

Auto-tuning procedures for distributed nonparametric regression algorithms

Damiano Varagnolo, Gianluigi Pillonetto, Luca Schenato

Abstract—We propose a distributed regression algorithm with the capability of automatically calibrating its parameters during its on-line functioning. The estimation procedure corresponds to a Regularization Network, i.e., the structural form of the estimator is a linear combination of basis functions which coefficients are computed by solving a linear system. The automatic tuning strategy instead constructs and then exploits opportune bounds on the distance between the distributed estimation results and the unknown centralized optimal estimate that would be computed processing the whole dataset at once. By numerical simulations we show how the proposed procedure allows the sensor networks to effectively self-tune the parameters of the distributed regression scheme by simple consensus strategies.

Index Terms—distributed regression, distributed calibration, self-organizing sensor networks, regularization networks, non-parametric estimation

I. INTRODUCTION

Applications like surveillance, monitoring, tracking and sensing, benefit of the distributed paradigm, where unmanned agents perform auxiliary and automatic operations. But to broaden the applicability of distributed paradigms, and to increase their robustness with respect to human error, algorithms should be self-configuring and self-tuning; these are indeed intermediate steps for implementing self-organizing and truly smart sensors and actuators networks.

Towards this vision we consider a specific class of distributed estimation strategies, more specifically nonparametric regression algorithms. Our interests in contributing to this field is indeed driven by some practical considerations, that make us believe in their technological possibilities: *i*) nonparametric strategies may be statistically more effective than parametric ones (e.g., identification of linear systems using Akaike Information Criterion plus Prediction Error Methods [1]); *ii*) nonparametric approaches may be consistent where parametric approaches fail to be [2], [3]; *iii*) nonparametric methods usually require the tuning of very few parameters, and this allows the implementation of fast calibration strategies [4]. We moreover specifically consider scenarios where agents have limited communication bandwidth, so that representations of the estimated quantities must be kept small.

Literature review: endowing nonparametric distributed estimators with self- and online-calibration capabilities is complicated by the fact that the regularization parameters (γ in the following Equation (5)), typical of nonparametric

strategies, combine with global quantities that are generally unknown to the single agents, such as the total number of measurements available in the whole network.

Up to now, and to the best of our knowledge, the problem of how to address this lack of information, and thus of how to tune regularization parameters of distributed nonparametric estimators in a online fashion, has not been treated. We recognize several implementations of ad-hoc distributed self-calibration / self-diagnosis strategies, e.g., [5], [6], [7], [8], [9], and literature on the calibration of centralized nonparametric estimators, e.g., [10, Chap. 5], [11, Chap. 7], but for distributed settings the usual approach is to assume the regularization parameter (or the parameters governing the sparsification rules) to be fixed and computed off-line [12], [13], [14], [15].

Statement of contributions: there are then two ways to overcome the lack of information on global quantities like the number of measurements in the network: either distributedly estimate this information, or bypass it and exploit some other structural property of the distributed nonparametric regression framework.

Here we consider the second approach, and devise on-line tuning procedures that are based on opportune Euclidean distances concepts. More specifically, we consider opportune a-posteriori probabilistic bounds on the distance between the outputs of the distributed regression strategy and the centralized optimal one. We notice that the proposed strategies do not follow iterative minimization procedures, but rather compare in parallel a set of different parameters and then choose the optimal one.

Organization of the manuscript: Sec. II describes the considered regression framework, while Secs. III and IV describe respectively a centralized nonparametric estimator and its distributed version. Secs. V-A and V-B introduce then the distributed procedures for the calibration of the parameters of the regression strategy. We conclude with numerical examples in Sec. VI and with some conclusions and indications of future works in Sec. VII. To improve the readability of the paper, the proofs have been collected in the appendix.

Notice that, to the best of our knowledge, strategies for the automatic tuning of the parameters of distributed nonparametric regression algorithms have never presented before. We are thus not able to offer comparative results with some other literature works.

II. REGRESSION FRAMEWORK

Let $f_\mu : \mathcal{X} \rightarrow \mathbb{R}$ denote an unknown function defined on the compact $\mathcal{X} \subset \mathbb{R}^d$. For brevity, and w.l.o.g. (the same

D. Varagnolo is with the Department of Computer Science, Electrical and Space Engineering, Luleå University of Technology, Luleå, Sweden. Email: damiano.varagnolo@ltu.se. G. Pillonetto and L. Schenato are with the Department of Information Engineering, University of Padova, Italy. Email: { giapi | schenato }@dei.unipd.it.

1 derivations could be made by letting the sensors collect more
 2 information), assume that there are S sensors, each collecting
 3 a single noisy measurement y_i , i.e.,

$$4 \quad y_i = f_\mu(x_i) + \nu_i, \quad i = 1, \dots, S \quad (1)$$

5 with ν_i white noise and i the sensor index. We assume that
 6 each input location x_i is known only to the i -th sensor and
 7 that it is independently drawn from a probability measure μ
 8 known to all the sensors.

9 Notice that hereafter we will use the following notation:
 10 • f_μ is the unknown function that has to be estimated; • f
 11 is a generic function; • f_c is a centralized estimate of f_μ ; •
 12 f_d is a distributed estimate of f_μ .

13 III. CENTRALIZED REGRESSION

14 Given the data set $\{x_i, y_i\}_{i=1}^S$, one of the most used ap-
 15 proaches to estimate f_μ relies upon the Tikhonov regular-
 16 ization theory [16], [17]. The hypothesis space is typically
 17 given by a reproducing kernel Hilbert space (RKHS) defined
 18 by a Mercer Kernel $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ [18], [19], [20] that
 19 is spanned by the eigenfunctions¹ ϕ_e of the positive integral
 20 operator

$$21 \quad \int_{\mathcal{X}} K(x, x') g(x') d\mu(x') \quad (2)$$

22 where the corresponding eigenvalues λ_e are s.t. $\lambda_1 \geq \lambda_2 \geq$
 23 $\dots \geq 0$. Under mild assumptions (see, e.g., [21]), the
 24 hypothesis space is given by the Hilbert space

$$25 \quad \mathcal{H}_K := \left\{ g \in \mathcal{L}^2(\mu) \text{ s.t. } g = \sum_{e=1}^{\infty} \alpha_e \phi_e \right. \\ \left. \text{with } \{\alpha_e\} \text{ s.t. } \sum_{e=1}^{\infty} \frac{\alpha_e^2}{\lambda_e} < +\infty \right\}. \quad (3)$$

26 Letting $g_1 = \sum_{e=1}^{\infty} \alpha_e \phi_e$ and $g_2 = \sum_{e=1}^{\infty} \beta_e \phi_e$, this implies
 27 that the inner product in \mathcal{H}_K is

$$28 \quad \langle g_1, g_2 \rangle_K := \sum_{e=1}^{\infty} \frac{\alpha_e \beta_e}{\lambda_e} \quad (4)$$

29 with the λ_e 's the eigenvalues of the kernel K .

