

Nonparametric Estimation & Mean-Shift Clustering



Prithwijit Guha
Dept. of EEE, IIT Guwahati

Unsupervised Learning



Parametric
Clustering
Algorithms

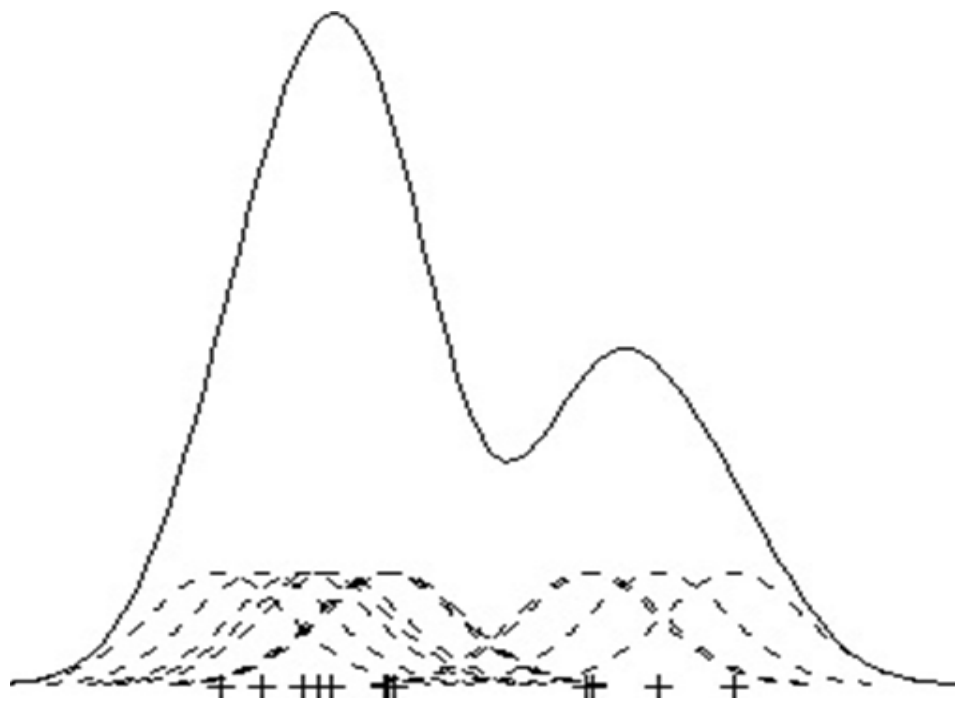
Generic
Clustering
Algorithms

**Estimation
Theory**

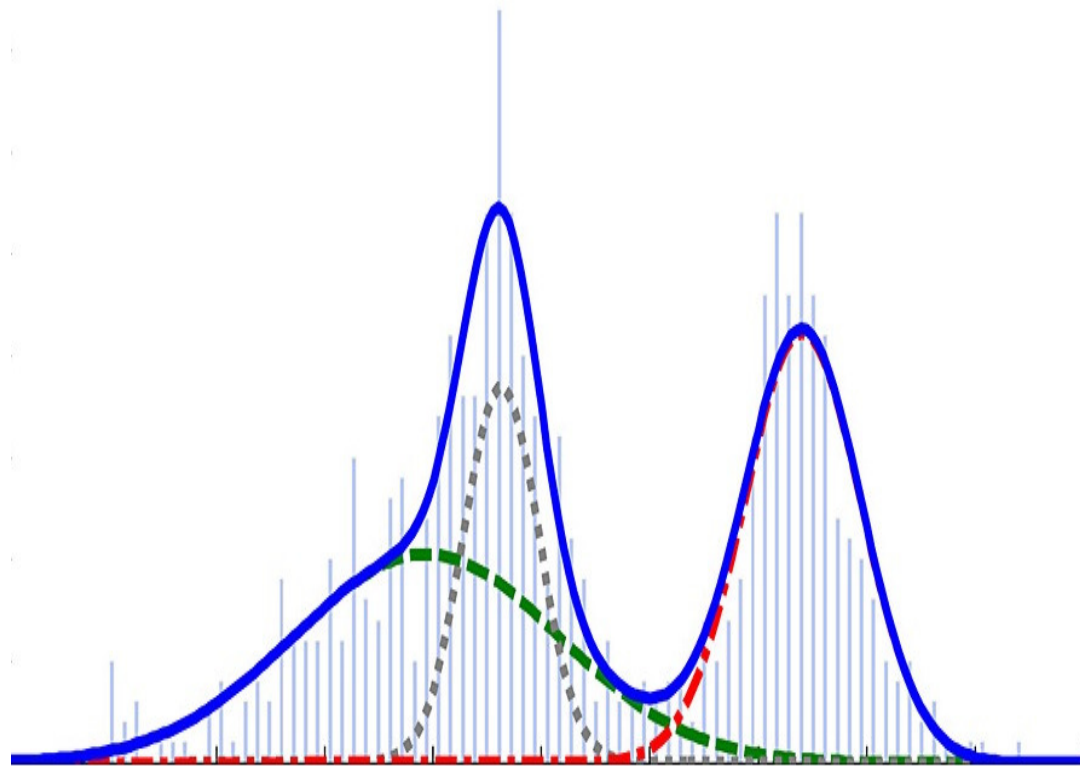
Generative
Models

Pattern
Mining

Estimation Theory



Non-Parametric Estimation



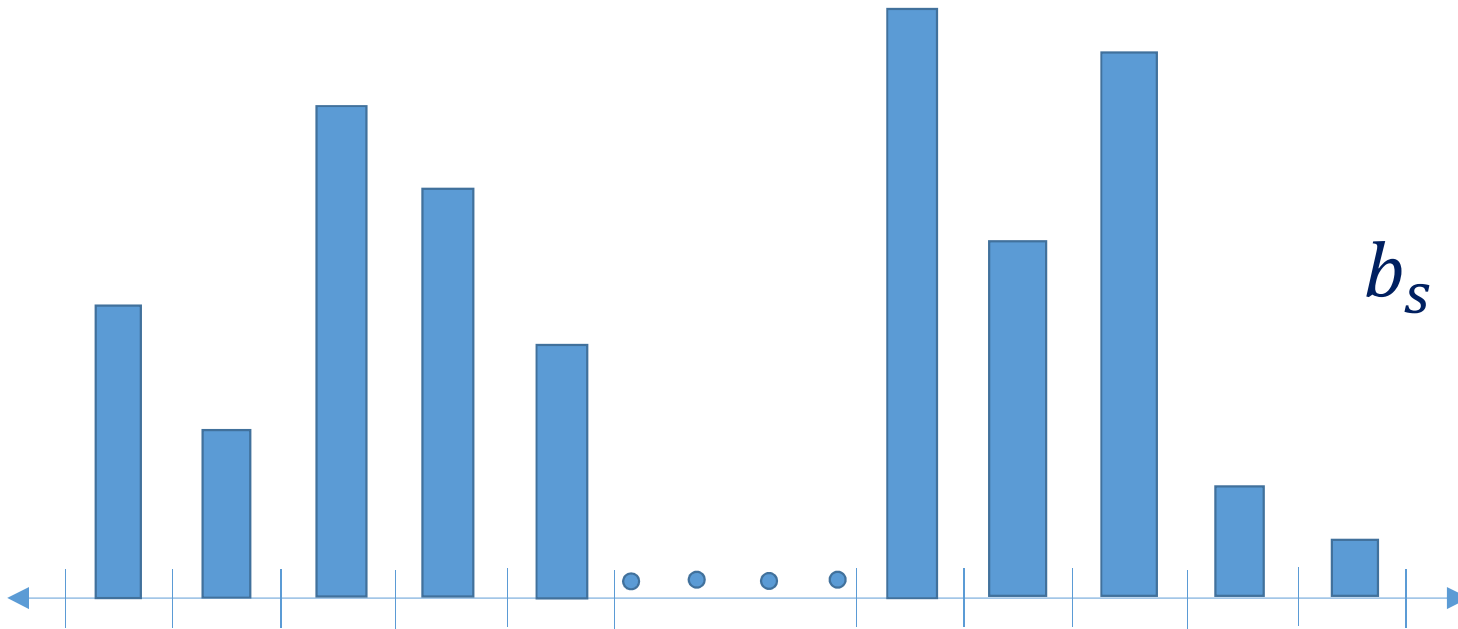
Parametric Estimation

Frequency Distribution

$$\mathbf{S} = \{x_1, \dots x_i, \dots x_n\}$$

$\min(\mathbf{S}) = x_{\min}$

$\max(\mathbf{S}) = x_{\max}$



$$b_s = \frac{x_{\max} - x_{\min}}{m}$$

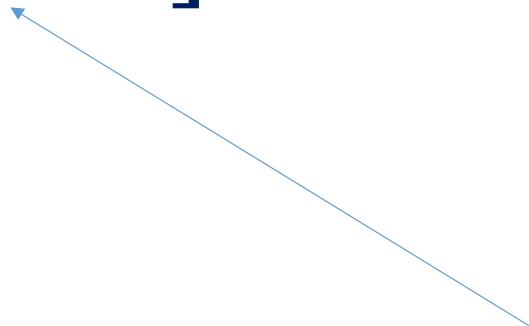
Frequency Distribution: Pseudo-code

- INPUT: $\mathbf{S} = \{x_1, \dots x_i, \dots x_n\}$
- Intervals: $x_{max} = \max(\mathbf{S}); x_{min} = \min(\mathbf{S})$
- Bin: $b_s = \frac{x_{max} - x_{min}}{m}$
- INITIALIZE: $\mathbf{H}[j] = 0; j = 1, \dots m$
- FOR $i = 1 \rightarrow n$
 1. Get Bin Index: $b_i = \left\lfloor \frac{x_i - x_{min}}{b_s} \right\rfloor + 1$
 2. UPDATE: $\mathbf{H}[b_i] = \mathbf{H}[b_i] + \frac{1}{n}$
- END FOR

