



ESC FALL Week 3

Spline

esl ch5

오천도

방상용

Contents

1 Spline Introduction

2 Smooth Spline

3 Reproducing Kernel Hilbert Space

4 Wavelet Smoothing

A photograph of a bronze statue of a man in a long coat, standing on a stone pedestal. He is positioned in front of a building with a stone facade and large windows, which is partially covered in ivy. The background is filled with trees displaying vibrant autumn colors. The image is split vertically by a thin white line.

Part I

Spline
Introducition

기저 함수 – Basis Function

비선형 데이터의 예측률을 높이는
비선형 모델을 만드려면 어떻게 해야 할까?

선형 모델 : 해석하기 쉬운 $f(x)$ 에 대한 일차 테일러 근사

기저 함수 : 설명변수를 함수 형태로 나타낸 것

기저 함수와 설명변수의 선형 결합 -> 모델 표현

기저 함수 – Basis Function

$$h_m(X) : R^p \rightarrow R, m = 1, \dots, M$$

$$f(X) = \sum_{m=1}^M \beta_m h_m(X)$$

기저함수 h 에 상관없이,
이전과 같이 선형

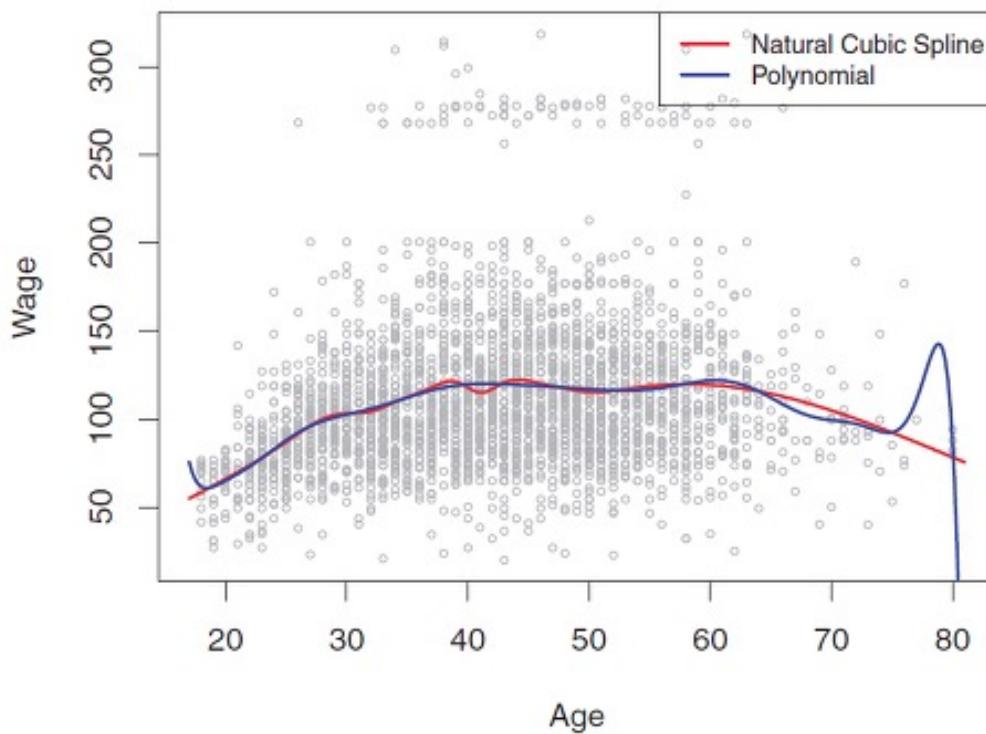
-> 계수 추정하기 위해서
least squares, 선형 모델에서
사용하던 도구들 적용이 가능함

기저함수 $h_m(X)$ 예시

- $h_m(X) = X_m, m = 1, \dots, p$
: 선형 모델
- $h_m(X) = X_j^2$
- $h_m(X) = X_j X_k$
: 고차 테일러 전개할 수 있도록 입력으로 다항식을 덧붙임
- $h_m(X) = \log(X_j)$
- $h_m(X) = \sqrt{X_j}$
: 비선형 변환

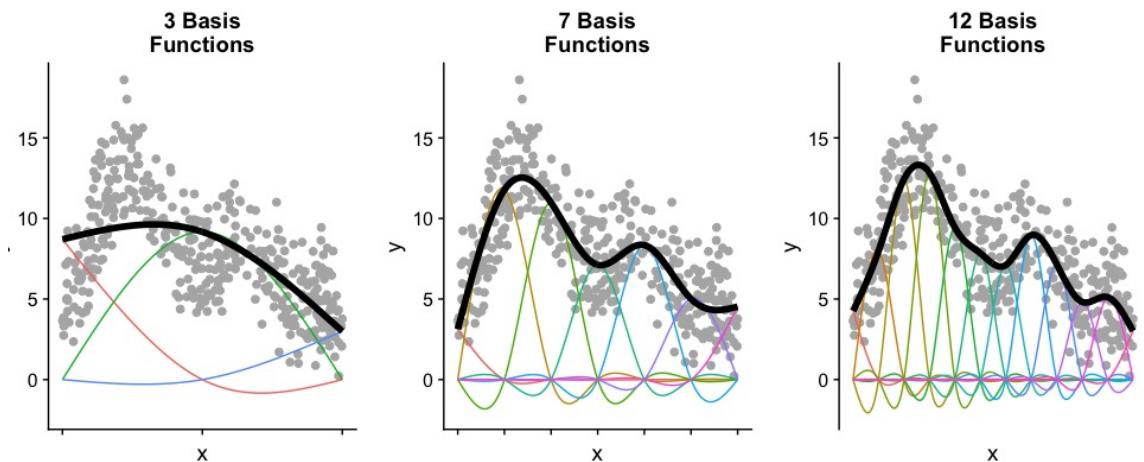
Polynomial Regression

비선형 데이터의 예측률을 높이는
비선형 모델을 만드려면 어떻게 해야 할까?



$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \epsilon_i,$$

$$y_i = \begin{cases} \beta_{01} + \beta_{11} x_i + \beta_{21} x_i^2 + \beta_{31} x_i^3 + \epsilon_i & \text{if } x_i < c; \\ \beta_{02} + \beta_{12} x_i + \beta_{22} x_i^2 + \beta_{32} x_i^3 + \epsilon_i & \text{if } x_i \geq c. \end{cases}$$



Spline

Knot : 매듭(점들)

Spline : 고정된 점들 사이를 잇는 방법

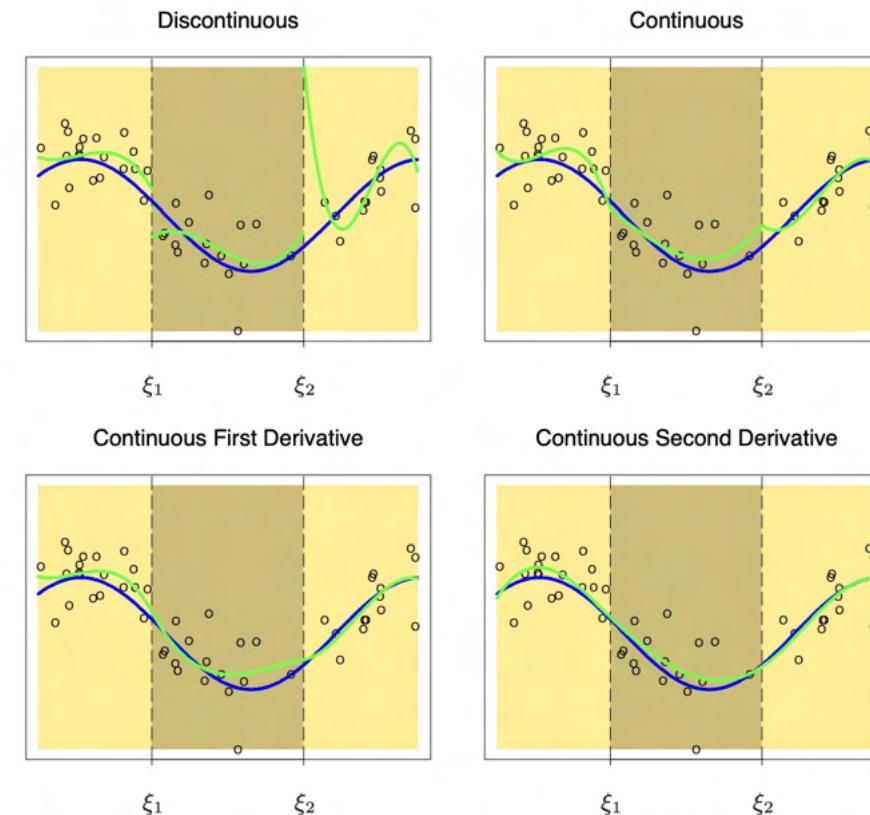
Piecewise polynomials : 고정된 점(knot) 사이를 부드럽게 연결



데이터의 문제에서는,
knot를 어디에, 몇 개 둬야 할까?

Degree-d Spline

각 knot에서 자기 자신과 도함수로부터 (d-1)계 도함수까지
모두 연속인 piecewise degree-d polynomial



Cubic Spline

자기자신(3차), 도함수, 이계도함수까지
모두 연속인 piecewise degree-3 polynomial

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \cdots + \beta_{K+3} b_{K+3}(x_i) + \epsilon_i$$

