

CHAPTER 8:

NONPARAMETRIC METHODS

Stella Grasshof

Overview

1. Repetition
2. Intro
3. Probabiliy Density Function (pdf)
4. Properties and Usages
5. Nonparametric Estimation
 - a) ... of pdf
 - b) ... of class-conditional pdf
 - c) ... for regression
6. Summary

What I miss / want repeated

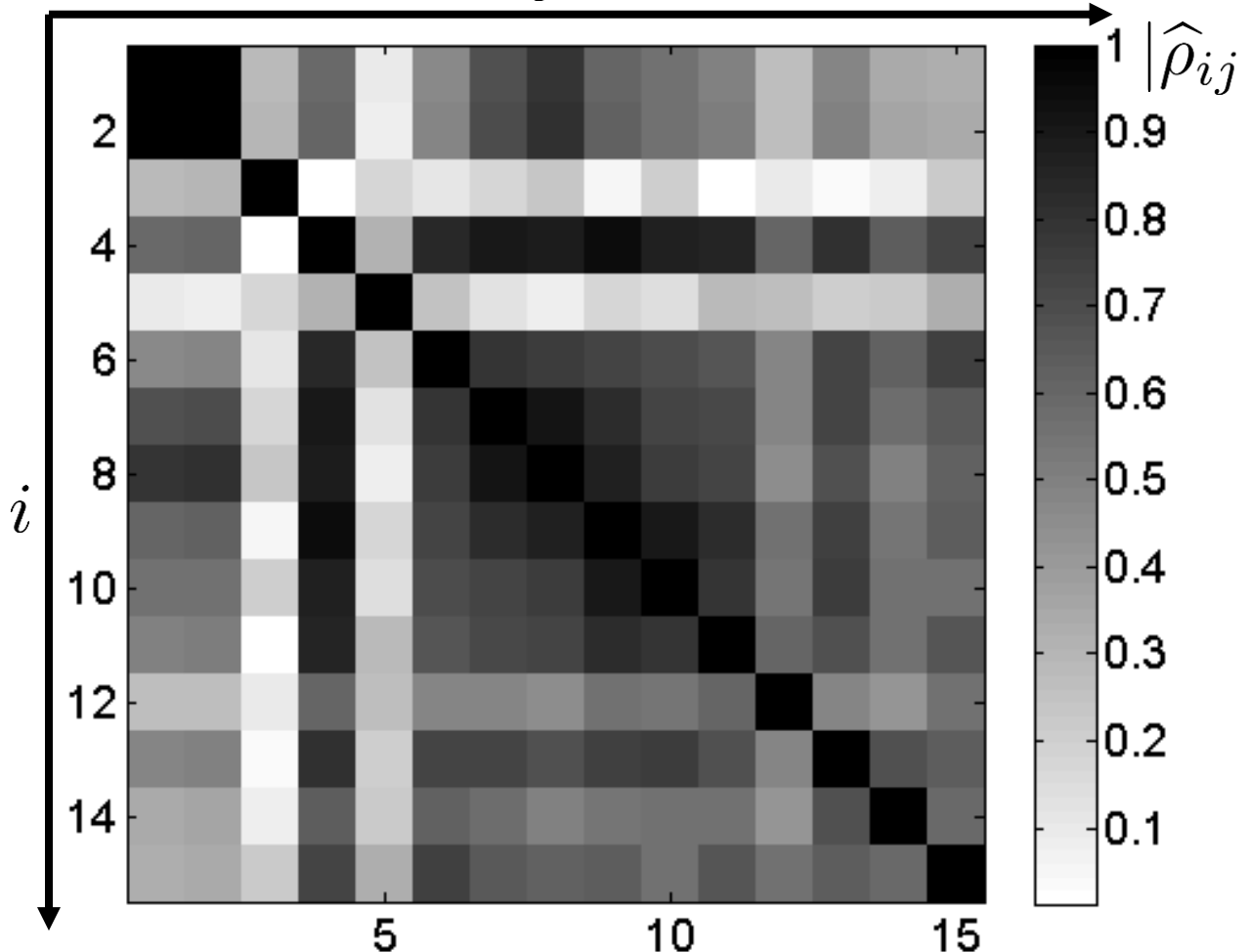
The word cloud contains the following terms:

- more applications
- covariance
- explain it to me like im
- im so lost
- classification
- overview of concepts
- parametric methods
- multivariate normal
- multivariate factor analysis
- exercise solutions
- notations datastructure
- concepts from beginning
- dimensionality reduction
- solutionssolutions
- a good book
- mle
- better schedule
- theoretical exercises
- multivariate methods
- pa
- illustrations
- everything
- solutions
- pca
- overview
- my dad
- explain it to me like im5
- animation for explanation
- covariance
- more graphs
- more time
- theoretical
- better book
- i'm so lost
- slower pace
- cca

3

Revisit: Correlation = scaled Cov

Example: absolute values of correlation matrix based on: $|\hat{\Sigma}|$
 $\hat{\Sigma} = \hat{\Sigma}^T$



Correlation matrix
consists of entries Σ
 $\rho_{ij} = \text{Corr}(\mathbf{x}_i, \mathbf{x}_j)$
pairwise correlation
between features

$\mathbf{x}_i, \mathbf{x}_j,$

$i, j = 1, \dots, D$

=> the darker
the pixel at (i,j) ,
the larger $|\rho_{ij}|$

Principal Components Analysis (PCA)

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

$$\mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_D]^T \in \mathbb{R}^D$$

1. Subtract mean $\mathbf{X} - \mathbf{M}$, $\mathbf{M} = (\hat{\boldsymbol{\mu}}, \dots, \hat{\boldsymbol{\mu}})^T \in \mathbb{R}^{N \times D}$
2. Compute covariance matrix $\hat{\boldsymbol{\Sigma}} = \frac{1}{N}(\mathbf{X} - \mathbf{M})^T(\mathbf{X} - \mathbf{M})$
3. Compute eigenvectors of covariance matrix

$$\hat{\boldsymbol{\Sigma}} \mathbf{w}_k = \lambda_k \mathbf{w}_k, \quad k = 1, \dots, D, \quad \lambda_i \geq \lambda_j, \quad i > j$$

$$\mathbf{w}_i^T \mathbf{w}_j = \begin{cases} 1 & , i = j \\ 0 & , i \neq j \end{cases} \quad [\mathbf{w}_1, \dots, \mathbf{w}_D] = \mathbf{W}^T$$

$$\mathbf{w}_k \text{ are principal components} \quad [\mathbf{w}_1, \dots, \mathbf{w}_K] = \mathbf{W}_K^T$$

4. New variables $\mathbf{z}_n = \mathbf{W}^T(\mathbf{x}_n - \hat{\boldsymbol{\mu}})$
5. Reconstruction $\hat{\mathbf{x}}_n = \mathbf{W}_K \mathbf{z}_n + \hat{\boldsymbol{\mu}}$

Principal Components Analysis (PCA)

1. Consider 2D data shall be mapped to 1D

$$\mathbf{x}_n \in \mathbb{R}^2 \mapsto \mathbf{z}_n \in \mathbb{R}^1$$

2. Covariance of mean centered data

$$\hat{\Sigma} = \frac{1}{N} (\mathbf{X} - \mathbf{M})^T (\mathbf{X} - \mathbf{M})$$

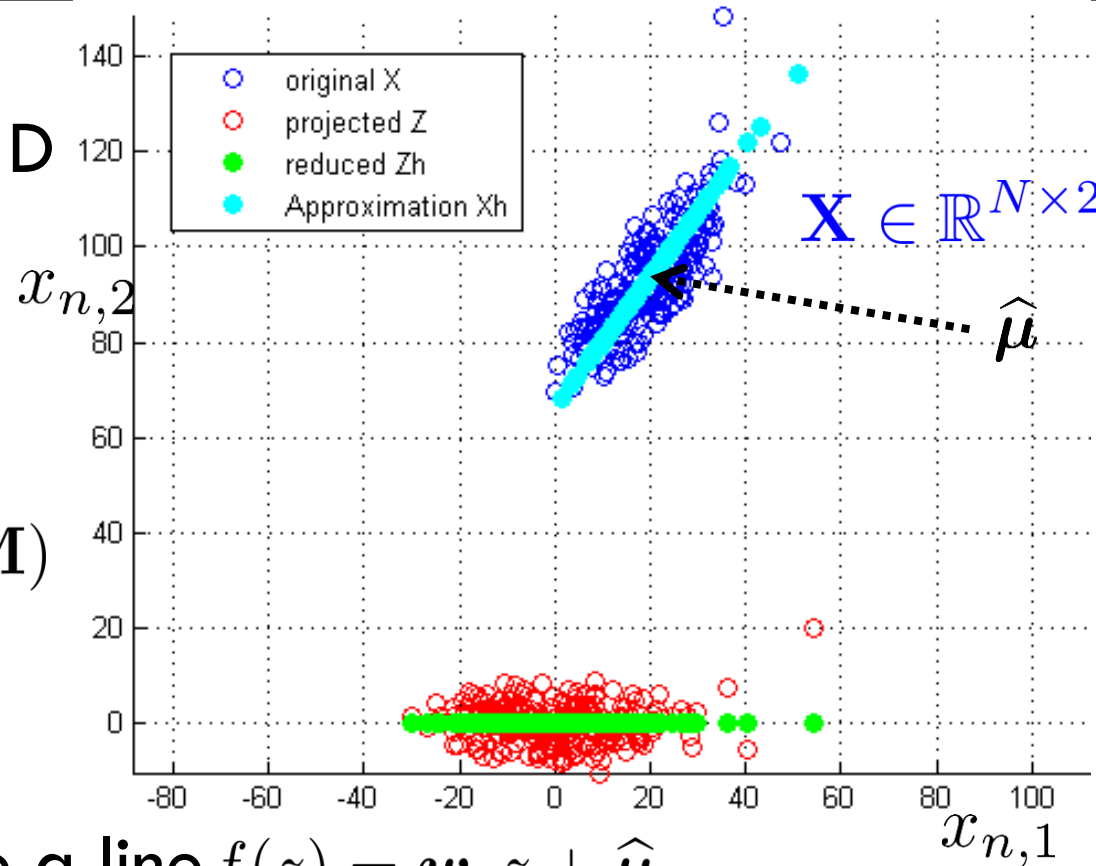
3. Eigenvectors

$$\hat{\Sigma} \mathbf{w}_k = \lambda_k \mathbf{w}_k, \quad k = 1, 2$$

4. Choose $k=1$ to define a line $f(z) = \mathbf{w}_1 z + \hat{\boldsymbol{\mu}}$ leads to

▣ Reduction $z_n = \mathbf{w}_1^T (\mathbf{x}_n - \hat{\boldsymbol{\mu}}) \in \mathbb{R}^1$

▣ Reconstruction $\hat{\mathbf{x}}_n = \mathbf{w}_1 z_n + \hat{\boldsymbol{\mu}} \in \mathbb{R}^2$



END Repetition

Nonparametric Estimation

- **Parametric:**
single global model
- **Semiparametric:**
small number of local models
- **Nonparametric:**
similar inputs have similar outputs
 - ▣ Keep the training data:
“let the data speak for itself”
 - ▣ Given x : interpolate from the closest training samples
 - ▣ Aka lazy/memory-based/case-based/instance-based learning

Probability Density Function (pdf)

