

# PROPAGATION OF UNCERTAINTY IN BAYESIAN KERNEL MODELS — APPLICATION TO MULTIPLE-STEP AHEAD FORECASTING

Joaquin Quiñonero Candela<sup>(1)</sup>, Agathe Girard<sup>(2)</sup>, Jan Larsen<sup>(1)</sup> & Carl Edward Rasmussen<sup>(3)</sup>

<sup>(1)</sup>Informatics and Math. Modelling  
Technical University of Denmark  
Richard Petersens Plads, Build. 321  
2800 Kongens Lyngby, Denmark  
{jqc,jl}@imm.dtu.dk

<sup>(2)</sup>Dept. of Computing Science  
Glasgow University  
17 Lilybank Gardens  
Glasgow G12 8QQ, Scotland  
agathe@dcs.gla.ac.uk

<sup>(3)</sup>Biological Cybernetics  
Max Planck Institute  
Spemannstraße 38  
72076 Tübingen, Germany  
carl@tuebingen.mpg.de

## ABSTRACT

The object of Bayesian modelling is the predictive distribution, which in a forecasting scenario enables evaluation of forecasted values and their uncertainties. In this paper we focus on reliably estimating the predictive mean and variance of forecasted values using Bayesian kernel based models such as the Gaussian Process and the Relevance Vector Machine. We derive novel analytic expressions for the predictive mean and variance for Gaussian kernel shapes under the assumption of a Gaussian input distribution in the static case, and of a recursive Gaussian predictive density in iterative forecasting. The capability of the method is demonstrated for forecasting of time-series and compared to approximate methods.

## 1. INTRODUCTION

The problem of nonlinear forecasting is relevant to numerous application domains e.g. in financial modelling and control. This paper focuses on providing better estimates of the forecasted value as well as its uncertainty. The object of interest in Bayesian modelling framework [1] is the predictive density which contains all information about the forecasted value given the history of known values. For many Bayesian models the predictive density can only be approximated using Monte-Carlo sampling, local expansions, or variational approaches. However, when using Bayesian Gaussian shaped kernel models such as the Gaussian Process (GP) with a Gaussian kernel [1, 2] or the Relevance Vector Machine (RVM) [3, 4] the predictive mean and variance are given by analytic expressions under mild assumptions. Moreover the Bayesian kernel methods have proven to be very efficient nonlinear models [2, 4], with flexible approximation capabilities and high generalization performance.

We focus on the nonlinear auto-regressive (NAR) model with Gaussian innovations although more flexible nonlinear time-series models [5] sometimes are more efficient. Multi-step ahead forecasting can be done as direct forecast or as iterative one-step ahead forecasting. In [6] it is concluded that iterative forecasting usually is superior to direct forecasting. Generally the complexity of the nonlinear mapping in direct forecasting increases with the forecast horizon and for a fixed length time-series the number of training

examples decreases with the forecast horizon. In iterative forecasting the complexity of the nonlinear mapping is much lower than in the direct case, the number of training samples higher, but the performance is diminished by the uncertainty of the forecasted values. Consequently the involved effects provide a delicate trade-off. We restrict this work to iterative forecasting, which offers the additional advantage that multi-step ahead forecasts can be obtained with only one trained model.

In classical iterative forecasting only the predictive mean is iterated, here we consider an improvement to the methods suggested in [7] which iterate both the predictive mean and variance. This corresponds to using the model in recall/test phase under uncertain input. We do not consider training the model under uncertain inputs, which has been addressed for nonlinear model in [8] and for linear models in e.g. [9].

In section 2 we introduce the Bayesian modelling framework. In section 3 we consider the evaluation of the prediction density with uncertain inputs, which is formulated for time-series forecasting in section 4. Finally section 5 provides numerical experiments, that demonstrate the capability of the proposed method.

