

# Nonparametric identification of LTIs: theory and practice

Gianluigi Pillonetto

Department of Information Engineering – University of Padova

May 8<sup>th</sup>, 2012



DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE



## Question

$$\begin{cases} y_1 &= \theta_1 + \nu_1 \\ &\vdots \\ y_N &= \theta_N + \nu_N \end{cases}$$

$$N \geq 3$$

$\theta_i$  deterministic

$$\nu_i \sim \mathcal{N}(0, \sigma^2)$$

$$\nu_i \perp \nu_j$$

Find  $\hat{\boldsymbol{\theta}}$  minimizing  $\mathbb{E} \left[ \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|_2^2 \right]$

## Question

$$\begin{cases} y_1 &= \theta_1 + \nu_1 \\ &\vdots \\ y_N &= \theta_N + \nu_N \end{cases} \quad \begin{array}{l} N \geq 3 \\ \theta_i \text{ deterministic} \\ \nu_i \sim \mathcal{N}(0, \sigma^2) \\ \nu_i \perp \nu_j \end{array}$$

Find  $\hat{\boldsymbol{\theta}}$  minimizing  $\mathbb{E} \left[ \left\| \boldsymbol{\theta} - \hat{\boldsymbol{\theta}} \right\|_2^2 \right]$

## Answer?

$$\hat{\boldsymbol{\theta}} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \quad \begin{array}{l} \text{ML} \\ \text{MVUE} \\ \text{efficient} \end{array}$$

# The James-Stein estimator (1956)

$$\hat{\boldsymbol{\theta}}_{JS} = \left(1 - \frac{(M-2)\sigma^2}{\|\mathbf{y}\|_2^2}\right) \mathbf{y}$$

- its MSEs always better than LSs ones, for every  $\boldsymbol{\theta} \in \mathbb{R}^M$
- its MSEs tend to LSs ones when  $\|\boldsymbol{\theta}\|_2$  is large

# Why?

$$\hat{\boldsymbol{\theta}}_{JS} = \left(1 - \frac{(M-2)\sigma^2}{\|\mathbf{y}\|_2^2}\right) \mathbf{y}$$

## The Bias – Variance dilemma

$$\mathbb{E} \left[ \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2^2 \right] = \left\| \mathbb{E} [\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}] \right\|_2^2 + \mathbb{E} \left[ \|\hat{\boldsymbol{\theta}} - \mathbb{E} [\hat{\boldsymbol{\theta}}]\|_2^2 \right]$$

# Why?

$$\hat{\boldsymbol{\theta}}_{JS} = \left( 1 - \frac{(M-2)\sigma^2}{\|\mathbf{y}\|_2^2} \right) \mathbf{y}$$

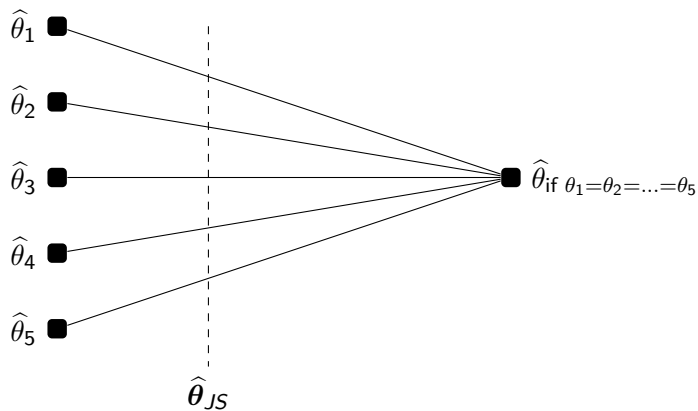
## The Bias – Variance dilemma

$$\mathbb{E} \left[ \|\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}\|_2^2 \right] = \left\| \mathbb{E} [\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}] \right\|_2^2 + \mathbb{E} \left[ \|\hat{\boldsymbol{\theta}} - \mathbb{E} [\hat{\boldsymbol{\theta}}]\|_2^2 \right]$$

