

# 2nd Congress of Polish Statistics



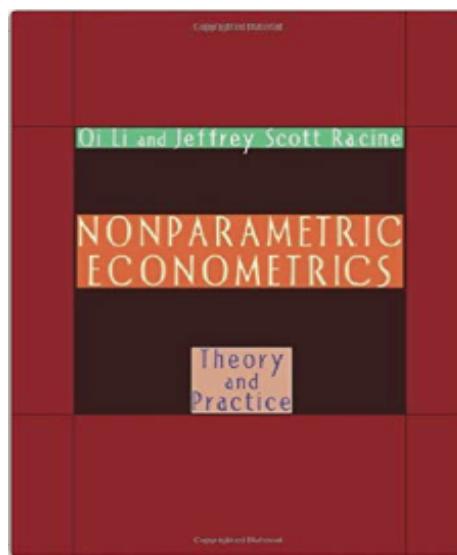
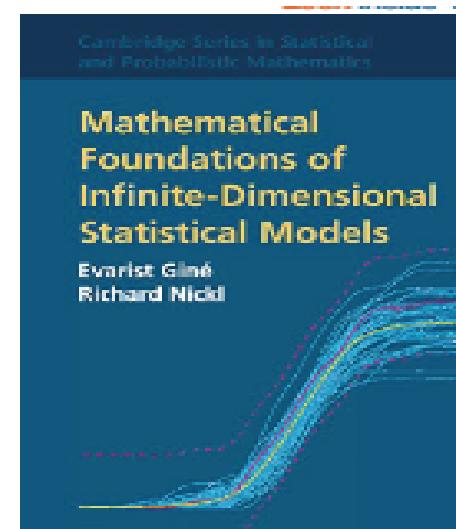
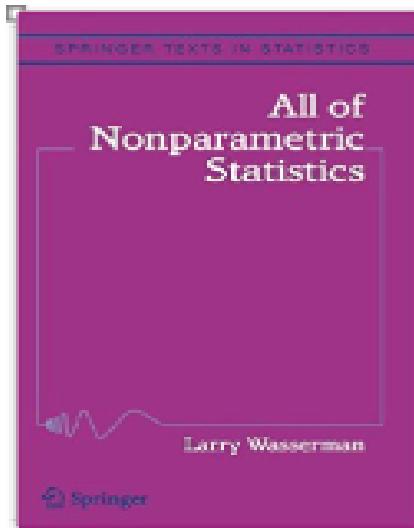
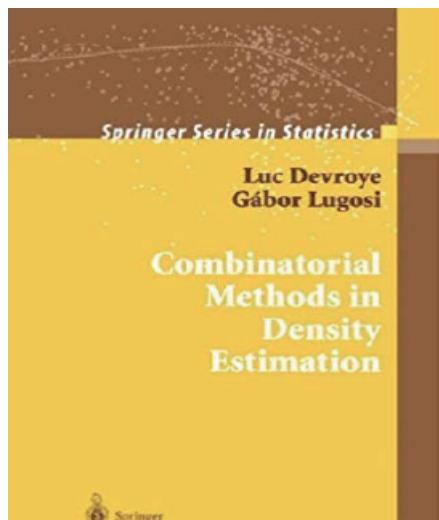
# **Nearest Neighbor Methods for Nonlinear Regression: Time Series Data Case**

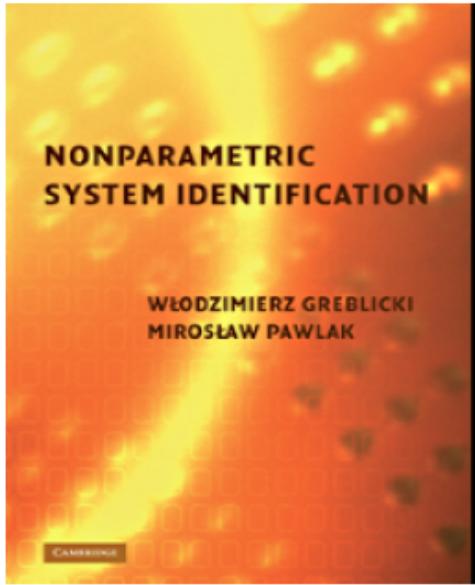
**Miroslaw Pawlak**

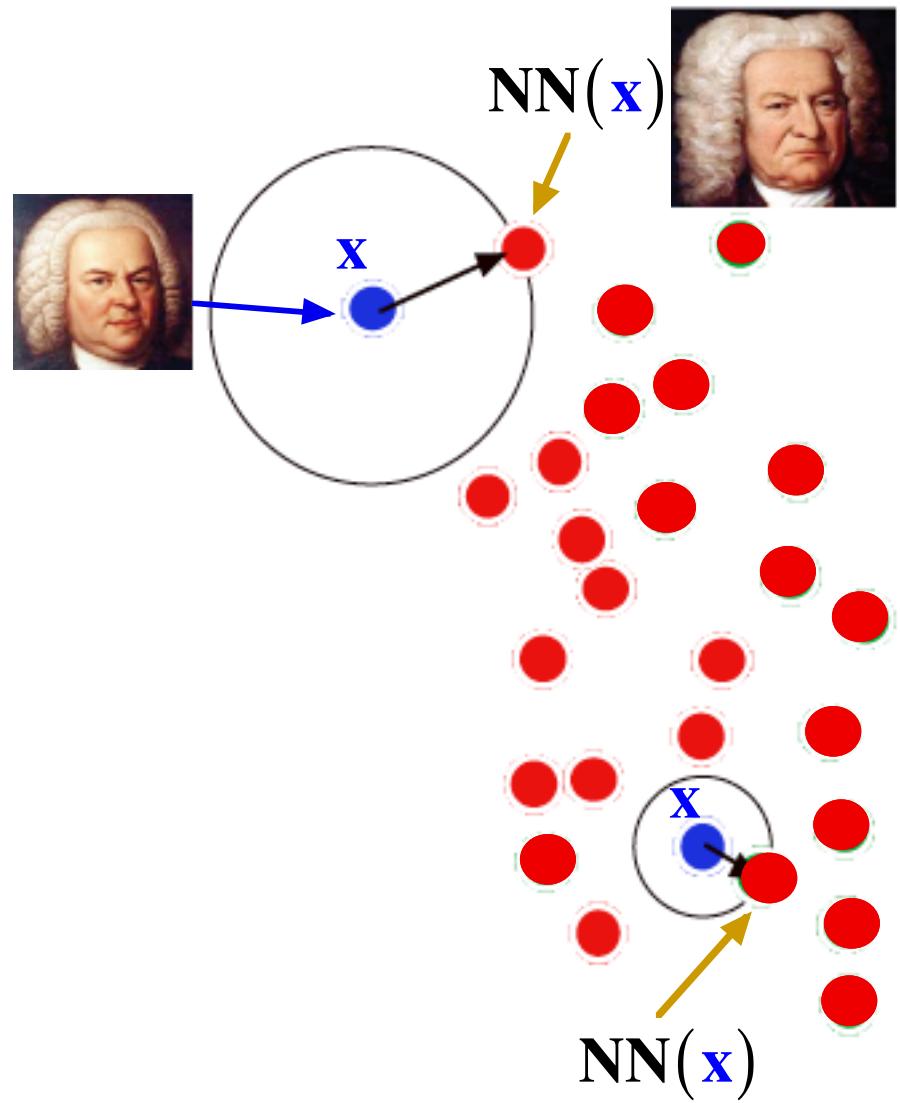
**U of M & AGH**

**Miroslaw.Pawlak@umanitoba.ca**

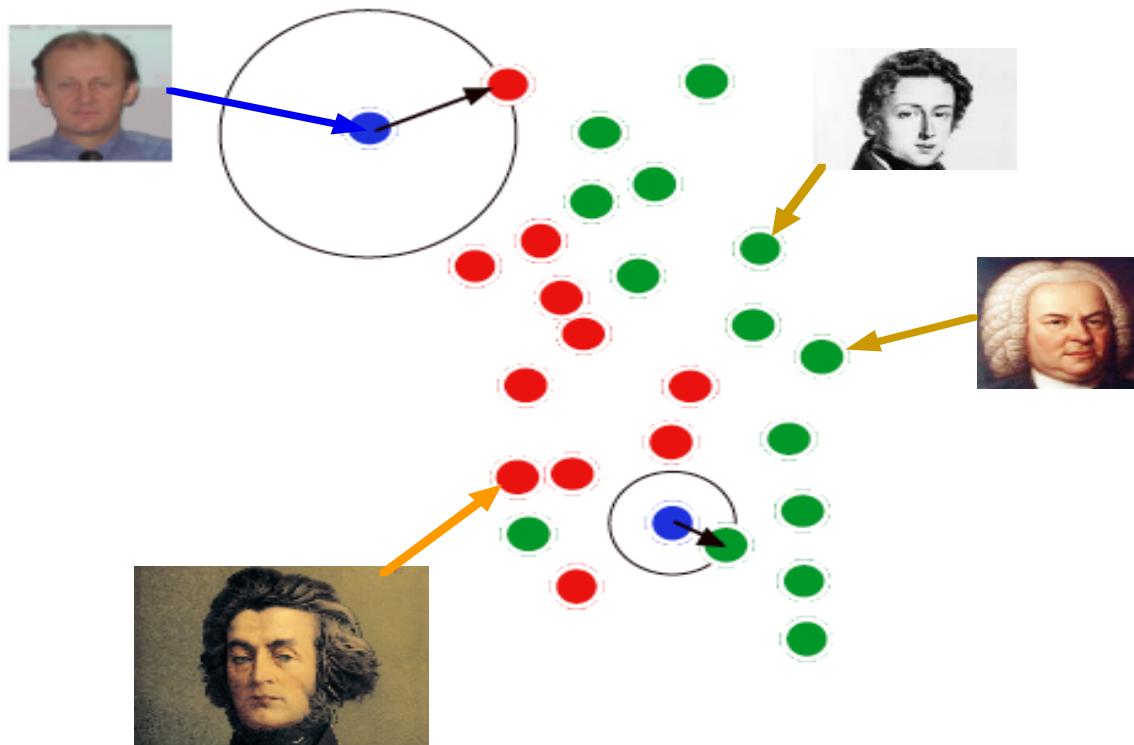
**mirek.pawlakk@gmail.com**







# 1-NN: Classification Problem



$$\mathbf{R}^* \leq \lim_{n \rightarrow \infty} \mathbf{P}(\text{NN error}) \leq 2\mathbf{R}^*$$

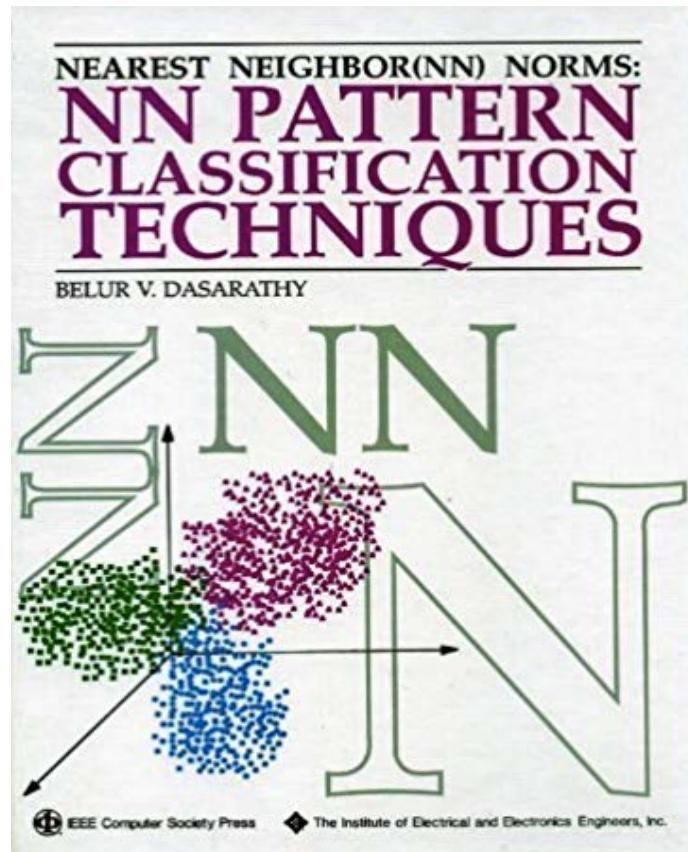
↑

$$\inf_{\text{all rules}} \mathbf{P}(\text{error})$$

## Nearest Neighbor Pattern Classification

T. M. COVER, MEMBER, IEEE, AND P. E. HART, MEMBER, IEEE

**IEEE Inf. Theory, 1967**



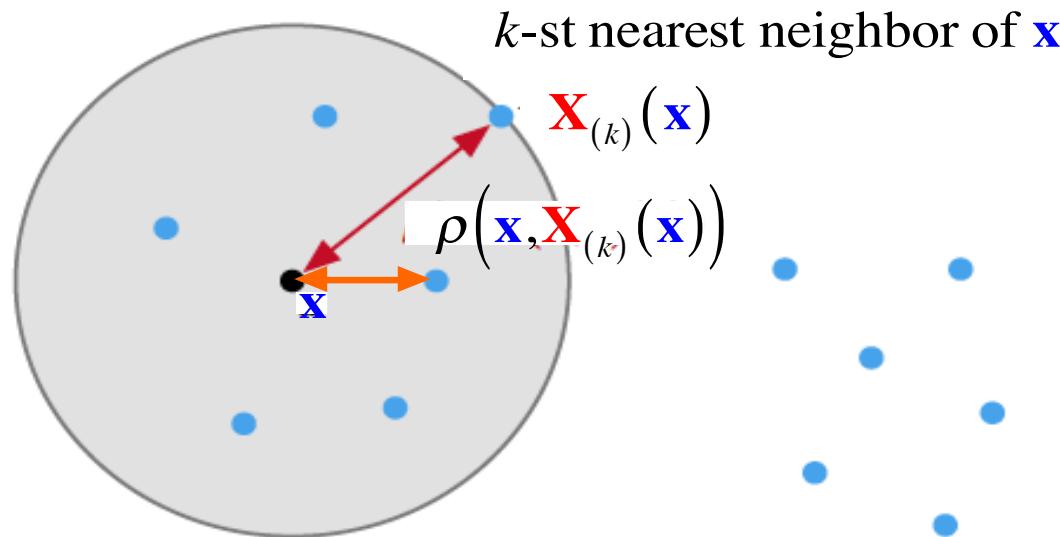
- E. Fix & J.L. Hodges, Discriminatory analysis. Nonparametric discrimination: Consistency properties. Technical Report 4, Project Number 21-49-004, USAF School of Aviation Medicine , Randolph Field, Texas, 1951.



INDY/PULSE

## **HAVING GOOD NEIGHBOURS CAN HELP CUT HEART ATTACK RISK, STUDY SHOWS**

**Estimation methods based on nearest neighbors are the simplest and most intuitively appealing of all nonparametric techniques**



- ➊ Density Estimation
- ➋ Classification
- ➌ Regression Estimation
- ➍ Time Series Analysis

# **OUTLINE**

**I NN Distance**

**II k-NN Regression Estimates: iid case**

**III k-NN Regression Estimates: Time Series Systems**

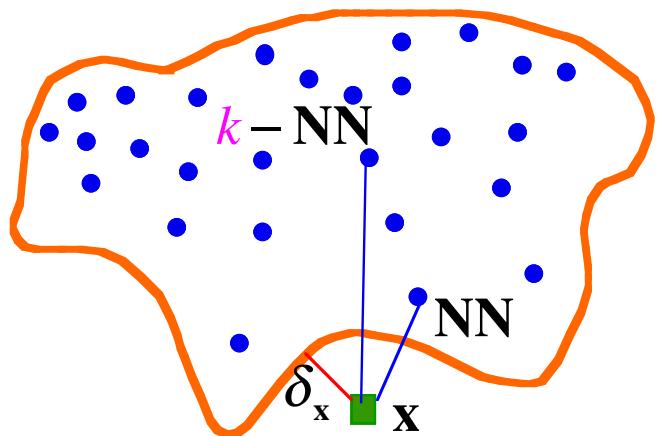
**IV k-NN Regression Estimates: Time Series Forecasting**

**V Concluding Remarks**

## I NN Distance

### Consistency

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \in \mathbf{R}^d \quad \text{iid}$$



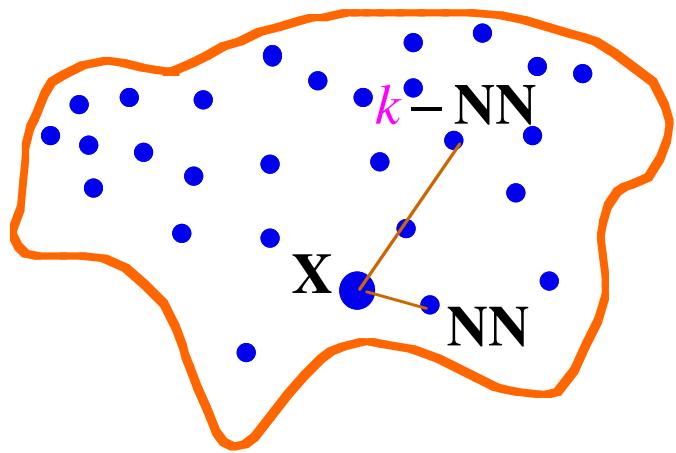
$$\|\mathbf{x} - \mathbf{X}_{(1)}(\mathbf{x})\| \leq \|\mathbf{x} - \mathbf{X}_{(2)}(\mathbf{x})\| \leq \dots \leq \|\mathbf{x} - \mathbf{X}_{(n)}(\mathbf{x})\|$$

If  $\frac{k}{n} \rightarrow 0 \Rightarrow \|\mathbf{x} - \mathbf{X}_{(k)}(\mathbf{x})\| \rightarrow \delta_{\mathbf{x}}$  (P, Pr.1)