## 2. BAYESIAN KERNEL MODELLING

Consider a  $D$ -dimensional column input vector  $\mathbf{x}$  and a single output  $y$ , then the nonlinear model is defined as<sup>1</sup>

$$y = f(\mathbf{x}) + \varepsilon, \quad (1)$$

where  $f(\cdot)$  is a nonlinear function implemented as a GP or a RVM, and  $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$  is additive i.i.d. Gaussian noise with variance  $\sigma_\varepsilon^2$ . Suppose that the training data set is  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ , where  $N$  is the number of training samples. When using a GP [1, 2] or a RVM [3, 4], the predictive distribution of the output,  $y$ , is Gaussian [10],

$$p(y|\mathbf{x}, \mathcal{D}) = \mathcal{N}(\mu(\mathbf{x}), \sigma^2(\mathbf{x})), \quad (2)$$

where  $\mathbf{x}$  is an arbitrary input at which we perform prediction. For a GP the mean and variance of the predictive distribution are given by

$$\mu(\mathbf{x}) = \mathbf{k}^\top(\mathbf{x})\mathbf{K}^{-1}\mathbf{y}, \quad \sigma_{\text{GP}}^2(\mathbf{x}) = 1 - \mathbf{k}^\top(\mathbf{x})\mathbf{K}^{-1}\mathbf{k}(\mathbf{x}), \quad (3)$$

<sup>1</sup>We tacitly assume that  $y$  has zero mean, although a bias term can be included, see further [10].

This work is supported by the Multi-Agent Control Research Training Network - EC TMR grant HPRN-CT-1999-00107. Roderick Murray-Smith is acknowledged for useful discussions.

where  $C_{\text{GP}}(\mathbf{x}_i, \mathbf{x}_j)$  is the kernel, which we set to the commonly used Gaussian form<sup>2</sup>. We have

$$C_{\text{GP}}(\mathbf{x}_i, \mathbf{x}_j) = \exp[-(\mathbf{x}_i - \mathbf{x}_j)^\top \mathbf{\Lambda}^{-1}(\mathbf{x}_i - \mathbf{x}_j)/2], \quad (4)$$

$$\mathbf{\Lambda} = \text{diag}[\lambda_1^2, \dots, \lambda_D^2], \quad (5)$$

$$\mathbf{K} = \{K_{ij}\} = \{C_{\text{GP}}(\mathbf{x}_i, \mathbf{x}_j) + \sigma_\varepsilon^2 \delta_{ij}\} \quad (6)$$

$$\mathbf{k}(\mathbf{x}) = [C_{\text{GP}}(\mathbf{x}, \mathbf{x}_1), \dots, C_{\text{GP}}(\mathbf{x}, \mathbf{x}_N)]^\top, \quad (7)$$

$$\mathbf{y} = [y_1, \dots, y_N]^\top. \quad (8)$$

The kernel width hyper-parameters,  $\lambda_p$ , are fitted by maximizing the evidence (ML-II) using conjugate gradient, see e.g. [2].

For the RVM, let  $\{\phi_j(\mathbf{x})\}$  and  $\{\alpha_j\}$  with  $j = 1, 2, \dots, M$  be respectively the basis functions and the weight hyper-parameters, where  $M$  is the number of relevance vectors. Since typically  $M < N$ , the RVM yields sparse kernels, spanned by a finite number of basis functions [3, 10]. For the RVM the predictive distribution (2) has mean and variance specified by

$$\mu(\mathbf{x}) = \phi^\top(\mathbf{x})\mathbf{w}_{\text{MP}}, \quad \sigma_{\text{RVM}}^2(\mathbf{x}) = \phi^\top(\mathbf{x})\mathbf{\Sigma}^{-1}\phi(\mathbf{x}), \quad (9)$$

where, choosing Gaussian basis functions, we have

$$\mathbf{w}_{\text{MP}} = \sigma_\varepsilon^{-2}\mathbf{\Sigma}\mathbf{\Phi}^\top\mathbf{y}, \quad (10)$$

$$\mathbf{\Sigma} = (\sigma_\varepsilon^{-2}\mathbf{\Phi}^\top\mathbf{\Phi} + \mathbf{A})^{-1}, \quad (11)$$

$$\mathbf{A} = \text{diag}[\alpha_1, \dots, \alpha_M], \quad (12)$$

$$\phi_j(\mathbf{x}) = \exp[-(\mathbf{x}_j - \mathbf{x})^\top \mathbf{\Lambda}^{-1}(\mathbf{x}_j - \mathbf{x})/2], \quad (13)$$

$$\phi(\mathbf{x}) = [\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^\top, \quad (14)$$

$$\mathbf{\Phi} = \{\Phi_{ij}\} = \{\phi_j(\mathbf{x}_i)\}, \quad i = [1; N], \quad j = [1; M]. \quad (15)$$

The details of training the RVM are described in [3, 4].