Frequency Distribution: Mathematical Expression

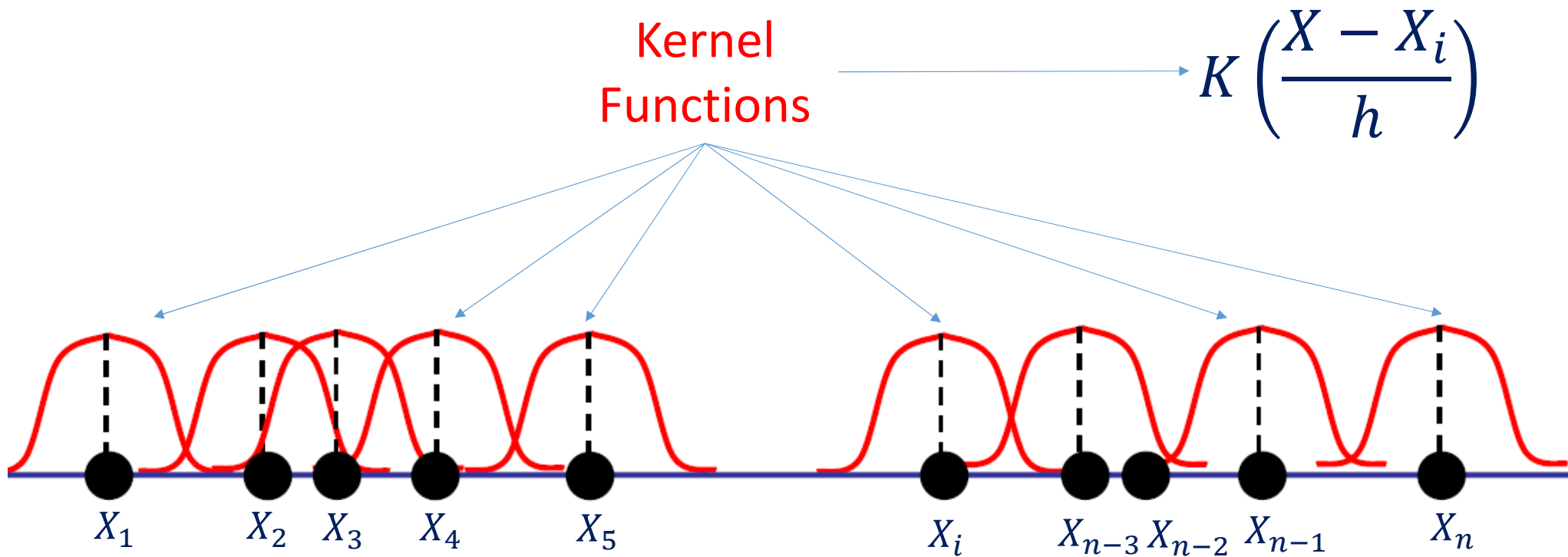
$$H[j] = \frac{1}{n} \sum_{i=1}^n \delta \left[\left\lfloor \frac{x_i - x_{min}}{b_s} \right\rfloor + 1 - j \right]$$


$$j = 1, \dots, m$$


$$b_s = \frac{x_{max} - x_{min}}{m}$$

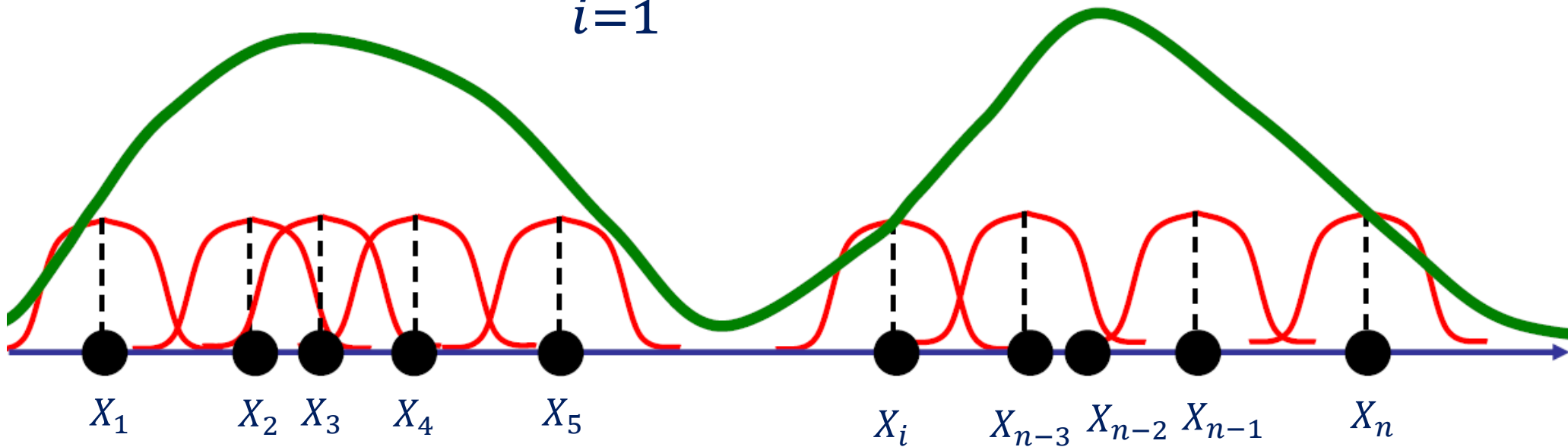
Distribution of Multivariate Data

$$\mathcal{S} = \{X_1, \dots, X_i, \dots, X_n\}$$

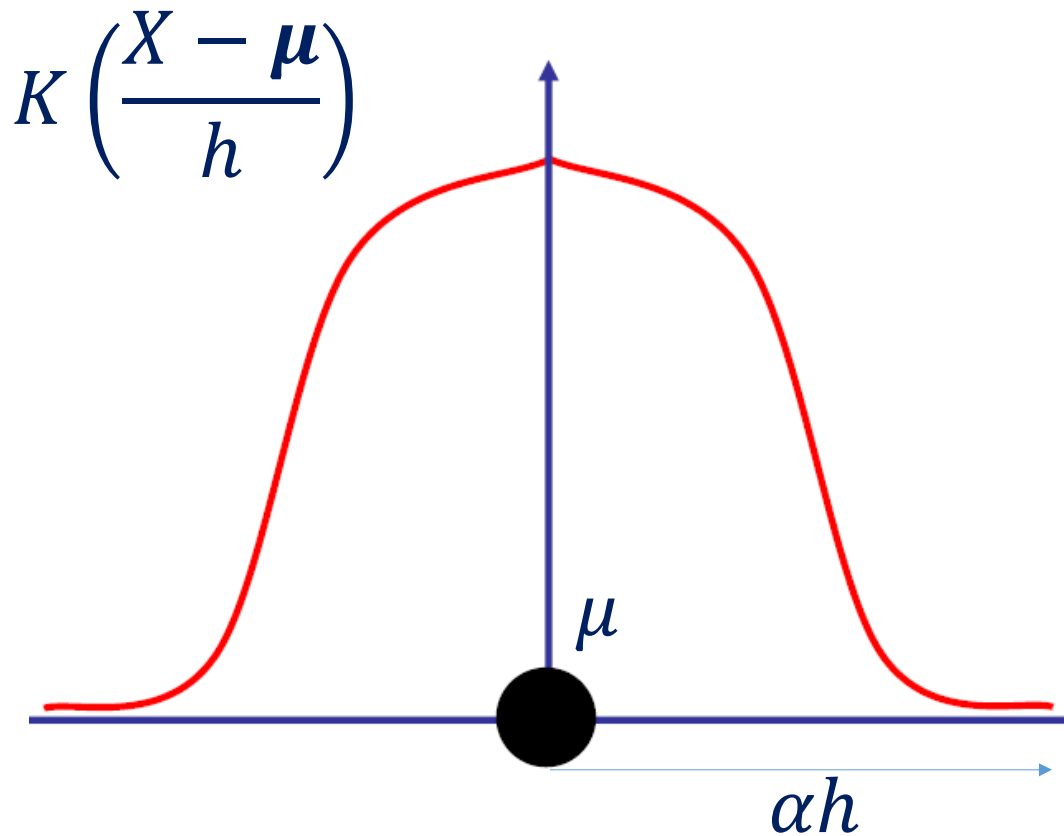


Distribution of Multivariate Data

$$P(X) = \frac{1}{n} \sum_{i=1}^n K \left(\frac{X - X_i}{h} \right)$$



Kernel Functions



Decays to ZERO As Moves
Away From μ

Rate of Decay Controlled by
Bandwidth h

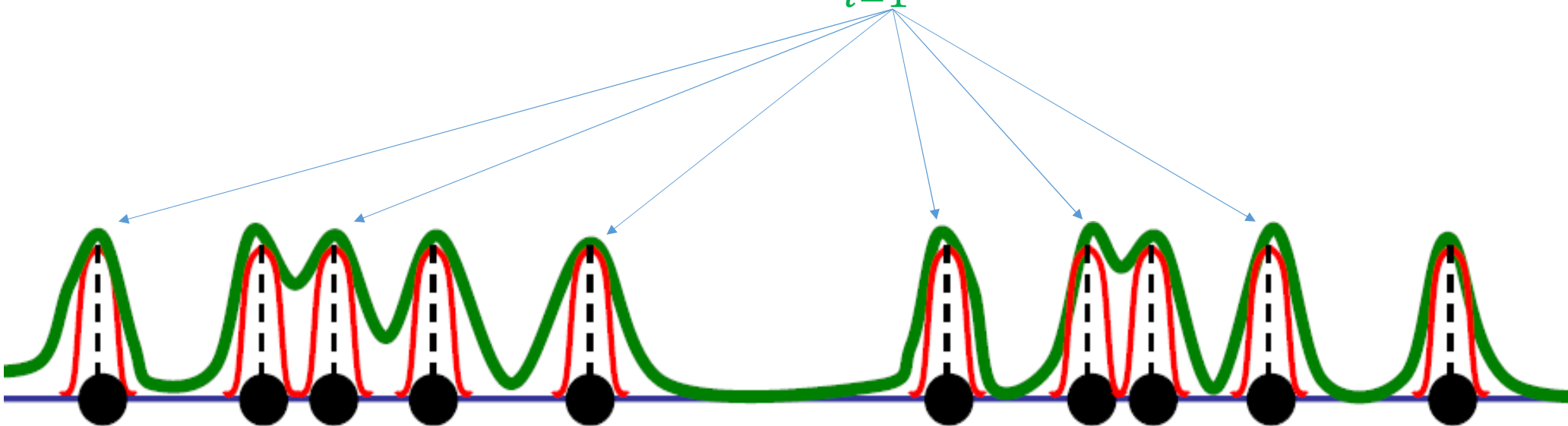
Maximum at μ

Symmetric around μ

$$\int_{\mathbb{R}^N} K\left(\frac{X - \mu}{h}\right) dX = 1$$

Effect of Low Bandwidth (h)

