

# Recursive Estimation for Sparse Gaussian Process Regression

Manuel Schürch<sup>a,b</sup>, Dario Azzimonti<sup>a</sup>, Alessio Benavoli<sup>c,a</sup>, Marco Zaffalon<sup>a</sup>

<sup>a</sup>*Istituto Dalle Molle di Studi sull'Intelligenza Artificiale (IDSIA); Manno, Switzerland*

<sup>b</sup>*Università della Svizzera italiana (USI); Lugano, Switzerland*

<sup>c</sup>*University of Limerick (UL); Limerick, Ireland*

---

## Abstract

Gaussian Processes (GPs) are powerful kernelized methods for non-parameteric regression used in many applications. However, their use is limited to a few thousand of training samples due to their cubic time complexity. In order to scale GPs to larger datasets, several sparse approximations based on so-called inducing points have been proposed in the literature. In this work we investigate the connection between a general class of sparse inducing point GP regression methods and Bayesian recursive estimation which enables Kalman Filter like updating for online learning. The majority of previous work has focused on the batch setting, in particular for learning the model parameters and the position of the inducing points, here instead we focus on training with mini-batches. By exploiting the Kalman filter formulation, we propose a novel approach that estimates such parameters by recursively propagating the analytical gradients of the posterior over mini-batches of the data. Compared to state of the art methods, our method keeps analytic updates for the mean and covariance of the posterior, thus reducing drastically the size of the optimization problem. We show that our method achieves faster convergence and superior performance compared to state of the art sequential Gaussian Process regression on synthetic GP as well as real-world data with up to a million of data samples.

*Key words:* Gaussian processes; Recursive estimation; Kalman filter; Non-parametric regression; Parameter estimation.

---

## 1 Introduction

*Gaussian process* (GPs) regression is used in many applications, ranging from machine learning, social sciences, natural sciences and engineering, due to its modeling flexibility, robustness to overfitting and availability of well-calibrated predictive uncertainty estimates. In control engineering, for example, GPs have been used in system identification for impulse response estimation [25,9,26,24], nonlinear ARX models [19,3], learning of ODEs [1,21], latent force modeling [42] and to learn the state space of a nonlinear dynamical system [12,23,35]. However, GPs do not scale to large data sets due to their  $\mathcal{O}(N^2)$  memory and  $\mathcal{O}(N^3)$  computational costs, where  $N$  is the number of training samples. For this reason several sparse GP approximations have been proposed in the literature. Often such approximations are based on *inducing points* methods, where the unknown function is represented by its values at a set of  $M \ll N$  pseudo-inputs, called inducing points. Among such methods, *Subset of Regressors* (SoR/DIC) approximations

[32,39,33] produces overconfident predictions when leaving the training data. On the other hand, *Deterministic Training Conditional* (DTC) [11,31], *Fully Independent Training Conditional* (FITC) [34], *Fully Independent Conditional* (FIC) [27] and *Partially Independent Training Conditional* (PITC) [27] all produce sensible uncertainty estimates. These models differ from each other in the definition of their joint prior over the latent function and test values. Titsias [36], instead, proposed to retain the exact prior but to perform approximate (variational) inference for the posterior, leading to the *Variational Free Energy* (VFE) method which converges to full GP as  $M$  increases. Bui et al. [7] introduced *Power Expectation Propagation* (PEP), based on the minimization of an  $\alpha$ -divergence, which unifies most of the previously mentioned models. Typically, inference is achieved in  $\mathcal{O}(M^2N)$  time and  $\mathcal{O}(MN)$  space. In order to find good parameters (inducing input points and kernel hyper-parameters), either the log marginal likelihood of the sparse models or a lower bound are numerically optimized. The previously mentioned approximations focus on the batch setting, i.e., all data is available at once and can be processed together. For big data, where the number of samples can be millions, keeping all data in memory is not possible, moreover the data might even arrive sequentially. Bui et al. [6] developed an algorithm to update hyper-parameters in

---

*Email addresses:* manuel@idsia.ch (Manuel Schürch),  
dario.azzimonti@idsia.ch (Dario Azzimonti),  
alessio.benavoli@ul.ie (Alessio Benavoli),  
zaffalon@idsia.ch (Marco Zaffalon).

an online fashion promising in a streaming setting, but with limited accuracy as each sample is considered only once.

We focus here on the setting where hyper-parameters are learned by reconsidering mini-batches several times. In order to speed up the optimization, we would like to update the parameters more frequently for a subset of data and update the posterior in a sequential way. In this setting, Hensman et al. [15] applied *Stochastic Variational Inference* (SVI, [17]) to an *uncollapsed* lower bound of the marginal likelihood. The resulting *Stochastic Variational Gaussian Process* (SVGP) method allows to optimize the parameters with mini-batches. Although showing high scalability and good accuracy, SVGP has two main drawbacks: i) the (variational) posterior is not given analytically, which leads to  $\mathcal{O}(M^2)$  additional many parameters; ii) the uncollapsed bounds are in practice often less tight than the corresponding collapsed VFE batch bounds because the (variational) posterior is not optimally eliminated. The large number of parameters ( $\approx MD + M^2$ , where  $D$  is the input space dimension) leads to a hard-to-tune optimization problem which requires appropriately decaying learning rates. Even for fixed parameters, each sample still needs to be reconsidered many times. An orthogonal direction was pursued by the authors in [14] and [30], where a connection between GPs and State Space models for particular kernels was established for spatio-temporal regression problems, which allows to apply sequential algorithm such as the Kalman Filter, see also [8,2]. Inspired by this line of research, the authors in [37] focused on efficient implementation and extended the methodology to varying sampling locations over time. These approaches can deal with sequential data and solve the problem of temporal time complexity, however the space complexity is still cubic in  $N$ . In addition, the hyper-parameters are usually fixed in advance.

In this work we propose a *recursive collapsed* lower bound to the log marginal likelihood which can be optimized stochastically with mini-batches.

