# Gaussian Processes for Nonlinear Signal Processing

Fernando Pérez-Cruz, Steven Van Vaerenbergh, Juan José Murillo-Fuentes, Miguel Lázaro-Gredilla and
Ignacio Santamaría

*Abstract*—**Gaussian processes (GPs) are versatile tools that have been successfully employed to solve nonlinear estimation problems in machine learning, but that are rarely used in signal processing. In this tutorial, we present GPs for regression as a natural nonlinear extension to optimal Wiener filtering. After establishing their basic formulation, we discuss several important aspects and extensions, including recursive and adaptive algorithms for dealing with non-stationarity, low-complexity solutions, non-Gaussian noise models and classification scenarios. Furthermore, we provide a selection of relevant applications to wireless digital communications.**

## I. INTRODUCTION

Gaussian processes (GPs) are Bayesian state-of-the-art tools for discriminative machine learning, i.e., regression [1], classification [2] and dimensionality reduction [3]. GPs were first proposed in statistics by Tony O'Hagan [4] and they are well-known to the geostatistics community as kriging. However, due to their high computational complexity they did not become widely applied tools in machine learning until the early XXI century [5]. GPs can be interpreted as a family of kernel methods with the additional advantage of providing a full conditional statistical description for the predicted variable, which can be primarily used to establish confidence intervals and to set hyper-parameters. In a nutshell, Gaussian processes assume that a Gaussian process prior governs the set of possible latent functions (which are unobserved), and the likelihood (of the latent function) and observations shape this prior to produce posterior probabilistic estimates. Consequently, the joint distribution of training and test data is a multidimensional Gaussian and the predicted

distribution is estimated by conditioning on the training data.

While GPs are well-established tools in machine learning, they are not as widely used by the signal processing community as neural networks or support vector machines (SVMs) are. In our opinion, there are several explanations for the limited use of GPs in signal processing problems. First, they do not have a simple intuition for classification problems. Second, their direct implementation is computationally demanding. Third, their plain vanilla version might seem uptight and not flexible enough. Fourth, to most signal processing experts Gaussian process merely stands for a noise model and not for a flexible algorithm that they should be using.

In this paper, we present an overview on Gaussian processes explained for and by signal processing practitioners. We introduce GPs as the natural nonlinear Bayesian extension to the linear minimum mean square error (MMSE) and Wiener filtering, which are central to many signal processing algorithms and applications. We believe that GPs provide the correct approach to solve an MMSE filter nonlinearly, because they naturally extend least squares to nonlinear solutions through the kernel trick; they use a simple yet flexible prior to control the nonlinearity; and, evidence sampling or maximization allows setting the hyper-parameters without overfitting. This last feature is most interesting: by avoiding cross-validation we are able to optimize over a larger number of hyperparameters, thus increasing the available kernel expressiveness. Additionally, GP provides a full statistical description of its predictions.

The tutorial is divided in three parts. We have summarized in Figure 1 the relationship between the regression techniques introduced throughout the different sections. In the first part, Section II provides a detailed overview of Gaussian processes for regression (GPR) [1]. We show that they are the natural nonlinear extension to MMSE/Wiener filtering and how they can be solved recursively. The second part of the paper focuses briefly on several key aspects of GP-based techniques. Consecutively, we review solutions to adjust the kernel function (Section III ), to tame the computational complexity of GPs (Section IV ), and to deal with non-Gaussian noise models (Section V ). In the third part, we cover

Fernando Pérez-Cruz and Miguel Lázaro-Gredilla are with the Dept. of Signal Theory and Communications, University Carlos III in Madrid, Spain, Email: {fernando,miguel}@tsc.uc3m.es.

Steven Van Vaerenbergh and Ignacio Santamaría are with the Dept. of Communications Engineering, University of Cantabria, Spain, Email: {steven,nacho}@gtas.dicom.unican.es.

Juan José Murillo-Fuentes is with the Dept. of Signal Theory and Communications, Spain, Email: murillo@etsi.us.es.
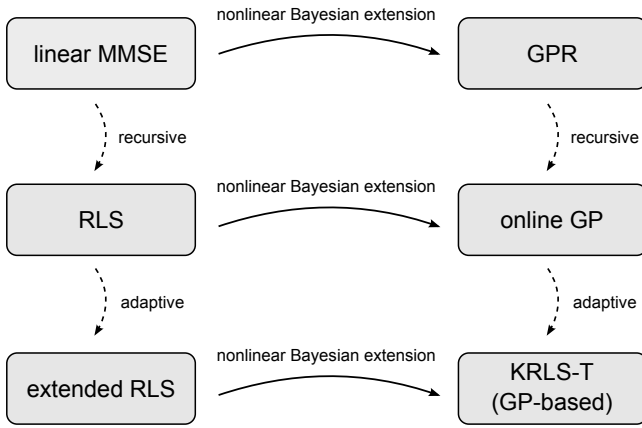
Fig. 1: Relationship between the regression techniques discussed in this tutorial.

additional extensions of interest to signal processing practitioners, in particular dealing with non-stationary scenarios (Section VI ) and classification problems (Section VII ). We illustrate them with relevant examples in signal processing for wireless communications. We conclude the paper with a discussion.

## II. Gaussian Processes for Machine Learning

### A. Minimum mean square error: a starting point

GPs can be introduced in a number of ways and we, as signal processing practitioners, find it particularly appealing to start from the MMSE solution. This is because the Wiener solution, which is obtained by minimizing the MSE criterion, is our first approach to most estimation problems and, as we show, GPs are its natural Bayesian extension.

Many signal processing problems reduce to estimating from an observed random process $\mathbf{x} \in \mathbb{R}^p$ another related process $y \in \mathbb{R}$. These two processes are related by a probabilistic, possibly unknown, model $p(\mathbf{x}|y)$. It is well known that the unconstrained MMSE estimate,

$$\underset{f(\mathbf{x})}{\operatorname{argmin}} E\left[\|y - f(\mathbf{x})\|^2\right], \qquad (1)$$

coincides with the conditional mean estimate of $y$ given $\mathbf{x}$

$$f_{mmse}(\mathbf{x}) = E[y|\mathbf{x}] = \int y p(y|\mathbf{x}) dy = \int y \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} dy. \qquad (2)$$

If $p(y, \mathbf{x})$ is jointly Gaussian, i.e. $p(y)$ and $p(\mathbf{x}|y)$ are Gaussians and $E[\mathbf{x}|y]$ is linear in $y$, this solution is linear. If $y$ and $\mathbf{x}$ are zero mean, the solution yields $E[y|\mathbf{x}] = \mathbf{w}^\top \mathbf{x}$, where

$$\mathbf{w}_{mmse} = \underset{\mathbf{w}}{\operatorname{argmin}} E\left[\left(y - \mathbf{w}^\top \mathbf{x}\right)^2\right] = \left(E\left[\mathbf{x}\mathbf{x}^\top\right]\right)^{-1} E\left[\mathbf{x}y\right]. \qquad (3)$$

Furthermore, these expectations can be easily estimated, using the sample mean, from independently and identically distributed (iid) samples drawn from $p(\mathbf{x}|y)$ and $p(y)$, namely $\mathcal{D}_n = \{\mathbf{x}_i, y_i\}_{i=1}^n$.

However, if $\mathbf{x}$ is not linearly related to $y$ (plus Gaussian noise) or $y$ is not Gaussian distributed, the conditional estimate of $y$ given $\mathbf{x}$ is no longer linear. Computing the nonlinear conditional mean estimate in (2) directly from $\mathcal{D}_n$ either leads to overfitted solutions, because there are no convergence guarantees for general density estimation [6], or to suboptimal solutions, if we restrict the density model to come from a narrow class of distributions. For instance, in channel equalization, although suboptimal, the sampled version of (3) is used due to its simplicity. One viable solution would be to minimize the sampled version of (1) with a restricted family of approximating functions to avoid overfitting. Kernel least squares (KLS) [7] and Gaussian process regression, among others, follow such approach.