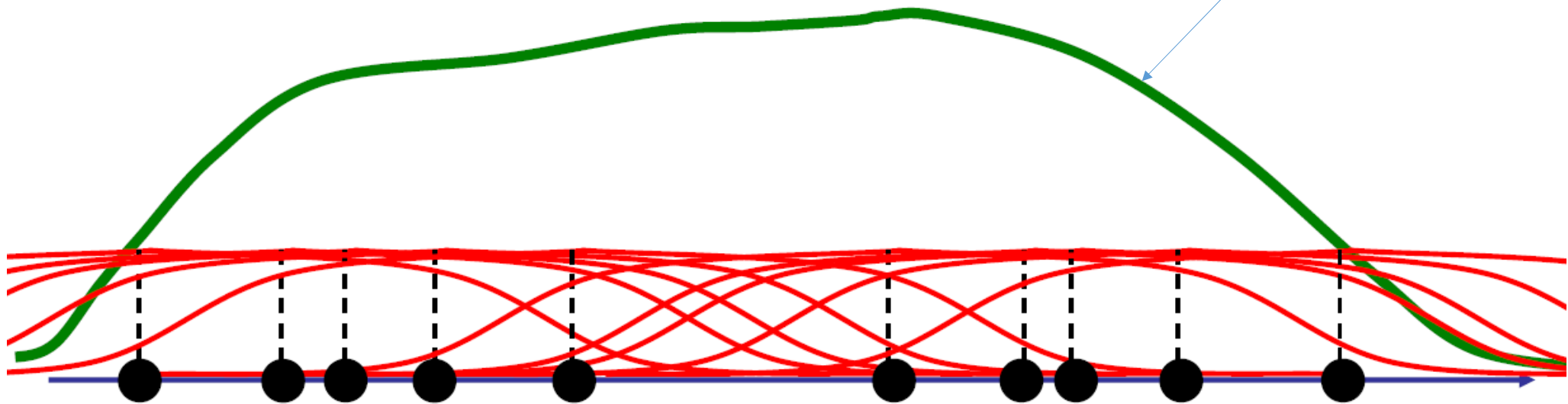
$$P(X) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)$$



Too Many Modes Lose the Interpolability

Effect of High Bandwidth (h)

$$P(X) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)$$



High Bandwidth Smoothens Out the Information

Kernel Profile Function

$$\mathbf{K} \left(\frac{\mathbf{X} - \boldsymbol{\mu}}{h} \right) = c \mathbf{k} \left(\left\| \frac{\mathbf{X} - \boldsymbol{\mu}}{h} \right\|^2 \right)$$

Diagram illustrating the relationship between the Kernel Function, Kernel Profile Function, and Normalization Constant:

- Kernel Function** (labeled \mathbf{K})
- Kernel Profile Function** (labeled $\mathbf{k}(x)$)
- Normalization Constant** (labeled c)

The equation shows that the Kernel Function \mathbf{K} is equal to the product of the Normalization Constant c and the Kernel Profile Function \mathbf{k} applied to the squared norm of the normalized difference between \mathbf{X} and $\boldsymbol{\mu}$.

Negative Derivative of Kernel Profile

$$\boldsymbol{g}(\boldsymbol{x}) = -\frac{d}{d\boldsymbol{x}}\{\boldsymbol{k}(\boldsymbol{x})\}$$

$$\nabla_X \boldsymbol{K}\left(\frac{X - \mu}{h}\right) = c \nabla_X \boldsymbol{k}\left(\left\|\frac{X - \mu}{h}\right\|^2\right)$$

Negative Derivative of Kernel Profile

$$\nabla_X \mathbf{K} \left(\frac{X - \mu}{h} \right) = c \nabla_X \mathbf{k} \left(\left\| \frac{X - \mu}{h} \right\|^2 \right)$$

$$= c \mathbf{k}' \left(\left\| \frac{X - \mu}{h} \right\|^2 \right) 2 \left(\frac{X - \mu}{h^2} \right)$$

$$= \frac{2c}{h^2} (\mu - X) \left\{ -\mathbf{k}' \left(\left\| \frac{X - \mu}{h} \right\|^2 \right) \right\}$$

Negative Derivative of Kernel Profile

$$\nabla_{\mathbf{X}} \mathbf{K} \left(\frac{X - \mu}{h} \right) = \frac{2c}{h^2} (\mu - X) \left\{ -\mathbf{k}' \left(\left\| \frac{X - \mu}{h} \right\|^2 \right) \right\}$$

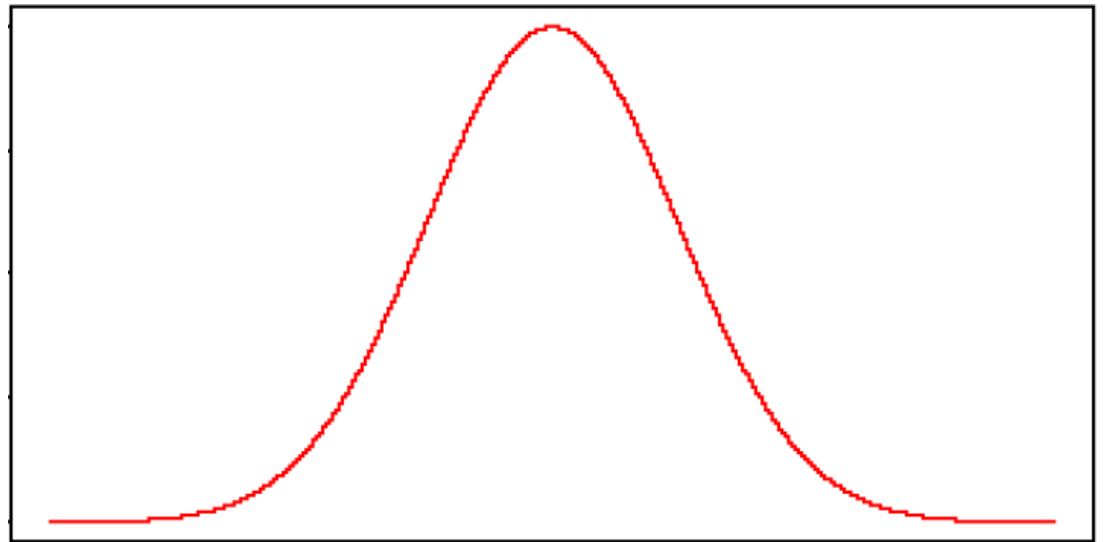
$$\nabla_{\mathbf{X}} \mathbf{K} \left(\frac{X - \mu}{h} \right) = \frac{2c}{h^2} (\mu - X) g \left(\left\| \frac{X - \mu}{h} \right\|^2 \right)$$

Gaussian Kernel Function

$$\mathbf{K}_N\left(\frac{X - \mu}{h}\right) = c_N \exp\left(-\frac{1}{2} \left\|\frac{X - \mu}{h}\right\|^2\right)$$