$$b_1(x_i) = x_i$$

$$b_2(x_i) = x_i^2$$

$$b_3(x_i) = x_i^3$$

$$b_{k+3}(x_i) = (x_i - \xi_k)_+^3, \quad k = 1, 2, \dots, K$$

4+K 개의 predictors를 Least Squares로 적합

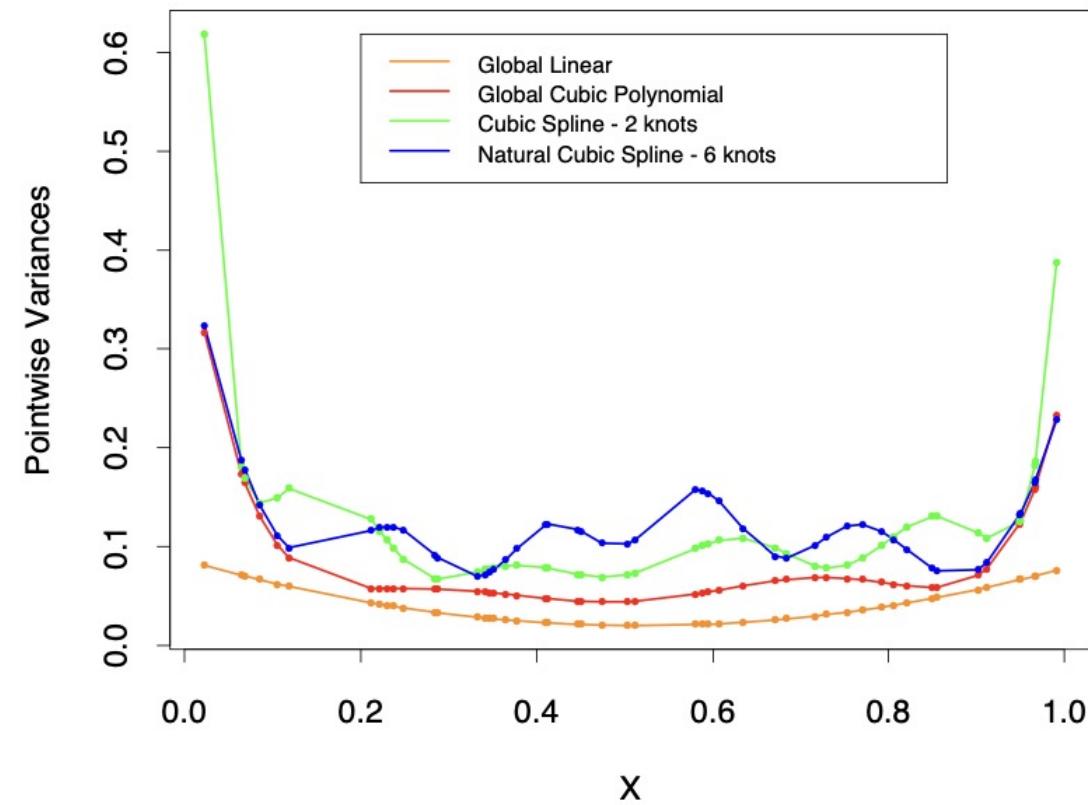
K : knot의 개수

4 : cubic function's unknowns

연속성 조건 -> knot 하나에 자유도 1 증가

Natural Cubic Spline

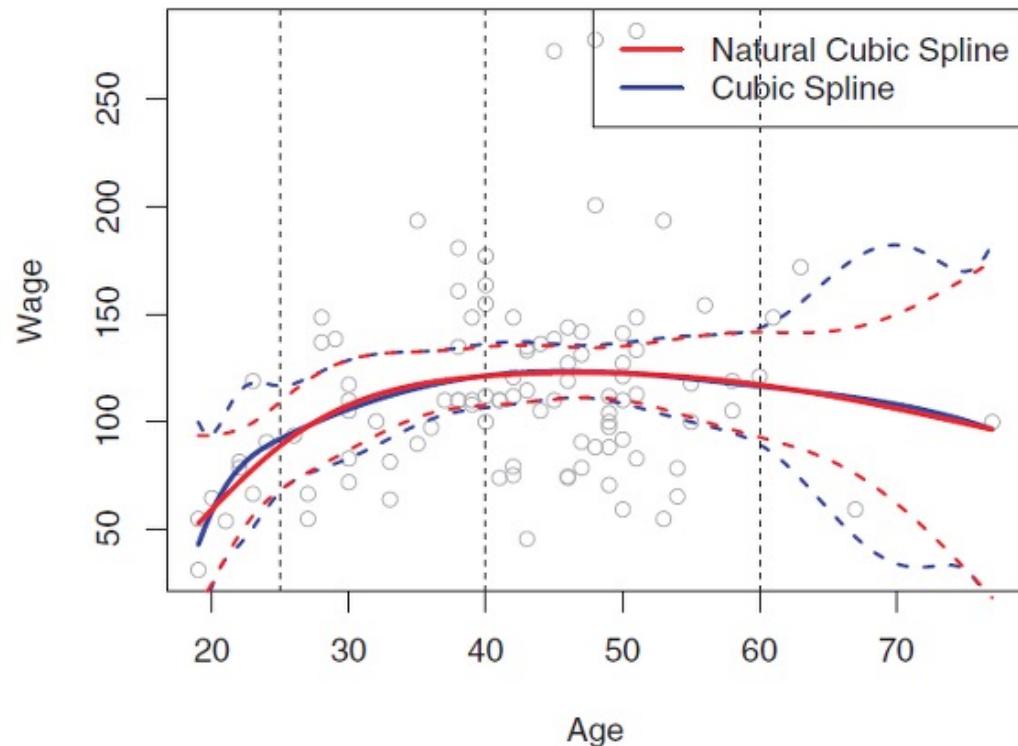
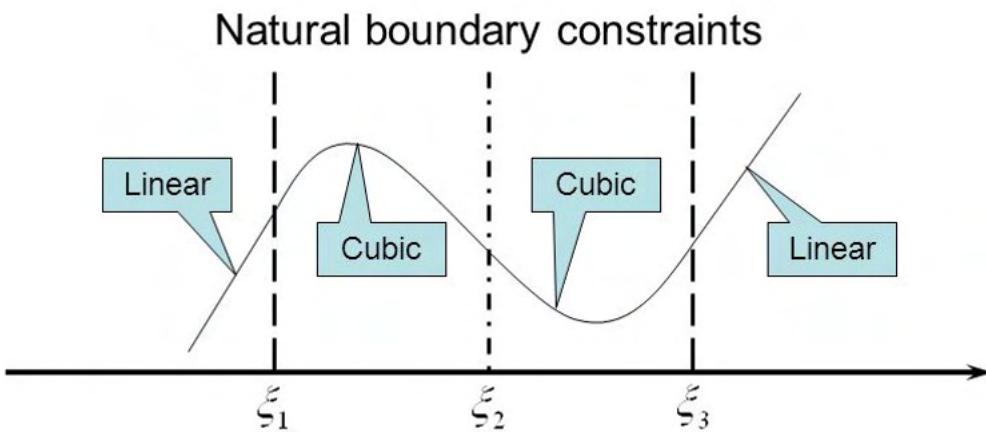
양 끝에서 부산이 튀는 문제는 어떻게 해결할까?



Natural Cubic Spline

Additional Boundary Constraints

: 가장 작은 knot보다 작거나, 가장 큰 knot보다 큰 부분에서는
함수가 직선이어야 한다는 조건 추가





Part 2

Smooth spline

Spline Method

Objective : find \hat{f} !



Basis function -> linear vs non-linear



Knot Selection -> Spline



Bias – Variance trade-off

Smooth Spline

Smooth Spline



Avoids the knot selection problem

→ Completely using a maximal set of knots



Penalty : smoothness of a function

Smooth Spline

$$\min_f RSS(f, \lambda) = \sum_i^N \{y_i - f(x_i)\}^2 + \lambda \int \{f''(t)\}^2 dt$$

Smooth Parameter

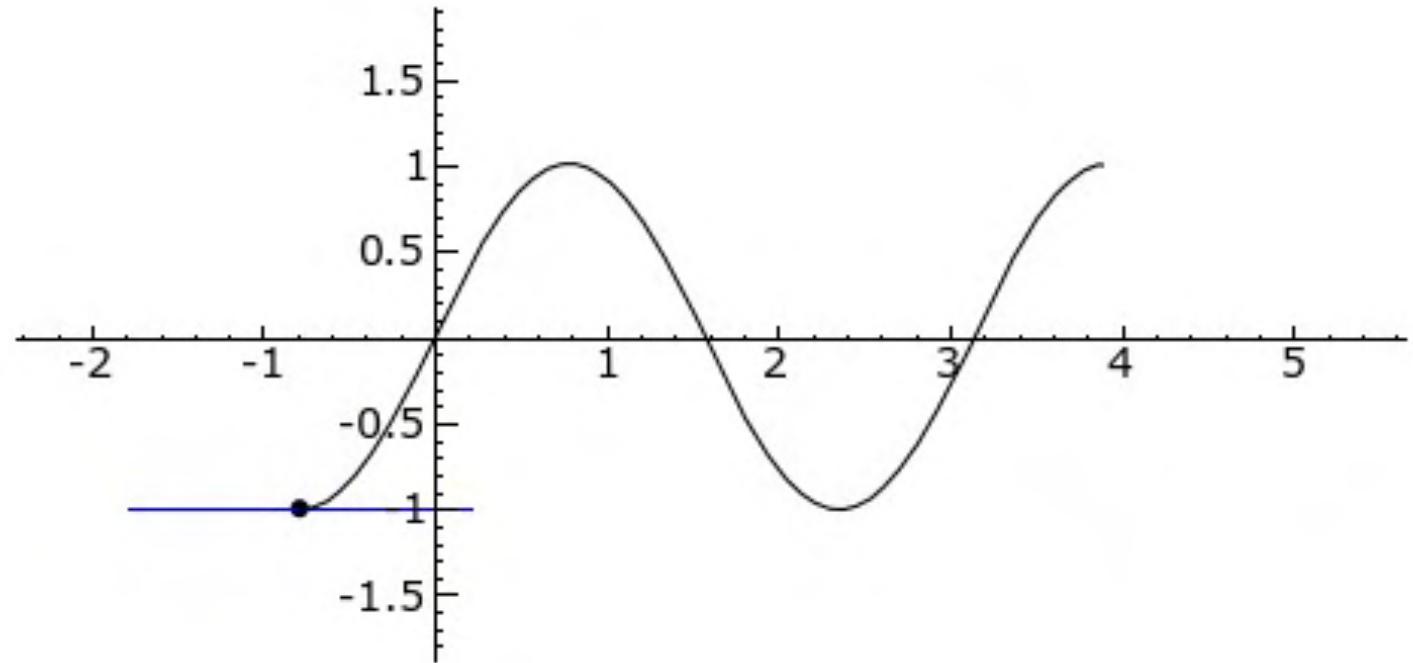
Measure of the smoothness

Second Derivative

이계도함수 f''

⇒ 함수의 오목성 측정!

⇒ 그래프의 **곡률** &
볼록성 지표



$$y = \sin(2x)$$

Smoothing parameter

$$\min_f RSS(f, \lambda) = \sum_i^N \{y_i - f(x_i)\}^2 + \lambda \int \{f''(t)\}^2 dt$$

Smooth Parameter

$\lambda \rightarrow 0$: no penalty : $f = \text{data interpolation}$

$\lambda \rightarrow \infty$: infinite penalty : $f = \text{simple least squares line fit}$

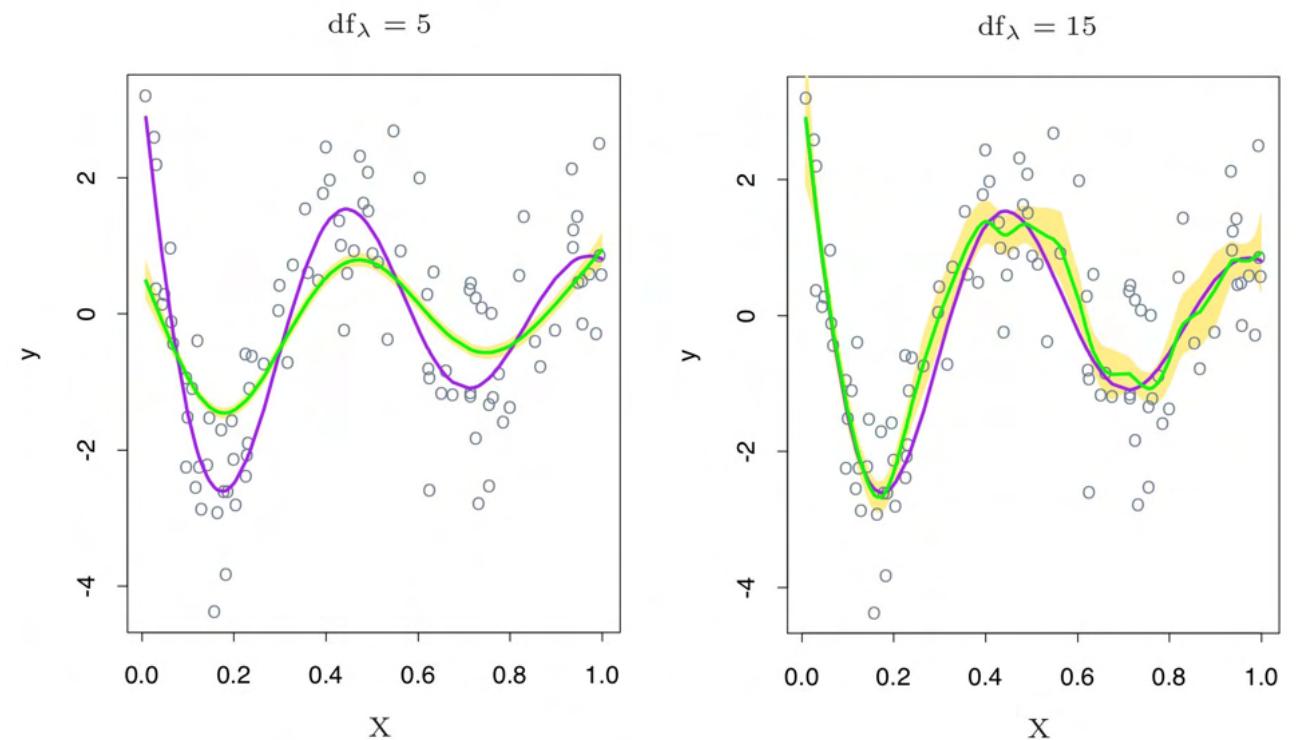
Effective degree of freedom

$$\hat{f} = S_\lambda y \quad : \quad df_\lambda = \text{tr}(S_y)$$

}

$\lambda \rightarrow 0 : df \rightarrow N, S_\lambda \rightarrow I$

$\lambda \rightarrow \infty : df \rightarrow 2, S_\lambda \rightarrow H$



Smooth Spline solution

$$\hat{\theta} = (N'N + \lambda\Omega_N)^{-1}N'y$$

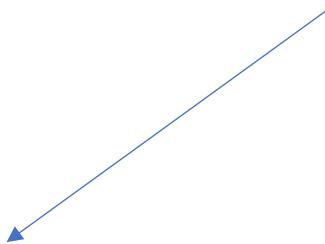
$$\{N\}_{ij} = N_j(x_i)$$

$$f(x) = \sum_j^N N_j(x)\theta_j$$

$$\{\Omega_N\}_{jk} = \int N_j''(t)N_k''(t)dt$$

Smooth Spline solution

$$f = N\theta \quad \longrightarrow \quad RSS(\theta, \lambda) = (y - N\theta)'(y - N\theta) + \lambda\theta'\Omega_N\theta$$



$$\hat{\theta} = (N'N + \lambda\Omega_N)^{-1}N'y$$

$$\hat{f} = N\hat{\theta} = N(N'N + \lambda\Omega_N)^{-1}N'y \coloneqq S_\lambda y$$