### B. Gaussian Processes for Regression

In its simplest form, GPR models the output nonlinearly according to

$$y = f(\mathbf{x}) + \nu, \qquad (4)$$

and it follows (1), without assuming that $\mathbf{x}$ and $y$ are linearly related or that $p(y)$ is Gaussian distributed. Nevertheless, it still considers that $p(y|\mathbf{x})$ is Gaussian distributed, i.e., $\nu$ is a zero-mean Gaussian[1]. In this way, GP can be understood as a natural nonlinear extension to MMSE estimation. Additionally, GPR does not only estimate (2) from $\mathcal{D}_n$, but it also provides a full statistical description of $y$ given $\mathbf{x}$, namely

$$p(y|\mathbf{x}, \mathcal{D}_n). \qquad (5)$$

GPs can be presented as a nonlinear regressor that expresses the input-output relation in (4) by assuming that a real-valued function $f(\mathbf{x})$, known as *latent function*, underlies the regression problem and that this function follows a Gaussian process. Before the labels are revealed, we assume this latent function has been drawn from a Gaussian process prior. GPs are characterized by their mean and covariance functions, denoted by $\mu(\mathbf{x})$ and $k(\mathbf{x}, \mathbf{x}')$, respectively. Even though nonzero mean priors might be of use, working with zero-mean priors typically represents a reasonable assumption and it simplifies the notation. The covariance function explains the correlation between each pair of points in the input space and characterizes the functions that

---

[1] A further relaxation to this condition is discussed in Section V.

can be described by the Gaussian process. For example, $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$ only yields linear latent functions and it is used to solve Bayesian linear regression problems, for which the mean of the posterior process coincides with the MMSE solution in (3), as shown in Section II-E. We cover the design of covariance functions in Section III.

For any finite set of inputs $\mathcal{D}_n$, a Gaussian process becomes a multidimensional Gaussian defined by its mean (zero in our case) and covariance matrix, $(\mathbf{K}_n)_{ij} = k(\mathbf{x}_i, \mathbf{x}_j), \forall \mathbf{x}_i, \mathbf{x}_j \in \mathcal{D}_n$. The Gaussian process prior becomes

$$p(\mathbf{f}_n | \mathbf{X}_n) = \mathcal{N}(\mathbf{0}, \mathbf{K}_n), \tag{6}$$

where $\mathbf{f}_n = [f(\mathbf{x}_1), f(\mathbf{x}_2), \ldots, f(\mathbf{x}_n)]^\top$ and $\mathbf{X}_n = [\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n]$. We want to compute the estimate for a general input $\mathbf{x}$, when the labels for the $n$ training examples, denoted by $\mathbf{y}_n = [y_1, y_2, \ldots, y_n]^\top$, are known. We can analytically compute (5) by using the standard tools of Bayesian statistics: Bayes' rule, marginalization and conditioning.

We first apply Bayes' rule to obtain the posterior density for the latent function

$$p(f(\mathbf{x}), \mathbf{f}_n | \mathbf{x}, \mathcal{D}_n) = \frac{p(\mathbf{y}_n | \mathbf{f}_n) p(f(\mathbf{x}), \mathbf{f}_n | \mathbf{x}, \mathbf{X}_n)}{p(\mathbf{y}_n | \mathbf{X}_n)}, \tag{7}$$

where $p(f(\mathbf{x}), \mathbf{f}_n | \mathbf{x}, \mathbf{X}_n)$ is the Gaussian process prior in (6) extended with a general input $\mathbf{x}$, $p(\mathbf{y}_n | \mathbf{f}_n)$ is the likelihood for the latent function at the training set, in which $\mathbf{y}_n$ is independent of $\mathbf{X}_n$ given the latent function $\mathbf{f}_n$, and $p(\mathbf{y}_n | \mathbf{X}_n)$ is the marginal likelihood or evidence of the model.

The likelihood function is given by a factorized model:

$$p(\mathbf{y}_n | \mathbf{f}_n) = \prod_{i=1}^{n} p(y_i | f(\mathbf{x}_i)), \tag{8}$$

because the samples in $\mathcal{D}_n$ are iid. In turn, for each pair $(f(\mathbf{x}_i), y_i)$, the likelihood is given by (4), therefore

$$p(y_i | f(\mathbf{x}_i)) \sim \mathcal{N}(f(\mathbf{x}_i), \sigma_\nu^2). \tag{9}$$

A Gaussian likelihood function is conjugate to the Gaussian prior and hence the posterior in (7) is also a multidimensional Gaussian, which simplifies the computations to obtain (5). If the observation model were not Gaussian, warped Gaussian processes (see Section V) could be used to estimate (5).

Finally, we can obtain the posterior density in (5) for a general input $\mathbf{x}$ by conditioning on the training set and $\mathbf{x}$, and by marginalizing the latent function:

$$p(y | \mathbf{x}, \mathcal{D}_n) = \int p(y | f(\mathbf{x})) p(f(\mathbf{x}) | \mathbf{x}, \mathcal{D}_n) df(\mathbf{x}), \tag{10}$$

where[2]

$$p(f(\mathbf{x}) | \mathcal{D}_n, \mathbf{x}) = \int p(f(\mathbf{x}), \mathbf{f}_n | \mathbf{x}, \mathcal{D}_n) d\mathbf{f}_n. \tag{11}$$

We have divided the marginalization in two separate equations to show the marginalization of the latent function over the training set in (11), and the marginalization of the latent function at a general input $\mathbf{x}$ in (10). As mentioned earlier, the likelihood and the prior are Gaussians and therefore the marginalization in (10) and (11) only involves Gaussian distributions. Thereby, we can analytically compute (10) and (11) by using Gaussian conditioning and marginalization properties, leading to the following Gaussian density for the output:

$$p(f(\mathbf{x}) | \mathbf{x}, \mathcal{D}_n) \sim \mathcal{N}\left(\mu_{f(\mathbf{x})}, \sigma_{f(\mathbf{x})}^2\right), \tag{12}$$

where

$$\mu_{f(\mathbf{x})} = \mathbf{k}^\top \mathbf{C}_n^{-1} \mathbf{y}_n, \tag{13a}$$
$$\sigma_{f(\mathbf{x})}^2 = k(\mathbf{x}, \mathbf{x}) - \mathbf{k}^\top \mathbf{C}_n^{-1} \mathbf{k}, \tag{13b}$$

with

$$\mathbf{k} = [k(\mathbf{x}_1, \mathbf{x}), k(\mathbf{x}_2, \mathbf{x}), \ldots, k(\mathbf{x}_n, \mathbf{x})]^\top, \tag{14}$$
$$\mathbf{C}_n = \mathbf{K}_n + \sigma_\nu^2 \mathbf{I}_n. \tag{15}$$

The mean for $p(y | \mathbf{x}, \mathcal{D}_n)$ is also given by (13a), i.e., $\mu_y = \mu_{f(\mathbf{x})}$, and its variance is

$$\sigma_y^2 = \sigma_{f(\mathbf{x})}^2 + \sigma_\nu^2, \tag{16}$$

which, as expected, also accounts for the noise in the observation model.

The mean prediction of GPR in (13a) is the solution provided by KLS, or kernel ridge regression (KRR) [7], in which the covariance function takes the place of the kernel. However, unlike standard kernel methods, GPR provides error bars for each estimate in (13b) or (16) and has a natural procedure for setting the covariance/kernel by evidence sampling or maximization, as detailed in Section III. In SVM or KRR the hyper-parameters are typically adjusted by cross-validation, needing to retrain the models for different settings of the hyper-parameters on a grid search. So, typically only one or two hyper-parameters can be fitted. GPs can actually learn tens of hyper-parameter, because either sampling or evidence maximization allows setting them by a hassle-free procedure.

