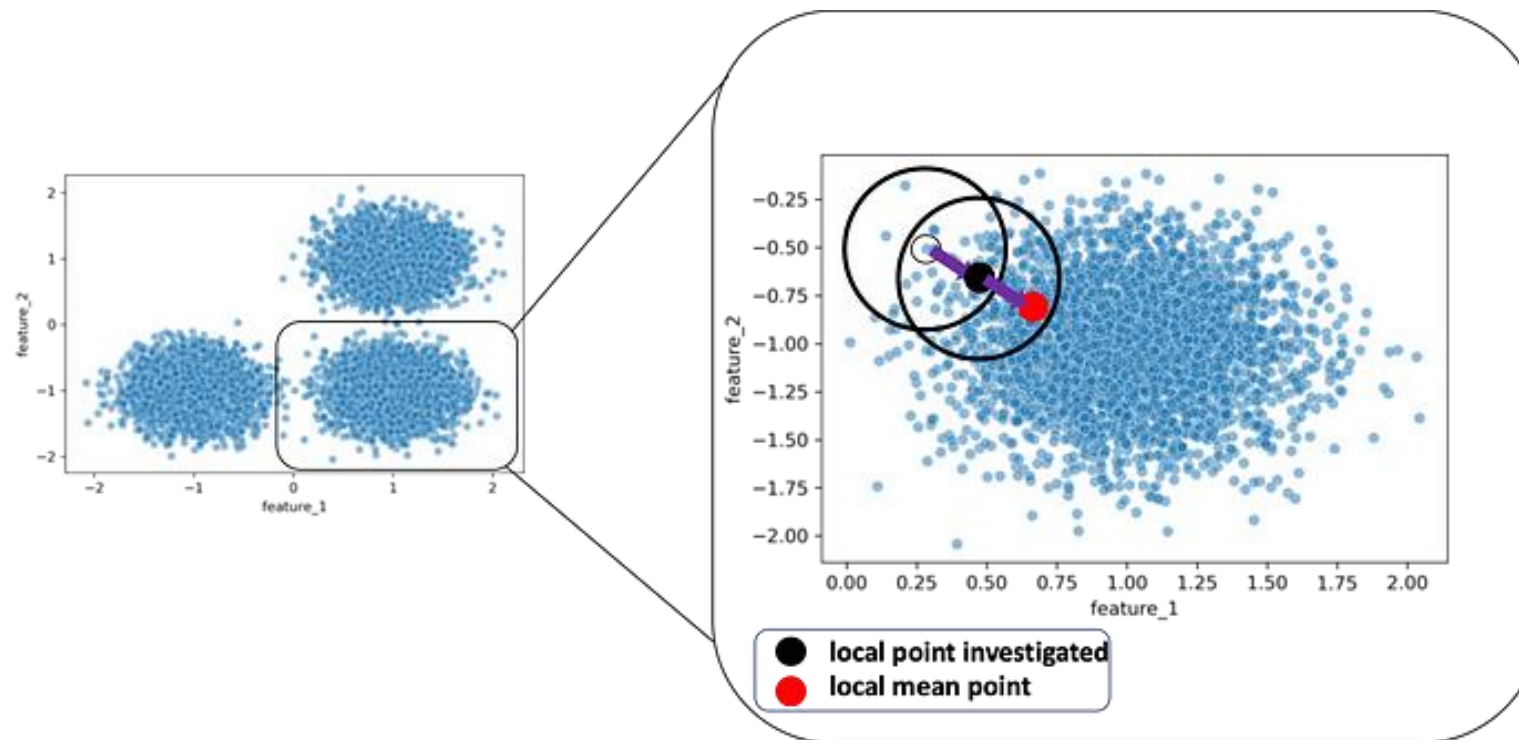
- Compute weighted mean based on area defined by kernel