►  $k = 1$       ►  $k = \sqrt{n}$       ►  $k = n/2$



## Consistency



$\mathbf{X}$  independent of the data  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$

If  $\frac{k}{n} \rightarrow 0 \Rightarrow \|\mathbf{X} - \mathbf{X}_{(k)}(\mathbf{X})\| \rightarrow 0$  (P, Pr.1)

►  $k = 1$

►  $k = \sqrt{n}$

►  $k = n / 2$



**Rate**  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \in \mathbf{A} \subset R^d; \quad \mathbf{X} \in \mathbf{A}$

↖ **compact set** with diameter = 1

	$E\ \mathbf{X} - \mathbf{X}_{(1)}(\mathbf{X})\ ^2$	$E\ \mathbf{X} - \mathbf{X}_{(k)}(\mathbf{X})\ ^2$
$d = 1$	$\frac{2}{n}$	$8\frac{k}{n}$
$d = 2$	$\frac{7.42...}{n}$	$c(2)\frac{k}{n}$
$d \geq 3$	$c(d)\left(\frac{1}{n}\right)^{2/d}$	$c(d)\left(\frac{k}{n}\right)^{2/d}$

## ▲ General distribution on $\mathbf{A}$



**Rate**  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \in \mathbf{A} \subset R^d$ ;  $\mathbf{X} \in \mathbf{A}$  - compact set

	$E\ \mathbf{X} - \mathbf{X}_{(1)}(\mathbf{X})\ ^2$	$E\ \mathbf{X} - \mathbf{X}_{(k)}(\mathbf{X})\ ^2$
$d = 1$	$c(1) \frac{1}{n^2}$ *	$c(1) \frac{k}{n^2}$ *
$d = 2$	$c(2) \frac{1}{n}$	$c(2) \frac{k}{n}$
$d \geq 3$	$c(d) \left(\frac{1}{n}\right)^{2/d}$	$c(d) \left(\frac{k}{n}\right)^{2/d}$

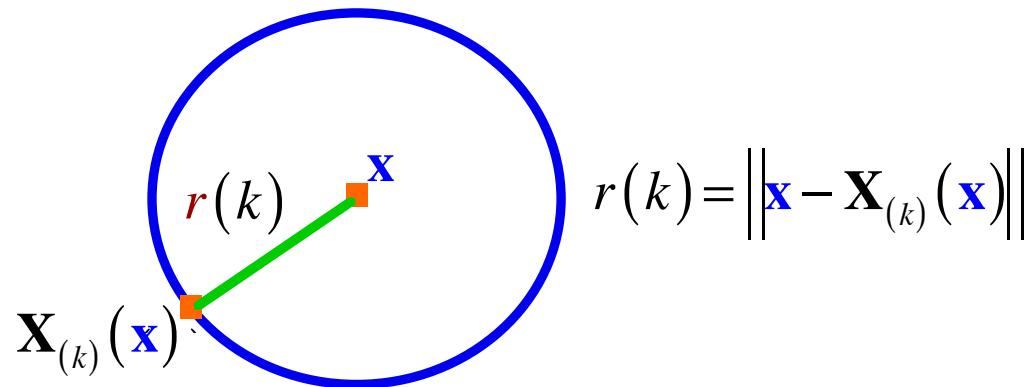


**Distribution with a density defined on A**



## Coverage Probability - Pr. Integral Transform

$$\int_{\mathbf{B}(\mathbf{x}; r(k))} f(\mathbf{v}) d\mathbf{v} \sim \text{Beta}(k, n-k+1)$$



- $\mathbf{E}\left[\int_{\mathbf{B}(\mathbf{x}; r(k))} f(\mathbf{v}) d\mathbf{v}\right] = \frac{k}{n+1}$
- $\mathbf{Var}\left[\int_{\mathbf{B}(\mathbf{x}; r(k))} f(\mathbf{v}) d\mathbf{v}\right] = \frac{k(n-k+1)}{(n+1)^2(n+2)}$

• **LLN**       $k \rightarrow \infty \Rightarrow \frac{\int_{\mathbf{B}(\mathbf{x}; r(k))} f(\mathbf{v}) d\mathbf{v}}{k/n} \rightarrow 1 \quad (\mathbf{P})$

## II k-NN Regression Estimates: iid data



### Regression

- $$\underbrace{\mathbf{E}[Y \mid \mathbf{X} = \mathbf{x}]}_{m(\mathbf{x})} = \arg \min_{g: \mathbf{R}^d \rightarrow \mathbf{R}} \mathbf{E}[Y - g(\mathbf{X})]^2$$
- $$\mathbf{E}[Y - g(\mathbf{X})]^2 = \underbrace{\mathbf{E}\left\{\|\mathbf{Y} - m(\mathbf{X})\|^2\right\}}_{\text{Bayes Risk}} + \mathbf{E}\left\{\|g(\mathbf{X}) - m(\mathbf{X})\|^2\right\}$$
- $$Y = m(\mathbf{X}) + \varepsilon \Rightarrow \mathbf{E}[Y - g(\mathbf{X})]^2 = \text{var}[\varepsilon] + \mathbf{E}\left\{\|g(\mathbf{X}) - m(\mathbf{X})\|^2\right\}$$



## k-NN Estimates

- Wish to estimate  $m(\mathbf{x}) = \mathbb{E}[Y | \mathbf{X} = \mathbf{x}]$  from

$$\mathfrak{I}_n = \left\{ (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \right\}; \quad (\mathbf{X}_i, Y_i) \in \mathbf{R}^d \times \mathbf{R}$$

- Ordered data set

$$\mathfrak{I}_n(\mathbf{x}) = \left\{ (\mathbf{X}_{(1)}, Y_{[1]}), \dots, (\mathbf{X}_{(n)}, Y_{[n]}) \right\} \leftarrow \mathbf{x} \in \mathbf{R}^d$$

$$\left\| \mathbf{X}_{(1)} - \mathbf{x} \right\| \leq \left\| \mathbf{X}_{(2)} - \mathbf{x} \right\| \leq \dots \leq \left\| \mathbf{X}_{(n)} - \mathbf{x} \right\|$$

$\downarrow$                      $\downarrow$                      $\downarrow$   
 $Y_{[1]}$                      $Y_{[2]}$                      $Y_{[n]}$       ← **concomitants**



## k-NN Estimates: iid data



### 1-NN Estimate

$$\tilde{m}_{NN}(\mathbf{x}) = Y_{[1]}(\mathbf{x})$$

$$Y = m(\mathbf{X}) + \varepsilon$$

$$Y_{[1]}(\mathbf{x}) = m(\mathbf{X}_{(1)}(\mathbf{x})) + \varepsilon_{[1]}$$

$$E \left[ m(\mathbf{X}_{(1)}(\mathbf{X})) - m(\mathbf{X}) \right]^2 \rightarrow 0$$



## 1-NN Estimate

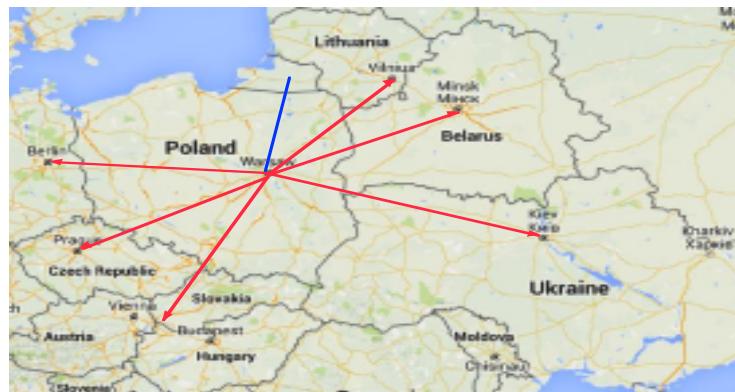
**Non-consistent estimate of  $m(\textcolor{blue}{u})$**

$$\tilde{\mathbf{m}}_{NN}(\mathbf{x}) \rightarrow m(\mathbf{x}) + \sqrt{\underbrace{\mathbf{E}\left\{\|\mathbf{Y} - m(\mathbf{X})\|^2\right\}}_{\text{var of noise}}} \quad (\mathbf{P})$$

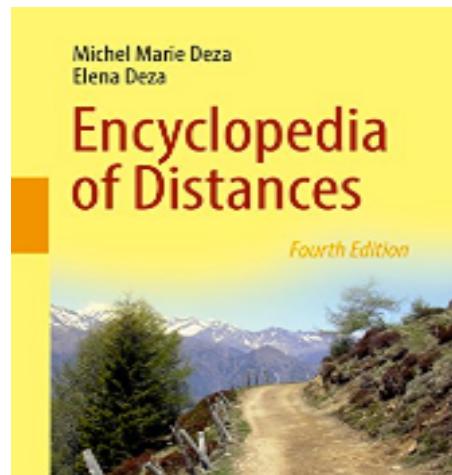
$$\mathbf{E}\left[\tilde{\mathbf{m}}_{NN}(\mathbf{X}) - m(\mathbf{X})\right]^2 \rightarrow \underbrace{\mathbf{E}\left[\|\mathbf{Y} - m(\mathbf{X})\|^2\right]}_{\sigma_e^2}$$

- The Standard k-NN Estimate:  $\tilde{m}(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k Y_{[i]}$

- $k$  — number of neighbors



- distance "metric"



□ **The Standard k-NN Estimate:**  $\tilde{m}(\mathbf{x}) = \frac{1}{k_n} \sum_{i=1}^{k_n} Y_{[i]}$

- $k_n \rightarrow \infty$  - controls variance

$\Leftarrow$  kernel method  $(h_n \rightarrow 0, nh_n^d \rightarrow \infty)$

- $\frac{k_n}{n} \rightarrow 0$  - controls bias

$$\hat{m}(\mathbf{x}) = \frac{\sum_{i=1}^n Y_i K\left(\frac{\|\mathbf{x} - \mathbf{X}_i\|}{h}\right)}{\sum_{i=1}^n K\left(\frac{\|\mathbf{x} - \mathbf{X}_i\|}{h}\right)}$$

□ **The Standard k-NN Estimate: Universal Consistency**

**Theorem 1** Let  $E|Y|^2 < \infty$ . Then

$$E\left\{\left|\tilde{m}(\mathbf{X}) - m(\mathbf{X})\right|^2\right\} \rightarrow 0 \text{ for all distributions of } (\mathbf{X}, Y)$$

iff

$$k_n \rightarrow \infty \text{ and } \frac{k_n}{n} \rightarrow 0$$

► Stone: Annals of Stat., 1977

► Distance tie-breaking      ► Geometric Proof  
Covering Numbers



## The Standard k-NN Estimate: Rate

$\Leftarrow d = 1$

**Theorem 2** Let  $m(\bullet) \in \text{Lip}(\mathbf{R})$  and

$$\tilde{k}_n^* = \lceil c_1 n^{2/3} \rceil .$$

Then

$$\mathbf{E} \left[ \tilde{m}(X) - m(X) \right]^2 = O(n^{-2/3})$$

This rate **cannot** be improved for smoother characteristics

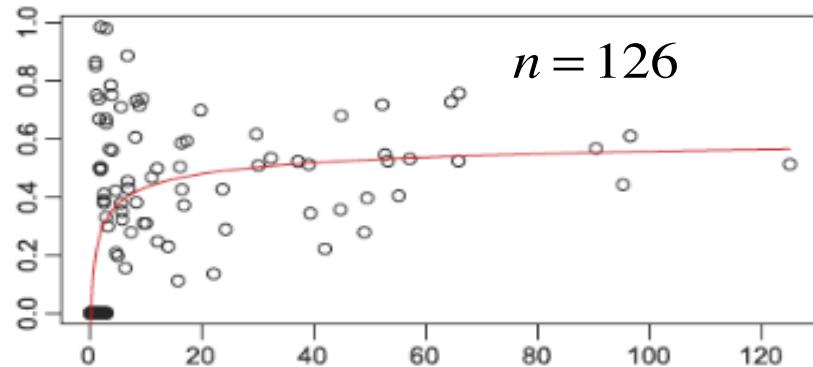
$$O(n^{-2/3}) \rightarrow O(n^{-4/5}) \rightarrow O(n^{-4/5}) \rightarrow \dots$$

- **The Standard k-NN Estimate: Useful Properties**
- The use of  $k$ -nearest neighbors weights is convenient in a number of nonparametric estimation applications because they avoid the **random denominator problem** intrinsic to commonly used **kernel weights**.

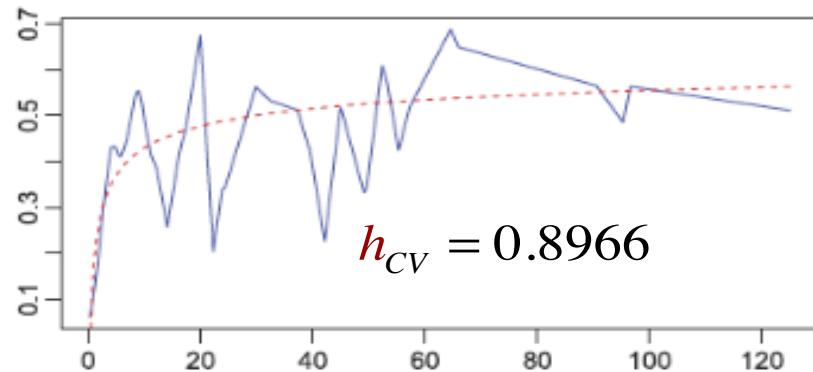
$$\tilde{m}(\mathbf{x}) = \frac{1}{k_n} \sum_{i=1}^{k_n} Y_{[i]} \leftrightarrow \hat{m}(\mathbf{x}) = \frac{\sum_{i=1}^n Y_i K\left(\frac{\|\mathbf{x} - \mathbf{X}_i\|}{h}\right)}{\sum_{i=1}^n K\left(\frac{\|\mathbf{x} - \mathbf{X}_i\|}{h}\right)}$$



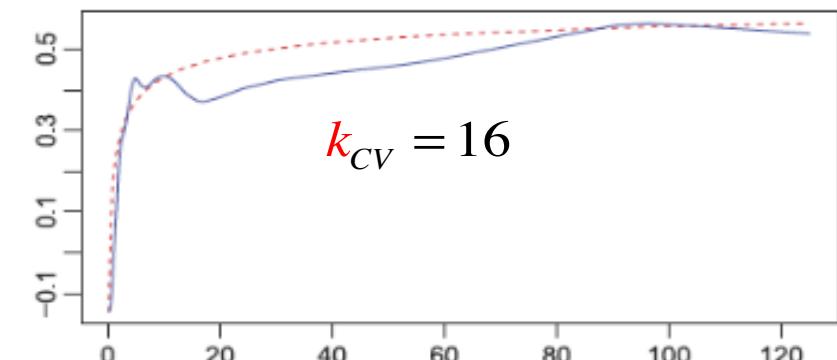
The  $k$ -NN estimate adapts to the non-homogeneity of data



LS for  $m(x; \theta) = \theta_1 + \theta_2 x^{\theta_3}$



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

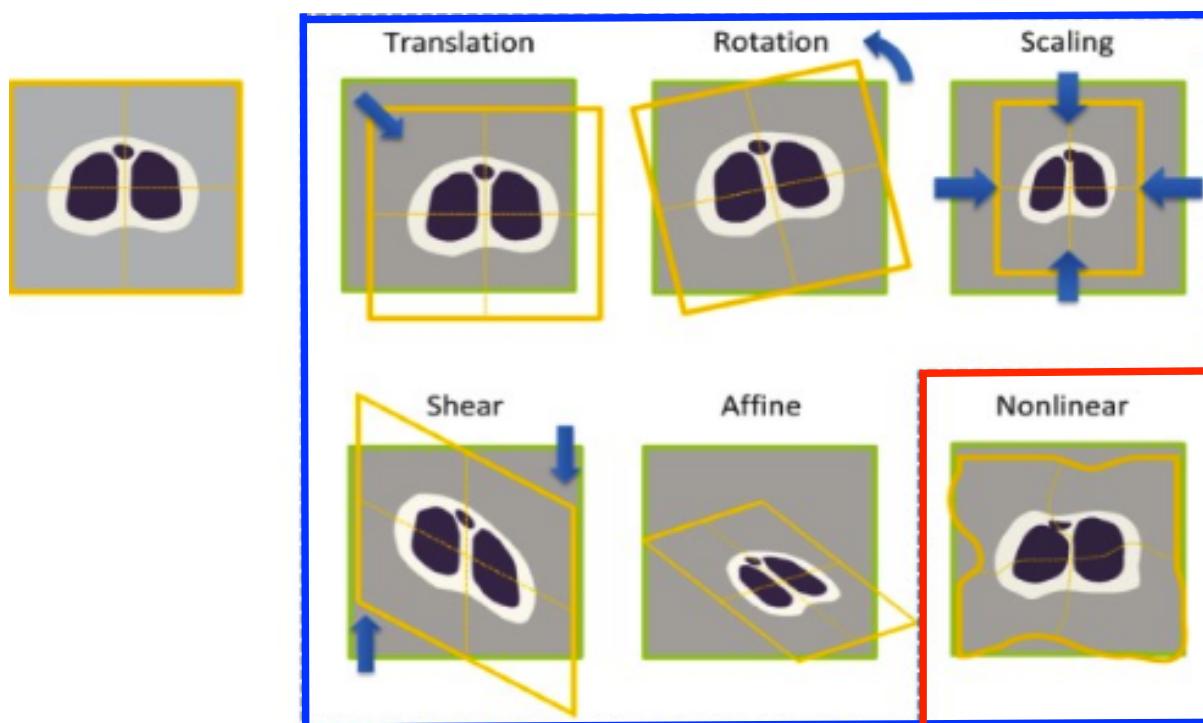


$$\frac{1}{k} \sum_{i=1}^k Y_{[i]}$$

- The  $k$ -NN estimate can incorporate invariance

$$\tilde{m}((\mathbf{Ax} + \mathbf{b}); (\mathbf{Ax} + \mathbf{b})\mathfrak{S}_n) = \tilde{m}(\mathbf{x}; \mathfrak{S}_n)$$