Selection of degree of freedom



Integrated squared prediction error (EPE)

$$\begin{aligned}
 EPE(\hat{f}_\lambda) &= E \left[(Y - \hat{f}_\lambda(X))^2 \right] \\
 &= Var(Y) + E \left[Bias^2(\hat{f}_\lambda(X)) + Var(\hat{f}_\lambda) \right] \\
 &= \sigma^2 + MSE(\hat{f}_\lambda)
 \end{aligned}$$

$$Bias(\hat{f}) = f - E(\hat{f}) = f - S_\lambda f$$

, when $Y = f(X) + \epsilon$,

$$Cov(\hat{f}) = S_\lambda Cov(y) S_\lambda' = S_\lambda S_\lambda'$$



$$f(X) = \frac{\sin(12(X + 0.2))}{x + 0.2}$$

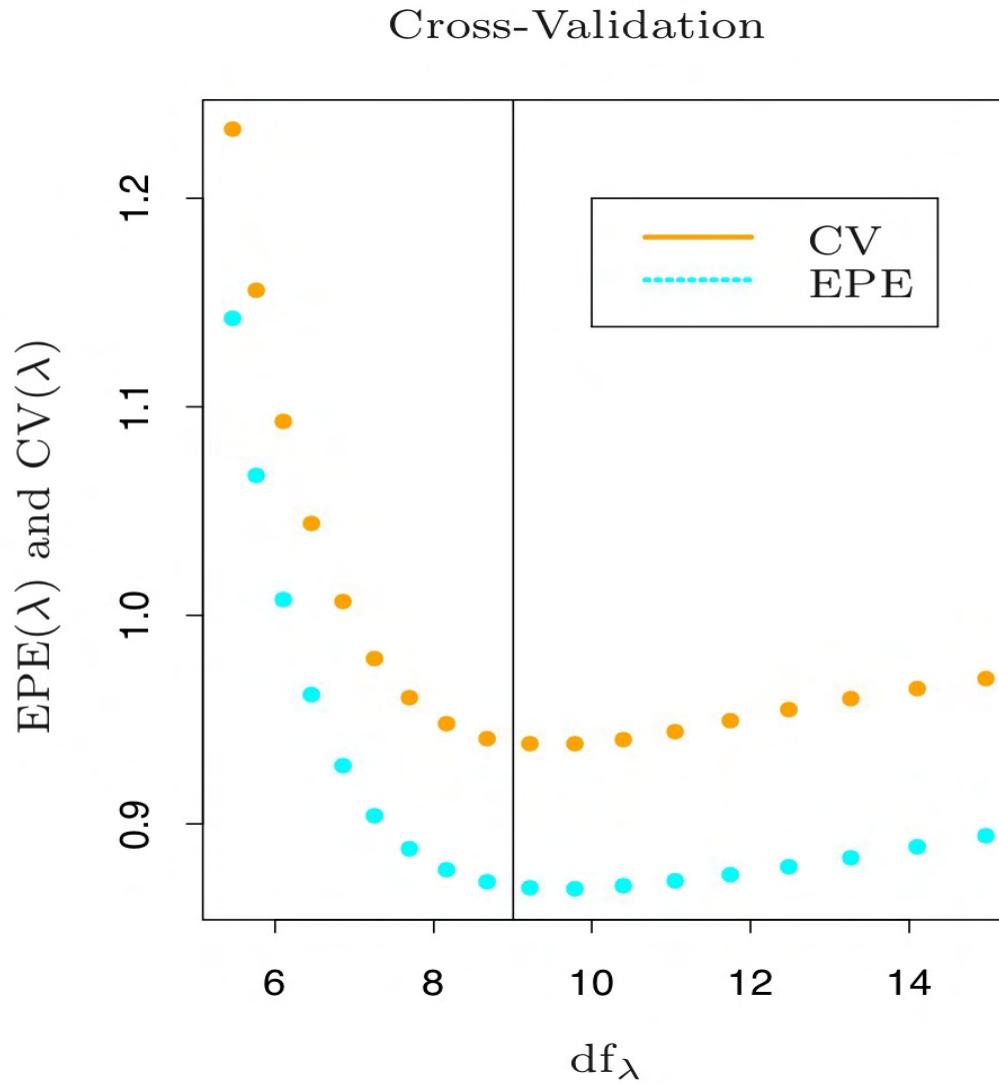
Selection of degree of freedom



Cross Validation

$$\text{CV}(\hat{f}_\lambda) = \frac{1}{N} \sum_i^N \left(y_i - \hat{f}_\lambda^{-i}(x_i) \right)^2 = \frac{1}{N} \sum_i^N \left(\frac{y_i - \hat{f}_\lambda(x_i)}{1 - S_\lambda(i, i)} \right)^2$$

Selection of degree of freedom





Part 3

Reproducing
Kernel Hilbert
Space(RKHS)

Motivation



Smooth Spline 문제의 solution



SVM Kernel trick

Background. Hilbert Space

Hilbert Space = “Complete” “Inner Product” “linear”

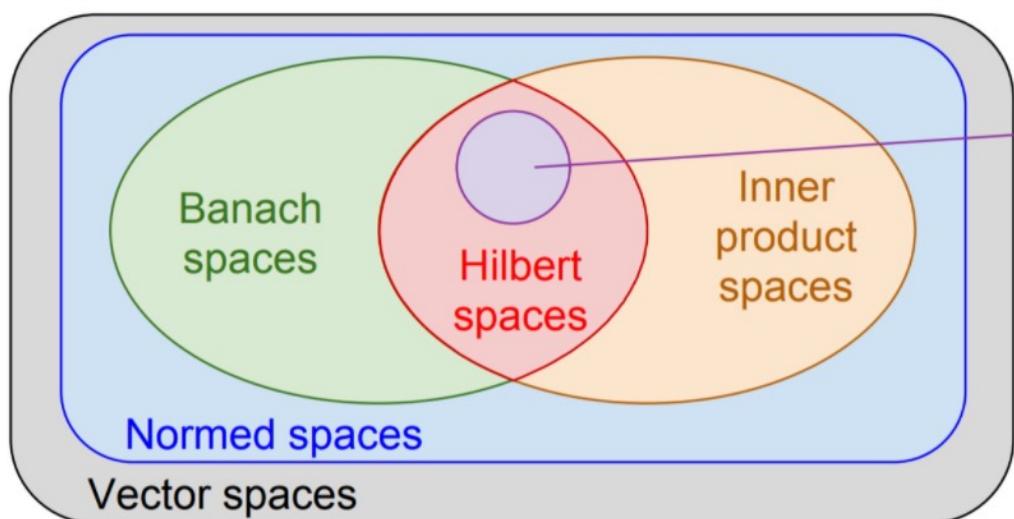


Cauchy Sequence

For any $\epsilon > 0, \exists n \in \mathbb{Z}^+ s.t. d(x_k, x_m) < \epsilon$ whenever $n \geq k, m$

Background. Hilbert Space

Hilbert Space = “Complete” “Inner Product” “linear”

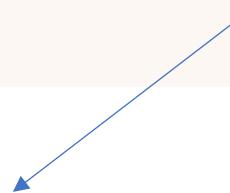


$$\int_{\mathcal{X}} f(x)d(x)dx = \langle f, g \rangle_H$$

$$\sqrt{\langle f, f \rangle_H} = \|f\|_H$$

Background. Hilbert Space

Hilbert Space = “Complete” “Inner Product” “linear”



For all $f, g \in H, f + g \in H$

For all $f, g \in H, fg \in H$

General Optimization Problem

$$\min_{f \in H} \left[\sum_i^N L(y_i, f(x_i)) + \lambda J(f) \right]$$

$$J(f) = \int \frac{|\widetilde{f(s)}|^2}{\widetilde{G(s)}} ds$$

$$f(x) = \sum_{k=1}^K \alpha_l \phi_k(X) + \sum_i^N \theta_i G(X - x_i)$$

Reproducing Kernel Hilbert Space (RKHS)

Produce RKHS



Feature map



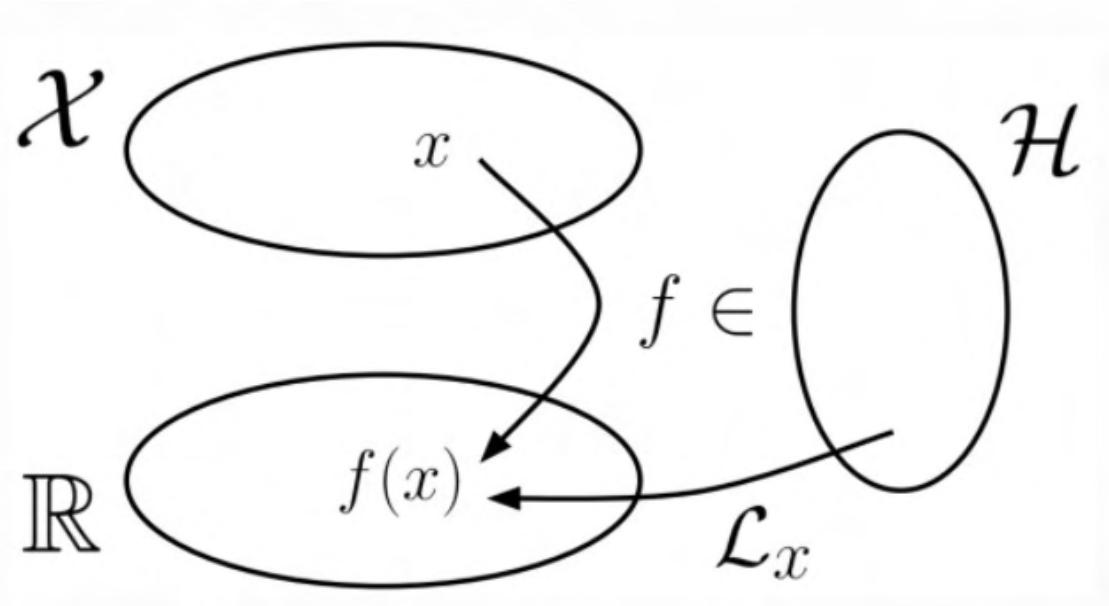
Integral as a linear operator

Reproducing Kernel Hilbert Space (RKHS) – Method 1

Let χ be an arbitrary set, and H a Hilbert space of all functions $f : \chi \rightarrow \mathbb{R}$. For each element $x \in \chi$, the evaluation map functional is a linear functional that evaluates each $f \in H$ at the point x , written

