

# Bayesian learning of probability density functions: a Markov chain Monte Carlo approach

Simone Del Favero, Damiano Varagnolo, Gianluigi Pillonetto

**Abstract**—The paper considers the problem of reconstructing a probability density function from a finite set of samples independently drawn from it. We cast the problem in a Bayesian setting where the unknown density is modeled via a nonlinear transformation of a Bayesian prior placed on a Reproducing Kernel Hilbert Space. The learning of the unknown density function is then formulated as a minimum variance estimation problem. Since this requires the solution of analytically intractable integrals, we solve this problem by proposing a novel algorithm based on the Markov chain Monte Carlo framework. Simulations are used to corroborate the goodness of the new approach.

**Index Terms**—stochastic regularization, regularization parameter, Reproducing Kernel Hilbert Spaces, Metropolis-Hastings algorithm, stochastic processes.

## I. INTRODUCTION

In many fields ranging from basic science to engineering one is often confronted with reconstructing the stochastic mechanism generating some observational data. Examples of applications abound and we cite e.g. pattern classification, clustering, time series prediction, characterization of materials, spatial modeling [1], [2], [3], [4].

Reconstructing a probability density function is in general intricate. The problem is in fact intrinsically nonlinear, since it includes nonnegative and unitary constraints. In addition, it is subject to the so-called bias/variance dilemma [5], [6], [7]. If the hypothesis space where the unknown function is searched is too large, the estimate may turn out close to the maximum likelihood one, i.e. a sum of delta function spikes centered at the observations. This estimate is in general poor, since a priori information about the smoothness of the function is typically available. On the other hand, if the hypothesis space is too narrow, the solution could turn out too few adherent to experimental data, with still a poor predictive capability on new data.

Parametric approaches, e.g., [8], tackle these difficulties assuming finite-dimensional hypothesis spaces. This is done by imposing a known parametric form to the unknown density, e.g. a mixture of Gaussians. Regularity and nonnegativity assumptions on the unknown function can thus

be easily included in the estimation process, and the problem can be solved by just fitting parameters against data, e.g., via standard nonlinear least squares algorithms. However, the model designer has often too few information to specify so strong a priori assumptions on the density shape. This represents the major drawback of parametric techniques.

A more powerful alternative is represented by nonparametric approaches, which have a wider range of applicability since they do not require to postulate a fixed-in-advance functional form. Examples of nonparametric techniques are penalized likelihood methods [5], [9], [10], [11], ad-hoc penalized likelihood methods, smoothed histograms [12], kernel methods [13], [14], regularized Gaussian Mixtures [15] and orthogonal series estimates [16]. In particular, among the most employed approaches, we cite Parzen’s window estimator, the  $k$ -nearest neighbor approach and smoothing spline density estimation. Although applied with success in many applications, all these methods have however some limitations as how they handle the bias/variance dilemma. In fact, key parameters controlling the complexity of the hypothesis space and having a major effect on the final estimate, e.g., the kernel width in Parzen’s approach, is in practice chosen empirically. Methods used to estimate the optimal values of such parameters are often asymptotic [17], thus prone to error when dealing with small data sets, or are based on cross validation techniques [18], [9], thus possibly subject to statistical error.

*In practice, because of the nonlinearity of the problem, it is hard to define rigorous statistical criteria determining the right amount of regularization to be included in the estimation problem. In this paper we propose a statistical modeling approach to overcome this problem.*

In particular, we embed density estimation within a stochastic framework, by interpreting Tikhonov regularization as placing an opportune Bayesian prior on a Reproducing Kernel Hilbert Space (RKHS) [19], [18], [7]. We then solve the resulting Bayesian estimation problem with a novel algorithm based on the Markov chain Monte Carlo (MCMC) framework [20], [21]. In particular *we jointly learn the regularization parameter and the unknown density function determining their minimum variance estimates.* The paper is organized as follows: in Section II we provide the statement of the problem, then formulate our statistical assumptions on the unknown density function and provide a brief overview on RKHS theory in Section III. In Section IV we connect our statistical model with Tikhonov regularization and compare our approach with some other literature. In Section V, after briefly

All the authors are with the Department of Information Engineering, University of Padova, via Gradenigo 6/b, 35131 Padova, Italy. Emails: { simone.delfavero | varagnolo | giapi } @dei.unipd.it.

The research leading to these results has received funding from the European Union Seventh Framework Programme [FP7/2007-2013] under grant agreement n°257462 HYCON2 Network of excellence and n°223866 FeedNetBack, by Progetto di Ateneo CPDA090135/09 funded by the University of Padova, and by the Italian PRIN Project “New Methods and Algorithms for Identification and Adaptive Control of Technological Systems”.

1 introducing the MCMC framework, we propose our novel  
2 numerical algorithm. In Section VI we use simulations to  
3 test the relative performance of the approach. Conclusions  
4 are finally offered in Section VII.

## 5 II. STATEMENT OF THE PROBLEM

6 In the following, given a vector  $w$ , we use  $w_i$  to refer  
7 to the  $i$ -th component of  $w$ . Moreover, all vectors will be  
8 column vectors.

9 We are given  $n$  random samples  $\{y_i\}$  collected in  
10 the vector  $y$  and independently drawn from an unknown  
11 probability density function  $f(x)$ . Such density is assumed  
12 to have support on the compact set  $X \subset \mathbf{R}^d$ . Our aim is to  
13 estimate  $f$  from  $y$ .

## 14 III. STATISTICAL ASSUMPTIONS ON THE UNKNOWN 15 DENSITY

16 Before specifying our statistical assumptions on  $f$ , we first  
17 briefly sketch some properties of RKHS which are relevant  
18 in the context of this paper.

