

# Auto-tuning procedures for distributed nonparametric regression algorithms

**Damiano Varagnolo**, Gianluigi Pillonetto, Luca Schenato

ECC 2015 - Linz

July 15, 2015



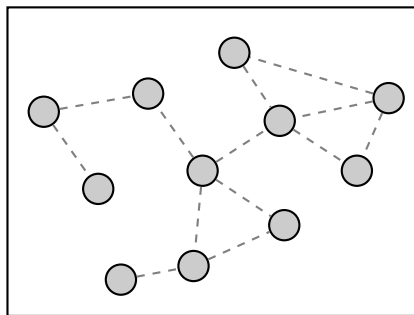
Thanks to...



Gianluigi Pillonetto  
Univ. of Padova



Luca Schenato  
Univ. of Padova



Agents:

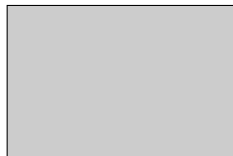
- noisily sample the same  $f$
- limited computational & communication capabilities
- 1 measurement  $\times$  agent (ease of notation)
- $M$  measurements in total

# Measurement model

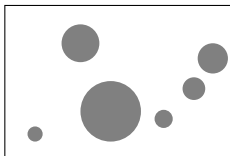
$$y_m = f(x_m) + \nu_m \quad (1)$$

- $f : \mathcal{X} \subset \mathbb{R}^d \rightarrow \mathbb{R}$  unknown ( $\mathcal{X}$  compact)
- $\nu_m \perp x_m$ , zero mean and variance  $\sigma^2$
- $x_m \sim \mu$  i.i.d. (*agents know  $\mu$ !!*)

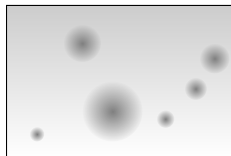
examples of  $\mu$ :



uniform

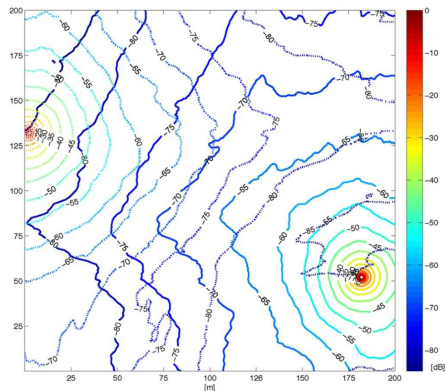


jitter



generic

# Example 1 - channel gains in geographical areas



$x \in \mathbb{R}^2$  : position  
 $t$  : time  
 $f(x, t)$  : channel gain

source: Dall'Anese et al., 2011

## Example 2 - waves power extraction



source: [www.graysharboroceanenergy.com](http://www.graysharboroceanenergy.com)

$x \in \mathbb{R}^2$  : position

$t$  : time

$f(x, t)$  : sea level

## Example 3 - multi robot exploration



source: <http://www-robotics.jpl.nasa.gov>

$x \in \mathbb{R}^2$  : position  
 $f(x)$  : ground level

## Considered cost function

$$Q(f) = \sum_{m=1}^M (y_m - f(x_m))^2 + \gamma \|f\|_K^2$$



## Considered cost function

$$Q(f) = \sum_{m=1}^M (y_m - f(x_m))^2 + \gamma \|f\|_K^2$$

↑  
lives in an infinite dimensional space

↑  
regularization factor,  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$   
Mercer kernel

## Considered cost function

$$Q(f) = \sum_{m=1}^M (y_m - f(x_m))^2 + \gamma \|f\|_K^2$$

↑  
lives in an infinite dimensional space

↑  
regularization factor,  $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$   
Mercer kernel

### Centralized optimal solution as a Regularization Network

$$f_c = \sum_{m=1}^M c_m K(x_m, \cdot) \quad \begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = \left( \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \gamma I \right)^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

## Drawbacks

$$f_c = \sum_{m=1}^M c_m K(x_m, \cdot) \quad \begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = \left( \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix} + \gamma I \right)^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