$L_x : H \rightarrow \mathbb{R}$, where $L_x(f) = f(x)$ for all $f \in H$

L_x is continuous at every $f \in H$



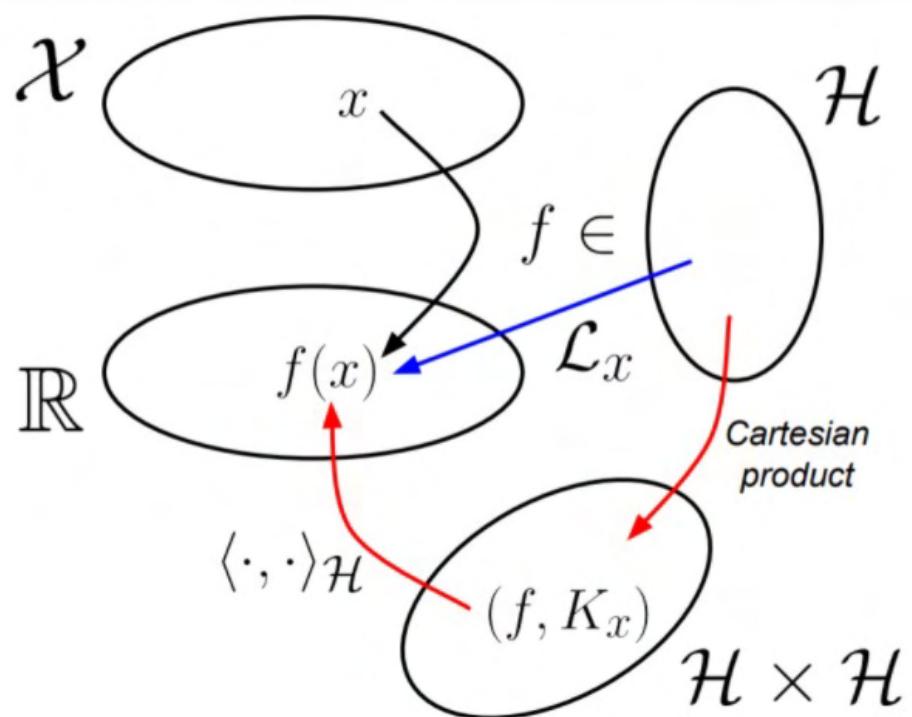
Q. evaluation map?

Reproducing Kernel Hilbert Space (RKHS) – Method 1

*If L_x is continuous at every $f \in H$, then, for each L_x ,
there is a unique function $K(\cdot, x) \in H$ s. t. for every $f \in H$,*

$$L_x(f) = f(x) = \langle f, K(\cdot, x) \rangle_H$$

This equation is known as the reproducing property.



Reproducing Kernel Hilbert Space (RKHS) – Method 1

*If L_x is continuous at every $f \in H$, then, for each L_x ,
there is a unique function $K(\cdot, x) \in H$ s. t. for every $f \in H$,*

$$L_x(f) = f(x) = \langle f, K(\cdot, x) \rangle_H$$

This equation is known as the reproducing property.

$$L_x(f) = f(x) = \langle f, K(\cdot, x) \rangle_H$$

$$L_x(K(\cdot, y)) = K(x, y) = \langle K(\cdot, x), K(\cdot, y) \rangle_H$$

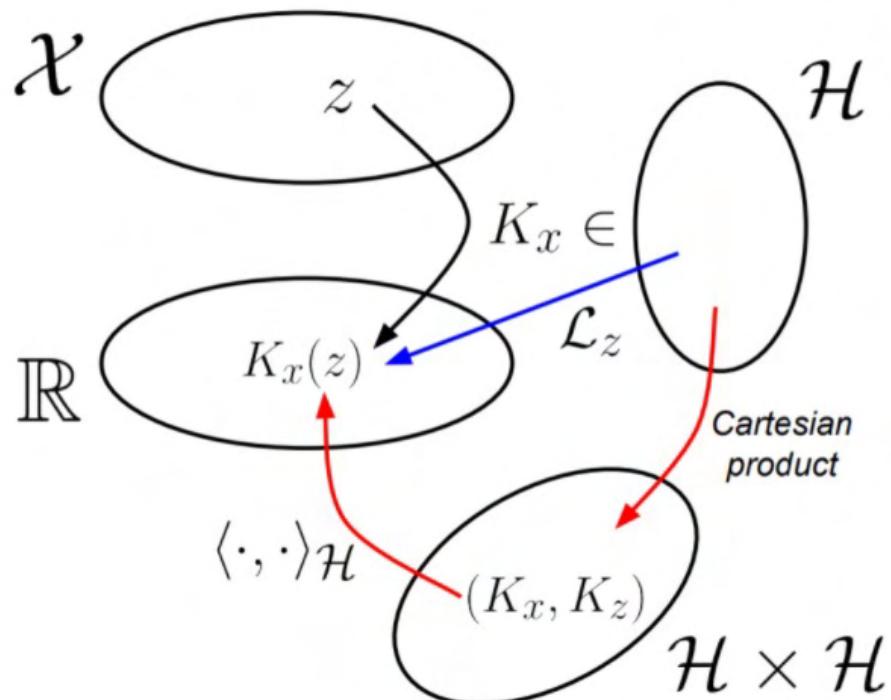
Reproducing Kernel Hilbert Space (RKHS) – Method 1

Let χ be an arbitrary set, and H a Hilbert space of all functions $f : \chi \rightarrow \mathbb{R}$.

If, for all $x \in \chi$, the linear evaluation functional $L_x : H \rightarrow \mathbb{R}$ is continuous at every $f \in H$, we can construct the reproducing kernel, which is a bivariate function $K : \chi \times \chi \rightarrow \mathbb{R}$ defined by

$$K(x, y) = \langle K(\cdot, x), K(\cdot, y) \rangle_H,$$

and the Hilbert space H is called a reproducing kernel Hilbert space (RKHS).

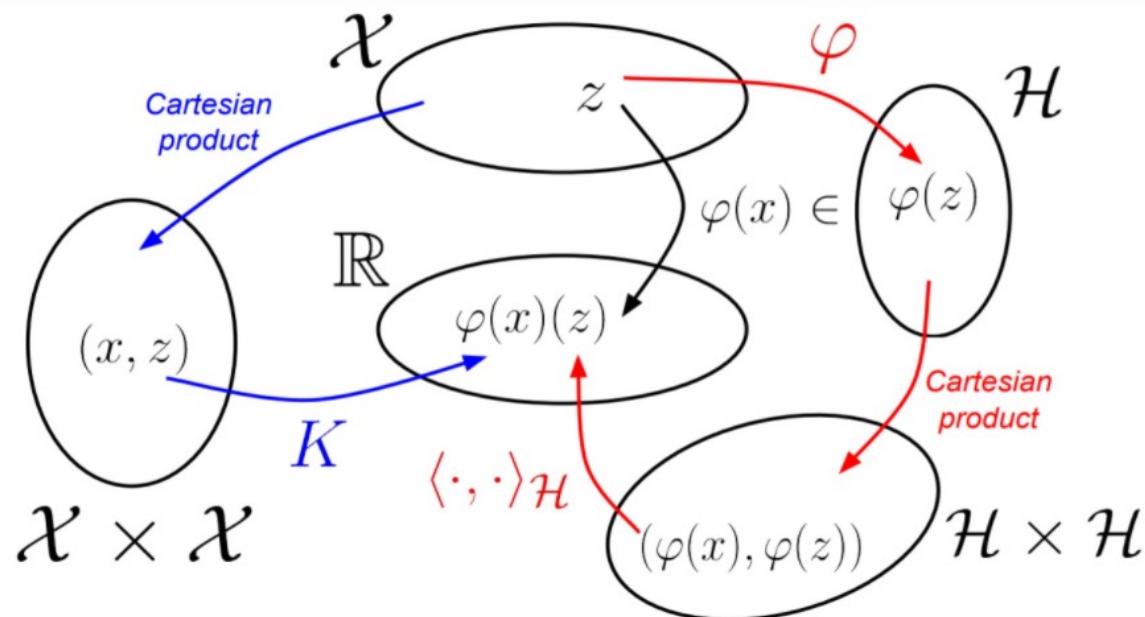


$$K(x, y) = \langle K(\cdot, x), K(\cdot, y) \rangle_H$$

Reproducing Kernel Hilbert Space (RKHS) – Method 1

Given the RKHS H , we can find a function $\phi: \chi \rightarrow H$ that links up χ and H , and defined as $\phi(x) = K(., x)$ for all $x \in \chi$.

In machine learning, ϕ is frequently referred to as the 'feature map', and H as the 'feature space'.



$$K(x, y) = \langle \phi(x), \phi(y) \rangle_H$$

$$\begin{aligned} H &:= \overline{\text{span}}\{\phi(x) : x \in \chi\} \\ &= \overline{\text{span}}\{K(x, .) : x \in \chi\} \end{aligned}$$

Properties of Kernel



Symmetry (assumption)

$$K(x, y) = K(y, x) \text{ for all } x, y \in \chi$$



Positive Semi-definite (from the reproducing property)

*Gram matrix $\{K\}_{ij}$ is positive semidefinite,
i.e. $a'Ka \geq 0$*

$$\int \int K(x, y) dx dy \geq 0$$

Reproducing Kernel Hilbert Space (RKHS) – Method 2

(Mercer's theorem)

Fix a symmetric kernel function K on a compact set χ

$$K : \chi^2 \rightarrow R$$
$$T_k : L_2(\chi) \rightarrow L_2(\chi)$$

$$T_k f(x) = \int K(x, y) f(y) dy$$

The integral operator T is characterized by K .

Reproducing Kernel Hilbert Space (RKHS) – Method 2

(Mercer's Theorem)

If K is continuous and T_k is positive semidefinite, then T has eigenfunction ψ_i , with eigenvalue λ_i s.t.

$$K(u, v) = \sum_j^{\infty} \lambda_j \psi_j(u) \psi_j(v)$$

That series converges uniformly.

$$\psi_i \in L_2(\chi), \quad \lambda \geq 0, \quad \|\psi_i\|_{L^2} = 1$$

$$(T_K \psi_a)(x_i) = \lambda_a \psi_a(x_i)$$

Reproducing Kernel Hilbert Space (RKHS) – Method 2

(Mercer's Theorem)

*For $\chi \in \mathbb{R}^d$ compact and $K: \chi^2 \rightarrow \mathbb{R}$ continuous and symmetric,
the followings are equivalent:*

1. Every Gram matrix is positive semidefinite.

2. The integral operator T_k is positive semidefinite.

3. We can express K as

$K(u, v) := \sum \lambda_i \psi_i(u) \psi_i(v)$ for fixed $\lambda_i \geq 0$ and orthonormal $\psi_i: \chi \rightarrow \mathbb{R}$

4. K is the reproducing kernel of an RKHS on χ

Reproducing Kernel Hilbert Space (RKHS) – Method 2

By the reproducing property, $\langle f, K(\cdot, x) \rangle_H = f(x)$

$$Thus, \langle \psi_j, \psi_j' \rangle_H = \begin{cases} \frac{1}{\lambda_j} & \text{if } j = j' \\ 0 & \text{otherwise} \end{cases}$$

$$\theta_j = \langle f, \psi_j \rangle, \text{ thus } f = \sum_j^{\infty} \theta_j \psi_j \quad \xrightarrow{\text{blue arrow}} \quad \|f\|_H^2 = \langle f, f \rangle_H = \sum_j^{\infty} \frac{\theta_j^2}{\lambda_j}$$

(Representer Theorem)

$$\min Q_{\lambda}(f) = \sum_i^n (y_i - f(x_i))^2 + \lambda \|f\|_H^2 \quad \Rightarrow \quad \text{solution } \hat{f}(x) = \sum_i^n \alpha_i K(x, x_i)$$

Non-uniqueness of basis function

$$K(x, y) = \langle K(\cdot, x), K(\cdot, y) \rangle_H$$

$$\Phi_1(x) = K(\cdot, x)$$

$$\Phi_2(x) : x \rightarrow (\sqrt{\lambda_1} \psi_1(x), \sqrt{\lambda_2} \psi_2(x))$$

$$K(x, y) = \langle \Phi_1(x), \Phi_1(y) \rangle_H = \langle \Phi_2(x), \Phi_2(y) \rangle_H$$

But... Kernel \Leftrightarrow RKHS : one-to-one by the theorem not in this ppt

Meaning of RKHS



무한 차원의 문제 => 유한 차원의 solution!



