

# Multi-agent map-building: Kalman Filtering meets Machine Learning



**Luca Schenato**

University of Padova  
Stuttgart, April 2017

**M.g.i.C.**  
Multi Agent Intelligent Control

DEPARTMENT OF  
INFORMATION  
ENGINEERING  
UNIVERSITY OF PADOVA



# Joint work with

## Colleagues at Univ. of Padova



Ruggero Carli



Gianluigi Pillonetto

## Current Ph.D/post-docs:



Marco Todescato



Nicoletta Bof

## Former Ph.D/post-docs:



Damiano Varagnolo  
Lulea Univ., Sweden



Guido Cavraro  
Virginia Tech, USA



Andrea Carron  
ETH, Switzerland



# Outline

---

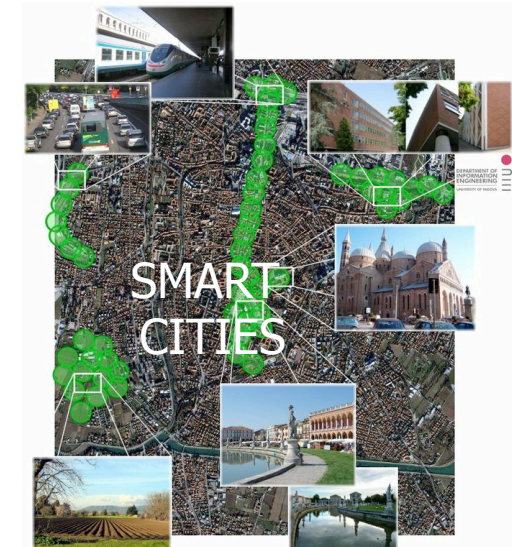
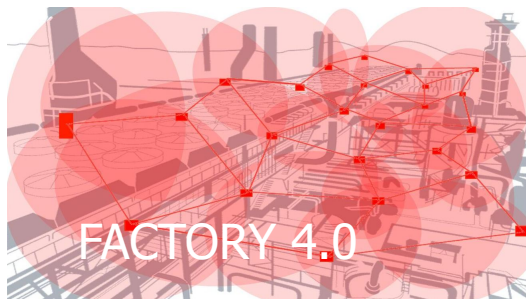
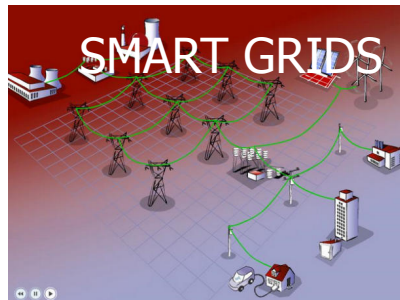
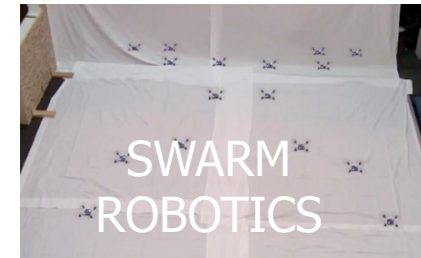
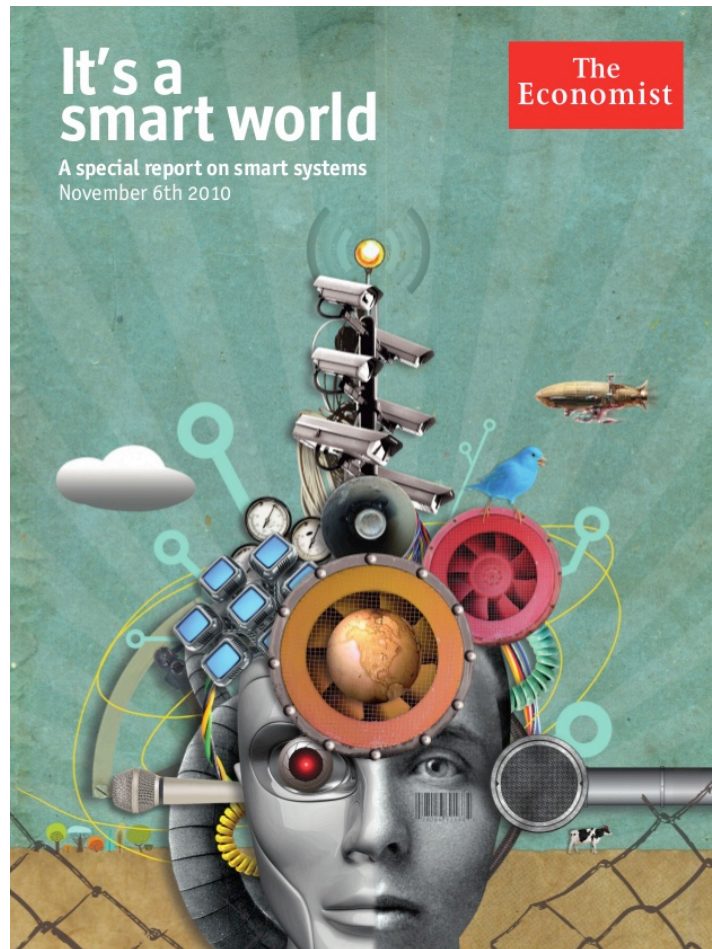
- Motivations, target applications & challenges
- Parametric regression
- Non-parametric regression
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems

# Outline

---

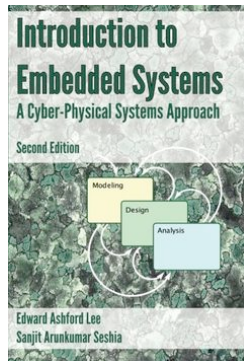
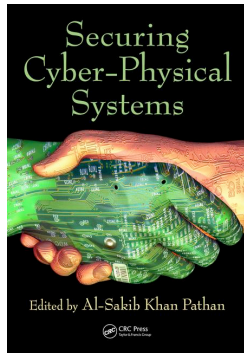
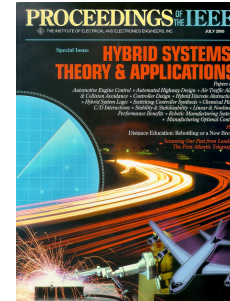
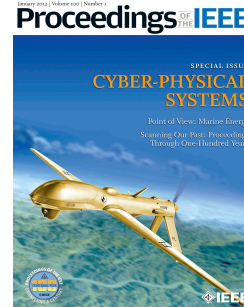
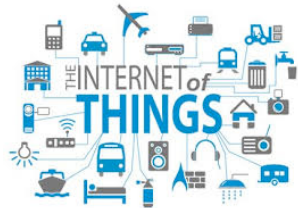
- Motivations, target applications & challenges
- Parametric regression
- Non-parametric regression
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems

# The XXI century: a Smart World

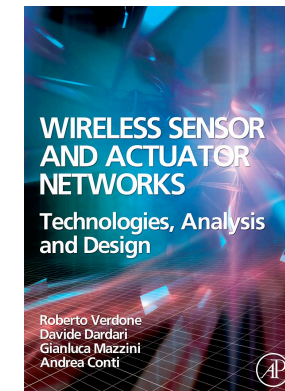
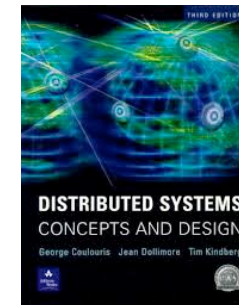
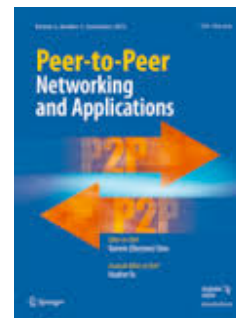
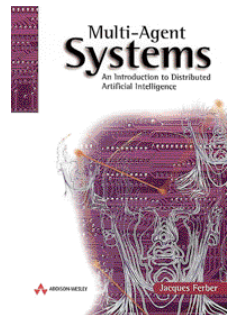
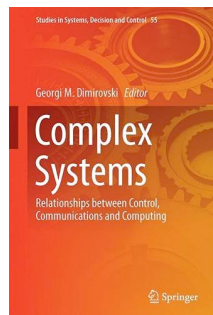




# The ICT scientific army



Before 2000's	Today	Tomorrow
Centralized Reliable comm.	Multi-agent Reliable comm.	Multi-agent Unrelia. Comm.
✓	✓	?



# Target applications: MAGIC Lab. at University of Padova

**Wireless Sensor  
Actuator Networks**

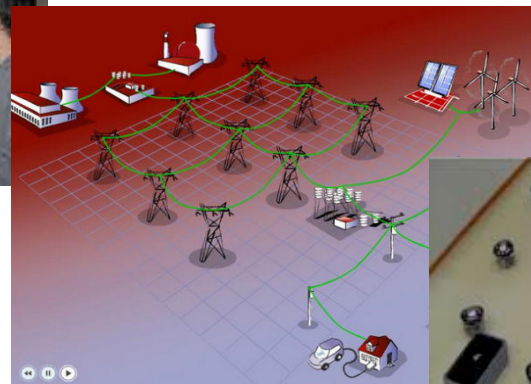


**Smart Camera  
Networks**



**M.  g. i. c.**  
Multi Agent Intelligent Control

**Smart Energy  
Grids**



**Robotic  
Networks**



**2002**

**2008**

**2012**



# Challenges

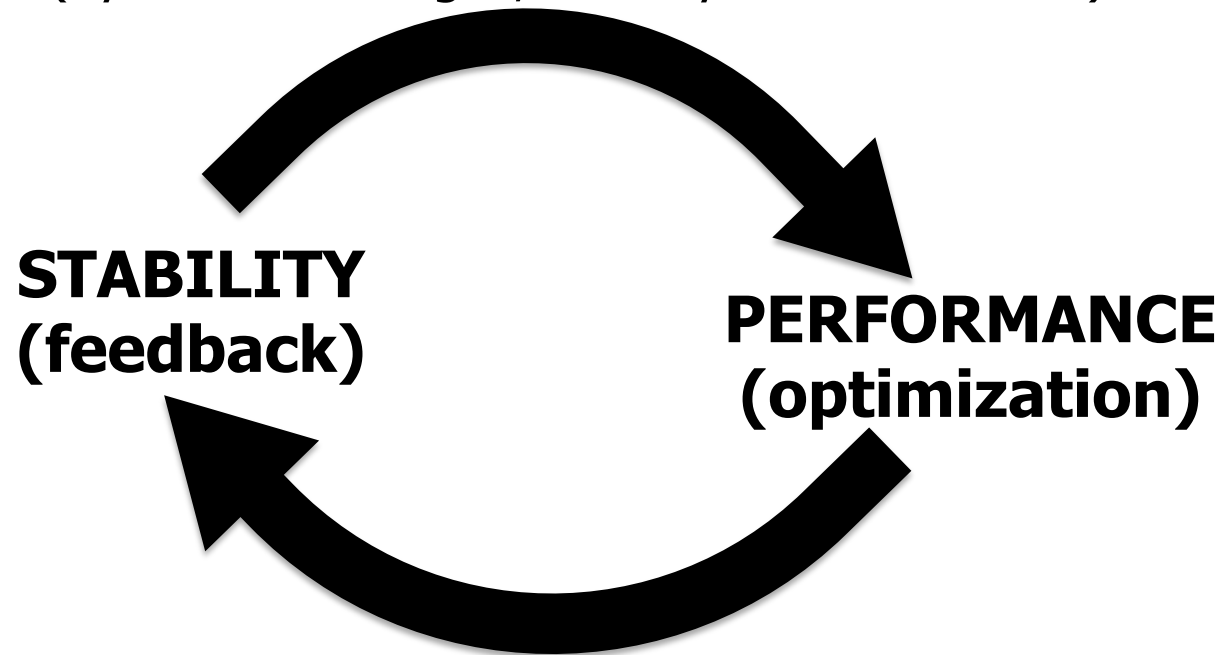
- Unreliable (wireless) communication:
  - Random delay, packet loss, limited communication range
- Scalability:
  - Complexity (CPU, memory, communication) per agent must be constant
- Robustness/resilience and adaptiveness/learning:
  - Mild performance degradation when local failures
  - Continuous environmental learning
- Architecture:
  - Centralized vs hierarchical vs distributed vs decentralized
  - Cooperative vs competitive