### 3. PREDICTION WITH UNCERTAIN INPUT

Assume that the test input  $\mathbf{x}$  can not be observed directly and the uncertainty is modeled as  $\mathbf{x} \sim p(\mathbf{x}) = \mathcal{N}(\mathbf{u}, \mathbf{S})$ , with mean  $\mathbf{u}$  and covariance matrix  $\mathbf{S}$ . The resulting predictive distribution is then obtained by marginalizing over the test input

$$p(y|\mathbf{u}, \mathbf{S}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathcal{D})p(\mathbf{x}) d\mathbf{x}. \quad (16)$$

The principle is shown in Figure 1. The marginalization can in most cases only be carried out using Monte-Carlo numerical approximation techniques, however, in the case of Gaussian kernels<sup>3</sup> it is possible to obtain exact analytical expressions for the mean and variance of the marginalized predictive distribution:

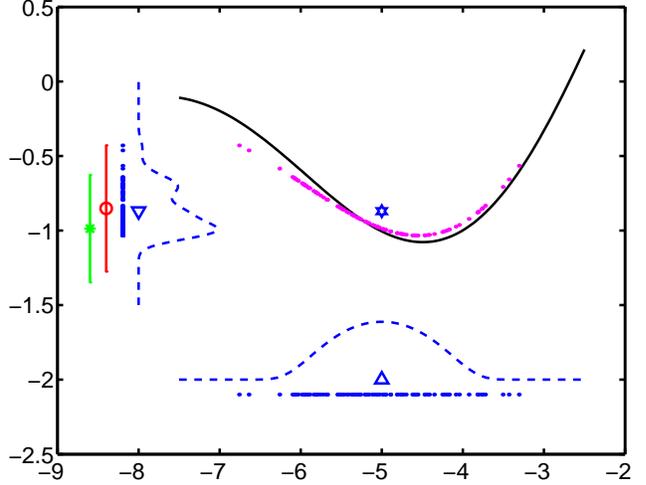
$$m(\mathbf{u}, \mathbf{S}) = \int y \cdot p(y|\mathbf{u}, \mathbf{S}, \mathcal{D}) dy, \quad \text{and} \quad (17)$$

$$v(\mathbf{u}, \mathbf{S}) = \int (y - m(\mathbf{u}, \mathbf{S}))^2 p(y|\mathbf{u}, \mathbf{S}, \mathcal{D}) dy. \quad (18)$$

The proposed method is an extension of the work presented in [7], which makes additional approximations, viz. Taylor series expansions of  $\mu(\mathbf{x})$  and  $\sigma^2(\mathbf{x})$  to first and second order around  $\mathbf{u}$  and

<sup>2</sup>The exponential in equation (4) is usually multiplied by an additional hyperparameter whose value is fitted during training. We here set it to 1 for clarity, which requires normalizing the data to unit variance.

<sup>3</sup>Exact analytical results can also be obtained for polynomial kernels,  $C(\mathbf{x}_p, \mathbf{x}_q) \propto |\mathbf{x}_p - \mathbf{x}_q|^r$ , e.g. a linear model.



**Fig. 1.** Prediction with uncertain input. *On the x-axis*, the dashed line represents the Gaussian input distribution, with mean located by the triangle, from which we draw 100 samples (dots under it). *In the middle of the figure*, the solid line represents the true underlying function. We fit a model to it, and propagate the 100 input samples through the model (dots close to the true function). *On the y-axis* we project the 100 predicted values (dots) and use them to estimate the predictive density (dashed line), with mean located by the triangle. The error bar with a circle and the error bar with a star show the mean and 95% confidence interval of the Gaussian approximation with exact computation of mean and variance and of the method with Taylor expansion respectively.