X : input space \rightarrow 모집단

$\{x_1, x_2, \dots, x_N\}$: sample, or training data

모집단 전체를 고려하는 최적화 문제

\Rightarrow sample data만을 가지고 solution 도출 가능

Bayesian interpretation : RKHS



A function f = a realization of a zero-mean stationary Gaussian process
With prior covariance function K



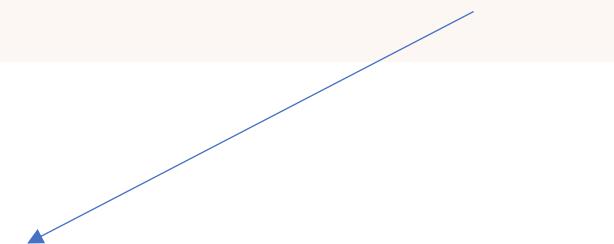
$$H = H_0 \oplus H_1$$

The direct sum of null space & orthogonal space

Back to Smooth Spline

$$\min_f RSS(f, \lambda) = \sum_i^N \{y_i - f(x_i)\}^2 + \lambda \int \underline{\{f''(t)\}^2 dt}$$

$$f(x) = \sum_j^N N_j(x) \theta_j$$



In “Sobolev Space”

Back to Smooth Spline

Banach Space
▷ *Sobolev Space*

Sobolev Space
vs **Hilbert Space**

Normed vector space
i.e. vector space + norm

Inner product space
e.v. + p.s.
(not complete : can have holes)

Euclidian space
finite dim. on \mathbb{R} + scal. p.

Hilbert space
i.e. complete inner product space
(Banach complete for its norm)

Banach space
complete for its norm

Sobolev space
Sobolev space + scal. p.

Sobolev Space

$$W^{m,p}(\Omega) := \{u \in L^p(\Omega) : D^\alpha u \in L^p(\Omega) \text{ for } 0 \leq |\alpha| \leq m\}$$



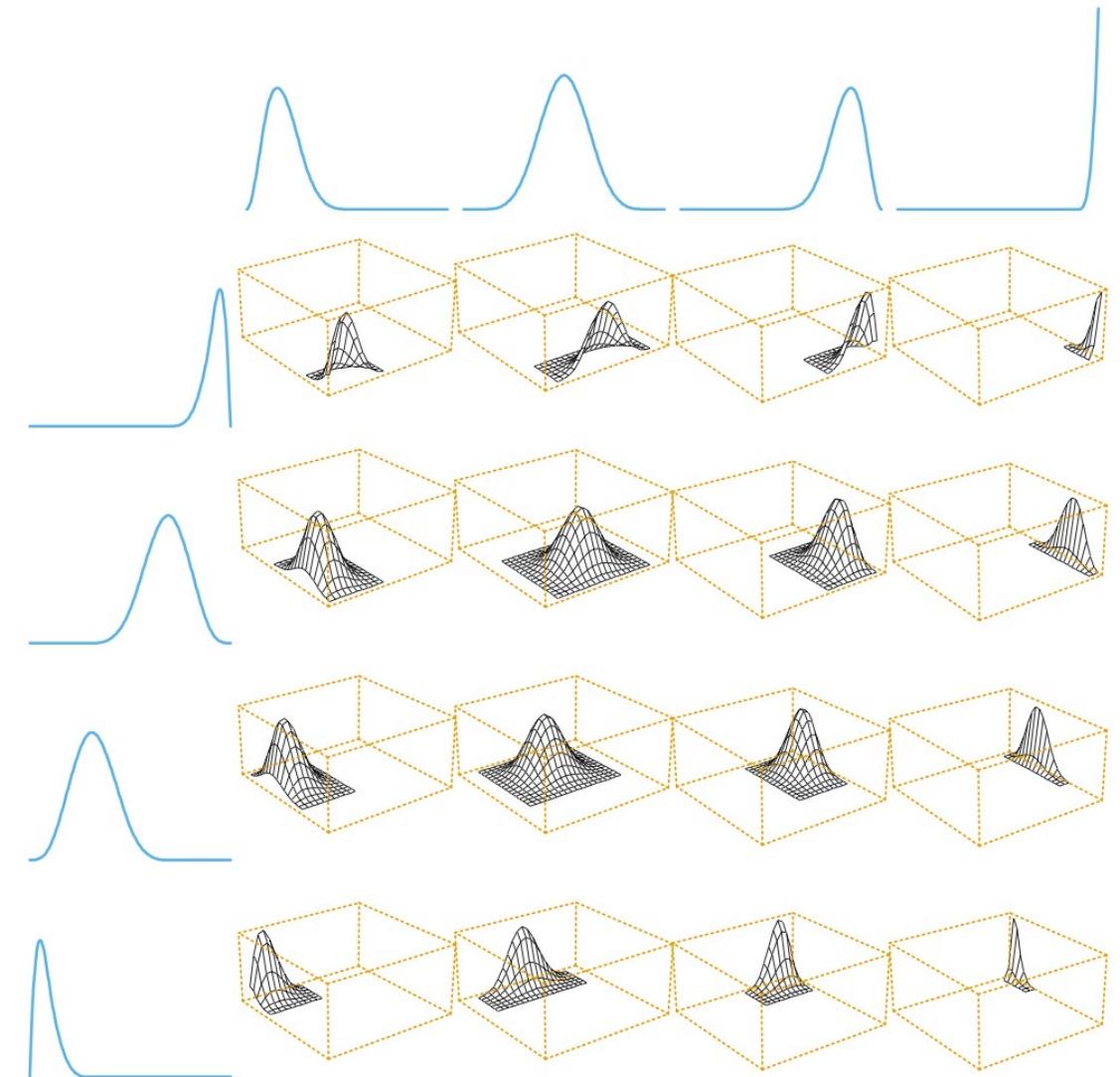
함수의 n차 미분
smoothness 가정!

When $p = 2$, $W = H$ (Hilbert Space)

Multivariate smooth spline – Tensor product basis

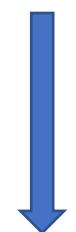
$$g_{jk}(X) = h_{1j}(X_1)h_{2k}(X_2), \\ j = 1, \dots, M_1, k = 1, \dots, M_2$$

, where $h_{1k}(X_1)$ is a basis function of X_1 ,
 $h_{2k}(X_2)$ is a basis function of X_2



Multivariate smooth spline

$$\min_f Q_\lambda(f) = \sum_i^N (y_i - f(x_i))^2 + \lambda J(f)$$



$$J(f) = \iint \left[\left(\frac{\partial^2 f(x)}{\partial x_1^2} + \frac{\partial^2 f(x)}{\partial x_2^2} + 2 \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \right)^2 \right] dx_1 x_2$$

$$f(x) = \beta_0 + \beta' x + \sum_j^N \alpha_j h_j(x)$$

$$h_j(x) = \left| |x - x_j| \right|^2 \log | |x - x_j| | \text{ : Radial Basis}$$



Part 4

Wavelet Spline

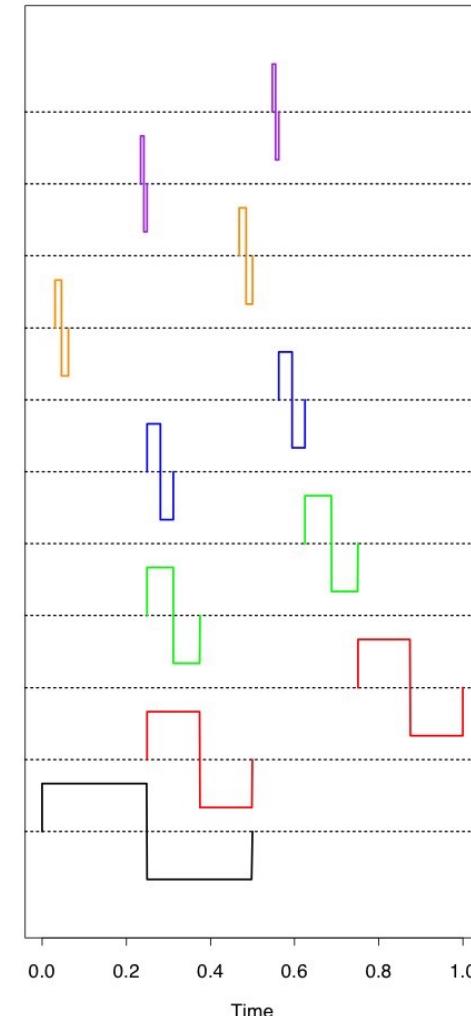
Wavelet Smoothing

“Signal Processing”