### 19 A. A brief overview on RKHS theory

20 In the sequel, let  $\mathbf{L}^2(X)$  the classical Lebesgue space of  
21 square integrable functions on  $X$ , equipped with the inner  
22 product  $\langle \cdot, \cdot \rangle_2$ , and let also  $K : X \times X \mapsto \mathbf{R}$ .

23 **Definition 1.** We say that  $K$  is definite positive if for  
24 all finite sets  $\{x_1, x_2, \dots, x_k\} \subset X$  the  $k \times k$  matrix  
25 whose  $(i, j)$ -th entry is  $K(x_i, x_j)$  is semi-definite positive.  
26 Moreover, we say that  $K$  is a Mercer kernel if it is  
27 continuous, symmetric and definite positive.

28 The following proposition can be obtained by combining  
29 the Spectral Theorem for compact operators and Mercer's  
30 theorem (see e.g. [22], [23]).

31 **Proposition 2.** Let  $K(s, t)$  a Mercer kernel. Then there exist  
32 a sequence  $\{\lambda_j \geq 0 : \lambda_{j+1} \geq \lambda_j, j = 1, \dots, \infty\}$  and a basis  
33 in  $\mathbf{L}^2(X)$  of continuous functions  $\{\phi_j : j = 1, \dots, \infty\}$  such  
34 that

$$\begin{aligned}
35 \quad \langle \phi_j, \phi_k \rangle_2 &= \begin{cases} 1 & \text{if } j = k \\ 0 & \text{otherwise} \end{cases} \\
36 \quad \int_X K(s, t) \phi_j(t) dt &= \lambda_j \phi_j(s) \\
37 \quad K(s, t) &= \sum_{j=1}^{\infty} \lambda_j \phi_j(s) \phi_j(t) \\
38
\end{aligned}$$

39 where the above convergence is uniform in  $X \times X$ .

40 The following proposition characterizes the Hilbert space  
41 associated to the Mercer Kernel  $K$  (see e.g. [24], [7]).

42 **Proposition 3.** Assigned a Mercer kernel  $K$  there exists a  
43 unique Hilbert space  $H$  such that

- 44 •  $K(x, y) \in H \quad \forall x \in X$ ;
- 45 • the span of the set  $\{K(x, \cdot), x \in X\}$  is dense in  $H$ ;
- 46 •  $f(x) = \langle f(y), K(y, x) \rangle_H \quad \forall f \in H$ .

In particular, the space  $H$  takes the following form

$$47 \quad H = \left\{ f \in \mathbf{L}^2(X) \mid f = \sum_{j=1}^{\infty} a_j \phi_j \text{ and } \sum_{j=1}^{\infty} \frac{a_j^2}{\lambda_j} < \infty \right\}$$

48 equipped with the inner product  $\langle \cdot, \cdot \rangle_H$  where, given  $f, g \in$   
49  $H$  with  $f = \sum_{j=1}^{\infty} a_j \phi_j$  and  $g = \sum_{j=1}^{\infty} b_j \phi_j$ , we have

$$50 \quad \langle f, g \rangle_H = \sum_{j=1}^{\infty} \frac{a_j b_j}{\lambda_j}.$$

51 The space  $H$  is also known in literature as the RKHS  
52 associated to the reproducing kernel  $K$ . Remarkably, the  
53 above Proposition enables us to interpret  $H$  as a certain  
54 subset of smooth functions in  $\mathbf{L}^2(X)$  determined by the  
55 eigenvalues and eigenvectors of  $K$ .  
56

57 **Example 4.** As an example of RKHS, let's define the  
58 Green's function  $G_{W_m}$  and the reproducing kernel  $K_{W_m}$  on  
59  $[0, T] \times [0, T]$ ,  $T \in \mathbf{R}$  as

$$60 \quad G_{W_m}(x, y) := \begin{cases} 0 & \text{if } x \leq y \\ 1 & \text{if } x > y \text{ and } m = 1 \\ \frac{(x-y)^{m-1}}{(m-1)!} & \text{otherwise} \end{cases}$$

$$61 \quad K_{W_m}(x, y) := \int_0^T G_{W_m}(x, \tau) G_{W_m}(y, \tau) d\tau.$$

Given a function  $f : [0, T] \mapsto \mathbf{R}$ , we use  $f^{(i)}$  to denote the  
62  $i$ -th derivative of  $f$ . The RKHS associated to  $K_{W_m}$  is then

$$63 \quad W_m = \left\{ f : [0, T] \mapsto \mathbf{R} \mid f^{(m)} \in \mathbf{L}^2[0, T], \right. \\
64 \quad \left. f^{(j)} \text{ absolutely continuous and } f^{(j)}(0) = 0 \text{ for } j = 0, \dots, m-1 \right\}$$

65 equipped with the inner product (see e.g. [18])

$$66 \quad \langle f, g \rangle_{W_m} = \left\langle f^{(m)}, g^{(m)} \right\rangle_2.$$

67 The eigenvalues and eigenvectors of  $K_{W_m}$  can be in  
68 general numerically computed (see e.g. Lemma 9 in [25]).  
69 In particular, if  $m$  equals 1,  $\phi_{W_1, j}$  and  $\lambda_{W_1, j}$  admit the  
70 following closed forms [26]:

$$\begin{aligned}
71 \quad \lambda_{W_1, j} &= T^2 / [(j-1)\pi + \pi/2]^2 \\
72 \quad \phi_{W_1, j}(t) &= \sqrt{2/T} \sin [(x/T)(j\pi - \pi/2)].
\end{aligned}$$

### 73 B. Stochastic modeling of the unknown probability density 74 function

75 We now cast our density estimation problem in a  
76 stochastic framework. We start defining a Bayesian prior  
77 for the unknown function  $f$ . In the sequel, we denote  
78 with  $\Psi$  a certain deterministic, continuous and nonlinear  
79 transformation mapping the space of continuous functions  
80 on  $X$  into itself.