# Mean Shift Algorithm: Shifting



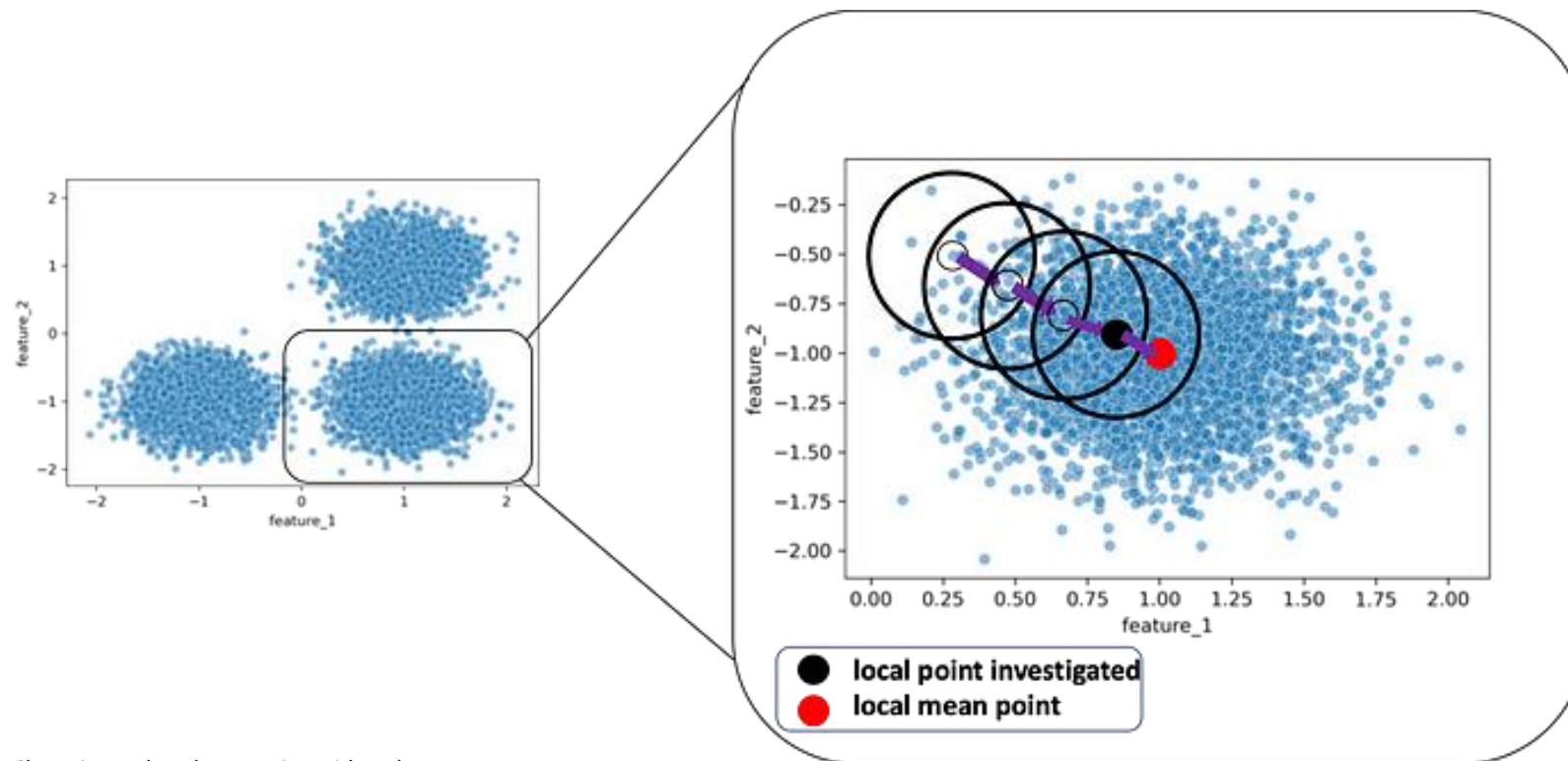
- Repeat procedure for new position (i.e., neighborhood)



# Mean Shift Algorithm: Shifting



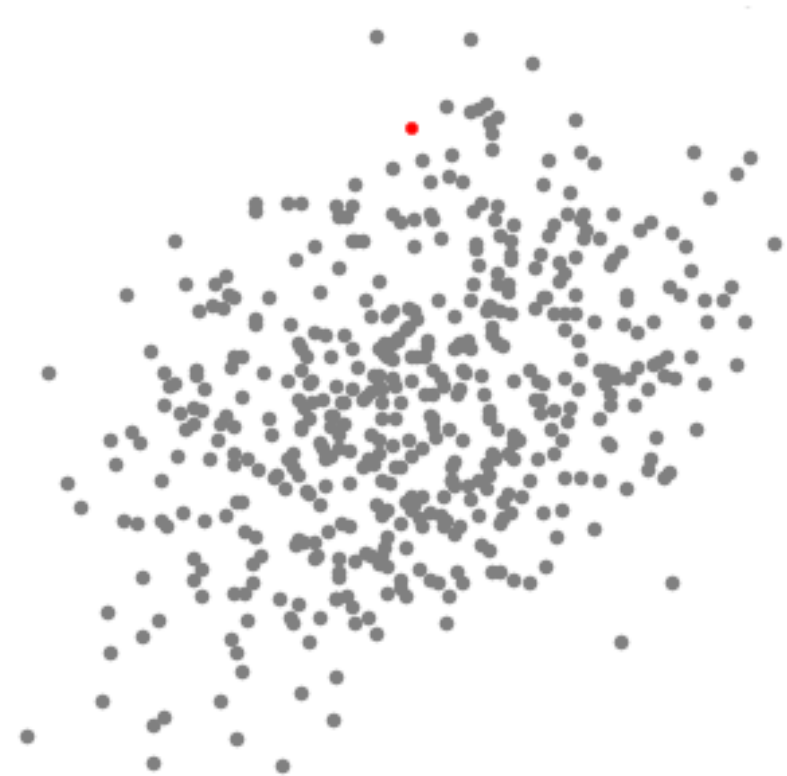
- Terminate after enough iterations or until the number of points within a cluster no longer increase (convergence)



