

Distributed non-parametric Gaussian regression: recent results and open problems

Damiano Varagnolo

Department of Information Engineering
University of Padova

May 27th, 2009



Part I

Introduction

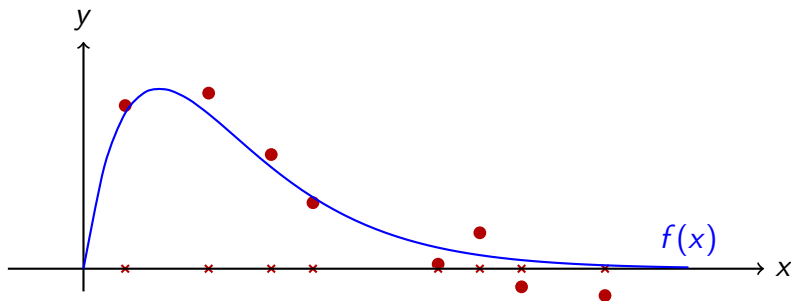


Problem statement - centralized scenario

inputs: set of noisy measurements of a certain signal:

$$y_m = f(x_m) + \nu_m \quad m = 1, \dots, M$$

goal: estimate $f(x)$

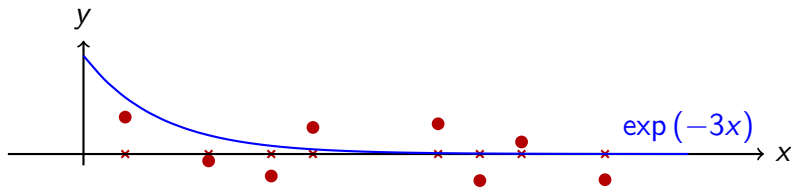


Parametric approach

assumption: known structure *but* unknown parameters

example: exponential:

$$f(x) = \exp(-\theta x) \quad \theta, x \in \mathbb{R}^+$$



goal: estimate θ starting from the data set $\{(x_m, y_m)\}$

\Rightarrow various approaches depending on the model on f :

- Maximum Likelihood
- Least Squares
- ...

Nonparametric approach

assumption: signal f lives in a certain functions space:

$$f \in \mathcal{H}_K$$

goal: search the estimate \hat{f} directly inside this space:

$$\hat{f} = \arg \min_{\tilde{f} \in \mathcal{H}_K} \left(\text{Loss function}(\tilde{f}, \{y_m\}) + \gamma \left\| \tilde{f} \right\|_{\mathcal{H}_K}^2 \right)$$

motivations: functional structure of f could be not easily managed with parametric structures

our approach: use Reproducing Kernel Hilbert Spaces



Part II

Centralized Learning



Introduction to the centralized learning scenario

- hypoteses:
- there is only one sensor
 - exist a certain and oportune “bidimensional” function

$$K(\cdot, \cdot) : \text{Input locations} \times \text{Input locations} \rightarrow \mathbb{R}$$

- working flow:
- $K(\cdot, \cdot)$ defines a function space \mathcal{H}_K
 - use $K(\cdot, \cdot)$ to construct the estimating function



Reproducing Kernel Hilbert Spaces - hypotheses

question: how \mathcal{H}_K is made?

first assumption: $K(\cdot, \cdot)$ is a Mercer Kernel:

- continuous
- symmetric
- definite positive

second assumption: input locations domain is compact



Reproducing Kernel Hilbert Spaces - implications

implication 1: $K(\cdot, \cdot)$ defines a compact linear positive definite integral operator:

$$(L_K f)(x_m) := \int_{\mathcal{X}} K(x_m, x') f(x') dx'$$

implication 2: L_K has at most a numerable set of eigenfunctions:

$$\phi_k(\cdot) = \lambda_k (L_K \phi_k)(\cdot) \quad k = 1, 2, \dots$$



Reproducing Kernel Hilbert Spaces

Theorem

with the previous hypotheses:

- $\{\lambda_k\}$ are real and non-negative: $\lambda_1 \geq \lambda_2 \geq \dots \geq 0$
- $\{\phi_k(\cdot)\}$ is an orthonormal basis for the space

$$\mathcal{H}_K = \left\{ f \in \mathcal{L}_2 \text{ s.t. } f = \sum_{k=1}^{\infty} a_k \phi_k \right. \\ \left. \text{with } \{a_k\} \text{ s.t. } \sum_{k=1}^{\infty} \frac{a_k \cdot a_k}{\lambda_k} < +\infty \right\}$$

$$\bullet f_1 = \sum_{k=1}^{\infty} a_k \phi_k, \quad f_2 = \sum_{k=1}^{\infty} b_k \phi_k \quad \Rightarrow$$

$$\langle f_1, f_2 \rangle_{\mathcal{H}_K} = \sum_{k=1}^{\infty} \frac{a_k \cdot b_k}{\lambda_k}$$

RKHS based learning

assumption: \mathcal{H}_K is defined via the kernel K

1° question: how we construct the estimate?

2° question: how can we interpret it?



RKHS based learning - cost-function interpretation

recall: estimation \hat{f} has to:

- fit \rightarrow loss functions
- not to overfit \rightarrow Tikhonov regularizer

our approach: loss functions = quadratic functions

$$\hat{f} = \arg \min_{\tilde{f} \in \mathcal{H}_K} \left(\sum_{\text{measurements } m} \frac{(\tilde{f}(x_m) - y_m)^2}{\sigma^2} + \gamma \|\tilde{f}\|_{\mathcal{H}_K}^2 \right)$$

why this choice?



RKHS based learning - Bayesian interpretation

goal: construct the Bayesian estimator of f :

$$\hat{f} = \text{cov}(f, \mathbf{y}) \text{var}(\mathbf{y})^{-1} \mathbf{y} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

Proposition

If:

- f is a **Gaussian** process with covariance K :

$$\text{cov} \left(f(x_m) f(x_n)^T \right) = K(x_m, x_n)$$

- loss functions = quadratic functions

then **cost-function regularization is equivalent to Bayes estimation**

RKHS based learning - numerical results

$$\hat{f}(\cdot) = \sum_{m=1}^M c_m K(x_m, \cdot) \quad \text{with} \quad \begin{bmatrix} c_1 \\ \vdots \\ c_M \end{bmatrix} = (K_{\mathbf{x}} + \gamma I_M)^{-1} \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix}$$

Example:

RKHS based learning - drawbacks

$$\hat{f}(\cdot) = \sum_{m=1}^M c_m K(x_m, \cdot)$$

1° feature: must invert $(K_x + \gamma I_M)^{-1}$

2° feature: must store $[c_1, \dots, c_M]$

caveat: M big \Rightarrow

- computationally hard to invert $M \times M$ matrices
- function representations pretty big

our requirements:

- compact representations (**necessary** for distributed algorithms)
- light computations (preferable)



Approximated RKHS learning

want compact representations & light computations?

⇒ **must approximate**

$$\text{model: } f = \sum_{k=1}^{+\infty} a_k \phi_k \quad \Rightarrow \quad y_m = \sum_{k=1}^{+\infty} a_k \phi_k(x_m) + \nu_m$$

new goal: estimate only the first E coefficients with $E \ll M$

definition: Reduced Hilbert Space:

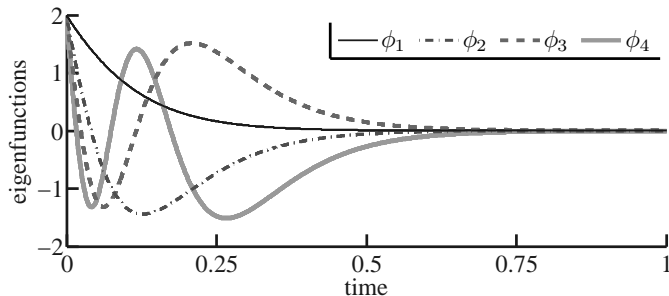
$$\mathcal{H}_K^E := \left\{ \tilde{f} \in \mathcal{L}_2 \text{ s.t. } \tilde{f} = \sum_{k=1}^E a_k \phi_k \right. \\ \left. \mathbf{a} := [a_1 \dots a_E]^T \in \mathbb{R}^E \right\}$$



Approximated learning - example

Kernel for BIBO stable linear time-invariant systems:

$$K(x, x'; \beta) = \begin{cases} \frac{\exp(-2\beta x)}{2} \left(\exp(-\beta x') - \frac{\exp(-\beta x)}{3} \right) & \text{if } x \leq x' \\ \frac{\exp(-2\beta x')}{2} \left(\exp(-\beta x) - \frac{\exp(-\beta x')}{3} \right) & \text{if } x \geq x' \end{cases}$$