# Challenges

- Unreliable (wireless) communication:
  - Random delay, packet loss, limited communication range
- Scalability:
  - Complexity (CPU, memory, communication) per agent must be constant
- Robustness/resilience and **adaptiveness/learning**:
  - Mild performance degradation when local failures
  - Continuous environmental learning
- Architecture:
  - Centralized vs hierarchical vs distributed vs decentralized
  - Cooperative vs competitive

# Dynamic learning and optimization

**Environment learning**  
(dynamical changes, "steady state" scenario)

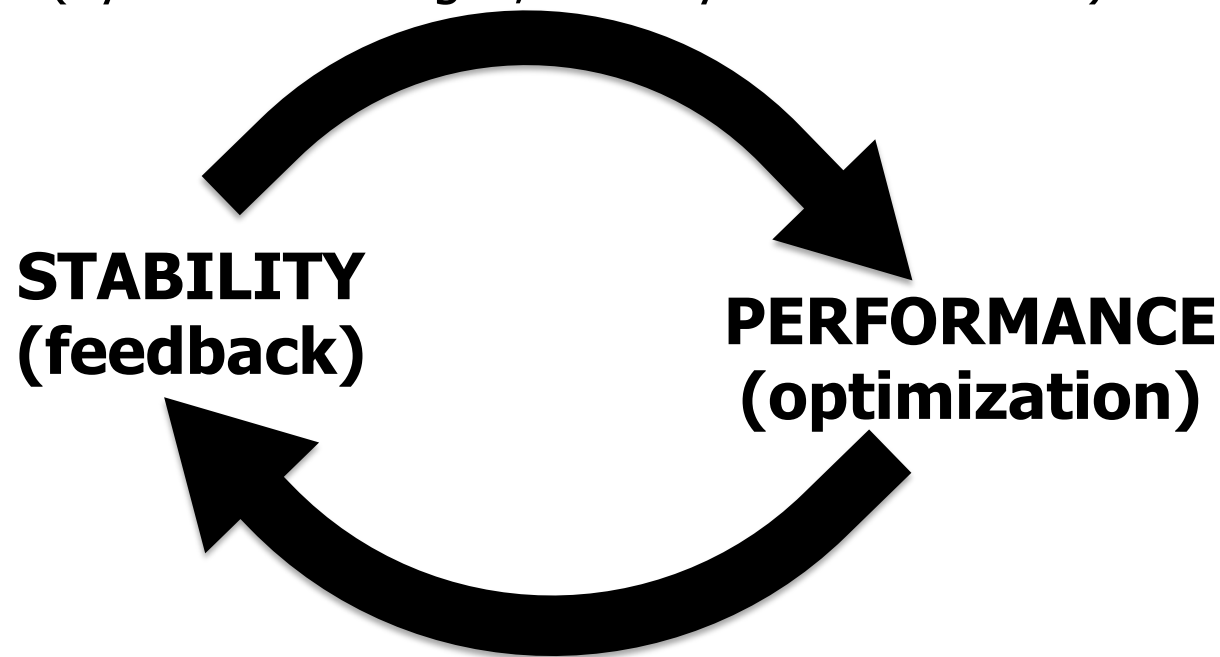


**Disruptive events detection**  
(local failures, communication blackouts, new agents integration, ....)

# Dynamic learning and optimization

## Environment learning

(dynamical changes, "steady state" scenario)



## Disruptive events detection

(local failures, communication blackouts, new agents integration, ....)

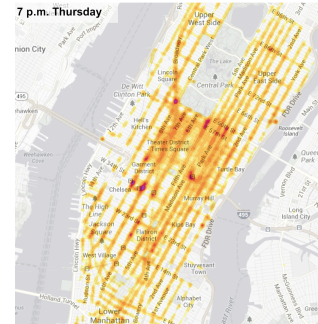
# Learning problems:

## Density estimation:

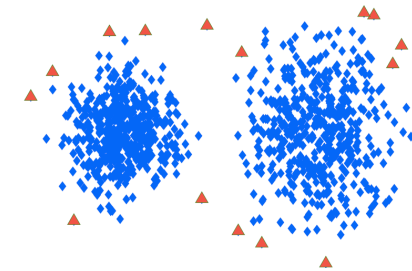
$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f(x) \geq 0, \quad \int f(x) = 1$$

$$\mathcal{D} = \{x_1, x_2, \dots\}: \text{events}$$



Taxi pick-up calls



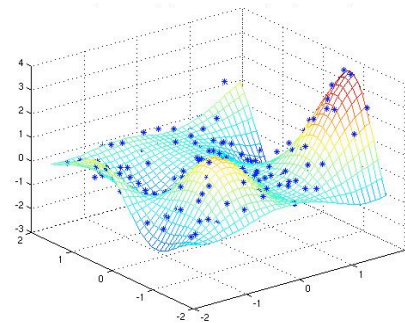
Anomaly detection

## Regression:

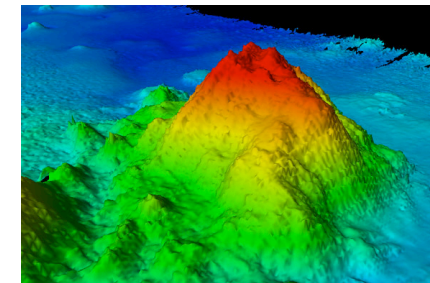
$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$y_i = f(x_i) + v_i$$

$v_i$  noise



Pollution level profile



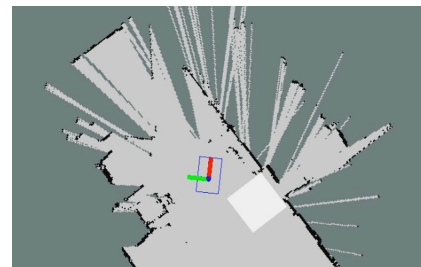
Seabed depth profile

## Classification

$$f(x) : \mathbb{R}^n \rightarrow \{0, 1\}$$

$$y_i = f(x_i) + v_i$$

$v_i$  noise



Obstacles map



Oil-spill boundary



# Multi-agent regression

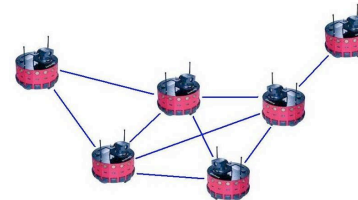
## ■ Parametric vs non-parametric

$$f(x) = \sum_{i=1}^m \theta_i g_i(x)$$

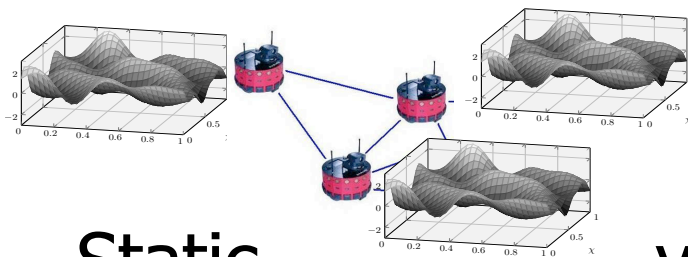
$\theta \in \mathbb{R}^m$ , unknown  
 $g_i(x)$  known

$f(x) \in RKHS$ , infinite dimensional  
 $f(x)$  defined via Kernel  $k(x, x')$   
 $k(x, x')$  known

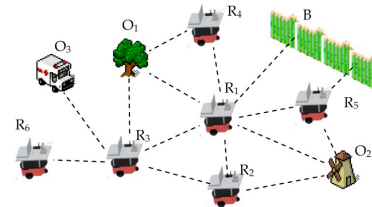
## ■ Cloud-based vs peer-to-peer



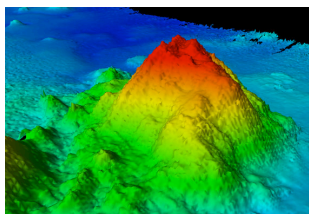
## ■ Global vs Local estimation



## Local estimation

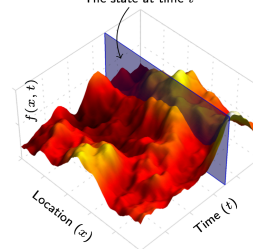


## ■ Static vs dynamic maps:



$$f(x)$$

## dynamic maps:



$$f(x, t)$$

# Multi-agent regression

## ■ Parametric

$$f(x) = \sum_{i=1}^m \theta_i g_i(x)$$

$\theta \in \mathbb{R}^m$ , unknown  
 $g_i(x)$  known

vs

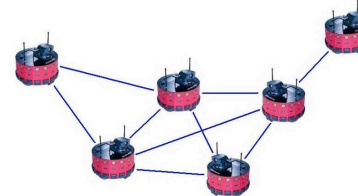
## non-parametric

$f(x) \in RKHS$ , infinite dimensional  
 $f(x)$  defined via Kernel  $k(x, x')$   
 $k(x, x')$  known

## ■ Cloud-based

vs

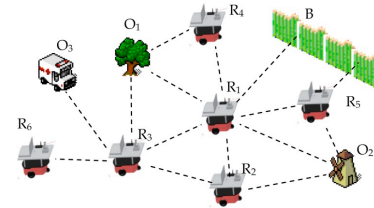
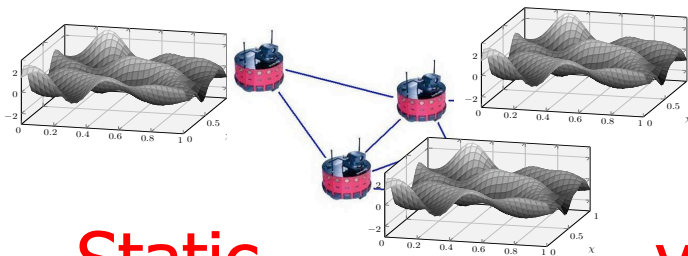
## peer-to-peer