30 To define the estimator of f_μ given the dataset
 31 $\{(x_i, y_i)\}_{i=1, \dots, S}$, a commonly used cost function is

$$32 \quad Q(f) := \sum_{i=1}^S (y_i - f(x_i))^2 + \gamma \|f\|_K^2 \quad (5)$$

33 where γ is the so called *regularization parameter* that trades
 34 off empirical evidence and smoothness information on f_μ .
 35 Assume w.l.o.g. γ to be known (cf. the discussion at the
 36 beginning of Sec. V). It is known that the optimal estimate

$$37 \quad f_c := \arg \min_{f \in \mathcal{H}_K} Q(f) \quad (6)$$

38 admits the structure of a Regularization Network, see [19],
 39 being the sum of S basis functions with expansion coeffi-
 40 cients obtainable by inverting a system of linear equations.

¹For numerical computation of eigenvalues and eigenfunctions see for
 example [10, Chap. 4.3.2].

41 IV. DISTRIBUTED REGRESSION

42 A potential strategy for computing f_c over networks is to
 43 route all the information to a specific unit, and let that unit
 44 perform the computations. Since this requires the processing
 45 unit to perform $O(S^3)$ operations and to store all the x_i 's,
 46 generally this strategy is impractical in distributed scenarios,
 47 where agents may have both limited computational and
 48 communication resources.

49 We thus aim at deriving an alternative approach, more
 50 suitable for distributed settings. To this aim we consider the
 51 following roadmap:

- rewrite the optimization problem (6) in an alternative
 52 but equivalent way, by exploiting the structure of \mathcal{H}_K ;
- change, thanks to Principal Components Analysis-like
 53 concepts, the hypothesis space from \mathcal{H}_K to an approx-
 54 imated one;
- derive the distributed estimator as an approximated
 55 version of the centralized one.

59 A. Rewriting optimization problem (6)

60 Let \mathbb{R}^∞ be the space of vectors with an infinite number
 61 of real scalar components. Introducing the map $T : \mathcal{H}_K \rightarrow$
 62 \mathbb{R}^∞ associating to a function $f(\cdot) = \sum_{e=1}^{+\infty} a_e \phi_e(\cdot)$ in \mathcal{H}_K
 63 the sequence $[a_1, a_2, \dots]$ of its eigenfunctions weights, it is
 64 possible to rewrite the estimand f_μ as the novel estimand
 65 $b_\mu = T[f_\mu]$. Of course b_μ and f_μ are equivalent.

66 Letting moreover

$$67 \quad C_i := [\phi_1(x_i) \ \phi_2(x_i) \ \dots], \quad (7)$$

68 it is possible to rewrite the measurement model (1) as

$$69 \quad y_i = C_i b_\mu + \nu_i, \quad i = 1, \dots, S, \quad (8)$$

70 and the cost function (5) as

$$71 \quad Q(b) := \sum_{i=1}^S (y_i - C_i b)^2 + \gamma \|b\|_K^2. \quad (9)$$

72 The optimal estimate $b_c := \arg \min_{b \in \mathbb{R}^\infty} Q(b)$ of the esti-
 73 mand b_μ is thus (see also [4])

$$74 \quad b_c = \left(\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \right)^{-1} \left(\sum_{i=1}^S C_i^T y_i \right) \quad (10)$$

75 with $\text{diag}(\alpha_e)$ indicating the matrix with diagonal elements
 76 given by $\alpha_1, \alpha_2, \dots$

77 B. Changing the hypothesis space

78 The optimal estimate b_c in (10) is infinite dimensional,
 79 and thus numerically intractable. To obtain a numerically
 80 tractable estimator, we consider the most natural finite-
 81 dimensional alternative of \mathcal{H}_K , i.e., the subspace \mathcal{H}_K^E gen-
 82 erated by the first E eigenfunctions ϕ_e , i.e.,

$$83 \quad \mathcal{H}_K^E := \left\{ g \in \mathcal{L}^2(\mu) \text{ s.t. } g = \sum_{e=1}^E \alpha_e \phi_e \right. \\ \left. \text{with } [\alpha_1, \dots, \alpha_E]^T \in \mathbb{R}^E \right\}. \quad (11)$$

84 Substituting \mathcal{H}_K with \mathcal{H}_K^E is then motivated by the presence
 85 of the penalty term $\|\cdot\|_K^2$ in (5): from Bayesian viewpoints,

1 \mathcal{H}_K^E represents the subspace that, before seeing the data,
 2 captures the biggest part of the signal variance among all the
 3 subspaces of dimension E [22], [10], in accordance with the
 4 Rayleigh's principle which underlies Principal Component
 5 Analysis [23].

6 C. Deriving the distributed estimator

7 Given the change from the hypothesis space \mathcal{H}_K to \mathcal{H}_K^E ,
 8 consider also the change from C_i in (7) to

$$9 \quad C_i^E = C^E(x_i) := [\phi_1(x_i), \dots, \phi_E(x_i), 0, 0, \dots], \quad (12)$$

10 and from the cost function (9) to

$$11 \quad Q^E(b) := \sum_{i=1}^S (y_i - C_i^E b)^2 + \gamma \|b\|_K^2. \quad (13)$$

12 In this case the optimal estimate of b_μ using \mathcal{H}_K^E as hypoth-
 13 esis space is then given by (see also [4])

$$14 \quad b_r := \arg \min_{b \in \mathcal{H}_K^E} Q(b) = \arg \min_{b \in \mathcal{H}_K^E} Q^E(b) \\
 15 = \left(\frac{1}{S} \text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E \right)^{-1} \left(\frac{1}{S} \sum_{i=1}^S (C_i^E)^T y_i \right) \quad (14)$$

16 Thus, if sensors know the number of measurements S
 17 and the regularization parameter γ , then b_r can be dis-
 18 tributedly computed through two parallel average consensus
 19 algorithms: one on $(C_i^E)^T C_i^E$ and one on $(C_i^E)^T y_i$, plus
 20 multiplications and inversions of $E \times E$ matrices and E -
 21 dimensional vectors.