[2]Given the training data set, $\mathbf{f}_n$ takes values in $\mathbb{R}^n$ as it is a vector of $n$ samples of a Gaussian process.
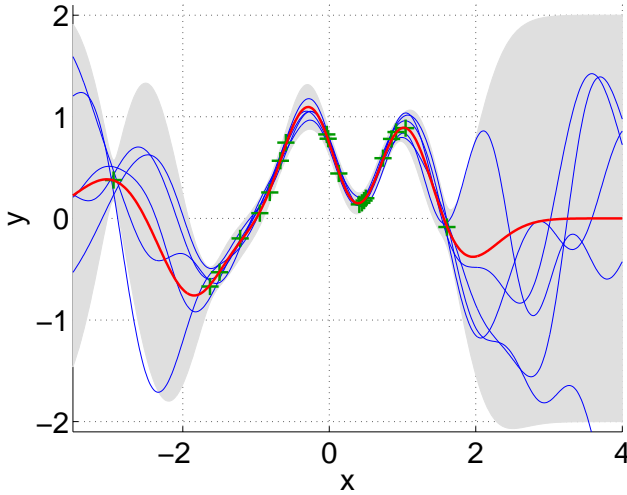
Fig. 2: Example of a Gaussian process posterior in (12) with 20 training samples, denoted by green +. Five instances of the posterior are plotted by thin blue lines and the mean of the posterior, $\mu_y$, by a red thick line. The shaded area denotes the error bars for the mean prediction: $\mu_y \pm 2\sigma_y$.

### C. An example

In Fig. 2 we include an illustrative example with 20 training points, in which we depict (12) for any $\mathbf{x}$ between $-3$ and $4$. We used standard functions from the GPML toolbox, available at http://www.gaussianprocess. org/gpml/, to generate the GP in this figure. We have chosen a Gaussian kernel that is fixed[3] as $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-2||\mathbf{x}_i - \mathbf{x}_j||^2)$ and $\sigma_\nu = 0.1$. In the plot, we show the mean of the process in red and the shaded area denotes the error bar for each prediction, i.e., $\mu_y \pm 2\sigma_y$. We also plot 5 samples from the posterior in thin blue lines.

We observe three different regions in the figure. On the right-hand side, we do not have samples and, for $\mathbf{x} > 3$, the GPR provides the solution given by the prior (zero mean and $\pm 2$). At the center, where most of the data points lie, we have a very accurate view of the latent function with small error bars (close to $\pm 2\sigma_\nu$). On the left hand side, we only have two samples and we notice the mixed effect of the prior widening the error bars and the data points constraining the values of the mean to lie close to the available samples. This is the typical behavior of GPR, which provides an accurate solution where the data lies and high error bars where we do not have available information and, consequently, we presume that the prediction in that area is not accurate.

---

[3]The kernel is typically expressed in a parametric form, see Section III.

### D. Recursive GPs

In many signal processing applications, the samples become available sequentially and estimation algorithms should obtain the new solution every time a new datum is received. In order to keep the computational complexity low, it is more interesting to perform inexpensive recursive updates rather than to recalculate the entire batch solution. *Online Gaussian Processes* [8] fulfill these requisites as follows.

Let us assume that we have observed the first $n$ samples and that at this point the new datum $\mathbf{x}_{n+1}$ is provided. We can readily compute the predicted distribution for $y_{n+1}$ using (13a), (13b) and (16). Furthermore, by using the formula for the inverse of a partitioned matrix and the Woodbury identity we update $\mathbf{C}_{n+1}^{-1}$ from $\mathbf{C}_n^{-1}$

$$\mathbf{C}_{n+1}^{-1} = \begin{bmatrix} \mathbf{C}_n^{-1} + \mathbf{C}_n^{-1}\mathbf{k}_{n+1}\mathbf{k}_{n+1}^\top\mathbf{C}_n^{-1}/\sigma_{y_{n+1}}^2 & -\mathbf{C}_n^{-1}\mathbf{k}_{n+1}/\sigma_{y_{n+1}}^2 \\ -\mathbf{k}_{n+1}^\top\mathbf{C}_n^{-1}/\sigma_{y_{n+1}}^2 & 1/\sigma_{y_{n+1}}^2 \end{bmatrix},$$

(17)

where $\sigma_{y_{n+1}}^2$ and $\mathbf{k}_{n+1}$ correspond to (16) and (14), respectively, for $\mathbf{x} = \mathbf{x}_{n+1}$.

Nevertheless, for online scenarios, it is more convenient to update the predicted mean and covariance matrix for all the available samples, as it is easier to interpret how the prediction changes with each new datum. Additionally, as we will show in Section VI, this formulation makes the adaptation to non-stationary scenarios straightforward. Let us denote by $\boldsymbol{\mu}_n$ and $\boldsymbol{\Sigma}_n$ the posterior mean and covariance matrix for the samples in $\mathcal{D}_n$. By applying (13a) and (13b) we obtain

$$\boldsymbol{\mu}_n = \mathbf{K}_n\mathbf{C}_n^{-1}\mathbf{y}_n, \tag{18a}$$

$$\boldsymbol{\Sigma}_n = \mathbf{K}_n - \mathbf{K}_n\mathbf{C}_n^{-1}\mathbf{K}_n, \tag{18b}$$

Once the new datum $(\mathbf{x}_{n+1}, y_{n+1})$ is observed, the updated mean and covariance matrix can be computed recursively as follows:

$$\boldsymbol{\mu}_{n+1} = \begin{bmatrix} \boldsymbol{\mu}_n \\ \mu_{f(\mathbf{x}_{n+1})} \end{bmatrix} - \frac{\mu_{f(\mathbf{x}_{n+1})} - y_{n+1}}{\sigma_{y_{n+1}}^2} \begin{bmatrix} \mathbf{h}_{n+1} \\ \sigma_{f(\mathbf{x}_{n+1})}^2 \end{bmatrix}$$

(19a)

$$\boldsymbol{\Sigma}_{n+1} = \begin{bmatrix} \boldsymbol{\Sigma}_n & \mathbf{h}_{n+1} \\ \mathbf{h}_{n+1}^\top & \sigma_{f(\mathbf{x}_{n+1})}^2 \end{bmatrix} - \frac{1}{\sigma_{y_{n+1}}^2} \begin{bmatrix} \mathbf{h}_{n+1} \\ \sigma_{f(\mathbf{x}_{n+1})}^2 \end{bmatrix} \begin{bmatrix} \mathbf{h}_{n+1}^\top & \sigma_{f(\mathbf{x}_{n+1})}^2 \end{bmatrix},$$

(19b)

where $\mathbf{h}_{n+1} = \boldsymbol{\Sigma}_n\mathbf{K}_n^{-1}\mathbf{k}_{n+1} = (\mathbf{I}_n - \mathbf{K}_n\mathbf{C}_n^{-1})\mathbf{k}_{n+1}$. As can be observed in (19a), the mean of the new process is obtained by applying a correction term to the previous mean, proportional to the estimation error, $\mu_{f(\mathbf{x}_{n+1})} - y_{n+1}$. Because of the relation between $\boldsymbol{\Sigma}_n$ and $\mathbf{C}_n^{-1}$ stated in (18b), only one of the two matrices needs to be stored and updated in an online formulation.

Some authors [8] prefer to rely on $\mathbf{C}_n^{-1}$, whereas others [9] store and update $\boldsymbol{\Sigma}_n$.

The recursive update of the mean in (19a) is equivalent to what is known as *kernel recursive least-squares* (KRLS) in the signal processing literature (see for instance [8]–[10] ). The unbounded growth of the involved matrices, visible in (19) and (17), is the main limitation in the KRLS formulation. Practical KRLS implementations typically either limit this growth [10,11] or even fix the matrix sizes [12]. Nevertheless, the solution of KRLS is limited to the mean only and it cannot estimate confidence intervals. By using a GP framework, though, an estimate of the entire posterior distribution is obtained, including the covariance in (19b).

*E. Connection to MMSE: GPR with a linear latent function*

If we replace $f(\mathbf{x})$ in (4) with a linear model

$$y = \mathbf{w}^\top \mathbf{x} + \nu,$$

the Gaussian process prior over $f(\mathbf{x})$ becomes a spherical-Gaussian prior distribution over $\mathbf{w}$, $p(\mathbf{w}) \sim \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{w}}^2 \mathbf{I})$.