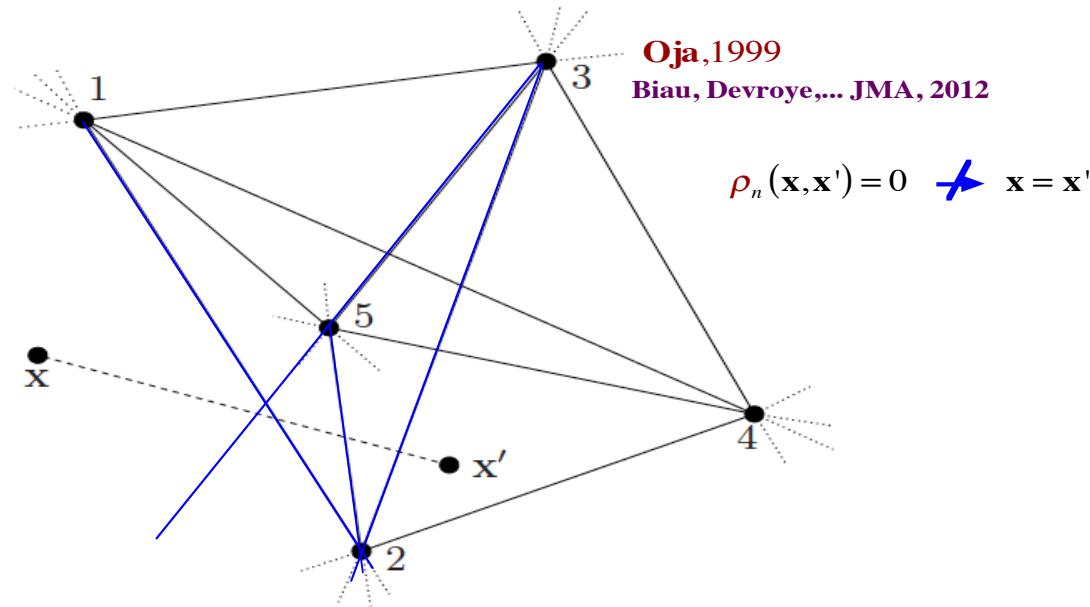
linear



9  $\longleftrightarrow ?$  6



## An Affine Invariant Distance



Oja, 1999  
Biau, Devroye,... JMA, 2012

$$\rho_n(\mathbf{x}, \mathbf{x}') = 0 \quad \not\rightarrow \quad \mathbf{x} = \mathbf{x}'$$

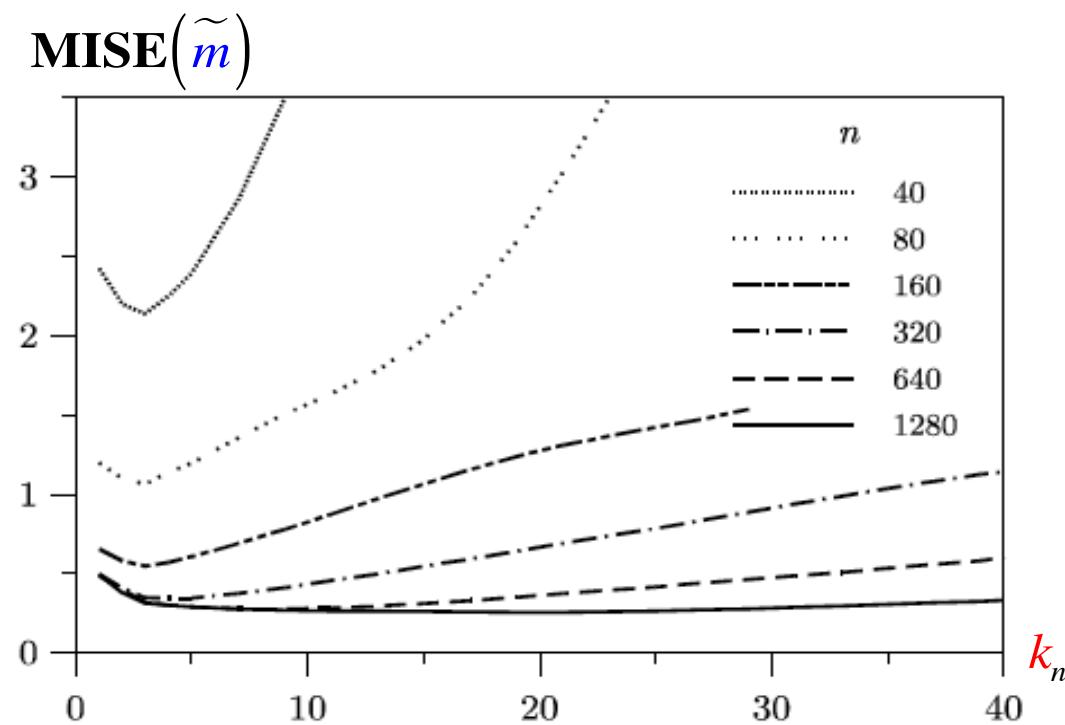
$\rho_n(\mathbf{x}, \mathbf{x}') = \# \text{ of hyperplanes passing through 2 out of the points}$

$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  that are separating  $\mathbf{x}$  and  $\mathbf{x}'$

$= \# \text{ of hyperplanes in } \mathbf{R}^d \text{ passing through } d \text{ out of the points}$   
 $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  that are separating  $\mathbf{x}$  and  $\mathbf{x}'$

- The **k -NN** estimate can be extended to infinite dimensional spaces
  - Separable Banach spaces
  - Hilbert spaces with a countable orthonormal basis
  - Reproducing Kernel Hilbert Space
- Biau,...; IEEE IT, 2010

- Choice of  $k_n$



► CV choice for the k-NN estimate: Li, Annals of Stat. 84

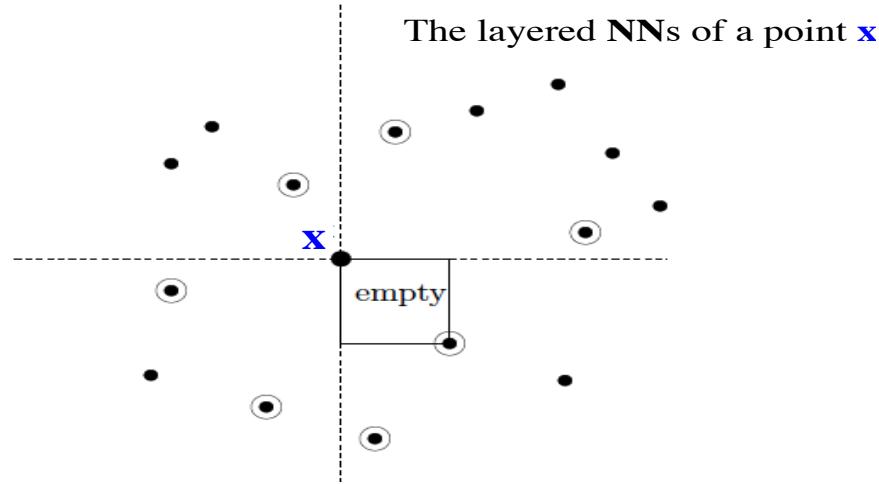
## ▀ Other k-NN Estimates

### □ Weighted k-NN Estimates

$$\hat{m}(\mathbf{x}) = \sum_{i=1}^{k_n} w_{ni} Y_{[i]}$$
$$\sum_{i=1}^{k_n} w_{ni} = 1$$

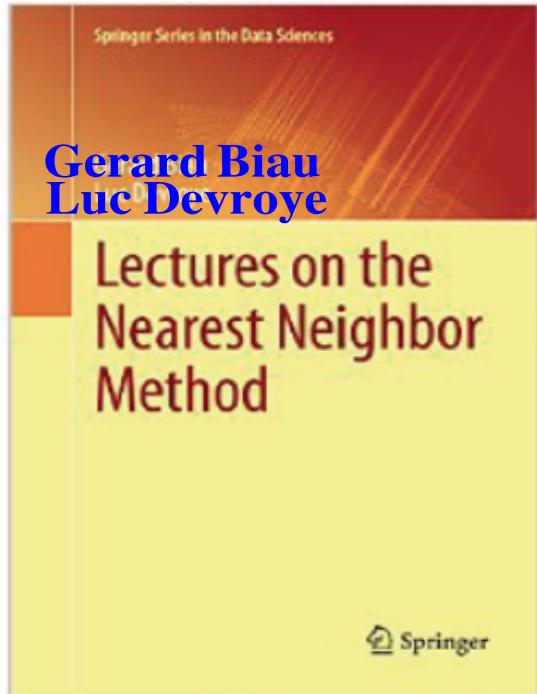
Stone's Theorem (**Stone: Annals of Stat., 1977**) gives general **iff** conditions on the weights in order to obtain the **universal  $L_p$  - consistency**.

## □ Layered k-NN estimate



$$\tilde{m}(\mathbf{x}) = \frac{1}{|l_n(\mathbf{x})|} \sum_{\{\mathbf{x}_i \in l_n(\mathbf{x})\}} Y_i$$

► Biau, Devroye, 2010



← iid data

### **III k-NN Regression Estimates: Dependent Data**

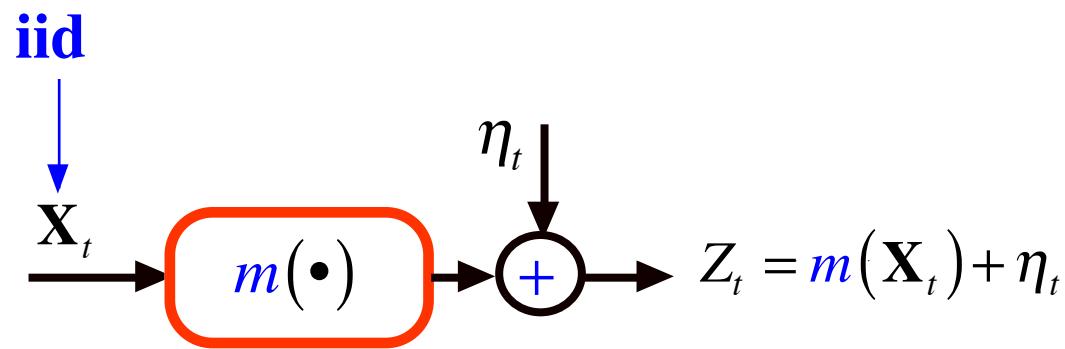
**S. Yakowitz.** Nearest-neighbour methods for time series.  
J. of Time Series, 1987.

**S. Kulkarni and S. Posner.** Rates of convergence of nearest neighbour estimation under arbitrary sampling.  
IEEE Trans. on Inf. Theory, 1995.

**G. Boente and R. Fraiman.** Robust nonparametric regression estimation for dependent observations. The Annals of Statistics, 1989.



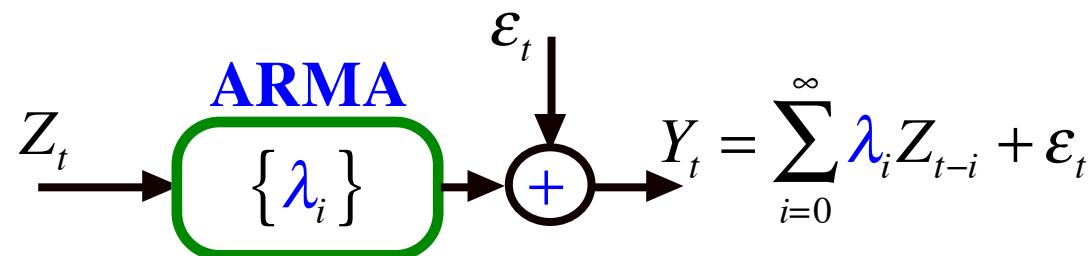
## Nonlinear Regression



$$m(\mathbf{x}) = \mathbf{E}[Z_t | \mathbf{X}_t = \mathbf{x}]$$



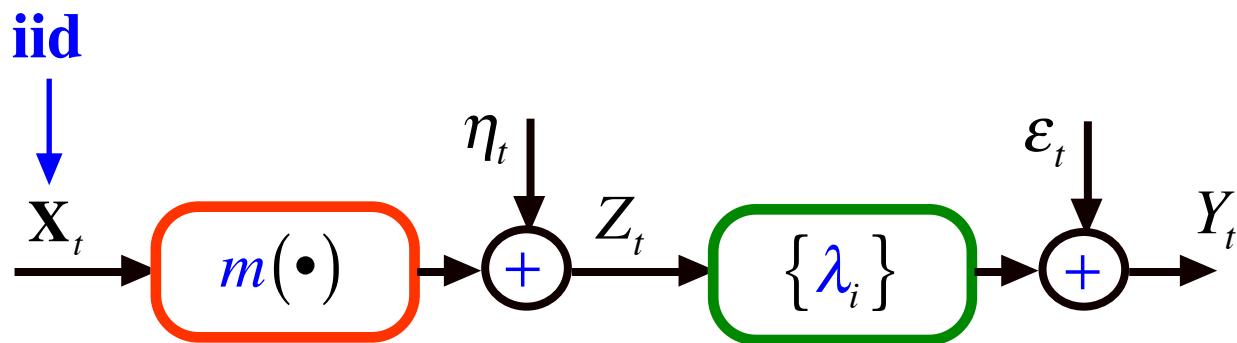
## Autoregressive (Linear) Process



$$\lambda_i = \text{cov}\{Y_t, Z_{t-i}\}$$



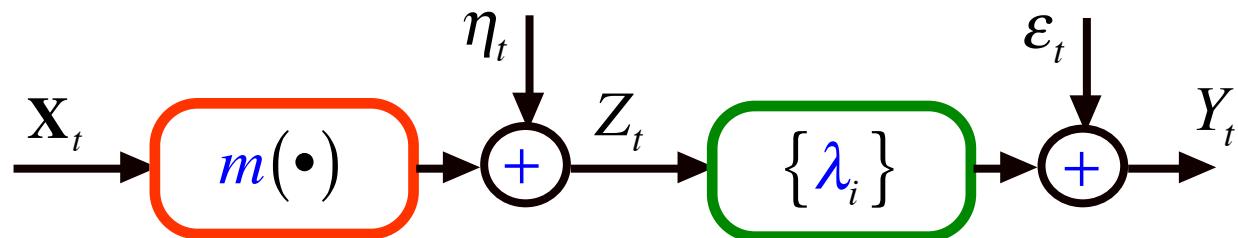
## Nonlinear Time-Series System: Output Dependence Model



$$Y_t = \sum_{i=0}^{\infty} \lambda_i m(X_{t-i}) + \sum_{i=0}^{\infty} \lambda_i \eta_{t-i} + \varepsilon_t$$



## Nonlinear System = Hammerstein System



$$x(t) + \int_0^T \underbrace{\Lambda}_{known} (s,t) \underbrace{m}_{known} (x(s)) ds = \underbrace{y}_{known} (t)$$

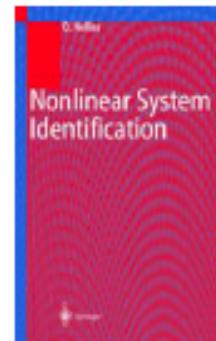
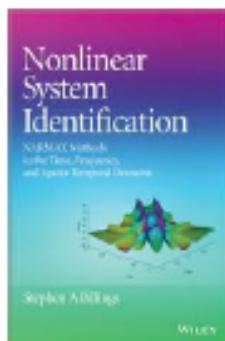
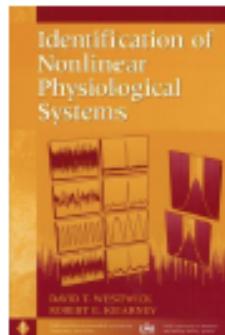
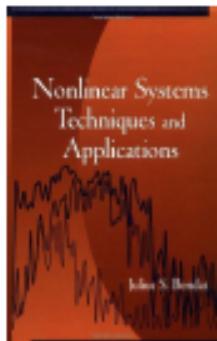
A. Hammerstein: Nichtlineare integralgleichungen nebst anwendungen,

Acta Mathematica, 1930.