# Mean Shift Algorithm Steps



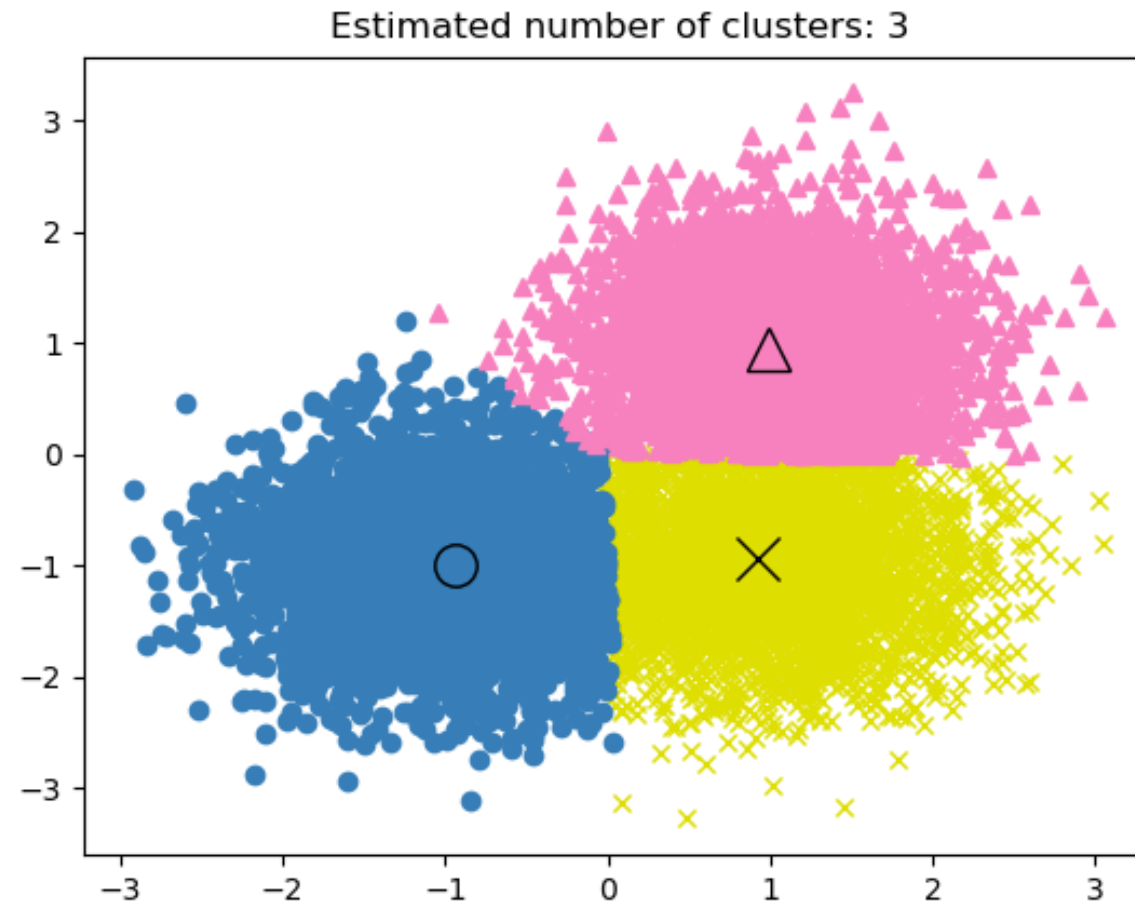
- Initialize data points as cluster centers
- Repeat until convergence
  - Compute the mean of all data points within a certain radius (kernel)
  - Shift the data to the mean
  - **Identify the cluster centroids as points that have not moved**
  - **Return cluster assignments**