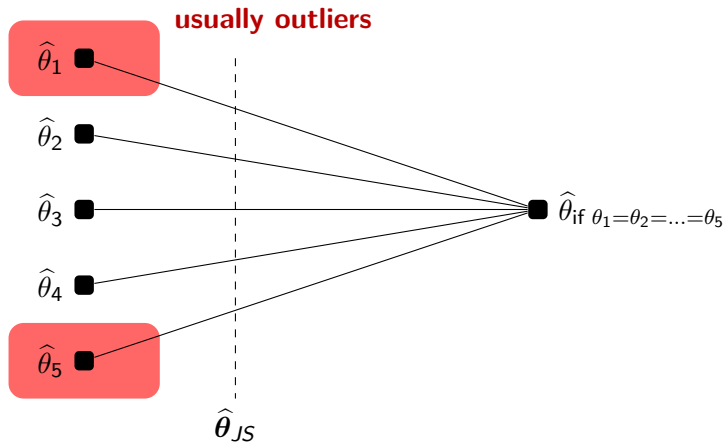
$\hat{\boldsymbol{\theta}}_{JS}$  is a regularized estimator

Rule-of-thumb: more regularization  $\Rightarrow \begin{cases} \text{more bias} \\ \text{less variance} \end{cases}$

# Graphical intuition



# Graphical intuition



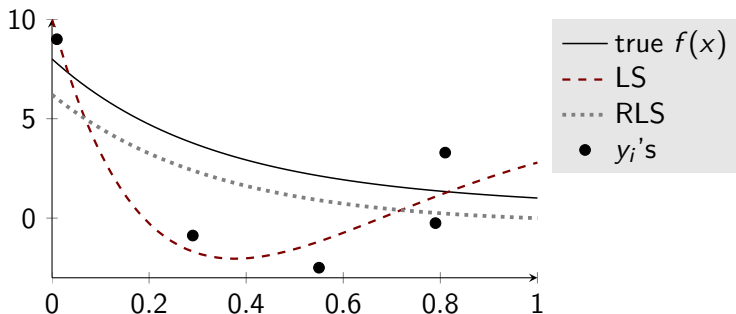


## Some comments

$$\hat{\theta}_{JS} = \left( 1 - \frac{(M-2)\sigma^2}{\|\mathbf{y}\|_2^2} \right) \mathbf{y}$$

- $\hat{\theta}_{JS}$  “learns” the mechanism generating the  $y_m$ ’s  
(connections with *Empirical Bayes*)
- improves the **total MSE** , not the MSEs of the single components
- very close to LS if  $y_m$ ’s very far apart

## Regularization helps: an example



$$f(x) = \theta_1 e^{-x} + \theta_2 e^{-2x} + \theta_3 e^{-3x} \quad \theta = \begin{bmatrix} 3 \\ -4 \\ 9 \end{bmatrix} \quad \nu_i \sim \mathcal{N}(0, 9)$$

$$LS = \arg \min_{\bar{\theta} \in \mathbb{R}^3} \sum_i \left( y_i - f_{\bar{\theta}}(x_i) \right)^2 \quad RLS = \arg \min_{\bar{\theta} \in \mathbb{R}^3} \sum_i \left( y_i - f_{\bar{\theta}}(x_i) \right)^2 + 5 \left\| \bar{\theta} \right\|_2^2$$

# Examples of regularization

- Ridge regression ( $\ell_2$  norms)
- LASSO ( $\ell_1$  norms)
- Elastic network (combination of  $\ell_1$  and  $\ell_2$  norms)
- ...