81 **Assumption 5.** There exist a positive real number  $\gamma$  and  
82 numerable collections of functions  $\{\phi_i(x)\}$  on  $X$  and non-  
83 negative real numbers  $\{\lambda_i\}$  such that

- 84 •  $K(x, y) = \sum_{j=1}^{\infty} \lambda_j \phi_j(x) \phi_j(y)$  is a Mercer kernel of a  
85 RKHS  $H$ ;

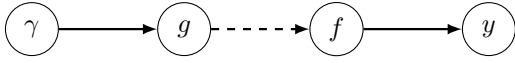


Fig. 1. Bayesian network describing the nonlinear stochastic model for density estimation.

- the function  $f$  is a random field of the form

$$f(x) = \frac{\Psi(g)}{\int_X \Psi(g) dx} \quad (1)$$

where

$$g(x) = \sum_{j=1}^{\infty} a_j \phi_j(x)$$

and  $\{a_j\}$  are independent Gaussian random variables, being the variance of  $a_j$  equal to  $\lambda_j/\gamma$ .

A graphical description of our model is depicted in Fig. 1 by using the formalism of Bayesian networks (see [27]). Here, random fields/vectors are represented by nodes, stochastic relationships by arrows and deterministic relationships by dashed arrows. One can thus note that in our framework also the regularization parameter is modeled as random variable. In particular, following a standard statistical choice (see e.g. [27]) we specify the entire network by assigning to  $\gamma$  a Gamma distribution  $\Gamma(\alpha, \beta)$  of mean  $\alpha/\beta$ .

#### IV. RELATIONSHIP WITH TIKHONOV REGULARIZATION THEORY

##### A. Connection with Tikhonov nonlinear regularization

In this sub-section we derive a connection between the probabilistic model of Fig. 1 and Tikhonov nonlinear regularization. Given a function  $g \in H$ , where  $g = \sum_{j=1}^{\infty} a_j \phi_j(x)$ , let  $a^N$  the vector containing the first  $N$  components of  $\{a_j\}$ . We define as  $g_a^N$  the following finite-dimensional approximation of  $g$ ,

$$g_a^N(x) := \sum_{j=1}^N a_j^N \phi_j(x),$$

where  $x \in X$ . We define the following prior distribution for  $a^N$ :

$$\begin{aligned} p_{a^N}(a^N) &= \frac{1}{(2\pi\gamma^{-1})^{N/2} \sqrt{\lambda_1 \cdots \lambda_N}} \exp\left(-\frac{\gamma}{2} \sum_{j=1}^N (a_j^N)^2 / \lambda_j\right) \\ &= \frac{1}{(2\pi\gamma^{-1})^{N/2} \sqrt{\lambda_1 \cdots \lambda_N}} \exp\left(-\frac{\gamma}{2} \|g_a^N(x)\|_H^2\right). \end{aligned}$$

Let  $\Psi(g_a^N)_{y_j}$  be the function  $\Psi(g_a^N)$  evaluated at  $y_j$ . Then the conditional density for  $y$  given  $a^N$  and  $\gamma$  is

$$p_{y|a^N, \gamma}(y|a^N, \gamma) = \frac{\prod_{j=1}^n \Psi(g_a^N)_{y_j}}{\int_X \Psi(g_a^N) dx}.$$

The corresponding negative log of the likelihood for a given  $y \in \mathbf{R}^n$  and  $a^N$  is

$$\begin{aligned} l^N(y, a^N | \gamma) &= \frac{1}{2} \sum_{j=1}^n \log\left(\frac{2\pi\lambda_j}{\gamma}\right) - \sum_{j=1}^n \log(\Psi(g_a^N)_{y_j}) \\ &\quad + \log\left(\int_X \Psi(g_a^N) dx\right) + \gamma \frac{\|g_a^N\|_H^2}{2}. \end{aligned}$$

We point out that  $H$ , being a RKHS, is a subset of the space of continuous functions and convergence in the topology induced by  $\|\cdot\|_H$  implies uniform convergence, see e.g. [7]. Then we easily have that

$$l^N(y, g_a^N | \gamma) - \frac{1}{2} \sum_{j=1}^n \log\left(\frac{2\pi\lambda_j}{\gamma}\right) \xrightarrow{N \rightarrow \infty} l(y, g | \gamma)$$

where

$$\begin{aligned} l(y, g | \gamma) &:= - \sum_{j=1}^n \log(\Psi(g)_{y_j}) \\ &\quad + \gamma \frac{\|g\|_H^2}{2} + \log\left(\int_X \Psi(g) dx\right). \end{aligned}$$

Given the model of Fig. 1, we can interpret

$$\hat{g} = \arg \min_{g \in H} l(y, g | \gamma) \quad (2)$$

as the maximum a posteriori (MAP) estimate of  $g$  given  $y$  and  $\gamma$ . MAP estimate is thus provided by a Tikhonov nonlinear variational problem which contains two contrasting terms. The first one, equal to  $-\sum_{j=1}^n \log(\Psi(g)_{y_j})$ , takes into account the experimental evidence, the second one, equal to  $\|g\|_H^2$ , the a priori information about the regularity of the solution. The trade-off between these two components is then established by the regularization parameter  $\gamma$ . Finally, a third term is also present in the estimator, equal to  $\log(\int_X \Psi(g) dx)$ , which is to enforce the nonnegative and unitary constraint on the unknown probability density function.