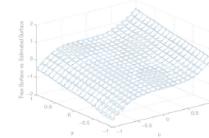


## Nonlinear System Identification

- hot research area with numerous applications
- mostly parametric methods
- block-oriented models are the most popular
- lack of the basic statistical theory ←



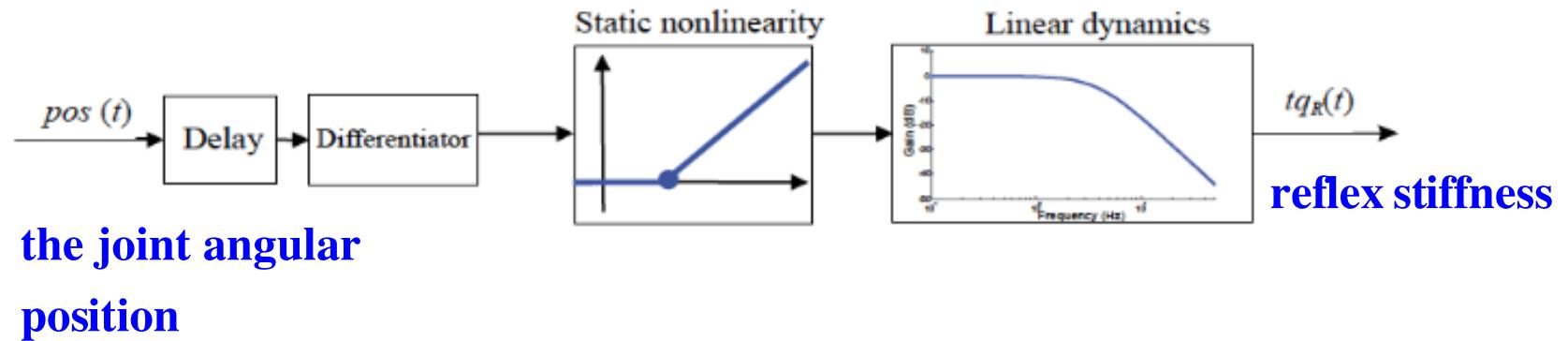
Recursive  
Identification and  
Parameter Estimation



Han-Fu Chen  
Wenxiao Zhao

## • Biomedical Engineering

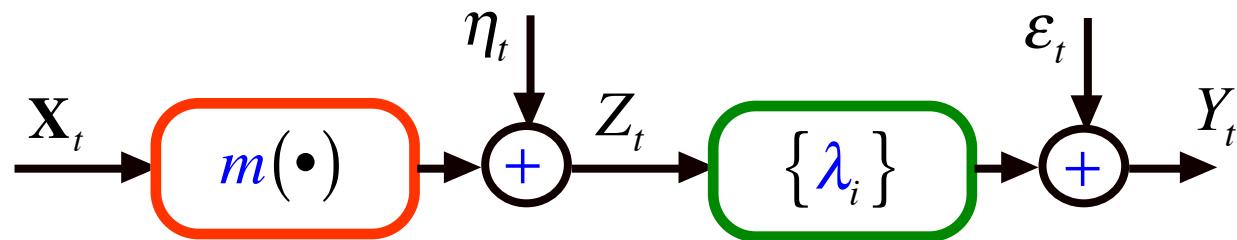
### Reflex Stiffness Pathway



- Kearny, ... : Identification of intrinsic and reflex contributions to human ankle stiffness dynamics," IEEE Trans. on Biomedical Engineering, 1997
- This **nonlinear model** has been widely used to characterize the **reflex stiffness** in the ankle and trunk joints of both normal and pathological subjects.



## Nonlinear System and Regression:



$$Y_t = \textcolor{blue}{m}(\mathbf{X}_t) + \textcolor{green}{e}_t \quad \leftarrow \quad \lambda_0 = 1$$

$$\textcolor{green}{e}_t = \sum_{i=1}^{\infty} \lambda_i \textcolor{blue}{m}(\mathbf{X}_{t-i}) + \sum_{i=0}^{\infty} \lambda_i \eta_{t-i} + \varepsilon_t$$

$$\text{cov}[\textcolor{green}{e}_t; \textcolor{green}{e}_{t+k}] = \text{var}[\textcolor{blue}{m}(\mathbf{X})] \sum_{i=1}^{\infty} \lambda_i \lambda_{i+k} + \sigma_{\eta}^2 \sum_{i=0}^{\infty} \lambda_i \lambda_{i+k} + \sigma_{\varepsilon}^2 \mathbf{1}(k=0)$$

$$\lambda_0 = 1$$

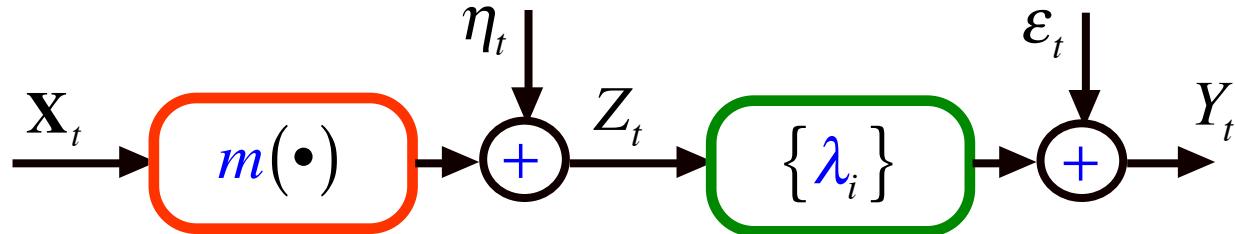
$$E\left[ \textcolor{blue}{m}\left(\mathbf{X}_t\right) \right] = 0 \quad \Rightarrow \quad \textcolor{red}{c} = 0$$



$$E\left\{ Y_t \mid \mathbf{X}_t = \mathbf{x} \right\} = \textcolor{blue}{m}\left(\mathbf{x}\right)$$



## Problem Statement



Given  $\mathfrak{S}_n = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$  estimate  $m(\mathbf{x})$

- A1** •  $\{\mathbf{X}_t\}$  - **iid** with pdf  $f(\cdot)$
- A2** •  $\bullet \mathbb{E}[m^2(\mathbf{X}_t)] < \infty, \quad \bullet \sum_{i=0}^{\infty} |\lambda_i| < \infty$
- A3** •  $\{\eta_t\}, \{\varepsilon_t\}$  - **iid**, zero mean, finite variance



## The System $\leftrightarrow$ Nonparametric Inference $\leftrightarrow$ Dependence Structure

*The Annals of Statistics*

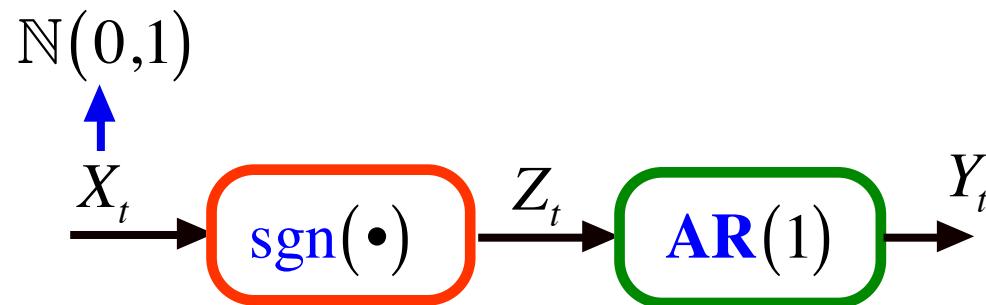
equivariant estimators are considered: (i) estimators based on kernel methods and (ii) estimators based on  $k$ -nearest neighbor kernel methods. Strong consistency of both families is proved under mild conditions.

We obtain strong pointwise consistency of both families of estimates, when the observations  $(X_t, Y_t)$  have a  $\varphi$ - or an  $\alpha$ -mixing dependence structure [see Billingsley (1968) and Rosenblatt (1956), respectively]. For kernel weights, the

In this paper we go far beyond this classical framework by showing that support vector machines (SVMs) essentially only require that the data-generating process satisfies a certain law of large numbers. We then consider the learnability of SVMs for  $\alpha$ -mixing (not necessarily stationary) processes for both classification and regression,....

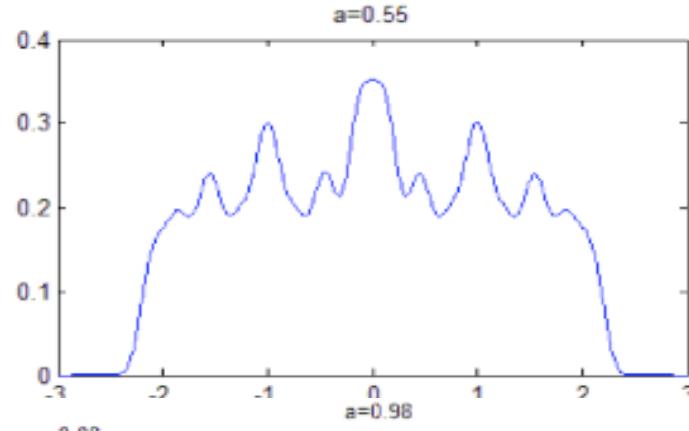
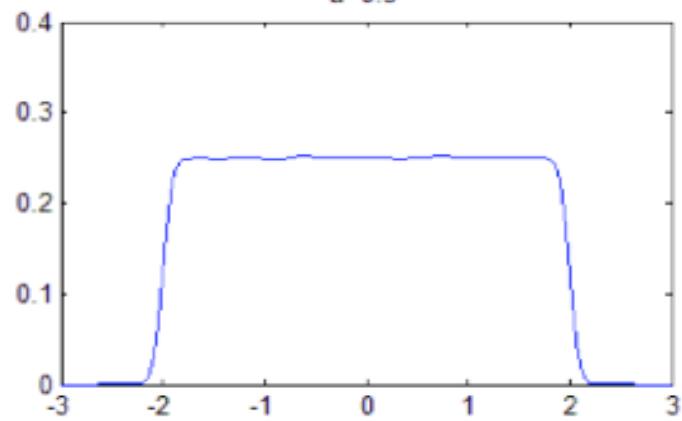
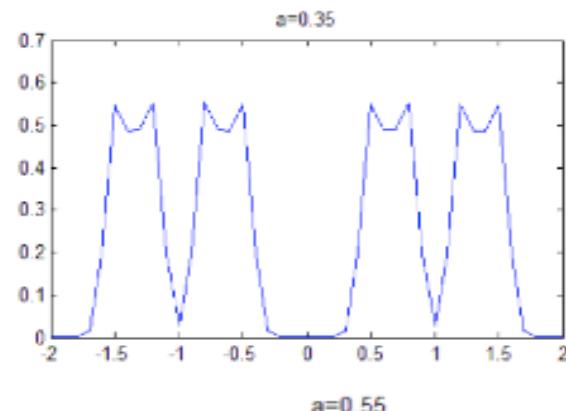
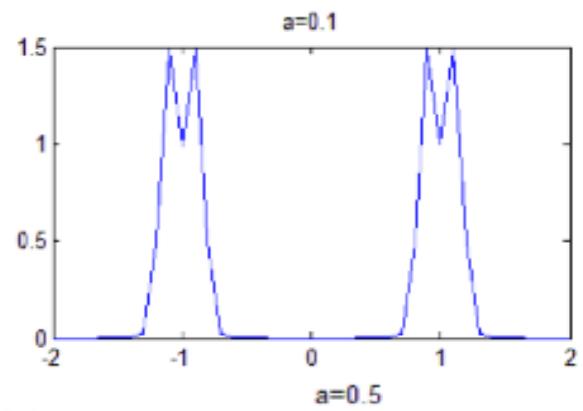


## The System $\leftrightarrow$ Dependence Structure

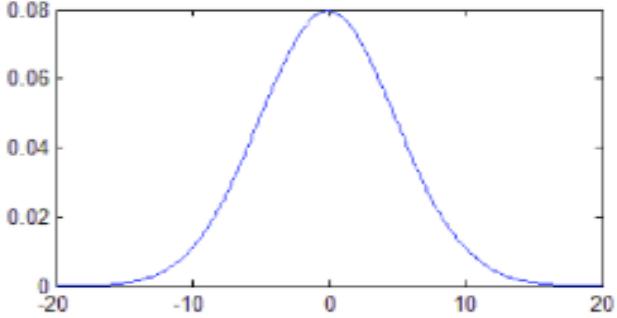


$$Y_t = \underbrace{\color{blue}{a} Y_{t-1} + \color{blue}{\text{sgn}}(X_t)}_{\pm 1 \in \text{pr}.1/2}$$

$$Y_t = \sum_{j=0}^{\infty} \color{blue}{a}^j (\pm 1)$$



$$\hat{f}_Y(y) = (nh)^{-1} \sum_{t=1}^n \mathbf{K}\left(\frac{y - Y_t}{h}\right)$$



$$Y_t = \sum_{j=0}^{\infty} \textcolor{blue}{a}^j (\pm 1) - \textbf{Bernoulli Convolution}$$

- **Paul Erdos.** On a family of symmetric Bernoulli convolutions. Amer. J. Math., 1939.
- **Boris Solomyak.** On the random series  $\sum \pm \lambda^n$  (an Erdos problem). Ann.of Math., 1995.
- **Peres, Schalg and Solomyak.** Sixty years of Bernoulli convolution, 2000.

$$Y_t = \sum_{j=0}^{\infty} a^j (\pm 1)$$

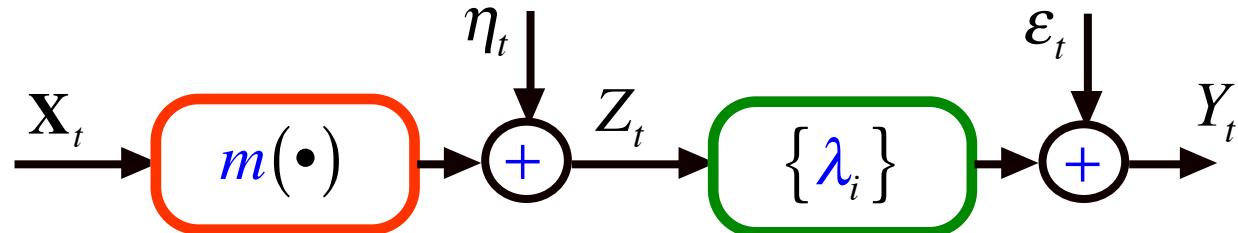
- $\{Y_t\}$  is **not mixing** for  $a \in (0, 1/2]$

**D. W. K. Andrews.** Non-strong mixing autoregressive processes.  
J.Applied Probability, 1984.

- $a \in (0, 1/2) \Rightarrow F_Y$  - **singular**
- $a = 1/2 \Rightarrow F_Y$  - **Uniform**[-2,2]
- $a \in (1/2, 1) \Rightarrow F_Y$  - either **absolutely cont.** or **singular** ?



## k-NN Estimates



Wish to estimate

$$m(\mathbf{x}) = E[Y_t | \mathbf{X}_t = \mathbf{x}]$$

from

$$\mathfrak{S}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$$



## Weighted k-NN Estimates

$$\hat{m}(\mathbf{x}) = \sum_{i=1}^{k_n} w_{ni} Y_{[i]}$$

$$\sum_{i=1}^{k_n} w_{ni} = 1$$



## k-NN Estimates for Dependent Data

**S. Yakowitz.** Nearest-neighbour methods for time series. *J. of Time Series*, 1987.

**S. Kulkarni and S. Posner.** Rates of convergence of nearest neighbour estimation under arbitrary sampling. *IEEE Trans. on Inf. Theory*, 1995.

**G. Boente and R. Fraiman.** Robust nonparametric regression estimation for dependent observations. *The Annals of Statistics*, 1989.



## Asymptotic Theory

$$\hat{m}(\mathbf{x}) = \sum_{i=1}^{k_n} w_{ni} Y_{[i]}$$

$$\sum_{i=1}^{k_n} w_{ni} = 1, \quad w_{ni} \geq 0$$

□ **Theorem 3 (Convergence)**

$$k_n \rightarrow \infty, \quad \sum_{i=1}^{k_n} w_{ni}^2 \rightarrow 0 \quad \Leftarrow \text{variance}$$

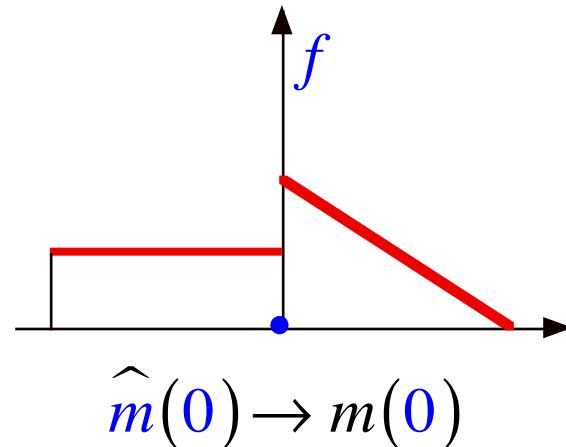
$$\frac{k_n}{n} \rightarrow 0 \quad \Leftarrow \text{Bias}$$

$$\hat{m}(\mathbf{x}) \rightarrow m(\mathbf{x}) \quad (\mathbf{P})$$

at every point  $\mathbf{x} \in S(f)$  at which  $m(\mathbf{x})$  is **continuous**

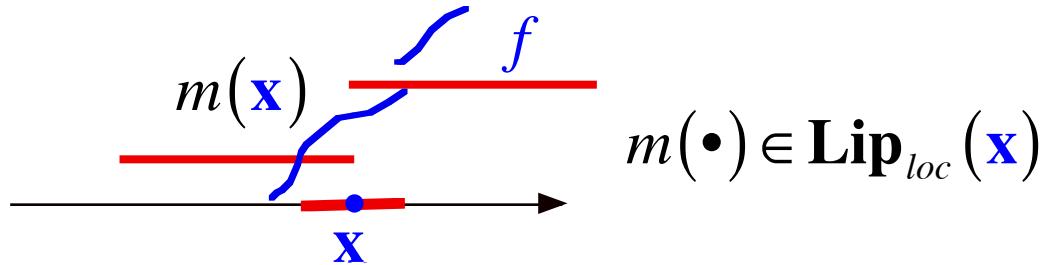
$$\hat{m}(\mathbf{x}) \rightarrow m(\mathbf{x}) \quad (\mathbf{P})$$

at every point  $\mathbf{x} \in S(f)$  at which  $m(\mathbf{x})$  is **continuous**



## Proof: Ranks

□ **Theorem 4 (MSE)**



$m(\bullet) \in \text{Lip}_{loc}(\mathbf{x})$

$$\begin{aligned} \mathbb{E}[\hat{m}(\mathbf{x}) - m(\mathbf{x})]^2 \\ \doteq \frac{\mathbf{D}}{n^2} \left[ \sum_{i=1}^{k_n} i w_{ni} \right]^2 &\quad \Leftarrow \text{Bias } (d=1) \end{aligned}$$

$$\begin{aligned} + \sum_{i=1}^{k_n} w_{ni}^2 \\ + O(n^{-1}) &\quad \Leftarrow \text{Variance} \end{aligned}$$