- computational cost:  $O(M^3)$  (inversion of  $M \times M$  matrix)
- transmission cost:  $O(M)$  (knowledge of whole  $\{x_m, y_m\}_{m=1}^M$ )



***need to find alternative solutions***

# Alternative centralized optimal solution (1<sup>st</sup> on 2)

Structure of  $K$  implies

- $K(x_1, x_2) = \sum_{e=1}^{+\infty} \lambda_e \phi_e(x_1) \phi_e(x_2)$   $\lambda_e = \text{eigenvalue}$   
 $\phi_e = \text{eigenfunction}$
- $f(x) = \sum_{e=1}^{+\infty} b_e \phi_e(x)$

$\Rightarrow$  measurement model can be rewritten as

$$y_m = \overbrace{[\phi_1(x_m), \phi_2(x_m), \dots]}^{C_m :=} \overbrace{\begin{bmatrix} b_1 \\ b_2 \\ \vdots \end{bmatrix}}^{b :=} + \nu_m \quad (2)$$

## Alternative centralized optimal solution (2<sup>nd</sup> on 2)

$$b_c = \left( \frac{1}{M} \text{diag} \left( \frac{\gamma}{\lambda_e} \right) + \frac{1}{M} \sum_{m=1}^M C_m^T C_m \right)^{-1} \left( \frac{1}{M} \sum_{m=1}^M C_m^T y_m \right) \quad (3)$$

involves infinite dimensional objects:

$$b_c = \begin{bmatrix} \bullet & \cdots & \cdots \\ \vdots & \ddots & \\ \vdots & & \ddots \end{bmatrix}^{-1} \begin{bmatrix} \bullet \\ \vdots \\ \vdots \end{bmatrix}$$

$\Rightarrow$  *cannot be computed exactly*

## Suboptimal finite dimensional solution

### New estimator

$$b_r = \left( \frac{1}{M} \text{diag} \left( \frac{\gamma}{\lambda_e} \right) + \frac{1}{M} \sum_{m=1}^M \left( C_m^E \right)^T C_m^E \right)^{-1} \left( \frac{1}{M} \sum_{m=1}^M \left( C_m^E \right)^T y_m \right)$$

- computable (involves  $E \times E$  matrices and  $E$ -dimensional vectors)
- minimizes  $Q^E(b) := \sum_{m=1}^M \left( y_m - C_m^E b \right)^2 + \gamma \sum_{e=1}^E \frac{b_e^2}{\lambda_e}$

# Suboptimal finite dimensional solution

## New estimator

$$b_r = \left( \frac{1}{M} \text{diag} \left( \frac{\gamma}{\lambda_e} \right) + \frac{1}{M} \sum_{m=1}^M (C_m^E)^T C_m^E \right)^{-1} \left( \frac{1}{M} \sum_{m=1}^M (C_m^E)^T y_m \right)$$

- computable (involves  $E \times E$  matrices and  $E$ -dimensional vectors)
- minimizes  $Q^E(b) := \sum_{m=1}^M (y_m - C_m^E b)^2 + \gamma \sum_{e=1}^E \frac{b_e^2}{\lambda_e}$

## Drawbacks

- ①  $O(E^3)$  computational effort
- ②  $O(E^2)$  transmission effort
- ③ must know  $M$

## Derivation of the distributed estimator

$$b_r = \left( \frac{1}{M} \text{diag} \left( \frac{\gamma}{\lambda_e} \right) + \frac{1}{M} \sum_{m=1}^M (C_m^E)^T C_m^E \right)^{-1} \left( \frac{1}{M} \sum_{m=1}^M (C_m^E)^T y_m \right)$$

### Consider the approximations

- $M \rightarrow M_g$  (guess)
- $\frac{1}{M} \sum_{m=1}^M (C_m^E)^T C_m^E \rightarrow \mathbb{E}_\mu \left[ (C_m^E)^T C_m^E \right] = I$