In this respect, the first contribution of this paper is the derivation of a *novel Kalman-filter-like* (KF) formulation for a generic sparse inducing point method. In particular we show that sparse inducing point models can be seen as a Bayesian kernelized linear regression model with input dependent observation noise, a particular choice of basis functions and noise covariance. Given the model hyper-parameters, KF allows to train sparse GP methods analytically and exactly in an online setting (considering each sample only once, as opposed to the work in [15] and [16]). In this formulation the posterior distribution obtained online is equivalent to full batch methods. This constitutes an interesting technique on its own for applications where hyper-parameters are given, however the analysis above provides a key insight for parameter estimation.

Our second main contribution is a *recursive approach to hyper-parameter estimation* based on the KF formulation. It is based on recursively exploiting the chain rule for derivatives by recursively propagating the analytical gradients of the posterior which enables us to compute the derivatives of the lower bound sequentially. We show that, when com-

puting the gradients of the recursive collapsed bound in a non-stochastic way, they exactly match the corresponding batch ones. This new *Stochastic Recursive Gradient Propagation* (SRGP)<sup>1</sup> approach constitutes an efficient method to train a very general class of sparse GP regression models with much fewer parameters to be estimated numerically ( $\approx MD$ ) than state of the art sequential GP regression methods ( $\approx MD + M^2$ ). Since the number  $M$  of inducing points determines the quality of the approximation to full GP, this reduction in number of parameters from  $M^2$  to  $M$  is crucial and results in more accurate and faster convergence than state of the art approaches such as SVGP. For example, in the application to learn the input output behavior of a non-linear plant presented in Sect. 5 the number of parameters estimated by SVGP is  $\approx 10500$  while our approach only estimates  $\approx 500$  parameters due to the analytical updates.

## 2 Background on GP Regression

Consider a training set  $\mathcal{D} = \{y_i, \mathbf{x}_i\}_{i=1}^N$  of  $N$  pairs of inputs  $\mathbf{x}_i \in \mathbb{R}^D$  and noisy scalar outputs  $y_i$  generated by adding independent Gaussian noise to a latent function  $f(\mathbf{x})$ , that is  $y_i = f(\mathbf{x}_i) + \varepsilon_i$ , where  $\varepsilon_i \sim \mathcal{N}(0, \sigma_n^2)$ . We denote by  $\mathbf{y} = [y_1, \dots, y_N]^T$  the vector of observations and by  $\mathbf{X} = [\mathbf{x}_1^T, \dots, \mathbf{x}_N^T]^T \in \mathbb{R}^{N \times D}$  the input points.

We model  $f$  with a *Gaussian Process* (GP), a stochastic process defined by its mean function  $m(\mathbf{x})$  and covariance kernel  $k(\mathbf{x}, \mathbf{x}')$ . The kernel  $k$  is a positive definite function [28], such as, for instance, the *squared exponential* (SE) kernel with individual lengthscales  $l_i$  for each dimension, that is  $k(\mathbf{x}, \mathbf{x}') = \sigma_0^2 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^T \text{Diag}[l_1^2, \dots, l_D^2](\mathbf{x} - \mathbf{x}')\right)$ . We assume  $m(\mathbf{x}) \equiv 0$  for the sake of simplicity and we use the SE kernel throughout this paper however all methods work with any positive definite kernel. Given the training values  $\mathbf{f} = f(\mathbf{X}) = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)]^T$  and a test latent function value  $f_* = f(\mathbf{x}_*)$  at a test point  $\mathbf{x}_* \in \mathbb{R}^D$ , then the joint distribution  $p(\mathbf{f}, f_*)$  is Gaussian. Our likelihood is Gaussian,  $p(\mathbf{y}|\mathbf{f}) = \mathcal{N}(\mathbf{y}|\mathbf{f}, \sigma_n^2 \mathbb{I})$ , and with Bayes theorem (see e.g. [28]) we obtain analytically the posterior predictive distribution  $p(f_*|\mathbf{y}) = \mathcal{N}(f_*|\boldsymbol{\mu}_*, \boldsymbol{\Sigma}_*)$  with

$$\begin{aligned} \boldsymbol{\mu}_* &= \mathbf{K}_{*\mathbf{X}} (\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{I})^{-1} \mathbf{y}, \\ \boldsymbol{\Sigma}_* &= \mathbf{K}_{**} - \mathbf{K}_{*\mathbf{X}} (\mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{I})^{-1} \mathbf{K}_{\mathbf{X}*}, \end{aligned} \quad (1)$$

where  $[\mathbf{K}_{\mathbf{A}\mathbf{B}}]_{ij} = k(\mathbf{a}_i, \mathbf{b}_j)$  for any  $\mathbf{A} \in \mathbb{R}^{M_1 \times D}$  and  $\mathbf{B} \in \mathbb{R}^{M_2 \times D}$  with the corresponding rows  $\mathbf{a}_i, \mathbf{b}_j$ . For brevity, we use  $*$  to indicate  $\mathbf{x}_*$ . The GP depends via the kernel matrices on the hyper-parameters  $\boldsymbol{\phi} = \{\sigma_0, l_1, \dots, l_D, \sigma_n\}$  typically estimated by maximizing the log marginal likelihood

$$\log p(\mathbf{y}|\boldsymbol{\phi}) = \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{K}_{\mathbf{X}\mathbf{X}} + \sigma_n^2 \mathbb{I}). \quad (2)$$

Note that the computations for inference require the inversion of the matrix in Eq. (1) which scales as  $\mathcal{O}(N^3)$  in time and  $\mathcal{O}(N^2)$  for memory (given  $\boldsymbol{\phi}$ ).

<sup>1</sup> Code is available at <https://github.com/manuelIDSIA/SRGP>.