Approximated learning - kind of approaches

cost-function based:

- data fitting \rightarrow loss functions
- not overfit \rightarrow Tikhonov regularizer

$$\hat{f} = \arg \min_{\tilde{f} \in \mathcal{H}_K^E} \left(\sum_{\text{measurements } m} \frac{(\tilde{f}(x_m) - y_m)^2}{\sigma^2} + \gamma \left\| \tilde{f} \right\|_{\mathcal{H}_K^E}^2 \right)$$

Bayesian approach:

- consider a prior
- find the best linear estimator

$$\hat{\mathbf{a}} = \text{cov}(\mathbf{a}, \mathbf{y}) \text{var}(\mathbf{y})^{-1} \mathbf{y}$$



Approximated learning - compact notation

notation: $y_m = \sum_{k=1}^{+\infty} a_k \phi_k(x_m) + \nu_m \quad \rightarrow \quad \mathbf{y} = \mathbf{C}\mathbf{a} + \mathbf{e} + \boldsymbol{\nu}$

definitions:

$$\mathbf{y} := \begin{bmatrix} y_1 \\ \vdots \\ y_M \end{bmatrix} \quad \mathbf{C} := \begin{bmatrix} \phi_1(x_1) & \dots & \phi_E(x_1) \\ \vdots & & \vdots \\ \phi_1(x_M) & \dots & \phi_E(x_M) \end{bmatrix}$$

$$\mathbf{a} := \begin{bmatrix} a_1 \\ \vdots \\ a_E \end{bmatrix} \quad \mathbf{e} := \begin{bmatrix} \sum_{k=E+1}^{+\infty} a_k \phi_k(x_1) \\ \vdots \\ \sum_{k=E+1}^{+\infty} a_k \phi_k(x_M) \end{bmatrix} \quad \boldsymbol{\nu} := \begin{bmatrix} \nu_1 \\ \vdots \\ \nu_E \end{bmatrix}$$



Approximated learning - numerical solutions

$$\left(\Sigma_{\mathbf{a}} := \text{diag}(\lambda_1, \dots, \lambda_E) \quad \Sigma_{\mathbf{e}} := \text{var}(\Sigma_{\mathbf{e}}) \right)$$

cost-function approach:

$$\hat{\mathbf{a}} = \left(\sigma^2 \Sigma_{\mathbf{a}} C^T C + \gamma I_E \right)^{-1} \Sigma_{\mathbf{a}} C^T \mathbf{y}$$

(computations load: $O(E^3 + E^2 M + EM^2)$ operations)

Bayesian approach:

$$\hat{\mathbf{a}} = \Sigma_{\mathbf{a}} C^T \left(C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \sigma^2 I_M \right)^{-1} \mathbf{y}$$

(computations load: $O(M^3)$ operations)



never equivalent!



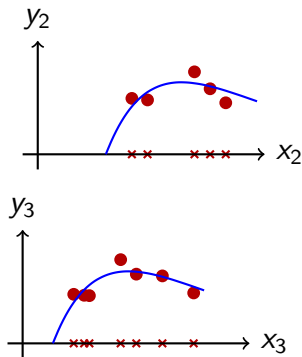
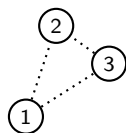
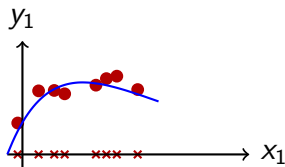
Part III

Distributed Approximated Learning



Distributed approximated learning - introduction

hypothesis: there are S sensors
all measuring the **same** function f



goal: distributely estimate f

constraints:

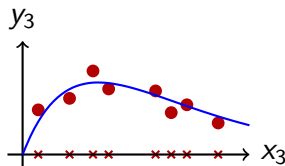
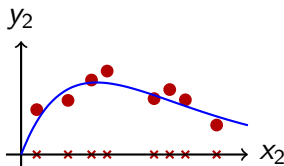
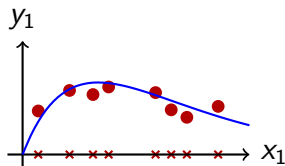
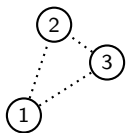
- limited amount of exchangeable information
- limited size of the representation of the estimated function

further hypothesis: consensus-based algorithms can be implemented

Distributed approx. learning - starting hypotheses

Simplifications:

- 1 each sensor knows **exactly** S (n° of sensors)
- 2 **no time-delay** between measured signals
- 3 **common input-locations grid** among sensors



Bayesian strategy: when distributed = centralized?