# Mean Shift Algorithm Cluster Assignments



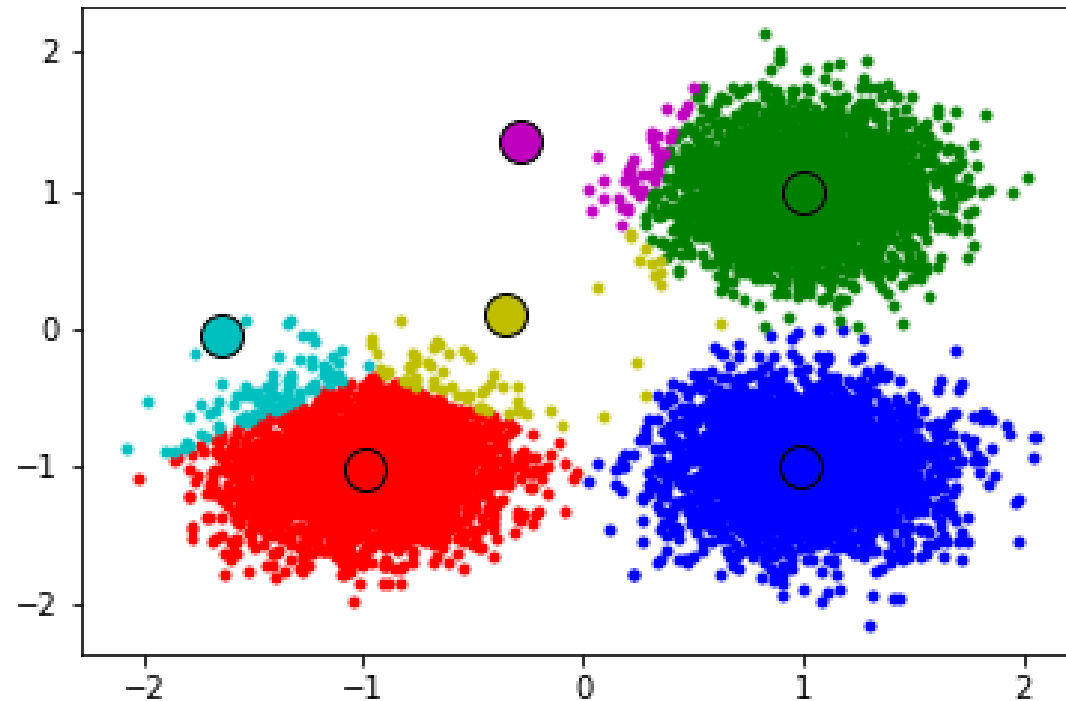
TEXAS A&M UNIVERSITY  
Engineering



# Mean Shift Algorithm: Bandwidth

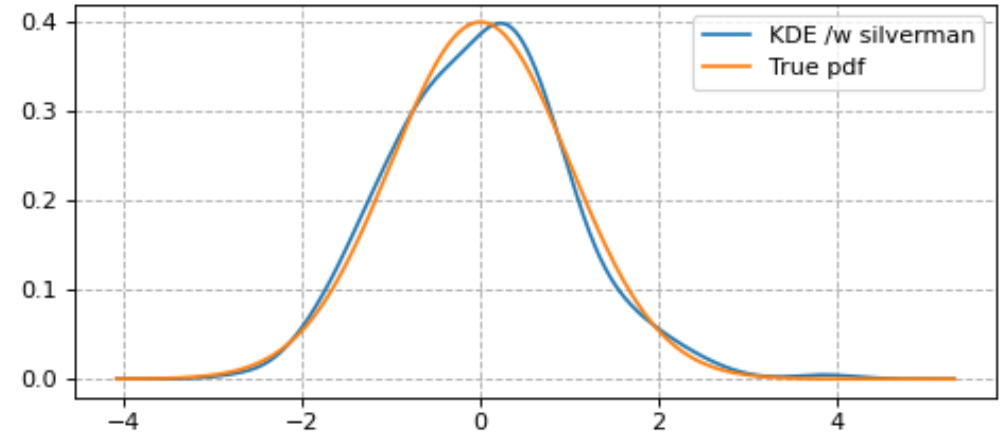


- Too small of a bandwidth can lead to “non-sense” clusters

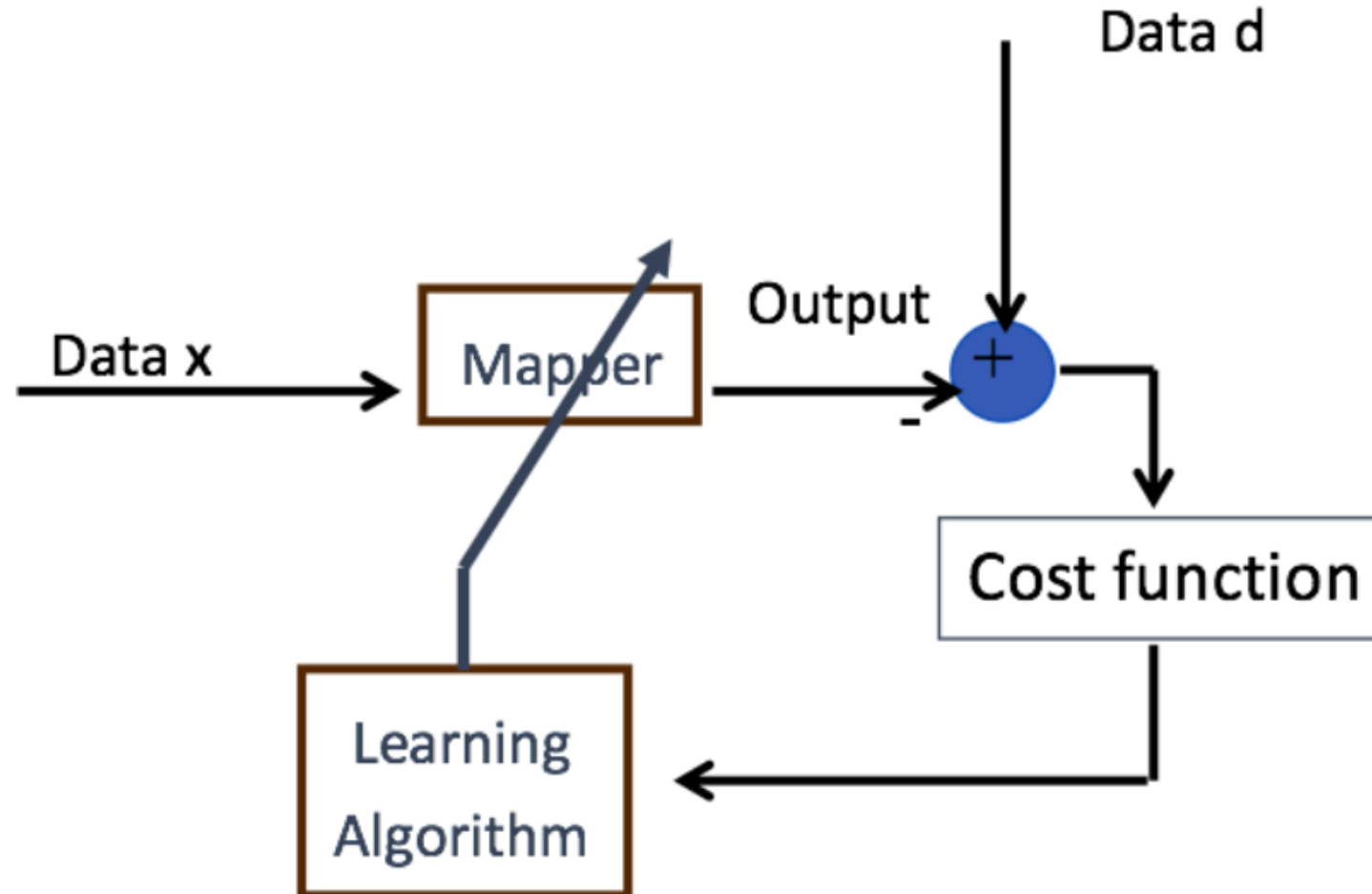