## 2.1 Batch Sparse GP Regression

Sparse GP regression methods based on *inducing points* approaches reduce the computational complexity by introducing  $M \ll N$  inducing points  $\mathbf{u} \in \mathbb{R}^M$  that optimally summarize the dependency of the whole training data. The inducing *inputs*  $\mathbf{R} \in \mathbb{R}^{M \times D}$  are in the  $D$ -dimensional input data space and the inducing *outputs*  $\mathbf{u} := f(\mathbf{R})$  are the corresponding GP-function values, see also Fig. 1. The GP prior over  $\mathbf{f}$  and  $f_*$  is augmented with the inducing outputs  $\mathbf{u}$ , leading to a joint  $p(\mathbf{f}, f_*, \mathbf{u})$  and marginal  $p(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \mathbf{K}_{RR})$  prior. By marginalizing out the inducing points, the original prior  $p(\mathbf{f}, f_*) = \int p(\mathbf{f}, f_*, \mathbf{u}) p(\mathbf{u}) d\mathbf{u}$  is recovered. The fundamental approximation in all sparse GP models is that given the inducing outputs  $\mathbf{u}$ ,  $\mathbf{f}$  and  $f_*$  are conditionally independent. Consequently, inference in these models can be done in  $\mathcal{O}(M^2N)$  time and  $\mathcal{O}(MN)$  space [34].

We briefly recall here the sparse predictive distribution and the variational lower bound to the log marginal likelihood for the *Power Expectation Propagation* (PEP) model [7] because it unifies the main sparse inducing points approaches. The variational lower bound is used for optimizing the parameters  $\theta := \{\phi, \mathbf{R}\}$ . In the following, we denote  $\mathbf{Q}_{AB} = \mathbf{K}_{AR} \mathbf{K}_{RR}^{-1} \mathbf{K}_{RB}$  and  $\mathbf{D}_A = \mathbf{K}_{AA} - \mathbf{Q}_{AA}$  for any  $A, B$ . The predictive distribution  $p(f_* | \mathbf{y}) = \mathcal{N}(f_* | \mu_*, \Sigma_*)$  of PEP is given by

$$\begin{aligned} \mu_* &= \mathbf{Q}_{*X} (\overline{\mathbf{K}}_{XX} + \sigma_n^2 \mathbb{I})^{-1} \mathbf{y}, \\ \Sigma_* &= \mathbf{K}_{**} - \mathbf{Q}_{*X} (\overline{\mathbf{K}}_{XX} + \sigma_n^2 \mathbb{I})^{-1} \mathbf{Q}_{X*}, \end{aligned} \quad (3)$$

where  $\overline{\mathbf{K}}_{XX} = \mathbf{Q}_{XX} + \alpha \text{Diag}[\mathbf{D}_X]$ . A lower bound to the sparse log marginal likelihood is analytically available

$$\begin{aligned} \mathcal{L}_{PEP}(\theta) &= \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \overline{\mathbf{K}}_{XX} + \sigma_n^2 \mathbb{I}) \\ &\quad - \frac{1-\alpha}{2\alpha} \sum_{i=1}^N \log \left( 1 + \frac{\alpha}{\sigma_n^2} [\mathbf{D}_X]_{ii} \right), \end{aligned} \quad (4)$$

where we omit the explicit dependency on  $\theta$  via  $\overline{\mathbf{K}}_{XX}$  and  $\mathbf{D}_X$  for the sake of brevity. This bound can be used to learn the parameters  $\theta$ , similarly to Eq. (2) for full GP.

The special case  $\alpha \rightarrow 0$  was originally introduced in [36] where the author proposed to maximize a variational lower bound to the true GP marginal likelihood, obtaining the *Variational Free Energy* (VFE) or the *collapsed lower bound*

$$\mathcal{L}_{VFE}(\theta) = \log \mathcal{N}(\mathbf{y} | \mathbf{0}, \mathbf{Q}_{XX} + \sigma_n^2 \mathbb{I}) - \frac{\text{Tr}[\mathbf{D}_X]}{2\sigma_n^2}. \quad (5)$$

In (5) the variational distribution over the inducing points is optimally eliminated and analytically available. The rightmost term in (5) acts as a regularizer that prevents overfitting and has the effect that the sparse GP predictive distribution (3) converges [36] to the exact GP predictive distribution (1) as the number of inducing points increases, when optimizing  $\theta$  with (5). See also [7,20,27,28] for recent reviews on the subject.

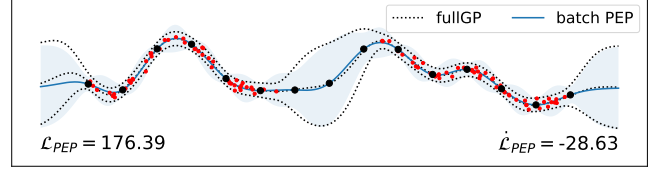


Fig. 1. Full GP and batch sparse GP regression with PEP model ( $\alpha = 0.5$ ).  $N = 100$  data samples are summarized with 15 equidistant inducing points (black dots). A slightly smaller than optimal lengthscale was selected and no parameters  $\theta$  were optimized. The numbers in the left and right corner indicate the lower bound to the log marginal likelihood in (4) and its derivative with respect to the lengthscale, respectively.

## 2.2 Sequential Sparse GP Regression

The optimization for  $\theta$  of the collapsed lower bound (5) requires to process the whole dataset, which is very inefficient and not feasible for large  $N$ . We would like to update the parameters more frequently, therefore, we split the data  $\mathcal{D} = \{\mathbf{y}_k, \mathbf{X}_k\}_{k=1}^K$  into  $K$  mini-batches of size  $B$  and denote  $\mathbf{f}_k$  the corresponding sparse GP value. *Stochastic Variational Gaussian Process* (SVGP) [15] achieves this result by applying stochastic optimization to an *uncollapsed lower bound* to the log marginal likelihood

$$\begin{aligned} \mathcal{L}_{SVGP}(\mu, \Sigma, \theta) &= -\text{KL}[q(\mathbf{u}) || p(\mathbf{u} | \theta)] \\ &\quad + K \sum_{k=1}^K \int q(\mathbf{u}) p(\mathbf{f}_k | \mathbf{u}, \theta) p(\mathbf{y}_k | \mathbf{f}_k, \theta) d\mathbf{u}, \end{aligned} \quad (6)$$