22 But even if sensors know the number of measurements
 23 S and the regularization parameter γ , as noticed in [24],
 24 the distributed implementation of (14) may still be problem-
 25 atic since it requires $O(E^2)$ -communication and $O(E^3)$ -
 26 computational costs, i.e., to exchange an amount of informa-
 27 tion that scales with the square of E , potentially too high.
 28 To this aim it is possible to consider that

$$29 \quad \frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E \approx \mathbb{E}_\mu \left[(C_i^E)^T C_i^E \right] = \text{diag}(I, 0) \quad (15)$$

30 where I is $E \times E$ -dimensional, and 0 is infinite dimensional.
 31 This equivalence is guaranteed by the fact that for $1 \leq$
 32 $m, n \leq E$

$$33 \quad \left[\frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E \right]_{mn} = \frac{1}{S} \sum_{i=1}^S \phi_m(x_i) \phi_n(x_i) \quad (16)$$

and, that, due to the orthogonality of the eigenfunctions of
 the kernel K in $\mathcal{L}^2(\mu)$ and the fact that the x_i 's are i.i.d.
 and extracted from μ ,

$$34 \quad \frac{1}{S} \sum_{i=1}^S \phi_m(x_i) \phi_n(x_i) \xrightarrow{S \rightarrow +\infty} \int_{\mathcal{X}} \phi_i(x) \phi_j(x) d\mu(x) = \delta_{ij}.$$

35 This means that b_r can be approximated with

$$36 \quad b_d := \text{diag} \left(\frac{\lambda_e}{\gamma/S + \lambda_e} \right) \left(\frac{1}{S} \sum_{i=1}^S (C_i^E)^T y_i \right), \quad (17)$$

an estimator that is particularly suitable for distributed es-
 37 timation purposes since it does neither require sensors to
 38 exchange information on their input locations x_i (i.e., the
 39 C_i^E) nor to compute matrix inversions; it only requires an
 40 average consensus on the E -dimensional vectors $(C_i^E)^T y_i$.
 41

42 V. AUTOTUNING PROCEDURES

43 Consider estimator b_d in (17). This estimator is parametrized
 44 in the number of eigenfunctions E , the regularization pa-
 45 rameter γ , and the total number of measurements in the
 46 network S . E drives the computational and communication
 47 requirements of the distributed strategy, but also the accuracy
 48 of the final estimate (as noticed in Sec. IV-B. The ratio γ/S ,
 49 instead, dictates how much the empirical evidence of the final
 50 solution should be traded off with its smoothness.

51 In practical situations, both E and γ/S should be chosen
 52 *a-posteriori*, i.e., after that sensors have collected their y_i .
 53 The aim of this paper is then the following: considering S
 54 and E as unknown (γ can instead w.l.o.g. be considered
 55 known, or arbitrarily be set to 1), develop in-line strategies so
 56 that sensors will find a guess S_g for S and for E maximizing
 57 in some sense the performance of b_d .

In other words we highlight this parametric dependency
 of b_d on S_g and E by writing

$$b_d = b_d(S_g, E),$$

58 and thus propose a distributed in-line self-calibration tech-
 59 nique that allows the sensors to opportunely select E and
 60 S_g assuming that the y_i 's are locally available. The details
 61 of this strategy are offered in the following sections, and are
 62 based on the following mild assumption:

Assumption 1 $S \in [S_{\min}, S_{\max}]$ and sensors have knowl-
 edge about S_{\min} and S_{\max} .

Remark 2 Even if γ and S are known, because of the
 additional noise coming from the approximation $I \approx$
 $\frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E$, it can be shown that in general, for
 any fixed E , implementing b_d with the exact S does not
 maximize the predictive capabilities of b_d . So, even if S is
 actually known, one may want to find on-line that S_g that
 maximizes the statistical performance of b_d .

A. Calibration of the Regularization Parameter

63 Assume for now E to be fixed, and write $b_d(S_g)$ instead
 64 of $b_d(S_g, E)$. Despite the fact that, for any finite number
 65 of measurements, it may happen that an opportunely tuned
 66 $b_d(S_g)$ has better predictive capabilities of the centralized
 67 optimal estimate b_c , usually b_c has bigger generalization
 68 capabilities of $b_d(S_g)$ for any $S_g \in \mathbb{R}_+$. It is then meaningful
 69 to consider $\|b_d(S_g) - b_c\|_2$ as a performance indicator, and
 70 try to tune S_g seeking to minimize this distance.
 71

72 Importantly, in actual distributed estimation scenarios it is
 73 impossible to compute

$$74 \quad S_g^* := \arg \min_{S_g \in \mathbb{R}_+} \|b_d(S_g) - b_c\|_2. \quad (18)$$

1 since b_c is unknown. It is thus necessary to proceed find-
 2 ing appropriate bounds for $\|b_d(S_g) - b_c\|_2$ that depend on
 3 S_g , and then find S_g^* minimizing these bounds. The first
 4 step is given by the following proposition, that bounds
 5 $\|b_d(S_g) - b_c\|_2$ with terms that can then be computed
 6 by agents independently. (The numerical validity of these
 7 bounds is analyzed in Sec. VI.)

Proposition 3 Let

$$C_i^{\setminus E} := [0, \dots, 0, \phi_{E+1}(x_i), \phi_{E+2}(x_i), \dots] \quad (19)$$

$$\gamma_a := \sup_{x \in \mathcal{X}} \left\| \text{diag} \left(\frac{\lambda_e}{\gamma} \right) (C_i^{\setminus E}(x))^T \right\|_2 \quad (20)$$

$$\gamma_b := \sup_{x \in \mathcal{X}} \left\| \text{diag} \left(\frac{\lambda_e}{\gamma} \right) (C_i^{\setminus E}(x))^T C_i^E(x) \right\|_2 \quad (21)$$

$$V_r := \left(\frac{1}{S} \text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E \right)^{-1} \quad (22)$$

$$V_d(S_g) := \left(\frac{1}{S_g} \text{diag} \left(\frac{\gamma}{\lambda_e} \right) + I \right)^{-1} \quad (23)$$

$$U_C := I - \frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E \quad (24)$$

$$U_S(S_g) := \left(\frac{1}{S_g} - \frac{1}{S} \right) \text{diag} \left(\frac{\gamma}{\lambda_e} \right). \quad (25)$$