## ■ Global

vs

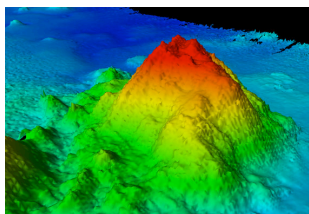
## Local estimation



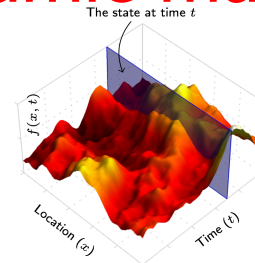
## ■ Static

vs

## dynamic maps:



$f(x)$



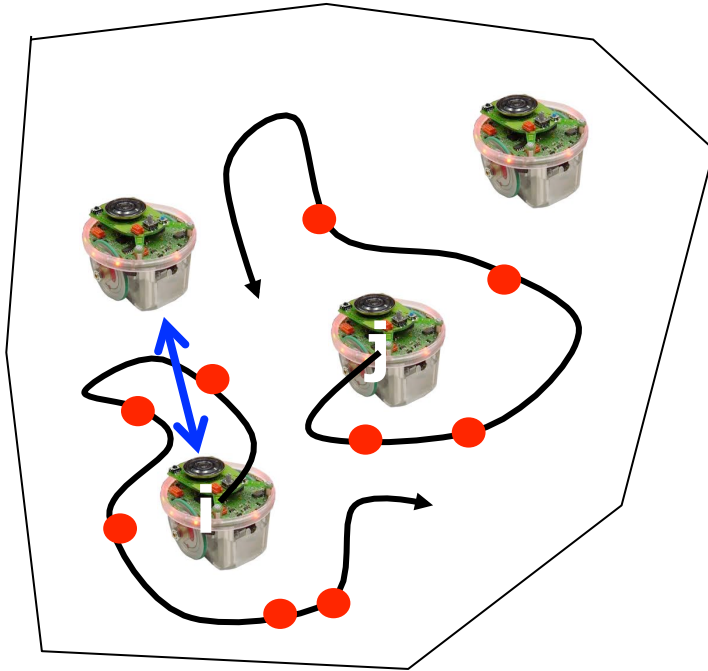
$f(x, t)$

# Outline

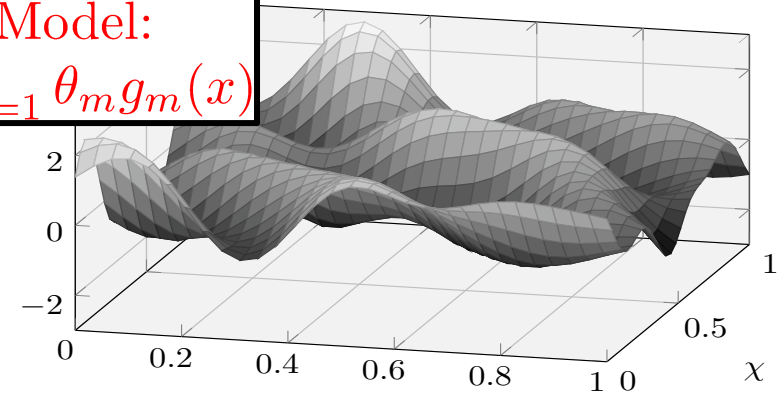
---

- Motivations, target applications & challenges
- **Parametric regression**
- Non-parametric regression
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems

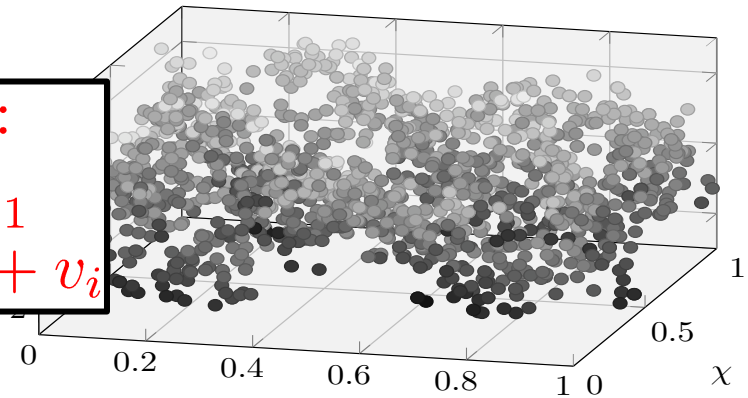
# Example: Map-building in robotic networks



Parametric Model:  
 $f(x) = \sum_{m=1}^M \theta_m g_m(x)$



Noisy data:  
 $\{(x_i, y_i)\}_{i=1}^N$   
 $y_i = f(x_i) + v_i$



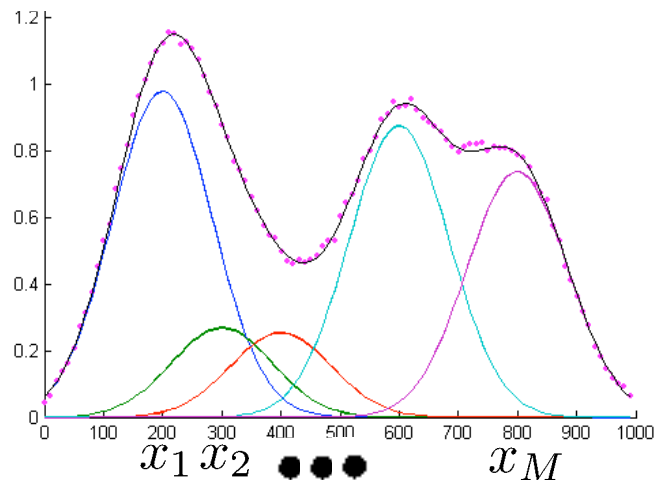
Goal:  
 $\min_{\theta} \sum_i v_i^2$

# Parametric model: linear vs non-linear

## Linear combination of radial basis functions

$$f(x) = \sum_{i=1}^M \theta_i g_i(x)$$

$$g_i(x) = e^{-\frac{|x-x_i|^2}{2\sigma^2}}$$

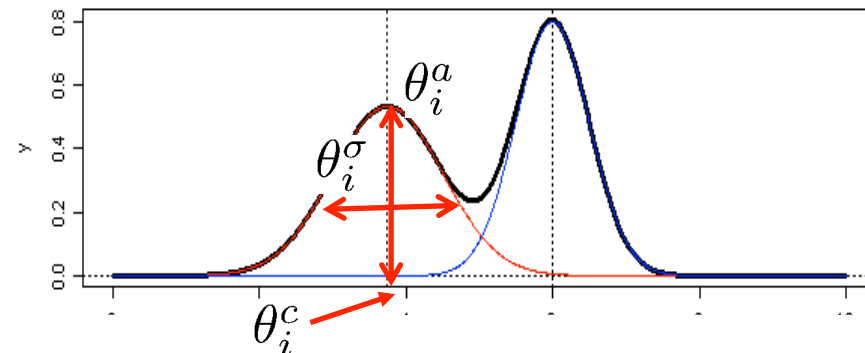


- Large number of basis
- Convex problem

## Mixture of Gaussians

$$f(x) = \sum_{i=1}^M g_i(x, \theta_i)$$

$$g_i(x, \theta_i) = \theta_i^a e^{-\frac{|x-\theta_i^c|^2}{2\theta_i^\sigma{}^2}}$$



- Needs fewer functions
- Non-linear problem



# Map-building as least-squares regression

- Model class:

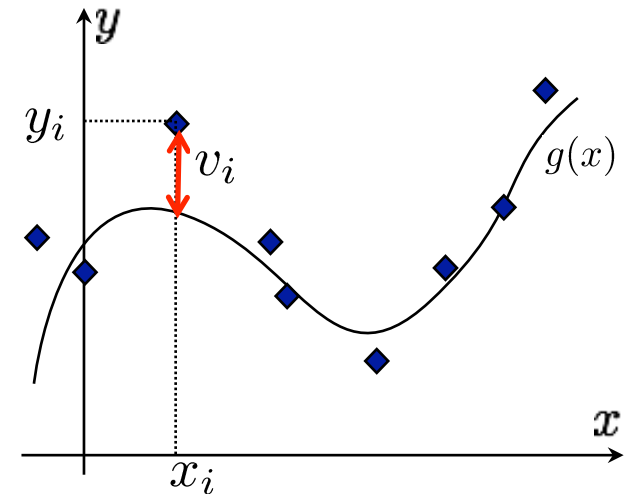
$$f(x) = \sum_{m=1}^M \theta_m g_m(x)$$

- Noisy measurements:

$$y_i = \sum_{m=1}^M \theta_m g_m(x_i) + v_i, \quad i = 1, \dots, N$$

$$\begin{bmatrix} y(x_1) \\ y(x_2) \\ \vdots \end{bmatrix} = \begin{bmatrix} g_1(x_1) & \dots & g_M(x_1) \\ \vdots & \vdots & \vdots \\ g_1(x_N) & \dots & g_M(x_N) \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_M \end{bmatrix} + \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix}$$

$$y = G\theta + v$$



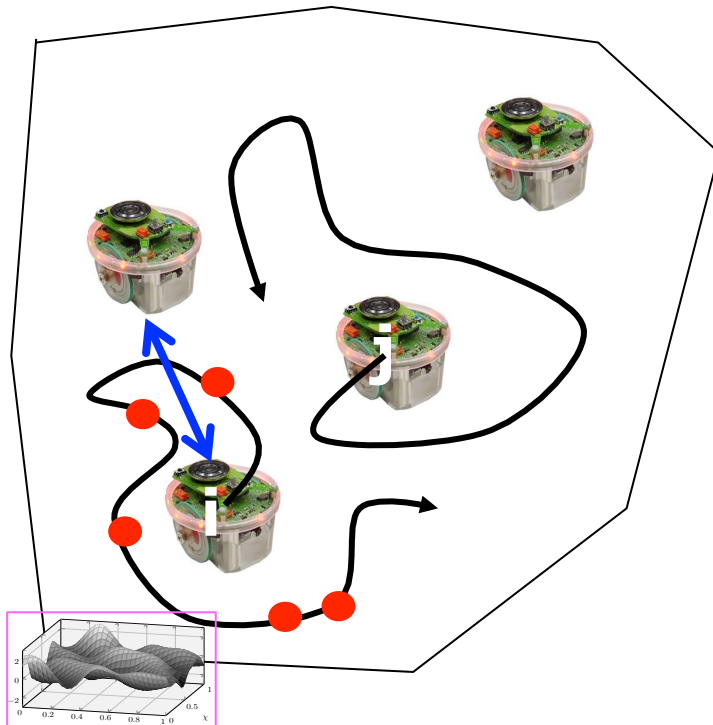
- Goal: minimize sum of squares of residues

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{i=1}^N v_i^2$$

$$\begin{aligned} \hat{\theta} &= (\sum_{i=1}^N G_i G_i^T)^{-1} (\sum_{i=1}^N G_i y_i) \\ &= (\frac{1}{N} \sum_{i=1}^N G_i G_i^T)^{-1} (\frac{1}{N} \sum_{i=1}^N G_i y_i) \end{aligned}$$

- Xiao-Boyd-Lall, 2005
- Bolognani-Del Favero-Schenato-Varagnolo, 2010

# Consensus-based Map-building: gossip communication



## ALGORITHM:

1) Initialize statistics:

$$Z_0^i = 0 \in R^{M \times M}$$

$$z_0^i = 0 \in R^M$$

2) Collect data and build local statistics:

$$Z_{t+1}^i = Z_t^i + G_t^i G_t^{iT}$$

$$z_{t+1}^i = z_t^i + G_t^i y_t^i$$

3) Choose neighbor  $j$  and do gossip consensus:

$$Z_{t+1}^j = Z_{t+1}^i = \frac{1}{2} Z_t^i + \frac{1}{2} Z_t^j$$

$$z_{t+1}^j = z_{t+1}^i = \frac{1}{2} z_t^i + \frac{1}{2} z_t^j$$

4) Estimate map:

$$\hat{\theta}_t^i = (Z_t^i)^{-1} z_t^i$$

5) Repeat steps 2,3,4 (non necessarily in order)