## Proof:

- $r(k) = \left\| \mathbf{X}_{(k)}(\mathbf{x}) - \mathbf{x} \right\|^2$
- $\mathbf{B}(k) = \int_{\mathbf{B}_{r(k)}(\mathbf{x})} f(\mathbf{v}) d\mathbf{v} \sim \text{Beta}(k, n-k+1)$
- $\sum_{i=1}^{k_n} \sum_{j=1}^{k_n} w_{ni} w_{nj} \mathbf{E}[\mathbf{B}(i)\mathbf{B}(j)]$



Optimal Estimates  $\leftarrow$  Optimal Weights +  $k_n^*$

$$E[\hat{m}(\mathbf{x}) - m(\mathbf{x})]^2$$

$$\simeq \frac{\mathbf{D}}{n^2} \left[ \sum_{i=1}^{k_n} i w_{ni} \right]^2 + \sum_{i=1}^{k_n} w_{ni}^2 \xrightarrow{\text{minimize}} \min_{\{w_{ni}\}, k_n} \text{s.t. } \sum_{j=1}^{k_n} w_{nj} = 1 \\ w_{ni} \geq 0$$

□ **Theorem 5 (Optimal Weights)**

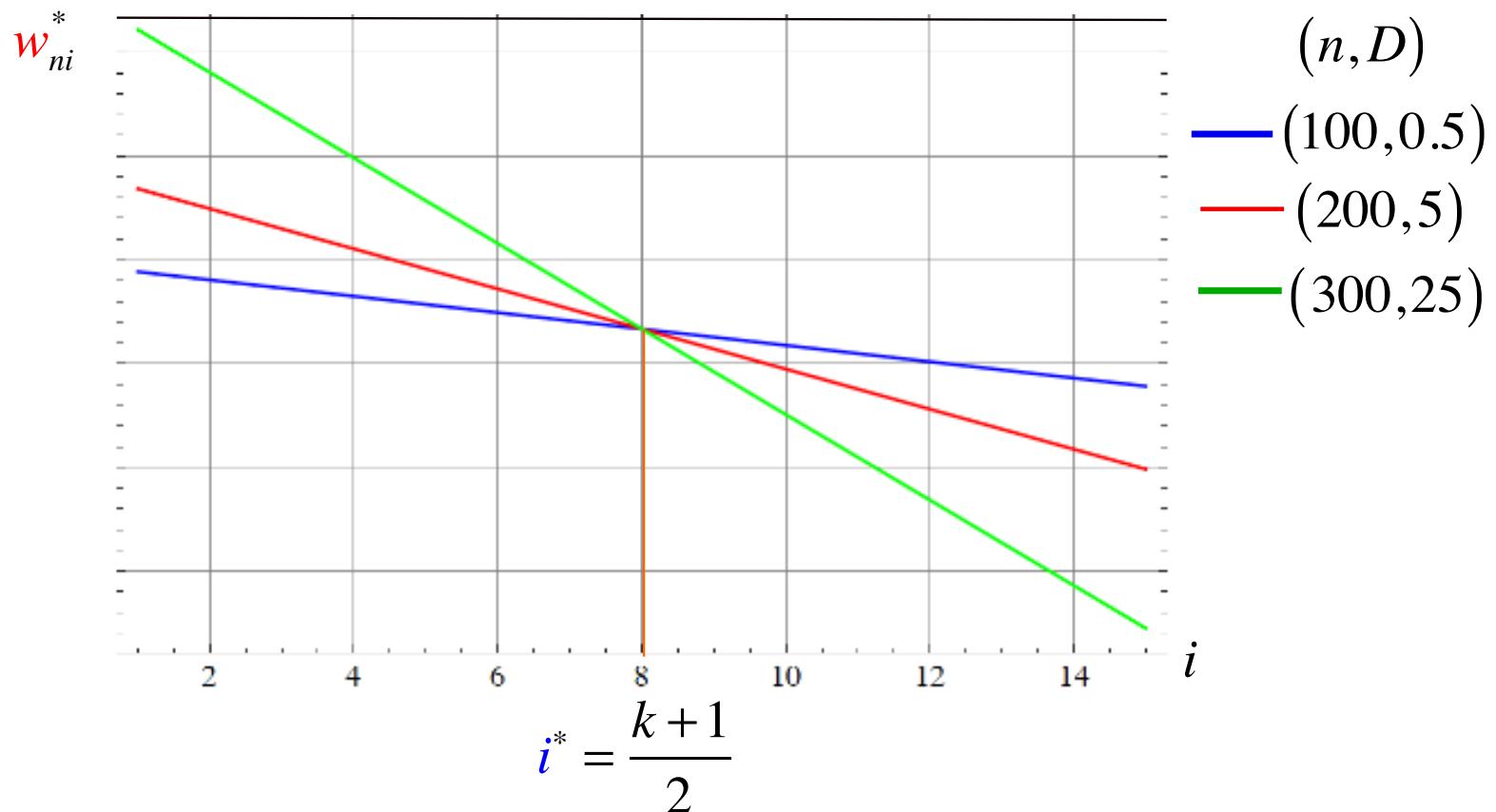
$$w_{ni}^* = \textcolor{red}{a}_n^* - i\textcolor{red}{b}_n^*, \quad i = 1, \dots, \textcolor{red}{k}_n^*$$

$$\textcolor{red}{a}_n^* = \frac{1}{\textcolor{red}{k}_n^*} + \frac{\mathbf{D}}{2} (\textcolor{red}{k}_n^* + 1) \Delta_n^* \quad , \quad \textcolor{red}{b}_n^* = \mathbf{D} \Delta_n^*$$

$$\Delta_n^* = \frac{6(\textcolor{red}{k}_n^* + 1)}{12n^2 + \mathbf{D}(\textcolor{red}{k}_n^* - 1)\textcolor{red}{k}_n^*(\textcolor{red}{k}_n^* + 1)}$$

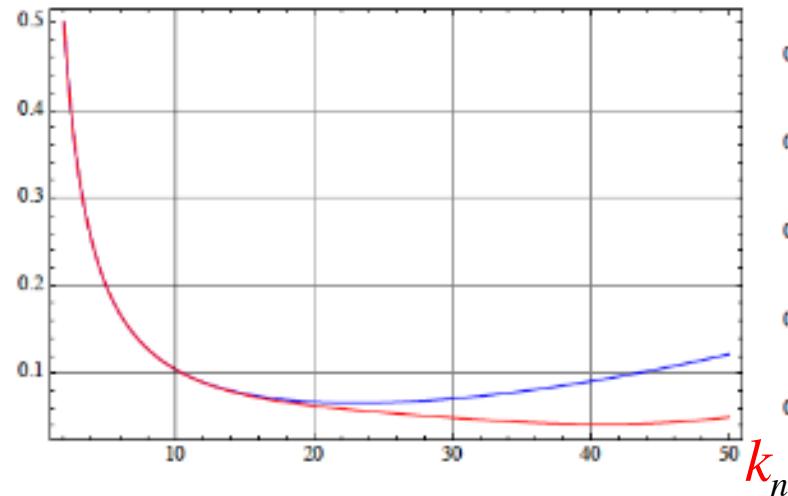
$$\textcolor{red}{k}_n^* = \left\lceil \left( \frac{6}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil$$

## □ Optimal Linear Weights

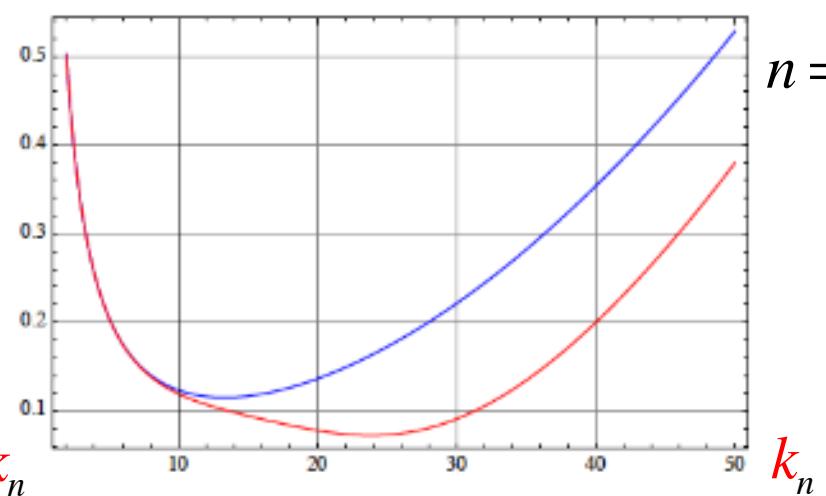


- Optimal Weights  $\leftrightarrow$  Uniform Weights

$\mathbf{D} = 1$



$\mathbf{D} = 5$



$n = 80$

$\text{MSE}\left(w_{ni}^*\right)$  —

$\text{MSE}\left(\frac{1}{k_n}\right)$  —

## □ Theorem 6 (Rate of Convergence)

Let  $k_n^* = \left\lceil \left( \frac{6}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil$  and  $w_{ni}^* \uparrow$

$$\mathbf{E}[\hat{m}(\mathbf{x}) - m(\mathbf{x})]^2 = O(n^{-2/3})$$

This rate can be improved for smoother characteristics  
(weights < 0)

$$O(n^{-2/3}) \rightarrow O(n^{-4/5}) \rightarrow O(n^{-6/9}) \rightarrow \dots \rightarrow O\left(n^{-\frac{2s}{2s+1}}\right) \rightarrow$$

smoothness →

□ **Theorem 6 (Rate of Convergence)**

Let  $k_n^* = \left\lceil \left( \frac{6}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil$  and  $w_{ni}^* \uparrow$

$$\mathbf{E} \left[ \hat{m}(\mathbf{x}) - m(\mathbf{x}) \right]^2 = O(n^{-2/3})$$

**The rate  $O(n^{-2/3})$  is optimal**

$$n^{2/3} \mathbf{E} \left[ \hat{m}(0) - m(0) \right]^2 \rightarrow c > 0$$

## Proof:

$$\bullet \quad \mathbf{E} \left[ \sum_{i=1}^{k_n} \mathbf{B}(i) \right]^2$$

$$\bullet \quad \mathbf{E} \left[ \sum_{i=1}^{k_n} \sum_{j=1}^{k_n} j \mathbf{B}(i) \mathbf{B}(j) \right]$$

$$\bullet \quad \mathbf{E} \left[ \sum_{j=1}^{k_n} i \mathbf{B}(i) \right]^2$$

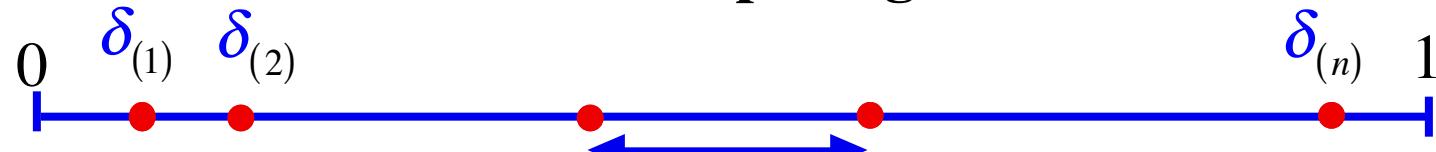
□ **Theorem 7 (Uniform Rate)**  $m(\bullet) \in \text{Lip}(\mathbb{R})$

Let  $k_n^* \simeq \left\lceil \frac{n^{2/3}}{\log^{2/3} n} \right\rceil$  and  $\{\mathbf{w}_{ni}^*\} \uparrow$

$$\sup_{x \in S(f)} |\hat{m}(x) - m(x)| = O_P\left(\frac{\log^{1/3} n}{n^{1/3}}\right)$$

**Proof:**

## Uniform Spacings



$$D_n = \max_{0 \leq i \leq n} \{ \delta_{(i+1)} - \delta_{(i)} \}$$

$$D_n = \frac{\log n}{n} + O\left(\frac{\log \log n}{n}\right) \text{ a.s}$$

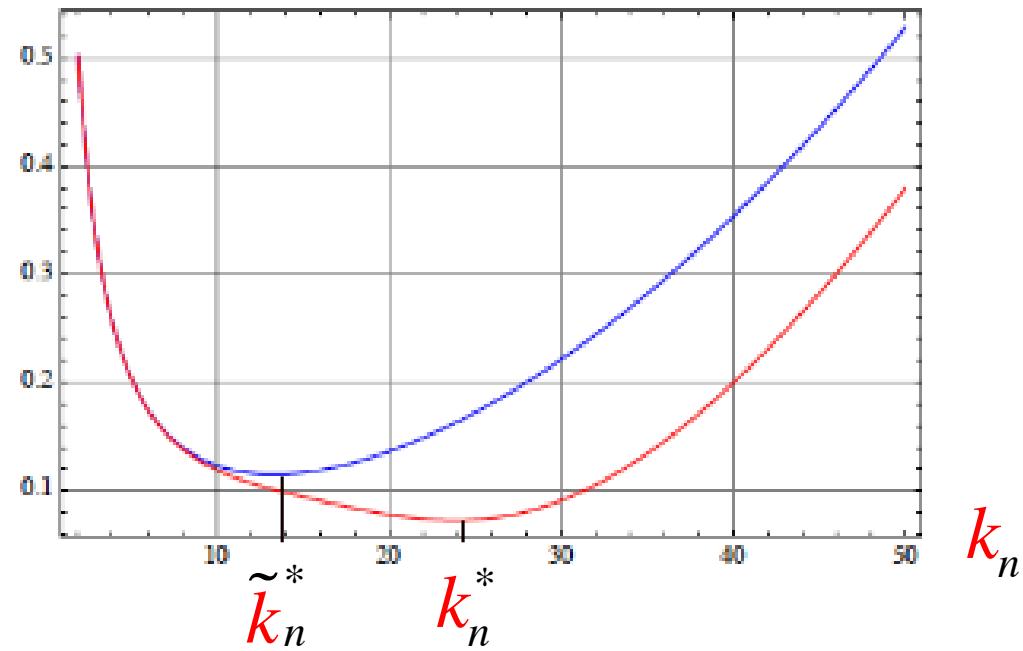
Slud, Devroye

□ Universal relationship for  $k_n^*$

**Optimal Linear Weights:**  $k_n^* = \left\lceil \left( \frac{6}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil$

**Uniform Weights:**  $\tilde{k}_n^* = \left\lceil \left( \frac{2}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil$

$$\frac{k_n^*}{\tilde{k}_n^*} = 3^{1/3} = 1.44\dots$$



Apply CV for  $\tilde{m}(\mathbf{x}) = \frac{1}{\tilde{k}_n} \sum_{i=1}^{\tilde{k}_n} Y_{[i]} \Rightarrow \tilde{k}_{CV} \Rightarrow k_{CV} = 3^{1/3} \tilde{k}_{CV}$

□ Kernel Weights:  $\hat{m}(\textcolor{blue}{x}) = \sum_{i=1}^{\textcolor{red}{k}_n} \textcolor{red}{v}_{ni} Y_{[i]}$

$$\textcolor{red}{v}_{ni} = \int_{(i-1)/\textcolor{red}{k}_n}^{i/\textcolor{red}{k}_n} \textcolor{blue}{p}(t) dt, \quad i = 1, \dots, \textcolor{red}{k}_n$$

↑ pdf on [0,1]

$$\mathbf{E} \left[ \hat{m}(\textcolor{blue}{x}) - m(\textcolor{blue}{x}) \right]^2 \quad \simeq \quad \frac{\mathbf{D}}{n^2} \left[ \sum_{i=1}^{\textcolor{red}{k}_n} i \textcolor{red}{v}_{ni} \right]^2 + \sum_{i=1}^{\textcolor{red}{k}_n} \textcolor{red}{v}_{ni}^2$$

## □ Kernel Weights

$$\bullet \quad \frac{1}{k_n} \sum_{i=1}^{k_n} i \textcolor{red}{v}_{ni} = \int_0^1 \textcolor{blue}{tp}(t) dt \left[ 1 + O\left(\frac{1}{k_n}\right) \right]$$

$$\bullet \quad k_n \sum_{i=1}^{k_n} \textcolor{red}{v}_{ni}^2 = \int_0^1 \textcolor{blue}{p}^2(t) dt \left[ 1 + O\left(\frac{1}{k_n}\right) \right]$$