## Derivation of the distributed estimator

obtain:

$$b_d = \left( \frac{1}{M_g} \text{diag} \left( \frac{\gamma}{\lambda_e} \right) + I \right)^{-1} \left( \frac{1}{M} \sum_{m=1}^M (C_m^E)^T y_m \right)$$

### Advantages

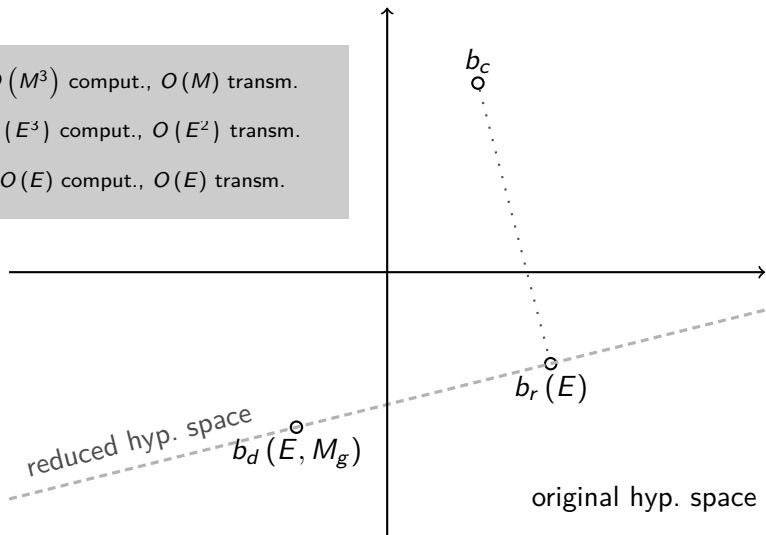
- ①  $O(E)$  computational effort
- ②  $O(E)$  transmission effort

# Summary of proposed estimation schemes

$b_c$ :  $O(M^3)$  comput.,  $O(M)$  transm.

$b_r$ :  $O(E^3)$  comput.,  $O(E^2)$  transm.

$b_d$ :  $O(E)$  comput.,  $O(E)$  transm.

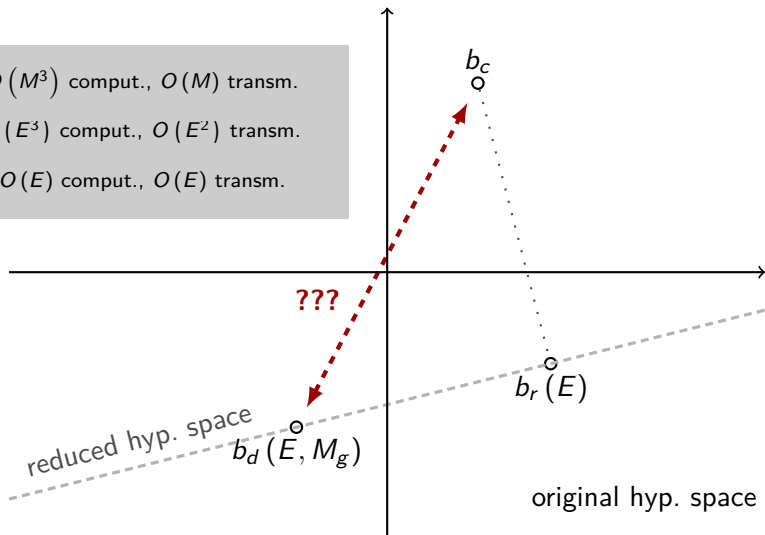


# Summary of proposed estimation schemes

$b_c$ :  $O(M^3)$  comput.,  $O(M)$  transm.

$b_r$ :  $O(E^3)$  comput.,  $O(E^2)$  transm.

$b_d$ :  $O(E)$  comput.,  $O(E)$  transm.