Suppose random variable X has pdf f

Properties of

- Non-negative $f(x) \geq 0, \quad \forall x \in \mathbb{R}$
- Normalized $\int_{-\infty}^{\infty} f(x) dx = 1$
- A pdf and a cumulative distribution function relate

$$F(x) = P(X \leq x) = \int_{-\infty}^x f(t) dt \qquad \frac{d}{dx} F(x) = f(x)$$

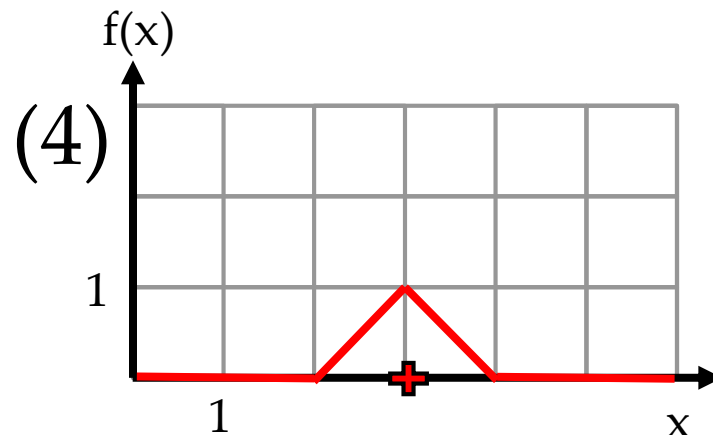
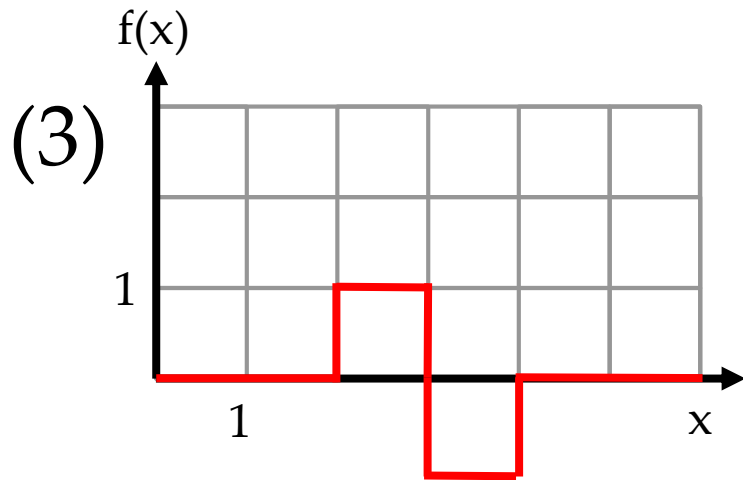
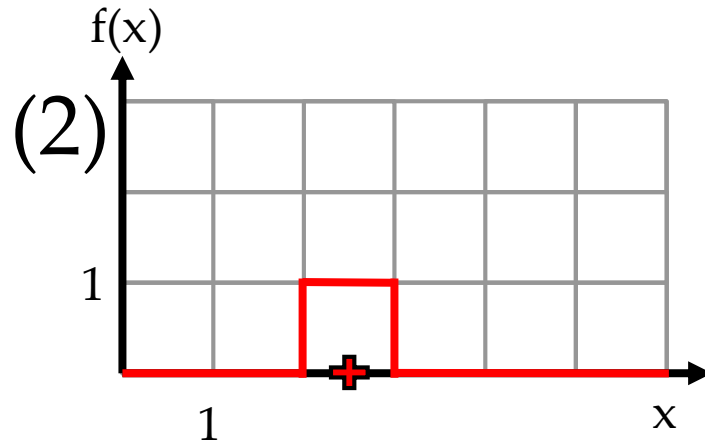
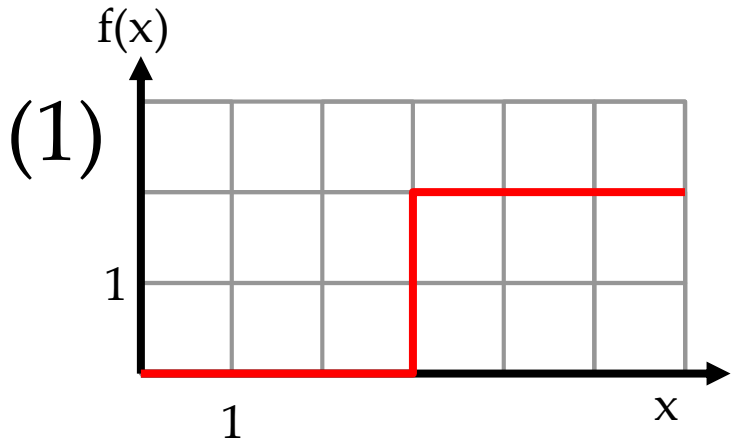
- Pdf relates area and probability

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

=> Usages?

Probability Density Function (pdf)

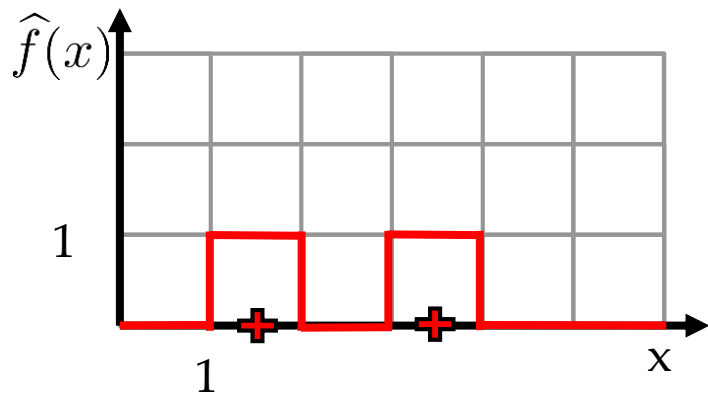
Which is a pdf?



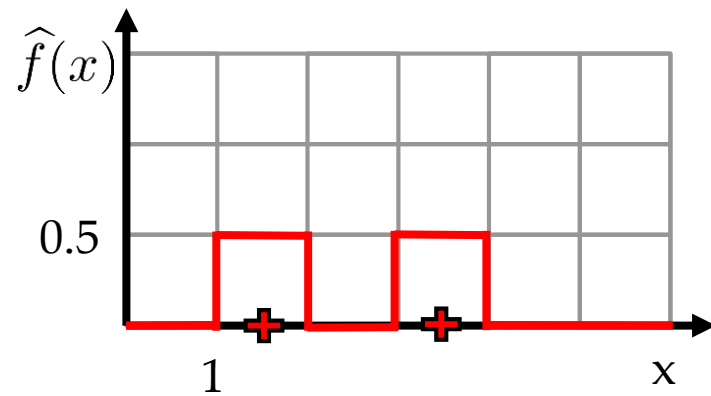
Probability Density Function (pdf)

Suppose random variable X has pdf f

it should be estimated from samples: $\text{+} = x_n, n = 1, \dots, N$

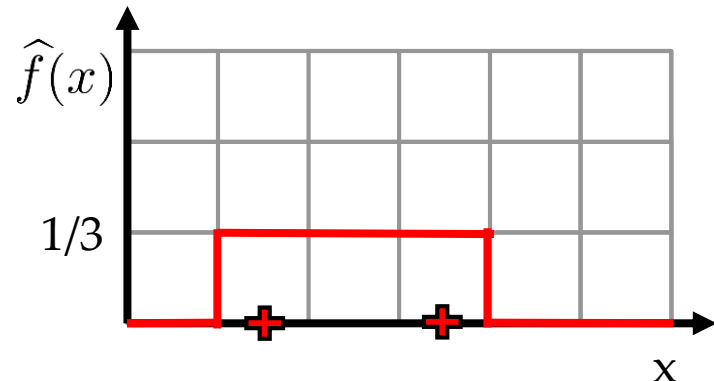


not a pdf



this is a pdf

Different possible estimates
for pdf from the samples



Density Estimation

- Given the training set x_n , $n = 1, \dots, N$ drawn iid from distribution described by pdf p
- \Rightarrow estimate pdf $p(x) \approx \hat{p}(x)$
- Divide data into bins of size h