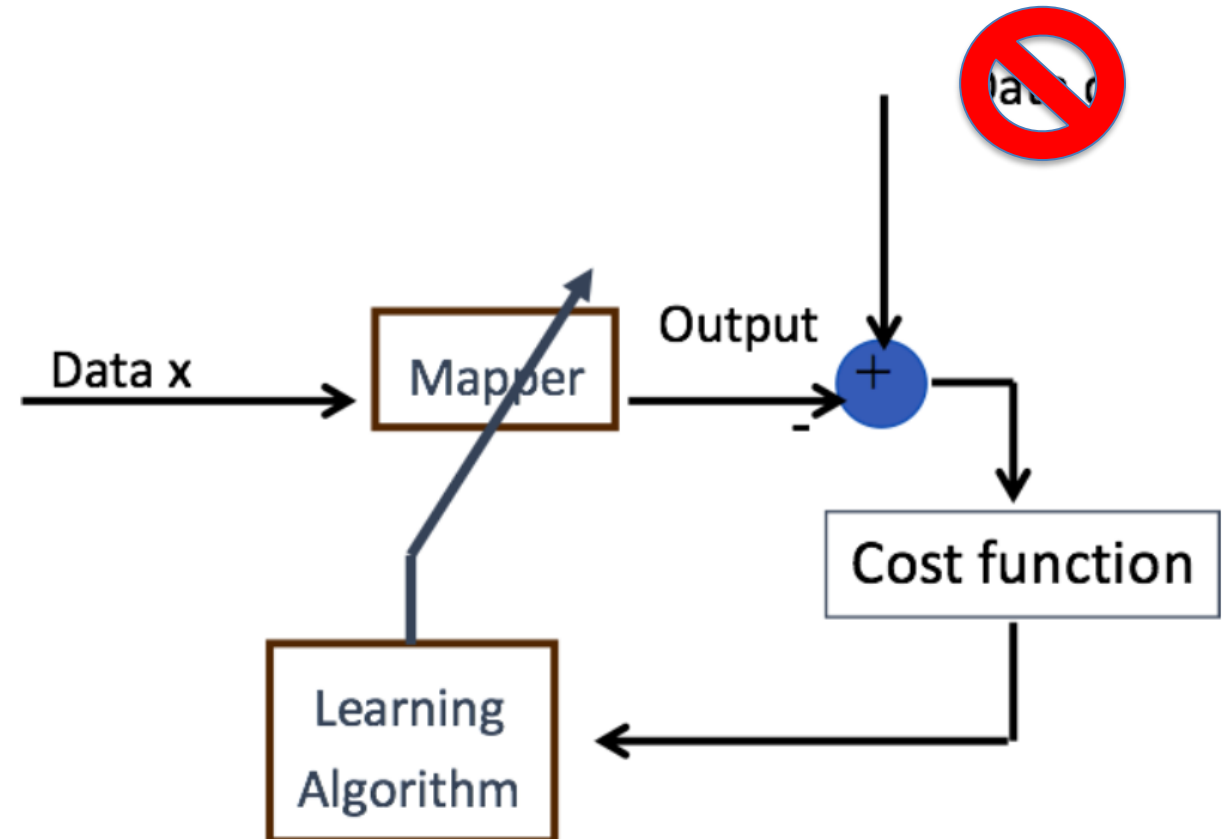
- Sklearn built-in function to estimate bandwidth based on data
- Silverman's rule
  - Assumes Gaussian distribution
- Empirically
- Etc.



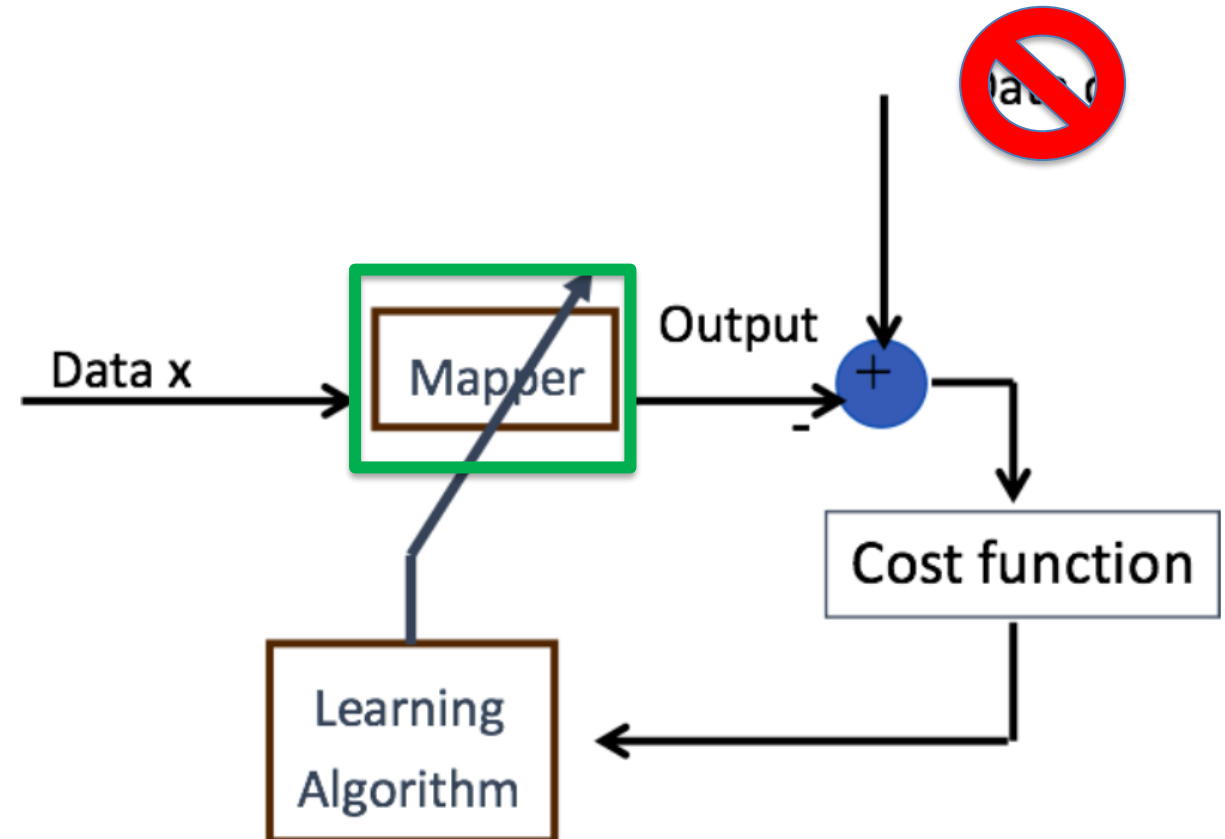
$$h = 0.9 \min \left( \hat{\sigma}, \frac{IQR}{1.34} \right) n^{-\frac{1}{5}}$$