where the variational distribution  $q(\mathbf{u})$  is part of the bound and explicitly parametrized as  $q(\mathbf{u}) = \mathcal{N}(\mathbf{u} | \mu, \Sigma)$ . This uncollapsed bound satisfies  $\mathcal{L}_{SVGP}(\mu, \Sigma, \theta) \leq \mathcal{L}_{VFE}(\theta)$  with equality when inserting the optimal mean and covariance of the variational distribution of VFE. The key property of this bound is that it can be written as a sum of  $K$  terms, which allows *Stochastic Variational Inference* (SVI, [17]). Note that collapsing the bound, i.e. inserting the optimal distribution, reintroduces dependencies between the observations, and eliminates the global parameter  $\mathbf{u}$  which is needed for SVI. For this reason, all variational parameters are numerically estimated by following the noisy gradients of a stochastic estimate of the lower bound  $\mathcal{L}_{SVGP}$ . By passing through the training data a sufficient number of times, the variational distribution converges to the batch solution of VFE method. This approach, however, requires a large number of parameters: in addition to the parameters  $\theta$ , all entries in the mean vector  $\mu$  and the covariance matrix  $\Sigma$  have to be estimated numerically, which is in order  $\mathcal{O}(M^2)$ .

## 3 Recursive Sparse GP Regression

In this section we establish the connection between Bayesian recursive estimation and sparse inducing point GP models. We recall the *weight-space view* for a large class of sparse inducing point GP models, which we present here as a particular kernelized version of a Bayesian linear regression model. See [28, Ch. 2.1], for an analogous discussion on the full GP model. Here, however, we show how to ex-

plot the KF to train many sparse methods analytically either in an online setting for fixed hyper-parameters. This allows us to introduce a recursive log marginal likelihood with a model specific regularization term for parameter estimation.

### 3.1 Weight-Space View of Generic Sparse GP

For a mini-batch  $\mathbf{X} \in \mathbb{R}^{B \times D}$  of size  $B$ , consider the generic sparse GP model

$$f(\mathbf{X}) = H(\mathbf{X})\mathbf{u} + \gamma(\mathbf{X}) \quad (7)$$

where the sparse GP value  $f(\mathbf{X})$  is modeled by a linear combination of basis-functions  $H(\mathbf{X}) \in \mathbb{R}^{B \times M}$ , (stochastic) weights  $\mathbf{u} \in \mathbb{R}^M$  with a prior  $p(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \Sigma_0)$  and an input dependent error term  $\gamma(\mathbf{X}) \sim \mathcal{N}(0, V(\mathbf{X}))$  that takes into account the sparse approximation. For  $k = 1, \dots, K$ , the noisy observations  $\mathbf{y}_k$  are obtained by adding independent noise  $\varepsilon_k \sim \mathcal{N}(0, \sigma_n^2 \mathbb{I})$  to  $f(\mathbf{X}_k)$ , yielding the model

$$\mathbf{y}_k = \mathbf{f}_k + \varepsilon_k; \quad (8)$$

$$\mathbf{f}_k = \mathbf{H}_k \mathbf{u} + \gamma_k \quad \text{and} \quad \mathbf{f}_* = \mathbf{H}_* \mathbf{u} + \gamma_*, \quad (9)$$

where we distinguish the training  $\mathbf{f}_k = f(\mathbf{X}_k)$  and test  $\mathbf{f}_* = f(\mathbf{X}_*)$  cases depending on the input  $\mathbf{X}_k$  and  $\mathbf{X}_*$ . Assuming  $\gamma_k, \gamma_*$  and  $\varepsilon_k$  are independent, by linearity and Gaussianity we can compactly write

$$p(\mathbf{y}_k | \mathbf{f}_k) = \mathcal{N}(\mathbf{y}_k | \mathbf{f}_k, \sigma_n^2 \mathbb{I}); \quad (10)$$

$$p(\mathbf{u}) = \mathcal{N}(\mathbf{0}, \Sigma_0); \quad (11)$$

$$p(\mathbf{f}_k | \mathbf{u}) = \mathcal{N}(\mathbf{f}_k | \mathbf{H}_k \mathbf{u}, \bar{\mathbf{V}}_k); \quad (12)$$

$$p(\mathbf{f}_* | \mathbf{u}) = \mathcal{N}(\mathbf{f}_* | \mathbf{H}_* \mathbf{u}, \mathbf{V}_*), \quad (13)$$

Combining (10) and (12) and by integrating out  $\mathbf{f}_k$  we obtain the likelihood  $p(\mathbf{y}_k | \mathbf{u}) = \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \mathbf{V}_k)$  where  $\mathbf{V}_k = \bar{\mathbf{V}}_k + \sigma_n^2 \mathbb{I}$ . This shows that a generic sparse GP regression model can be seen as a Bayesian non-linear regression model with additional input dependent observation noise, a particular choice of basis functions  $\mathbf{H}_k$  and covariance structures  $\bar{\mathbf{V}}_k, \mathbf{V}_*$ . For inducing inputs  $\mathbf{R} \in \mathbb{R}^{M \times D}$ , we have  $\Sigma_0 = \mathbf{K}_{RR}$ ,  $\mathbf{H}_k = \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1}$  and  $\mathbf{H}_* = \mathbf{K}_{*R} \mathbf{K}_{RR}^{-1}$ . Different choices of the quantities  $\bar{\mathbf{V}}_k$  and  $\mathbf{V}_*$  lead to a range of sparse GP models summarized in the bottom table in Fig. 2.

### 3.2 Training

Given the prior  $p(\mathbf{u})$  and the likelihood  $p(\mathbf{y}_k | \mathbf{u})$ , the posterior over the weights  $\mathbf{u}$  conditioned on the data  $\mathbf{y}_{1:k}$  can be computed either in a batch or in a recursive manner.