# Examples of regularization

- Ridge regression ( $\ell_2$  norms)
- LASSO ( $\ell_1$  norms)
- Elastic network (combination of  $\ell_1$  and  $\ell_2$  norms)
- ...

regularization is most useful  
when problems are ill-conditioned

# Examples of regularization

- Ridge regression ( $\ell_2$  norms)
- LASSO ( $\ell_1$  norms)
- Elastic network (combination of  $\ell_1$  and  $\ell_2$  norms)
- ...

regularization is most useful  
when problems are ill-conditioned

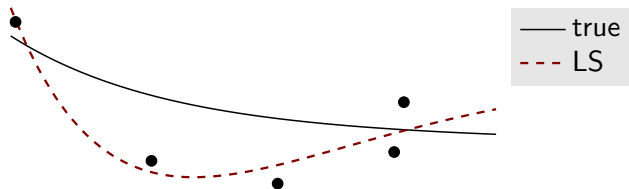
here we focus on Ridge regression

## From parametric to nonparametric

previous aim: introduce regularization

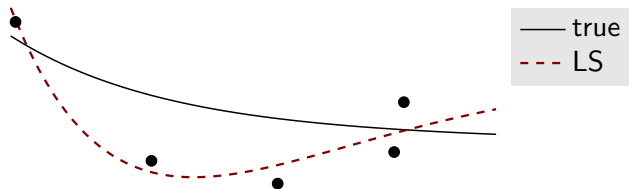
next aim: introduce nonparametric regression

# Nonparametric approaches – motivations 1



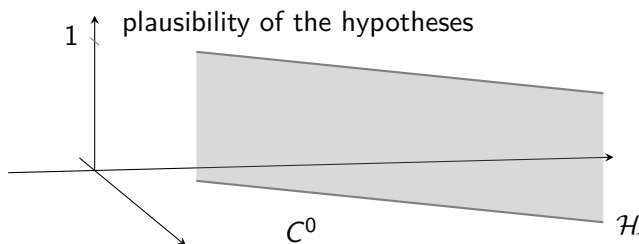
Key fact in previous example:  $\mathcal{H} = \text{span} \{e^{-x}, e^{-2x}, e^{-3x}\}$

# Nonparametric approaches – motivations 1



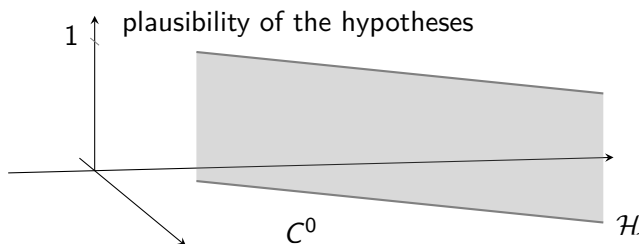
Key fact in previous example:  $\mathcal{H} = \text{span} \{e^{-x}, e^{-2x}, e^{-3x}\}$

peculiarity of parametric approaches!





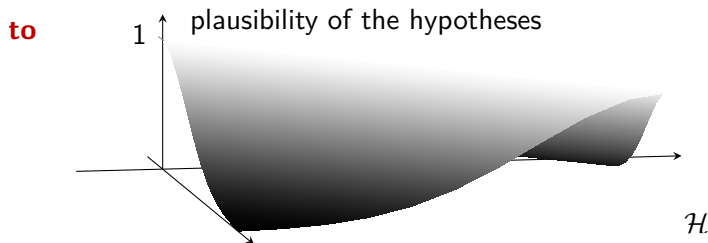
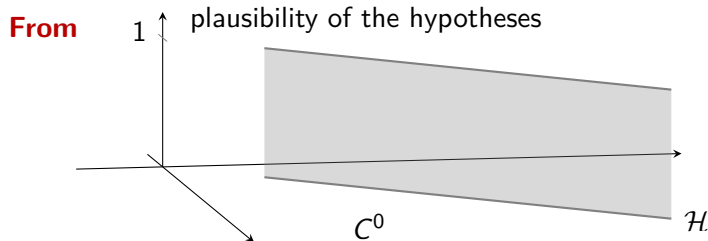
# Nonparametric approaches – motivations 2



## (some) drawbacks of parametric approaches

- require high levels of prior knowledge
- complex systems may lead to proliferation of parameters (*high variance!*)