##### B. Connection with other density estimation approaches

We notice that Estimator (2) has already been proposed in the literature with different choices of  $\Psi$ . For instance, [28] introduces a penalty on the second derivative of the squared root of the function of interest. Under our framework, this corresponds to define  $\Psi : f \mapsto f^{(2)}$ , embedding the problem in  $W_2$ . [9] instead assumes  $\Psi$  to be exponential, i.e.  $\Psi : f \mapsto e^f$ , and proposes efficient iterative algorithms to implement this model (asymptotic properties are described in [29]).

However, rigorous statistical criteria to determine  $\gamma$  in (2) have not been so far proposed. In addition it is worth pointing out that, even if  $\gamma$  were known, the MAP estimate of  $g$  given  $y$  and  $\gamma$  is in general less robust than its a posteriori expected value. In the next Section we then show how to compute the minimum variance estimates of  $f$  and  $\gamma$  via a MCMC approach.

We finally remark that density estimation through infinite Gaussian mixtures [30] correspond to MCMC approaches using more vague priors: in (2) designers can directly

1 encode assumptions on the regularity of the density through  
 2 appropriate  $K(\cdot, \cdot)$ 's, while in infinite Gaussian mixtures this  
 3 design opportunity is missing.

#### 4 V. NUMERICAL ALGORITHMS

5 Recovering minimum variance estimates and confidence  
 6 intervals from the a posteriori probability density function of  
 7  $f$  and  $\gamma$  graphically described in Fig. 1 requires the solution  
 8 of analytically intractable integrals. Here, we derive a novel  
 9 MCMC algorithm for density estimation which circumvents  
 10 this difficulty. We start providing a brief overview on the  
 11 Metropolis-Hastings algorithm [20] on which our numerical  
 12 procedure relies on.

##### 13 A. A brief overview on the MCMC framework

The goal of a MCMC algorithm is to simulate realizations  
 from a certain posterior distribution so that empirical  
 estimates for any statistics of interest can be determined.  
 A MCMC procedure thus consists of two steps. Firstly,  
 a Markov process with limiting quantities that follow the  
 invariant distribution of interest, in our case the a posteriori  
 probability density function of  $f$  and  $\gamma$ , is designed. This  
 first step is used to recover the target distribution of interest  
 in sampled form. Secondly, a Monte Carlo integration  
 is done to obtain the integrals of interest. The common  
 mechanism by which the first step can be performed is the  
 Metropolis/Hastings algorithm. A variant of this procedure,  
 named the single-component Metropolis/Hastings algorithm,  
 will be in particular used in this paper. To describe it,  
 let  $\pi(\theta)$  the target density, being  $\theta$  a finite-dimensional  
 vector containing the parameters of interest. We denote with  
 $q_i(Z_{t+1}|Z_t)$  a proposal density from which a candidate  
 value  $Z_{t+1}$  is drawn when the current state of the chain  
 is  $Z_t$ . The scheme suggests to divide  $Z$  into  $h$  portions  
 of desired dimension, then specifying  $h$  proposal density  
 functions  $q_i(\cdot|\cdot)$  with  $i = 1, 2, \dots, h$ . Every iteration of the  
 algorithm is composed by  $h$  distinct phases where Metropolis  
 Hastings updates are employed to explore the parameter  
 space by proposing moves which are subsequently either  
 accepted or rejected. To be specific, at step  $i$  of iteration  
 $t+1$  the  $i$ -th portion of  $Z$ , denoted with  $Z_{t+1,i}$ , is drawn by  
 the kernel  $q_i(Z_{t+1,i}|Z_{t,i}, Z_{t,-i})$ , where

$$Z_{t,-i} := \{Z_{t+1,1}, \dots, Z_{t+1,i-1}, Z_{t,i+1}, \dots, Z_{t,h}\}.$$

In practice, the first  $i-1$  components of  $Z_{t,-i}$  come from  
 the first  $i-1$  steps computed at instant  $t+1$ . The candidate  
 is then accepted with probability

$$\delta(Z_{t,-i}, Z_{t,i}, Z_{t+1,i}) = \min \left( 1, \frac{\pi(Z_{t+1,i}|Z_{t,-i})q_i(Z_{t,i}|Z_{t+1,i}, Z_{t,-i})}{\pi(Z_{t,i}|Z_{t,-i})q_i(Z_{t+1,i}|Z_{t,i}, Z_{t,-i})} \right).$$

14 If the proposal is rejected then the chain remains in  
 15 the current state. This scheme guarantees, under mild  
 16 additional conditions,  $\pi$  to be the limiting distribution of the  
 17 Markov chain generated (see e.g. [21]). This virtually holds  
 18 independently of the particular proposal densities  $\{q_i(\cdot|\cdot)\}$   
 19 employed. Even if their choice is essentially arbitrary, it has

20 however a crucial influence on the rate of convergence of  
 21 the algorithm. In other words, it can be often problematic to  
 22 design an efficient MCMC scheme which obtains an accurate  
 23 reconstruction (in sampled form) of  $\pi$  after a reasonable  
 24 number of iterations.