$$\mathbf{k}_N(x) = \exp\left(-\frac{x}{2}\right)$$

$$\mathbf{g}_N(x) = \frac{1}{2} \exp\left(-\frac{x}{2}\right)$$

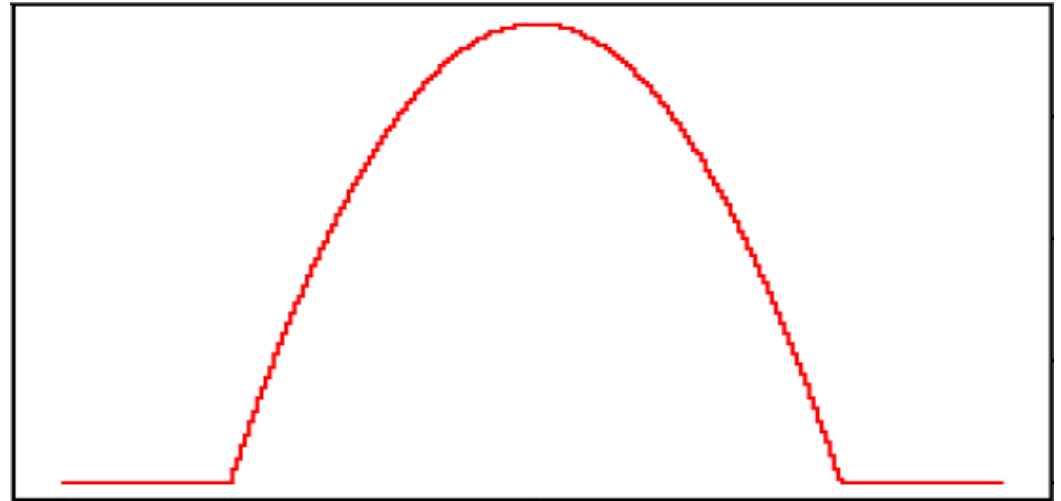


Epanechnikov Kernel Function

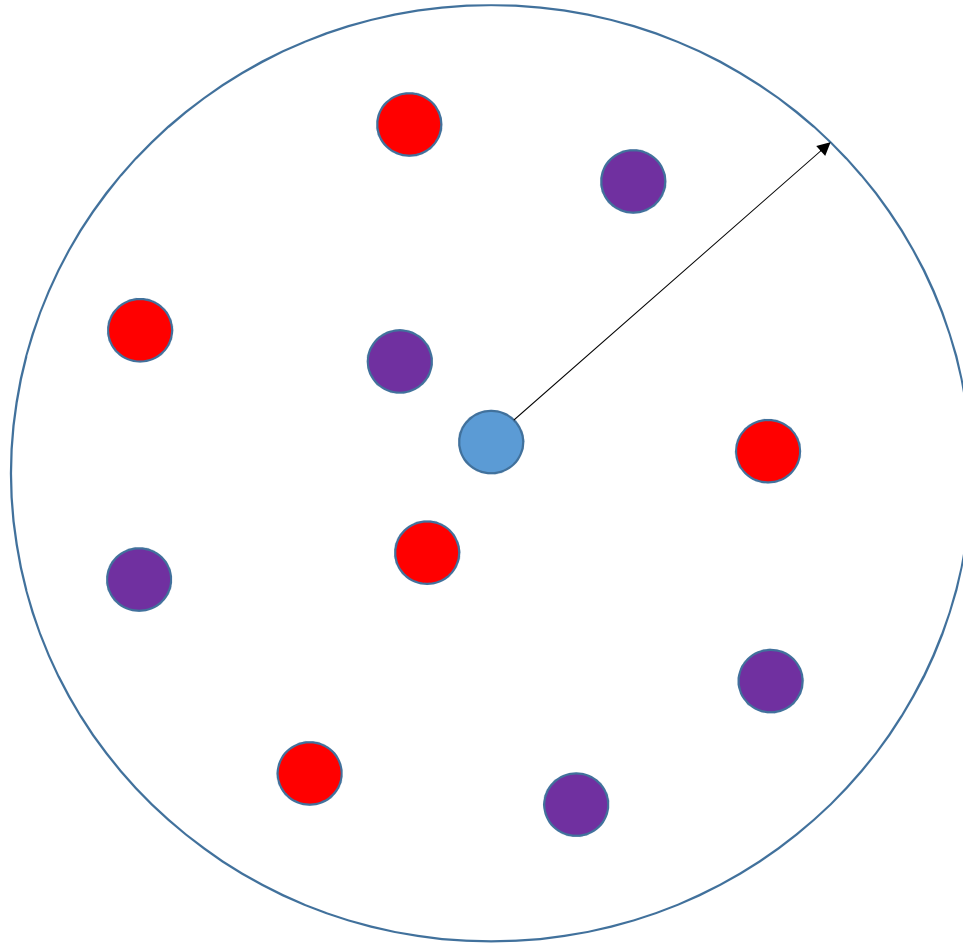
$$\mathbf{K}_E\left(\frac{X - \mu}{h}\right) = \begin{cases} c_E \left(1 - \left\|\frac{X - \mu}{h}\right\|^2\right), & \|X - \mu\| \leq h \\ 0, & \|X - \mu\| > h \end{cases}$$

$$\mathbf{k}_E(x) = \begin{cases} 1 - x^2, & |x| \leq 1 \\ 0, & |x| > 1 \end{cases}$$

$$\mathbf{g}_E(x) = \begin{cases} 1, & |x| \leq 1 \\ 0, & |x| > 1 \end{cases}$$



Classification & Regression in the Neighborhood



Application: Regression

$$\mathcal{S} = \{(X_i, Y_i); i = 1, \dots, n\}$$

$$\mathbf{X} \in \mathbb{R}^N$$

$$\mathbf{Y} \in \mathbb{R}^M$$

$$Y = \frac{\sum_{i=1}^n Y_i K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

Application: Classification

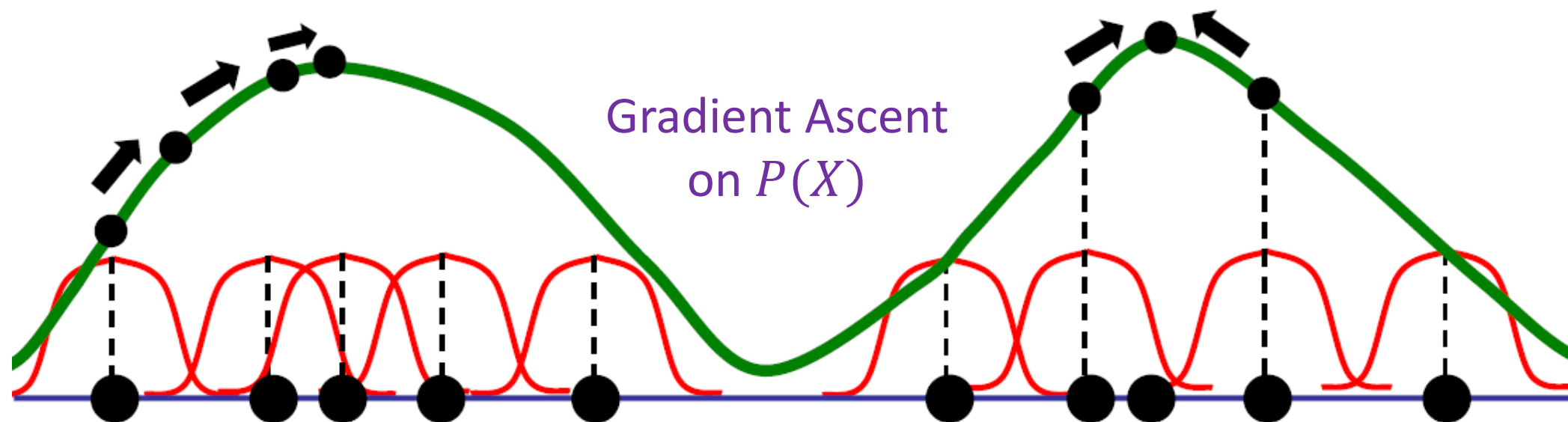
$$\mathcal{S} = \{(X_i, Y_i); i = 1, \dots, n\} \quad \begin{array}{l} X \in \mathbb{R}^N \\ Y \in \{0,1\}^M \end{array}$$

$$V = \frac{\sum_{i=1}^n Y_i K\left(\frac{X - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{X - X_i}{h}\right)}$$

$$Y[j] = \frac{V[j]}{\sum_{r=1}^m V[r]}$$

Seeking the Modes of $P(X)$

$$X^{(t+1)} = X^{(t)} + \eta_t \nabla_X P(X^{(t)})$$



Evaluating Gradient of $P(X)$

$$P(X) = \frac{1}{n} \sum_{i=1}^n \mathbf{K} \left(\frac{X - X_i}{h} \right) = \frac{1}{n} \sum_{i=1}^n c \mathbf{k} \left(\left\| \frac{X - X_i}{h} \right\|^2 \right)$$

$$\nabla_X P(X) = \frac{1}{n} \sum_{i=1}^n c \nabla_X \mathbf{k} \left(\left\| \frac{X - X_i}{h} \right\|^2 \right)$$

$$\nabla_X P(X) = \frac{1}{n} \sum_{i=1}^n c \mathbf{k}' \left(\left\| \frac{X - X_i}{h} \right\|^2 \right) 2 \left(\frac{X - X_i}{h^2} \right)$$

Evaluating Gradient of $P(X)$

$$\nabla_X P(X) = \frac{1}{n} \sum_{i=1}^n c \mathbf{k}' \left(\left\| \frac{X - X_i}{h} \right\|^2 \right) 2 \left(\frac{X - X_i}{h^2} \right)$$

$$\nabla_X P(X) = \frac{2c}{nh^2} \sum_{i=1}^n (X_i - X) \left\{ -\mathbf{k}' \left(\left\| \frac{X - X_i}{h} \right\|^2 \right) \right\}$$

$$\nabla_X P(X) = \frac{2c}{nh^2} \sum_{i=1}^n (X_i - X) \mathbf{g} \left(\left\| \frac{X - X_i}{h} \right\|^2 \right)$$

Seeking Modes of $P(X)$: Gradient Ascent

$$X^{(t+1)} = X^{(t)} + \eta_t \nabla_X P(X^{(t)})$$

$$X^{(t+1)} - X^{(t)} = \eta_t \left\{ \frac{2c}{nh^2} \sum_{i=1}^n (X_i - X^{(t)}) \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\}$$

$$= \frac{2c\eta_t}{nh^2} \left\{ \sum_{i=1}^n X_i \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) - \sum_{i=1}^n X^{(t)} \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\}$$

Seeking Modes of $P(X)$: Mean Shift Vector

$$X^{(t+1)} - X^{(t)} = \frac{2c\eta_t}{nh^2} \left\{ \sum_{i=1}^n X_i \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) - \sum_{i=1}^n X^{(t)} \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\}$$

$$X^{(t+1)} - X^{(t)} = \underbrace{\frac{2c\eta_t}{nh^2} \left\{ \sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\}}_{\text{scalar}} \left\{ \frac{\sum_{i=1}^n X_i \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)}{\sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)} - X^{(t)} \right\}$$

Mean Shift Vector

Seeking Modes of $P(X)$: Mean Shift Iteration

$$X^{(t+1)} - X^{(t)} = \frac{2c\eta_t}{nh^2} \left\{ \sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\} \left\{ \frac{\sum_{i=1}^n X_i \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)}{\sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)} - X^{(t)} \right\}$$

$$\eta_t = \frac{nh^2}{2c} \left\{ \sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\}^{-1} \Rightarrow \frac{2c\eta_t}{nh^2} \left\{ \sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right) \right\} = 1$$

$$X^{(t+1)} - X^{(t)} = \frac{\sum_{i=1}^n X_i \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)}{\sum_{i=1}^n \mathbf{g} \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)} - X^{(t)}$$