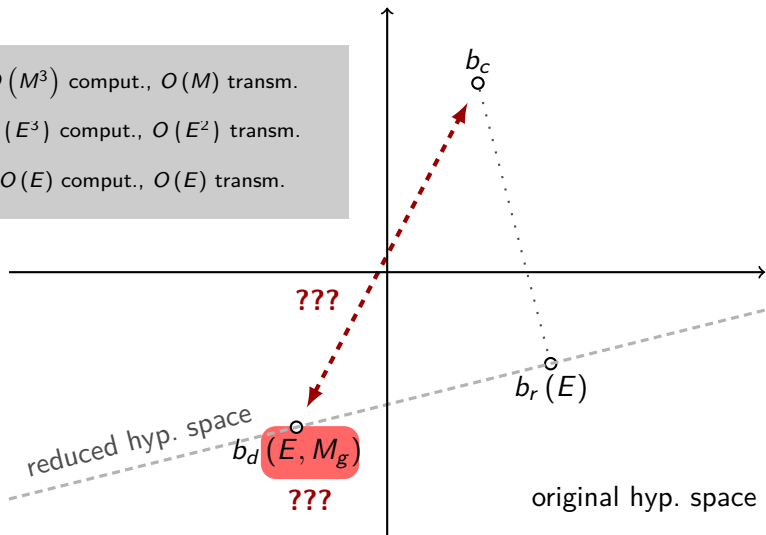


# Summary of proposed estimation schemes

$b_c$ :  $O(M^3)$  comput.,  $O(M)$  transm.

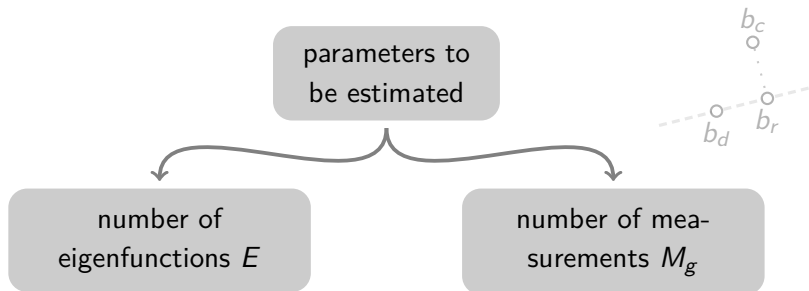
$b_r$ :  $O(E^3)$  comput.,  $O(E^2)$  transm.

$b_d$ :  $O(E)$  comput.,  $O(E)$  transm.