##### B. MCMC algorithm for density estimation

25 Following [31], we define a MCMC procedure which  
 26 relies upon a representation of  $g$  in terms of a finite subset  
 27 of eigenvectors  $\{\phi_j\}$ . To be specific, we assume  
 28

$$g(x) = \sum_{j=1}^N a_j^N \phi_j(x) \quad (3) \quad 29$$

30 where  $N$  depends on the specific problem and has to  
 31 be chosen large enough so as to provide an accurate  
 32 approximation of the original infinite-dimensional model.

We then block the parameter space into two groups, i.e.  
 $\gamma$  and  $a^N$ . As concerns the updating of  $\gamma$ , let  $\Lambda$  the  $N \times N$   
 diagonal matrix with  $(j, j)$ -th entry equal to  $\lambda_j$ . After some  
 computations we obtain

$$p_{\gamma|a^N, y}(\gamma | a^N, y) = \Gamma \left( \frac{N}{2} + \alpha, \frac{1}{2} (a^N)' \Lambda^{-1} (a^N) + \beta \right).$$

33 We then choose as the proposal density for  $\gamma$  its conditional  
 34 distribution given  $a^N$  and  $y$ , thus defining a Gibbs sampler  
 35 update, see [21].

36 As concerns  $a^N$ , the same kind of strategy can no more  
 37 be exploited, since the corresponding posterior does not take  
 38 a standard form. We then propose a move in which these  
 39 parameters are updated by sampling a value in a symmetric  
 40 interval around the current position, in accordance with a  
 41 Gaussian distribution with covariance matrix  $\Sigma$ . In particular,  
 42 we adapt the proposal scales as follows. We begin a pilot-  
 43 tuning run from some arbitrary values  $\gamma_0$  and  $a_0^N$ . We  
 44 preliminarily set  $\Sigma$  to a matrix proportional to  $\gamma_0^{-1} \Lambda$ , with the  
 45 scale factor chosen in order to make the acceptance ratio for  
 46  $a^N$  to be around 0.2–0.4. This preliminary stage is used in  
 47 order to let the algorithm approximately learn the a posteriori  
 48 correlation of the components of  $a^N$  given  $y$ . Then, after  
 49 a certain number of iterations, we set once and for all  $\Sigma$   
 50 as proportional to the covariance matrix of the generated  
 51 samples of  $a^N$ , still choosing a scale factor which ensures  
 52 an acceptance ratio of the proposed moves around 0.3 [21].

53 We summarize the MCMC procedure by means of the  
 54 following Algorithm 1. The symbol  $\mathcal{N}(\mu, \Sigma)$  is therein used  
 55 to denote a Gaussian probability density function having  
 56 mean  $\mu$  and covariance matrix  $\Sigma$ .

57 **Remark 6.**  $\Psi$ , as well as the RKHS  $H$ , should be chosen  
 58 in accordance with available information on  $f$ . Moreover,  
 59 as concerns the transformation  $\Psi(g)/\int \Psi(g)dx$ , it appears  
 60 important to define it as injective. Otherwise, the posterior  
 61 probability of  $a^N$  given  $y$  could turn out multimodal with  
 62 many peaks (or also improper), thus making difficult the  
 63 convergence of the generated Markov chain.

---

**Algorithm 1**


---

1: (*initialization*) set  $(\gamma_0, a_0^N)$  and  $k = 1$

---

2: **for**  $k = 1, 2, \dots$  **do**

3:   sample  $\gamma_k$  from

$$\Gamma\left(\frac{N}{2} + \alpha, \frac{1}{2}(a_{k-1}^N)' \Lambda^{-1}(a_{k-1}^N) + \beta\right)$$

4:   sample  $s$  from  $\mathcal{N}(a_{k-1}^N, \Sigma)$

5:   accept  $s$  with probability

$$\delta(s, a_{k-1}^N, \gamma_k) = \min\left(1, \frac{\rho(s, \gamma_k)}{\rho(a_{k-1}^N, \gamma_k)}\right)$$

where

$$\rho(a^N, \gamma) := \frac{\prod_{i=1}^n \Psi(g_a^N)_{y_i}}{\int_X \Psi(g_a^N) dx} \exp\left(-\frac{\gamma}{2}(a^N)' \Lambda^{-1}(a^N)\right)$$

6:   if  $s$  is accepted then set  $a_k^N = s$ , otherwise  $a_k^N = a_{k-1}^N$

---

## VI. NUMERICAL EXPERIMENTS

We present two set of univariate simulations on  $[0, 1]$  to examine the relative effectiveness of the proposed methodology. Before introducing the examples in detail, it is worth pointing out that in every case study the binary control of Raftery and Lewis has been used in order to assess the convergence of the generated Markov chains, see [32]. In particular we have always required to estimate the quantiles 0.025, 0.25, 0.5, 0.75, 0.975 with precision respectively 0.005, 0.01, 0.01, 0.01, 0.005, and with probability 0.95.

### A. Bayesian learning of an exponential probability density function

We start considering the reconstruction of an exponential density, with mean 0.15, from 100 samples independently drawn from it. Data are displayed by means of a histogram in Fig. 2 (top panel).

In (1) we include information about the smoothness of the unknown function  $f$  by setting  $\Psi$  to an exponential transform and assuming that  $g$  belongs to  $W_1$ , with support on<sup>1</sup>  $[0, 1]$ . We recall that this exponential transformation is the nonparametric counterpart of the classic log-normal choice in parametric frameworks.