## ■ PROS:

- Can be distributed
- Gradient-based implementation: ADMM, gradient-consensus,
- Extension to robust costs, e.g.  $\| \cdot \|_1$

## ■ CONS:

- How to select basis functions
- No estimate unless at least M data
- Gradient-based implementations require step-size design



# Simulations:

## broadcast based map building

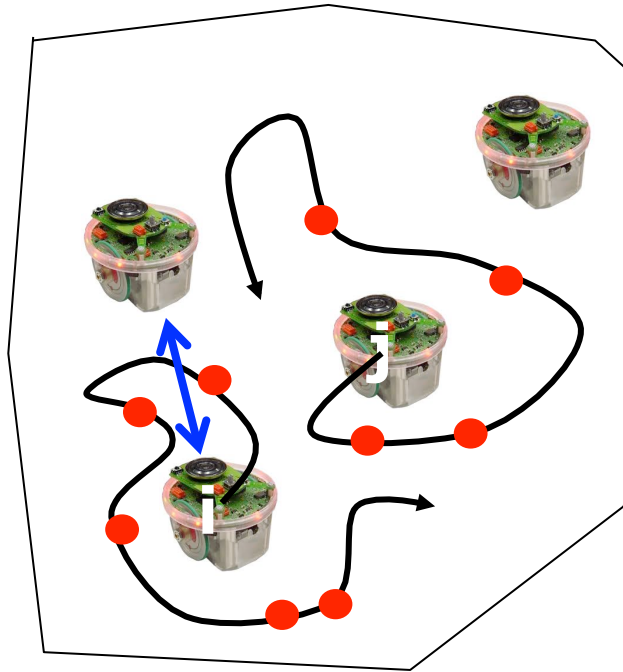


# Outline

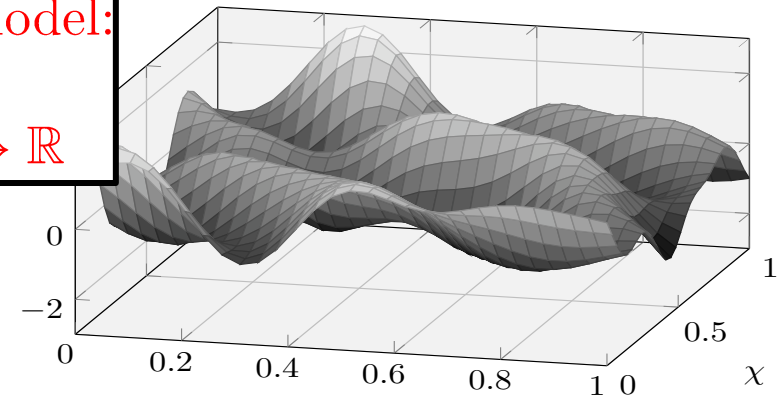
---

- Motivations, target applications & challenges
- Parametric regression
- **Non-parametric regression**
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems

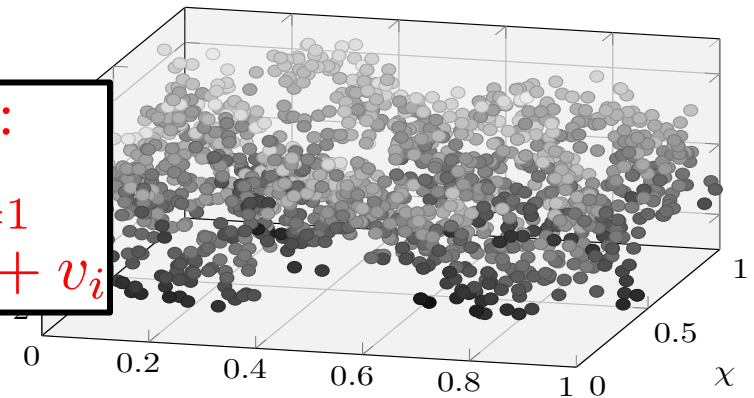
# Gaussian regression (non parametric)



Non-parametric Model:  
 $f(x) \in \text{RKHS}$   
 $k(x, x') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$



Noisy data:  
 $\{(x_i, y_i)\}_{i=1}^N$   
 $y_i = f(x_i) + v_i$



Goal:  
 $\min_{\theta} \sum_i v_i^2$



# Reproducing Kernel Hilbert Spaces (RKHS) (con't)

$k(x, x') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ : Mercer Kernel

1)  $k(\cdot, \cdot)$  continuous,  $\mathcal{X}$ : compact

2) symmetric:  $k(x, x') = k(x', x)$

3) positive semidefinite:  $K \in \mathbb{R}^{N \times N} \geq 0$ ,  $[K]_{i,j} = k(x_i, x_j), \forall x_i, \forall N$ ,

## Bayesian Interpretation:

$\mathbb{E}[f(x)] = 0$ ,  $\mathbb{E}[f(x)f(x')] = k(x, x')$ : zero-mean gaussian process

$\{(x_i, y_i)\}_{i=1}^N$ : Noisy data:

$$y_i = f(x_i) + v_i, \quad v_i \sim \mathcal{N}(0, \sigma^2)$$

$$\hat{f}(x) = \mathbb{E}[f(x) | \{x_i, y_i\}_{i=1}^M] = \sum_{i=1}^M c_i k(x_i, x)$$

$$\begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = (K + \sigma^2 I)^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}, \quad K := \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_M) \\ \vdots & & \vdots \\ k(x_M, x_1) & \cdots & k(x_M, x_M) \end{bmatrix}$$

# Parametric vs non-parametric

$\{(x_i, y_i)\}_{i=1}^N$ : Noisy data:

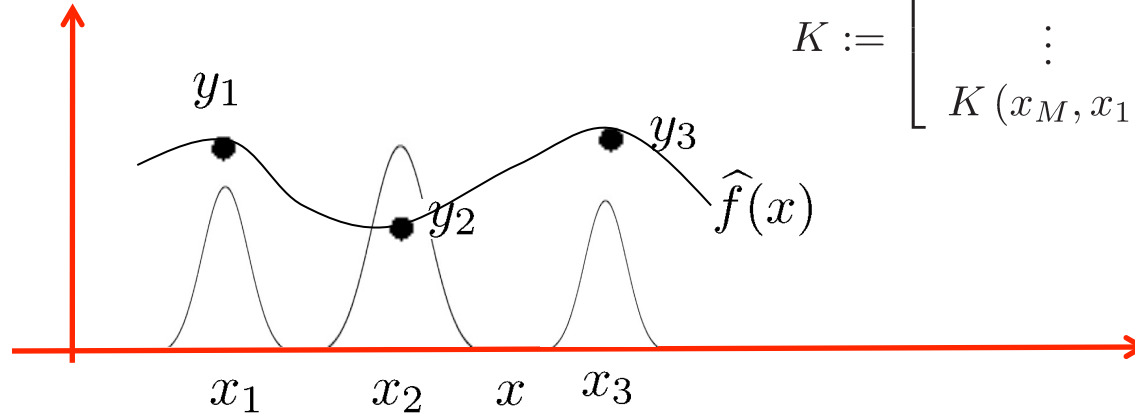
$$y_i = f(x_i) + v_i, \quad v_i \sim \mathcal{N}(0, \sigma^2)$$

$$f(x) = \sum_{i=1}^M \theta_i g_i(x)$$

$$g_i(x) = e^{-\frac{|x-x_i|^2}{2\sigma^2}}$$

$$k(x, x') = e^{-\frac{|x-x'|^2}{2\sigma^2}}$$

$$K := \begin{bmatrix} K(x_1, x_1) & \cdots & K(x_1, x_M) \\ \vdots & & \vdots \\ K(x_M, x_1) & \cdots & K(x_M, x_M) \end{bmatrix}$$



$$\hat{f}(x) = \sum_{i=1}^M \hat{\theta}_i g_i(x)$$

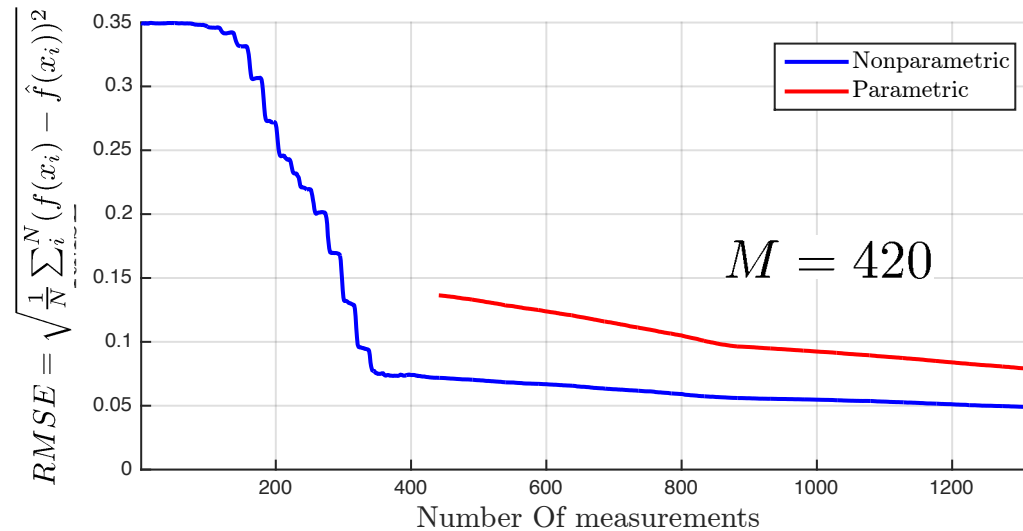
$$\hat{f}(x) = \sum_{i=1}^M \hat{c}_i k(x_i, x) = \sum_{i=1}^M \hat{c}_i g_i(x)$$

$$\hat{\theta} = K^{-1}y$$

$$\hat{c} = (K + \sigma^2 I)^{-1}y$$

regularization term

# Parametric vs non-parametric



	PARAMETRIC	$N$ NON-PARAMETRIC
PROS	<ul style="list-style-type: none"> <li>• Distributed (consensus)</li> <li>• Bounded complexity <math>O(M^3)</math></li> </ul>	<ul style="list-style-type: none"> <li>• Better performance</li> <li>• Adaptable resolution</li> </ul>
CONS	<ul style="list-style-type: none"> <li>• What <math>g_i(x)</math> ?</li> <li>• Need <math>N &gt; M</math> points</li> <li>• Over-fitting &amp; ill-conditioned</li> </ul>	<ul style="list-style-type: none"> <li>• Regularization factor design</li> <li>• Data-limited complexity <math>O(N^3)</math></li> </ul>

# Outline

---

- Motivations, target applications & challenges
- Parametric regression
- Non-parametric regression
- **Semi non-parametric regression**
- Non-parametric regression for dynamical systems
- Conclusion and open problems

# Representer theorem

$k(x, x') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ : Mercer Kernel

1)  $k(\cdot, \cdot)$  continuous,  $\mathcal{X}$ : compact

2) symmetric:  $k(x, x') = k(x', x)$

3) positive semidefinite:  $K \in \mathbb{R}^{N \times N} \geq 0$ ,  $[K]_{i,j} = k(x_i, x_j), \forall x_i, \forall N$ ,

$\mu : \mathcal{X} \rightarrow \mathbb{R}^+$  : measure function (sampling density)