- Unsupervised: No labels,  $d$

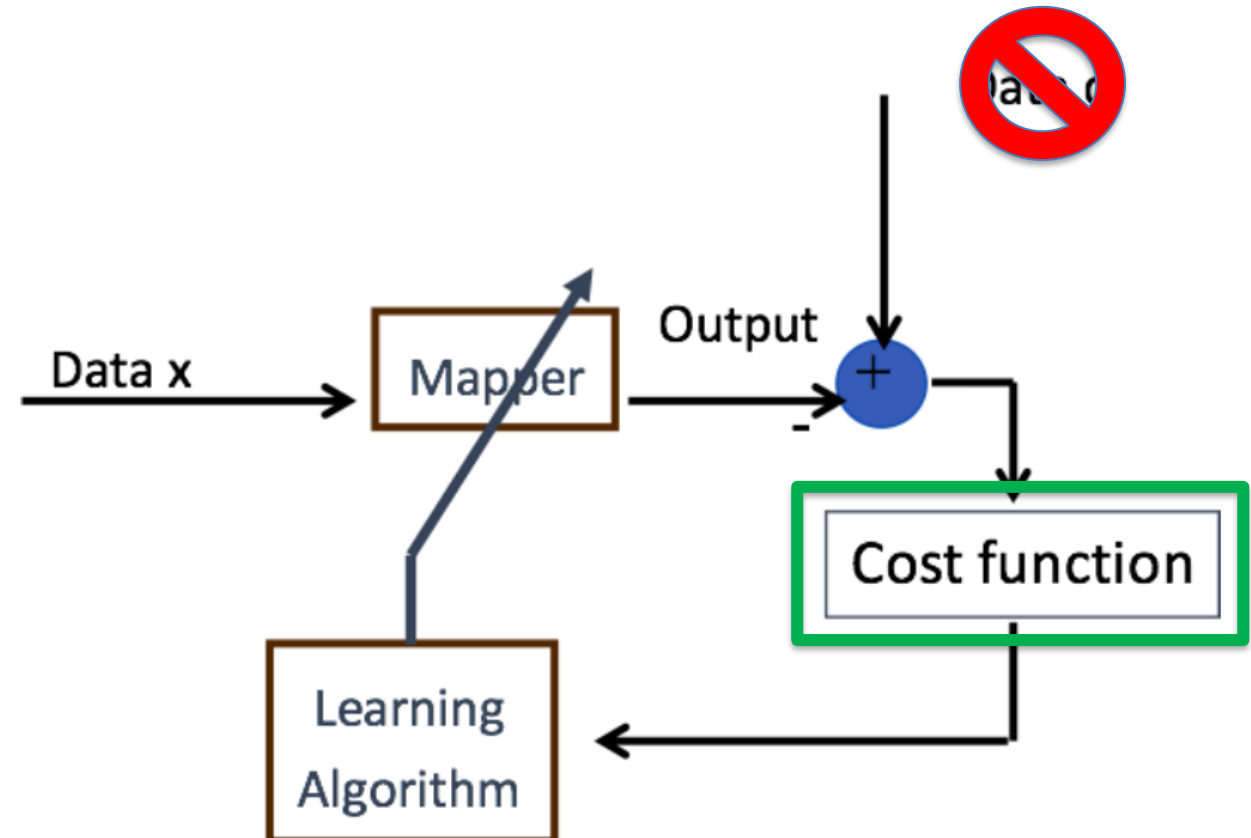


- Unsupervised: No labels,  $d$
- **Mapper:**
  - Density-based Clustering
  - Takes input data and groups into  $k$  clusters