*S.* Using properties of the conditional mean and variance

$$m(\mathbf{u}, \mathbf{S}) = E_x[E_y[y|\mathbf{x}]] = E_x[\mu(\mathbf{x})], \quad (19)$$

$$v(\mathbf{u}, \mathbf{S}) = E_x[V_y[y|\mathbf{x}]] + V_x[E_y[y|\mathbf{x}]] \\ = E_x[\sigma^2(\mathbf{x})] + V_x[\mu(\mathbf{x})], \quad (20)$$

where  $E_x[\cdot]$ ,  $V_x[\cdot]$  denote the expectation and variance wrt.  $\mathbf{x}$ . When using Gaussian kernels in GPs and Gaussian basis functions in RVMs, the expressions for  $\mu(\mathbf{x})$  in eq. (3) and (9) are Gaussian shaped functions of  $\mathbf{x}$  and the expressions for  $\sigma^2(\mathbf{x})$  are products of Gaussian shaped functions in  $\mathbf{x}$ . Therefore the integrands involved in determining  $m(\mathbf{u}, \mathbf{S})$  and  $v(\mathbf{u}, \mathbf{S})$  are products of Gaussian shaped functions, which allows an analytical calculation. In [10] it is shown that

$$m(\mathbf{u}, \mathbf{S}) = \beta^\top \mathbf{l}. \quad (21)$$

For the GP  $\beta = \{\beta_1, \dots, \beta_N\} = \mathbf{K}^{-1}\mathbf{y}$  and for the RVM  $\beta = \{\beta_1, \dots, \beta_M\} = \mathbf{w}_{\text{MP}}$ . Vector  $\mathbf{l} = \{l_1, \dots, l_N\}$  is given by

$$l_j = |\mathbf{\Lambda}^{-1}\mathbf{S} + \mathbf{I}|^{-\frac{1}{2}} \\ \cdot \exp\left[-\frac{1}{2}(\mathbf{u} - \mathbf{x}_j)^\top (\mathbf{\Lambda} + \mathbf{S})^{-1}(\mathbf{u} - \mathbf{x}_j)\right], \quad (22)$$

where  $\mathbf{I}$  is the identity matrix. Note that if  $\mathbf{S}$  is the zero matrix, then  $\mathbf{l} = \mathbf{k}(\mathbf{u})$  and  $m(\mathbf{u}, \mathbf{S}) = \mu(\mathbf{u})$  as would be expected.

Further, for the GP

$$v(\mathbf{u}, \mathbf{S}) = \sigma_{\text{GP}}^2(\mathbf{u}) + \text{Tr}\left(\tilde{\mathbf{L}}(\beta\beta^\top - \mathbf{K}^{-1})\right) \\ + \text{Tr}\left((\mathbf{k}(\mathbf{u})\mathbf{k}(\mathbf{u})^\top - \mathbf{l}\mathbf{l}^\top)\beta\beta^\top\right), \quad (23)$$

where  $\tilde{\mathbf{L}} = \mathbf{L} - \mathbf{k}(\mathbf{u})\mathbf{k}(\mathbf{u})^\top$  and the elements of matrix  $\mathbf{L}$  are

$$L_{ij} = \mathbf{k}_i(\mathbf{u})\mathbf{k}_j(\mathbf{u}) \cdot |2\boldsymbol{\Lambda}^{-1}\mathbf{S} + \mathbf{I}|^{-\frac{1}{2}} \cdot \exp\left[2(\mathbf{u} - \mathbf{x}_d)^\top \boldsymbol{\Lambda}^{-1}(2\boldsymbol{\Lambda}^{-1} + \mathbf{S}^{-1})^{-1}\boldsymbol{\Lambda}^{-1}(\mathbf{u} - \mathbf{x}_d)\right], \quad (24)$$

and where  $\mathbf{x}_d = \frac{1}{2}(\mathbf{x}_i + \mathbf{x}_j)$ . For the RVM

$$v(\mathbf{u}, \mathbf{S}) = \sigma_{\text{RVM}}^2(\mathbf{u}) + \text{Tr}\left(\tilde{\mathbf{L}}(\boldsymbol{\beta}\boldsymbol{\beta}^\top + \boldsymbol{\Sigma}^{-1})\right) + \text{Tr}\left([\mathbf{k}(\mathbf{u})\mathbf{k}(\mathbf{u})^\top - \mathbf{U}^\top]\boldsymbol{\beta}\boldsymbol{\beta}^\top\right). \quad (25)$$

Notice that both for GPs and RVMs, when  $\mathbf{S}$  is the zero matrix,  $\tilde{\mathbf{L}}$  is also the zero matrix, again  $\mathbf{l} = \mathbf{k}(\mathbf{u})$ , and  $v(\mathbf{u}, \mathbf{S}) = \sigma^2(\mathbf{u})$ .