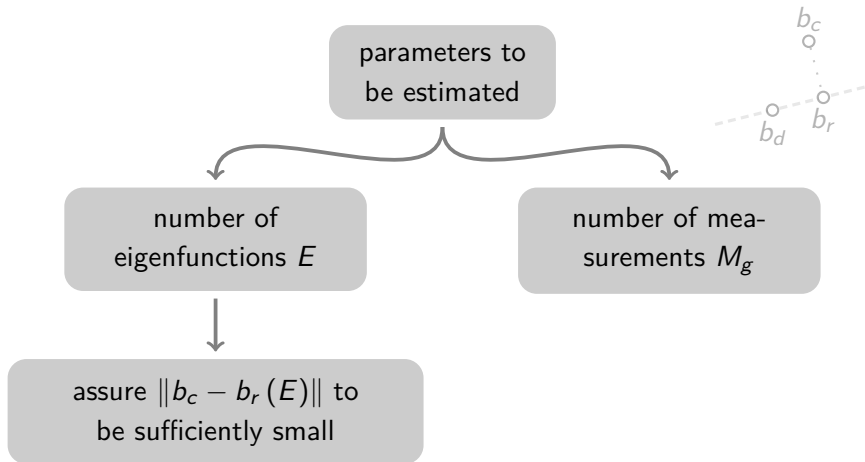
# Tuning of the parameters - key ideas

Assumption: have some information on the energy of  $f$



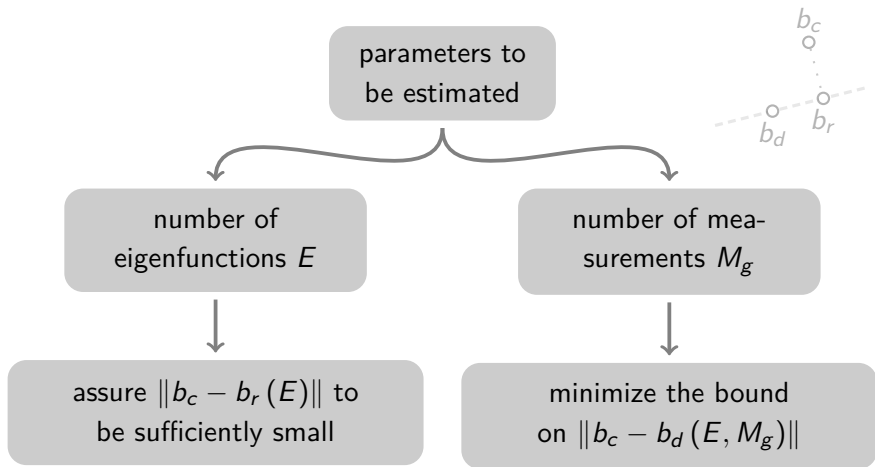
# Tuning of the parameters - key ideas

Assumption: have some information on the energy of  $f$

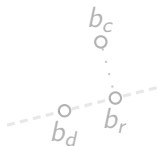


# Tuning of the parameters - key ideas

Assumption: have some information on the energy of  $f$



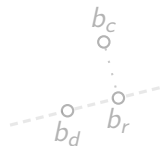
## Tuning of the parameters - in practice



$$\|b_c - b_d\|_2 \leq \frac{1}{M} \sum_{m=1}^M |r_m| + \|U_M b_d\|_2 + \|U_C b_d\|_2$$



# Tuning of the parameters - in practice

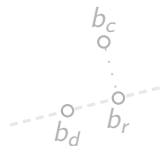


$$\|b_c - b_d\|_2 \leq \frac{1}{M} \sum_{m=1}^M |r_m| + \|U_M b_d\|_2 + \|U_C b_d\|_2$$

local residuals, func. of  $M_g$

# Tuning of the parameters - in practice

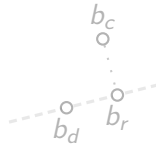
func. of  $M_g$  and  $\propto \frac{1}{M_{\min}} - \frac{1}{M_{\max}}$



$$\|b_c - b_d\|_2 \leq \frac{1}{M} \sum_{m=1}^M |r_m| + \|U_M b_d\|_2 + \|U_C b_d\|_2$$

local residuals, func. of  $M_g$

# Tuning of the parameters - in practice


$$\|b_c - b_d\|_2 \leq \frac{1}{M} \sum_{m=1}^M |r_m| + \|U_M b_d\|_2 + \|U_C b_d\|_2$$

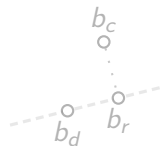
func. of  $M_g$  and  $\propto \frac{1}{M_{\min}} - \frac{1}{M_{\max}}$

local residuals, func. of  $M_g$

func. of  $M_g$  and  $\propto I - \frac{1}{M} \sum_{m=1}^M (C_m^E)^T C_m^E$

# Tuning of the parameters - in practice

func. of  $M_g$  and  $\propto \frac{1}{M_{\min}} - \frac{1}{M_{\max}}$



$$\|b_c - b_d\|_2 \leq \frac{1}{M} \sum_{m=1}^M |r_m| + \|U_M b_d\|_2 + \|U_C b_d\|_2$$