($\mathbf{y}_i :=$ set of measurements of sensor i)

$$\begin{aligned} \hat{\mathbf{a}}_{\text{cent}} &= \text{cov} \left(\mathbf{a}, \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_S \end{bmatrix} \right) \text{var} \left(\begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_S \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_S \end{bmatrix} \\ &= (\dots \text{some messages} \dots) \\ &= \frac{1}{S} \sum_{i=1}^S \left(\Sigma_{\mathbf{a}} C^T \left(C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \frac{\sigma^2}{S} I_M \right)^{-1} \mathbf{y}_i \right) \end{aligned}$$

\Rightarrow equivalent to an **average consensus** on **locally computable** quantities!



Bayesian strategies: distributed vs local

local Bayesian strategy:

$$\hat{\mathbf{a}}_{\text{loc},s} = \text{cov}(\mathbf{a}, \mathbf{y}_i) \text{var}(\mathbf{y}_i)^{-1} \mathbf{y}_i = \Sigma_{\mathbf{a}} C^T (C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \sigma^2 I_M)^{-1} \mathbf{y}_s$$

distributed strategy: (equivalent to centralized)

- 1 initial local estimation:

$$\hat{\mathbf{a}}_s(0) = \Sigma_{\mathbf{a}} C^T \left(C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \frac{\sigma^2}{S} I_M \right)^{-1} \mathbf{y}_s$$

- 2 average consensus on the various $\hat{\mathbf{a}}_s(0)$

difference = how to weight the measurement noise!



Gussed distributed strategy

hypothesis removal: sensors do **not** know S (n° of sensors)



all sensors make the same guess: \mathbf{S}_g (“g” = guess)

how distributed estimator changes?

distributed strategy:

- 1 initial local estimation:

$$\hat{\mathbf{a}}_s(0) (\mathbf{S}_g) = \Sigma_a C^T \left(C \Sigma_a C^T + \Sigma_e + \frac{\sigma^2}{\mathbf{S}_g} I_M \right)^{-1} \mathbf{y}_s$$

- 2 average consensus on the various $\hat{\mathbf{a}}_s(0) (\mathbf{S}_g)$



Comparison: centralized vs guessed distributed

centralized (or distributed) strategy:

$$\hat{\mathbf{a}}_{\text{cent}} = \Sigma_{\mathbf{a}} C^T \left(C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \frac{\sigma^2}{S} I_M \right)^{-1} \bar{\mathbf{y}}$$

guessed distributed strategy:

$$\hat{\mathbf{a}}_{\text{dist}}(\mathbf{S}_{\mathbf{g}}) = \Sigma_{\mathbf{a}} C^T \left(C \Sigma_{\mathbf{a}} C^T + \Sigma_{\mathbf{e}} + \frac{\sigma^2}{\mathbf{S}_{\mathbf{g}}} I_M \right)^{-1} \bar{\mathbf{y}}$$

centralized strategy use the correct measurements variance!



Comparisons between estimators performances

performance “=” estimation error variance

centralized vs local: centralized is always better than local

centralized vs guessed distributed: centralized is always better than guessed distributed (equal iff $S = \mathbf{S}_g$, (guess is correct))

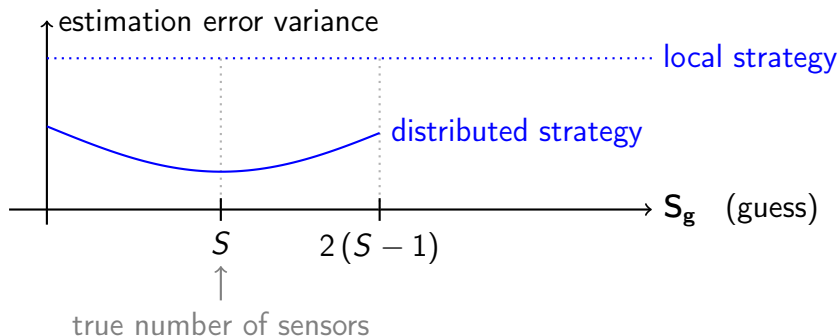
guessed distributed vs local: depends!!