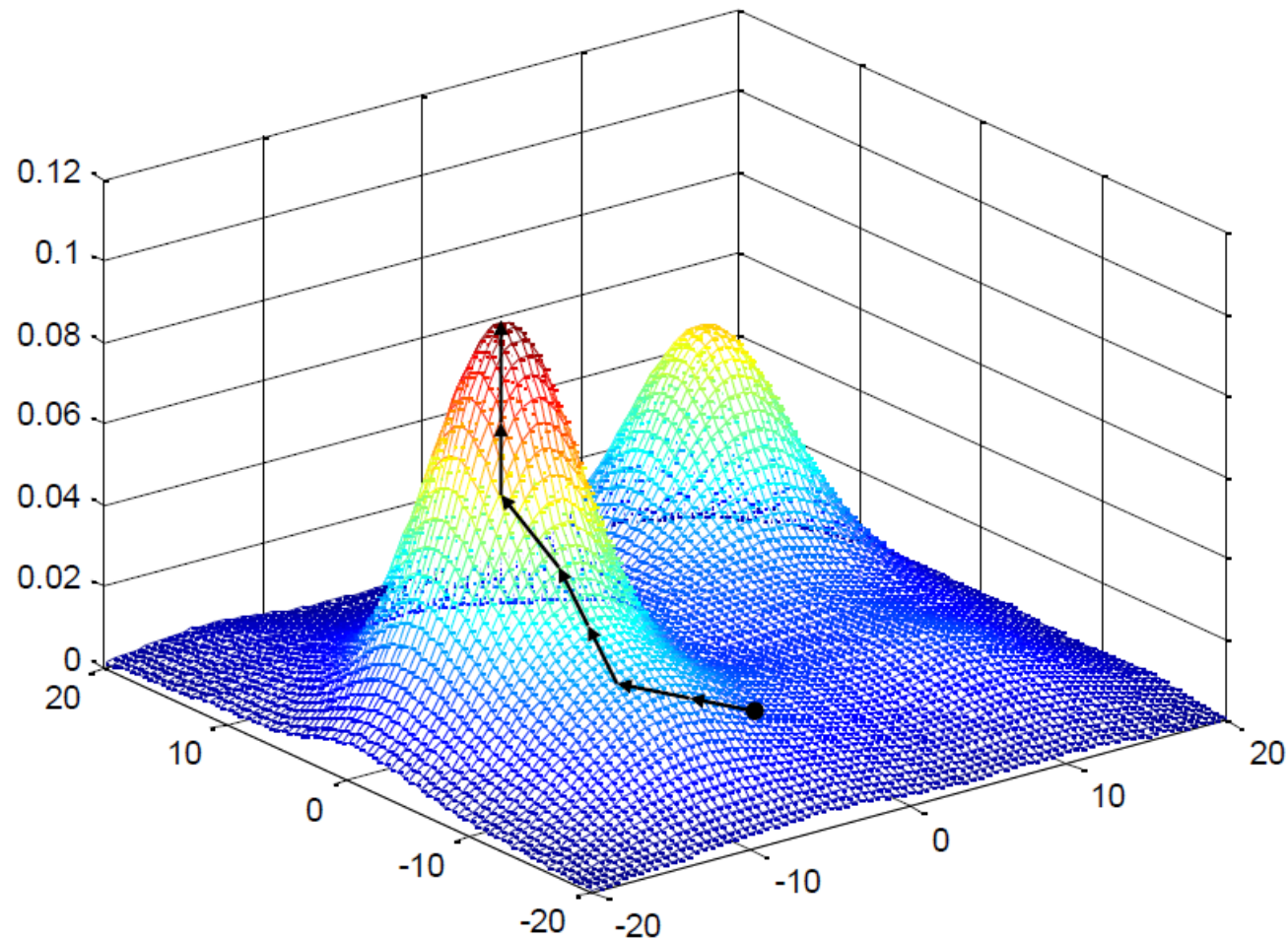
Seeking Modes of $P(X)$: Mean Shift Iteration

$$X^{(t+1)} = \frac{\sum_{i=1}^n X_i g\left(\left\|\frac{X^{(t)} - X_i}{h}\right\|^2\right)}{\sum_{i=1}^n g\left(\left\|\frac{X^{(t)} - X_i}{h}\right\|^2\right)}$$