#### 3.2.1 Batch Estimation

The batch likelihood is  $p(\mathbf{y} | \mathbf{u}) = \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{u}) = \mathcal{N}(\mathbf{y} | \mathbf{H}\mathbf{u}, \mathbf{V})$  with  $\mathbf{H} = [\mathbf{H}_1^T, \dots, \mathbf{H}_K^T]^T \in \mathbb{R}^{N \times M}$  and  $\mathbf{V}$  a block-diagonal matrix with blocks  $\mathbf{V}_k$ . The posterior over  $\mathbf{u}$  given the data  $\mathbf{y}$  can be obtained by Bayes' rule, i.e.

$$p(\mathbf{u} | \mathbf{y}) \propto p(\mathbf{y} | \mathbf{u}) p(\mathbf{u}) \propto \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}_K, \boldsymbol{\Sigma}_K), \quad (14)$$

with  $\boldsymbol{\Sigma}_K = (\boldsymbol{\Sigma}_0^{-1} + \mathbf{H}^T \mathbf{V}^{-1} \mathbf{H})^{-1}$  and  $\boldsymbol{\mu}_K = \boldsymbol{\Sigma}_K \mathbf{H}^T \mathbf{V}^{-1} \mathbf{y}$ .

a) parameters			b) transformed parameters		
$\Sigma_0$	$H_k$	$H_*$	$\tilde{\Sigma}_0$	$\tilde{H}_k$	$\tilde{H}_*$
$K_{RR}$	$K_{X_k R} K_{RR}^{-1}$	$K_{*R} K_{RR}^{-1}$	$K_{RR}^{-1}$	$K_{X_k R}$	$K_{*R}$

	I) training	II) prediction	III) optimization
	$\bar{\mathbf{V}}_k$	$\mathbf{V}_*$	$a_k$
DIC [31]	0	0	0
DTC [9]	0	$D_*$	0
FITC [32]	$Diag[D_{X_k}]$	$D_*$	0
FIC [26]	$Diag[D_{X_k}]$	$Diag[D_*]$	0
PITC [26]	$D_{X_k}$	$D_*$	0
VFE [34]	0	$D_*$	$\frac{1}{2\sigma_n^2} Tr[D_{X_k}]$
PEP [6]	$\alpha Diag[D_{X_k}]$	$D_*$	$\frac{1-\alpha}{2\alpha} \sum_i \log(1 + \frac{\alpha}{\sigma_n^2} [D_{X_k}]_{ii})$
PEPB	$\alpha D_{X_k}$	$D_*$	$\frac{1-\alpha}{2\alpha} \log  I + \frac{\alpha}{\sigma_n^2} D_{X_k} $

Fig. 2. Summary of parameters for sparse GP models for recursive estimation. For all models, we have  $\boldsymbol{\mu}_0 = \mathbf{0}$ ,  $\mathbf{V}_k = \bar{\mathbf{V}}_k + \sigma_n^2 \mathbb{I}$  and  $\Sigma_0, \mathbf{H}_k, \mathbf{H}_*$  from the table a) or a transformed version b). Using the model specific quantities for the observation noise  $\bar{\mathbf{V}}_k$ , the prediction covariance  $\mathbf{V}_*$  and the regularization term  $a_k$  from the bottom table allows the training with the recursive approaches.

#### 3.2.2 Recursive Estimation

An equivalent solution can be obtained by propagating recursively  $p(\mathbf{u} | \mathbf{y}_{1:k-1})$ . By interpreting this previous posterior as the prior, the updated posterior can be recursively computed by

$$p(\mathbf{u} | \mathbf{y}_{1:k}) \propto p(\mathbf{y}_k | \mathbf{u}) p(\mathbf{u} | \mathbf{y}_{1:k-1}) \propto \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k), \quad (15)$$

where  $\boldsymbol{\mu}_k = \boldsymbol{\Sigma}_k (\mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{y}_k + \boldsymbol{\Sigma}_{k-1}^{-1} \boldsymbol{\mu}_{k-1})$  and  $\boldsymbol{\Sigma}_k = (\boldsymbol{\Sigma}_{k-1}^{-1} + \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{H}_k)^{-1}$ .

**Kalman Filter like updating:** the KF constitutes an efficient way to update the mean and covariance of  $p(\mathbf{u} | \mathbf{y}_{1:k})$ . Applying the Woodbury identity to  $\boldsymbol{\Sigma}_k$  in Eq. (15) and introducing temporary variables yields

$$\begin{aligned} \mathbf{r}_k &= \mathbf{y}_k - \mathbf{H}_k \boldsymbol{\mu}_{k-1}; \\ \mathbf{S}_k &= \mathbf{H}_k \boldsymbol{\Sigma}_{k-1} \mathbf{H}_k^T + \mathbf{V}_k; & \boldsymbol{\mu}_k &= \boldsymbol{\mu}_{k-1} + \mathbf{G}_k \mathbf{r}_k; \\ \mathbf{G}_k &= \boldsymbol{\Sigma}_{k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1}; & \boldsymbol{\Sigma}_k &= \boldsymbol{\Sigma}_{k-1} - \mathbf{G}_k \mathbf{S}_k \mathbf{G}_k^T. \end{aligned} \quad (16)$$

Starting the recursion with  $\boldsymbol{\mu}_0 = \mathbf{0}$  and  $\boldsymbol{\Sigma}_0$ , the posterior distribution at step  $K$  is equivalent to (14) independent of the order of the data. We want to emphasize that the only difference in the estimation part between the sparse GP models is the form of the additional noise  $\mathbf{V}_k = \bar{\mathbf{V}}_k + \sigma_n^2 \mathbb{I}$ .