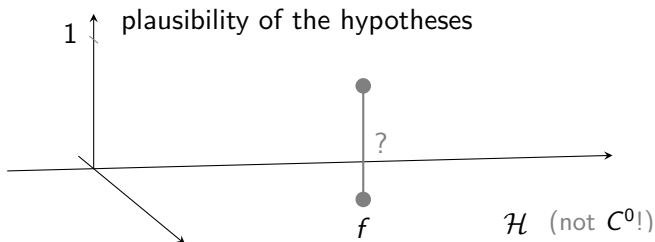
# Towards nonparametric approaches



how should we define the novel plausibility?

plausibility = regularity

plausibility = regularity



concept of regularity depends on prior assumptions!

intuitive examples:

$$\|f\|_2^2 = \int_{\mathcal{X}} f(x)^2 dx$$

$$\|f\|_{\text{cub.sp.}}^2 = \int_{\mathcal{X}} \ddot{f}(x)^2 dx$$

**assumption:**  $f \in \mathcal{H}_\star$  and regularity of  $f := \|f\|_\star$

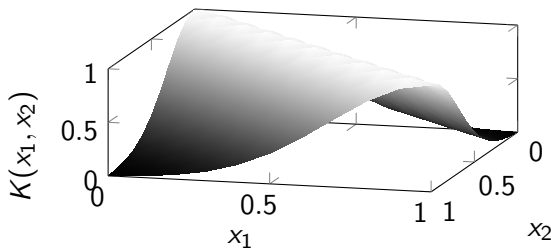
**question:** what can  $\star$  be?

# RKHSs – the key ingredient

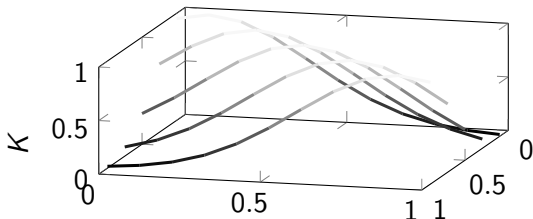
$$\star = K \quad K : \mathcal{X} \times \mathcal{X} \mapsto \mathbb{R} \quad K : \begin{cases} \text{symmetric} \\ \text{continuous} \\ \text{positive definite} \end{cases}$$

$(\mathcal{X} = \text{domain of } f)$

Example: Gaussian Kernel  $K(x_1, x_2) = \exp\left(-\frac{(x_1 - x_2)^2}{2\sigma^2}\right)$



$$K \leftrightarrow \mathcal{H}_K := \overline{\text{span}\{K(x, \cdot) \text{ s.t. } x \in \mathcal{X}\}}$$



RKHSs = spaces of functions where the evaluation functional is bounded and linear

Remarks:

- $\mathcal{H}_K \subset C^0$
- $f(\cdot) = \sum_i a_i K(x_i, \cdot) \Rightarrow \|f\|_K^2 = \sum_{i,j} a_i a_j K(x_i, x_j)$

## Kernels are the goggles with which one sees the world

$$f(\cdot) = \sum_i a_i K(x_i, \cdot) \Rightarrow \|f\|_K^2 = \sum_{i,j} a_i a_j K(x_i, x_j)$$

same  $f$ , different  $K$ 's, different  $\|f\|_K$ 's



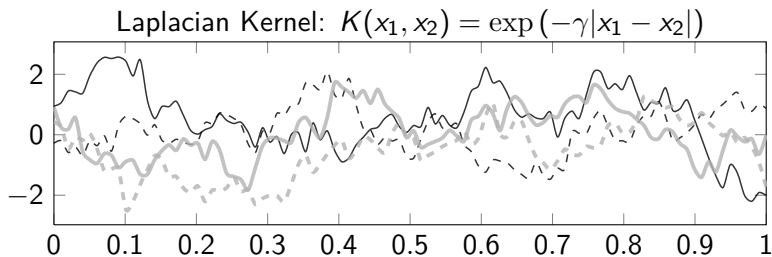
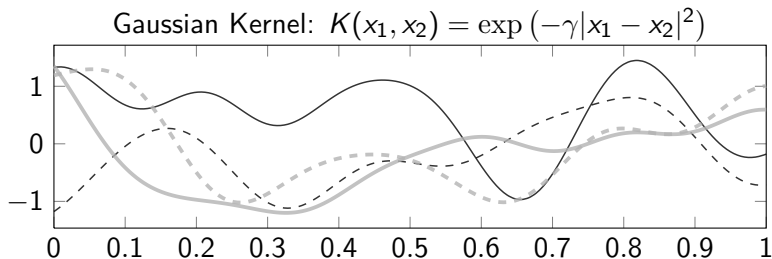
# Kernels are the goggles with which one sees the world