Input: $\mathcal{S} = \{X_i; i = 1, \dots, n\}$

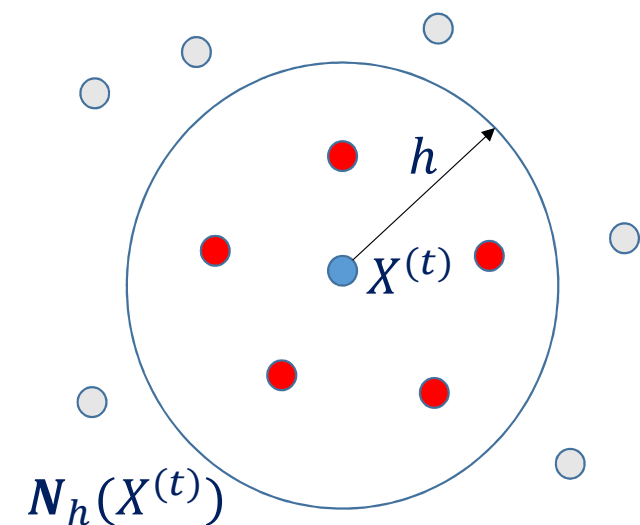
$X^{(0)} \in \mathcal{S}$

Seeking Modes of $P(X)$: Visualization

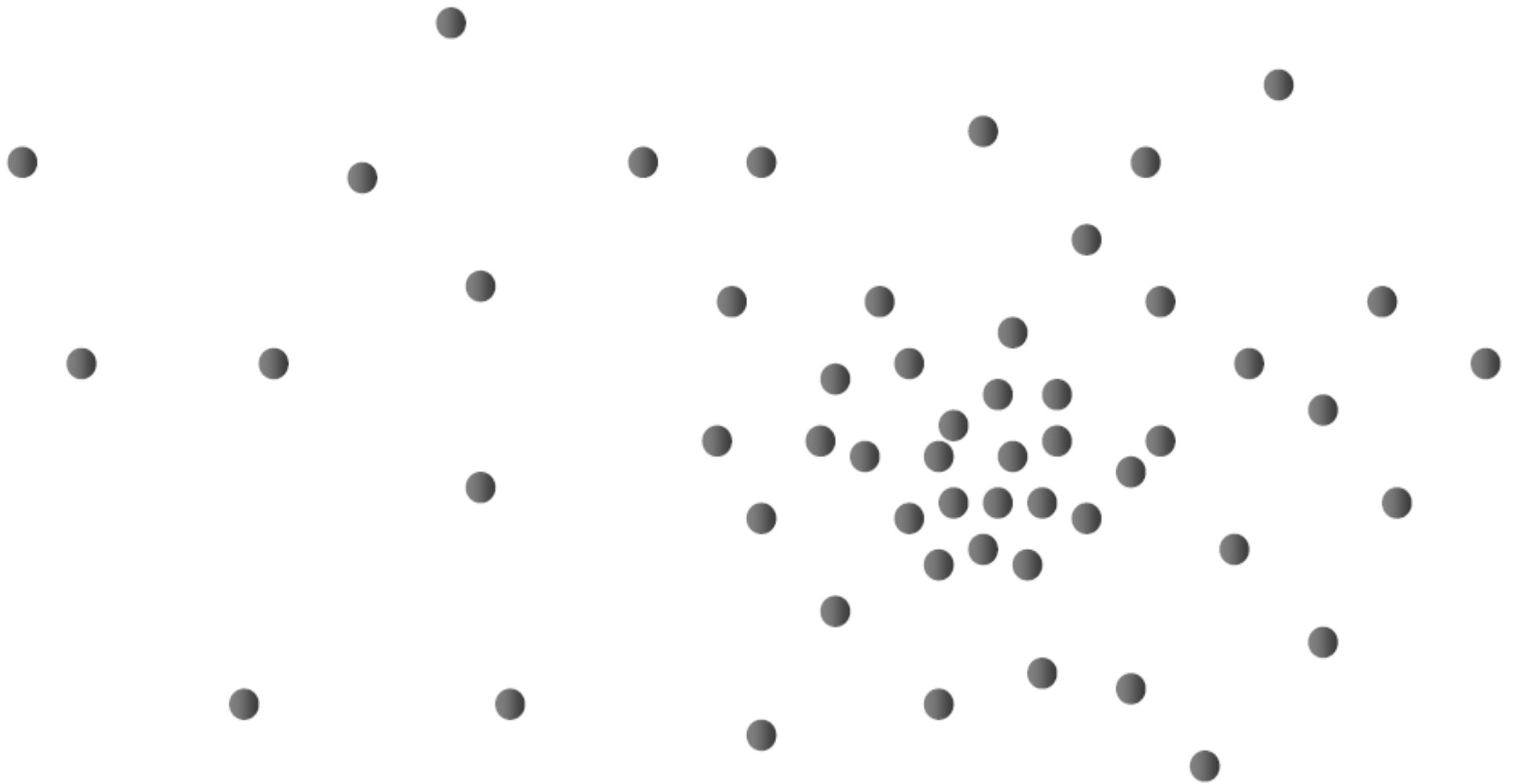


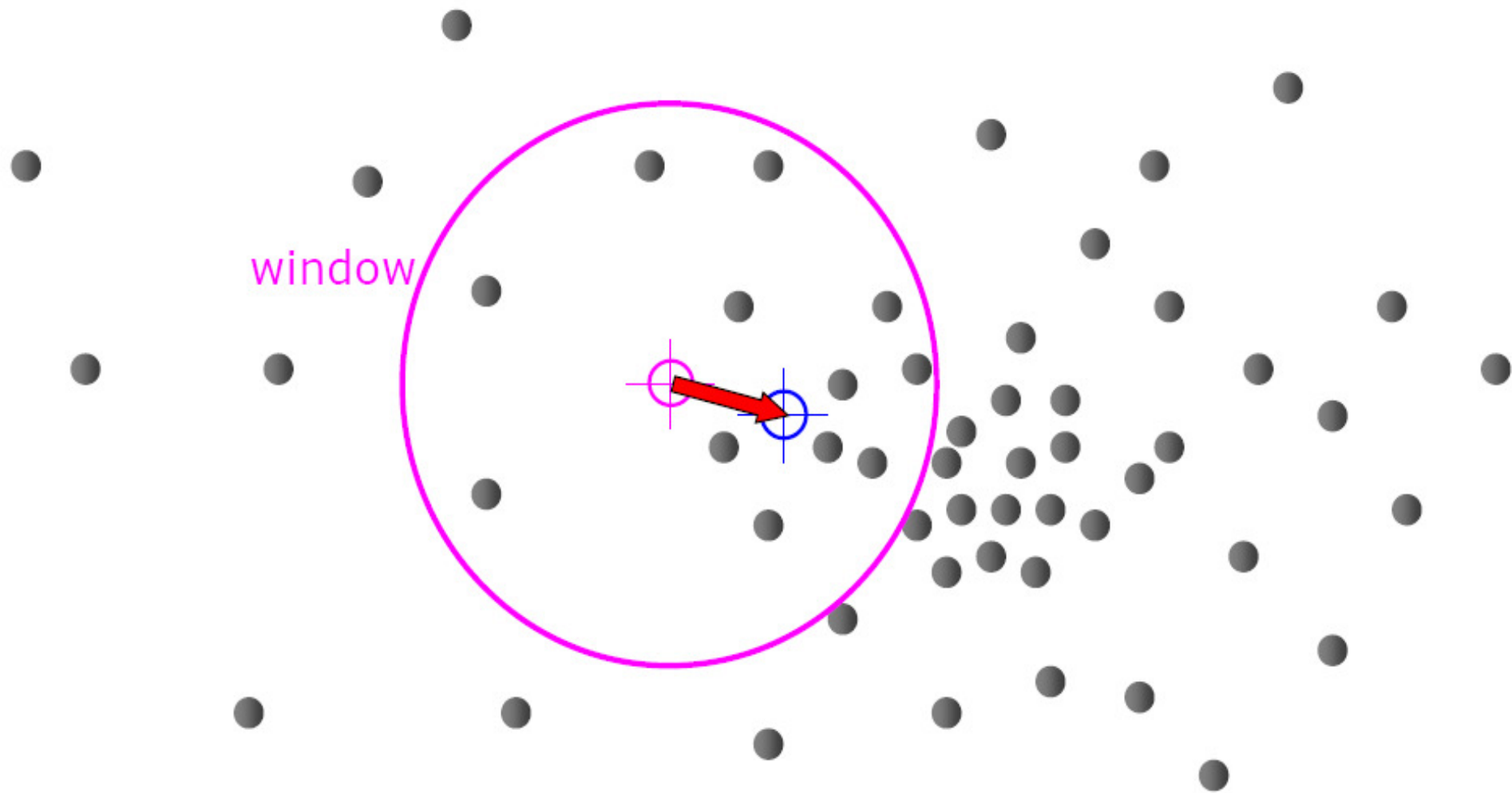
Seeking Modes of $P(X)$: Choice of Kernel

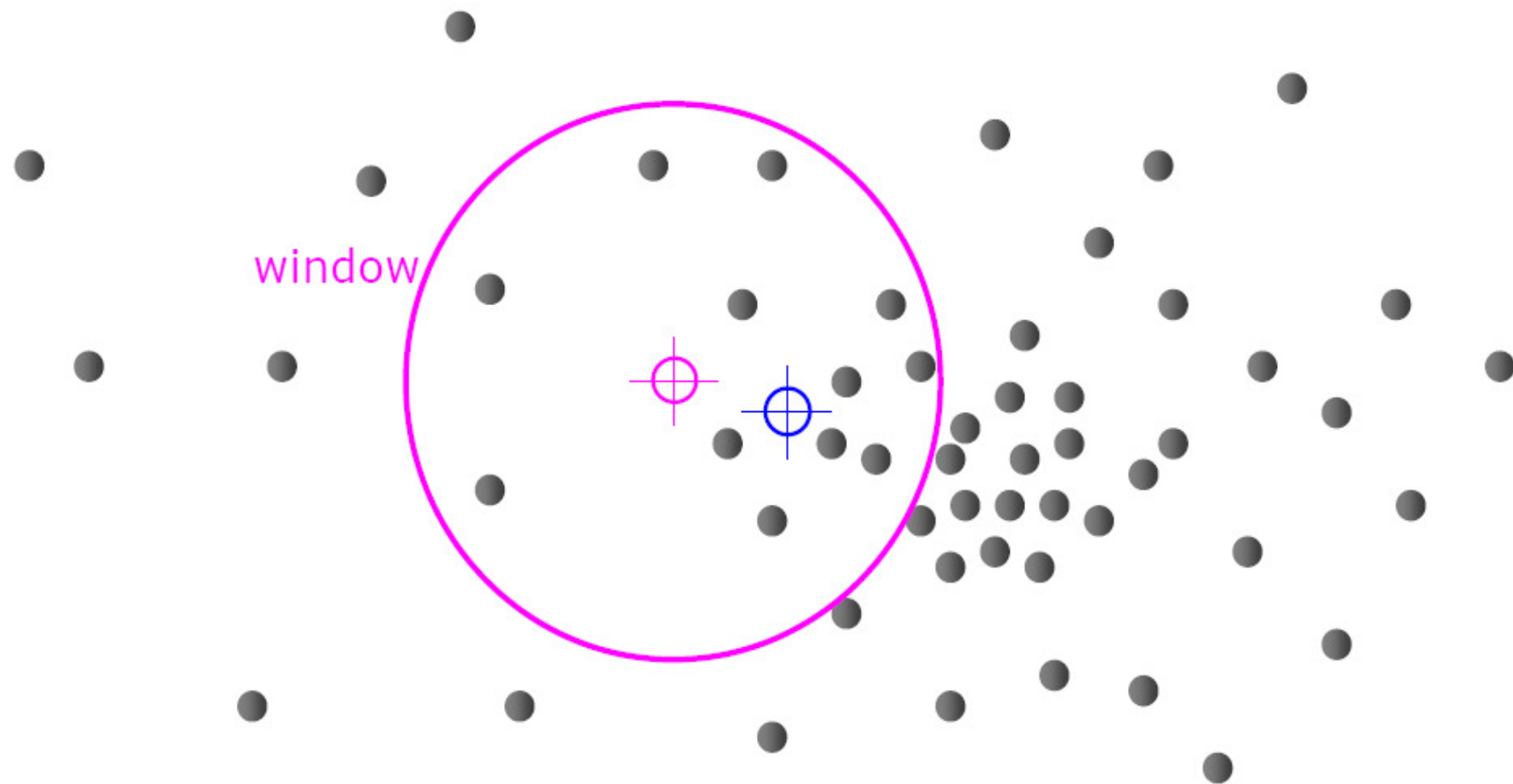
$$X^{(t+1)} = \frac{\sum_{i=1}^n X_i \mathbf{g}_E \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)}{\sum_{i=1}^n \mathbf{g}_E \left(\left\| \frac{X^{(t)} - X_i}{h} \right\|^2 \right)}$$