As concerns  $\gamma$ ,  $\alpha$  and  $\beta$  are chosen so as to define a poorly informative prior on this hyper-parameter. After setting  $N$  to 100 in (3), the model in Fig. 1 has been solved by reconstructing the joint posterior density of  $\gamma$  and  $f$  via a Monte Carlo run where 4000 samples were generated. In Fig. 3 (top panel) the (unnormalized) posterior of  $\gamma$  is depicted. In Fig. 4 (top panel) we report the minimum variance estimate of  $f$  (solid line), together with its 95% confidence interval (shaded area) and the true function

<sup>1</sup>It is easy to assess that this kind of choice makes  $\Psi(g)/\int \Psi(g)dx$  injective. This in particular holds thanks to the side condition at 0 present in the definition of  $W_1$ .

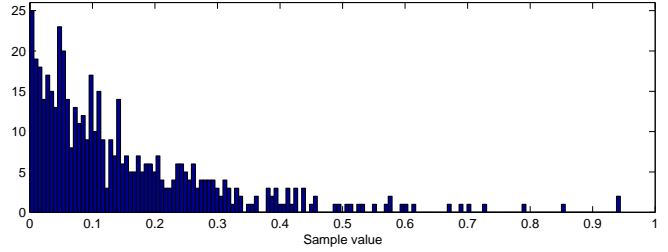
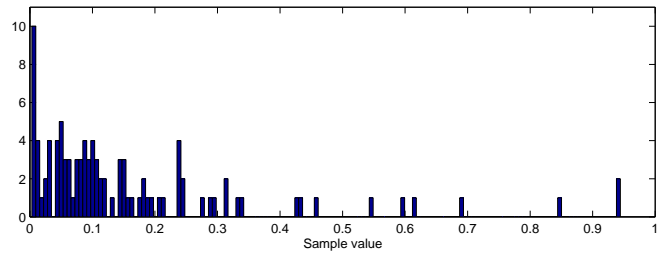


Fig. 2. Reconstruction of an exponential density function. *Top*: histogram of 100 samples independently drawn from the unknown density. *Bottom*: histogram of 500 samples independently drawn from the unknown density.

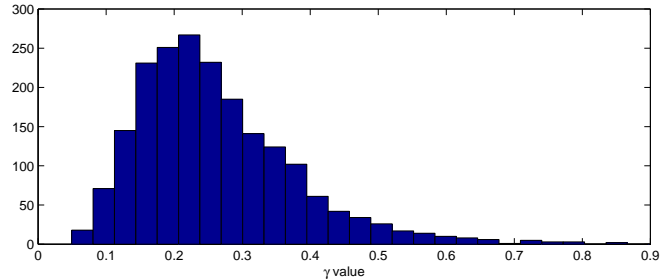
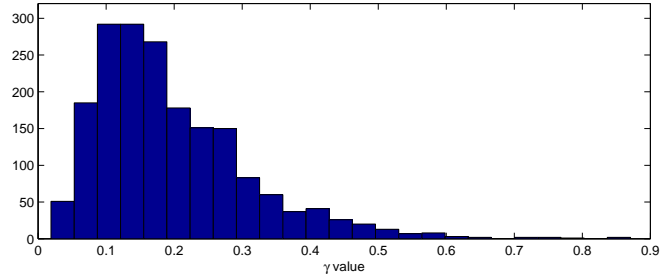


Fig. 3. Reconstruction of an exponential density function. *Top*: posterior of  $\gamma$  obtained (in sampled form) by MCMC using the training set of 100 samples. *Bottom*: posterior of  $\gamma$  obtained (in sampled form) by MCMC using the training set of 500 samples.

(dashed line). Even though the training set size is small, the density estimate is somewhat close to the true one. It thus appears that a suitable amount of regularization has been introduced by the nonlinear estimator.

We then repeated the entire estimation process by adding to the training set other 400 samples independently drawn from the exponential density. Data are displayed by means of a histogram in Fig. 2 (bottom panel). After setting  $N$  to 100 in (3), a Monte Carlo run of 3000 samples were generated. The posterior of  $\gamma$  (in sampled form) is visible in Fig. 3 (bottom panel). Results regarding the estimate of  $f$  are then

32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42

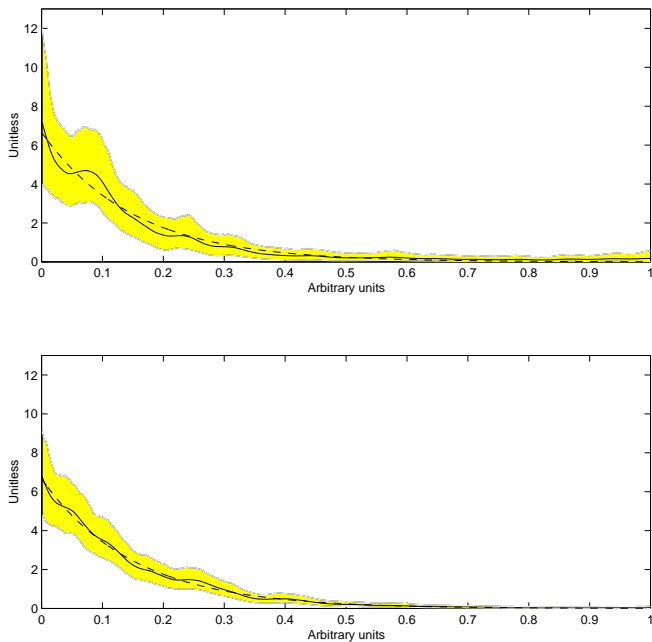


Fig. 4. Reconstruction of an exponential density function. *Top*: minimum variance estimate (continuous line) with 95% confidence interval (shaded area) and true density (dashed line) using the training set of 100 samples. *Bottom*: minimum variance estimate (continuous line) with 95% confidence interval (shaded area) and true density (dashed line) using the training set of 500 samples.