## □ Kernel Weights

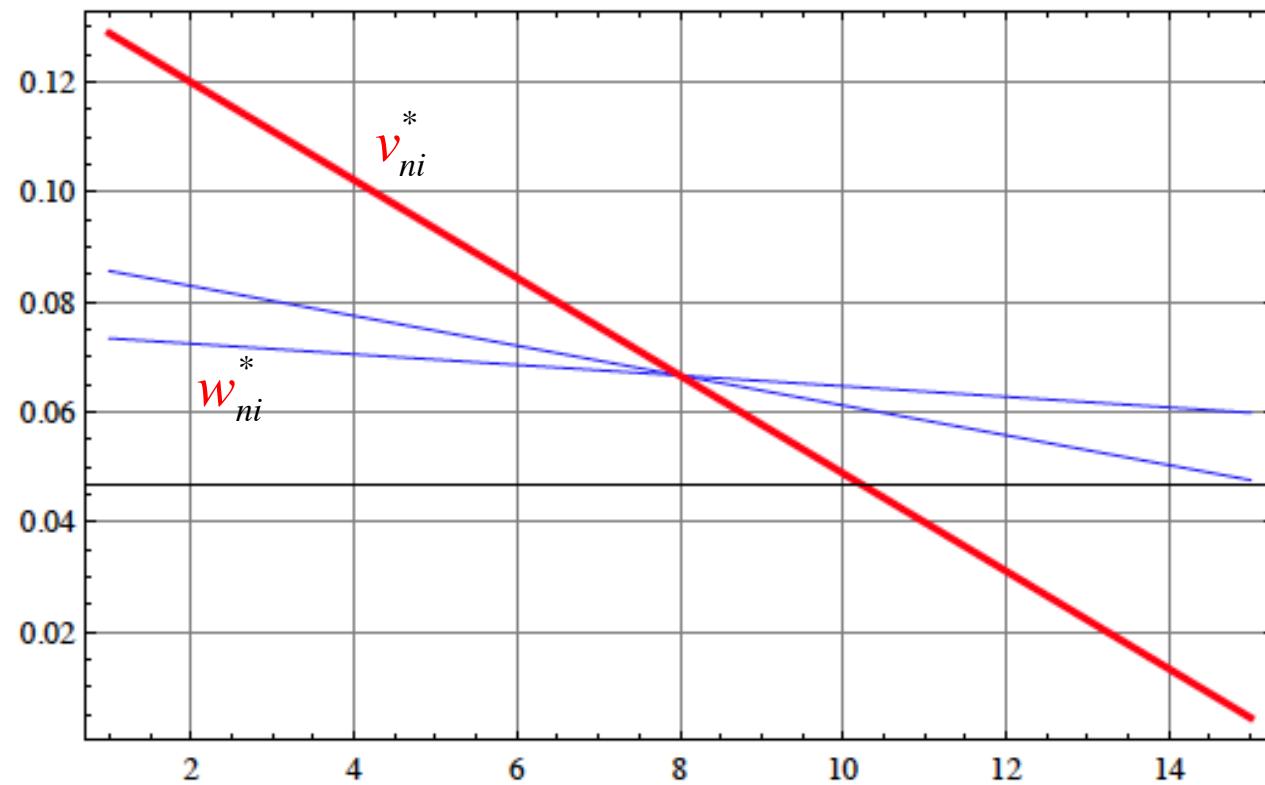
$$\mathbf{E} \left[ \widehat{\mathbf{m}}(\mathbf{x}) - \mathbf{m}(\mathbf{x}) \right]^2 \simeq \frac{\mathbf{D}}{n^2} \left[ \sum_{i=1}^{k_n} i \mathbf{v}_{ni} \right]^2 + \sum_{i=1}^{k_n} \mathbf{v}_{ni}^2 \xrightarrow{\quad} \min_{\mathbf{p}(\bullet), k_n}$$

$$\mathbf{D} \left[ \frac{k_n}{n} \int_0^1 t \mathbf{p}(t) dt \right]^2 + \frac{1}{k_n} \int_0^1 \mathbf{p}^2(t) dt \xrightarrow{\quad} \min_{\mathbf{p}(\bullet), k_n}$$

## □ Optimal Kernel Weights

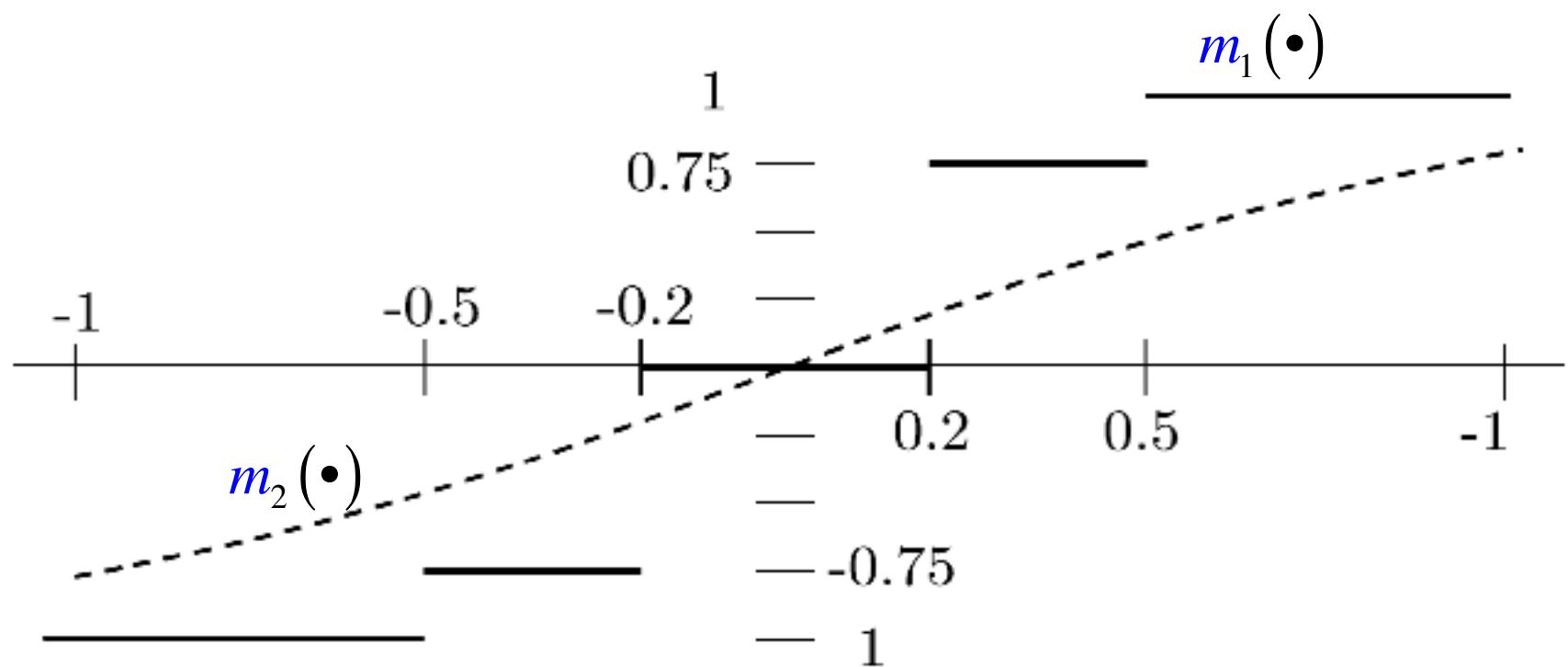
$$\nu_{ni}^* = \frac{1}{k_n^*} + \frac{k_n^* + 1}{k_n^{*2}} - i \frac{2}{k_n^{*2}},$$

$$k_n^* = \left\lceil \left( \frac{6}{\mathbf{D}} \right)^{1/3} n^{2/3} \right\rceil \quad \text{blue arrow} \rightarrow \text{the same as for } w_{ni}^*$$



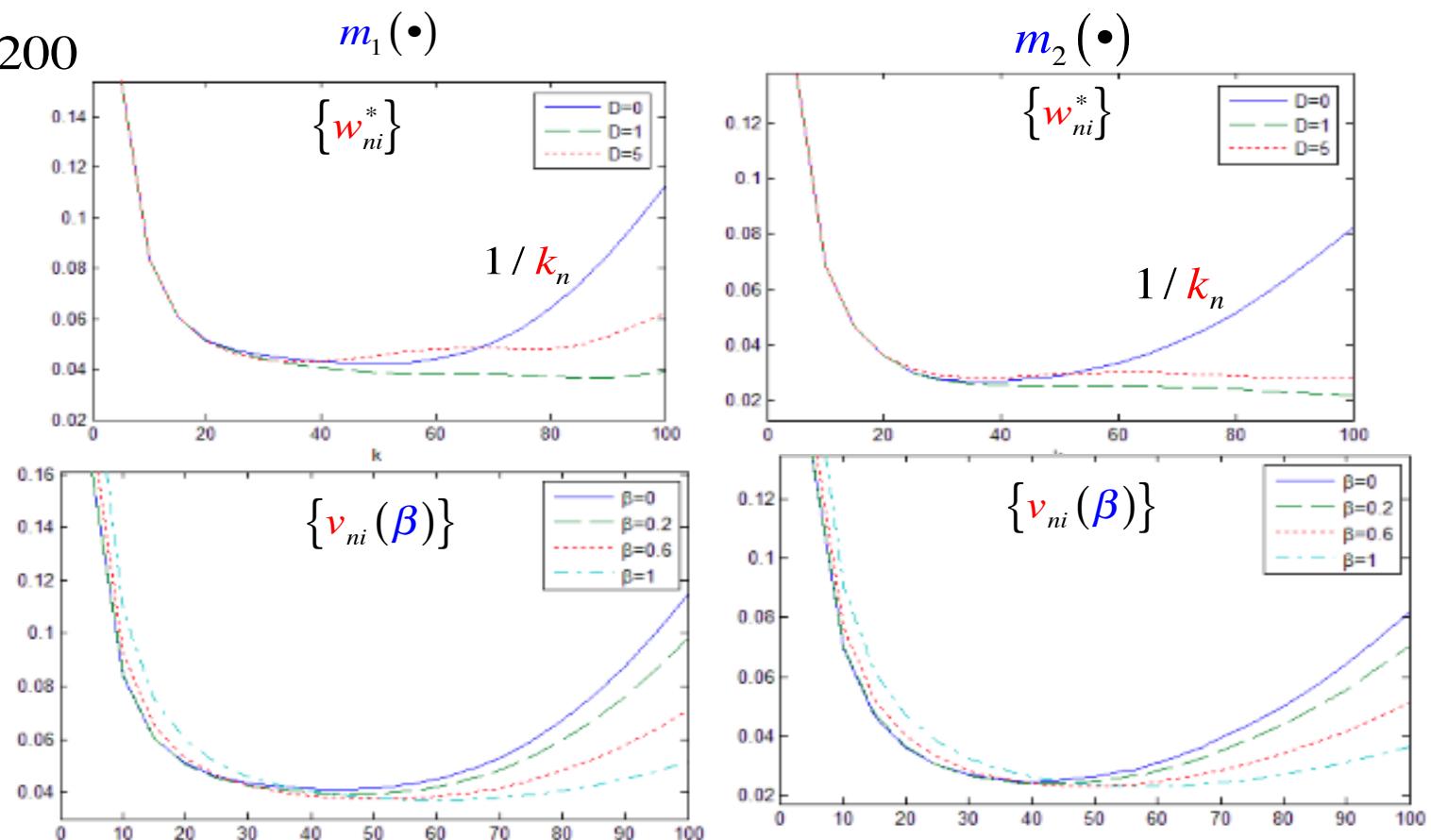


## Finite n



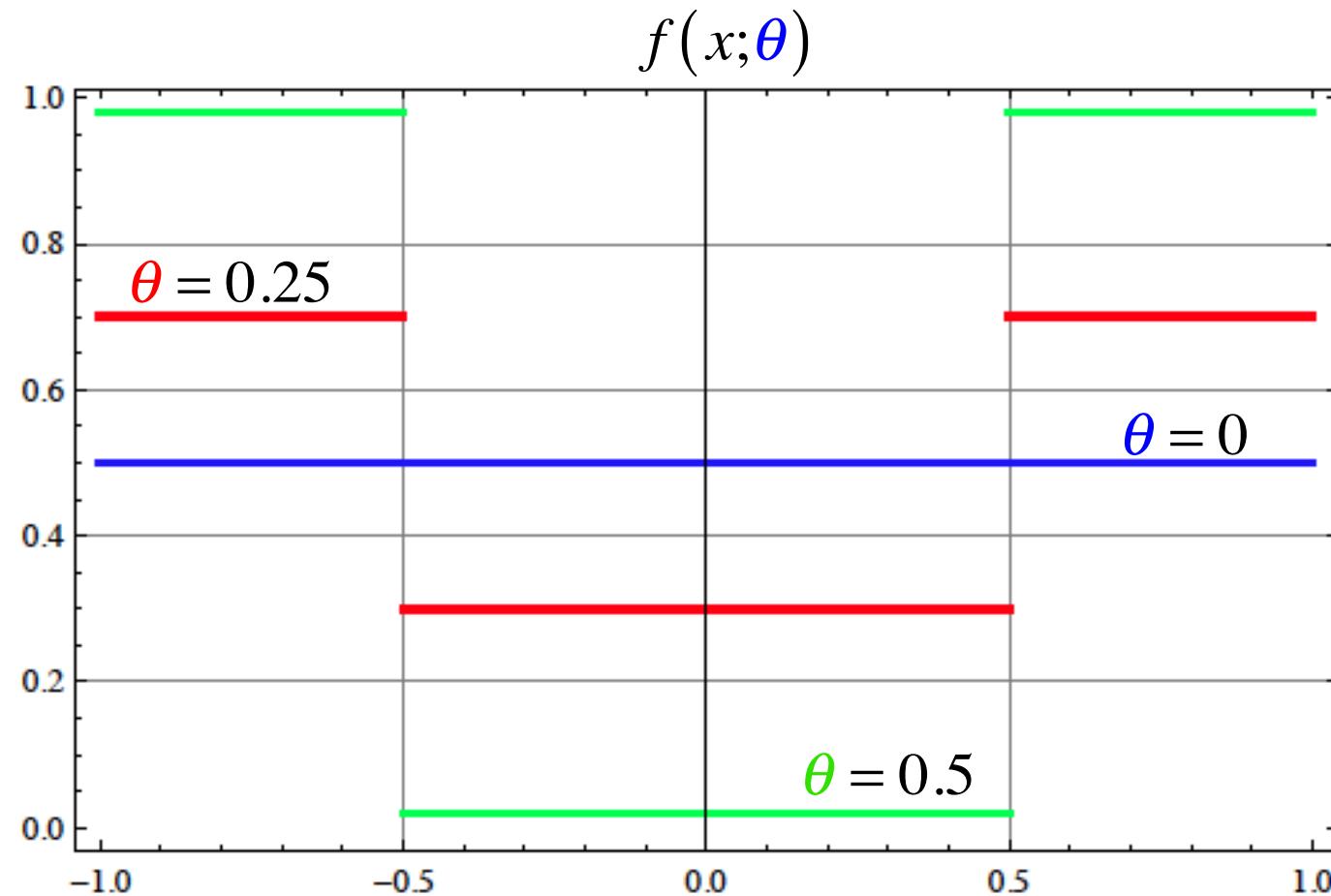
## Ex. 1: MISE versus $k_n$

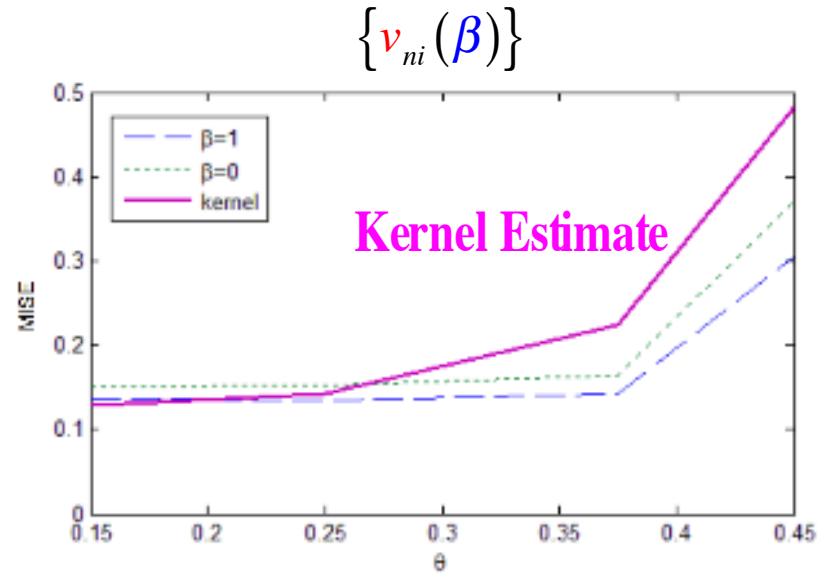
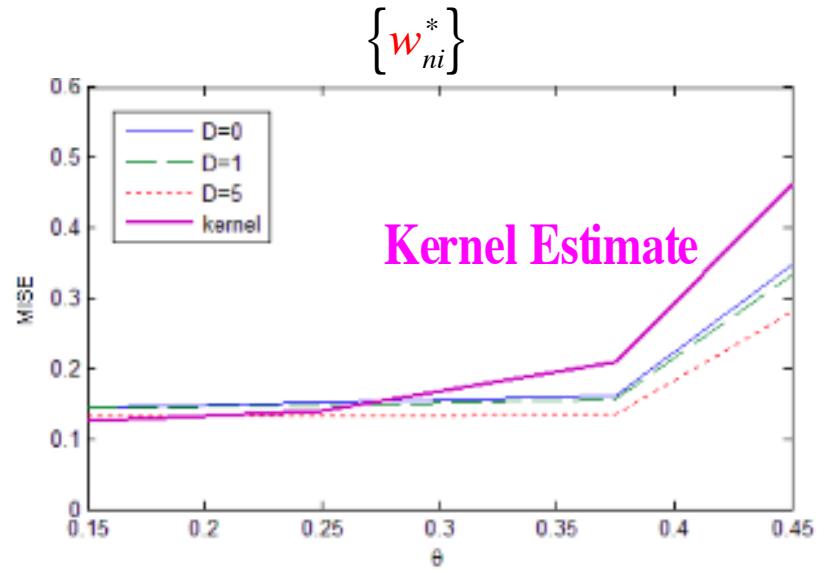
$n = 200$



$$v_{ni}(\beta) = \frac{1}{k_n^*} + \beta \frac{k_n^* + 1}{k_n^{*2}} - i \frac{2\beta}{k_n^{*2}} \Rightarrow v_{ni}(1) = v_{ni}^* \quad v_{ni}(0) = 1/k_n$$