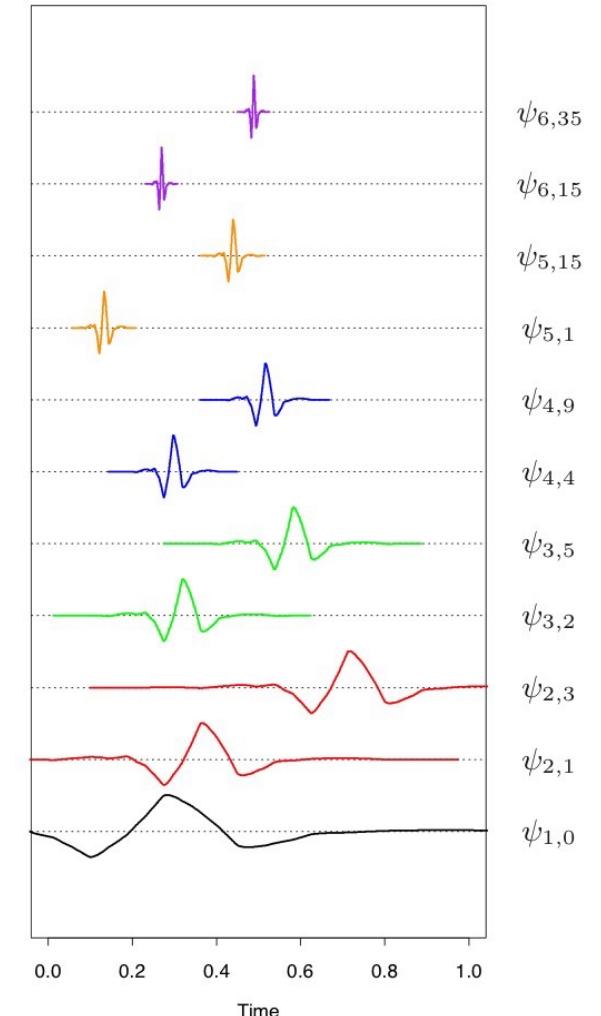


- Wavelet Bases
- Wavelet smoothing

Haar Wavelets



Symmlet-8 Wavelets



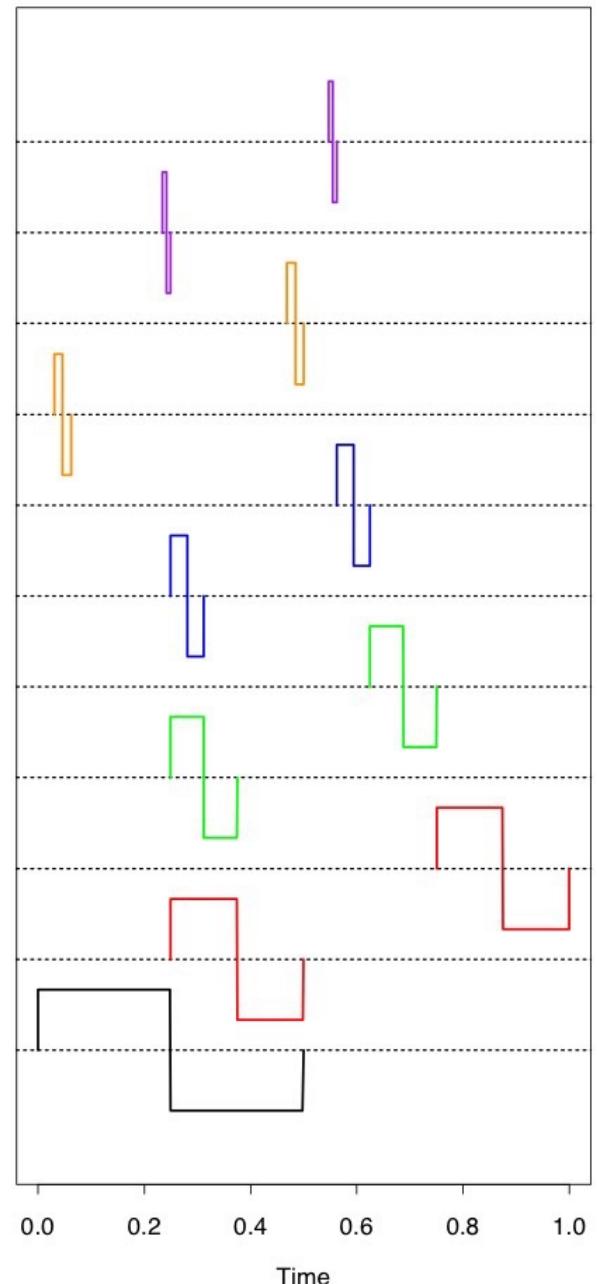
Wavelet Bases

Let $\phi(x) = I(x \in [0,1])$.
 $\phi_{0,k} = \phi(x - k)$
: orthonormal basis of reference space V_0

The dilations $\phi_{1,k}(x) = \sqrt{2} \phi(2x - k)$: orthonormal basis of $V_1 \supset V_0$

$\dots \supset V_1 \supset V_0 \supset V_{-1} \supset \dots$
each $V_i = \text{span}\{\phi_{i,k} 2^{\frac{i}{2}} \phi(2^i x - k)\}$

Haar Wavelets



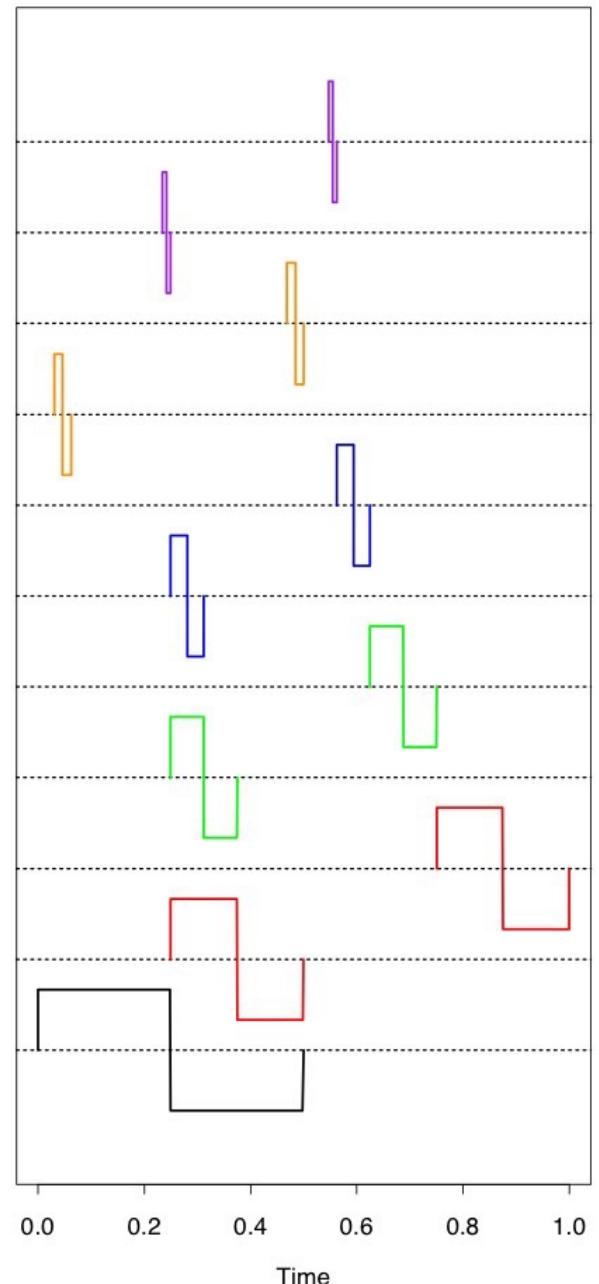
Wavelet Bases

$\psi(x) = \phi(2x) - \phi(2x - 1)$
: an orthonormal basis of W_0

inductively : $\psi_{j,k}$
 $= 2^{\frac{j}{2}}\psi(2^j x - k)$: *a basis for W_j*

W : “Detail”
 $V_J = V_0 \oplus W_0 \oplus W_1 \oplus \dots \oplus W_{J-1}$

Haar Wavelets



Adaptive Wavelet Filtering

Wavelet Transform

$y^* = W'y$, where W is the $N * N$ orthonormal basis

Inverse Wavelet Transform $\hat{f} = W\hat{\theta}$

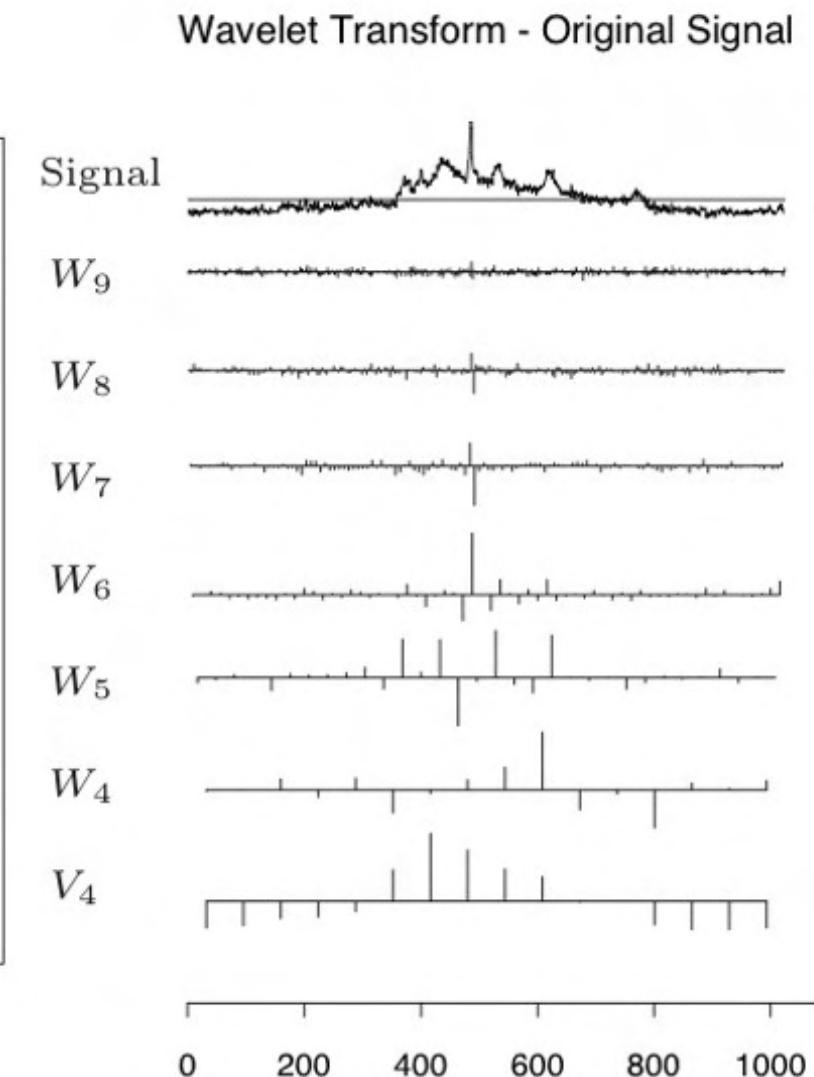
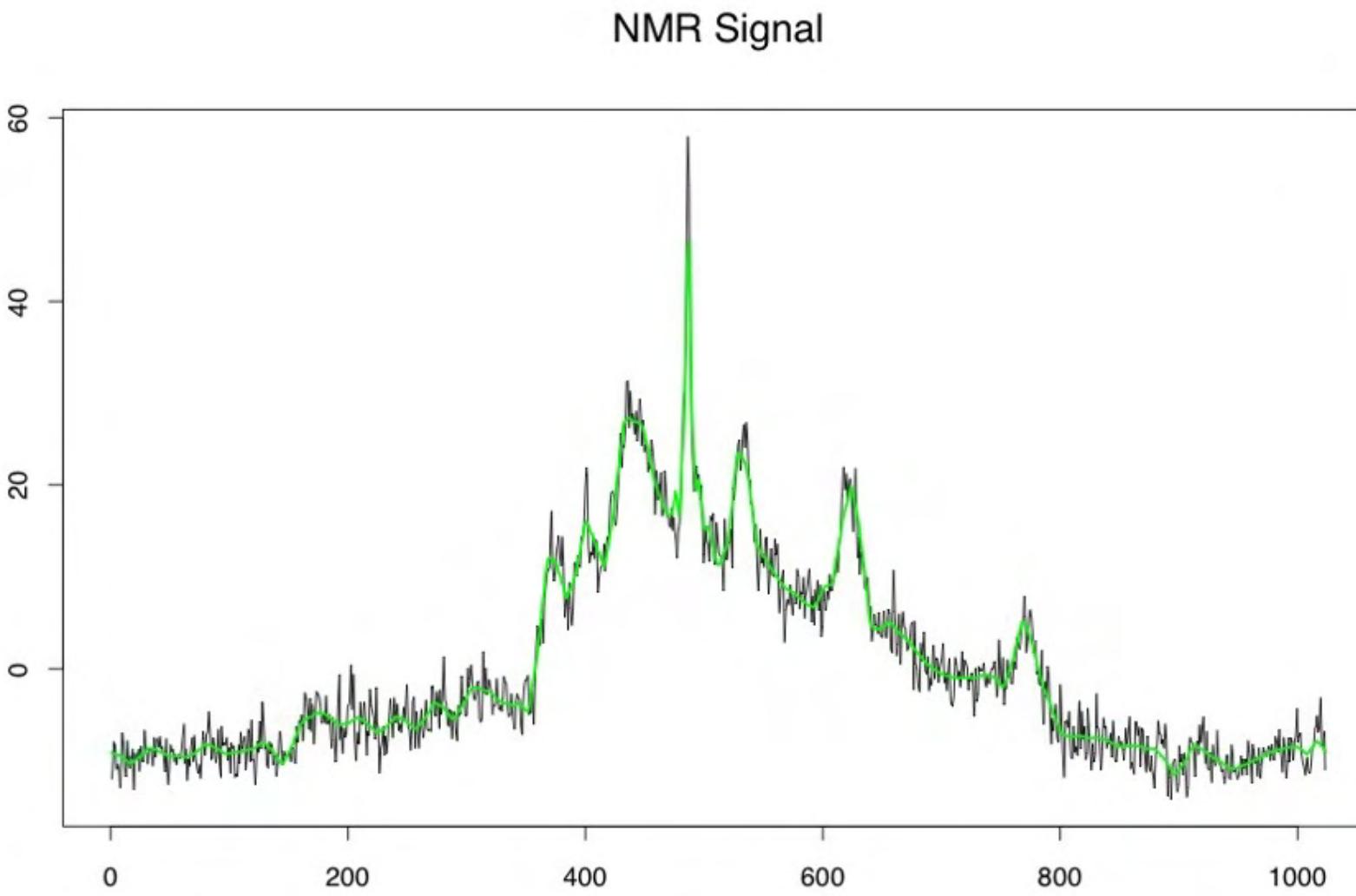