Comparison: guessed distributed vs local (1)

Proposition

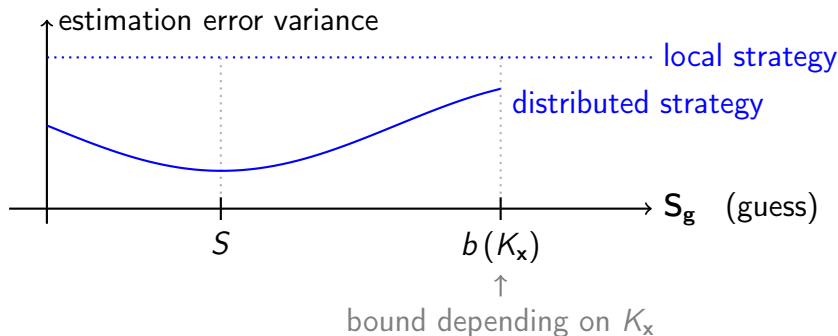
If $S_g \in [1, 2(S - 1)]$ then guessed distributed strategy is better than local independently of the kernel, noise power, number of measurements, etc.



Comparison: guessed distributed vs local (2)

Proposition

If we consider kernel then the previous bound can be enlarged



- 1° note: the bound depends on the eigenvalues of K_x
- 2° note: the bound is conservative

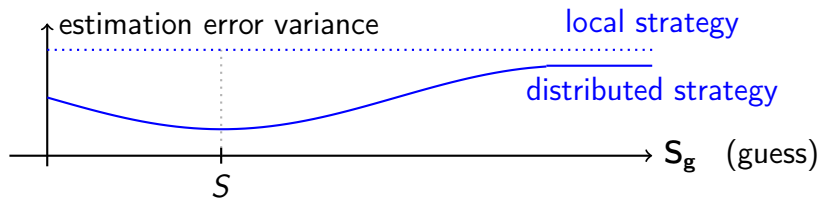
Comparison: guessed distributed vs local (3)

Proposition

If

$$\min(\text{eig}(K_x)) \geq \frac{\sigma^2}{S-1}$$

then the guessed distributed strategy is always better than the local one, for all guesses $S_g \in [1, +\infty)$

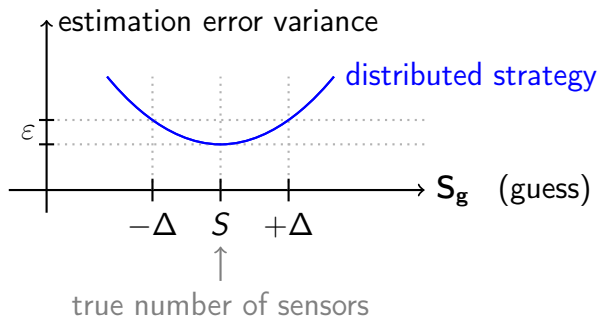


implication: in this case communications always improve estimation

Loss of performances: TODO

TODO: characterize the loss of performance when making wrong guesses (i.e. $S_g \neq S$)

desired propositions: small error in guessing (say Δ) $\stackrel{?}{\Rightarrow}$ loss of performance is small (say ε)? When? How much is ε ? Does it depend on K_x ? ...



Distributed estimation without initial guessings (1)

Initial guessing \mathbf{S}_g could be undesirable \Rightarrow look for estimators without required initializations

requested features:

- must not require guessings
- must be a linear transformation of the measurements \mathbf{y}_i
- must lead to the smallest possible estimation error variance
- dimension of exchanged vectors must be at most E



Distributed estimation without initial guessings (2)

Proposed algorithm:

- 1 whiten the noise $\mathbf{e} + \nu$
- 2 compress information using an SVD decomposition
- 3 run an average consensus algorithm

$$\rightarrow \hat{\mathbf{a}} = A(\mathbf{y}_1, \dots, \mathbf{y}_S)$$

“Corollary algorithm”: maximum likelihood estimator for the number of sensors S :

$$\hat{S}_{\text{ML}} := \arg \max_S P(A(\mathbf{y}_1, \dots, \mathbf{y}_S) \mid S)$$