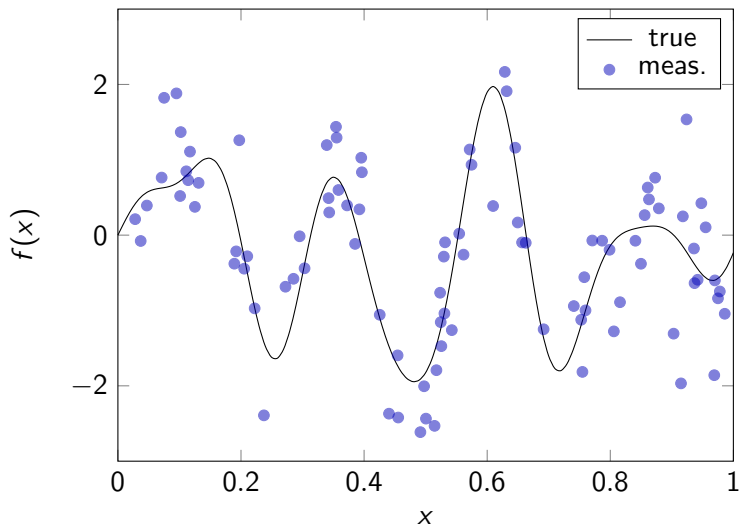
local residuals, func. of  $M_g$

$$\text{func. of } M_g \text{ and } \propto I - \frac{1}{M} \sum_{m=1}^M (C_m^E)^T C_m^E$$

computable through distributed MC

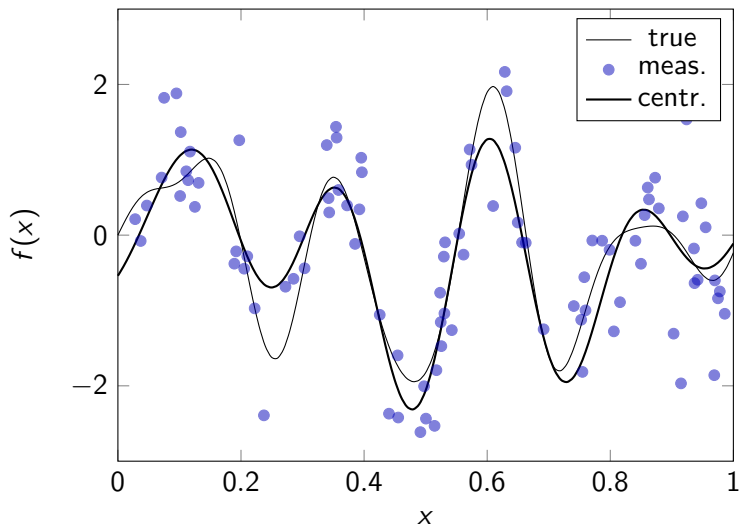
## Regression strategy effectiveness example

$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$ ,  $\text{SNR} \approx 2.5$



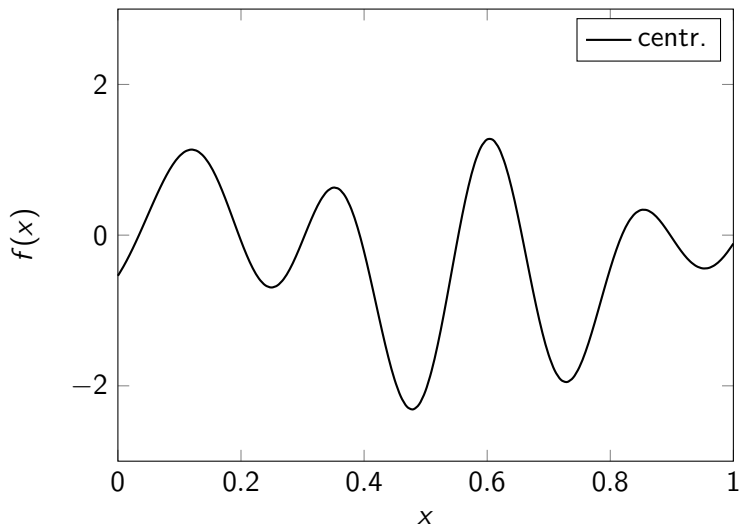
## Regression strategy effectiveness example

$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$ ,  $\text{SNR} \approx 2.5$