$$N = 8$$

$$h = 1$$

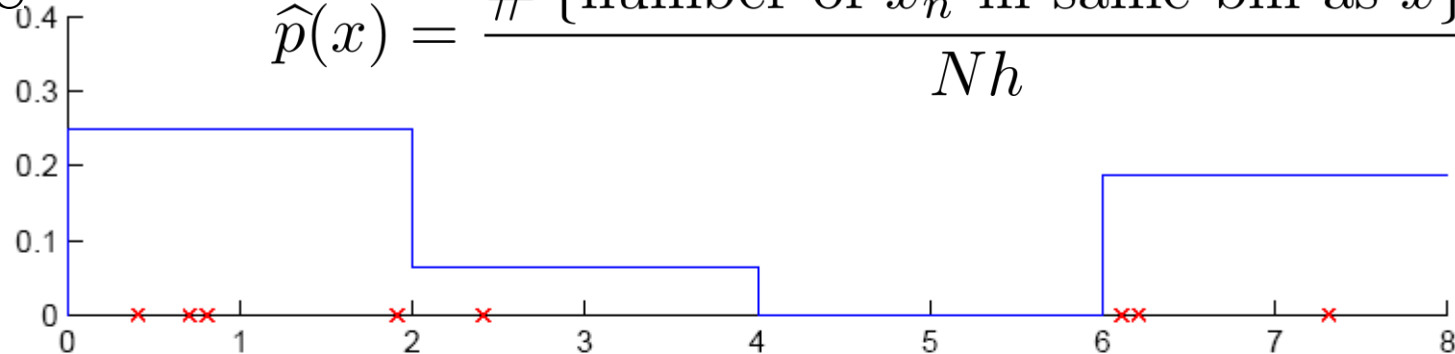


Histogram Estimator

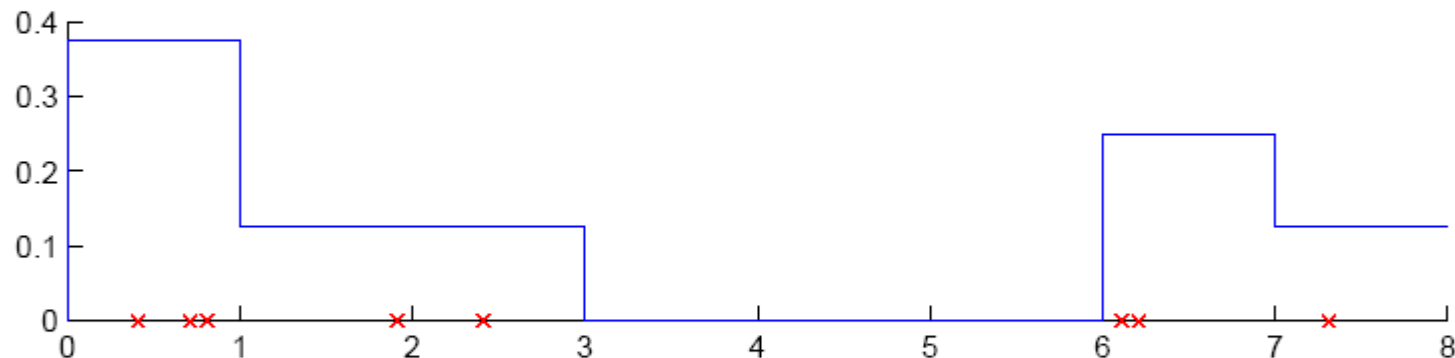
$$N = 8$$

$$h = 2$$

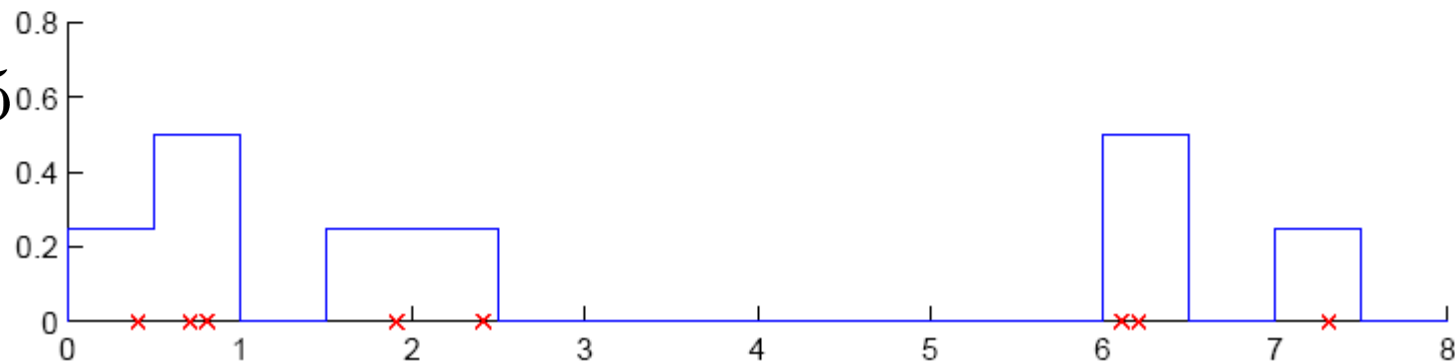
$$\hat{p}(x) = \frac{\# \{ \text{number of } x_n \text{ in same bin as } x \}}{Nh}$$



$$h = 1$$



$$h = 0.5$$

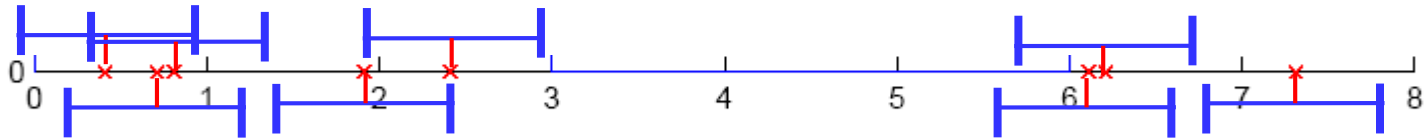


Density Estimation

- Given the training set x_n , $n = 1, \dots, N$ drawn iid from distribution described by pdf p
- \Rightarrow estimate pdf $p(x) \approx \hat{p}(x)$
- Divide data into bins of size h
- *Center at data points*

$$N = 8$$

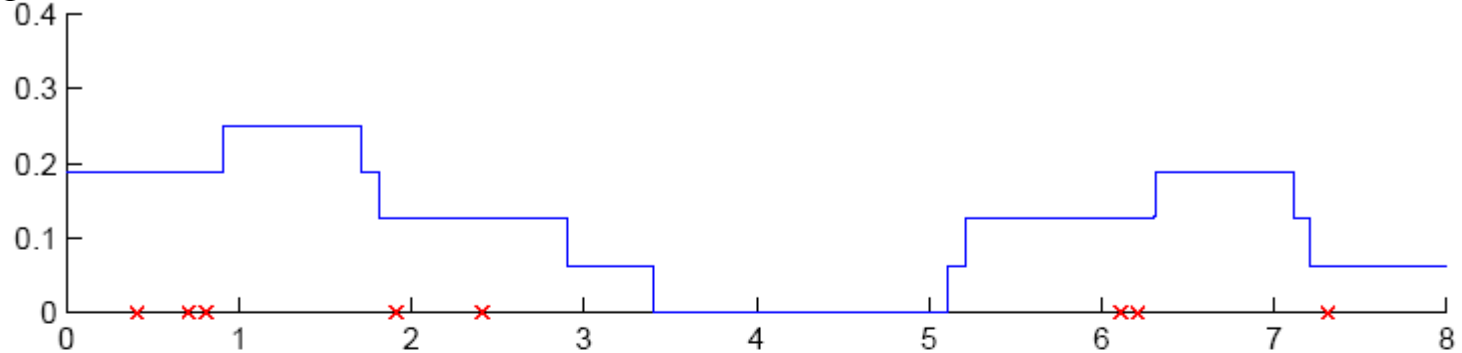
$$h = 1$$



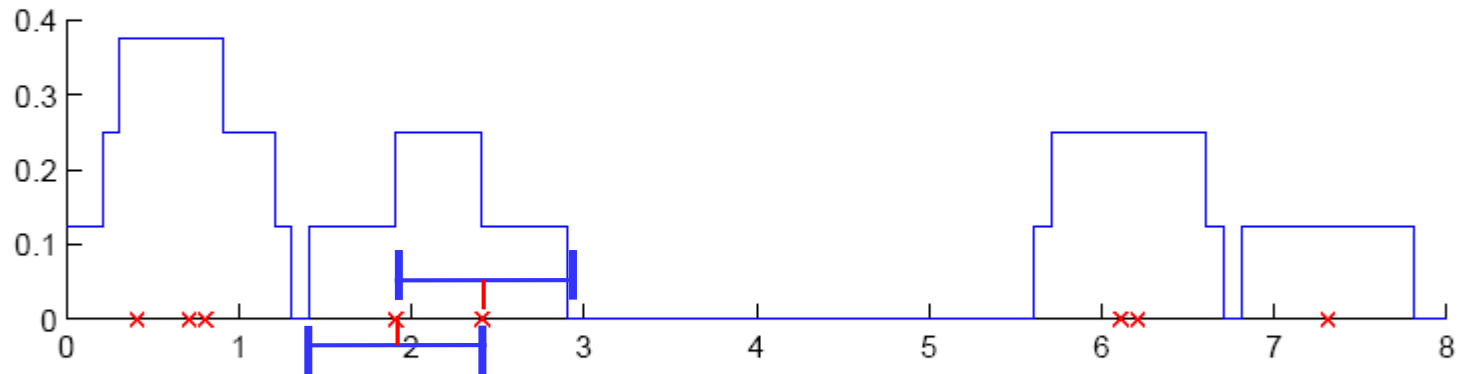
Naive Estimator

$$N = 8$$

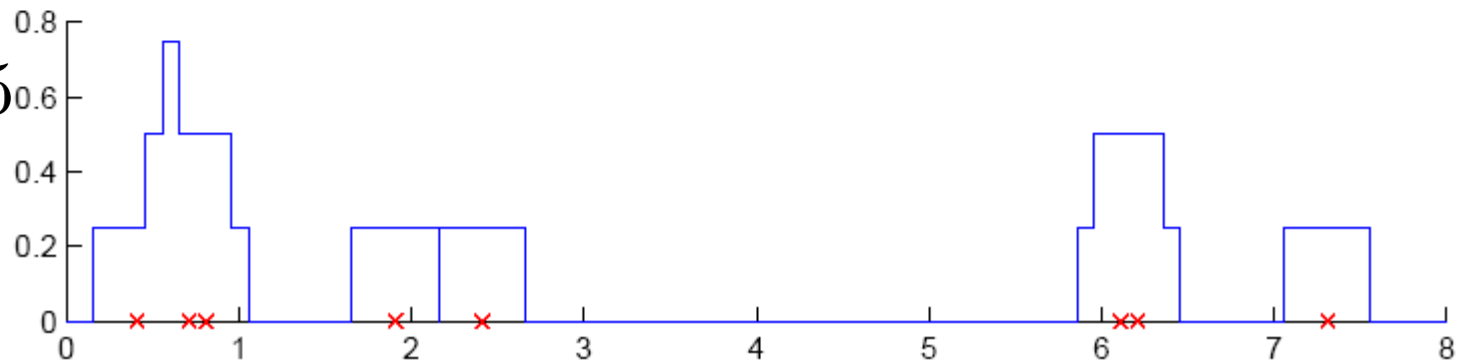
$$h = 2$$



$$h = 1$$



$$h = 0.5$$



Density Estimation

- Given the training set x_n , $n = 1, \dots, N$ drawn iid from $p(x)$

- Divide data into bins of size h

- Histogram:

$$\hat{p}(x) = \frac{\# \{ \text{number of } x_n \text{ in same bin as } x \}}{Nh}$$

- Naive estimator:

$$\hat{p}(x) = \frac{\# \{ x - h/2 < x_n \leq x + h/2 \}}{Nh}$$

or

$$\hat{p}(x) = \frac{1}{Nh} \sum_{n=1}^N w \left(\frac{x - x_n}{h} \right), \quad w(u) = \begin{cases} 1 & , |u| < 1/2 \\ 0 & , \text{else} \end{cases}$$

Drawback: Discontinuities

Kernel Estimator

- Replace:

$$w(u) = \begin{cases} 1 & , |u| < 1/2 \\ 0 & , \text{ else} \end{cases}$$

- by another kernel function, e.g. Gaussian kernel:

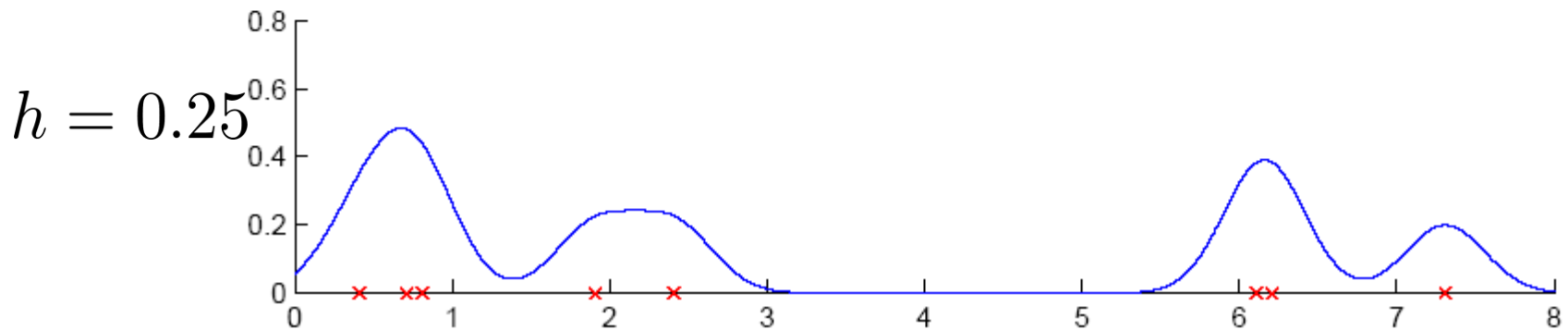
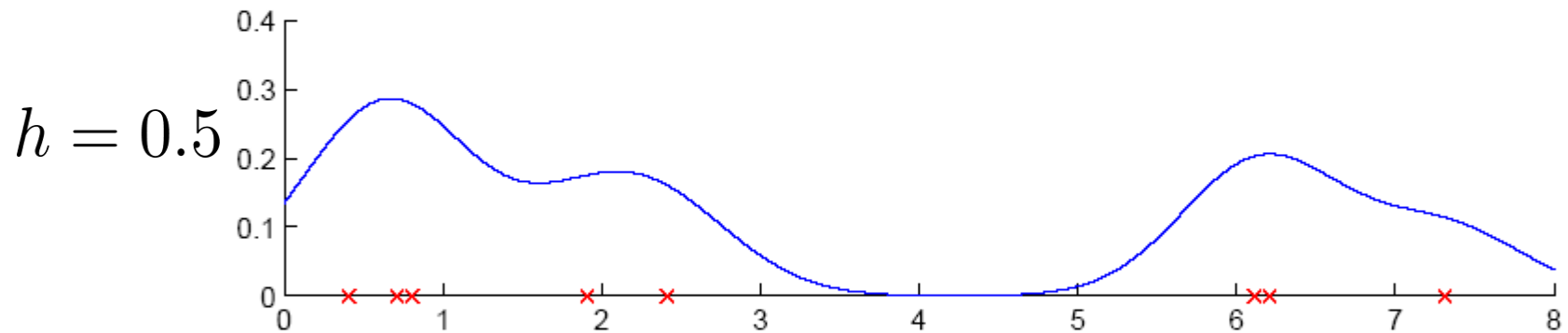
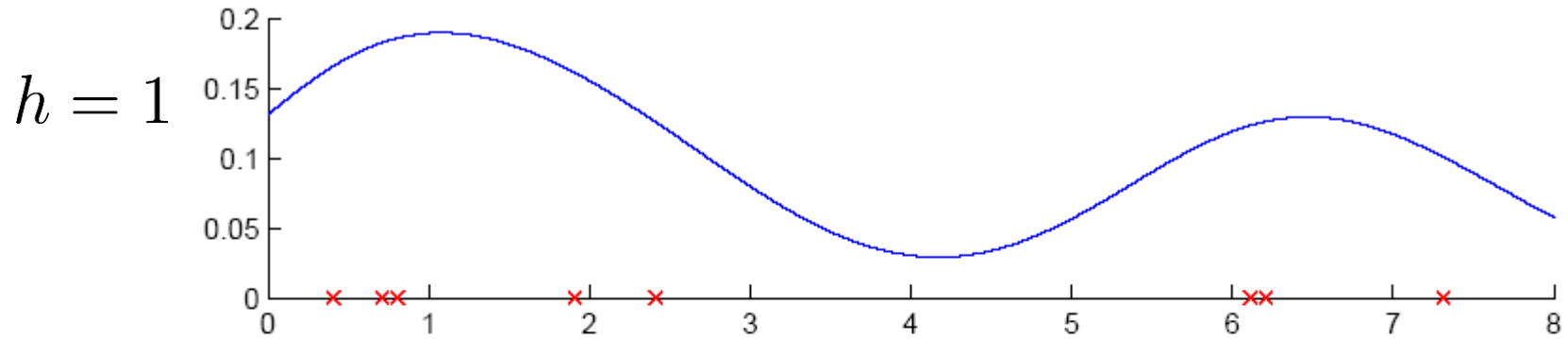
$$\varphi(x) = \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{x^2}{2} \right] \quad \text{this is: } \mathcal{N}(0, 1)$$

- Gives kernel estimator (Parzen windows)

$$\hat{p}(x) = \frac{1}{Nh} \sum_{n=1}^N \varphi \left(\frac{x - x_n}{h} \right)$$

... is continuous, differentiable

Kernel Estimator

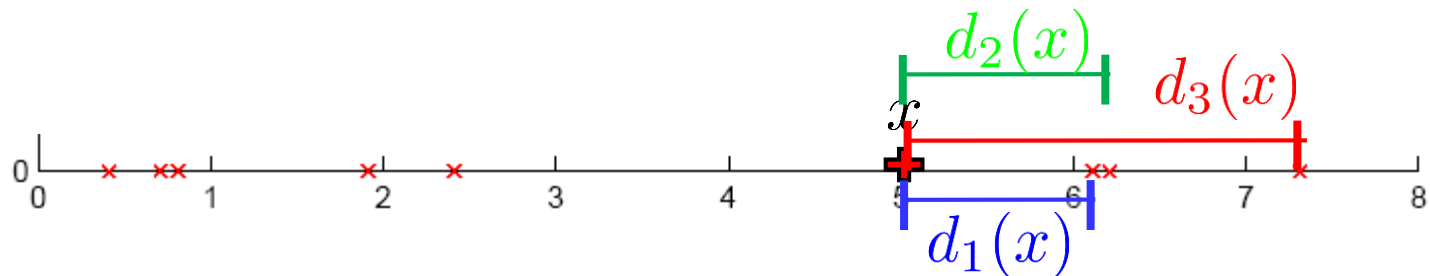


k-Nearest Neighbor Estimator

Instead of fixing bin width h
and counting the number of instances
fix the instances (neighbors) k and check bin width

$\Rightarrow d_k(x)$ distance to k th closest instance to x

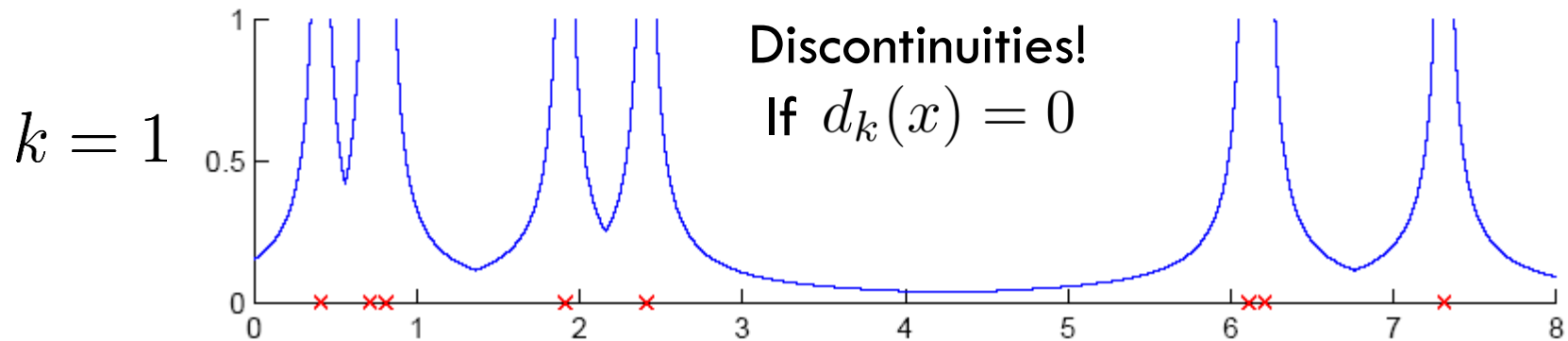
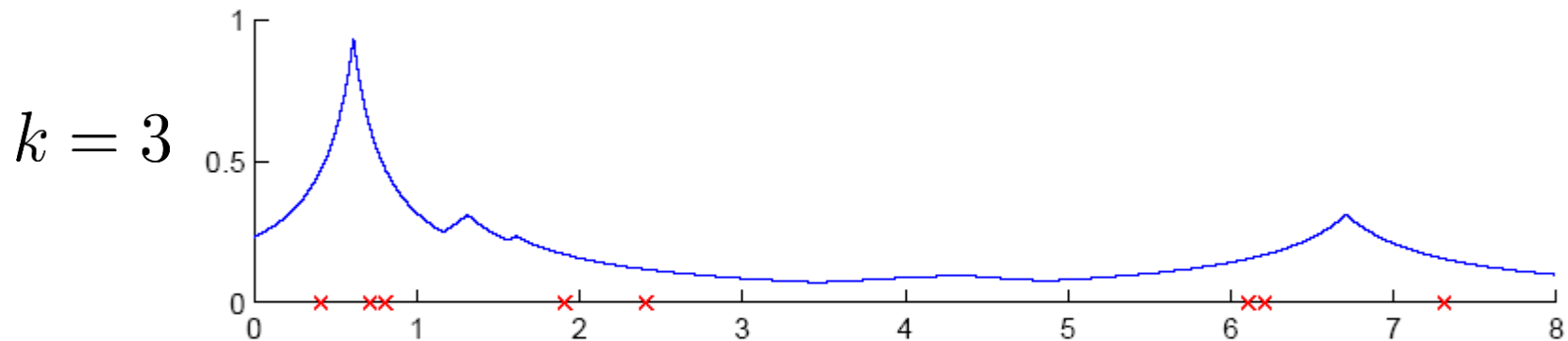
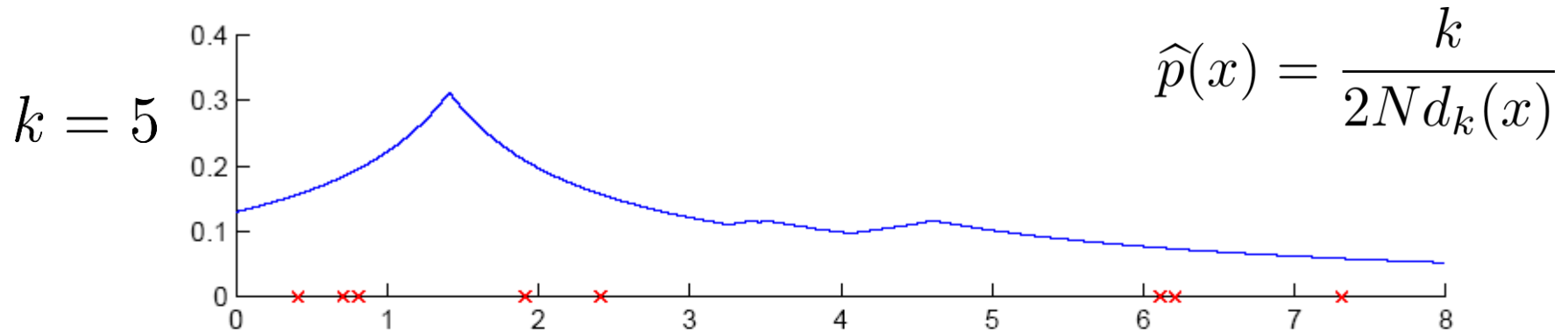
aim: small distance = high pdf



$$\hat{p}(x) = \frac{k}{2Nd_k(x)}$$

Drawback: Discontinuities and is NOT pdf
improvement by kernel function

k-Nearest Neighbor Estimator



Multivariate Data

$$\mathbf{x}, \mathbf{x}_n \in \mathbb{R}^D$$

Kernel density estimator

$$\hat{p}(\mathbf{x}) = \frac{1}{Nh^D} \sum_{n=1}^N \varphi \left(\frac{\mathbf{x} - \mathbf{x}_n}{h} \right)$$