Then

$$\|b_d(S_g) - b_r\|_2 \leq \|V_r U_S(S_g) b_d(S_g)\|_2 + \|V_r U_C b_d(S_g)\|_2 \quad (26)$$

and

$$\|b_d(S_g) - b_c\|_2 \leq (\gamma_b S_{\max} + 1) \|b_d(S_g) - b_r\|_2 + \sum_{i=1}^S \gamma_a \|y_i - C_i^E b_d(S_g)\|_2 \quad (27)$$

8 The terms involved in Prop. 3 have the following interpre-
 9 tations:

- 10 • $C_i^{\setminus E}$ is the part of the transformation expressed in (8)
 11 corresponding to the discarded eigenfunctions;
- 12 • γ_a and γ_b respectively bound how much the residuals
 13 $y_i - C_i^E b_d(S_g)$ and $b_d(S_g) - b_r$ will influence the overall
 14 approximation error $b_d(S_g) - b_c$;
- 15 • V_r is s.t. $\frac{1}{S} V_r^{-1}$ is an approximation of the true covari-
 16 ance of the set of measurements $\{y_i\}$. More precisely,
 17 $\frac{1}{S} V_r^{-1}$ would be the actual covariance if $\lambda_{E+1} =$
 18 $\lambda_{E+2} = \dots = 0$. The smaller these eigenvalues are,
 19 the better $\frac{1}{S} V_r^{-1}$ is an approximation of the actual
 20 covariance;
- 21 • $V_d(S_g)$ corresponds to an opportune approximation of
 22 V_r ;
- 23 • U_C corresponds to the approximation error encountered
 24 replacing $\frac{1}{S} \sum_{i=1}^S (C_i^E)^T C_i^E$ with $\mathbb{E}_\mu \left[(C_i^E)^T C_i^E \right]$;

- 25 • $U_S(S_g)$ modulates how the error on the regularization
 26 parameter affect the regularization properties of the
 27 proposed distributed estimator.

The usefulness of Prop. 3 is that it is possible to
 28 build on top of it to construct the following bound for
 29 $\|b_d(S_g) - b_c\|_2$:
 30

$$\mathcal{B}(S_g) := (\gamma_b S_{\max} + 1) \left(\|V_r U_S(S_g) b_d(S_g)\|_2 + \right. \\ \left. + \|V_r U_C b_d(S_g)\|_2 \right) + \sum_{i=1}^S \gamma_a \|y_i - C_i^E b_d(S_g)\|_2. \quad (28)$$

One would then want to optimize on-line the unknown
 32 parameter S_g through
 33

$$S_g^* := \arg \min_{S_g \in \mathbb{R}_+} \mathcal{B}(S_g); \quad (29)$$

nonetheless $\mathcal{B}(S_g)$ cannot be directly used for computing S_g
 35 since the quantities V_r , $U_S(S_g)$, U_C and S are unknown to
 36 the various sensors.
 37

To cope with this lack of information we propose thus to:

- 1) majorize $U_S^*(S_g)$ with $U_S^*(S_g)$, defined as
 39

$$U_S^*(S_g) := \max \left(\left| \frac{1}{S_g} - \frac{1}{S_{\max}} \right|, \left| \frac{1}{S_g} - \frac{1}{S_{\min}} \right| \right) \cdot \text{diag} \left(\frac{\gamma}{\lambda_e} \right) \quad (30)$$

and exploiting Assumption 1. Indeed it is immediate to
 check that

$$U_S^*(S_g) \geq U_S(S_g) \quad \forall S_g \in \mathbb{R}_+$$

where the inequality is in a matricial positive definite
 41 sense.
 42

- 2) majorize V_r and U_C with quantities that are generated
 43 locally by each sensor i as follows: a) locally simulate a
 44 particular scenario of the network by locally generating
 45 S_{\min} independent virtual input locations $x_{i,j}$ by means
 46 of density μ , i.e., each i generates
 47

$$x_{i,j} \sim \mu \quad \text{where} \quad j = 1, \dots, S_{\min}. \quad (31)$$

- b) then each i locally computes

$$C_{i,j}^E := [\phi_1(x_{i,j}), \dots, \phi_E(x_{i,j})],$$

$$V_{r,i}^* := \left(\frac{1}{S_{\max}} \text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \frac{1}{S_{\max}} \sum_{j=1}^{S_{\min}} (C_{i,j}^E)^T C_{i,j}^E \right)^{-1} \quad (32)$$

$$U_{C,i}^* := \left(I - \frac{1}{S_{\min}} \sum_{j=1}^{S_{\min}} (C_{i,j}^E)^T C_{i,j}^E \right), \quad (33)$$

i.e., from probabilistic viewpoints, generate $V_{r,i}^*$ and
 53 $U_{C,i}^*$ as pessimistic but informative versions of the true
 54 and unknown V_r and U_C .
 55

By means of the previous scheme, optimization of S_g is
 56 now then possible through solving
 57

$$S_g^* := \arg \min_{S_g \in \mathbb{R}_+} \mathcal{B}^*(S_g) \quad (34)$$

1 where

$$\begin{aligned}
\mathcal{B}^*(S_g) := & (\gamma_b S_{\max} + 1) \cdot \frac{1}{S} \sum_{i=1}^S \left(\|V_{r,i}^* U_S^*(S_g) b_d(S_g)\|_2 \right. \\
& \left. + \|V_{r,i}^* U_{C,i}^* b_d(S_g)\|_2 \right) \\
& + (\gamma_a S_{\max}) \cdot \frac{1}{S} \sum_{i=1}^S \|y_i - C_i^E b_d(S_g)\|_2.
\end{aligned} \tag{35}$$

2
3 Intuitively, thus, agents try to minimize a pessimistic
4 estimate $\mathcal{B}^*(S_g)$ of $\mathcal{B}(S_g)$ instead of $\mathcal{B}(S_g)$ itself. The
5 complete algorithm is then reported in Alg. 1, solving
6 problem (35) by gridding, i.e., selecting the best S_g from
7 a set of candidates $S_g^{(1)}, \dots, S_g^{(P)}$.

Algorithm 1 Distributed calibration of the regularization parameter

Off-line work: Sensors are given $S_{\min}, S_{\max}, \mu, E, \gamma_a, \gamma_b$, a set of R different candidates $S_g^{(1)}, \dots, S_g^{(P)}$ and relative matrices $U_S^*(S_g^{(1)}), \dots, U_S^*(S_g^{(P)})$. In addition, each sensor i locally generates S_{\min} independent virtual input locations $x_{i,j}, j = 1, \dots, S_{\min}$ by means of density μ , from which it computes $C_{i,j}^E, V_{r,i}^*$ and $U_{C,i}^*$.