#### 4. APPLICATION TO TIME-SERIES FORECASTING

Suppose that  $\{y_t\}$  are the ordered samples of a time-series, where  $t$  is an integer time index. We wish to make time-series forecasting using a NAR model (1), where the inputs are formed by a collection of previous output values,  $\mathbf{x}_t = [y_{t-1}, y_{t-2}, \dots, y_{t-L}]$ , where the integer  $L$  is the size of the lag space.

Given that we have observed the values  $y^T \equiv \{y_t\}_{t=1}^T$ ,  $T$  being the number of observed samples, computing the predictive density of the value  $y_{T+1}$  is readily given by the model from (2) as

$$p(y_{T+1}|\mathbf{x}_{T+1}) = \mathcal{N}(\mu(\mathbf{x}_{T+1}), \sigma^2(\mathbf{x}_{T+1}))$$

The predictive density of the value  $y_{T+2}$  (two steps ahead) depends on  $\mathbf{x}_{T+1}$ , which now contains a stochastic element. In general, the predictive distribution of  $y_{T+k}$ , with  $k \geq 2$ , requires integrating out the uncertainty of the input:

$$p(y_{T+k}|y^T) = \int p(y_{T+k}|\mathbf{x}_{T+k})p(\mathbf{x}_{T+k}|y^T)d\mathbf{x}_{T+k}. \quad (26)$$

It is straightforward that this scheme leads to a recursive density estimation. The integral in (26) has no analytical solution. A naïve approach to the recursion is to ignore the uncertainty in the distribution of the input by setting  $p(\mathbf{x}_{T+k}) = \delta(\mathbf{x} - [\mu(\mathbf{x}_{T+k-1}), \dots, \mu(\mathbf{x}_{T+k-L})]^\top)$ <sup>4</sup>, thus propagating only the mean predictions. This method yields very poor error-bars, since it in some way only considers one step ahead predictions, treating the previous predicted values as exact, and is therefore overconfident, [7]. Alternatively, one can approximate the predictive density of  $y_{T+k}$  by a Gaussian density and compute only the mean and variance of  $p(y_{T+k}|y^T)$ . By doing this one ensures that the input distribution  $p(\mathbf{x}_{T+k}|y^T)$  is always Gaussian, which allows to use the results described in section 3 for computing the mean and variance of  $y_{T+k}$ , see eq. (26). This can be done exactly (for Gaussian or polynomial kernels) or in an approximate fashion, [7]. The recursive mechanism works because the predictive distribution of  $y_{T+1}$  at the first step is Gaussian (26), and therefore the input distribution of  $\mathbf{x}_{T+1}$  is also Gaussian. We call this procedure of recursively approximating the predictive density by a Gaussian the Recursive Gaussian Predictive Density (RGPD), and distinguish between exact-RGPD for the case of exact computation of mean and variance and approximate-RGPD for the case where the model is approximated by a Taylor expansion, [7].

<sup>4</sup>Where  $\delta(x)$  is 1 for  $x = 0$  and 0 otherwise. If  $k < L$ , we have simply  $\mu(\mathbf{x}_n) = y_n$  for  $n < T$ .

In the RGPD scheme, the input distribution is given by<sup>5</sup>

$$p(\mathbf{x}_{T+k}|y^T) = \mathcal{N}(\mathbf{u}_{T+k}, \mathbf{S}_{T+k}), \quad (27)$$

where

$$\mathbf{u}_{T+k} = [m(\mathbf{u}_{T+k-1}, \mathbf{S}_{T+k-1}), \dots, m(\mathbf{u}_{T+k-L}, \mathbf{S}_{T+k-L})],$$

and where  $\mathbf{S}_{T+k}$  is iteratively computed by using the fact that its first column is given by

$$(\mathbf{S}_{T+k})_{1:L,1} = \text{cov}(y_{T+k}, \mathbf{x}_{T+k}) = \sum_j \beta_j L_j (\mathbf{c}_j - \mathbf{u}_{T+k}), \quad (28)$$

where  $\mathbf{c}_j = (\boldsymbol{\Lambda}^{-1} + \mathbf{S}^{-1})^{-1}(\boldsymbol{\Lambda}^{-1}\mathbf{x}_j + \mathbf{S}^{-1}\mathbf{u})$ , refer to [10].