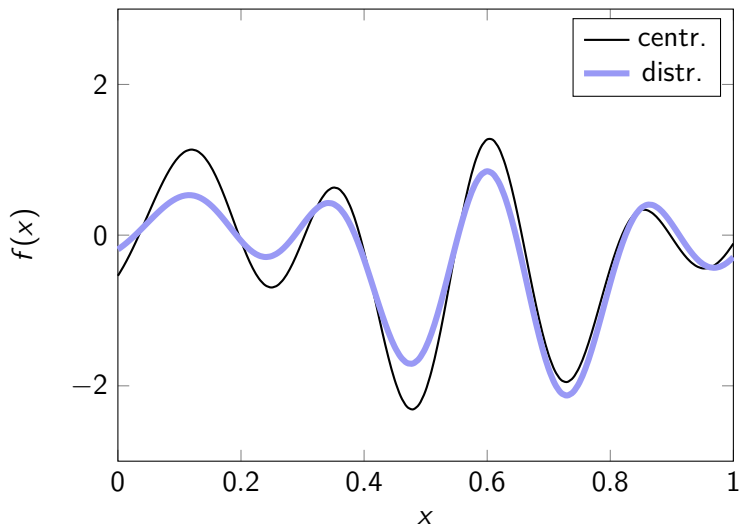
## Regression strategy effectiveness example

$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$ ,  $\text{SNR} \approx 2.5$



## Regression strategy effectiveness example

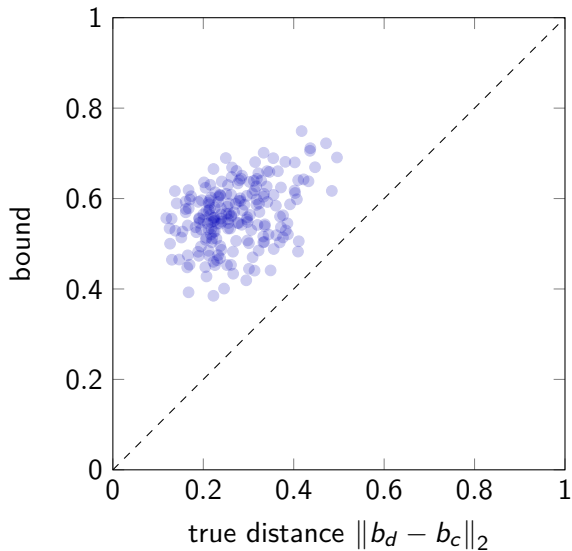
$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$ ,  $\text{SNR} \approx 2.5$





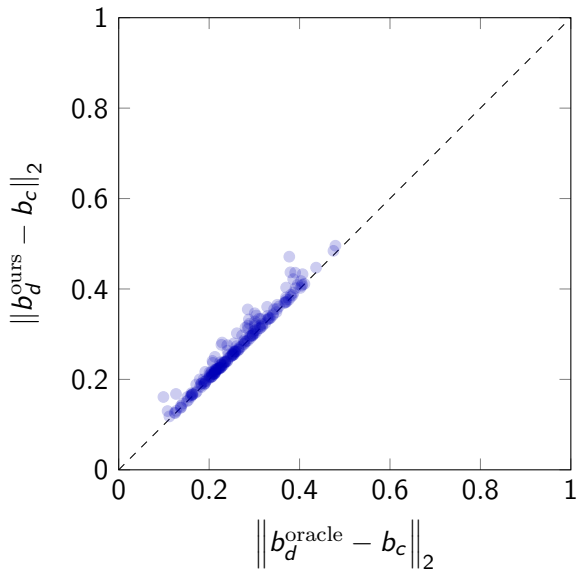
## Accuracy of the computed bound

$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$



# Comparison with oracle

$M = 100$ ,  $E = 20$ ,  $M_{\min} = 90$ ,  $M_{\max} = 110$



## Conclusions

Strategy:

- is effective and easy to be implemented
- has self-evaluation capabilities
- has self-tuning capabilities

## Future works

- exploit statistical knowledge about  $M$
- incorporate effects of finite number of steps in consensus algorithms
- extend to dynamic scenarios (long term objective)

# Auto-tuning procedures for distributed nonparametric regression algorithms

**Damiano Varagnolo**, Gianluigi Pillonetto, Luca Schenato

ECC 2015 - Linz

July 15, 2015

[damiano.varagnolo@ltu.se](mailto:damiano.varagnolo@ltu.se)