On-line and distributed work:

- 1: (distributed step) sensors distributedly compute, by means of average consensus protocols, the E -dimensional vector

$$\mathcal{Z} := \frac{1}{S} \sum_{i=1}^S (C_i^E)^T y_i \tag{36}$$

- 2: (local step) each sensor i computes the P versions of the estimator (17), namely $b_d(S_g^{(p)}) = V_d(S_g^{(p)}) \mathcal{Z}$, for $p = 1, \dots, P$.
- 3: (local step) each sensor i computes the local P auxiliary scalars, for $p = 1, \dots, P$

$$\begin{aligned}
\mathcal{B}_i^*(S_g^{(p)}) := & (\gamma_b S_{\max} + 1) \left\| V_{r,i}^* U_S^*(S_g^{(p)}) b_d(S_g^{(p)}) \right\|_2 \\
& + (\gamma_b S_{\max} + 1) \left\| V_{r,i}^* U_{C,i}^* b_d(S_g^{(p)}) \right\|_2 \\
& + (\gamma_a S_{\max}) \cdot \left\| y_i - C_i^E b_d(S_g^{(p)}) \right\|_2
\end{aligned}$$

- 4: (distributed step) sensors distributedly compute, by means of average consensus protocols, the P scalars, for $p = 1, \dots, P$

$$\mathcal{B}^*(S_g^{(p)}) := \frac{1}{S} \sum_{i=1}^S \mathcal{B}_i^*(S_g^{(p)}) \tag{37}$$

- 5: (local step) each sensor i computes $S_g^* = S_g^{(p^*)}$ where

$$(p^*) = \arg \min_{(p)} \mathcal{B}^*(S_g^{(p)}) \tag{38}$$

B. Calibration of the Number of Eigenfunctions

The maximum admissible value for E is upper bounded by computational complexity and transmission capability constraints. Assuming \bar{E} to be this maximum value, the usage of a naïve strategy like $E = \bar{E}$ could lead to resource wasting. In the following algorithm 2 we offer a practical and general guideline for the choice of E exploiting pessimistic bounds on the approximation error $\|b_c - b_r\|_2$.

Algorithm 2 Calibration of the number of eigenfunctions

- 1: assume the knowledge of a lower bound on the energy of the unknown signal f_μ , indicated with $\min \|f_\mu\|_2$
- 2: choose a threshold δ for the maximal tolerable error $\frac{\|b_c - b_r\|_2}{\|f_\mu\|_2}$
- 3: compute the minimal value of E s.t.

$$\frac{3\sigma S_{\max} \gamma_a(E)}{\min \|f_\mu\|_2} \leq \delta \tag{39}$$

where we highlighted the dependence of γ_a on E .

From a practical point of view, the algorithm returns a number E assuring the operator that the normalized approximation error $\frac{\|b_c - b_r\|_2}{\|f_\mu\|_2}$ is smaller than a certain threshold. The algorithm is derived from the consideration that inequality (55) in the proof of Prop. 3 implies

$$\|b_c - b_r\|_2 \leq \gamma_a \sum_{i=1}^S \|y_i - C_i b_r\|_2 \tag{40}$$

and the consideration that, in general, residuals $\|y_i - C_i b_r\|_2$ are far smaller than 3 times the standard deviation of the measurement noise. We notice that this choice is arbitrary and relies on the assumption that the estimation result will have a certain minimum level of generalization capabilities. Pessimistic considerations can lead to increase the number of standard deviations, with the limit case of no approximation capabilities of b_r corresponding to set $b_r = 0$ in (40) and to substitute 3σ with $\max_i \|y_i\|_2$ in (39).

We notice that, substituting $\min \|f_\mu\|_2$ with $\max_i \|y_i\|_2$ in (39), algorithm 2 can be used in a-posteriori scenarios, where sensors decide E by means of a max consensus on $\|y_i\|_2$ before computing (36). We also notice that high uncertainties on S lead to overestimations of E because of the approximation S_{\max} .

VI. NUMERICAL EXAMPLES

In this section we show the effectiveness of the proposed strategies through some numerical examples. We consider $f_\mu : \mathcal{X} = [0, 1] \rightarrow \mathbb{R}$ to be given by

$$f_\mu(x) = \sum_{n=1}^{100} \alpha_n \sin(\omega_n x) \tag{41}$$

with $\alpha_n \sim \mathcal{N}(0, 0.01)$ i.i.d., $\omega_n \sim \mathcal{U}[0, 25]$ i.i.d., $\mu \sim \mathcal{U}[0, 1]$ and a measurement noise standard deviation $\sigma =$

1 0.75 s.t., on average, $\text{SNR} := \frac{\text{var}(f_\mu)}{\sigma^2} \approx 2.5$. Moreover we
 2 consider the Gaussian kernel

$$3 \quad K(x, x') = \exp\left(-\frac{(x - x')^2}{0.02}\right) \quad (42)$$

4 with the estimators (10) and (17) defined by $\gamma = 0.3$.

5 To show the effectiveness of the estimation strategy (17), a
 6 randomly generated realization of f_μ is sampled by $S = 100$
 7 sensors and estimated using $E = 20$ eigenfunctions² under
 8 two different uncertainty levels on S , namely case (a), where
 9 $S_{\min} = 90$ and $S_{\max} = 110$, and case (b), where $S_{\min} = 20$
 10 and $S_{\max} = 2000$. In Fig. 1 we plot then the actual realization
 11 (solid line), its estimates reconstructed from b_c (dashed line)
 12 and $b_d^{(\cdot)}(S_g^*)$ with S_g^* chosen by Algorithm 1 among 20
 13 candidates logarithmically spaced inside $[1, S_{\max}]$, and $(\cdot) =$
 14 (a) or (b) accordingly to the level of uncertainty on S (dotted
 and dashed-dotted lines, respectively). We claim an overall

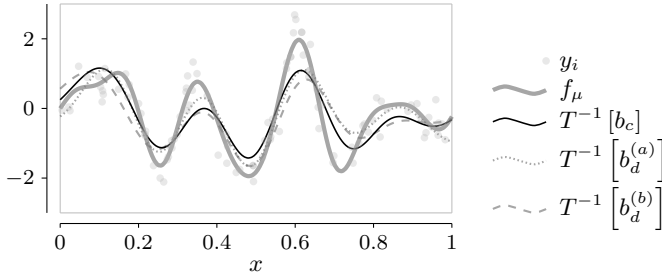


Fig. 1. Effectiveness of the estimation strategy (17) on a randomly generated f_μ , for various levels of uncertainty on S .