**Transformation:** instead of running a KF with  $\Sigma_0 = \mathbf{K}_{RR}$ ,  $\mathbf{H}_k = \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1}$  and  $\mathbf{H}_* = \mathbf{K}_{*R} \mathbf{K}_{RR}^{-1}$ , an equivalent predictive distribution is also obtained when using  $\tilde{\Sigma}_0 = \mathbf{K}_{RR}^{-1}$  and  $\tilde{\mathbf{H}}_k = \mathbf{K}_{X_k R}$  together with  $\tilde{\mathbf{H}}_* = \mathbf{K}_{*R}$ . For any  $k$  we then propagate a transformed posterior distribution  $\tilde{\boldsymbol{\mu}}_k = \mathbf{K}_{RR}^{-1} \boldsymbol{\mu}_k$ ,  $\tilde{\boldsymbol{\Sigma}}_k = \mathbf{K}_{RR}^{-1} \boldsymbol{\Sigma}_k \mathbf{K}_{RR}^{-1}$  and  $\boldsymbol{\mu}_k = \mathbf{K}_{RR} \tilde{\boldsymbol{\mu}}_k$ ,  $\boldsymbol{\Sigma}_k =$

$\mathbf{K}_{RR}\tilde{\Sigma}_k\mathbf{K}_{RR}$ , respectively. This parametrization constitutes a computational shortcut, since the basis functions are very easy to interpret and do not include any matrix multiplication. Note that also the log marginal likelihood discussed below is not affected by this transformation.

### 3.3 Prediction

Given a new  $\mathbf{X}_* \in \mathbb{R}^{A \times D}$ , the predictive distribution after seeing  $\mathbf{y}_{1:k}$  of the sparse GP methods can be obtained by

$$\begin{aligned} p(\mathbf{f}_*|\mathbf{y}_{1:k}) &= \int p(\mathbf{f}_*|\mathbf{u})p(\mathbf{u}|\mathbf{y}_{1:k})d\mathbf{u} \\ &= \mathcal{N}(\mathbf{f}_*|\mathbf{H}_*\boldsymbol{\mu}_k, \mathbf{H}_*\boldsymbol{\Sigma}_k\mathbf{H}_*^T + \mathbf{V}_*) \end{aligned} \quad (17)$$

using  $p(\mathbf{f}_*|\mathbf{u}) = \mathcal{N}(\mathbf{f}_*|\mathbf{H}_*\mathbf{u}, \mathbf{V}_*)$  with  $\mathbf{H}_* = \mathbf{K}_{*R}\mathbf{K}_{RR}^{-1}$  and  $\mathbf{V}_*$  the model specific prediction covariance. The predictions for  $\mathbf{y}^*$  are obtained by adding  $\sigma_n^2\mathbb{I}$  to the covariance of  $\mathbf{f}_*|\mathbf{y}_{1:k}$ . At step  $K$ , by applying the Woodbury identity to the batch covariance  $\boldsymbol{\Sigma}_K$  in (14), we get for the predictive distribution in (17)

$$\begin{aligned} \boldsymbol{\mu}_K^* &= \mathbf{H}_*\boldsymbol{\Sigma}_0\mathbf{H}^T\boldsymbol{\Sigma}\mathbf{y}, \\ \boldsymbol{\Sigma}_K^* &= \mathbf{H}_*\boldsymbol{\Sigma}_0\mathbf{H}_*^T - \mathbf{H}_*\boldsymbol{\Sigma}_0\mathbf{H}^T\boldsymbol{\Sigma}\mathbf{H}\boldsymbol{\Sigma}_0\mathbf{H}_*^T + \mathbf{V}_*, \end{aligned} \quad (18)$$

where  $\boldsymbol{\Sigma} = (\mathbf{H}\boldsymbol{\Sigma}_0\mathbf{H}^T + \mathbf{V})^{-1}$ . Inserting the particular choices for  $\boldsymbol{\Sigma}_0$ ,  $\mathbf{H}$  and  $\mathbf{H}_*$  yields the usual formulation for the sparse predictive distribution

$$\begin{aligned} \boldsymbol{\mu}_K^* &= \mathbf{Q}_{*X}(\mathbf{Q}_{XX} + \mathbf{V})^{-1}\mathbf{y}; \\ \boldsymbol{\Sigma}_K^* &= \mathbf{Q}_{**} - \mathbf{Q}_{*X}(\mathbf{Q}_{XX} + \mathbf{V})^{-1}\mathbf{Q}_{X*} + \mathbf{V}_*. \end{aligned} \quad (19)$$

Depending on the choice of the covariances  $\bar{\mathbf{V}}$  and  $\mathbf{V}_*$ , we obtain for instance (3) for PEP, or the analogous predictions for VFE and FITC, respectively.

### 3.4 Online learning

This connection between sparse GP models and recursive estimation allows us to train the sparse GP models analytically online for streaming data for fixed  $\boldsymbol{\theta}$ .

As an illustrative example consider  $N = 100$  data samples in  $D = 1$ , we are interested in training a PEP model with  $\alpha = 0.5$  with  $M = 15$  inducing points with fixed  $\boldsymbol{\theta}$ . Let's assume that the data samples  $\{\mathbf{x}_k, y_k\}$  arrive sequentially in a stream. Thus we have  $B = 1$  and  $K = 100$ . Here we use the transformation and we apply, for each data sample  $k$ , the recursion in (16). We note that here  $\tilde{\mathbf{r}}_k$ ,  $\tilde{\mathbf{V}}_k$  and  $\tilde{\mathbf{S}}_k$  are numbers as  $B = 1$  and  $\tilde{\mathbf{H}}_k, \tilde{\mathbf{G}}_k \in \mathbb{R}^{15}$ .

We obtain predictions for new data  $\mathbf{X}_*$ , by applying (17) with  $\tilde{\mathbf{H}}_* = \mathbf{K}_{*R}$  and  $\mathbf{V}_* = \mathbf{K}_{*R}\mathbf{K}_{RR}^{-1}\mathbf{K}_{R*}$ . Note that there is no need to transform back the posterior over the inducing points, since it is already taken into account in the prediction step. After processing all  $N$  samples, the predictive distribution and the cumulative bound of log marginal likelihood correspond to the batch version, as shown in Fig. 3.