We can now compute the posterior for $\mathbf{w}$, as we did for the latent function in (7)

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w}) p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})} = \frac{p(\mathbf{w})}{p(\mathbf{y}|\mathbf{X})} \prod_{i=1}^n p(y_i|\mathbf{x}_i, \mathbf{w}),$$

where $p(y_i|\mathbf{x}_i, \mathbf{w})$ is the likelihood. Since the prior and likelihood are Gaussians, so it is the posterior, and its mean and covariance are given by

$$\mu_{\mathbf{w}} = \frac{1}{\sigma_\nu^2} \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{X}^\top \mathbf{y}, \tag{20a}$$

$$\boldsymbol{\Sigma}_{\mathbf{w}} = \left( \mathbf{X}^\top \mathbf{X}/\sigma_\nu^2 + \mathbf{I}/\sigma_{\mathbf{w}}^2 \right)^{-1}. \tag{20b}$$

We can readily notice that (20a) is the sampled version of (3), when the prior variance $\sigma_{\mathbf{w}}^2$ tends to infinity (i.e., the prior has no effect of the solution). The precision matrix (the inverse covariance) is composed of two terms: the first depends on the data and the other one on the prior over $\mathbf{w}$. The effect of the prior in the mean and covariance fades away, as we have more available data. The estimate for a general input $\mathbf{x}$ is computed as in (10)

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathbf{w}) p(\mathbf{w}|\mathcal{D}) d\mathbf{w}, \tag{21}$$

which is a Gaussian distribution with mean and variance given by:

$$\mu_y = \mathbf{x}^\top \mu_{\mathbf{w}} = \frac{1}{\sigma_\nu^2} \mathbf{x}^\top \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{X}^\top \mathbf{y} \tag{22}$$

$$\sigma_y^2 = \mathbf{x}^\top \boldsymbol{\Sigma}_{\mathbf{w}} \mathbf{x} + \sigma_\nu^2 \tag{23}$$

Equations (22) and (23) can be, respectively, rewritten as (13a) and (16), if we use the inner product between the $\mathbf{x}_i$ multiplied by the width of the prior over $\mathbf{w}$, i.e. the kernel matrix is given by: $\mathbf{K}_n = \mathbf{X} \sigma_{\mathbf{w}}^2 \mathbf{I} \mathbf{X}^\top$. The kernel matrix must include the width of the prior over $\mathbf{w}$, because the kernel matrix represents the prior of the Gaussian process and $\sigma_{\mathbf{w}}^2$ is the prior of the linear Bayesian estimator. By using the Woodbury's identity, it follows that

$$\boldsymbol{\Sigma}_{\mathbf{w}} = \sigma_{\mathbf{w}}^2 \mathbf{I} - \sigma_{\mathbf{w}}^2 \mathbf{I} \mathbf{X}^\top \left( \sigma_\nu^2 \mathbf{I} + \mathbf{K}_n \right)^{-1} \mathbf{X} \sigma_{\mathbf{w}}^2 \mathbf{I}. \tag{24}$$

Now, by replacing (24) in (22) and (23), we, respectively, recover (13a) and (16). These steps connect the estimation of a Bayesian linear model and the nonlinear estimation using a kernel or covariance function without needing to explicitly indicate the nonlinear mapping.

### III. COVARIANCE FUNCTIONS

In the previous section, we have assumed that the covariance functions $k(\mathbf{x}, \mathbf{x}')$ are known, which is not typically the case. In fact, the design of a good covariance function is crucial for GPs to provide accurate nonlinear solutions. The covariance function plays the same role as the kernel function in SVMs or KLS [7]. It describes the relation between the inputs and its form determines the possible solutions of the GPR. It controls how fast the function can change or how the samples in one part of the input space affect the latent function everywhere else. For most problems, we can specify a parametric kernel function that captures any available information about the problem at hand. As already discussed, unlike kernel methods, GPs can infer these parameters, the so-called *hyper-parameters*, from the samples in $\mathcal{D}_n$ using the Bayesian framework. Instead of relying on computational intensive procedures as cross-validation [13] or learning the kernel matrix [14], as kernel methods need to.

The covariance function must be positive semidefinite, as it represents the covariance matrix of a multidimensional Gaussian distribution. The covariance can be built by adding simpler covariance matrices, weighted by a positive hyper-parameter, or by multiplying them together, as the addition and multiplication of positive definite matrices yields a positive definite matrix. In general, the design of the kernel should rely on the information that we have for each estimation problem and should be designed to get the most accurate solution with the least amount of samples. Nevertheless, the following kernel in (25) often works well in signal

processing applications

$$k(\mathbf{x}_i, \mathbf{x}_j) = \alpha_1 \exp\left(-\sum_{\ell=1}^{d} \gamma_\ell ||x_{i\ell} - x_{j\ell}||^2\right) + \alpha_2 \mathbf{x}_i^\top \mathbf{x}_j + \alpha_3 \delta_{ij} \tag{25}$$

where $\boldsymbol{\theta} = [\alpha_1, \gamma_1, \gamma_2, \ldots, \gamma_d, \alpha_2, \alpha_3]^\top$ are the hyper-parameters. The first term is a radial basis kernel, also denoted as RBF or Gaussian, with a different length-scale for each input dimension. This term is universal and allows constructing a generic nonlinear regressor. If we have symmetries in our problem, we can use the same length-scale for all dimensions: $\gamma_\ell = \gamma$ for $\ell = 1, \ldots, d$. The second term is the linear covariance function. The last term represents the noise variance $\alpha_3 = \sigma_\nu^2$, which can be treated as an additional hyper-parameter to be learned from the data. We can add other terms or other covariance functions that allow for faster transitions, like the Matérn kernel among others [5].

If the hyper-parameters, $\boldsymbol{\theta}$, are unknown, the likelihood in (8) and the prior in (6) can, respectively, be expressed[4] as $p(\mathbf{y}|\mathbf{f}, \boldsymbol{\theta})$ and $p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta})$, and we can proceed to integrate out $\boldsymbol{\theta}$ as we did for the latent function, $\mathbf{f}$, in Section II-B. First, we compute the *marginal likelihood* of the hyper-parameters of the kernel given the training dataset

$$p(\mathbf{y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f}, \boldsymbol{\theta}) p(\mathbf{f}|\mathbf{X}, \boldsymbol{\theta}) d\mathbf{f}. \tag{26}$$

Second, we can define a prior for the hyper-parameters, $p(\boldsymbol{\theta})$, that can be used to construct its posterior. Third, we integrate out the hyper-parameters to obtain the predictions. However, in this case, the marginal likelihood does not have a conjugate prior and the posterior cannot be obtained in closed form. Hence, the integration has to be done either by sampling or approximations. Although this approach is well principled, it is computational intensive and it may be not feasible for some applications. For example, Markov-Chain Monte Carlo (MCMC) methods require several hundred to several thousand samples from the posterior of $\boldsymbol{\theta}$ to integrate it out. Interested readers can find further details in [5].

Alternatively, we can maximize the marginal likelihood in (26) to obtain its optimal setting [1]. Although setting the hyper-parameters by maximum likelihood (ML) is not a purely Bayesian solution, it is fairly standard in the community and it allows using Bayesian solutions in time sensitive applications. This optimization is nonconvex [15], but, as we increase the number of training samples, the likelihood becomes a unimodal distribution around the ML hyper-parameters and the

---

[4]We have dropped the subindex $n$, as it is inconsequential and unnecessarily clutters the notation.

solution can be found using gradient ascent techniques. See [5] for further details.

## IV. SPARSE GPs: DEALING WITH LARGE-SCALE DATA SETS

To perform inference under any GP model, the inverse of the covariance matrix must be computed. This is a costly operation, $\mathcal{O}(n^3)$, that becomes prohibitive for large enough $n$. Given the ever-increasing availability of large-scale databases, a lot of effort has been devoted over the last decade to the development of approximate methods that allow inference in GPs to scale linearly with the number of data points. These approximate methods are referred to as "sparse GPs", since they approximate the full GP model using a finite-basis-set expansion. This set of bases is usually spawned by using a common functional form with different parametrizations. For instance, it is common to use bases of the type $\{k(\mathbf{z}_b, \mathbf{x})\}_{b=1}^{m}$, where $\{\mathbf{z}_b\}_{b=1}^{m}$ —known as the *active set*— is a subset of the input samples parametrizing the bases.

Under the unifying framework of [16], it can be shown that most relevant sparse GP proposals [17,18], which were initially thought of as entirely different low-cost approximations, can be expressed as exact inference under different modifications of the original GP prior. This modified prior induces a rank-$m$ ($m \ll n$) covariance matrix —plus optional (block) diagonal correcting terms—, clarifying how the reduced $\mathcal{O}(m^2 n)$ cost of exact inference arises.