$$\lambda = \sigma\sqrt{2\log N} : Simple$$

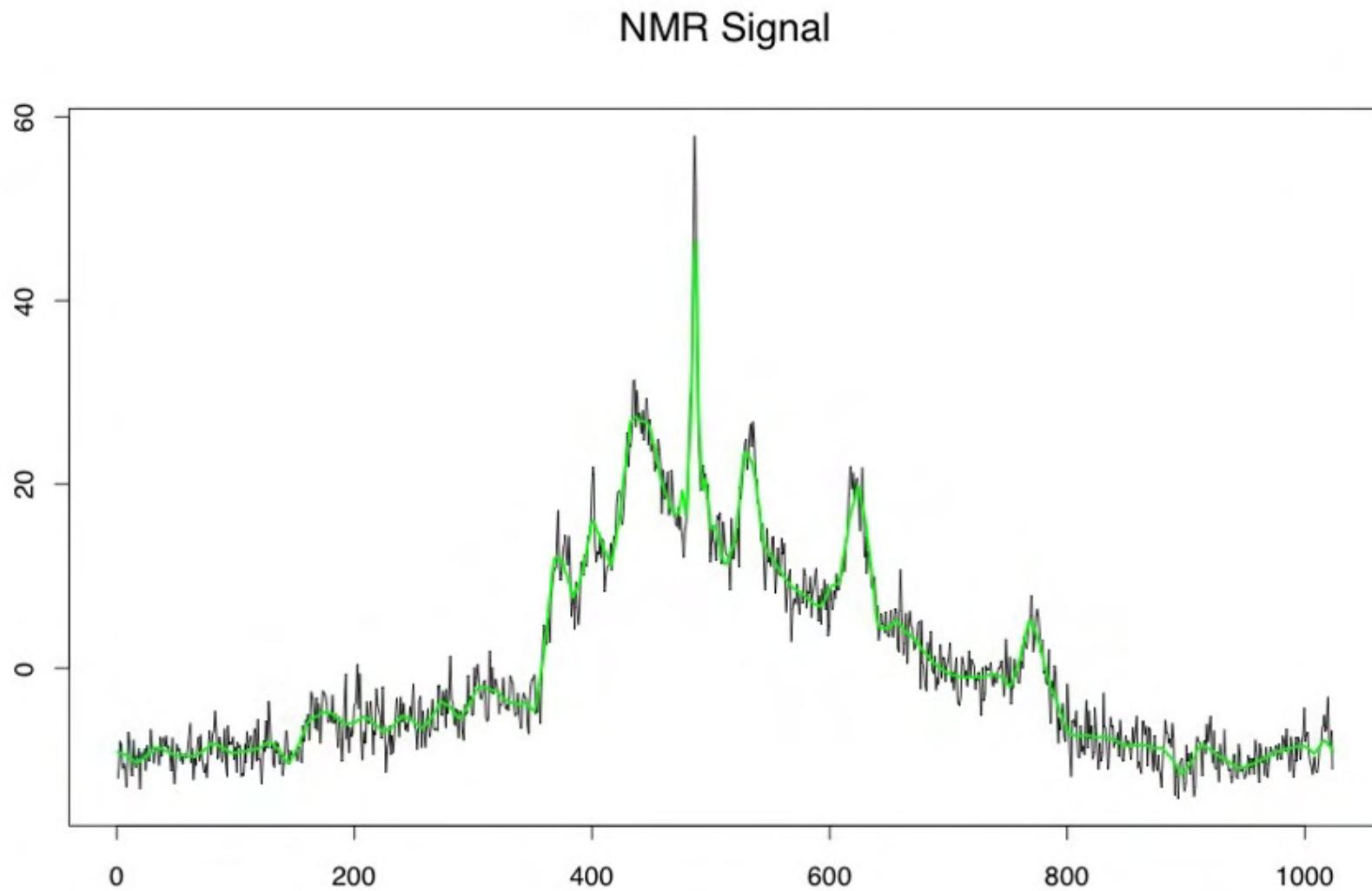
SURE shrinkage : Stein Unbiased Risk Estimation

$$\min_{\theta} ||y - W\theta||_2^2 + 2\lambda||\theta||_1 \longrightarrow \hat{\theta}_j = sign(y_j^*)(|y_j^*| - \lambda)_+$$