#### 5. EXPERIMENTS

We examine the comparative performance of the exact and approximate-RGPD on a hard prediction problem, the Mackey-Glass chaotic time series [11], which is well-known for its strong non-linearity. In [4] we showed that non-linear models, in particular RVMs, have a prediction error four orders of magnitude lower than optimized linear models. The Mackey-Glass attractor is a non-linear chaotic system described by the following equation:

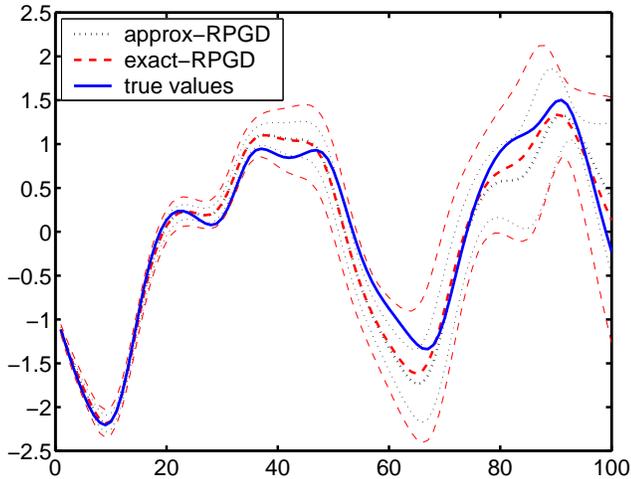
$$\frac{dz(t)}{dt} = -bz(t) + a \frac{z(t-\tau)}{1+z(t-\tau)^{10}} \quad (29)$$

where the constants are set to  $a = 0.2$ ,  $b = 0.1$  and  $\tau = 17$ . The series is re-sampled with period 1 according to standard practice. The inputs are formed by  $L = 16$  samples spaced 1 periods from each other  $\mathbf{x}_k = [z_{k-1}, z_{k-2}, \dots, z_{k-L}]$  and the targets are chosen to be  $y_k = z_k$ .

We train a GP model with Gaussian kernel on only 100 examples — enough to obtain a 1-step ahead normalized mean squared error on the order of  $10^{-4}$ . Besides, we normalize the data and contaminate it with a small amount of Gaussian noise with variance  $10^{-3}$ . Figure 2 shows the result of making 100 iterative predictions using a GP model, both for the exact-RGPD and the approximate-RGPD methods. By informal visual inspection, the error-bars of the exact-RGPD seem to be better than those of the approximate-RGPD. Consequently the exact-RGPD produces a better predictive density, which we show in figure 3. The mean value of the predictions seems also to be a slightly closer to the true target values for the exact-RGPD than for the approximate-RGPD.

In order to better evaluate the performance of the proposed methods, for a given prediction horizon, we compute the negative log predictive density, the squared error and the absolute error. While the two last measures only take into consideration the mean of the Gaussian predictive distribution, the first one also takes into account its variance. We average over 200 repetitions with different starting points (chosen at random from the series), and represent averages of the three loss measures for prediction horizons ranging from 1 to 100. Figure 3 shows the results. The means are slightly better for the exact-RGPD, but the predictive distribution is much improved. The better error-bars obtained by the exact-RGPD result in a lower value of the negative log predictive density for all values of the prediction horizon. The performance of the naïve iterative method is identical to that of the approximate-RGPD in terms of absolute and squared error. In terms of predictive density (since it produces unrealistic small error-bars) its performance is so poor that it is not worth reporting.

<sup>5</sup>If  $k < L$ , we have simply  $\mu(\mathbf{x}_n) = y_n$  for  $n < T$ .



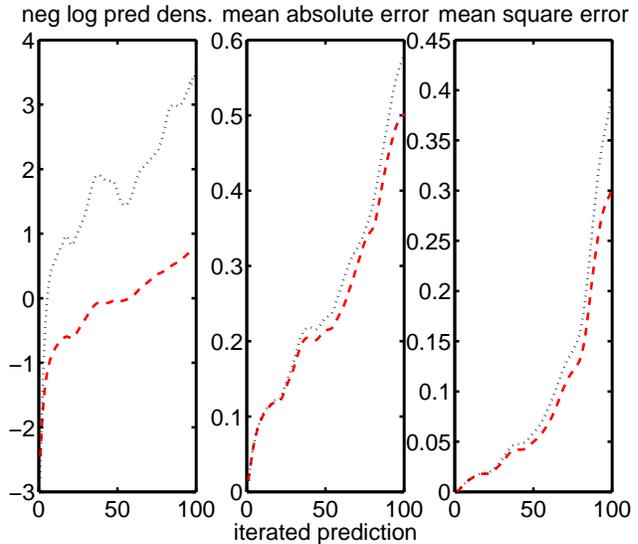
**Fig. 2.** 100 iterated predictions for the exact-RPGD (dashed) and approximate-RPGD (dotted): for each the thicker lines represent the mean of the predictive distributions and the two thinner lines around represent the upper and lower bounds of the 95% confidence interval of the Gaussian predictive distributions. The solid line shows the true target values.