$$f(\cdot) = \sum_i a_i K(x_i, \cdot) \Rightarrow \|f\|_K^2 = \sum_{i,j} a_i a_j K(x_i, x_j)$$

same  $f$ , different  $K$ 's, different  $\|f\|_K$ 's

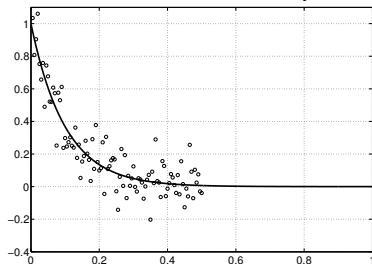
linear $x_1 x_2 + c$	polynomial $(\alpha x_1 x_2 + c)^d$	Gaussian $\exp(-\gamma(x_1 - x_2)^2)$
Laplacian $\exp(-\gamma(x_1 - x_2))$	hyperbolic $\tanh((\alpha x_1 x_2 + c)^d)$	rational quadratic $1 - \frac{ x_1 - x_2 ^2}{ x_1 - x_2 ^2 + c}$
multiquadratic $\sqrt{ x_1 - x_2 ^2 + c}$	inv. multiquadratic $(\sqrt{ x_1 - x_2 ^2 + c})^{-1}$	wave $\frac{\theta}{ x_1 - x_2 } \sin\left(\frac{ x_1 - x_2 }{\theta}\right)$
$\vdots$	$\vdots$	$\vdots$

## Examples of typical elements of $\mathcal{H}_K$

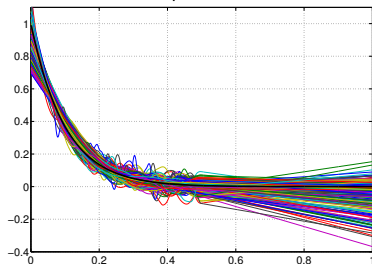


# Examples of different approximation capabilities

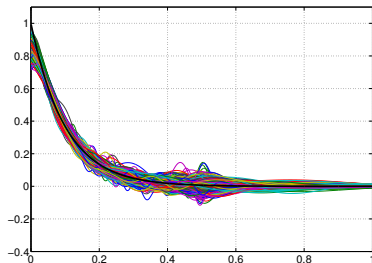
True function and a set of noisy data



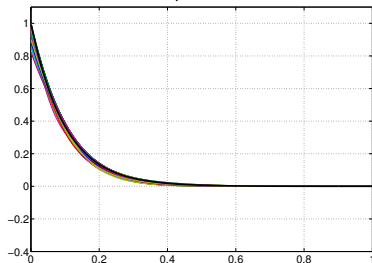
Cubic spline estimator



Gaussian kernel-based estimator



Stable spline estimator



## Example of regression with RKHSs

$$f^* = \arg \min_{f \in \mathcal{H}_K} \sum_i \left( y_i - L_i[f] \right)^2 + \gamma \|f\|_K^2$$

*tradeoff between fitting and regularity  
– same as before!!*

## Example of regression with RKHSs

$$f^* = \arg \min_{f \in \mathcal{H}_K} \sum_i \left( y_i - L_i[f] \right)^2 + \gamma \|f\|_K^2$$

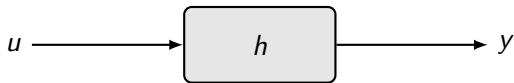
*tradeoff between fitting and regularity  
– same as before!!*

This case = Regularization Network

$$f^*(\cdot) = \sum_i c_i L_i[K(\cdot, \cdot)]$$

$$\begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = \left( \begin{bmatrix} L_1[L_1[K]] & \cdots & L_1[L_M[K]] \\ \vdots & & \vdots \\ L_M[L_1[K]] & \cdots & L_M[L_M[K]] \end{bmatrix} + \gamma I \right)^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