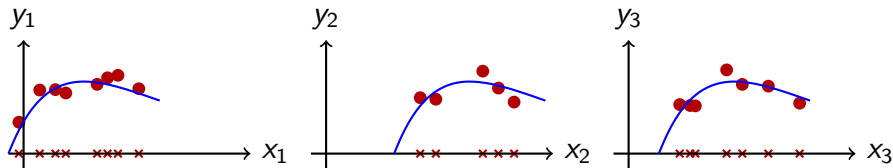
TODO: comparisons between this and “guessed distributed”:

- estimation error variance
- computational requirements



Adding temporal shifts to the measured function

hypothesis removal: there are unknown **time-delays** between measured signals



implication: unknown delays \Rightarrow no common sampling grid

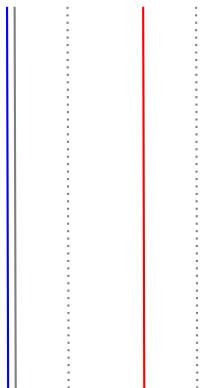
\Rightarrow much more difficult scenario!

Classic Time Delay Estimation

notation: $f_1, f_2 =$ noisy delayed versions of the same f

classic TDE: maximization of \mathcal{L}_2 's inner product:

$$\tau_{\text{optimal}} = \arg \max_{\tau} \langle f_1(x), f_2(x - \tau) \rangle_{\mathcal{L}_2}$$



Time Delay Estimation in RKHS framework

RKHS based TDE: maximization of \mathcal{H}_K 's inner product:

$$\tau_{\text{optimal}} = \arg \max_{\tau} \langle f_1(x), f_2(x - \tau) \rangle_{\mathcal{H}_K} = \arg \max_{\tau} \sum_{k=1}^{\infty} \frac{a_k \cdot b_k(\tau)}{\lambda_k}$$

Note: requires $f_1(x)$ and $f_2(x - \tau)$ in the same reference system



Joint “function and TD” Estimation - centralized scenario

proposed solution: Maximum likelihood strategy:

$$\mathcal{L}(\tilde{f}) := -J(\tilde{f}) = -\sum_{s=1}^S \sum_{m=1}^{M_s} \frac{(\tilde{f}(x_{s,m} - \tau_s) - y_{s,m})^2}{\sigma^2} - \gamma \|\tilde{f}\|_{\mathcal{H}_K}^2$$

implies:

$$\hat{f}_{ML} := \arg \max_{\tilde{f} \in \mathcal{H}_K^E} \mathcal{L}(\tilde{f})$$

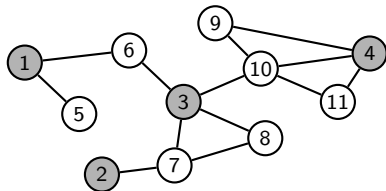
⇒ maximization via descent algorithms

Caveat: initialization strongly affects results!

Joint “function and TDE” - distributed scenario

proposed solution: distributed minimization of the centralized likelihood:

- 1 find the minimal *bridged sensor network* topology:



- 2 introduce some constraints in the centralized likelihood (one for each bridge)
- 3 construct a Lagrangian from the constrained likelihood
- 4 solve the Lagrangian via a distributed minimization algorithm

Caveat: initialization strongly affects results!

Part IV

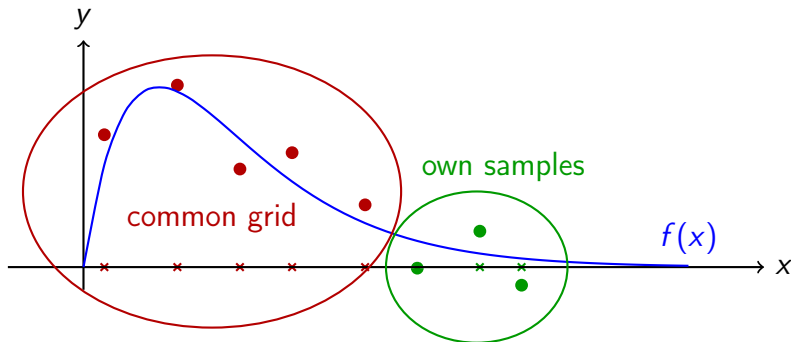
Other Open Problems



Grid based distributed estimation (1)

hypotheses:

- no time delays between measured functions
- sensors share a **common sampling grid**
- sensors have also some **own sampling locations**



Grid based distributed estimation (2)

proposed algorithm:

- 1 run distributed estimation on the grid (guessed / without guessing) $\rightarrow \hat{\mathbf{a}}_{\text{dist}}(\theta)$
- 2 fuse distributed estimation & data not previously used:
 $\hat{\mathbb{E}}[f \mid \hat{\mathbf{a}}_{\text{dist}}(\theta) \mathbf{y}_s]$

TODOs:

- characterize the combined estimation error variance (always better than pure local? no? when? why?)
- find suboptimal combination strategies of local estimates + distributed estimates (no recomputing everything)



Recursively updated approximated estimators

problematic issue: all current implementations of approximated estimators work “offline”:

single new measurement arrive \rightarrow recompute everything

desired strategy: find recursive equations:

$$\hat{\mathbf{a}}(t+1) = \Theta(\hat{\mathbf{a}}(t), (x_{t+1}, y_{t+1}))$$

TODO: everything (first step: find at least some suboptimal equations)



No-common-grid suboptimal strategy

previously proposed strategy: distributed lagrangian minimization →
slow convergence + local minima

suboptimal alternative idea:

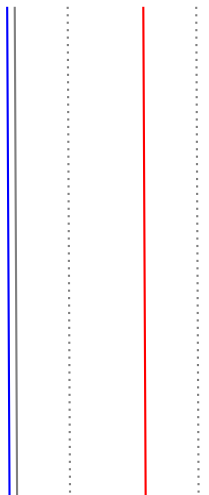
cycle the following

- 1 estimate the delays between the functions
- 2 construct an artificial grid
- 3 *create* some measurements on this grid
- 4 run distributed strategy on this grid
- 5 update the local estimations using the distributed one

TODO: everything (numerical equations, convergence, stability, error bounds, ...)

Time Delay Estimation using RKHS techniques

TDE via RKHS-based approximated representations: in our knowledge **never been proposed**



Distributed Fault Detection (1)

hypotheses:

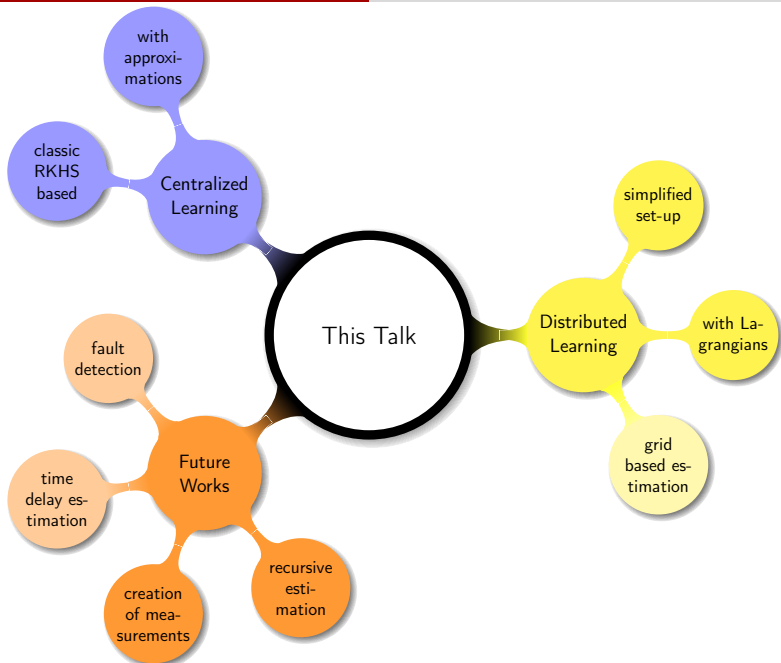
- sensors measure the same (or quite the same) function
- sensors can be faulty

Distributed Fault Detection (2)

proposed solution: each sensor:

- do distributed estimation
- do local estimation
- compare local and distributed estimations







Tack

damiano.varagnolo@dei.unipd.it

www.dei.unipd.it/~varagnolo

entirely written in **ETEX-2_ε** using **Beamer** and **Tik Z**