Among the mentioned approximations, the sparse pseudo-input GP (SPGP) [18] is generally regarded as the most efficient. Unlike other alternatives, it does not require the active set to be a subset of the training data. Instead, $\{\mathbf{z}_b\}_{b=1}^{m}$ can be selected to lie anywhere in the input space, thus increasing the flexibility of the finite set expansion. This selection is typically performed by evidence maximization. An even more flexible option, which does not require the active set to even lie in the input domain, is presented in [19].

Despite the success of SPGP, it is worth mentioning that increasing the number of bases in this algorithm does not yield, in general, convergence to the full GP solution because the active set $\{\mathbf{z}_b\}_{b=1}^{m}$ is not constrained to be a subset of input data. This might lead to overfitting in some pathological cases. A recent variational sparse GP proposal that guarantees convergence to the full GP solution while still allowing the active set to be unconstrained is presented in [20].

Further approaches yielding reduced computational cost involve numerical approximations to accelerate matrix-vector multiplications and compactly supported

covariance functions which set most entries of the co-variance matrix to zero [21].

Sparsity is often seen in online signal processing in the form of *pruning*, which restricts the active set to a subset of input data. The success of SPGP and its variational counterpart suggests that advanced forms of pruning may result in increased efficiency for a given sparsity level.

## V. WARPED GPs: BEYOND THE STANDARD NOISE MODEL

Even though GPs are very flexible priors for the latent function, they might not be suitable to model all types of data. It is often the case that applying a logarithmic transformation to the target variable of some regression task (e.g., those involving stock prices, measured sound power, etc) can enhance the ability of GPs to model it.

In [22] it is shown that it is possible to include a non-linear preprocessing of output data $h(y)$ (called *warping function* in this context) as part of the modeling process and learn it. In more detail, a parametric form for $z = h(y)$ is selected, then $z$ (which depends on the parameters of $h(y)$) is regarded as a GP, and finally, the parameters of $h(y)$ are selected by maximizing the evidence of such GP (i.e., a ML approach). The authors suggest using $h(y) = \sum_{i=1}^{l} a_i \tanh(b_i(y + c_i))$ as the parametric form of the warping function, but any option resulting in a monotonic function is valid. A non-parametric version of warped GPs using a variational approximation has been proposed in [23].

## VI. TRACKING NON-STATIONARY SCENARIOS: LEARNING TO FORGET

KRLS algorithms, discussed in Section II-D, traditionally consider that the mapping function $f(\cdot)$ is constant throughout the whole learning process [10,24]. However, in the signal processing domain this function (which might represent, for instance, a fading channel) is often subject to changes and the model must account for this non-stationarity. Some kernel-based algorithms have been proposed to deal with non-stationary scenarios. They include a kernelized version of the extended RLS filter [24], a sliding-window KRLS approach [12] and a family of projection-based algorithms [25,26].

In order to add adaptivity to the online GP algorithm described in Section II-D, it is necessary to make it "forget" the information contained in old samples. This becomes possible by including a "forgetting" step after each update

$$\boldsymbol{\mu} \leftarrow \sqrt{\lambda}\boldsymbol{\mu} \tag{27a}$$

$$\boldsymbol{\Sigma} \leftarrow \lambda\boldsymbol{\Sigma} + (1 - \lambda)\mathbf{K}. \tag{27b}$$

to shift the posterior distribution towards the prior (for $0 < \lambda < 1$), thus effectively reducing the influence of older samples. Note that when using this formulation there is no need to store or update $\mathbf{C}^{-1}$, see [9] for further details. The adaptive, GP-based algorithm obtained in this manner is known as KRLS-T.

Equations (27) might seem like an *ad-hoc* step to enable forgetting. However, it can be shown that the whole learning procedure —including the mentioned controlled forgetting step— corresponds exactly to a principled non-stationary scheme within the GP framework, as described in [27]. It is sufficient to consider an augmented input space that includes the time stamp $t$ of each sample and define a *spatio-temporal* covariance function:

$$k_{\text{st}}([t \ \mathbf{x}^\top]^\top, [t' \ \mathbf{x}'^\top]^\top) = k_{\text{t}}(t, t')k_{\text{s}}(\mathbf{x}, \mathbf{x}'), \tag{28}$$

where $k_{\text{s}}(\mathbf{x}, \mathbf{x}')$ is the already-known spatial covariance function and $k_{\text{t}}(t, t')$ is a temporal covariance function giving more weight to samples that are closer in time. Inference on this augmented model effectively accounts for non-stationarity in $f(\cdot)$ and recent samples have more impact in predictions for the current time instant. It is fairly simple to include this augmented model in the online learning process described in the previous section. When the temporal covariance is set to $k_{\text{t}}(t, t') = \lambda^{\frac{|t-t'|}{2}}, \ \lambda \in (0, 1]$, inference in the augmented spatio-temporal GP model is exactly equivalent to using (27) after each update (19) in the algorithm of Section II-D, which has the added benefit of being inexpensive and online. See [9,27,28] for further details.

Observe that $\lambda$ is used here to model the speed at which $f(\cdot)$ varies, playing a similar rôle to that of the forgetting factor in linear adaptive filtering algorithms. When used with a linear spatial covariance, the above model reduces to linear extended RLS filtering. The selection of this parameter is usually rather *ad-hoc*. However, using the GP framework, we can select it in a principled manner using Type-II ML, see [27].

In Fig. 3 we take the example of Fig. 2 and we apply a forgetting factor $\lambda = 0.8$. The red continuous line indicates the original mean function before forgetting. After applying one forgetting update, this mean function is displaced toward zero, as indicated by the the blue dashed line. The shaded gray area represents the error bars prior to forgetting. The forgetting update expands this area into the shaded red area, which tends to the prior variance of 1.
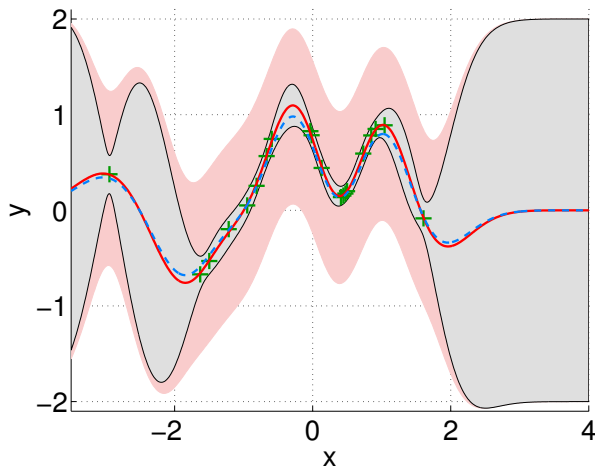
Fig. 3: Illustration of forgetting step (27) on the GP of Fig. 2: the dashed line represents the predictive mean that is pulled towards the prior mean, while the shaded red area represents the region $\mu_y \pm 2\sigma_y$ after forgetting.

### A. Tracking a time-selective nonlinear communication channel

To illustrate the validity of the adaptive filtering algorithm, we focus on the problem of tracking a nonlinear Rayleigh fading channel [29, Chapter 7]. The used model consists of a memoryless saturating nonlinearity followed by a time-varying linear channel, as shown in Fig. 4. This model appears for instance in broadcast or satellite communications when the amplifier operates close to saturation regime [30].

In a first, simulated setup, the time-varying linear fading channel consists of 5 randomly generated paths, and the saturating nonlinearity is chosen as $y = \tanh(x)$. We fix the symbol rate at $T = 1\mu s$, and we simulate two scenarios: one with a normalized Doppler frequency of $f_dT = 10^{-4}$ (where $f_d$ denotes the Doppler spread), representing a slow-fading channel, and another one with $f_dT = 10^{-3}$, corresponding to a fast time-varying channel. Note that a higher Doppler frequency yields a more difficult tracking problem, as it corresponds to a channel that changes faster in time. We consider a Gaussian source signal, and we add 30 dB of additive white Gaussian noise to the output signal. Given one input-output data pair per time instant, the tracking problem consists in estimating the received signal that corresponds to a new channel input.