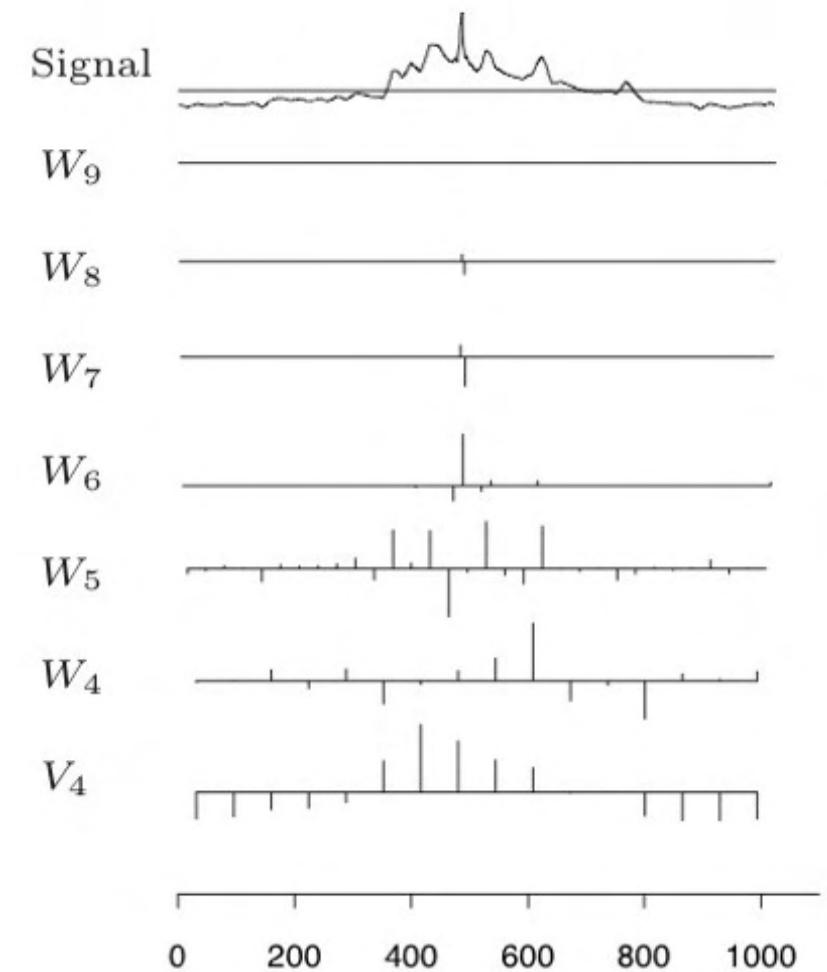
Wavelet Smoothing



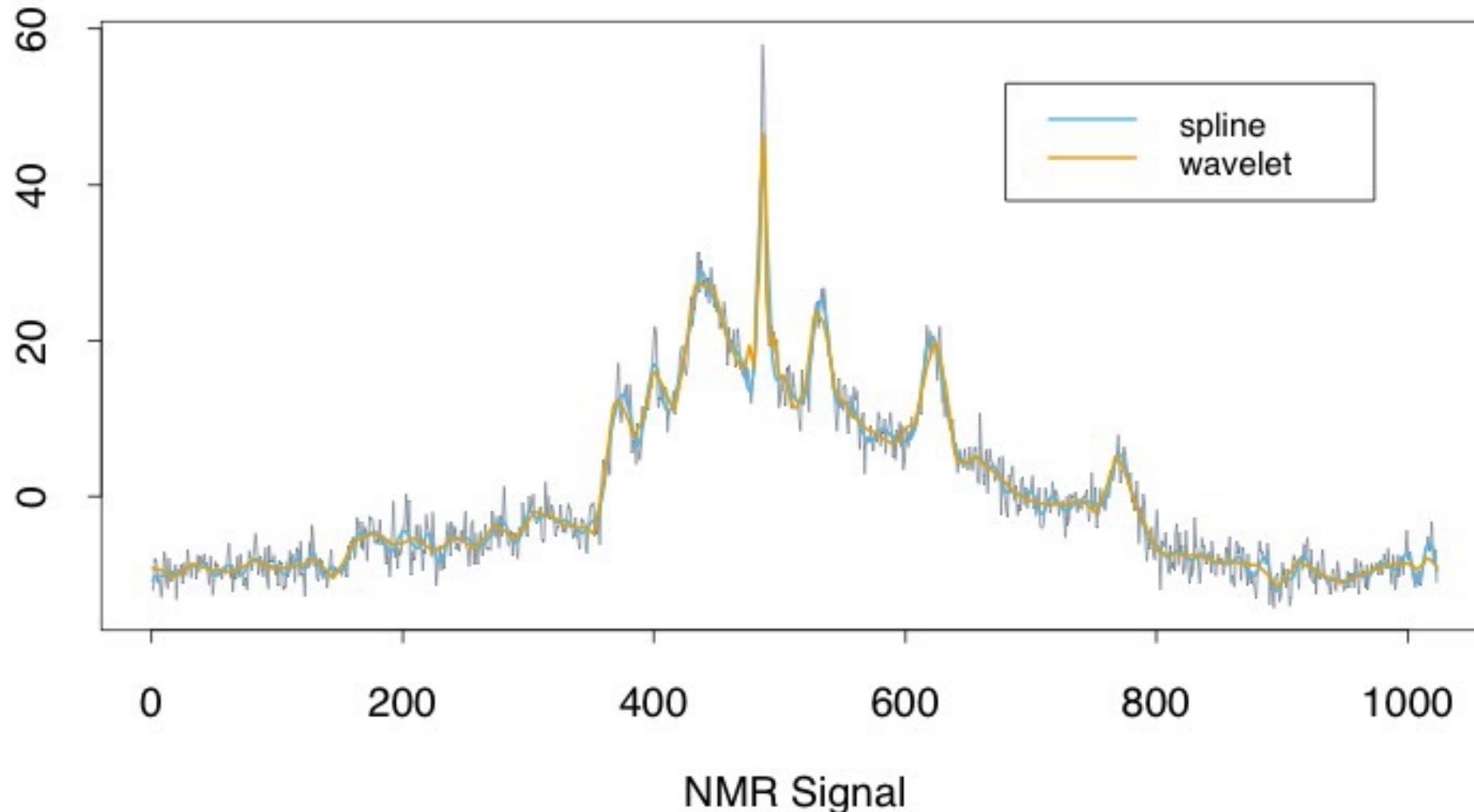
Wavelet Smoothing



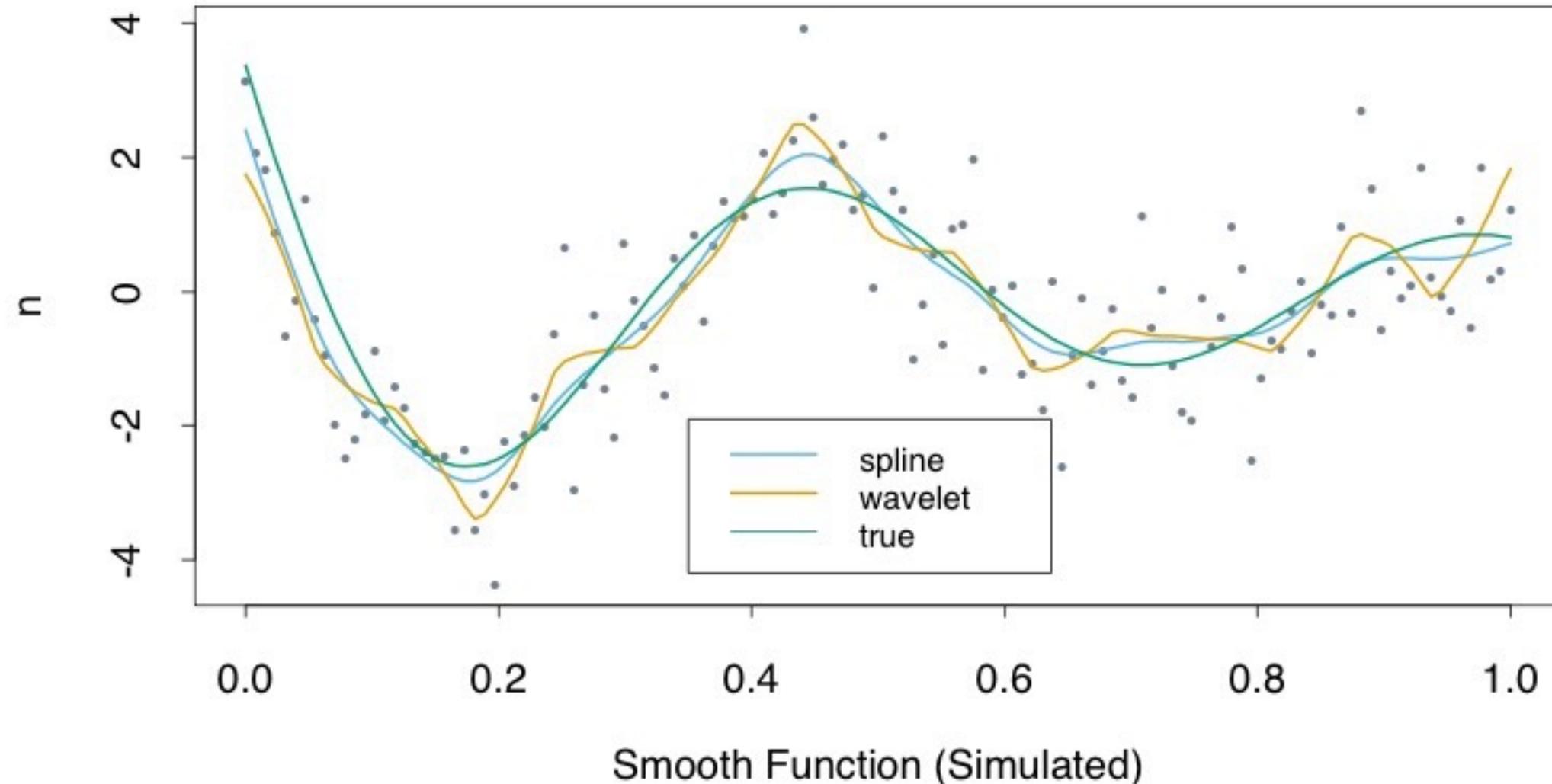
Wavelet Transform - WaveShrunk Signal



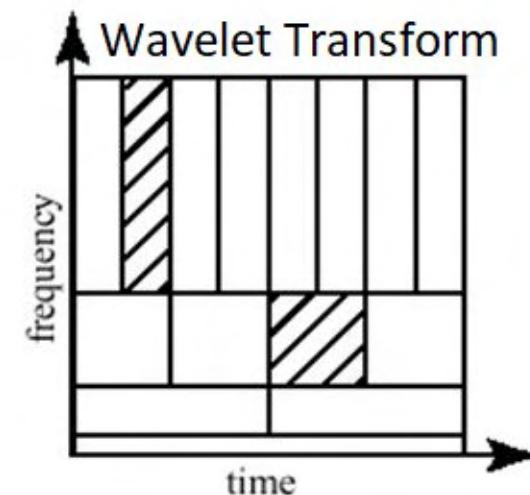
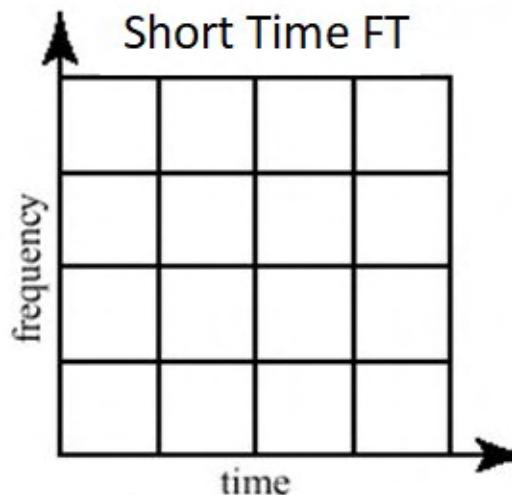
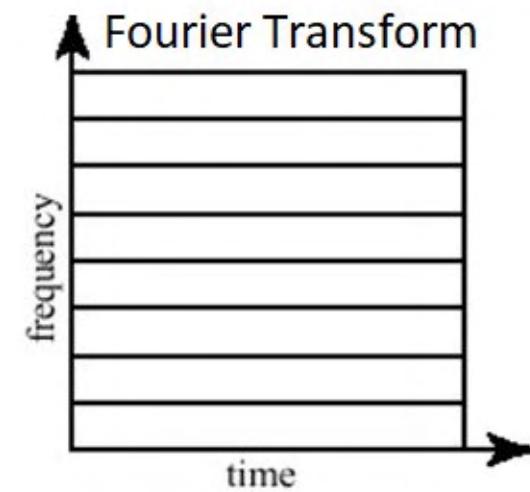
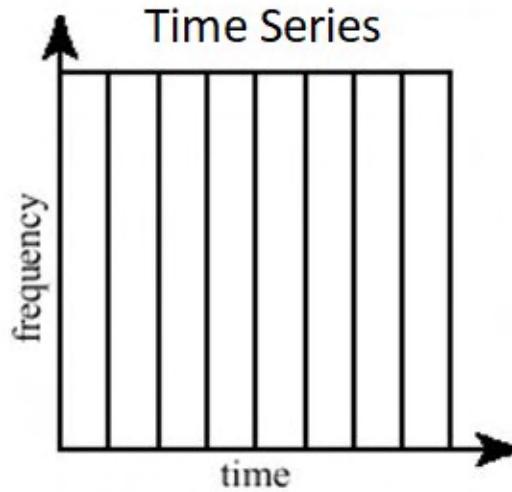
Comparison : Smooth Spline vs Wavelet Spline



Comparison : Smooth Spline vs Wavelet Spline



Comparison : FFT vs Wavelet Spline



Time Complexity Difference

FFT : $O(N \log N)$

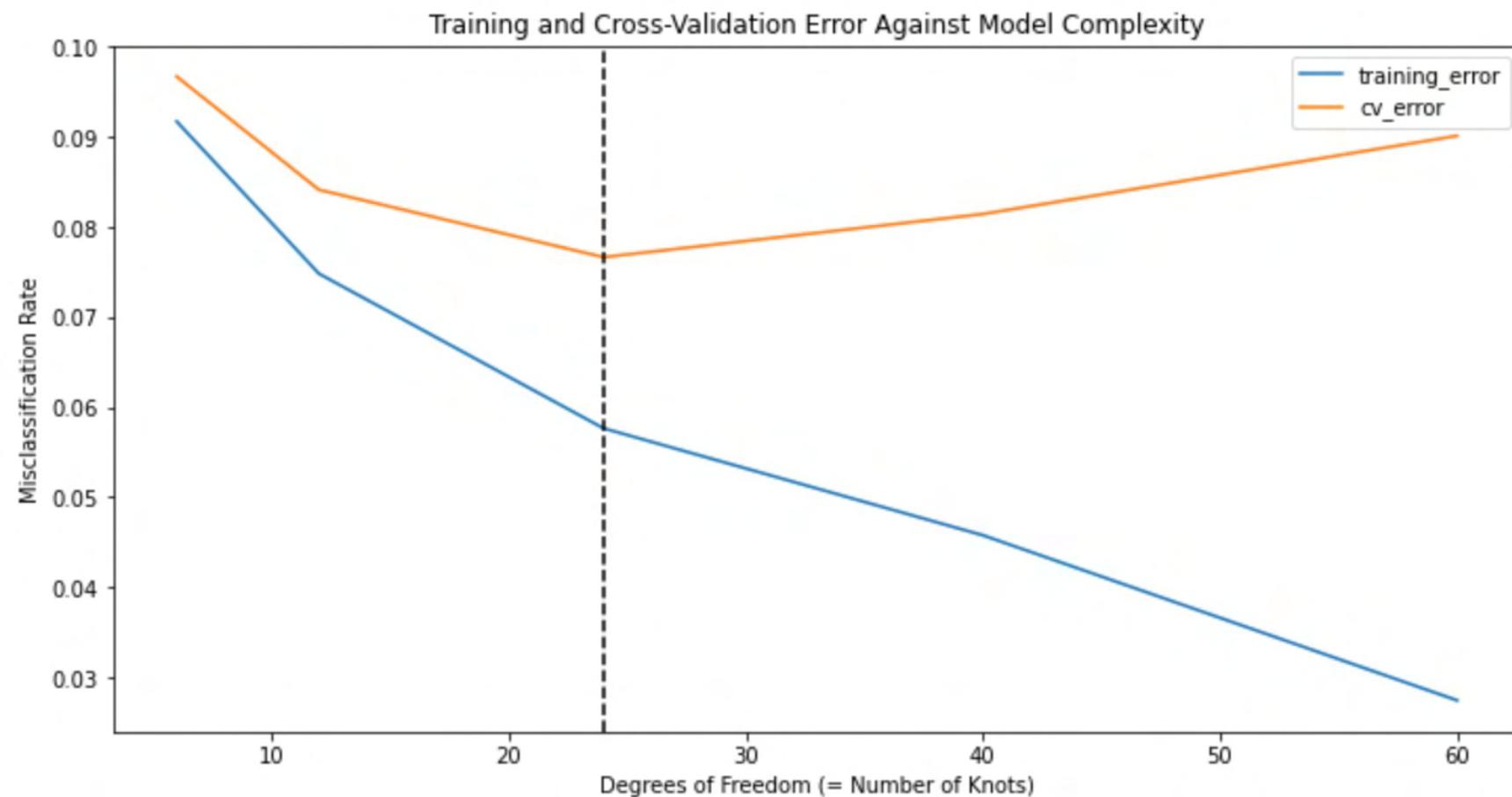
VS

Wavelet Transform : $O(N)$
Inverse Wavelet Transform : $O(N)$

Homework - ESL

Phoneme data

자세한 문제 설명은
깃허브에!



Homework - ESL

Ex. 5.12 Characterize the solution to the following problem,

$$\min_f \text{RSS}(f, \lambda) = \sum_{i=1}^N w_i \{y_i - f(x_i)\}^2 + \lambda \int \{f''(t)\}^2 dt, \quad (5.73)$$

where the $w_i \geq 0$ are observation weights.

Characterize the solution to the smoothing spline problem (5.9) when the training data have ties in X .

References

- <https://www.youtube.com/watch?v=EsOaXmCshm4&index=13&list=PLUXWLrdMaU9vK0mvcRRbg83RqwFvTjHJq&t=0s>
 - <https://jiminsun.github.io/2018-05-27/RKHS/>
 - <https://ngilshie.github.io/jekyll/update/2018/02/01/RKHS.html>
 - <https://ngilshie.github.io/jekyll/update/2018/02/01/RKHS.html>
 - <https://www.youtube.com/watch?v=lgvcyCPxNDQ>
 - <https://www.youtube.com/watch?v=0aZHkrF55cc>



Q&A