15 insensitivity of b_d on the uncertainty on S considering that
 16 both $T^{-1}[b_d^{(a)}]$ and $T^{-1}[b_d^{(b)}]$ are close to the centralized
 17 estimate $T^{-1}[b_c]$.

18 Despite this valuable property, bounds \mathcal{B}^* are good indi-
 19 cators about the actual distance $\|b_d(S_g^*) - b_c\|_2$ only for
 20 the case (a) (low uncertainty on S), as Fig. 2 indicates. In
 21 this figure we generate 200 independent realizations of f_μ ,
 22 then estimate each f_μ as before, and finally plot the actual
 23 distances $\|b_d^{(\cdot)} - b_c\|_2$ versus the obtained bounds \mathcal{B}^* . It is
 24 immediate to see that the bound provides, for the case (b),
 25 meaningless information on the actual distance. This lack of
 26 meaningfulness is caused by the presence in the bound of the
 27 multiplicative factor S_{\max} . This implies that in general the
 28 accuracy of the bound is tightly connected with the accuracy
 29 of the knowledge on S .

30 For sake of completeness, we show in Fig. 3 the values of
 31 the bounds $\mathcal{B}^*(S_g^{(p)})$ defined in (37) associated to the exper-
 32 iment of Fig. 1, and the relative distances $\|b_d(S_g^{(p)}) - b_c\|_2$.
 33 It is possible to see how the qualitative behavior of curve
 34 $\mathcal{B}^*(S_g^{(p)})$ is similar to the one of curve $\|b_d(S_g^{(p)}) - b_c\|_2$.
 35
 36

²This particular choice will be motivated later.

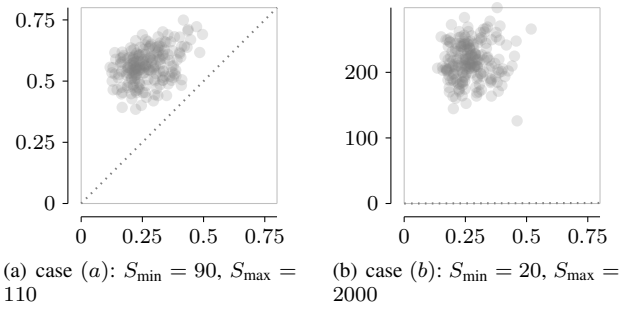


Fig. 2. Actual distances $\|b_d(S_g^*) - b_c\|_2$ vs. bounds values \mathcal{B}^* for different levels of uncertainty on S .

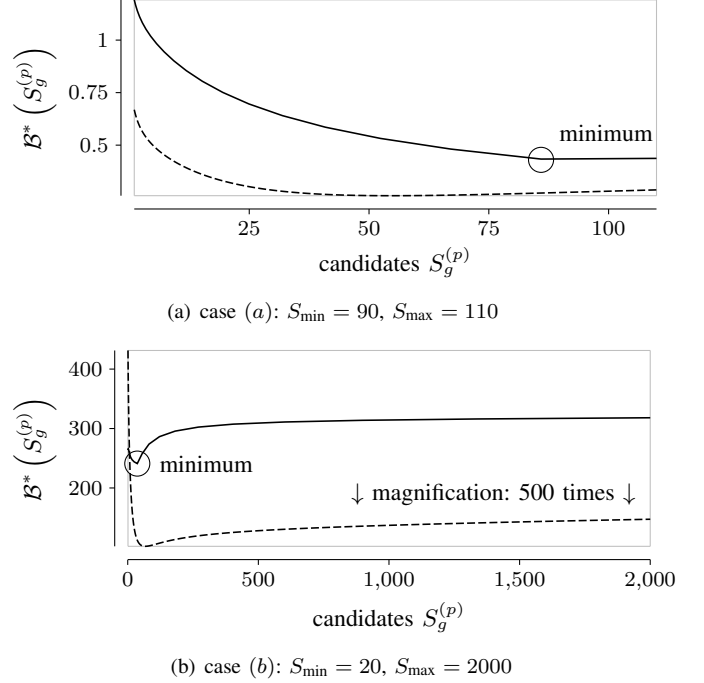


Fig. 3. Values of the bounds $\mathcal{B}^*(S_g^{(p)})$ under different uncertainty levels on S for the experiment of Fig. (1) (solid lines), and relative values of the distances $\|b_d(S_g^{(p)}) - b_c\|_2$ (dashed lines). Circles on the solid lines indicate the optimal values \mathcal{B}^* . The dashed line in panel (b) has been magnified 300 times.

37 We then aim to check if it is better to use Alg. 1 or to try
 38 to directly try to estimate S . We thus compare the estimation
 39 performance obtainable with three different naïve strategies
 40 for the choice of S_g , namely $S_g^* = S_{\min}$, $S_g^* = S_{\max}$, $S_g^* =$
 41 $S_{\text{ave}} := \frac{S_{\min} + S_{\max}}{2}$. Considering panels (a) of Figs. 2 and 3
 42 it is possible to infer that:

- in case of low uncertainty levels, Alg. 1 will not lead to big improvements w.r.t. to naïve strategies, but will give accurate descriptions of the actual distance with the centralized estimate;
- in case of high uncertainty levels, Alg. 1 will not give accurate descriptions of the actual distance with the centralized estimate but its usage will lead to improvements w.r.t. to naïve strategies.

To numerically prove the last statement, we consider the previously generated 200 independent realizations of f_μ and the case $S_{\min} = 20$ and $S_{\max} = 2000$. We then plot in Fig. 4 the 100 points

$$\left(\|b_d(S_{\min}) - b_c\|_2, \|b_d(S_g^*) - b_c\|_2 \right) \quad (43)$$

$$\left(\|b_d(S_{\text{ave}}) - b_c\|_2, \|b_d(S_g^*) - b_c\|_2 \right) \quad (44)$$

$$\left(\|b_d(S_{\max}) - b_c\|_2, \|b_d(S_g^*) - b_c\|_2 \right). \quad (45)$$

in panels (a), (b) and (c) respectively. Since these points generally lie below the bisector of the first quadrant, the distributed estimators b_d with S_g chosen with Alg. 1 are generally closer to the centralized estimates b_c than the ones with naïvely chosen S_g s. Finally, to check the level of suboptimality of the results of Alg. 1, in panel (d) of the same figure we plot also the points

$$\left(\|b_d(S_g^{\text{ora}}) - b_c\|_2, \|b_d(S_g^*) - b_c\|_2 \right) \quad (46)$$

where³ S_g^{ora} are the optimal S_g s obtained exactly solving problem (18). Since the distance of these points from the bisector is small, we can conclude that the level of suboptimality of Alg. 1 is also small.