## 6. CONCLUSIONS

We have derived analytical expressions for the exact computation of the mean and variance of the marginalized predictive distribution for uncertain Gaussian test inputs. This analytical expressions are valid for Gaussian processes and the Relevance Vector Machine (extended linear models) with Gaussian or polynomial kernels or basis functions. Our results extend the approximate method presented in [7], where the mean prediction was unaffected by the input uncertainty. In our case the input uncertainty biases the mean prediction, by smoothing, which is interesting in itself for predictions on noisy inputs. Furthermore, in the context of iterated time-series forecasting, our exact-RPGD not only gives much better error-bars, but the mean predictions are closer to the true values, both in terms of absolute and squared error. We are currently investigating efficient Monte Carlo methods to avoid the Gaussian approximation of the recursive predictive density.

## 7. REFERENCES

- [1] Radford M. Neal, *Bayesian Learning for Neural Networks*, Lecture Notes in Statistics, no. 118. Springer, New York, 1996.
- [2] Carl E. Rasmussen, *Evaluation of Gaussian Processes and Other Methods for Non-linear Regression*, Ph.D. thesis, Dept. of Computer Science, University of Toronto, 1996.
- [3] Michael E. Tipping, “Sparse bayesian learning and the relevance vector machine,” *Journal of Machine Learning Research*, vol. 1, pp. 211–244, 2001.
- [4] Joaquin Quiñero-Candela and Lars Kai Hansen, “Time series prediction based on the relevance vector machine with adaptive kernels,” in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2002, pp. 985–988.



**Fig. 3.** Negative log predictive density, mean absolute error and mean squared error as a function of the iterative prediction horizon for the exact-RPGD method (dashed) and for the approximate-RPGD (dotted). Averages over 200 repetitions.

- [5] I.J. Leontaritis and S.A. Billings, “Input-output parametric models for non-linear systems, part 1: Deterministic non-linear systems, part 2: Stochastic non-linear systems,” *International Journal of Control*, vol. 41, pp. 303–344, 1985.
- [6] J. Doyne Farmer and John J. Sidorowich, “Exploiting chaos to predict the future and reduce noise,” Tech. Rep. LA-UR-88, Los Alamos National Laboratory, 1988.
- [7] Agathe Girard, Carl Edward Rasmussen, and Roderick Murray-Smith, “Gaussian process with uncertain input - application to multiple-step ahead time-series forecasting,” in *Advances in Neural Information Processing Systems 15*, Suzanna Becker, Sebastian Thrun, and Klaus Obermayer, Eds. 2003, MIT Press.
- [8] Volker Tresp, Subutai Ahmad, and Ralph Neuneier, “Training neural networks with deficient data,” in *Advances in Neural Information Processing Systems 6*, Jack D. Cowan, Gerald Tesauero, and Joshua Alspector, Eds. 1994, pp. 128–135, Morgan Kaufmann Publishers, Inc.
- [9] Oscar Netares, David J. Fleet, and David J. Heeger, “Likelihood functions and confidence bounds for total-least-squares problems,” in *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition*, 2000, vol. 1, pp. 523–530.
- [10] Joaquin Quiñero-Candela and Agathe Girard, “Prediction at an uncertain input for gaussian processes and relevance vector machines - application to multiple-step ahead time-series forecasting,” Tech. Rep., IMM, DTU, 2002, [http://isp.imm.dtu.dk/staff/jqc/prop\\_uncert.ps.gz](http://isp.imm.dtu.dk/staff/jqc/prop_uncert.ps.gz).
- [11] M.C. Mackey and L. Glass, “Oscillation and chaos in physiological control systems,” *Science*, vol. 197, pp. 287–289, July 1977.