Multivariate Gaussian kernel

□ spheric $\mathcal{N}(\mathbf{0}, \mathbf{I})$

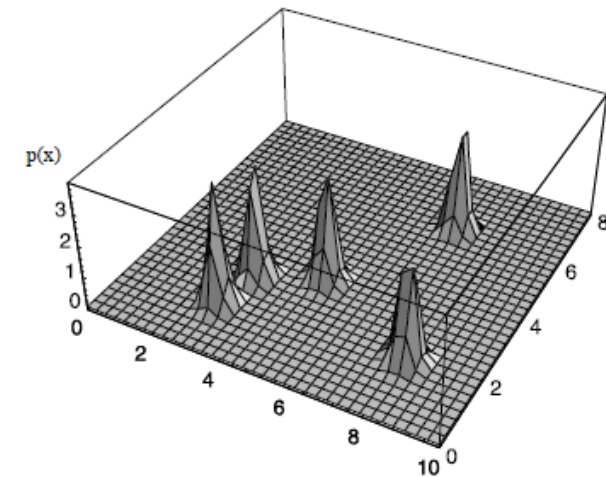
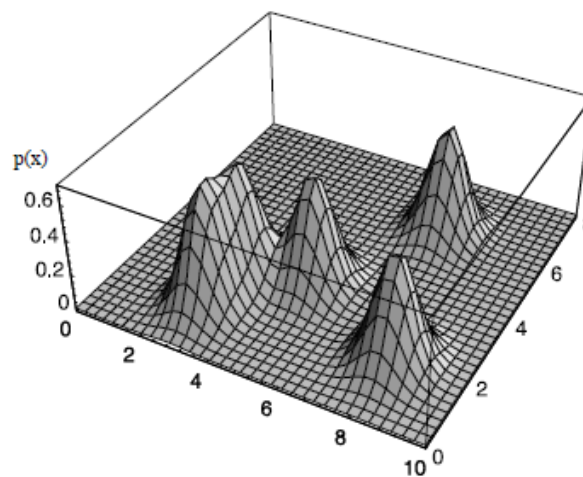
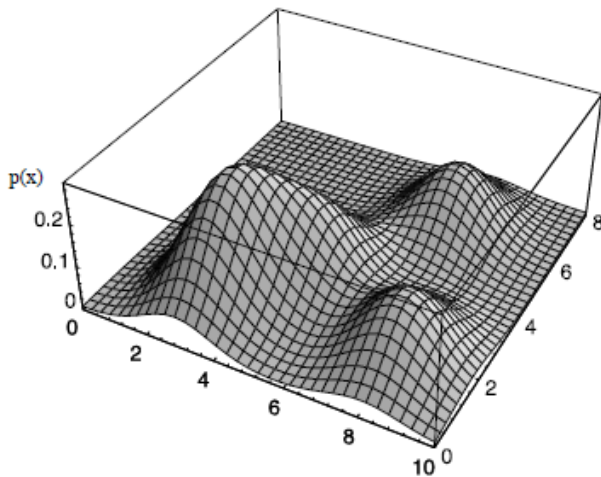
$$\varphi(\mathbf{x}) = \frac{1}{(2\pi)^{D/2}} \exp \left[-\frac{\|\mathbf{x}\|^2}{2} \right]$$

□ ellipsoid $\mathcal{N}(\mathbf{0}, \mathbf{S})$

$$\varphi(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} |\mathbf{S}|^{1/2}} \exp \left[-\frac{1}{2} \mathbf{x}^T \mathbf{S}^{-1} \mathbf{x} \right]$$

Kernel Estimator

Estimation from 5 points with varying width



Nonparametric Classification

- Discriminant for class k : $\hat{g}_i(\mathbf{x}) = \hat{P}(C_i)\hat{p}(\mathbf{x}|C_i)$
- Class probability as before: $\hat{P}(C_i) = \frac{N_i}{N}$
- Estimate conditional pdf $p(\mathbf{x}|C_i) \approx \hat{p}(\mathbf{x}|C_i)$
 - ▣ Kernel estimator (same as before, but select cases)

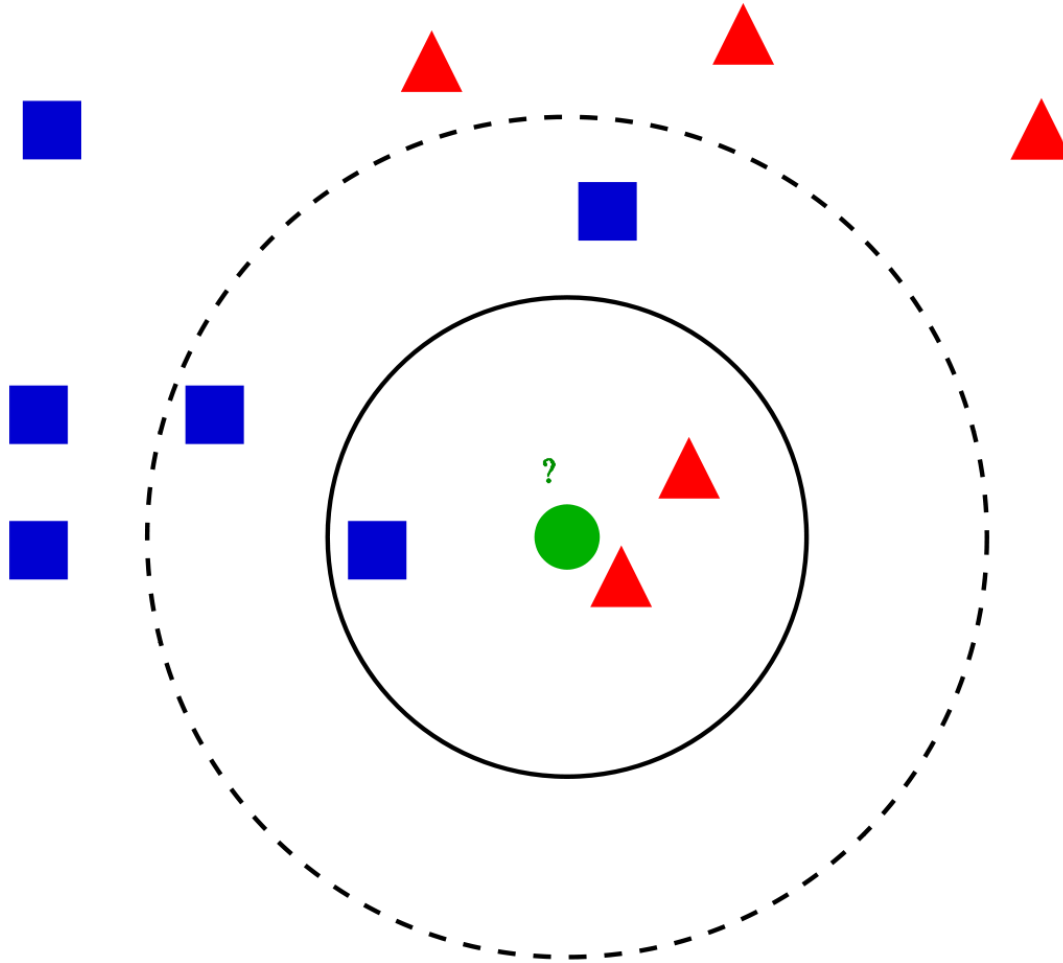
$$\hat{p}(\mathbf{x}|C_i) = \frac{1}{N_i h^D} \sum_{n=1}^N 1(\mathbf{x}_n \in C_i) \varphi\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

- ▣ k -NN estimator: get k nearest neighbors of \mathbf{x}
 - k_i of k neighbors belong to class C_i
 - $V^k(\mathbf{x})$ volume of D -dimensional hypersphere centered at \mathbf{x}


then: $\hat{p}(\mathbf{x}|C_i) = \frac{k_i}{N_i V^k(\mathbf{x})}$

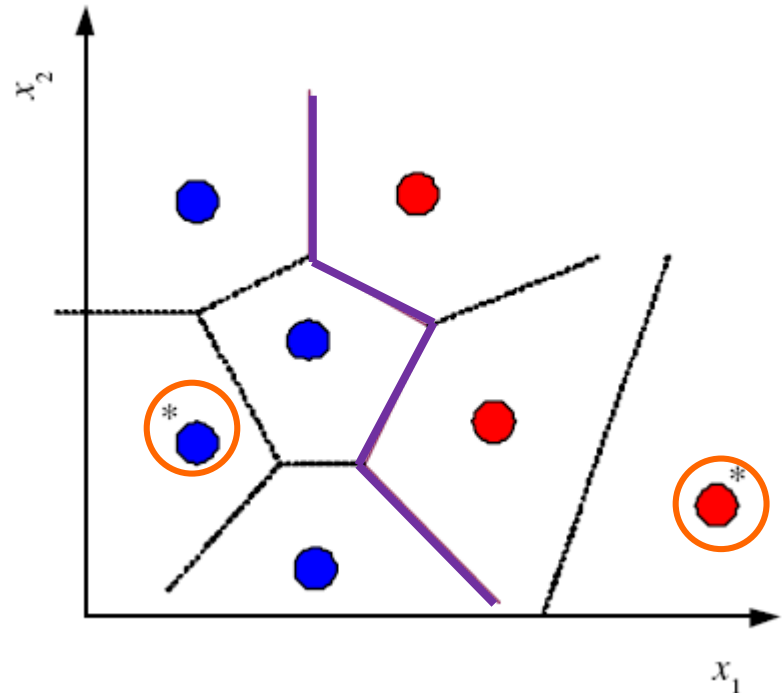
Decision is majority vote: $\hat{P}(C_i|\mathbf{x}) = \frac{k_i}{k}$

k-Nearest Neighbor Classifier



Condensed Nearest Neighbor

- Find a subset Z of X that is small and accurate in classifying X (Hart, 1968)
- The points  do not contribute to the decision boundary
⇒ delete them from X
receive smaller set Z



Condensed Nearest Neighbor

Incremental algorithm: Add instance if needed

$\mathcal{Z} \leftarrow \emptyset$

Repeat

For all $\mathbf{x} \in \mathcal{X}$ (in random order)

Find $\mathbf{x}' \in \mathcal{Z}$ s.t. $\|\mathbf{x} - \mathbf{x}'\| = \min_{\mathbf{x}^j \in \mathcal{Z}} \|\mathbf{x} - \mathbf{x}^j\|$

If $\text{class}(\mathbf{x}) \neq \text{class}(\mathbf{x}')$ add \mathbf{x} to \mathcal{Z}

Until \mathcal{Z} does not change

Distance-based Classification

- Find a distance function $\mathcal{D}(x_i, x_j)$ such that:
 - if x_i, x_j belong to the same class: $\mathcal{D}(x_i, x_j)$ small
 - if x_i, x_j belong to different classes: $\mathcal{D}(x_i, x_j)$ large
- Assume a parametric model and learn its parameters using data, e.g.,

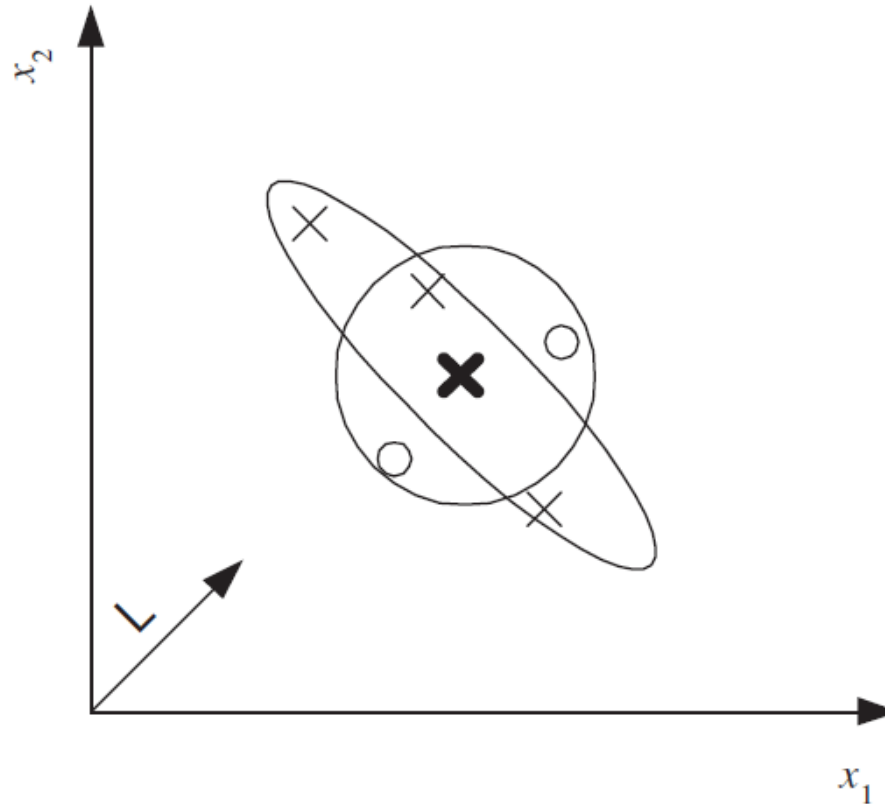
$$\mathcal{D}(x_i, x_j | M) = (x_i - x_j)^T M (x_i - x_j)$$