## Ex. 2: MISE $\leftrightarrow$ The Input Density



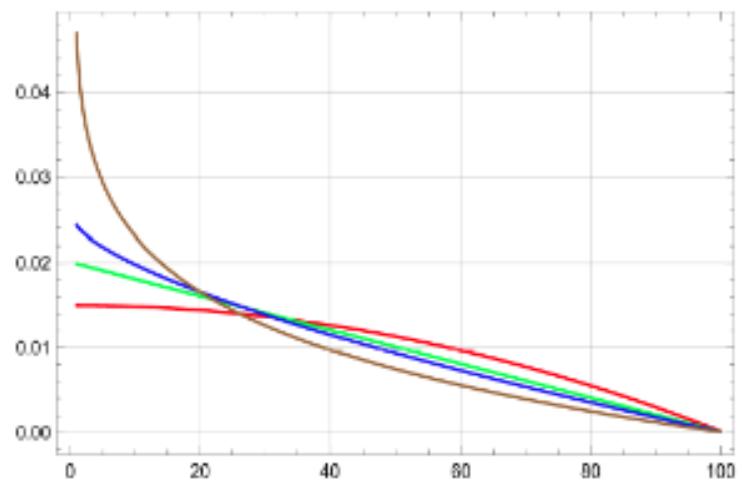
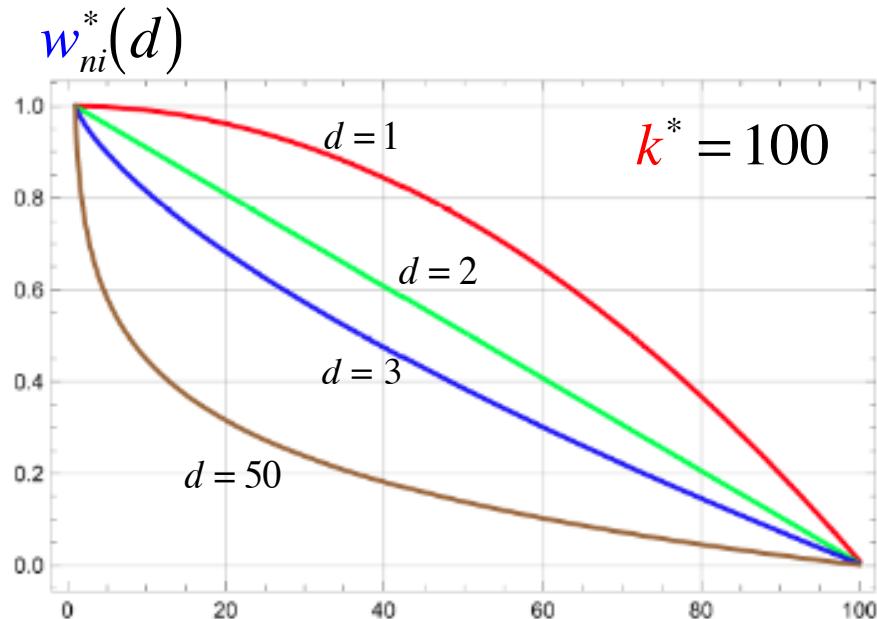


$$v_{ni}(\beta) = \frac{1}{k_n^*} + \beta \frac{k_n^* + 1}{k_n^{*2}} - i \frac{2\beta}{k_n^{*2}}$$

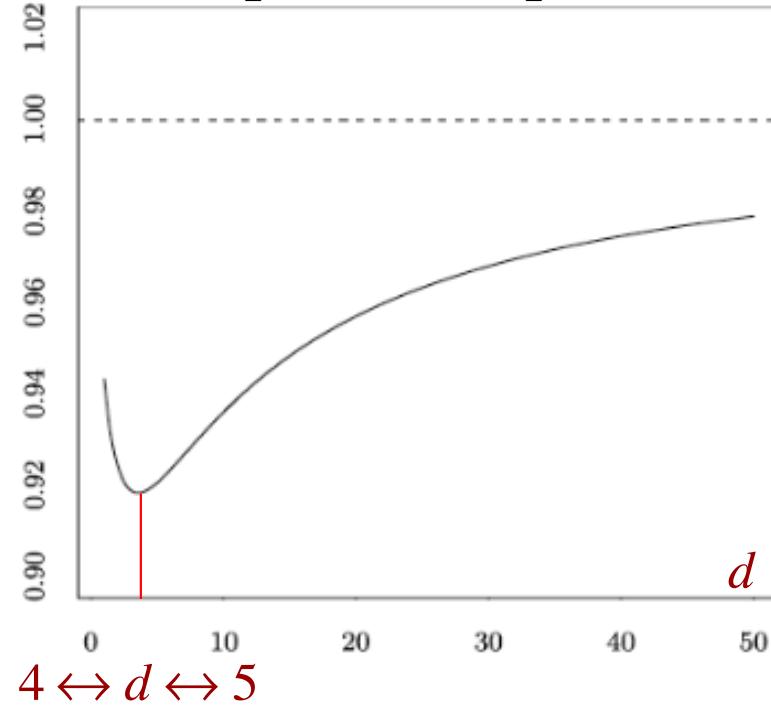
$$\Rightarrow v_{ni}(1) = v_{ni}^*$$

$$v_{ni}(0) = 1/k_n^*$$

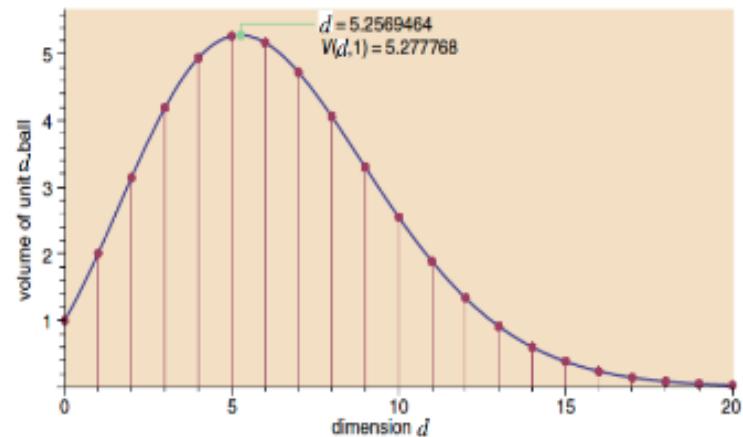
• Optimal Positive Weights in  $\mathbf{R}^d$ :  $w_{ni}^*(d) \Leftarrow m(\bullet) \in \mathbf{C}^2(\mathbf{R}^d)$



unweighted  $\rightarrow \mathbf{E}[\tilde{m}(\mathbf{X}) - m(\mathbf{X})]^2$   
 weighted  $\rightarrow \frac{\mathbf{E}[\hat{m}(\mathbf{X}) - m(\mathbf{X})]^2}{\mathbf{E}[\hat{m}(\mathbf{X})]^2}$



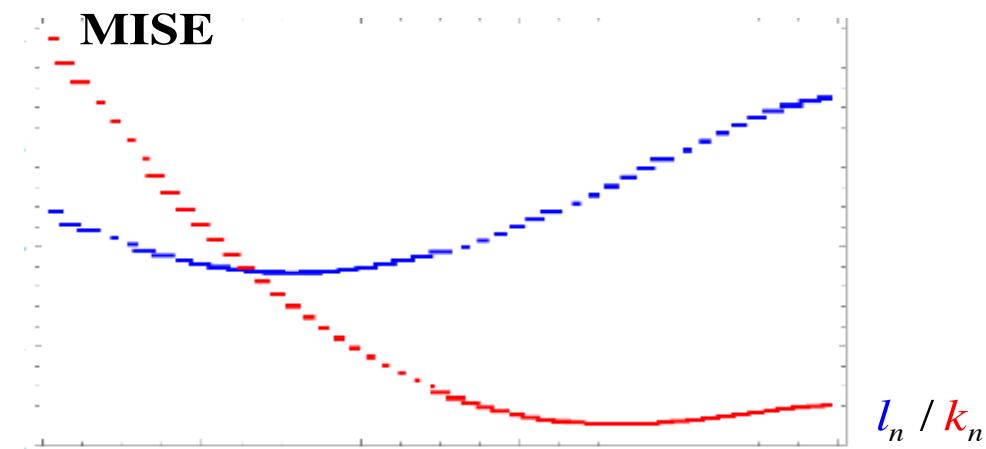
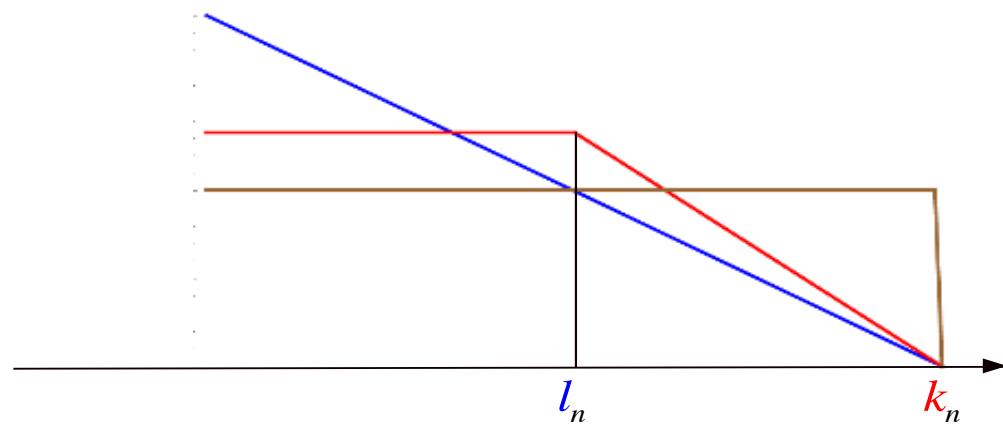
Volume of Unit Ball



- R. J. Samworth. Optimal weighted nearest neighbour classifiers. Ann. of Stat., 2012.



$L_\infty$  Weights  $\longrightarrow (k_n, l_n)$  - NN rules





**Negative Weights  $\Leftarrow$  Smooth  $m(\bullet)$**

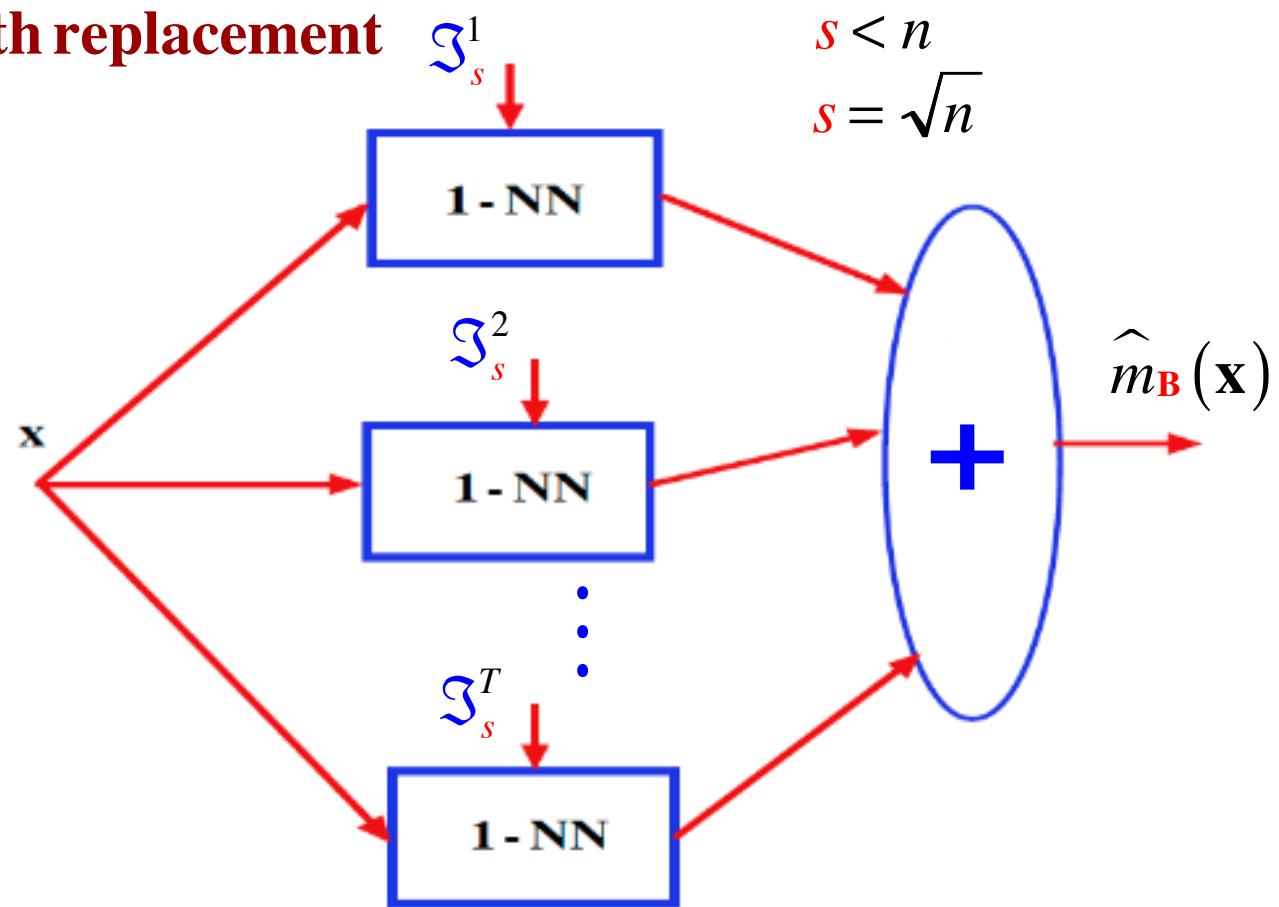
$$\sum_{i=1}^{k_n} i^2 \mathbf{w}_{ni} = 0 \quad \longrightarrow \text{better rates}$$

- The rates can be improved for smoother characteristics
- No smoothness of the input density is needed



## k-NN via 1-NN estimate Aggregation: Bagging Weights

resampling with replacement



- $\hat{m}_B(\mathbf{x}; \mathbf{s}) = \frac{m_{NN}(\mathbf{x}; \mathfrak{I}_{\mathbf{s}}^1) + \dots + m_{NN}(\mathbf{x}; \mathfrak{I}_{\mathbf{s}}^T)}{T}$
- For large  $T$  the average estimate  $\hat{m}_B(\mathbf{x}; \mathbf{s})$  is equivalent to the weighted NN estimate

$$\hat{m}_B(\mathbf{x}; \mathbf{s}) \approx \sum_{i=1}^n b_{ni} Y_i$$

$$b_{ni} = \left(1 - \frac{i-1}{n}\right)^{\mathbf{s}} - \left(1 - \frac{i}{n}\right)^{\mathbf{s}} \approx \tau (1 - \tau)^{i-1}, \quad \tau = \frac{\mathbf{s}}{n}$$

- If

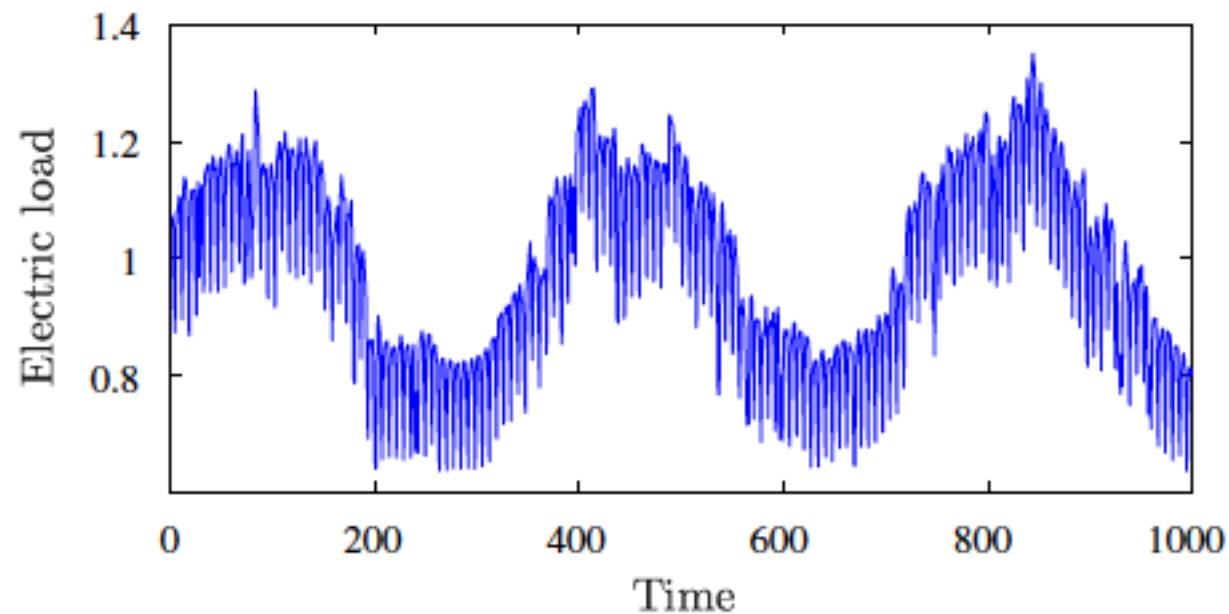
$$s_n \rightarrow \infty \text{ and } \frac{s_n}{n} \rightarrow 0$$

then  $\hat{m}_B(\mathbf{x}; s)$  is a consistent estimate of  $m(\mathbf{x})$

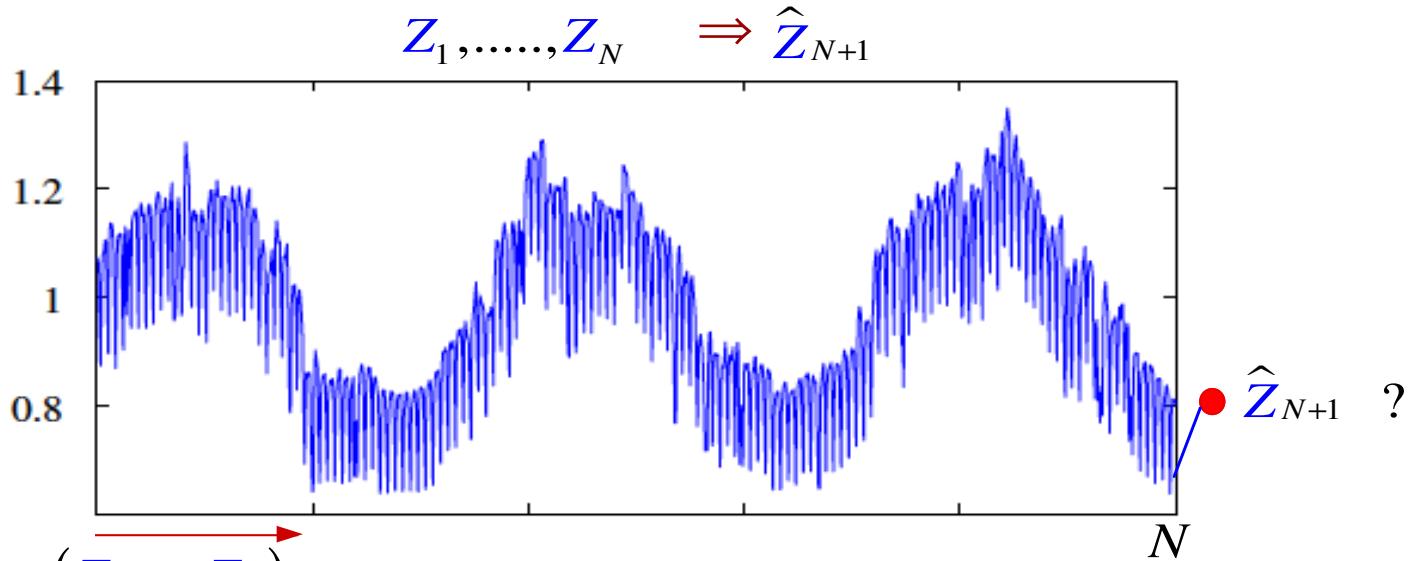
- Hall & Samworth. Properties of bagged NN classifiers  
J.Royal. Stat. Soc. 2005.
- Biau,..., Journal of Machine Learning Res. , 2010