1 reported in Fig. 4 (bottom panel), with the same rationale  
 2 followed in the top panel of the same figure. One can note  
 3 that the minimum variance estimate is very close to the truth.  
 4 Moreover, comparing the confidence intervals depicted in the  
 5 first and top panel, one can appreciate how the uncertainty  
 6 related to the estimate has been reduced by augmenting the  
 7 size of the data set. This illustrates the capability of the  
 8 proposed approach in clearly making the investigator assess  
 9 the amount of information that different training sets provide.

#### 10 B. Bayesian learning of a mixture of Gaussians

As a second example, we consider a benchmark problem  
 proposed in [9] which consists of reconstructing a density  
 on  $[0, 1]$  proportional to

$$\frac{1}{3}e^{-50(x-0.3)^2} + \frac{2}{3}e^{-50(x-0.7)^2}.$$

11 This function is depicted in Fig. 5 (dashed line), and it  
 12 can be noticed that it virtually corresponds to a mixture of  
 13 Gaussians.

14 In (1) we include information about the smoothness of  
 15 the unknown function  $f$  by setting  $\Psi$  as in Sec. VI-A and  
 16 assuming that  $g$  belongs to  $W_2$ , with support on  $[0, 1]$ .  
 17 As concerns  $\alpha$  and  $\beta$ , they are chosen so as to define a  
 18 poorly informative prior on  $\gamma$ . After setting  $N$  to 100 in  
 19 (3), we consider 300 replicates of this problem. For each  
 20 of these 300 simulations we generate a new training set of  
 21 200 samples and then use our MCMC scheme to obtain the  
 22 minimum variance estimate of  $f$ , in accordance with the  
 23 model depicted in Fig. 1. In Fig. 5 we report the mean of

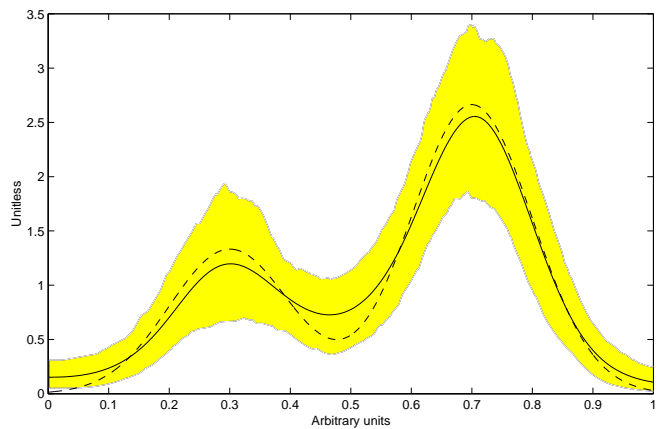


Fig. 5. Reconstruction of a mixture of Gaussians through Monte Carlo simulations: estimates mean (continuous line), with 95% bands of variability (shaded area), and true density (dashed line).

the 300 estimates (continuous line) together with the 95%  
 variability band (shaded area), i.e. the interval between the  
 2.5 and 97.5 percentiles of the Monte Carlo distribution of  
 the density estimates at each point of their support. Even  
 though the size of the data set is really small, one can  
 note that the mean is quite close to the true profile. This  
 is particularly evident near the second peak since, on average,  
 the most part of the samples is generated from it at every  
 of the 300 Monte Carlo runs. Moreover, 95% variability  
 bands show that the variance of the error is not high.

## VII. CONCLUSIONS

The choice of the regularization parameter is a crucial  
 issue in learning theory and, more in general, when dealing  
 with ill-posed problems [33], [18], [34], [7]. In particular,  
 when reconstructing a probability density function, its  
 determination presents formidable difficulties due to the  
 nonlinear nature of the problem. In this paper, we have  
 proposed a new technique which tackles this difficulty by  
 casting the density estimation problem within a Bayesian  
 framework. An MCMC approach is then used to implement  
 the resulting stochastic model.

The power of our method consists in providing estimates  
 that take into account all the sources of uncertainty present  
 in the problem. In particular, the proposed algorithm is able  
 to return minimum variance estimates of both the unknown  
 density and the regularization parameter. In addition, our  
 scheme can associate to the estimate a confidence interval,  
 thus allowing the investigator to assess how informative the  
 training set is.

In future work, we will test the methodology to reconstruct  
 multivariate densities, also improving its computational  
 efficiency.

## REFERENCES

- [1] A. Gersho and R. M. Gray, *Vector Quantization and Signal Compression*. Boston: Kluwer Academic Publishers, 1992.
- [2] A. S. Weigend and A. N. Srivastava, "Predicting conditional probability distributions: a connectionist approach," *International Journal of Neural Systems*, vol. 6, no. 2, pp. 109 – 118, June 1995.