$h(x) := T_{k,\mu}[g](x) := \int_{\mathcal{X}} g(x')k(x, x')d\mu(x')$ : Hilbert-Schmidt integral operator  
 $h(x), g(x) \in \mathcal{L}^2(\mu)$

Since T is a linear operator  $\rightarrow$  eigenvalues and eigenfunctions

$$T_{k,\mu}[\phi(x)] = \lambda \phi(x), \lambda \geq 0$$

**Representer Theorem:** Let  $k(\cdot, \cdot)$  be a Mercer kernel on  $\mathcal{X} \times \mathcal{X}$ ,  $\lambda_\ell > 0 \quad \forall \ell$  and  $\mu$  a non-degenerate measure. Then,  $\{\phi_\ell\}_{\ell=1}^{+\infty}$  is an orthonormal basis in  $\mathcal{L}^2(\mu)$  while the associated RKHS is

$$\mathcal{H}_K := \left\{ f(x) \in \mathcal{L}^2(\mu) \text{ s.t. } f(x) = \sum_{\ell=1}^{\infty} \alpha_\ell \phi_\ell(x) \text{ and } \sum_{e=1}^{\infty} \frac{\alpha_e^2}{\lambda_e} < +\infty \right\}$$

# Map-building as least-squares regression

- Model class:

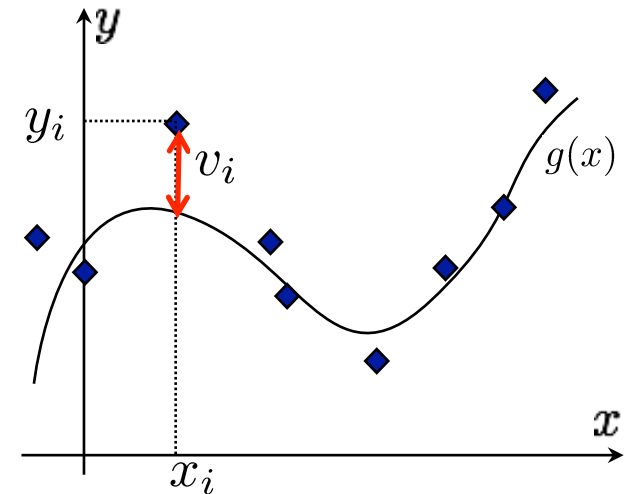
$$f(x) = \sum_{m=1}^M \theta_m g_m(x)$$

- Noisy measurements:

$$y_i = \sum_{m=1}^M \theta_m g_m(x_i) + v_i, \quad i = 1, \dots, N$$

$$\begin{bmatrix} y(x_1) \\ y(x_2) \\ \vdots \end{bmatrix} = \begin{bmatrix} g_1(x_1) & \dots & g_M(x_1) \\ \vdots & \vdots & \vdots \\ g_1(x_N) & \dots & g_M(x_N) \end{bmatrix} \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_M \end{bmatrix} + \begin{bmatrix} v_1 \\ \vdots \\ v_N \end{bmatrix}$$

$$y = G\theta + v$$



- Goal: minimize sum of squares of residues

$$\hat{\theta} = \operatorname{argmin}_{\theta} \sum_{i=1}^N v_i^2$$

$$\begin{aligned} \hat{\theta} &= (\sum_{i=1}^N G_i G_i^T)^{-1} (\sum_{i=1}^N G_i y_i) \\ &= (\frac{1}{N} \sum_{i=1}^N G_i G_i^T)^{-1} (\frac{1}{N} \sum_{i=1}^N G_i y_i) \end{aligned}$$

- Xiao-Boyd-Lall, 2005
- Bolognani-Del Favero-Schenato-Varagnolo, 2010

# Semi-parametric estimation

1<sup>st</sup> IDEA: Use first eigenfunctions as basis function for parametric estimation

$$f(x) = \sum_{\ell=1}^{+\infty} \alpha_{\ell} \phi_{\ell}(x)$$

$$y_i = \sum_{\ell=1}^{+\infty} \alpha_{\ell} \phi_{\ell}(x_i) = \underbrace{[\phi_1(x_i) \ \phi_2(x_i) \ \dots]}_{G_i^T} \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \end{bmatrix} + v_i$$

$$\hat{\alpha} = \left( \text{diag}\left(\frac{\sigma^2}{\lambda_{\ell}}\right) + \sum_{i=1}^N G_i G_i^T \right)^{-1} \left( \sum_{i=1}^N G_i y_i \right)$$

$$\hat{\alpha}^E = \left( \text{diag}\left(\frac{\sigma^2}{\lambda_{\ell}}\right) + \sum_{i=1}^N G_i^E (G_i^E)^T \right)^{-1} \left( \sum_{i=1}^N G_i^E y_i \right)$$

$$G_i^E = [\phi_1(x_i) \cdots \phi_E(x_i)]$$

(intuition:  $\alpha_i \approx 0$ , for  $i > E$ , therefore  $\hat{f}(x) \approx \hat{f}^E(x)$ )

## Semi-parametric estimation (cont'd)

2<sup>st</sup> IDEA: Use orthonormality of eigenfunctions  $\phi_n$  and i.i.d. sampling of  $x_i$

$$\hat{\alpha}^E = \left( \text{diag}\left(\frac{\sigma^2}{N\lambda_\ell}\right) + \frac{1}{N} \sum_{i=1}^N G_i^E (G_i^E)^T \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N G_i^E y_i \right)$$

$$\left[ \frac{1}{N} \sum_{i=1}^N G_i^E (G_i^E)^T \right]_{mn} = \frac{1}{N} \sum_{i=1}^N \phi_m(x_i) \phi_n(x_i)$$

$$\left[ \frac{1}{N} \sum_{i=1}^N \phi_m(x_i) \phi_n(x_i) \right] \xrightarrow{N \rightarrow +\infty, x_i \sim \mu(x)} \int \phi_m(x) \phi_n(x) d\mu(x) = \delta_{mn}$$

$$\hat{\alpha}^I(x) = \left( \text{diag}\left(\frac{\sigma^2}{N\lambda_\ell}\right) + I \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N G_i^E y_i \right)$$



# Complexity of semi-parametric approaches

$$\hat{f}(x) = \sum_{i=1}^N c_i k(x_i, x), \quad c = (K + I)^{-1} y, \quad [K]_{mn} = k(x_m, x_n)$$

$$\hat{f}^E(x) = \sum_{i=1}^E \alpha_i^E \phi_i(x) \quad \hat{\alpha}^E = \left( \text{diag}\left(\frac{\sigma^2}{N\lambda_\ell}\right) + \frac{1}{N} \sum_{i=1}^N G_i^E (G_i^E)^T \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N G_i^E y_i \right)$$

$$\hat{f}^I(x) = \sum_{i=1}^I \alpha_i^E \phi_i(x) \quad \hat{\alpha}^I = \left( \text{diag}\left(\frac{\sigma^2}{N\lambda_\ell}\right) + I \right)^{-1} \left( \frac{1}{N} \sum_{i=1}^N G_i^E y_i \right)$$

<i>estimator</i>	<i>comput. cost</i>	<i>commun. cost</i>	<i>memory cost</i>
$\hat{f}(x)$	$O(N^3)$	$O(N)$	$O(N)$
$\hat{f}^E(x)$	$O(E^3)$	$O(E^2)$	$O(E^2)$
$\hat{f}^I(x)$	$O(E)$	$O(E)$	$O(E)$

# Performance of semi-parametric approaches

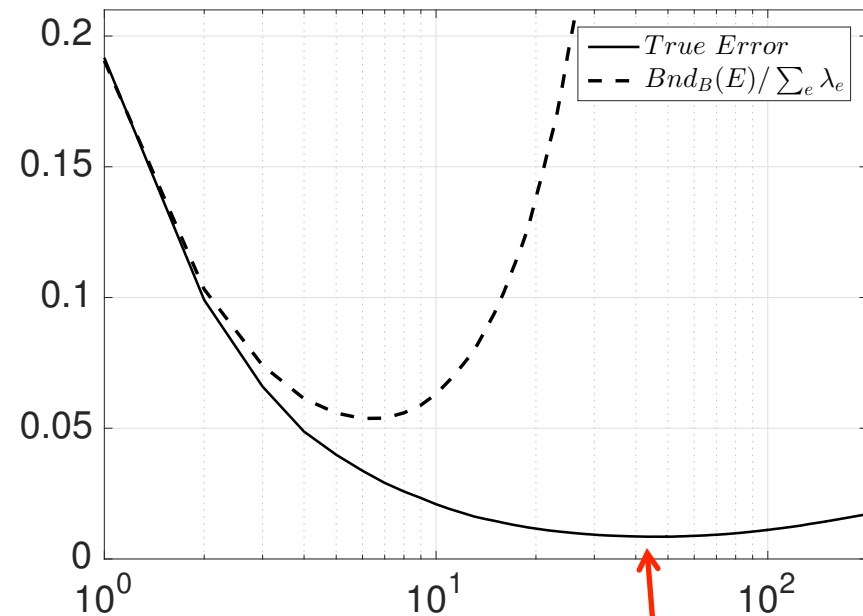
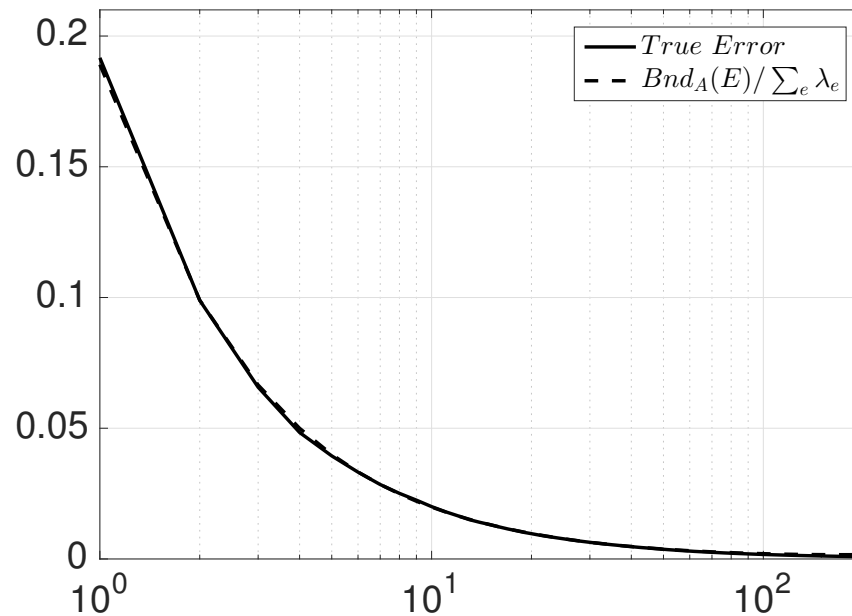
$N = 10000$