## An important example



$$y_t = L_t[f] = (u * h)_t$$

convolution is a linear functional

## important differences

Parametric:  $\theta^* = \arg \min_{\theta \in \Theta} \sum_i \left( y_i - L_i[f_\theta] \right)^2$

Nonparametric:  $f^* = \arg \min_{f \in \mathcal{H}_K} \sum_i \left( y_i - L_i[f] \right)^2 + \gamma \|f\|_K^2$

# System Identification as Nonparametric regression

1<sup>st</sup> requirement: select the most appropriate  $K(\cdot, \cdot)$

*assumption: our system is LTI BIBO stable*



# System Identification as Nonparametric regression

1<sup>st</sup> requirement: select the most appropriate  $K(\cdot, \cdot)$

***assumption: our system is LTI BIBO stable***

Translating the available a-priori information

- impulse response  $g \in \mathcal{L}^1$ , i.e.,  $\int_{\mathbb{R}^+} |g(x)| dx < +\infty$

# System Identification as Nonparametric regression

1<sup>st</sup> requirement: select the most appropriate  $K(\cdot, \cdot)$

***assumption: our system is LTI BIBO stable***

Translating the available a-priori information

- impulse response  $g \in \mathcal{L}^1$ , i.e.,  $\int_{\mathbb{R}^+} |g(x)| dx < +\infty$

**Definition**

$K(\cdot, \cdot)$  is said a **stable kernel** if  $\mathcal{H}_K \subset \mathcal{L}^1$

# Characterization of the stable kernels (1)

## Proposition

$$\int_{\mathbb{R}^+} \int_{\mathbb{R}^+} |K(x_1, x_2)| dx_1 dx_2 < +\infty \quad \Rightarrow \quad \mathcal{H}_K \subset \mathcal{L}^1$$

## Proposition

If  $K(x_1, x_2) \geq 0$  for all  $x_1, x_2 \in \mathbb{R}^+$  then

$$\int_{\mathbb{R}^+} \int_{\mathbb{R}^+} |K(x_1, x_2)| dx_1 dx_2 < +\infty \quad \Leftrightarrow \quad \mathcal{H}_K \subset \mathcal{L}^1$$

## Characterization of the stable kernels (2)

### Definition

Let  $q = \frac{p}{p-1}$ .  $K(\cdot, \cdot)$  is said  **$q$ -bounded** if

- $K(x, \cdot) \in \mathcal{L}^p$  for almost all  $x \in \mathbb{R}^+$
- $f \in \mathcal{L}^q \Rightarrow g(x) := \int_{\mathbb{R}^+} K(x, a)f(a)da \in \mathcal{L}^p$

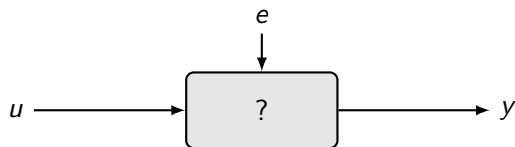
### Proposition

$$K \text{ } q\text{-bounded} \Leftrightarrow \mathcal{H}_K \subset \mathcal{L}^p$$

thus

$$\begin{aligned} K \text{ } \infty\text{-bounded} &\Leftrightarrow \int_{\mathbb{R}^+} \left| \int_{\mathbb{R}^+} K(x, a)f(a)da \right| dx < +\infty \quad \forall f \in \mathcal{L}^\infty \\ &\Leftrightarrow \mathcal{H}_K \subset \mathcal{L}^1 \end{aligned}$$

# Nonparametric identification of LTIs



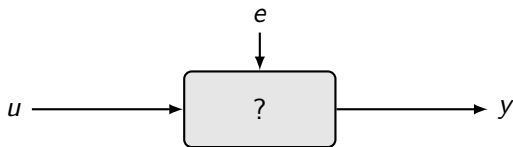
$$y_t = \sum_{i=1}^{\infty} f_i u_{t-i} + \sum_{i=1}^{\infty} g_i e_{t-i},$$

$$\text{dataset} = \{u_t\}, \{y_t\}$$

PEM approach:

$$\hat{y}_{t|t-1} = \sum_{i=1}^{\infty} h_i^u u_{t-i} + \sum_{i=1}^{\infty} h_i^y y_{t-i}$$