### 3.5 Marginal Likelihood

In the batch setting, the log marginal likelihood  $\log p(\mathbf{y})$  can be computed by marginalizing out  $\mathbf{u}$ , that is

$$\begin{aligned} \log p(\mathbf{y}) &= \log \int p(\mathbf{y}|\mathbf{u})p(\mathbf{u})d\mathbf{u} \\ &= \log \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{H}\boldsymbol{\Sigma}_0\mathbf{H}^T + \mathbf{V}). \end{aligned} \quad (20)$$

In the recursive setting,  $p(\mathbf{y})$  can be factorized into  $\prod_{k=1}^K p(\mathbf{y}_k|\mathbf{y}_{1:k-1})$ , where

$$\begin{aligned} p(\mathbf{y}_k|\mathbf{y}_{1:k-1}) &= \mathcal{N}(\mathbf{y}_k|\mathbf{H}_k\boldsymbol{\mu}_{k-1}, \mathbf{H}_k\boldsymbol{\Sigma}_{k-1}\mathbf{H}_k^T + \mathbf{V}_k) \\ &= \mathcal{N}(\mathbf{r}_k|\mathbf{0}, \mathbf{S}_k). \end{aligned} \quad (21)$$

The log of the joint marginal likelihood involving all terms of (21) can be explicitly written as

$$\begin{aligned} \log \prod_{k=1}^K p(\mathbf{y}_k|\mathbf{y}_{1:k-1}) &= \sum_{k=1}^K \log \mathcal{N}(\mathbf{r}_k|\mathbf{0}, \mathbf{S}_k) \\ &= -\frac{N}{2} \log 2\pi - \frac{1}{2} \sum_{k=1}^K \log |\mathbf{S}_k| + \mathbf{r}_k^T \mathbf{S}_k^{-1} \mathbf{r}_k. \end{aligned} \quad (22)$$

The iterative maximization of a lower bound of the recursive factorized marginal likelihood in (21) leads to the recursive KF updates in (16) for the posterior and to the lower bound

$$\psi(\boldsymbol{\theta}) = \sum_{k=1}^K \log \mathcal{N}(\mathbf{r}_k^\theta|\mathbf{0}, \mathbf{S}_k^\theta) - a_k(\boldsymbol{\theta}) \quad (23)$$

which includes a model specific regularization term  $a_k$  (see the right-most column in Fig. 2). We refer to this as the *recursive collapsed bound* and a detailed derivation for the VFE model is given in App. B. Using the model specific quantities  $\bar{\mathbf{V}}_k$  and  $a_k$ , this recursive computation of the lower bound of the marginal likelihood are equivalent to the batch counterparts for all sparse models, for instance (5) and (4) for VFE and PEP, respectively.

## 4 Hyper-parameters Estimation

The previous section presented an online procedure for training sparse GP models at fixed hyper-parameters  $\boldsymbol{\theta}$ . Here we show that, by exploiting the connections highlighted before, we can optimize  $\boldsymbol{\theta}$  sequentially. The *recursive collapsed bound* in (23) decomposes into a recursive sum over the mini-batches which allows to optimize the hyper-parameters  $\boldsymbol{\theta}$  sequentially as opposed to the collapsed bound (5). Our bound (23) enables the application of stochastic optimization without needing to estimate all entries in the posterior mean vector and covariance matrix as in SVGP. Compared to the uncollapsed bound in (6), the variational distribution is recursively and analytically eliminated, thus reducing the number of parameters to be numerically estimated drastically from  $\mathcal{O}(MD + M^2)$  to  $\mathcal{O}(MD)$ .

Finding a maximizer  $\boldsymbol{\theta} \in \Theta$  of an objective function  $\Psi(\boldsymbol{\theta}) = \sum_{k=1}^K \psi_k(\boldsymbol{\theta})$  can be achieved by applying *Stochastic*

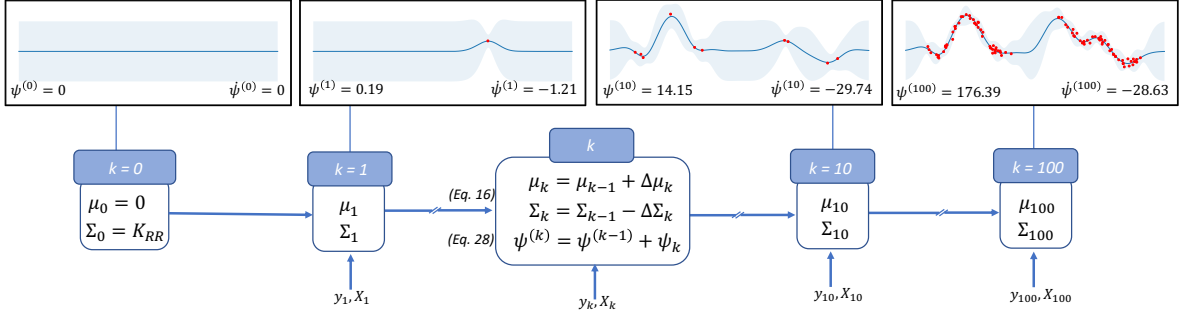


Fig. 3. Online learning for the toy example in Sect. 3.4 for fixed  $\theta$  with batch size  $B = 1$ . In each step  $k$ , the sample  $\mathbf{y}_k, \mathbf{X}_k$  is updated to the current posterior represented by  $\mu_k$  and  $\Sigma_k$  according to Eq. (16) together with the cumulative recursive bound  $\psi^{(k)}$  (Eq. (25) and the numbers in the left corners of the plots). For recursive parameter estimation as discussed in Sect. 4.1, in addition to the posterior and the cumulative bound, the recursive derivatives are propagated for each  $k$ . The derivatives of this bound w.r.t. to the lengthscale are indicated in the right corners of the plots. The cumulative bound and its derivative at step  $k = 100$  are equal to the corresponding batch version in Fig. 1.

Gradient Descent(SGD), with the update

$$\theta^{(t)} = \theta^{(t-1)} - \gamma^{(t-1)} \frac{\partial \psi_k}{\partial \theta} \Big|_{\theta=\theta^{(t-1)}}, \quad (24)$$

where  $\gamma^{(t-1)}$  might be a sophisticated function of  $\theta^{(0)}, \dots, \theta^{(t-1)}$  (for instance using ADAM [18], where also a bias correction term is included). We call one pass over the  $K$  mini-batches an epoch. We denote  $\theta^{(e,k)} \in \Theta^{(e,k)}$  the estimate of  $\theta$  in epoch  $e \in E$  for mini-batch  $k$ .