$\hat{f}^E(x)$

$\hat{f}^I(x)$

$K = \text{Splines}$

$K = \text{Splines}$



$E_{opt} = 50$

# Outline

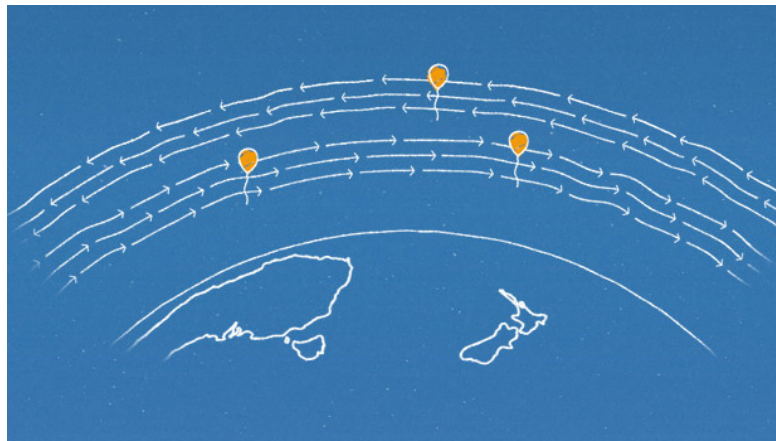
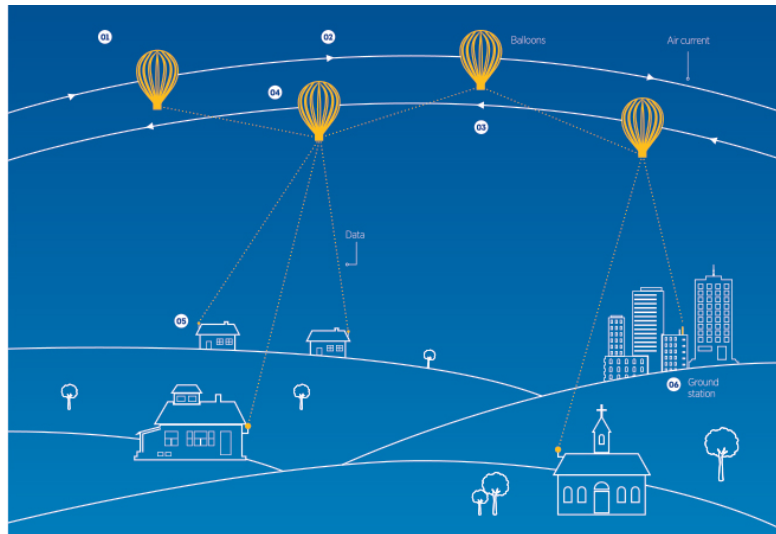
---

- Motivations, target applications & challenges
- Parametric regression
- Non-parametric regression
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems

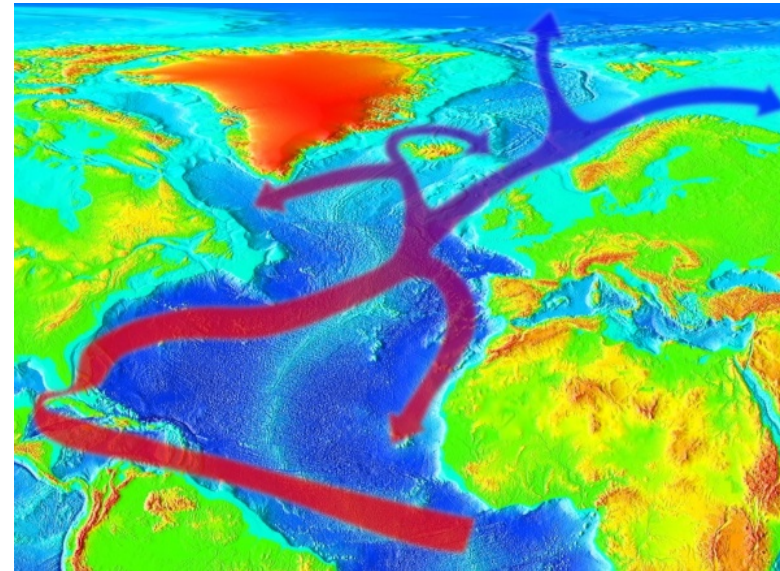


# Non-parametric regression for dynamical systems

Project Loon<sup>1</sup>



Wind/Ocean Current 4 Energy/Air Vehicles



# Time-varying regression $f(x,t)$ : parametric vs non-parametric

Single location:  $f(\bar{x}, t) : \mathbb{R} \rightarrow \mathbb{R}$

$$s_{k+1} = As_k + w_k, \quad w_k \sim \mathcal{N}(0, Q)$$

$$y_k = Cs_k + v_k, \quad v_k \sim \mathcal{N}(0, R)$$

$$f_k = Cs_k, \quad s_0 \sim \mathcal{N}(0, \Sigma_0)$$

Data:  $\{(t_k, y_k)\}_{i=1}^N, t_k = kT$

$$\hat{s}_{k+1|k} = A\hat{s}_{k|k}$$

$$\Sigma_{k+1|k} = A\Sigma_{k|k}A^T + Q$$

$$\hat{s}_{k+1|k+1} = \hat{s}_{k+1|k} + L_{k+1}(y_{k+1} - C_k\hat{s}_{k+1|k})$$

$$\Sigma_{k+1|k+1} = \Sigma_{k+1|k} - L_{k+1}C_k\Sigma_{k+1|k}$$

$$L_{k+1} = \Sigma_{k+1|k}C_k^T (C_k\Sigma_{k+1|k}C_k^T + R)^{-1}$$

$$\hat{s}_{k|k} = \mathbb{E}[s_k | y_k, \dots, y_0] \Rightarrow \hat{f}(\bar{x}, kT) = C\hat{s}_{k|k}$$

$$\hat{s}_{k|k} = \mathbb{E}[s_k | y_k, \dots, y_0] = \mathbb{E}[s_k | y_k, \hat{s}_{k-1|k-1}]$$

Numerical efficient but requires to  
know the exact model (A,C,Q,R)

Single instant:  $f(x, \bar{t}) : \mathbb{R}^p \rightarrow \mathbb{R}$

$$y_i = f(x_i) + v_i, \quad v_i \sim \mathcal{N}(0, \sigma^2)$$

$$f(x) \in \text{RKHS} \leftrightarrow k(x, x')$$

Data:  $\{(x_i, y_i)\}_{i=1}^N$

$$\hat{f}(x) = \sum_{i=1}^N c_i k(x_i, x)$$

$$c = (K + I)^{-1}y, \quad [K]_{mn} = k(x_m, x_n)$$

Model-free and best performance but  
high computational complexity  $O(N^3)$

# Time-varying regression $f(x,t)$ : parametric vs non-parametric

General case:  $f(x, t) : \mathbb{R}^p \times \mathbb{R} \rightarrow \mathbb{R}$

Stochastic PDEs:

$$\frac{\partial^2 f(x, t)}{\partial x^2} + \frac{\partial^2 f(x, t)}{\partial t^2} - \lambda^2 f(x, t) = w(x, t),$$

space-time white  
Gaussian noise



Data:  $\{(t_i, x_i, y_i)\}_{i=1}^N$

Discretize in time and space  
and run finite-element methods

Approximated solutions and requires  
a good physical model

Define  $\xi = (x, t)$  and kernel  $k(\xi, \xi')$ :

$$y_i = f(\xi_i) + v_i, \quad v_i \sim \mathcal{N}(0, \sigma^2)$$

$$f(\xi) \in \text{RKHS} \leftrightarrow k(\xi, \xi')$$

Data:  $\{(\xi_i, y_i)\}_{i=1}^N$

$$\hat{f}(\xi) = \sum_{i=1}^N c_i k(\xi_i, \xi)$$

$$c = (K + I)^{-1} y, \quad [K]_{mn} = k(\xi_m, \xi_n)$$

Time space treated equally and  
unbounded complexity  $O(t^3)$



# Combining parametric and non-parametric: Kalman filtering meets Machine Learning (1)

General case:  $f(x, t) : \mathbb{R}^p \times \mathbb{R} \rightarrow \mathbb{R}$

Two fundamental assumptions:

1. finite number of measurement locations  $\mathcal{X}_{meas} = \{x_1, \dots, x_M\}$
2. separable kernels:  $k(x, t, x', t') = k_s(x, x')h(\tau)$ ,  $\tau = t - t'$

$$\mathbf{f}_t = [f(x_1, t) \ f(x_2, t) \ \cdots \ f(x_M, t)]^T \quad \mathbf{f}_t : \mathcal{X}_{meas} \times \mathbb{R} \rightarrow \mathbb{R}^M$$

Fourier Transform  $S(\omega) = \mathcal{F}[h(\tau)]$

Spectral Factorization  $S(\omega) = W(i\omega)W(-i\omega)$

$$W(i\omega) = \frac{b_{r-1}(i\omega)^{r-1} + b_{r-2}(i\omega)^{r-2} + \cdots + b_0}{(i\omega)^r + a_{r-1}(i\omega)^{r-1} + \cdots + a_0}$$

State Space representation:

$$\dot{s}_t = F s_t + G w_t$$

$$z_t = H s_t$$

$$F = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ & & & \ddots & \\ 0 & 0 & 0 & \cdots & 1 \\ -a_0 & -a_1 & -a_2 & \cdots & -a_{r-1} \end{bmatrix}, \quad G = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix},$$

$$H = [b_0 \ b_1 \ b_2 \ \cdots \ b_{r-1}],$$

# Combining parametric and non-parametric: Kalman filtering meets Machine Learning (2)

$$\mathbf{f}_t = [f(x_1, t) \ f(x_2, t) \ \cdots \ f(x_M, t)]^T$$

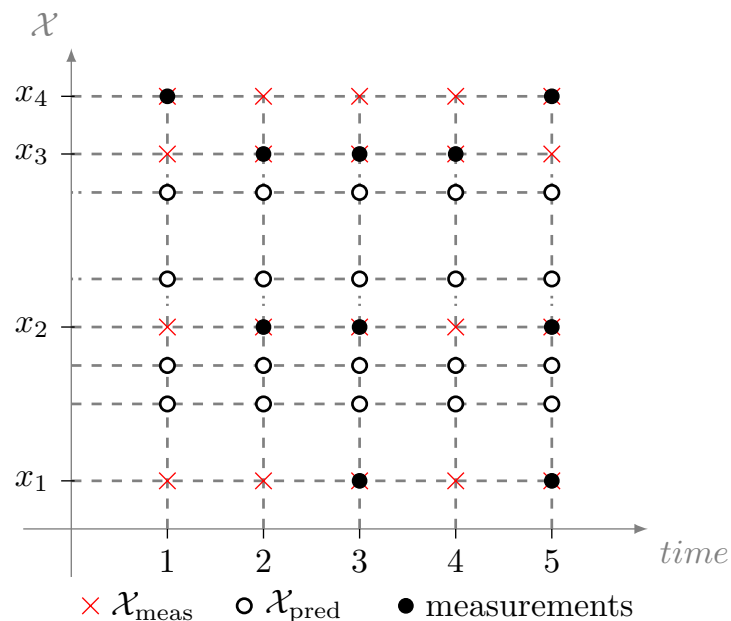
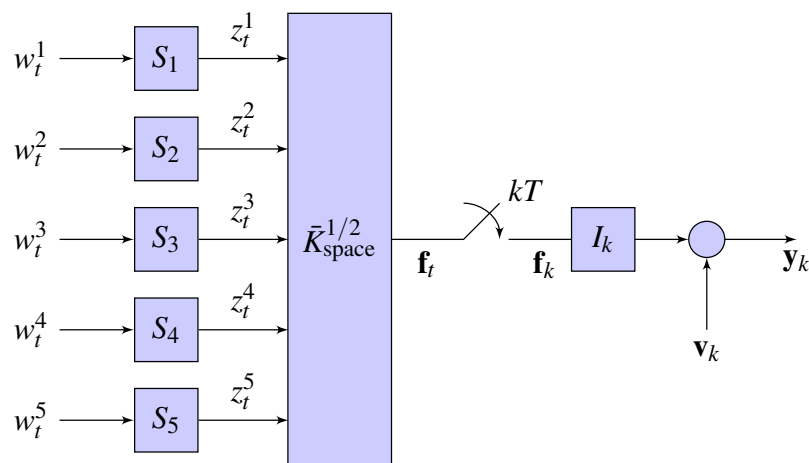
Dynamics of  $\mathbf{f}_t$ :