# Nonparametric identification of LTIs



$$y_t = \sum_{i=1}^{\infty} f_i u_{t-i} + \sum_{i=1}^{\infty} g_i e_{t-i},$$

$$\text{dataset} = \{u_t\}, \{y_t\}$$

PEM approach:

$$\hat{y}_{t|t-1} = \sum_{i=1}^{\infty} h_i^u u_{t-i} + \sum_{i=1}^{\infty} h_i^y y_{t-i}$$

SysId with Regularization Networks??

$$h^* = \arg \min_{h \in \star} \sum_t \left( y_t - \hat{y}_{t|t-1} \right)^2 + \gamma \|h\|^2 \star$$

# The Stable Splines Kernel



G. Pillonetto, G. De Nicolao (Automatica 2010)

A new kernel-based approach for linear system identification

$$\text{cubic splines: } W(s, t) = \begin{cases} \frac{s^2}{2} \left( t - \frac{s}{3} \right) & \text{if } s \leq t \\ \frac{t^2}{2} \left( s - \frac{t}{3} \right) & \text{if } s > t \end{cases}$$

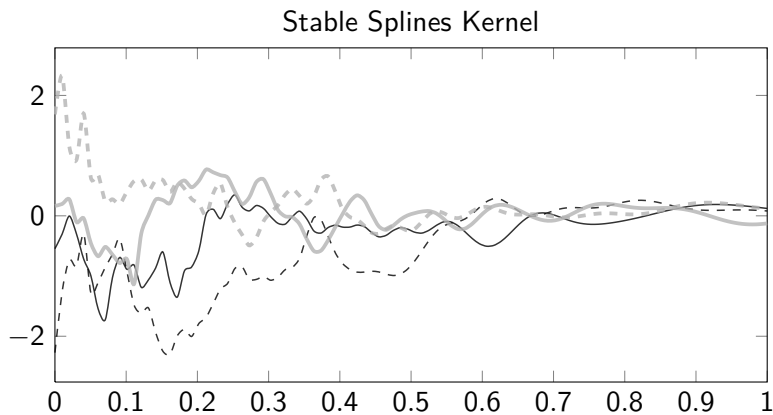
$$\text{stable splines: } K(x_1, x_2) = W(e^{-\beta x_1}, e^{-\beta x_2})$$

## Bayesian interpretation

Let  $f \sim \mathcal{GP}(0, K)$ . Then

$$\mathbb{P}[f = \text{imp. resp. of LTI BIBO stable system}] = 1$$

## Examples of typical elements of $\mathcal{H}_K$





*how to actually perform the identification*

# The model & the hyperparameters

$$A(z)y(t) = B(z)u(t) + C(z)e(t)$$

MISO  $\Rightarrow$  more than one impulse response!

$i :=$  impulse response index

# The model & the hyperparameters

$$A(z)y(t) = B(z)u(t) + C(z)e(t)$$

MISO  $\Rightarrow$  more than one impulse response!

$i :=$  impulse response index

$$h^i \sim \mathcal{GP}(0, \lambda_i^2 K(\cdot, \cdot; \beta))$$

# The model & the hyperparameters

$$A(z)y(t) = B(z)u(t) + C(z)e(t)$$

MISO  $\Rightarrow$  more than one impulse response!