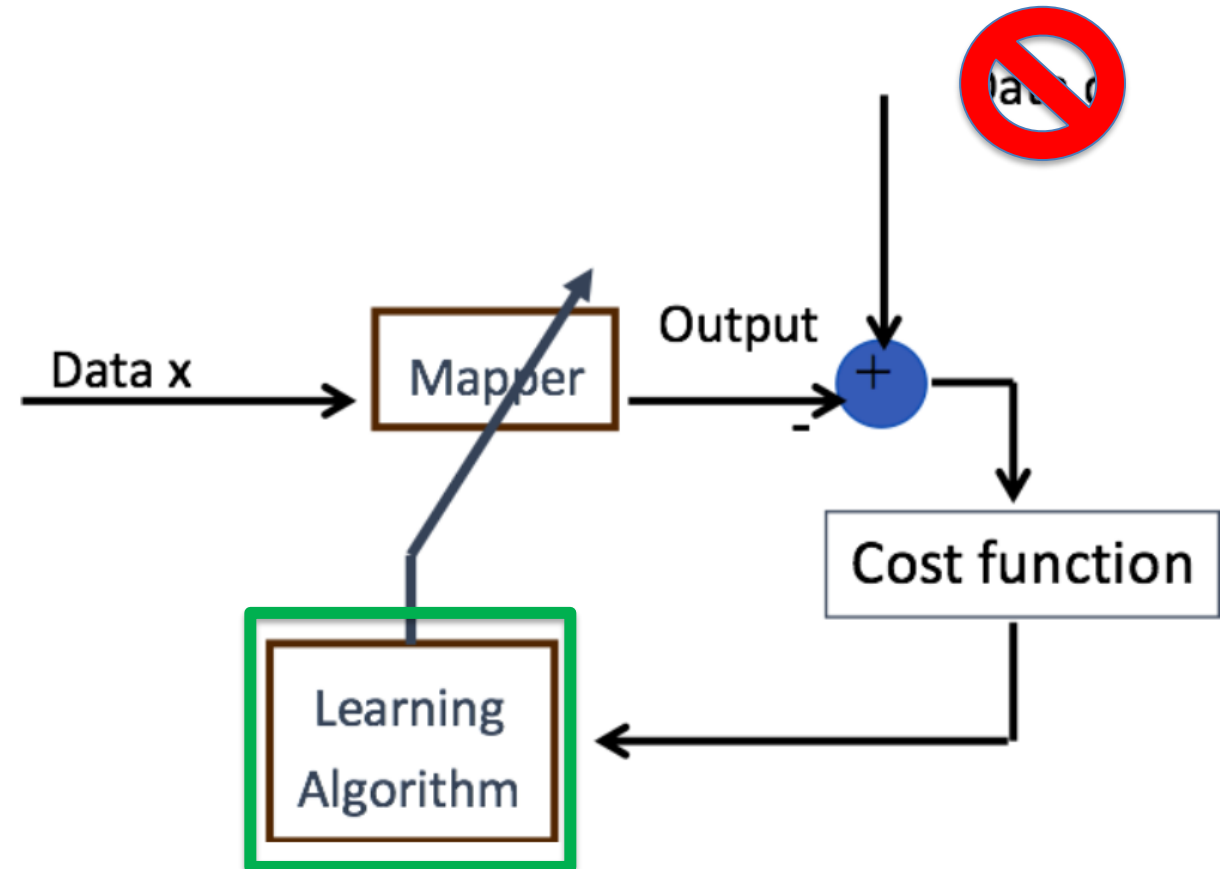


- Unsupervised: No labels,  $d$
- Mapper:
  - Density-based Clustering
  - Takes input data and groups into  $k$  clusters
- **Cost function:**
  - Kernel function

$$\frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right)$$



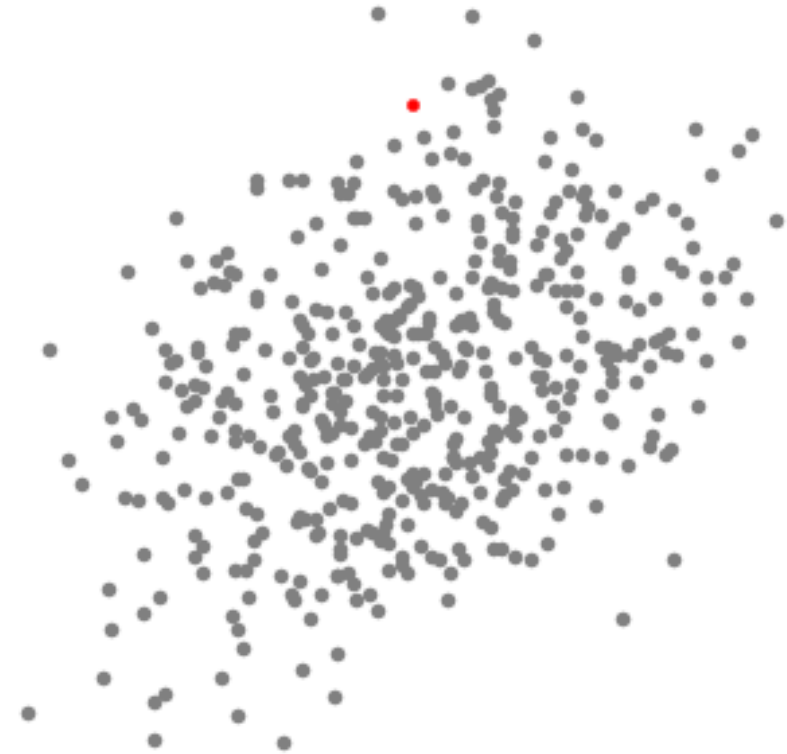
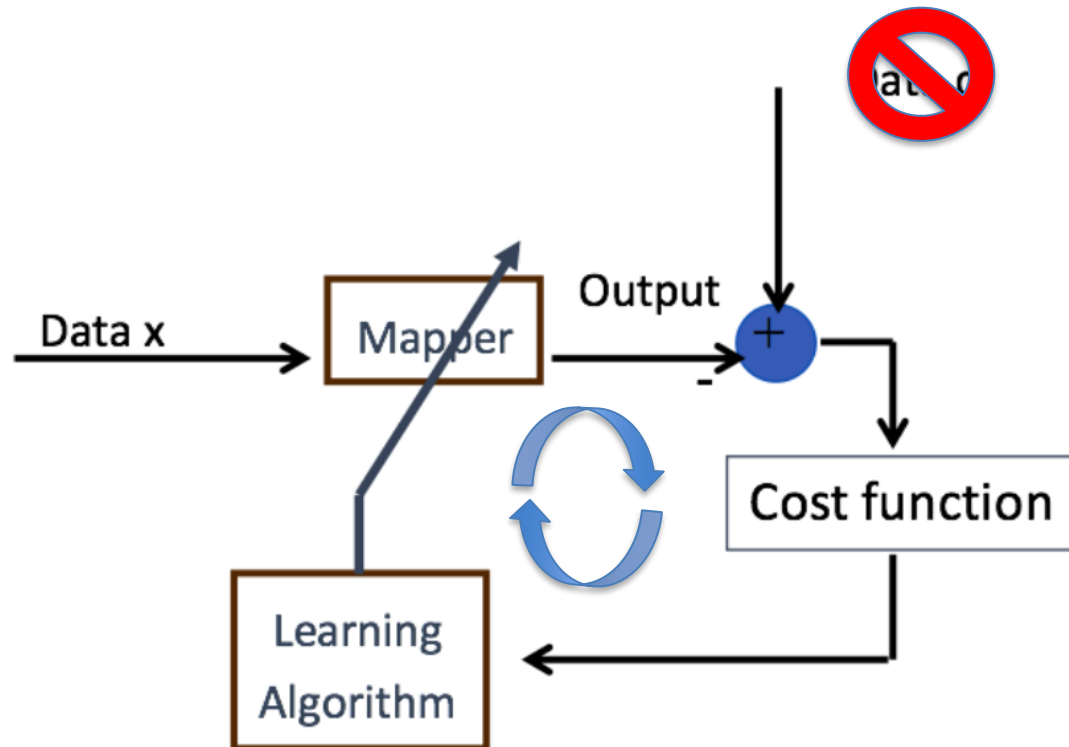
- Unsupervised: No labels,  $d$
- Mapper:
  - Density-based
  - Takes input data and groups into  $k$  clusters
- Cost function:
  - Kernel function
- **Learning algorithm**
  - Maximal density
    - DBSCAN: Density connected points
    - MeanShift: Modes