$$\dot{s}_t^j = F s_t^j + G w_t^j$$

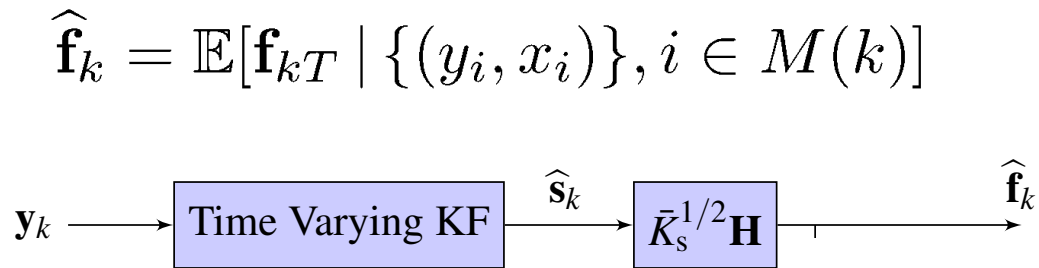
$$z_t^j = H s_t^j, \quad j = 1, \dots, M$$

$$\mathbf{f}_t = K_s^{1/2} \mathbf{z}_t, \quad [K_s]_{mn} = k_s(x_m, x_n)$$

$$y_k^j = f_k(x_j) + v_k^j \quad x_j \in M(k) \subset \mathcal{X}_{meas}, \quad v_k^j \sim \mathcal{N}(0, \sigma^2)$$



# Combining parametric and non-parametric: Kalman filtering meets Machine Learning (3)



$$\begin{aligned} \hat{s}_{k+1|k} &= A\hat{s}_{k|k} \\ \Sigma_{k+1|k} &= A\Sigma_{k|k}A^T + Q \\ \hat{s}_{k+1|k+1} &= \hat{s}_{k+1|k} + L_{k+1}(y_{k+1} - C_k\hat{s}_{k+1|k}) \\ \Sigma_{k+1|k+1} &= \Sigma_{k+1|k} - L_{k+1}C_k\Sigma_{k+1|k} \\ L_{k+1} &= \Sigma_{k+1|k}C_k^T (C_k\Sigma_{k+1|k}C_k^T + R)^{-1} \end{aligned}$$

$$A = \text{blkdiag}(\bar{F}, \dots, \bar{F}), \quad Q = \text{blkdiag}(\bar{Q}, \dots, \bar{Q})$$

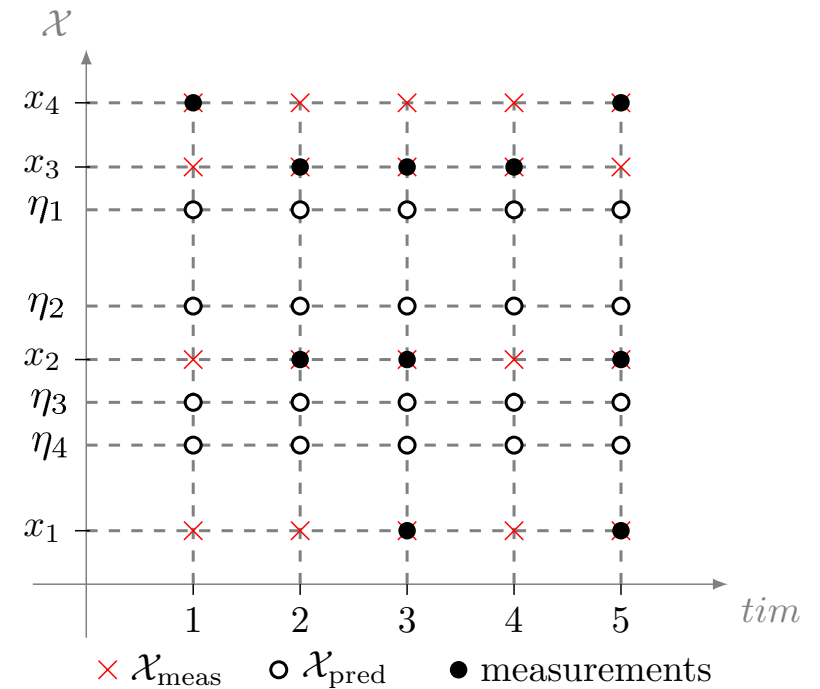
$$\bar{F} = e^{FT}, \quad Q = \int_0^T e^{F\tau}GG^T(e^{F\tau})^T d\tau, \quad R = \sigma^2 I$$

$$C_k := I_k K_s^{1/2} \mathbf{H}, \quad \mathbf{H} = \text{blkdiag}(H, \dots, H), \quad I_k \in \{0, 1\}^{M_k \times M}$$

$$\mathcal{X}_{pred} = \{x_{M+1}, \dots, x_{M+S}\}, \quad \mathcal{X}_{pred} \cap \mathcal{X}_{meas} = \emptyset$$

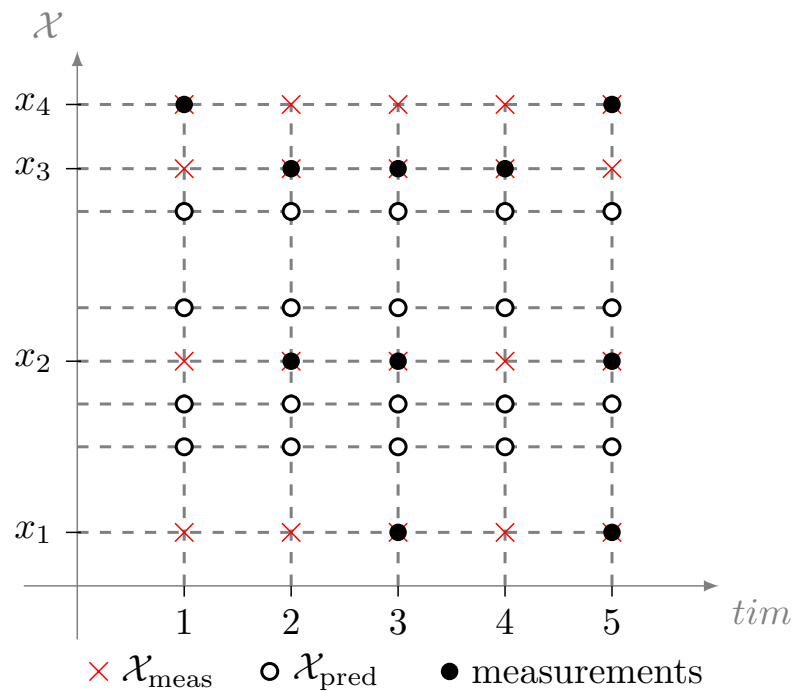
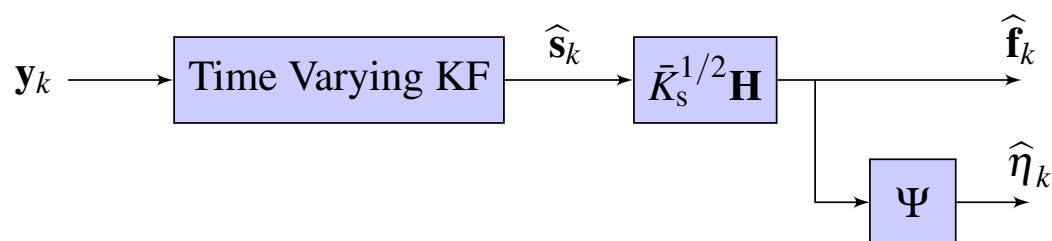
$$\eta_t = [f(x_{M+1}, t) \cdots f(x_{M+S}, t)]^T$$

$$\hat{\eta}_k = \Psi \hat{\mathbf{f}}_k, \quad \Psi \in \mathbb{R}^{S \times M}, \quad \Psi = \Psi(\mathcal{X}_{meas}, \mathcal{X}_{pred})$$



# Combining parametric and non-parametric: Kalman filtering meets Machine Learning (4)

$$\hat{\mathbf{f}}_k = \mathbb{E}[\mathbf{f}_{kT} \mid \{(y_i, x_i)\}, i \in M(k)]$$



# Truncated Gaussian regression vs Kalman-based Gaussian regression

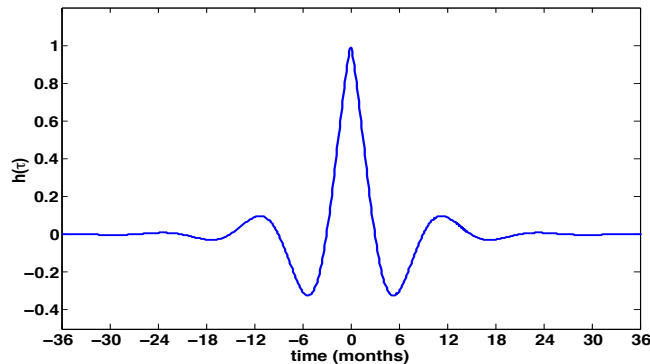
Colorado Weather Dataset: 365 stations, 100 years, monthly rain precipitation

$$K_S(x, x') = e^{-\sigma_s \|x - x'\|}, \quad \sigma_s = 0.5,$$

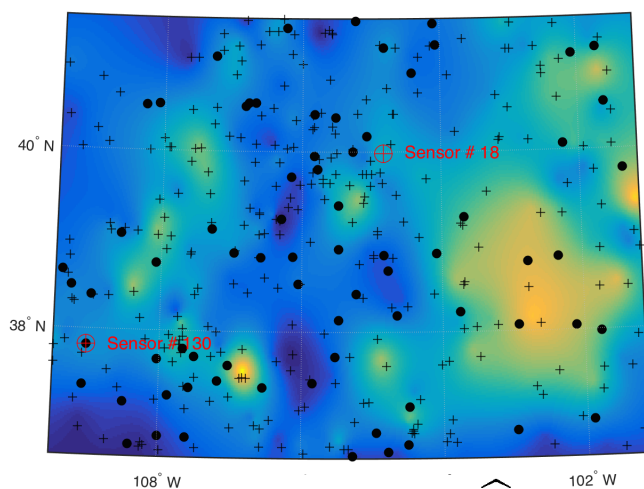
$$h(\tau) = \lambda \cos(2\pi f|\tau|)e^{-\sigma_t|\tau|}, \quad \lambda = 2 \times 10^3, \quad \sigma_t = 0.2,$$