- 2 [3] S. Fiori, "Nonsymmetric pdf estimation by artificial neurons: Application to statistical characterization of reinforced composites," *IEEE Transactions on Neural Networks*, vol. 14, no. 4, pp. 959 – 962, July 2003.
- 3
- 4 [4] F. Liang, "Continuous contour monte carlo for marginal density estimation with an application to spatial statistical model," *Journal of Computational and Graphical Statistics*, vol. 16, no. 3, pp. 608 – 632, September 2007.
- 5
- 6 [5] B. Silverman, "On the estimation of a probability density function by the maximum penalized likelihood method," *Annals of Statistics*, vol. 10, pp. 795 – 810, 1982.
- 7
- 8 [6] V. N. Vapnik, *Statistical learning theory*. New York: Wiley, 1998.
- 9
- 10 [7] F. Cucker and S. Smale, "On the mathematical foundations of learning," *Bulletin of the American mathematical society*, vol. 39, pp. 1–49, 2001.
- 11
- 12 [8] V. Krylov, G. Moser, S. B. Serpico, and J. Zerubia, "On the method of logarithmic cumulants for parametric probability density function estimation," INRIA Sophia Antipolis, Tech. Rep. RR-7666, 2011.
- 13
- 14 [9] C. Gu, "Smoothing spline density estimation: A dimensionless automatic algorithm," *Journal of the American Statistical Association*, vol. 88, pp. 495 – 504, June 1993.
- 15
- 16 [10] A. Komárek and E. Lesaffre, "Generalized linear mixed model with a penalized gaussian mixture as a random effects distribution," *Computational Statistics & Data Analysis*, vol. 52, no. 7, pp. 3441 – 3458, March 2008.
- 17
- 18 [11] C. Schellhase and G. Kauermann, "Density estimation and comparison with a penalized mixture approach," *Computational Statistics*, vol. –, pp. 1 – 21, 2011.
- 19
- 20 [12] G. Wahba, "Histosplines with knots which are order statistics," *Journal of the Royal Statistical Association*, vol. 38, pp. 140 – 151, 1976.
- 21
- 22 [13] E. Parzen, "On the estimation of a probability density function and mode," *The Annals of Mathematical Statistics*, vol. 33, pp. 1065–1076, 1962.
- 23
- 24 [14] E. López-Rubio and J. M. O. de Lazcano-Lobato, "Soft clustering for nonparametric probability density function estimation," *Pattern Recognition Letters*, vol. 29, pp. 2085 – 2091, 2008.
- 25
- 26 [15] C. Archambeau and M. Verleysen, "Fully nonparametric probability density function estimation with finite gaussian mixture models," in *International Conference on Advances in Pattern Recognition*, 2003.
- 27
- 28 [16] M. Girolami, "Orthogonal series density estimation and the kernel eigenvalue problem," *Neural Computation*, vol. 14, pp. 669 – 688, 2002.
- 29
- 30 [17] Y. Yang, "Penalized semiparametric density estimation," *Statistics and Computing*, vol. 19, no. 4, pp. 355 – 366, 2009.
- 31
- 32 [18] G. Wahba, *Spline models for observational data*. SIAM, Philadelphia, 1990.
- 33
- 34 [19] A. N. Tikhonov and V. Y. Arsenin, *Solutions of Ill-Posed Problems*. Washington, D.C.: Winston/Wiley, 1977.
- 35
- 36 [20] W. Hastings, "Monte Carlo sampling methods using Markov chain and their applications," *Biometrika*, vol. 57, pp. 97–109, 1970.
- 37
- 38 [21] W. R. Gilks, S. Richardson, and D. J. Spiegelhalter, *Markov chain Monte Carlo in Practice*. London: Chapman and Hall, 1996.
- 39
- 40 [22] J. Mercer, "Functions of positive and negative type and their connection with the theory of integral equations," *Phil. Trans. Roy. Soc. London Ser.*, vol. 209, pp. 415–446, 1909.
- 41
- 42 [23] R. Courant and D. Hilbert, *Methods of mathematical physics*. Interscience, 1962.
- 43
- 44 [24] N. Aronszajn, "Theory of reproducing kernels," *Transactions of the American Mathematical Society*, vol. 68, pp. 337 – 404, 1950.
- 45
- 46 [25] G. De Nicolao and G. Ferrari-Trecate, "Consistent identification of NARX models via regularization networks," *IEEE Transactions on Automatic Control*, vol. 44, pp. 2045–2049, 1999.
- 47
- 48 [26] A. M. Yaglom, *Correlation theory of stationary and related random functions*. Springer-Verlag, New York, 1987, vol. 1.
- 49
- 50 [27] P. Magni, R. Bellazzi, and G. De Nicolao, "Bayesian function learning using MCMC methods," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 20, pp. 1219–1331, 1998.
- 51
- 52 [28] L. Liu, M. Levine, and Y. Zhu, "A functional em algorithm for mixing density estimation via nonparametric penalized likelihood maximization," *Journal of Computational and Graphical Statistics*, vol. 18, no. 2, pp. 481 – 504, 2009.
- 53
- 54 [29] C. Gu and C. Qiu, "Smoothing spline density estimation: theory," *Annals of Statistics*, vol. 21, pp. 227–234, 1993.
- 55
- 56 [30] C. E. Rasmussen, "The infinite Gaussian mixture model." in *Advances in Neural Information Processing Systems (NIPS)*, vol. 12, 2000, pp. 554–560.
- 57
- 58 [31] G. Pillonetto and B. Bell., "Bayes and empirical Bayes semi-blind deconvolution using eigenfunctions of a prior covariance," *Automatica*, vol. 43, no. 10, pp. 1698–1712, 2007.
- 59
- 60 [32] A. E. Raftery and S. M. Lewis, *Implementing MCMC*. Markov Chain Monte Carlo in Practice. W.R. Gilks, S. Richardson, and D.J. Spiegelhalter, eds. London: Chapman and Hall, 1996, pp. 115–130.
- 61
- 62 [33] J. A. Rice, "Choice of smoothing parameter in deconvolution problems," *Contemporary Mathematics*, vol. 59, pp. 137–151, 1986.
- 63
- 64 [34] P. Hansen, "Analysis of discrete ill-posed problems by means of the L-curve," *SIAM Review*, vol. 34, pp. 561 – 580, 1992.
- 65
- 66
- 67
- 68
- 69
- 70
- 71
- 72
- 73
- 74
- 75