# Density-based Clustering Machine Learning Model



TEXAS A&M UNIVERSITY  
Engineering



- Bayesian and Nearest Neighbor Classification



INTEGRITY  
EXCELLENCE LEADERSHIP



TEXAS A&M UNIVERSITY  
Engineering

**Thank You! Questions?  
Joshua Peeples, Ph.D.**

**<https://www.joshpeeples.com/>**  
**[jpeeples@tamu.edu](mailto:jpeeples@tamu.edu)**





TEXAS A&M UNIVERSITY  
Engineering

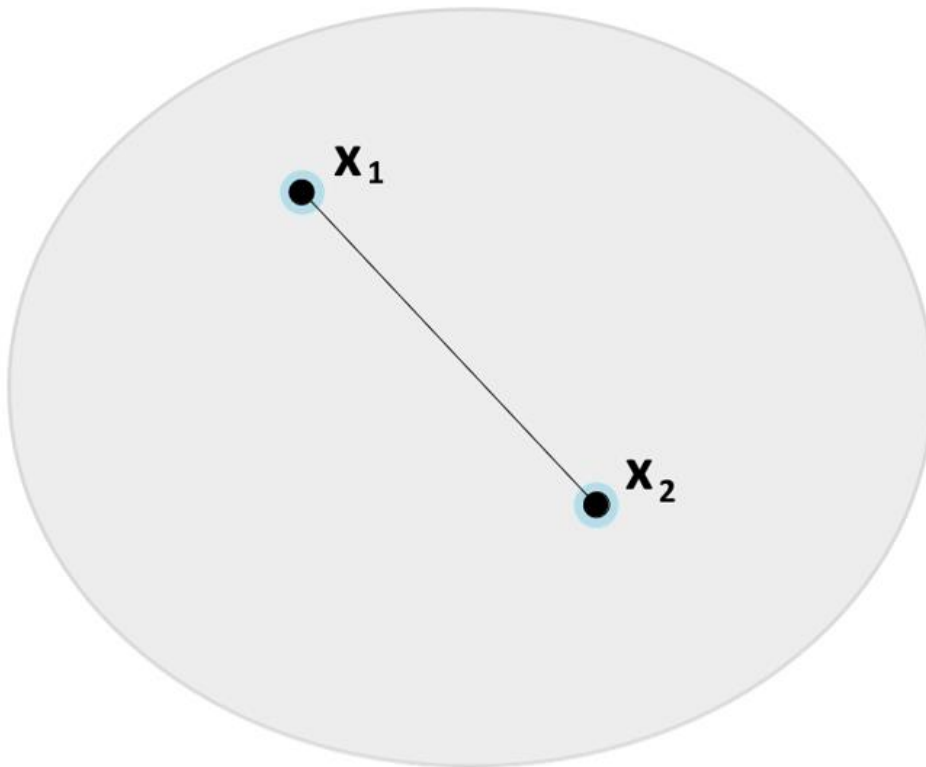
# Supplemental Slides

# Convex vs Non-Convex Clusters

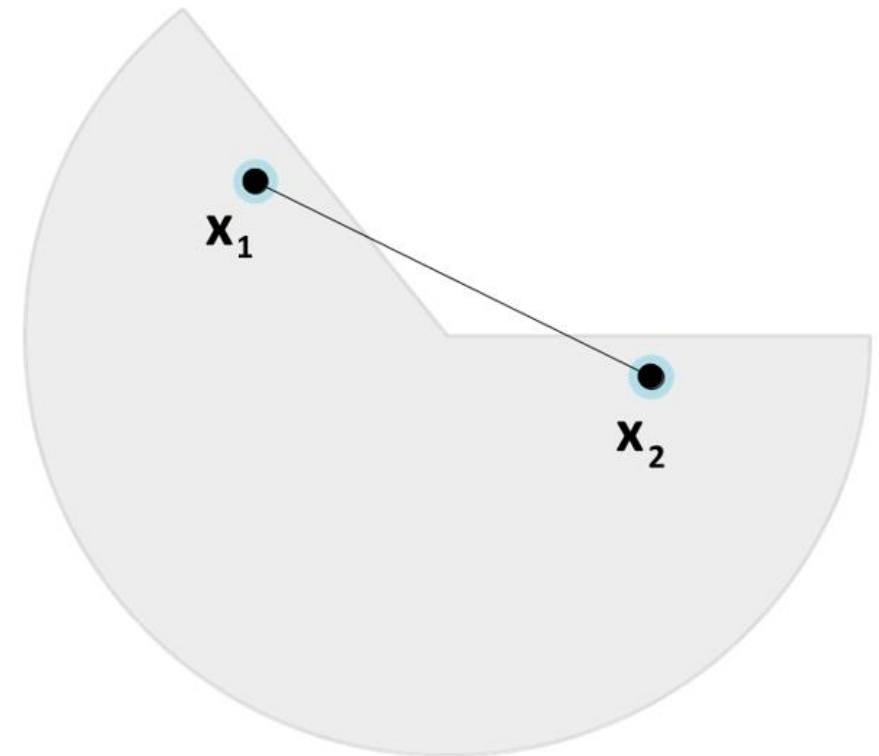


TEXAS A&M UNIVERSITY  
Engineering

Convex cluster



Non-Convex (or Concave) cluster



- [StatQuest: DBSCAN](#)
- [Mean Shift Clustering Summary](#)
- [Kernel Density Estimation](#)
- [HDBSCAN](#)
- [OPTICS](#)