Space kernel

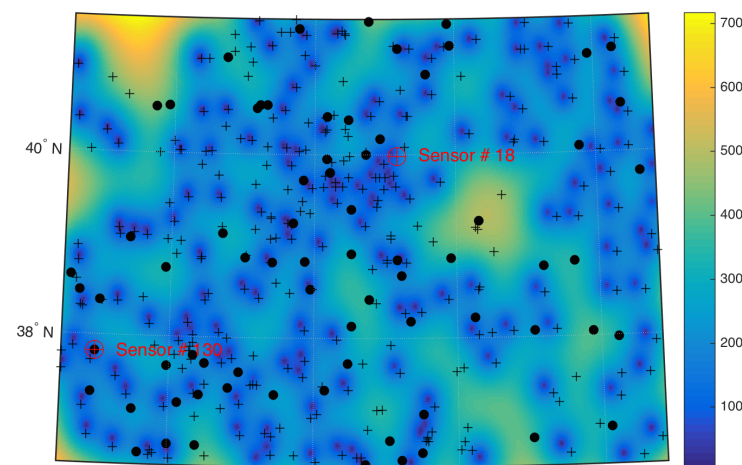
Time-Kernel



	Memory [MB]	CPU time [sec.]
Kalman-based Alg.1	4	0.02
Classical GP (all data)	$1.5 \cdot 10^6$	NA
Truncated GP (1 year data)	150	15
Truncated GP (2 years data)	600	120
Truncated GP (3 years data)	1300	410



Estimate  $\hat{f}$



Variance  $Var(\hat{f})$

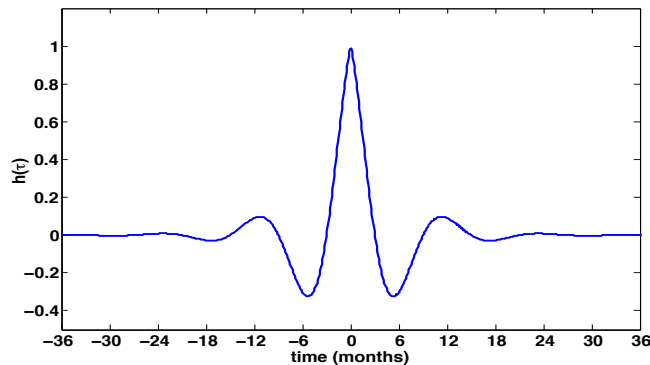
# Truncated Gaussian regression vs Kalman-based Gaussian regression

$$K_S(x, x') = e^{-\sigma_s \|x - x'\|}, \quad \sigma_s = 0.5,$$

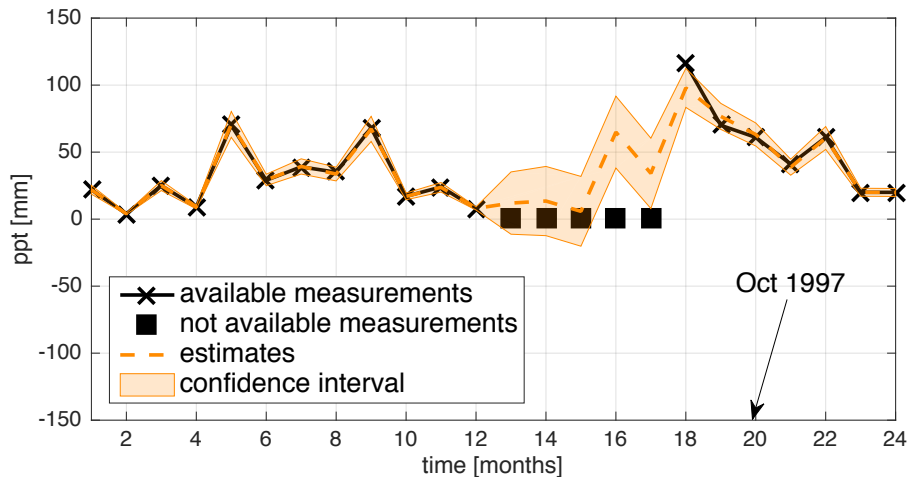
$$h(\tau) = \lambda \cos(2\pi f|\tau|) e^{-\sigma_t |\tau|}, \quad \lambda = 2 \times 10^3, \quad \sigma_t = 0.2,$$

Space kernel

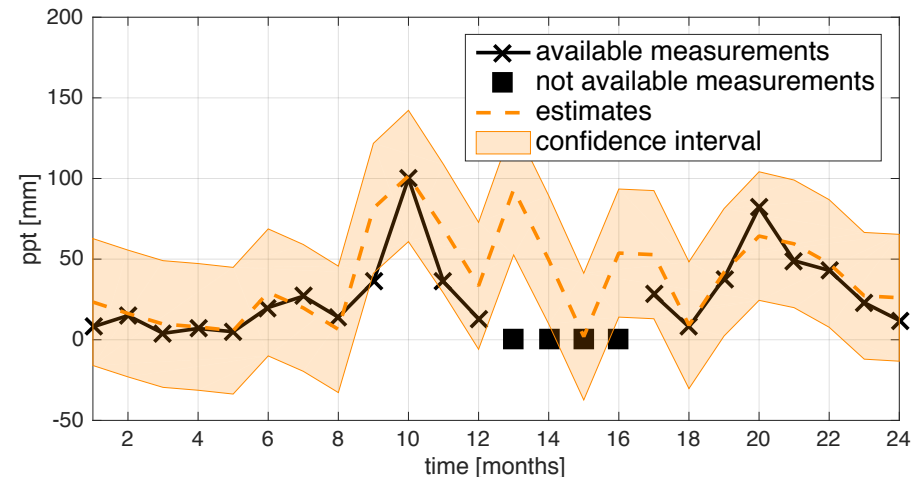
Time-Kernel



	Memory [MB]	CPU time [sec.]
Kalman-based Alg.1	4	0.02
Classical GP (all data)	$1.5 \cdot 10^6$	NA
Truncated GP (1 year data)	150	15
Truncated GP (2 years data)	600	120
Truncated GP (3 years data)	1300	410

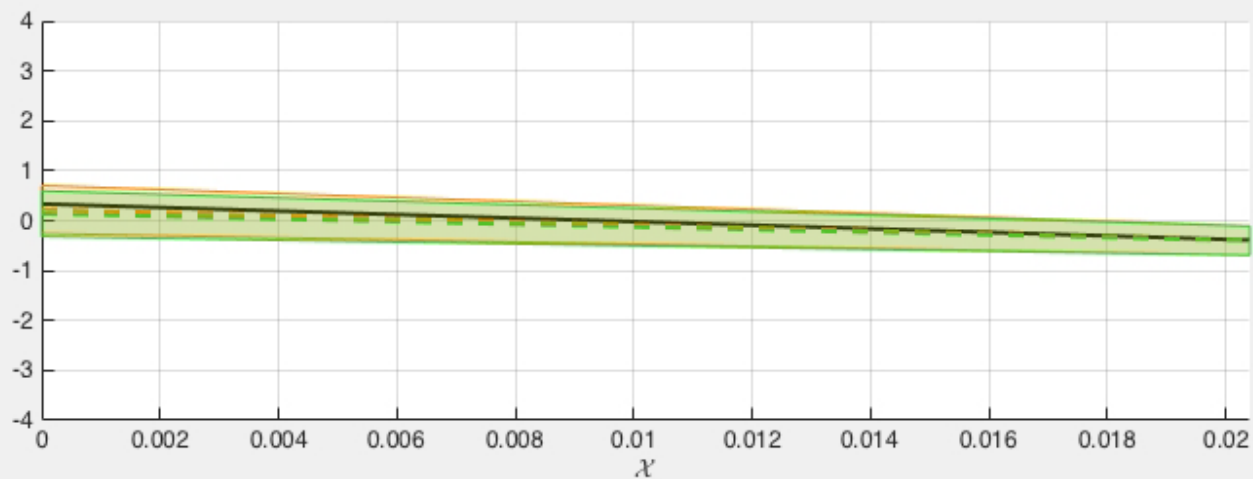
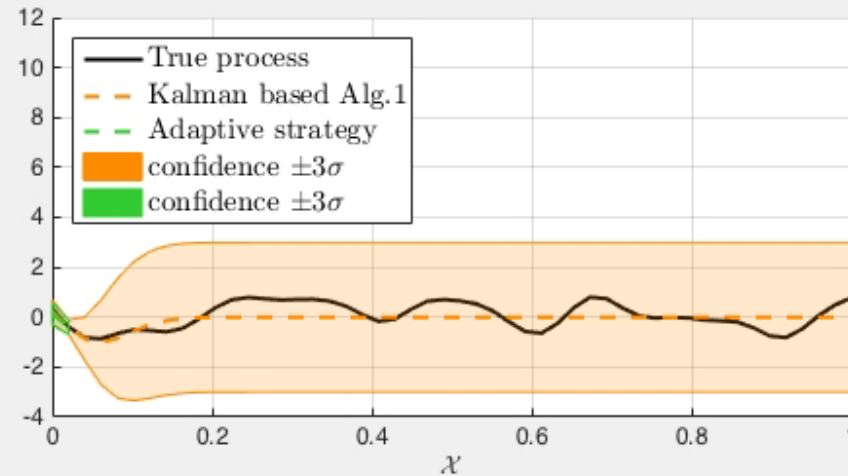


Estimate on measured location



Estimate on non-measured location

# Dynamic Grid (suboptimal solution)





# Outline

---

- Motivations, target applications & challenges
- Parametric regression
- Non-parametric regression
- Semi non-parametric regression
- Non-parametric regression for dynamical systems
- Conclusion and open problems



# Conclusions & open problems

- Non-parametric approach has great potential but it is unclear how to
  - make it distributed
  - incorporate time
  - adaptively design the sampling density, i.e.  $\mu(x) \propto \hat{f}(x)$
- Many details swept under the carpet:
  - Real-time and distributed design of regularization parameter for non-parametric approaches
  - Packet loss & asynchronous computation
  - Computation of eigenfunctions
- Integration of learning with control & optimization

# References

- Parametric vs non-parametric
  - D. Varagnolo, G. Pillonetto, L. Schenato. **Distributed parametric and nonparametric regression with on-line performance bounds computation.** *Automatica*, vol. 48(10), pp. 2468 -- 2481, 2012
- Cloud-based vs peer-to-peer
  - M. Todescato, A. Carron, R. Carli, G. Pillonetto, L. Schenato. **Multi-Robots Gaussian Estimation and Coverage Control: from Server-based to Distributed Architecture.** *Automatica [to appear]*
- Global vs Local estimation
  - D. Varagnolo, F. Zanella, A. Cenedese, G. Pillonetto, L. Schenato. **Newton-Raphson Consensus for Distributed Convex Optimization.** *IEEE Transactions on Automatic Control*, vol. 61(4), pp. 994--1009, 2016
  - A. Carron, M. Todescato, R. Carli, L. Schenato. **An asynchronous consensus-based algorithm for estimation from noisy relative measurements.** *IEEE Transactions on Control of Network Systems*, vol. 1(3), pp. 283 - 295, 2014
  - N. Bof, M. Todescato, R. Carli, L. Schenato. **Robust Distributed Estimation for Localization in Lossy Sensor Networks.** *6th IFAC Workshop on Distributed Estimation and control in Networked Systems (NecSys16)*, 2016
- Static vs dynamic maps:
  - M. Todescato, A. Carron, R. Carli, L. Schenato, G. Pillonetto. **Machine Learning meets Kalman Filtering.** *55th IEEE Conference on Decision and Control (CDC16)*

# Thank you

WC IFAC '17 (Toulouse)  
**Open Invited Track**

Multi-agent distributed learning and optimization of  
dynamical systems

Proponents: Ruggero Carli (Univ. Padova), Jongeun Choi (Yonsei University Seoul), Hideaki Ishii (Tokyo Institute of Technology), Jerome Le Ny (Polytechnique Montreal), Luca Schenato (Univ. Padova)