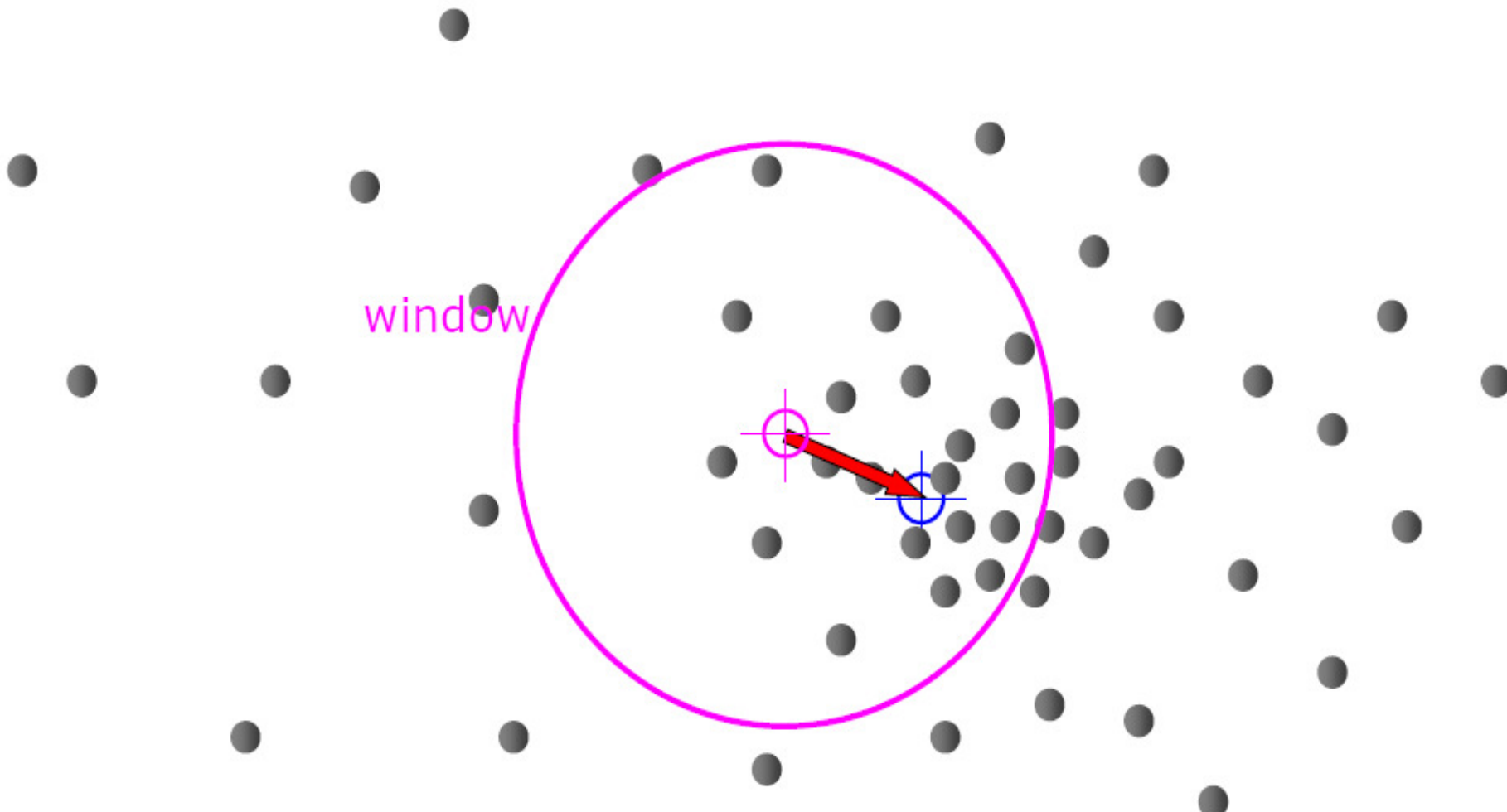


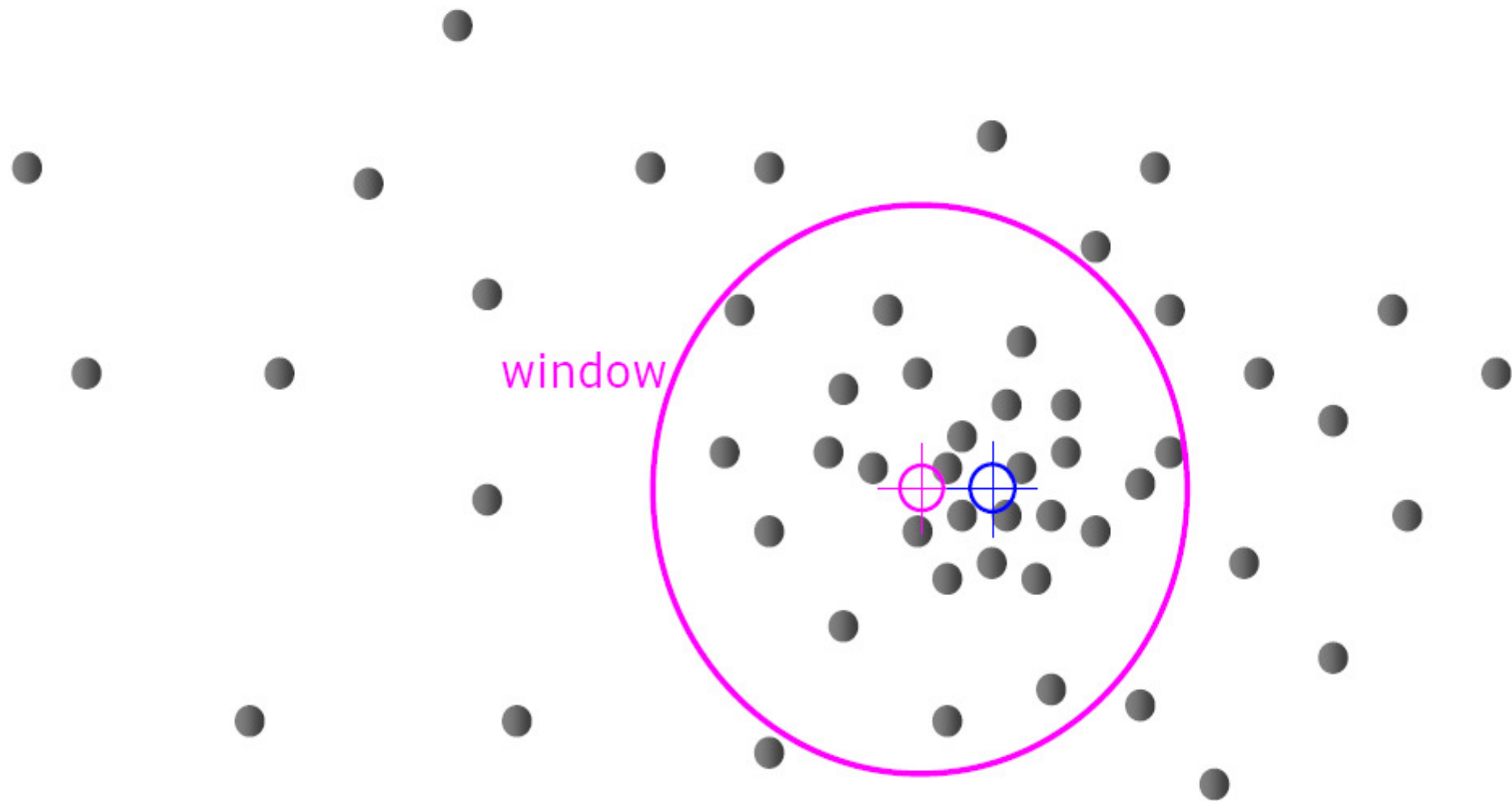
$$X^{(t+1)} = \frac{\sum_{X \in N_h(X^{(t)})} X \cdot 1}{\sum_{X \in N_h(X^{(t)})} 1}$$