Learning a Distance Function

- Distances in high dimensions, can be represented as euclidean distance in lower dimensions
- Consider dimensionality reduction: $z_i = Lx_i$
- $M = L^T L$ is $D \times D$ and L is $K \times D$

$$\begin{aligned}\mathcal{D}(x_i, x_j | M) &= (x_i - x_j)^T M (x_i - x_j) = \dots \\ &= \|z_i - z_j\|_2^2\end{aligned}$$

- Similarity-based representation using similarity scores



Euclidean distance (circle) is not suitable,
Mahalanobis distance using an \mathbf{M} (ellipse) is suitable.
After the data is projected along \mathbf{L} , Euclidean distance can be used.

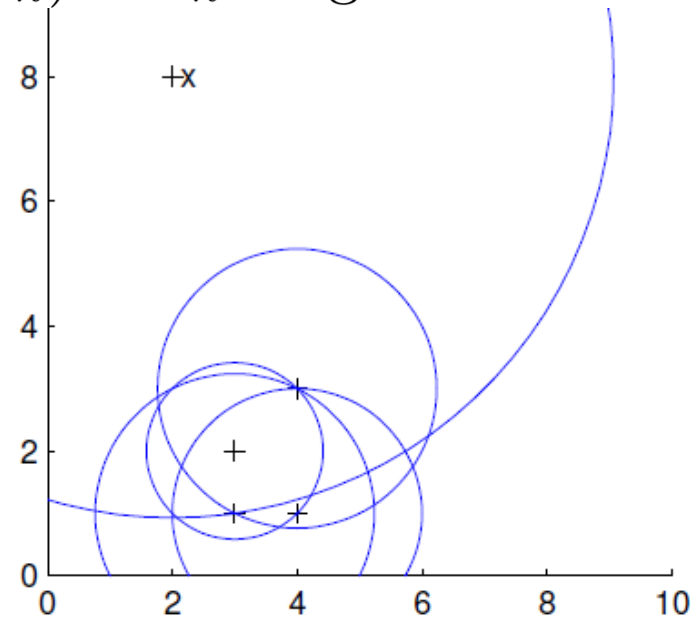
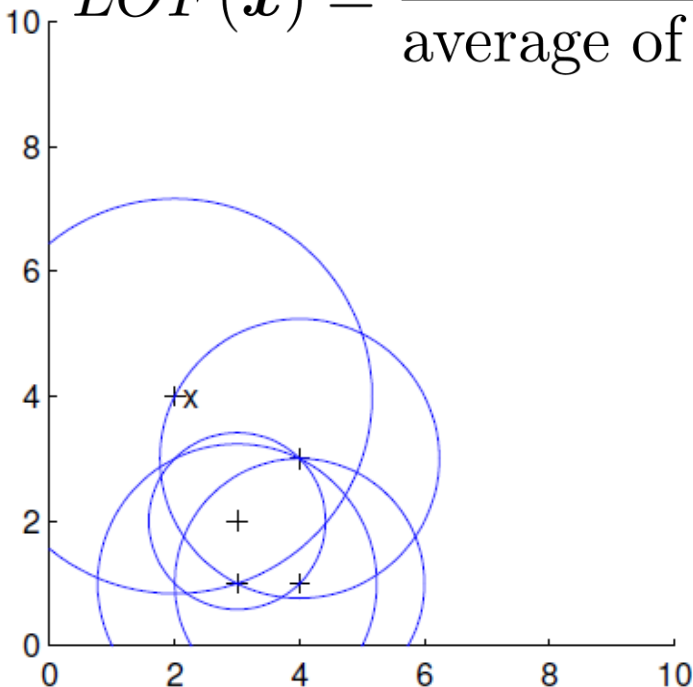
Outlier Detection

- Find outlier/novelty points
- Not a two-class problem!
 - ▣ Outliers vary, are seldom and unlabeled
- Instead: find instances with low probability
- In nonparametric case:
Find instances far away from other instances

Local Outlier Factor (LOF)

- Circle around each point, containing $k=3$ next neighbors
- Is the distance to the 3rd nearest neighbor similar to other points in that area?
 - ▣ Yes: LOF near 1
 - ▣ No: LOF bigger than 1 \Rightarrow Outlier

$$LOF(x) = \frac{d_k(x)}{\text{average of } d_k(x_n) : x_n \text{ neighbors of } x}$$



Nonparametric Regression

- Regression:

- ▣ Given training data: (\mathbf{x}_n, y_n)

- ▣ Estimate function such that $y_n = g(\mathbf{x}_n) + \epsilon$

$$y_n \approx g(\mathbf{x}_n)$$

- Nonparametric regression = smoothing models, e.g.

$g(\mathbf{x}) :=$ average of y_n from \mathbf{x}_n in the bin of \mathbf{x}

- Regressogram:

- ▣ Divide data into bins of width h

- ▣ Define value for each bin as:

- average of points inside of bin training data

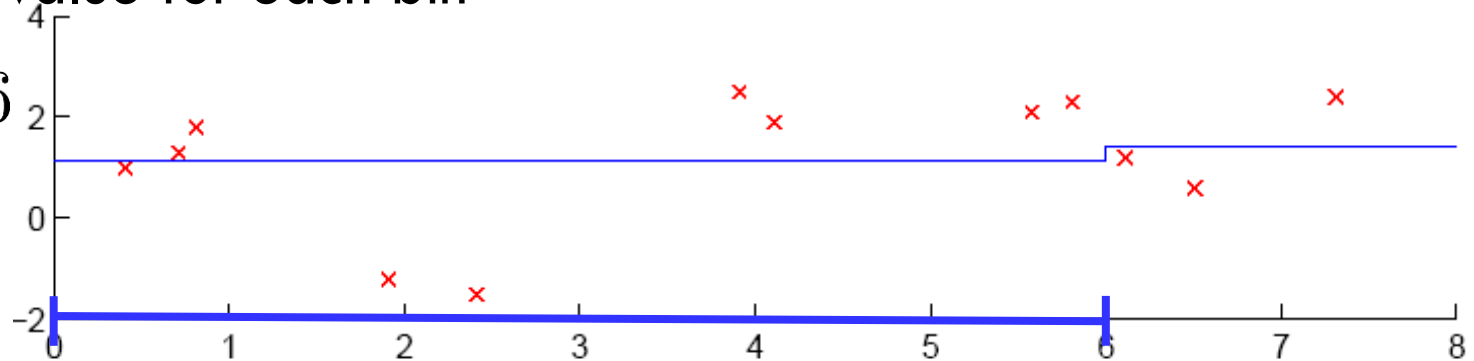
- Line through points

- ...