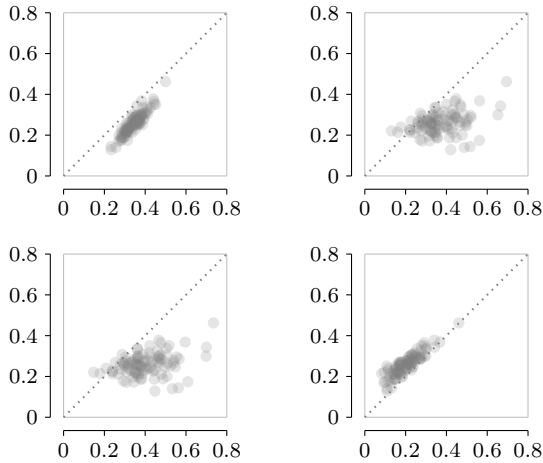


Fig. 4. Scatter plots to test the effectiveness of Alg. 1. Left-up panel: scatter plots of the points defined in (43). Right-up panel: points defined in (44). Left-down panel: points defined in (45). Right-down panel: points defined in (46). $S_{\min} = 20$ and $S_{\max} = 2000$.

To test the effectiveness of Alg. 2 and motivate the previous choice $E = 20$, we plot in Fig. 5 the values of E returned by the on-line version of Alg. 2, applied to the experiment of Fig. 1 with $S_{\min} = 20$ and $S_{\max} = 2000$, and fed with various values for the threshold δ . We notice that the exponential decay of the bound is inherited by the exponential decay of eigenvalues λ_e associated to the Gaussian kernel. Different kernels would lead to different outputs. Notice that if we let $\delta = 10^{-3}$ we obtain $E = 20$ and thus motivate the previous choice.

³We use the superscript “ora” pretending that solution has been provided by an *oracle*.

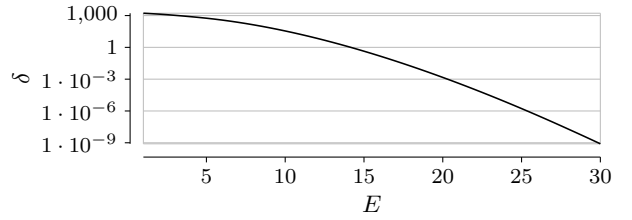


Fig. 5. Values of E returned by the on-line version of Alg. 2 fed with various choices of the threshold δ and applied to the experiment of Fig. 1 with $S_{\min} = 20$ and $S_{\max} = 2000$.

VII. CONCLUSIONS

In this paper we analyze how to endow distributed non-parametric regression strategies with self-tuning capabilities. The considered estimator is characterized by two parameters: the first one, the regularization parameter, that trades off the empirical evidence and the smoothness information on the true function. The second one, the number of eigenfunctions to be used, determines the size of the hypothesis space. Here we constructed a novel distributed and on-line parameters self-calibration strategy exploiting opportune a-posteriori probabilistic bounds on the distance between the parametrized distributed estimator and the unknown estimate that would be computed in a centralized scenario.

We also analyzed the performances of this distributed parameters calibration strategy through numerical experiments, and shown that under highly uncertain topological knowledge, the strategy leads to improvements with respect to naïve calibration strategies. On the contrary, in case of accurate knowledge on the number of sensors in the network, the computed probabilistic bounds constitute an accurate description of the distance between the distributed regression strategy and an optimal centralized one.

As examples of future works, we notice that the proposed strategy can be ameliorated exploiting statistical knowledge about the number of sensors in the network. Moreover, the strategy can be extended in order to compute on the fly the minimal number of eigenfunctions guaranteeing a certain regression quality.

APPENDIX

Proof (Proof of Prop. 3) We rewrite (14) as

$$V_r^{-1} b_r = \mathcal{Z} \quad (47)$$

and (17) as

$$(V_r^{-1} + V_d^{-1}(S_g) - V_r^{-1}) b_d(S_g) = \mathcal{Z}. \quad (48)$$

Subtracting (48) to (47) we then obtain

$$b_r - b_d(S_g) = V_r (V_d^{-1}(S_g) - V_r^{-1}) b_d(S_g) \quad (49)$$

from which it immediately follows that

$$\|b_d - b_r(S_g)\|_2 = \|V_r (V_d^{-1}(S_g) - V_r^{-1}) b_d(S_g)\|_2. \quad (50)$$

1 Defining then U_C and U_S by means of (24) and (25), it
 2 is immediate to check that $V_d^{-1}(S_g) - V_r^{-1} = U_S(S_g) +$
 3 U_C from which inequality (26) immediately follows.

4 To prove (27), we rewrite (14) as

$$\begin{aligned} & \left(\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \right) b_r + \left(\sum_{i=1}^S (C_i^E)^T C_i^E - \sum_{i=1}^S C_i^T C_i \right) b_r \\ & = \sum_{i=1}^S C_i^T y_i - \sum_{i=1}^S (C_i^{\setminus E})^T y_i \end{aligned} \quad (51)$$

5 and (10) as

$$\left(\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \right) b_c = \sum_{i=1}^S C_i^T y_i. \quad (52)$$

8 After subtracting (52) to (51), we obtain

$$\begin{aligned} & \left(\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \right) (b_c - b_r) = \\ & = \left(\sum_{i=1}^S (C_i^E)^T C_i^E - \sum_{i=1}^S C_i^T C_i \right) b_r + \sum_{i=1}^S (C_i^{\setminus E})^T y_i. \end{aligned} \quad (53)$$

9 Substituting now each C_i in the right side of (53) with $C_i^E +$
 10 $C_i^{\setminus E}$, exploiting the fact that $C_i^{\setminus E} b_r = 0$ (where 0 is an
 11 infinite dimensional vector of zeros), and properly collecting
 12 the various terms, we obtain

$$\begin{aligned} b_c - b_r = & \\ & \left(\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \right)^{-1} \sum_{i=1}^S (C_i^{\setminus E})^T (y_i - C_i b_r). \end{aligned} \quad (54)$$