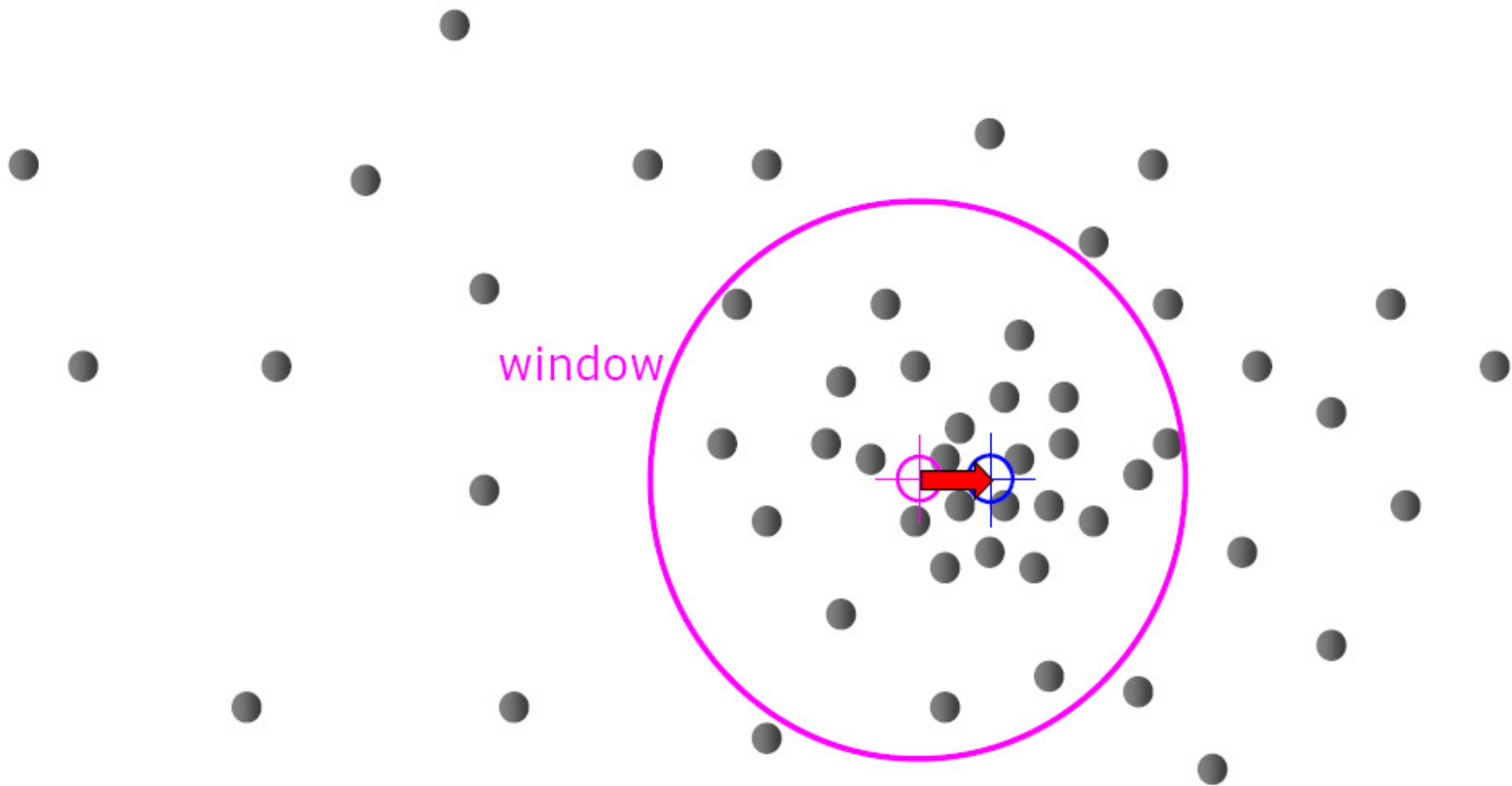


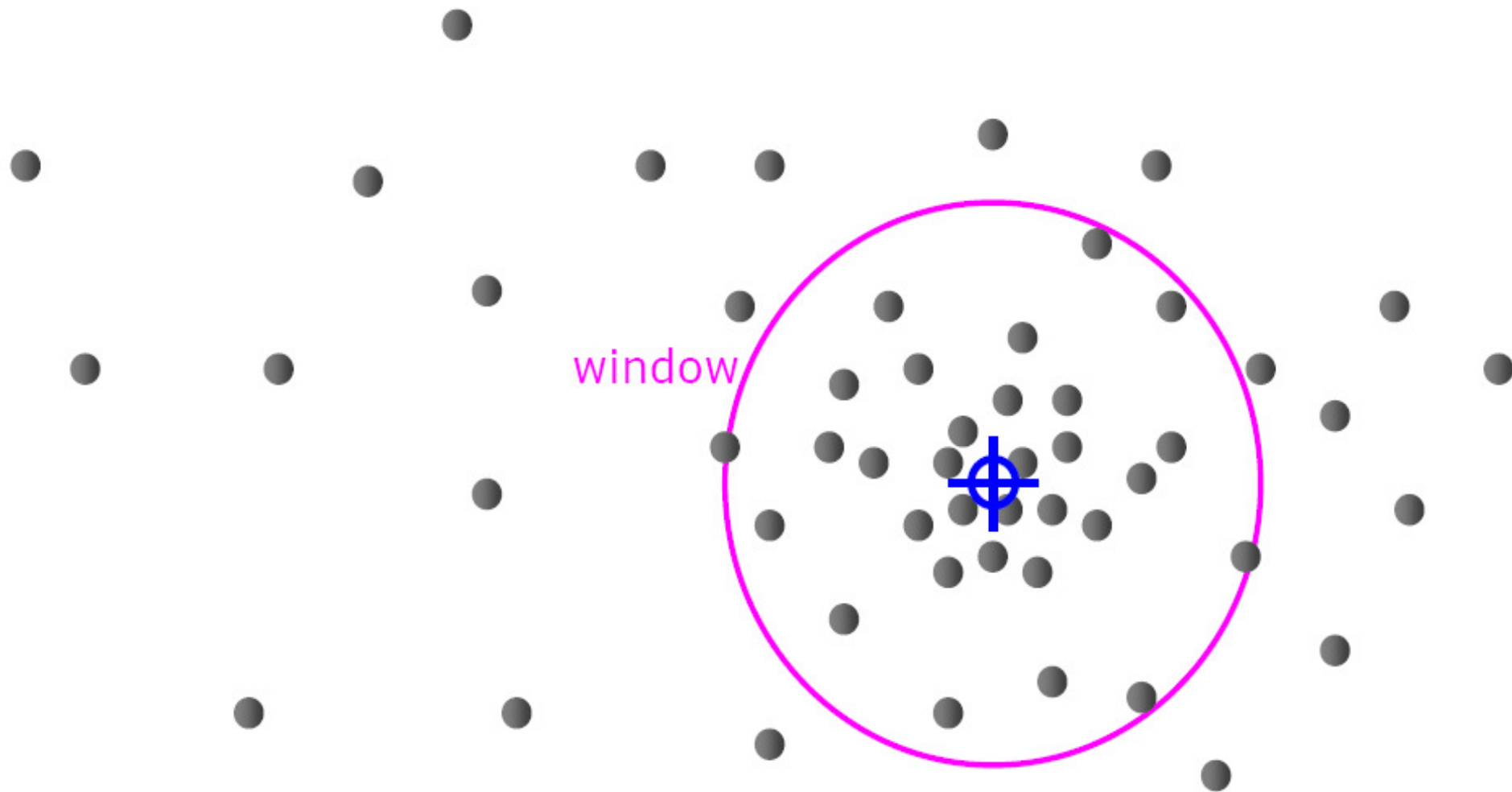












Mean-Shift Clustering: Problem

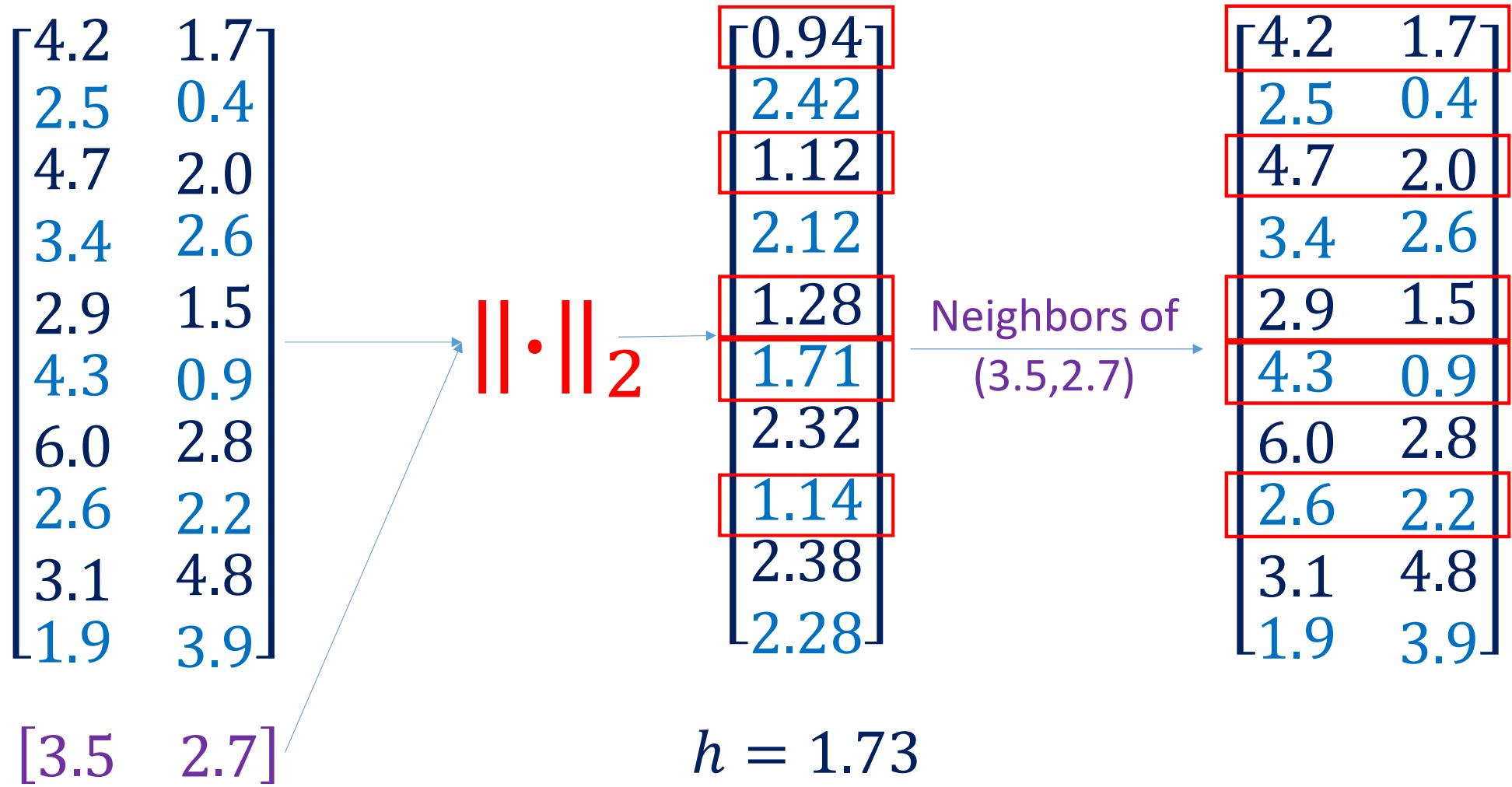
X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
4.2	2.5	4.7	3.4	2.9	4.3	6.0	2.6	3.1	1.9
1.7	0.4	2.0	4.6	1.5	0.9	2.8	2.2	4.8	3.9

$$X^{(0)} = (3.7, 2.5)$$

$$X^{(1)} = (\quad , \quad)$$

$$h = 1.73$$

Mean-Shift Clustering: Solution



Mean-Shift Clustering: Solution

$$X^{(1)} = \frac{\sum_{X \in N_{1.73}(3.5, 2.7)} X \cdot 1}{\sum_{X \in N_{1.73}(3.5, 2.7)} 1}$$

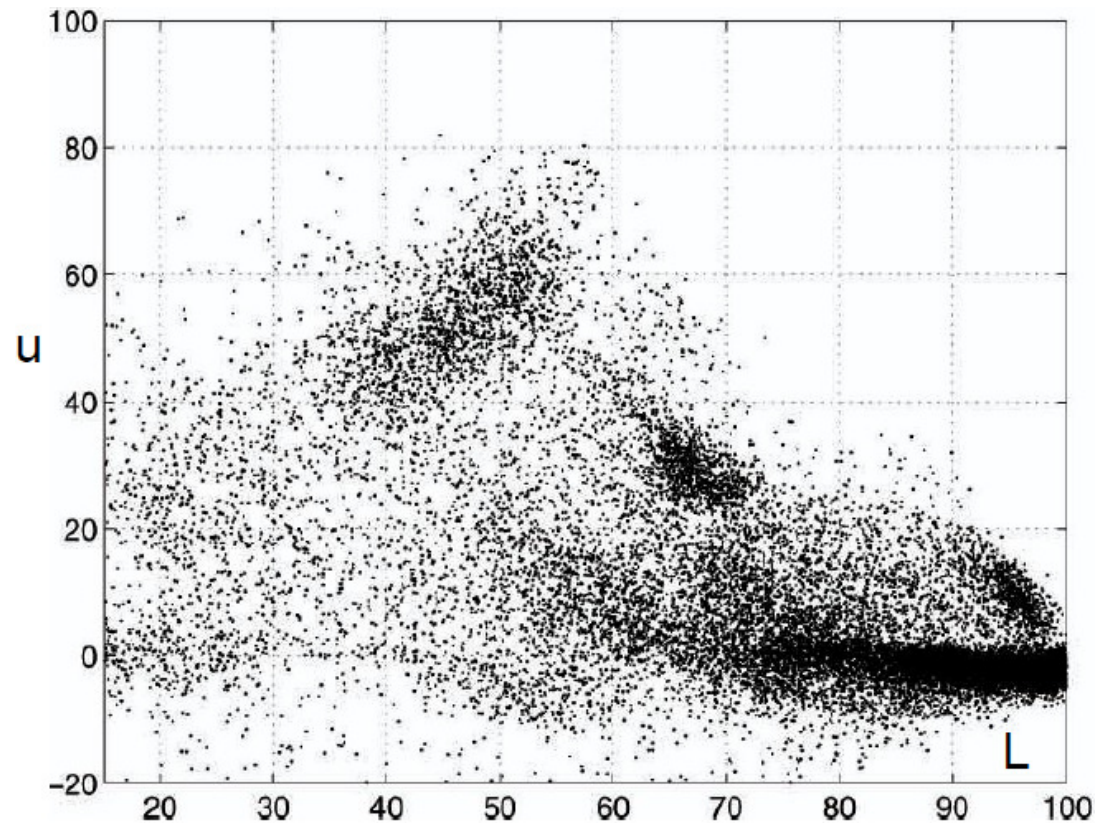
$$X^{(1)} = \frac{\begin{bmatrix} 4.2 \\ 1.7 \end{bmatrix} + \begin{bmatrix} 4.7 \\ 2.0 \end{bmatrix} + \begin{bmatrix} 2.9 \\ 1.5 \end{bmatrix} + \begin{bmatrix} 4.3 \\ 0.9 \end{bmatrix} + \begin{bmatrix} 2.6 \\ 2.2 \end{bmatrix}}{5} = \begin{bmatrix} 3.74 \\ 1.66 \end{bmatrix}$$

Application: Image Segmentation



Data Vector Formed
by Using (L, u, v, x, y)

Separate Bandwidth Along Each Dimension





Advantages

- Seeks the Mode of Distributions
- Ideal for Most Practical Cases
- Provides Model Free Nonlinear Regression
- Minimal Parameter Tuning

Disadvantages

- Needs All Data Points
- Often Computation Intensive

Summary

- Nonparametric Representation of Distributions
- Model Free Regression
- Nonparametric Classification
- Mean-Shift Clustering



Thank You