Figs. 5(a,b) illustrate the tracking results obtained by KRLS-T in these scenarios. As a reference, we include the performance of several state-of-the-art adaptive filtering algorithms, whose Matlab implementations are taken from the Kernel Adaptive Filtering Toolbox, available at http://sourceforge.net/projects/kafbox/. In particular, we
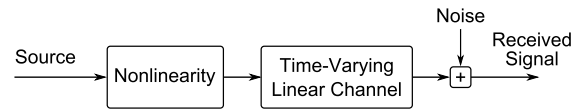


Fig. 4: The nonlinear channel used in the example consists of a nonlinearity followed by a linear channel.

compare KRLS-T with normalized least mean squares (NLMS), extended RLS (EX-RLS), both of which are linear algorithms, see [29], and quantized kernel LMS (QKLMS) [31], which is an efficient, kernelized version of the LMS algorithm. A Gaussian kernel $k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2)$ is used for QKLMS and KRLS-T. In each scenario the optimal hyperparameters of KRLS-T are obtained by performing Type-II ML optimization (see Section III ) on a separate data set of 500 test samples. The optimal parameters of the other algorithms are obtained by performing cross-validation on the test data set. To avoid an unbounded growth of the matrices involved in KRLS-T, its memory is limited to 100 bases which are selected by pruning the least relevant bases (see [9] for details on the pruning mechanism). The quantization parameter of QKLMS is set to yield similar memory sizes. As can be seen in Figs. 5(a,b), KRLS-T outperforms the other algorithms with a significant margin in both scenarios. By being kernel-based it is capable to deal with nonlinear identification problems, in contrast to the classical EX-RLS and NLMS algorithms. Furthermore, it shows excellent convergence speed and steady-state performance when compared to QKLMS. Additional experimental comparisons to other kernel adaptive filters can be found in [9].

In a second setup we used a wireless communication test bed that allows to evaluate the performance of digital communication systems in realistic indoor environments. This platform is composed of several transmit and receive nodes, each one including a radio-frequency front-end and baseband hardware for signal generation and acquisition. The front-end also incorporates a programmable variable attenuator to control the transmit power value and therefore the signal saturation. A more detailed description of the test bed can be found in [32]. Using the hardware platform, we reproduced the model corresponding to Fig. 4 by transmitting clipped orthogonal frequency-division multiplexing (OFDM) signals centered at 5.4 GHz over real frequency-selective and time-varying channels. Notice that, unlike the simulated setup, several parameters such as the noise level and the variation of the channel coefficients are unknown. To have an idea about the channel characteristics, we first measured the indoor channel using the procedure

(a) $f_d T = 10^{-4}$, simulated data

(b) $f_d T = 10^{-3}$, simulated data

(c) $f_d T = 10^{-3}$, real data
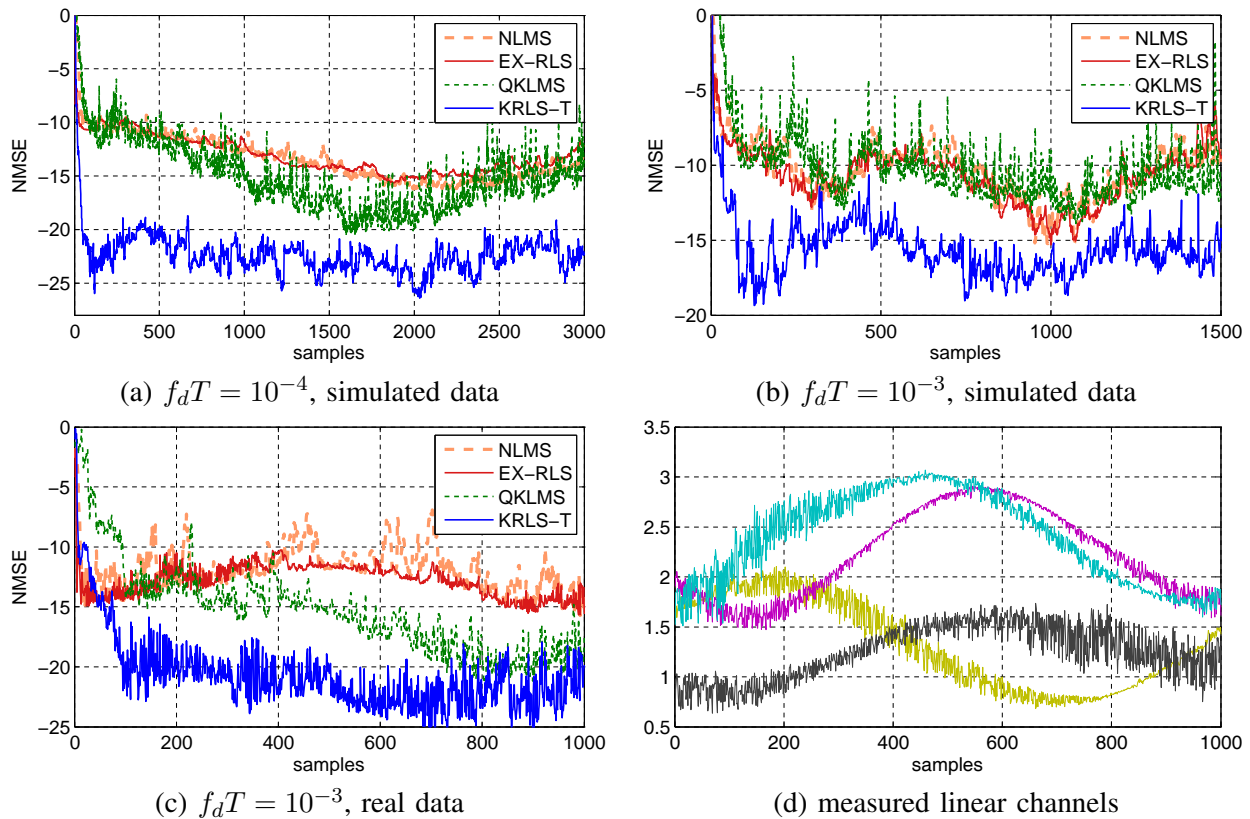
(d) measured linear channels

Fig. 5: Tracking results on a nonlinear Rayleigh fading channel: (a) simulation results for a slow-fading scenario; (b) simulation results for a fast time-varying scenario; (c) tracking results on data measured on the test bed with fast time-varying channels; (d) channel taps of the noisy linear channels, measured on the test bed setup.

TABLE I: Steady-state NMSE performance for Fig. 5.

|  | NLMS | EX-RLS | QKLMS | KRLS-T |
|---|---|---|---|---|
| $f_d T = 10^{-4}$, simulated | $-13.3$ dB | $-13.0$ dB | $-14.6$ dB | $-22.3$ dB |
| $f_d T = 10^{-3}$, simulated | $-10.6$ dB | $-11.0$ dB | $-9.9$ dB | $-15.3$ dB |
| $f_d T = 10^{-3}$, real data | $-11.5$ dB | $-12.5$ dB | $-15.8$ dB | $-21.3$ dB |

described in [32]. As an example, the variation of the four main channel coefficients is depicted in Fig. 5(d), indicating a normalized Doppler frequency around $f_d T = 10^{-3}$. We then transmitted periodically OFDM signals with the transmit amplifier operating close to saturation and acquired the received signals. The transmitted and received signals were used to track the nonlinear channel variations as in the simulated setup. The results, shown in Fig. 5(c), are similar to those of the simulated setup. Finally, the steady-state NMSE performances of all three scenarios, Figs. 5(a,b,c), are summarized in Table I.

## VII. GAUSSIAN PROCESSES FOR CLASSIFICATION

For classification problems, the labels are drawn from a finite set and GPs return a probabilistic prediction for each label in the finite set, i.e., how certain is the classifier about its prediction. In this tutorial, we limit our presentation of GPs for classification (GPC) for binary classification problems, i.e., $y_i \in \{0, 1\}$. For GPC, we change the likelihood model for the latent function at $\mathbf{x}$ using a *response function* $\Phi(\cdot)$:

$$p(y = 1 | f(\mathbf{x})) = \Phi(f(\mathbf{x})). \tag{29}$$

The response function "squashes" the real-valued latent function to an $(0, 1)$-interval that represents the posterior probability for $y$ [5]. Standard choices for the response function are $\Phi(a) = 1/(1 + \exp(-a))$ and the cumulative density function of a standard normal distribution, used in logistic and probit regression respectively.