## IV k-NN Regression Estimates: Time Series Forecasting

### ● Time Series Forecasting



**Poland Electricity Load: 1000 days**



$$1 \begin{cases} \mathbf{X}_1 = (\overrightarrow{Z_1, \dots, Z_d}) \\ Y_1 = Z_{d+1} \end{cases}$$

$$2 \begin{cases} \mathbf{X}_2 = (\overrightarrow{Z_2, \dots, Z_{d+1}}) \\ Y_2 = Z_{d+2} \end{cases}$$

$$n = N - d \begin{cases} \mathbf{X}_n = (\overrightarrow{Z_n, \dots, Z_{n+d-1}}) \\ Y_n = Z_{n+d} \end{cases}$$

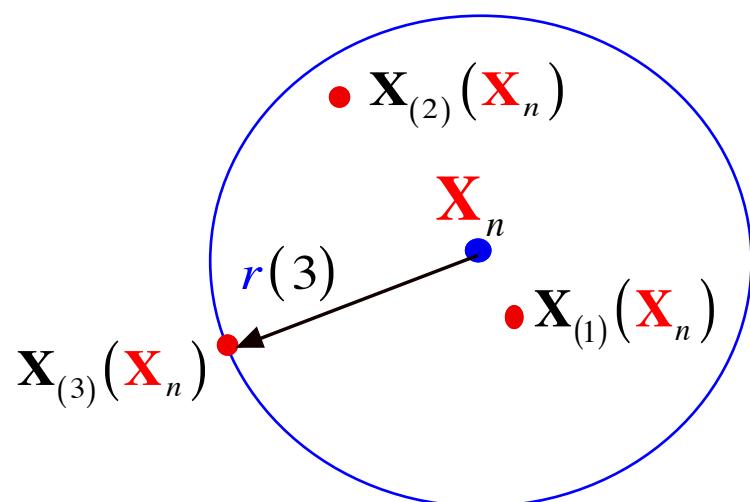
- The original time series is transformed into the  $d$ -dimensional regression type data for prediction of 1-unit of time in the future

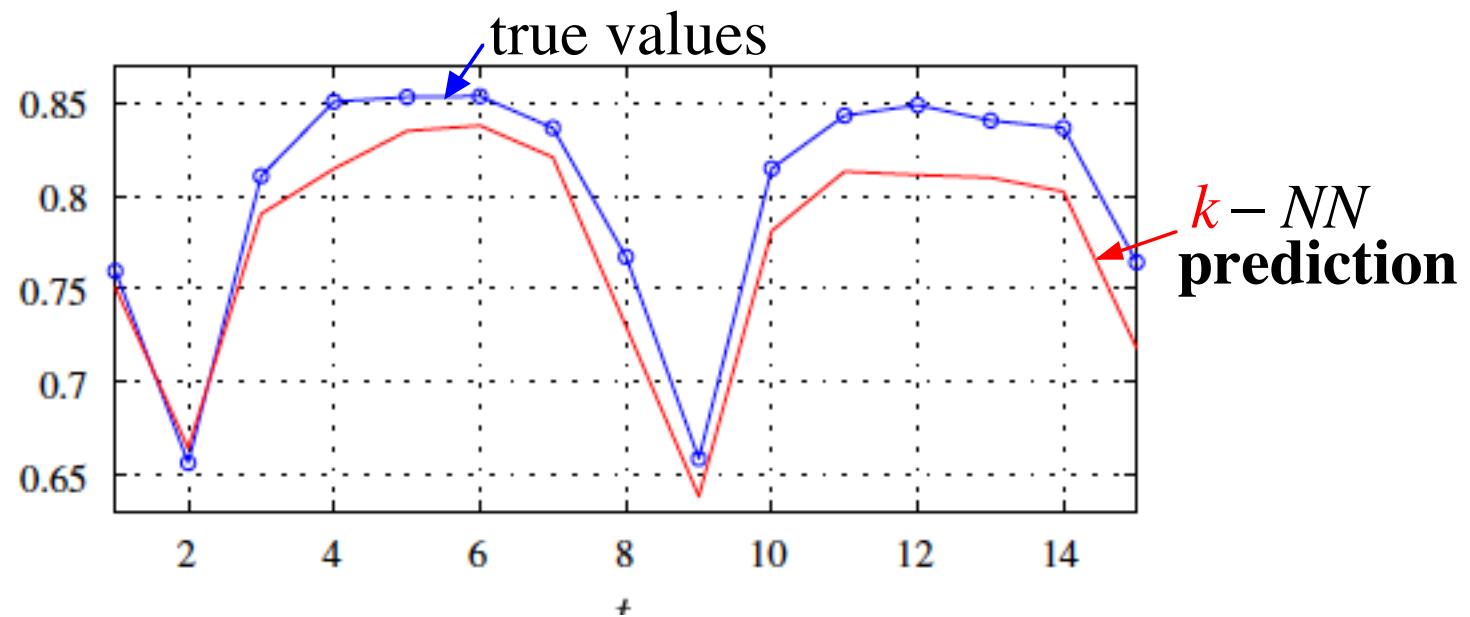
$$\mathbf{Z}_1, \dots, \mathbf{Z}_N \quad \xrightarrow{\hspace{1cm}} \quad \left\{ (\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n) \right\}$$

$$(\mathbf{X}_t, Y_t) \in \mathbf{R}^d \times \mathbf{R}$$

- To predict  $\hat{Z}_{N+1}$  we evaluate the  $k$ -NN regression estimate at the most recent point  $\mathbf{X}_n$

$$\hat{Z}_{N+1} = \frac{1}{k} \sum_{\{t < n : \mathbf{X}_t \in B_{r(k)}(\mathbf{X}_n)\}} Y_t , \quad \hat{Z}_{N+1} = \sum_{\{t < n : \mathbf{X}_t \in B_{r(k)}(\mathbf{X}_n)\}} w_t Y_t$$







## Asymptotic Theory

Let  $\{\mathbf{Z}_t\}$  be a stationary ( $\alpha$ - mixing ) stochastic process.

Then for  $\hat{\mathbf{Z}}_{N+1} = \frac{1}{k} \sum_{\{i: \mathbf{X}_i \in B_{(k)}(\mathbf{X}_n)\}} Y_i$  with  $k_n = cn^{2/(2+d)}$  we have

$$\mathbf{E}[\hat{\mathbf{Z}}_{N+1} - Z_{N+1}^*]^2 = O\left(n^{-\frac{2}{2+d}}\right) \leftarrow \text{optimal rate}$$



$Z_{N+1}^* = \mathbf{E}[\mathbf{Z}_{N+1} | \mathbf{Z}_N, \dots, \mathbf{Z}_{N-d+1}]$ - optimal predictor

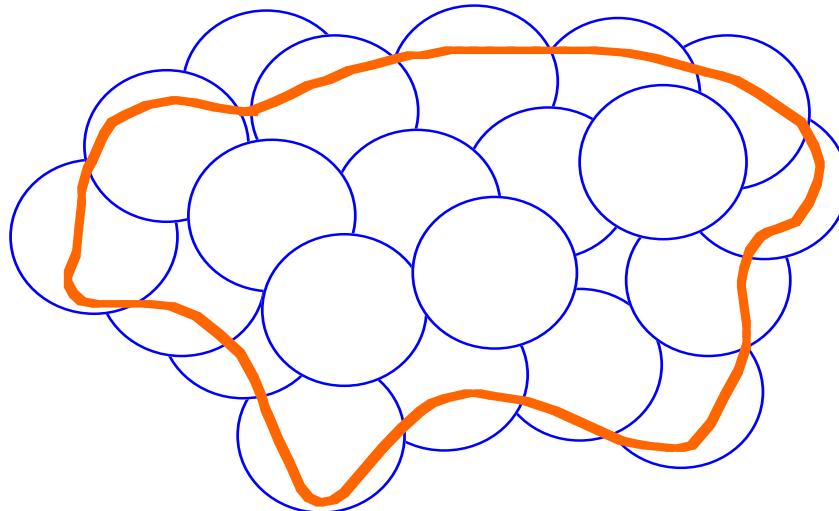
- **1-NN distance for stationary process**

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \in \mathbf{A} \subset R^d; \quad \mathbf{X}_{n+1} \in \mathbf{A}$$

	$E\left\ \mathbf{X}_{n+1} - \mathbf{X}_{(1)}(\mathbf{X}_{n+1})\right\ ^2$	$E\left\ \mathbf{X}_{n+1} - \mathbf{X}_{(k)}(\mathbf{X}_{n+1})\right\ ^2$
$d = 1$	$c(1) \frac{1 + \log(n)}{n}$	$c(1) \frac{k}{n} \left[ 1 + \log\left(\frac{n}{k}\right) \right]$
$d = 2$	$c(2) \frac{1}{n}$	$c(2) \frac{k}{n}$
$d \geq 3$	$c(d) \left(\frac{1}{n}\right)^{2/d}$	$c(d) \left(\frac{k}{n}\right)^{2/d}$

## Proof by the theory of covering numbers

$\varepsilon_{\min}(r; \mathbf{A})$  = the smallest radius such there exist  $r$  balls of this radius which cover the set  $\mathbf{A}$



$$\mathbf{E} \left\| \mathbf{X}_{n+1} - \mathbf{X}_{(1)}(\mathbf{X}_{n+1}) \right\|^2 \leq \frac{8}{n} \sum_{j=1}^n (\varepsilon_{\min}(j; \mathbf{A}))^2$$

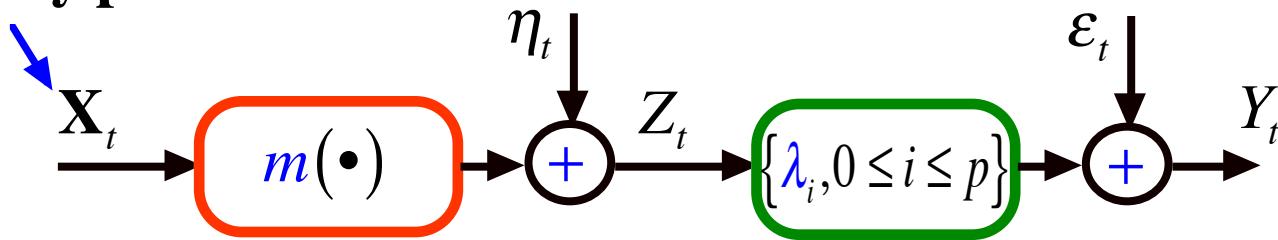
If  $\mathbf{A} \subset \mathbf{R}^d$  then  $\varepsilon_{\min}(r; \mathbf{A}) \leq \varepsilon_{\min}(1; \mathbf{A}) r^{-1/d}$

► Kulkarni & Posner - IT - 1995

## V Extensions and Concluding Remarks

### Input-Output Dependent System

stationary process

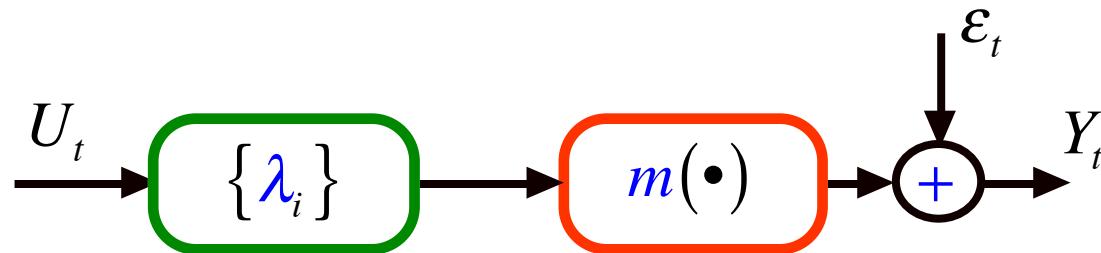


$$m(\mathbf{x}) = \sum_{l=0}^p a_l(\lambda) \underbrace{r_l(\mathbf{x})}_{\mathbb{E}[Y_t | \mathbf{X}_{t-l} = \mathbf{x}]}$$

Additive Solution in Terms of Lagged Regression Functions



## Other Time Series Systems $\Rightarrow$ Wiener System



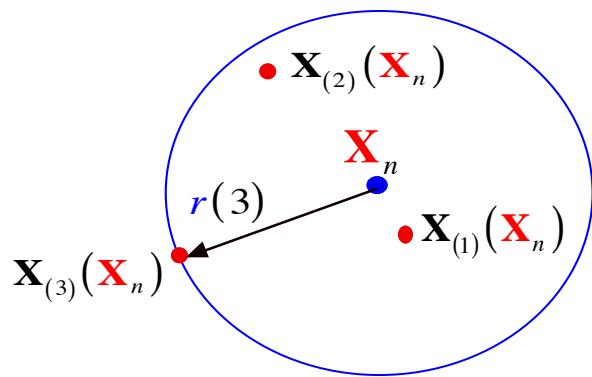
$$Y_t = m \left( \sum_{i=0}^{\infty} \lambda_i U_{t-i} \right) + \varepsilon_t$$

Dynamical Version of Single Index Model



## **$k$ -NN predictor with a distance correlation approach**

$$\hat{\mathbf{Z}}_{N+1} = \sum_{\{t < n : \mathbf{X}_t \in B_{r(k)}(\mathbf{X}_n)\}} w_t Y_t$$

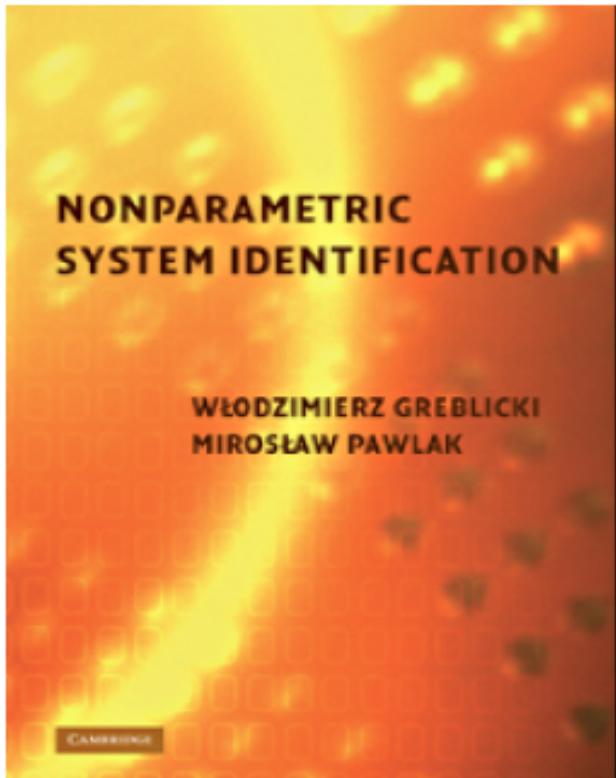


$\rho(\mathbf{X}_n; \mathbf{X}_t)$  = correlation distance between  $\mathbf{X}_n$  and  $\mathbf{X}_t$

Szekely, Rizzo,... Measuring and testing independence by correlation of distances; Annals of Stat. 2007.

## References

- W. Greblicki & M.Pawlak: IEEE Trans. Inf. Theory, 2017  
⇒ k-NN for nonlinear time series systems
- W. Greblicki & M.Pawlak: IEEE Trans. Automatic Control, 2018  
⇒ weighted k-NN for nonlinear time series systems
- M.Pawlak: IEEE Trans. Inf. Theory, 2018  
⇒ nonparametric model checks for nonlinear time series systems



**Thank you all very much  
for your time!**



**“Remember, the other team  
is using Machine Learning on your  
games to predict your play.  
So, kick the ball with your other foot!”**