15 Since $\text{diag} \left(\frac{\gamma}{\lambda_e} \right) + \sum_{i=1}^S C_i^T C_i \geq \text{diag} \left(\frac{\gamma}{\lambda_e} \right)$ (in a matricial
 16 positive definite sense), we obtain

$$\|b_c - b_r\|_2 \leq \sum_{i=1}^S \left\| \text{diag} \left(\frac{\lambda_e}{\gamma} \right) (C_i^{\setminus E})^T (y_i - C_i b_r) \right\|_2. \quad (55)$$

18 Rewriting $y_i - C_i b_r$ as $y_i - C_i^E b_d(S_g) + C_i^E b_d(S_g) - C_i^E b_r$
 19 and using definitions (20) and (21), it follows immediately
 20 that

$$\begin{aligned} \|b_c - b_r\|_2 & \leq \gamma_a \sum_{i=1}^S \|y_i - C_i b_d(S_g)\|_2 + \gamma_b \sum_{i=1}^S \|b_r - b_d(S_g)\|_2 \\ & \leq \gamma_a \sum_{i=1}^S \|y_i - C_i b_d(S_g)\|_2 + \gamma_b S_{\max} \|b_r - b_d(S_g)\|_2. \end{aligned} \quad (56)$$

22 Notice that γ_a is finite since for every $x \in \mathcal{X}$ it holds that

$$\left\| \text{diag} \left(\frac{\lambda_e}{\gamma} \right) C^{\setminus E}(x) \right\|_2^2 \leq \sup_{x \in \mathcal{X}, e \in \mathbb{N}_+} \phi_e(x) \cdot \sum_{e=E+1}^{+\infty} \frac{\lambda_e}{\gamma} \quad (57)$$

24 with $\sup_{x \in \mathcal{X}, e \in \mathbb{N}_+} \phi_e(x) < +\infty$ because eigenfunctions are
 25 continuous on a compact, and also with $\sum_{e=E+1}^{+\infty} \frac{\lambda_e}{\gamma} < +\infty$
 26 since K is Mercer. In the same way it is possible to show
 27 that also γ_b is finite.

28 (27) can then be proved substituting (56) in

$$\|b_c - b_d(S_g)\|_2 \leq \|b_c - b_r\|_2 + \|b_r - b_d(S_g)\|_2. \quad (58)$$

- [1] G. Pillonetto, A. Chiuso, and G. De Nicolao, "Prediction error identification of linear systems: A nonparametric gaussian regression approach," *Automatica*, vol. 47, no. 2, pp. 291 – 305, February 2011.
- [2] S. Smale and D.-X. Zhou, "Learning theory estimates via integral operators and their approximations," *Constructive approximation*, vol. 26, pp. 153–172, 2007.
- [3] G. De Nicolao and G. Ferrari-Trecate, "Consistent identification of NARX models via Regularization Networks," *IEEE Transactions on Automatic Control*, vol. 44, no. 11, pp. 2045 – 2049, November 1999.
- [4] G. Pillonetto and B. M. Bell, "Bayes and empirical Bayes semi-blind deconvolution using eigenfunctions of a prior covariance," *Automatica*, vol. 43, no. 10, pp. 1698–1712, October 2007.
- [5] S. M. S. Rezeki, W. Chan, M. R. Haskard, D. E. Mulcahy, and D. E. Davey, "Realization of self-diagnosis and self-calibration strategies using conventional signal processing and fuzzy approach for distributed intelligent sensor systems," in *SPIE conference on Smart Structures and Materials 1999: Smart Electronics and MEMS*, 1999.
- [6] M. Gopinathan, G. A. Pajunen, P. S. Neelakanta, , and M. Arockiasamy, "Linear quadratic distributed self-tuning control of vibration in a cantilever beam," in *SPIE conference on Smart Structures and Materials 1995: Smart Structures and Integrated Systems*, 1995.
- [7] A. Karnik, A. Kumar, and V. Borkar, "Distributed self-tuning of sensor networks," *Wireless Networks*, vol. 12, pp. 531 – 544, 2006.
- [8] Y. Li, J. Yu, M. Zhao, and K. Han, "Self-tuning distributed measurement fusion kalman filter," in *IEEE International Conference on Information and Automation*, 2010.
- [9] G.-L. Tao, W. Wei, and Z.-L. Deng, "The self-tuning distributed information fusion wiener filter for the ARMA signals," in *8th World Congress on Intelligent Control and Automation*, 2010.
- [10] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*. New York: Springer, 2001.
- [12] F. Pérez-Cruz and S. R. Kulkarni, "Robust and low complexity distributed kernel least squares learning in sensor networks," *IEEE Signal Processing Letters*, vol. 17, no. 4, pp. 355 – 358, April 2010.
- [13] J. B. Predd, S. R. Kulkarni, and H. V. Poor, "A collaborative training algorithm for distributed learning," *IEEE Transactions on Information Theory*, vol. 55, no. 4, pp. 1856 – 1871, April 2009.
- [14] —, "Distributed kernel regression: An algorithm for training collaboratively," in *Proceedings of the IEEE Information Theory Workshop*, March 2006, pp. 332 – 336.
- [15] P. Honeine, C. Richard, J. Bermudez, H. Snoussi, M. Essoloh, and F. Vincent, "Functional estimation in Hilbert space for distributed learning in wireless sensor networks," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, April 2009, pp. 2861–2864.
- [16] A. N. Tikhonov and V. Y. Arsenin, *Solution of Ill-posed Problems*. Wiston, 1977.
- [17] B. Schölkopf and A. J. Smola, *Learning with Kernels*. The MIT Press, 2002.
- [18] K. Yosida, *Functional Analysis*. Springer-Verlag, 1965, vol. 123.
- [19] T. Poggio and F. Girosi, "Networks for approximation and learning," *Proceedings of the IEEE*, vol. 78, no. 9, pp. 1481 – 1497, September 1990.
- [20] G. Wahba, *Spline models for observational data*. SIAM, 1990.
- [21] F. Cucker and S. Smale, "On the mathematical foundations of learning," *Bulletin of the American Mathematical Society*, vol. 39, pp. 1 – 49, 2002.
- [22] H. Zhu, C. K. I. Williams, R. Rohwer, and M. Morciniec, "Gaussian regression and optimal finite dimensional linear models," in *Neural Networks and Machine Learning*. Springer-Verlag, 1998.
- [23] W. Nef, *Linear Algebra*. McGraw-Hill, 1967.
- [24] D. Varagnolo, G. Pillonetto, and L. Schenato, "Distributed parametric and nonparametric regression with on-line performance bounds computation," *Automatica*, vol. 48, no. 10, pp. 2468 – 2481, 2012.