The integrals in (10) and (11) are now analytically intractable, because the likelihood and the prior are not conjugated. Therefore, we have to resort to numerical

methods or approximations to solve them. The posterior distribution in (7) is typically single-mode and the standard methods approximate it with a Gaussian [5]. Using a Gaussian approximation for (7) allows exact marginalization in (11) and we can use numerical integration for solving (10), as it involves marginalizing a single real-valued quantity. The two standard approximations are the Laplace method or expectation propagation (EP) [33]. In [2], EP is shown to be a more accurate approximation.

### A. Probabilistic channel equalization

GPC predictive performance is similar to other nonlinear discriminative methods, such as SVMs. However, if the probabilistic output is of importance, then GPC outperforms other kernel algorithms, because it naturally incorporates the confidence interval in its predictions. In digital communication, channel decoders follow equalizers, which work optimally when accurate posterior estimates are given for each symbol. To illustrate that GPC provide accurate posterior probability estimates, we equalize a dispersive channel model like the one in Fig. 4 using GPC and SVM with a probabilistic output. These outputs are subsequently fed to a low-density parity-check (LDPC) belief-propagation based channel decoder to assess the quality of the estimated posterior probabilities. Details for the experimental set up can be found in [34] in which linear and nonlinear channel models are tested. We now summarize the results for the linear channel model in that paper.

In Fig. 6, we depict the posterior probability estimates versus the true posterior probability, in (a) for the GPC-based equalizer and in (b) for SVM-based equalizer, to emphasize the differences between the equalizers we use a highly noisy scenario with normalized signal-to-noise ratio of 2 dB. If we threshold at 0.5, both equalizers provide similar error rates and we cannot tell if there is an advantage from using GPC. However, if we consider the whole probability space, GPC predictions are significantly closer to the main diagonal that represents a perfect match, hence GPC provides more accurate predictions to the channel decoder.

To further quantify the gain from using a GPC-based equalizer with accurate posterior probability estimates, we plot the bit error rate (BER) in Fig. 7 after the probabilistic channel encoder, in which the GPC-based equalizer clearly outperforms the SVM-based equalizer and is close to the optimal solution (known channel and forward-backward (BCJR) equalizer). This example is illustrative of the results that can be expected from GPC when a probabilistic output is needed to perform optimally.
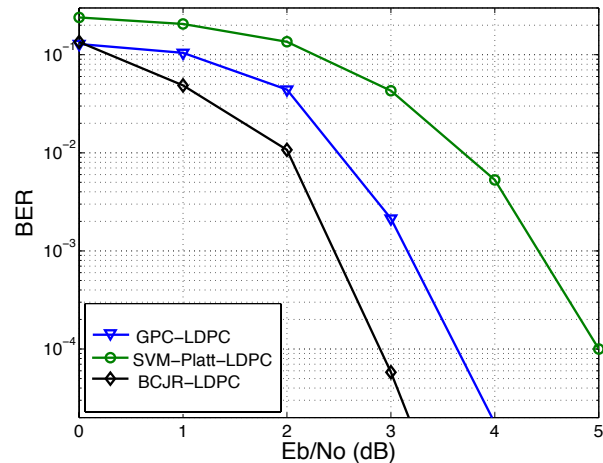


Fig. 7: GPC and SVM as probabilistic channel equalizer in channel LDPC decoding: BER for the GPC-LDPC ($\triangledown$), the SVM-LDPC ($\circ$) and the optimal solution ($\diamond$).

## VIII. DISCUSSION

In this tutorial, we have presented Gaussian Processes for Regression in detail from the point of view of MMSE/Wiener filtering, so it is amenable to signal processing practitioners. GPR provides the same mean estimate as KLS or KRR for the same kernel matrix. On the plus side, GPR provides error bars that take into account the approximation error and the error from the likelihood model, so we know the uncertainty of our model for any input point (see Fig. 2 ), while KLS assumes the error bars are given by the likelihood function (i.e., constant for the whole input space). Additionally, GPR naturally allows computing the hyper-parameters of the kernel function by sampling or maximizing the marginal likelihood, being able to set tens of hyper-parameters, while KLS or SVM need to rely on cross-validation, in which only one or two parameters can be easily tuned. On the minus side, the GP prior imposes a strong assumption on the error bars that might not be accurate, if the latent variable model does not follow a Gaussian process. Although, in any case, it is better than not having error bars.

We have also shown that some of the limitations of the standard GPR can be eased. GPs can be extended to non-Gaussian noise models and classification problems, in which GPC provides an accurate a posteriori probability estimate. The computational complexity of GPs can be reduced considerably, from cubic to linear in the number of training examples, without significantly affecting the mean and error bars prediction. Finally, we have shown the GP can be solved iteratively, with an RLS formulation that can be adapted to non-stationary environments efficiently.
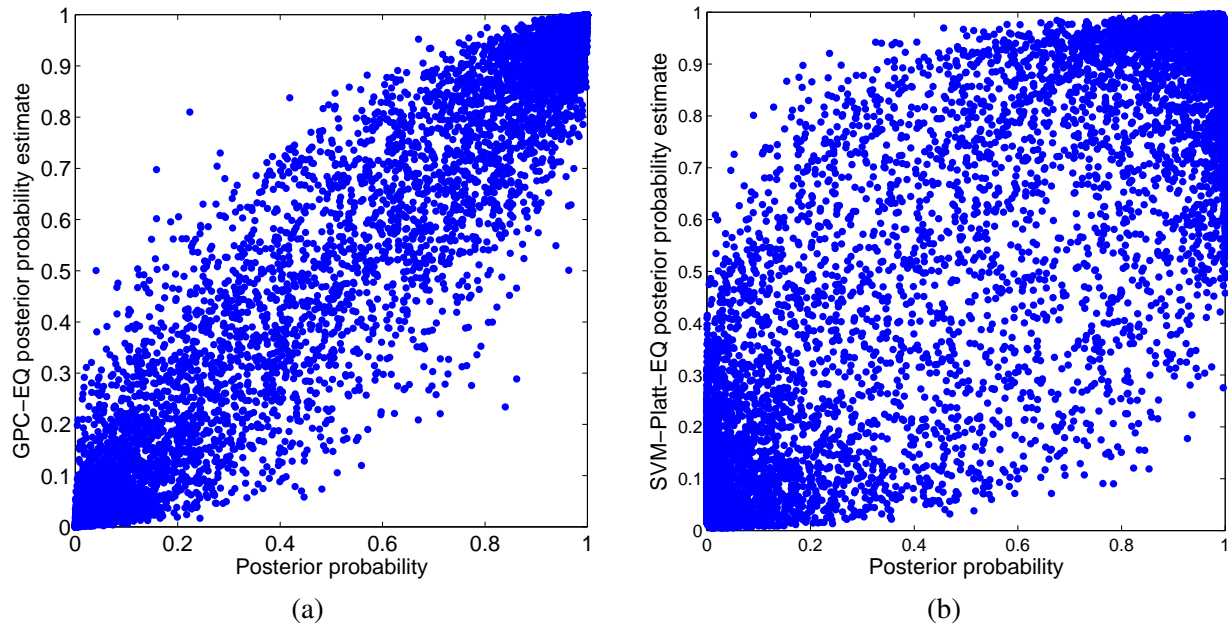
Fig. 6: GPC as probabilistic channel equalizer: (a) calibration curve for the GPC and (b) calibration curve for the SVM.

Instead of covering more methods and applications in detail, our intention was to provide a tutorial paper on how to use GPs in signal processing, with a number of illustrative examples. Nevertheless, since we assume that there are several other methods and applications that are relevant to the reader, we finish with a brief list of further topics. In particular, GPs have also been applied to problems including modeling human motion [35], source separation [36], estimating chlorophyll concentration [37], approximating stochastic differential equations [38] and multi-user detection [39], among others.