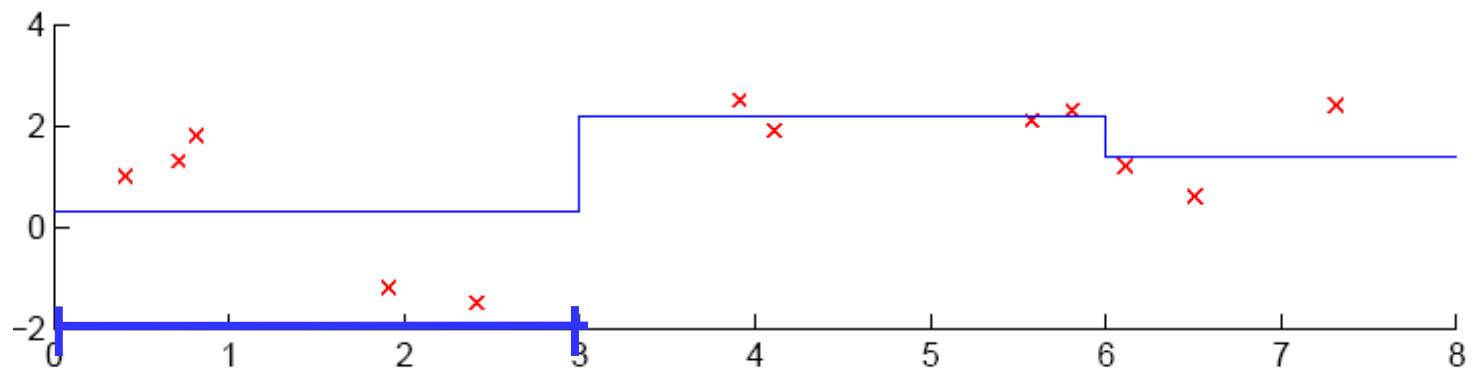
Regressogram smoother

Mean value for each bin

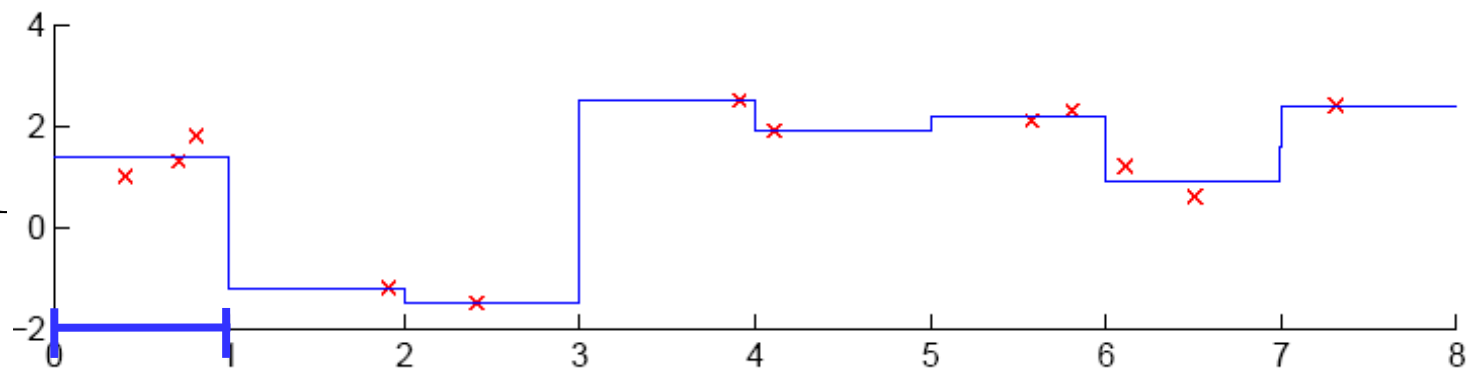
$h = 6$



$h = 3$



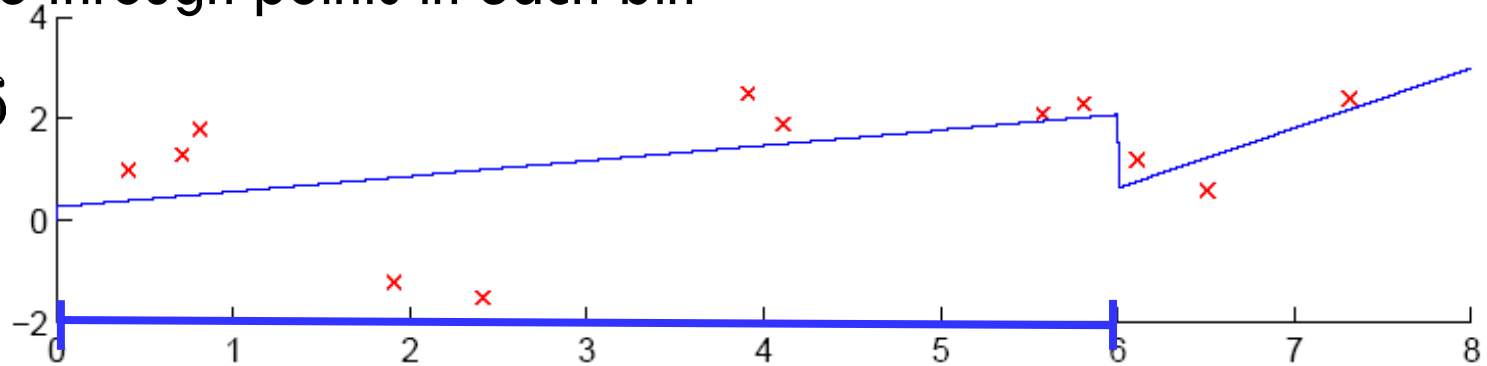
$h = 1$



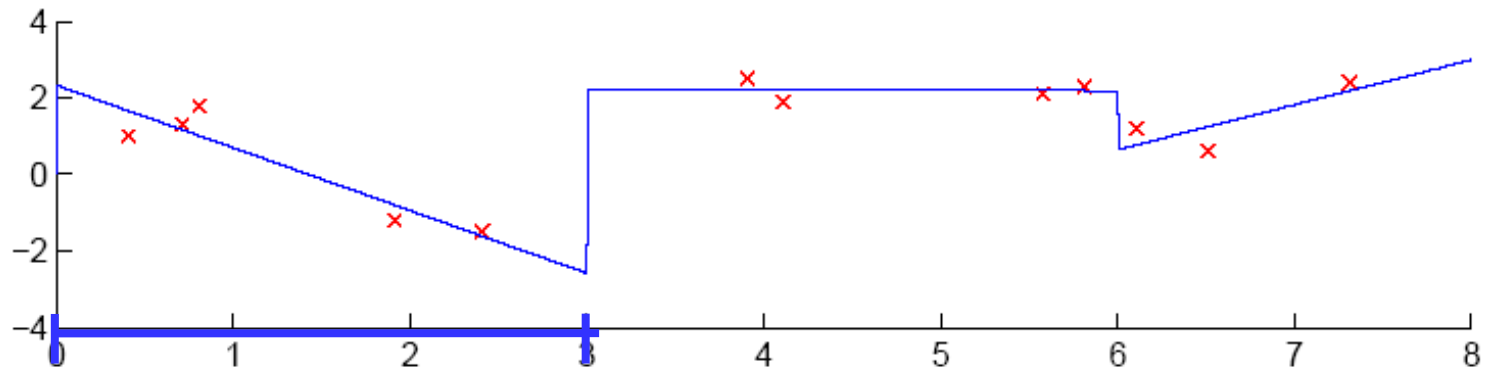
Nonparametric Regression

Fit line through points in each bin

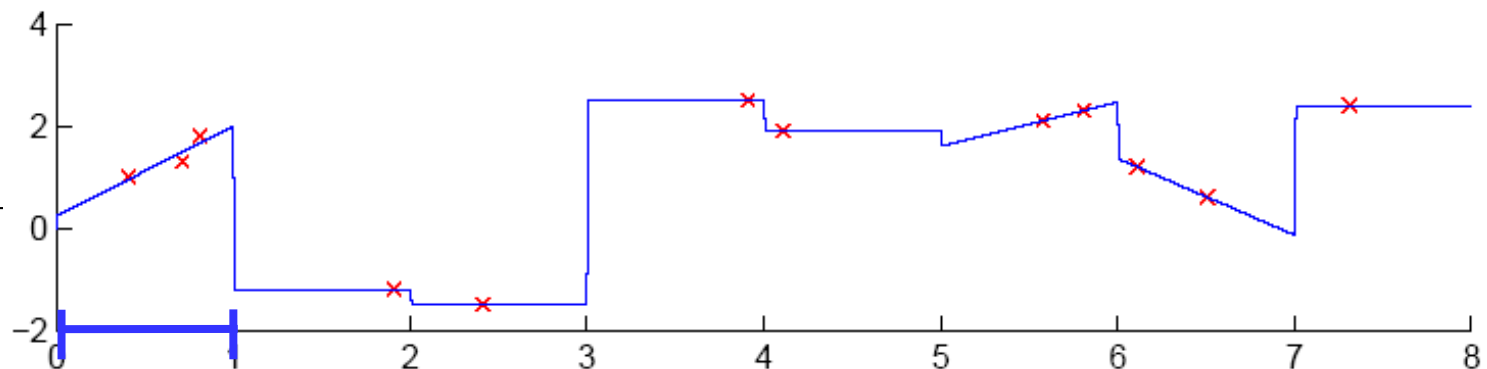
$h = 6$



$h = 3$



$h = 1$



Running Mean/Kernel Smoother

$$\hat{g}(x) = \frac{\sum_{n=1}^N \varphi\left(\frac{x-x_n}{h}\right) y_n}{\sum_{n=1}^N \varphi\left(\frac{x-x_n}{h}\right)}$$

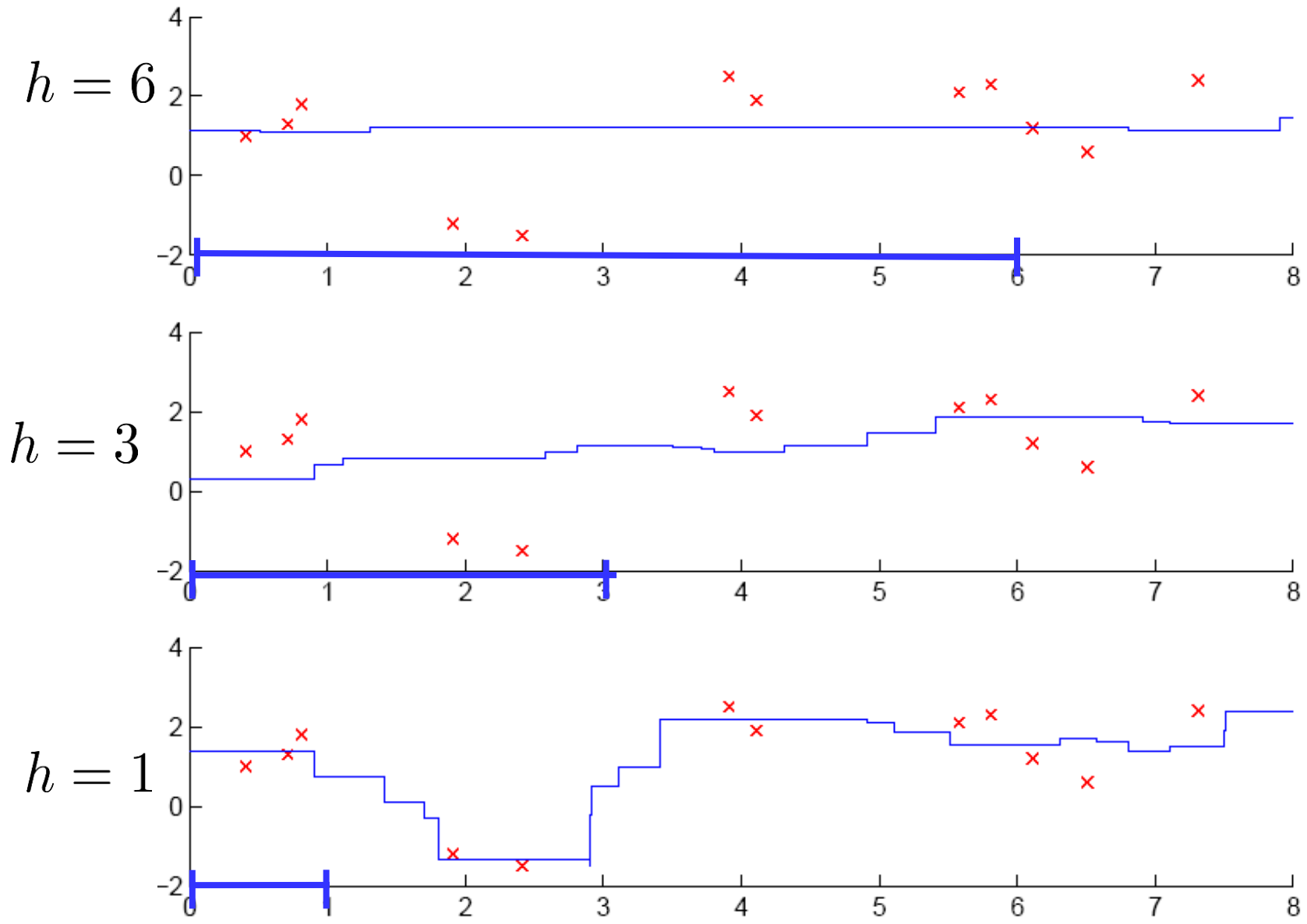
- Running mean smoother

$$\varphi(x) = \begin{cases} 1 & , \text{ if } |x| < 1 \\ 0 & , \text{ else} \end{cases}$$

- Line smoother by local regression
locally weighted line fitting
- Kernel smoother:

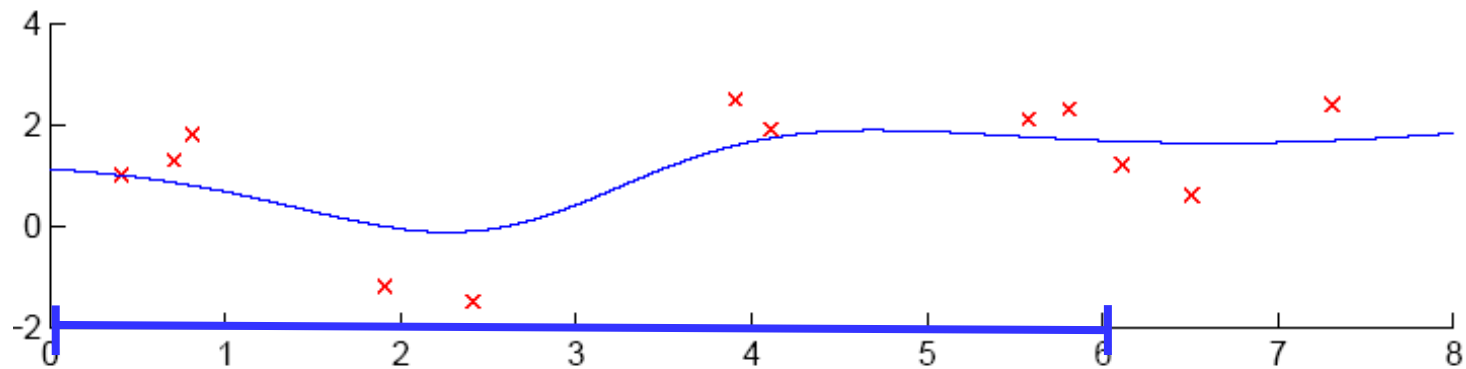
$$\varphi(x) = \mathcal{N}(0, 1)$$

Running mean smoother

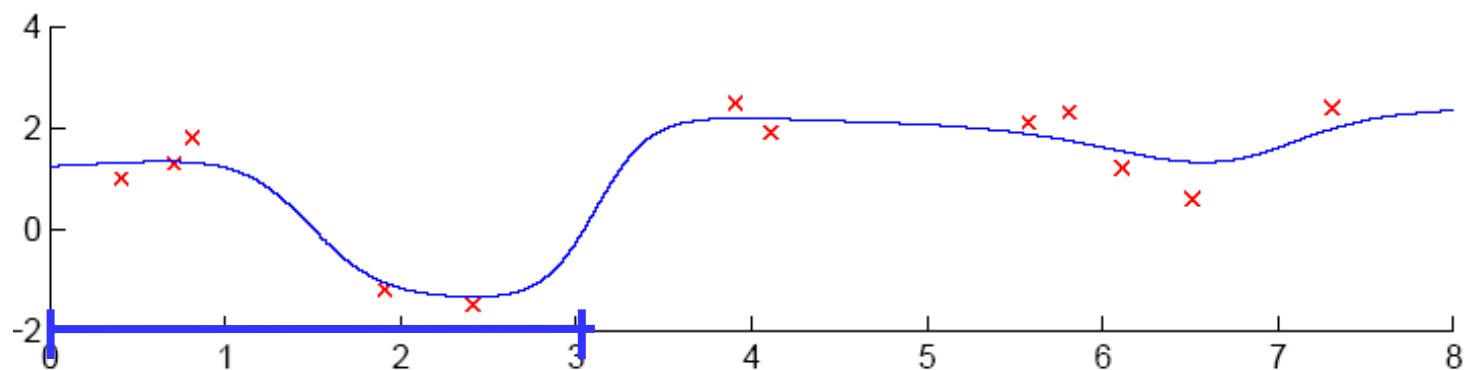


Kernel Smoother

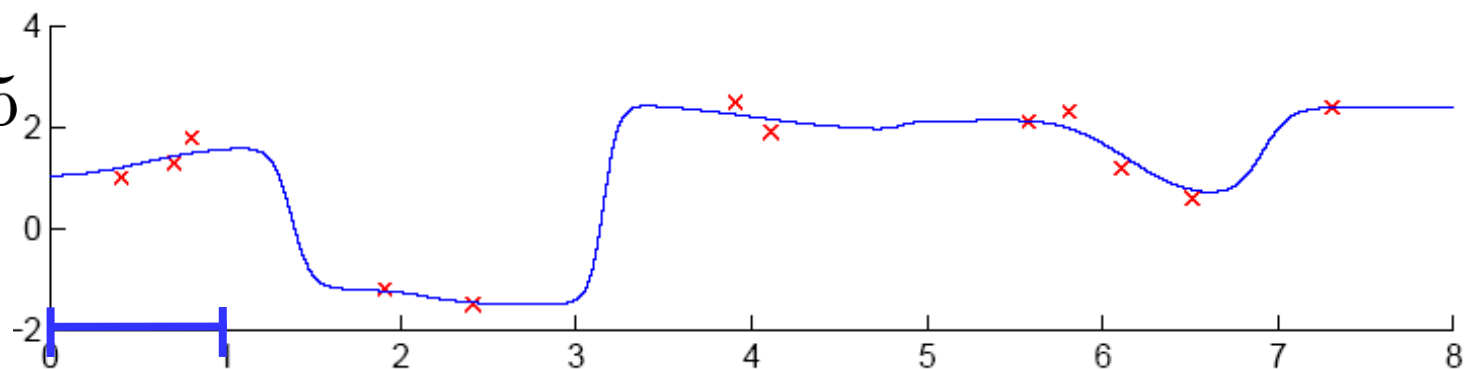
$h = 1$



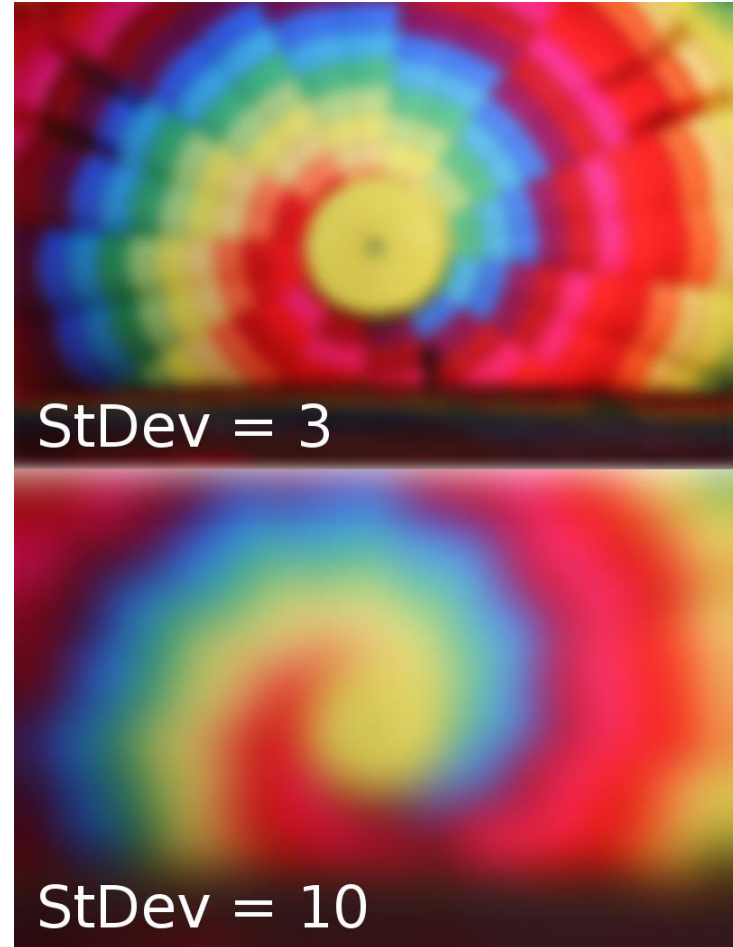
$h = 0.5$



$h = 0.25$



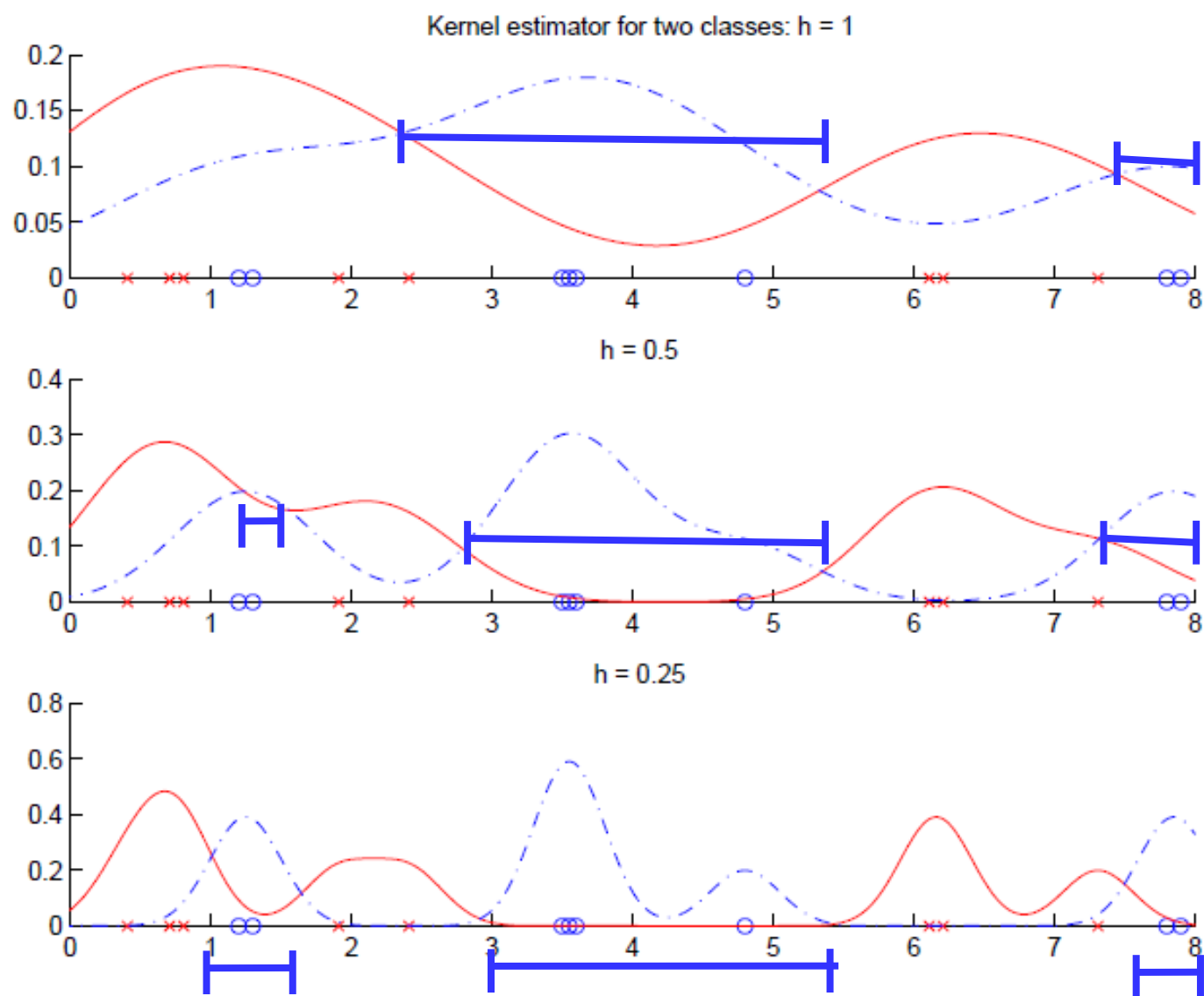
Gaussian Filtered Image



[https://en.wikipedia.org/wiki/Gaussian_blur#/media/
File:Cappadocia_Gaussian_Blur.svg](https://en.wikipedia.org/wiki/Gaussian_blur#/media/File:Cappadocia_Gaussian_Blur.svg)

How to Choose k or h ?

- k or h is small:
 - ▣ single instances matter
 - \Rightarrow bias is small, variance is large
 - \Rightarrow undersmoothing: High complexity
- k or h increases:
 - ▣ average over more instances
 - \Rightarrow variance decreases but bias increases
 - \Rightarrow oversmoothing: Low complexity
- Cross-validation is used to finetune k or h



Summary

Nonparametric

- No explicit model
- Use training data
 - ▣ “Smart” subset of training data
- More robust to outliers
- Different possibilities for pdf, regression ...

APPENDIX