$i :=$  impulse response index

$$h^i \sim \mathcal{GP}(0, \lambda_i^2 K(\cdot, \cdot; \beta))$$

$\lambda_i^2$ : “amplitude” of the  $i$ -th impulse response

$\beta$ : decay ratio

# Estimation of the hyperparameters

## Empirical Bayes in theory

- ① assume existence of prior distribution with *unknown hyperparameters*

## Empirical Bayes in practice (*with some abuses of notation*)

- ①  $p(y|h)$ ,  $p(h|\lambda)$  are known,  $\lambda$  unknown

# Estimation of the hyperparameters

## Empirical Bayes in theory

- 1 assume existence of prior distribution with *unknown hyperparameters*
- 2 compute the marginal likelihood

## Empirical Bayes in practice *(with some abuses of notation)*

- 1  $p(y|h)$ ,  $p(h|\lambda)$  are known,  $\lambda$  unknown
- 2 exploiting  $p(y, h|\lambda) = p(y|h, \lambda) p(h|\lambda)$  compute

$$p(y|\lambda) = \int p(y, h|\lambda) p(h|\lambda) dh$$

# Estimation of the hyperparameters

## Empirical Bayes in theory

- 1 assume existence of prior distribution with *unknown hyperparameters*
- 2 compute the marginal likelihood
- 3 estimate the hyperparameters maximizing the marginal likelihood

## Empirical Bayes in practice *(with some abuses of notation)*

- 1  $p(y|h)$ ,  $p(h|\lambda)$  are known,  $\lambda$  unknown
- 2 exploiting  $p(y, h|\lambda) = p(y|h, \lambda) p(h|\lambda)$  compute

$$p(y|\lambda) = \int p(y, h|\lambda) p(h|\lambda) dh$$

- 3  $\lambda^* = \arg \max_{\lambda} p(y|\lambda)$

# SSpline.m: a matlab toolbox

... even if not yet publicly available

```
[M,ip,Ak]=SSpline(y,U,p,l,mv,mb,cn,r,LP,LP2,ips)
```

**M:** estimated tf (idpoly object)

**ip:** estimated hyperparameters

**Ak:** AIC of the estimated tf

**y:** measured outputs

**U:** measured inputs

**p:** max. length of the to-be estimated impulse responses

**l:** type of tf to be estimated (ARMAX / ARX / etc.)

**mv:** one  $\lambda$  in common for all the  $h^i$ 's or not

**mb:** one  $\beta$  in common for all the  $h^i$ 's or not

**cn:** identify high frequencies components

**r:** number of data to be used for estimating the hyperparameters

**LP:** obtain sparse solutions

**LP2:** obtain approximated solutions

**ips:** start optimization from the assigned initial point



the routine returns an idpoly and has a  
pre-fixed impulse responses max. length!

**Question:** at the end of the day, we obtain an object of the same  
kind of classical PEM approaches. So, why should this be  
better?

the routine returns an idpoly and has a pre-fixed impulse responses max. length!

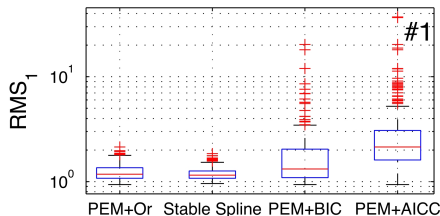
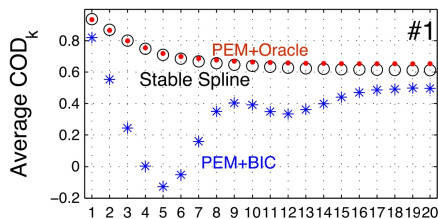
**Question:** at the end of the day, we obtain an object of the same kind of classical PEM approaches. So, why should this be better?

**Answer:** at the end of the day, ***because of the bias / variance tradeoff***

# Some comparisons

1000 MC runs, for each run:

- ARMAX order  $\in \{1, \dots, 30\}$
- ARMAX model  $\sim \text{drmodel}$
- no delays
- inputs  $\sim \text{idinput}$
- training set = 200 samples
- test set = 1000 samples



# And some conclusions

regularization usually has beneficial effects

nonparametric approaches are specially suited for complex situations

regularization in nonparametric approaches can be seen as using prior smoothness assumptions

regularization in nonparametric approaches leads to an extremely efficient LTIs sysid technique

# Nonparametric identification of LTIs: theory and practice

Gianluigi Pillonetto

Department of Information Engineering – University of Padova

May 8<sup>th</sup>, 2012