## IX. ACKNOWLEDGMENTS

## REFERENCES

[1] C. K. I. Williams and C. E. Rasmussen, "Gaussian processes for regression," in *Neural Information Processing Systems 8*. MIT Press, 1996, pp. 598–604.

[2] M. Kuss and C. Rasmussen, "Assessing approximate inference for binary Gaussian process classification," *Machine learning research*, vol. 6, pp. 1679–1704, Oct. 2005.

[3] N. Lawrence, "Probabilistic non-linear principal component analysis with Gaussian process latent variable models," *Machine learning research*, vol. 6, pp. 1783–1816, Nov. 2005.

[4] A. O'Hagan and J. F. Kingman, "Curve fitting and optimal design for prediction," *Journal of the Royal Statistical Society. Series B*, vol. 40, no. 1, pp. 1783–1816, 1978.

[5] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. MIT Press, 2006.

[6] V. N. Vapnik, *Statistical Learning Theory*. New York: John Wiley & Sons, 1998.

[7] F. Pérez-Cruz and O. Bousquet, "Kernel methods and their potential use in signal processing," *Signal Processing Magazine*, vol. 21, no. 3, pp. 57–65, 2004.

[8] L. Csató and M. Opper, "Sparse representation for Gaussian process models," in *Neural Information Processing Systems 13*. MIT Press, 2001, pp. 444–450.

[9] S. Van Vaerenbergh, M. Lázaro-Gredilla, and I. Santamaría, "Kernel recursive least-squares tracker for time-varying regression," *IEEE Trans. on Neural Networks and Learning Systems*, vol. 23, no. 8, pp. 1313–1326, Aug. 2012.

[10] Y. Engel, S. Mannor, and R. Meir, "The kernel recursive least squares algorithm," *IEEE Trans. on Signal Processing*, vol. 52, no. 8, pp. 2275–2285, Aug. 2004.

[11] W. Liu, I. Park, and J. C. Príncipe, "An information theoretic approach of designing sparse kernel adaptive filters," *IEEE Trans. on Neural Networks*, vol. 20, no. 12, pp. 1950–1961, 2009.

[12] S. Van Vaerenbergh, J. Vía, and I. Santamaría, "A sliding-window kernel RLS algorithm and its application to nonlinear channel identification," in *Proc. of IEEE Int. Conference on Acoustics, Speech and Signal Processing (ICASSP)*, vol. 5, Toulouse, France, May 2006, pp. 789–792.

[13] G. S. Kimeldorf and G. Wahba, "Some results in Tchebycheffian spline functions," *Journal of Mathematical Analysis and Applications*, vol. 33, pp. 82–95, 1971.

[14] O. Bousquet and D. J. L. Herrmann, "On the complexity of learning the kernel matrix," in *In Advances in Neural Information Processing Systems 15*. MIT Press, 2003, pp. 399–406.

[15] D. J. C. MacKay, *Information Theory, Inference and Learning Algorithms*. Cambridge, UK: Cambridge University Press, 2003.

[16] J. Quiñonero-Candela and C. E. Rasmussen, "A unifying view of sparse approximate Gaussian process regression," *Machine learning research*, vol. 6, pp. 1939–1959, Dec. 2005.

[17] M. Seeger, C. K. I. Williams, and N. D. Lawrence, "Fast for-

ward selection to speed up sparse Gaussian process regression," in *Proc. of 9th Int. Workshop on Artificial Intelligence and Statistics*, 2003.

[18] E. Snelson and Z. Ghahramani, "Sparse Gaussian processes using pseudo-inputs," in *Advances in Neural Information Processing Systems 18*. MIT Press, 2006, pp. 1259–1266.

[19] M. Lázaro-Gredilla and A. Figueiras-Vidal, "Inter-domain Gaussian processes for sparse inference using inducing features," in *Advances in Neural Information Processing Systems 22*. MIT Press, 2010, pp. 1087–1095.

[20] M. K. Titsias, "Variational learning of inducing variables in sparse Gaussian processes," in *Proc. of 12th Int. Workshop on Artificial Intelligence and Statistics*, 2009, pp. 567–574.

[21] D. Gu, "Spatial Gaussian process regression with mobile sensor networks," *IEEE Trans. on Neural Networks and Learning Systems*, vol. 23, pp. 1279–1290, 2012.

[22] E. Snelson, Z. Ghahramani, and C. Rasmussen, "Warped Gaussian processes," in *Advances in Neural Information Processing Systems 16*. MIT Press, 2003.

[23] M. Lázaro-Gredilla, "Bayesian warped Gaussian processes," in *Advances in Neural Information Processing Systems 26*. MIT Press, 2013.

[24] W. Liu, J. C. Príncipe, and S. Haykin, *Kernel Adaptive Filtering: A Comprehensive Introduction*. Wiley, 2010.

[25] K. Slavakis, S. Theodoridis, and I. Yamada, "Online kernel-based classification using adaptive projection algorithms," *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 2781 – 2796, july 2008.

[26] S. Theodoridis, K. Slavakis, and I. Yamada, "Adaptive learning in a world of projections," *IEEE Signal Processing Magazine*, vol. 28, no. 1, pp. 97 –123, Jan. 2011.

[27] S. Van Vaerenbergh, I. Santamaría, and M. Lázaro-Gredilla, "Estimation of the forgetting factor in kernel recursive least squares," in *Proc. of IEEE Int. Workshop on Machine Learning for Signal Processing (MLSP)*, Sept. 2012, pp. 1–6.

[28] M. Lázaro-Gredilla, S. Van Vaerenbergh, and I. Santamaría, "A Bayesian approach to tracking with kernel recursive least-squares," in *Proc. of IEEE Int. Workshop on Machine Learning for Signal Processing (MLSP)*, Sept. 2011, pp. 1 –6.

[29] A. Sayed, *Fundamentals of adaptive filtering*. Wiley-IEEE Press, 2003.

[30] K. Feher, *Digital Communications: Satellite/Earth Station Engineering*. Englewood Cliffs, N.J.: Prentice-Hall, 1983.

[31] B. Chen, S. Zhao, P. Zhu, and J. C. Principe, "Quantized kernel least mean square algorithm," *IEEE Trans. on Neural Networks and Learning Systems*, vol. 23, no. 1, pp. 22–32, Jan. 2012.

[32] J. Gutiérrez, Ó. González, J. Pérez, D. Ramírez, L. Vielva, J. Ibáñez, and I. Santamaría, "Frequency-domain methodology for measuring MIMO channels using a generic test bed," *IEEE Trans. on Instrumentation and Measurement*, vol. 60, no. 3, pp. 827–838, Mar. 2011.

[33] T. Minka, "Expectation propagation for approximate Bayesian inference," in *Proc. of 17th Conference in Uncertainty in Artificial Intelligence*, University of Washington, Seattle, Washington, USA, 2001, pp. 362–369.

[34] P. Olmos, J. Murillo-Fuentes, and F. Pérez-Cruz, "Joint nonlinear channel equalization and soft LDPC decoding with Gaussian processes," *IEEE Trans. on Signal Processing*, vol. 58, no. 3, pp. 1183–1192, 2010.

[35] J. Wang, D. Fleet, and A. Hertzmann, "Gaussian process dynamical models for human motion," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 30, no. 2, pp. 283 – 298, Feb. 2008.

[36] A. Liutkus, R. Badeau, and G. Richard, "Gaussian processes for underdetermined source separation," *IEEE Trans. on Signal Processing*, vol. 59, no. 7, pp. 3155–3167, July 2011.

[37] L. Pasolli, F. Melgani, and E. Blanzieri, "Gaussian process regression for estimating chlorophyll concentration in subsurface waters from remote sensing data," *IEEE Geoscience and Remote Sensing Letters*, vol. 7, no. 3, pp. 464 – 468, Mar. 2010.

[38] C. Archambeau, D. Cornford, M. Opper, and J. Shawe-Taylor, "Gaussian process approximations of stochastic differential equations," *Journal of Machine Learning Research*, vol. 1, pp. 1–16, 2007.

[39] J. J. Murillo-Fuentes and F. Pérez-Cruz, "Gaussian process regressors for multiuser detection in DS-CDMA systems," *IEEE Trans. on Communications*, vol. 57, no. 8, pp. 2339–2347, Aug. 2009.