#### 4.1 Recursive Gradient Propagation (RGP)

We rewrite the recursive collapsed bound in (23) as

$$\psi^{(K)}(\theta) = \sum_{k=1}^K d_k(\theta) - a_k(\theta) = \sum_{k=1}^K \psi_k(\theta) \quad (25)$$

where  $d_k(\theta) = \log \mathcal{N}(\mathbf{r}_k^\theta | 0, \mathbf{S}_k^\theta)$  and  $a_k(\theta)$  the model specific regularization term. Since  $\psi^{(K)}(\theta)$  decomposes into a (recursive) sum over the mini-batches, we directly compute the derivative of  $\psi_k(\theta)$  w.r.t.  $\theta \in \Theta$ . The derivative of  $a_k$  is straightforward, for  $d_k$  we have

$$\frac{\partial d_k(\theta)}{\partial \theta} = -\frac{1}{2} \frac{\partial \log |\mathbf{S}_k^\theta|}{\partial \theta} - \frac{1}{2} \frac{\partial (\mathbf{r}_k^\theta)^T (\mathbf{S}_k^\theta)^{-1} \mathbf{r}_k^\theta}{\partial \theta} \quad (26)$$

with  $\mathbf{r}_k^\theta = \mathbf{y}_k - \mathbf{H}_k^\theta \mu_{k-1}^\theta$  and  $\mathbf{S}_k^\theta = \mathbf{H}_k^\theta \Sigma_{k-1}^\theta (\mathbf{H}_k^\theta)^T + \mathbf{V}_k^\theta$ . It is important to note that ignoring naively the dependency of  $\theta$  through  $\mu_{k-1}^\theta$  and  $\Sigma_{k-1}^\theta$  completely forgets the past and thus results in overfitting the current mini-batch. In order to compute the derivatives of  $\mu_k^\theta$  and  $\Sigma_k^\theta$ , we exploit the chain rule for derivatives and recursively propagate the gradients of the mean and the covariance over time, that is

$$\begin{aligned} \frac{\partial \mathbf{u}_k}{\partial \theta} &= \frac{\partial \mathbf{u}_{k-1}}{\partial \theta} + \frac{\partial \mathbf{G}_k}{\partial \theta} \mathbf{r}_k + \mathbf{G}_k \frac{\partial \mathbf{r}_k}{\partial \theta} \\ \frac{\partial \Sigma_k}{\partial \theta} &= \frac{\partial \Sigma_{k-1}}{\partial \theta} - \frac{\partial \mathbf{G}_k}{\partial \theta} \mathbf{S}_k \mathbf{G}_k^T - \mathbf{G}_k \frac{\partial \mathbf{S}_k}{\partial \theta} \mathbf{G}_k - \mathbf{G}_k \mathbf{S}_k \frac{\partial \mathbf{G}_k^T}{\partial \theta}, \end{aligned} \quad (27)$$

where  $\frac{\partial \mathbf{G}_k}{\partial \theta}$ ,  $\frac{\partial \mathbf{r}_k}{\partial \theta}$  and  $\frac{\partial \mathbf{S}_k}{\partial \theta}$  are computed recursively according to (16). Computing the derivatives of  $d_k$  as explained in (26) and (27), the stochastic gradient

$$\frac{\partial \psi_k(\theta)}{\partial \theta} = \frac{\partial d_k(\theta)}{\partial \theta} - \frac{\partial a_k(\theta)}{\partial \theta} \quad (28)$$

can be computed for each mini-batch  $k$ .

**Proposition 1** Consider a fixed parameter vector  $\bar{\theta} \in \Theta$ , the gradients of the full batch lower bound  $\mathcal{L}_{PEP}(\bar{\theta})$  in (4) with respect to  $\theta$  are equal to the cumulative partial derivatives  $\frac{\partial \psi^{(K)}}{\partial \theta}$  of the recursive collapsed bound with the above learning procedure. That is, it holds  $\frac{\partial \mathcal{L}_{PEP}(\bar{\theta})}{\partial \theta} = \frac{\partial \psi^{(K)}(\bar{\theta})}{\partial \theta}$ .

This shows that, when the gradients are cumulated over all data samples, each gradient step of our recursive procedure is equivalent to a gradient step for  $\mathcal{L}_{PEP}$  and, therefore, follows the gradients of an optimal collapsed lower bound.

The toy example in Fig. 3 also shows this equivalence. The numbers in the bottom left and right corners show the cumulative recursive collapsed bound  $\psi^{(k)}$  and its cumulative derivative  $\frac{\partial \psi^{(k)}}{\partial \theta}$  (abbreviated as  $\dot{\psi}^{(k)}$ ) with respect to the lengthscale. The lower bound of the marginal likelihood as well as its derivatives are exactly the same value as the corresponding batch counterpart in Fig. 1.

For Stochastic Recursive Gradient Propagation (SRGP), in each epoch  $e$  and mini-batch  $k$ , we interleave the update step of the inducing points in Eq. (15) with the SGD update (24) of the parameters  $\theta^{(e,k)}$ , i.e.

$$p(\mathbf{u} | \mathbf{y}_{1:k}, \theta^{(e,k)}) \approx p(\mathbf{y}_k | \mathbf{u}, \theta^{(e,k)}) p(\mathbf{u} | \mathbf{y}_{1:k-1}, \theta^{(e,k-1)}).$$

More concretely, we update after each mini-batch  $k$  the parameters  $\theta^{(e,k)}$  with (28),(24) and propagate recursively the posterior with (16) and its derivative (27). In order to compute all the derivatives with respect to  $\theta \in \Theta$ , we exploit several matrix derivative rules which simplify the computation significantly, see App. A. Finally note that the form



of the gradients in eq. (26) and (28) implies that the noise in the stochastic gradient at step  $k$  depends on the noise at step  $k-1$ , thus excluding standard convergence proofs such as [5]. Recent results on non-convex optimization problems [10,40,41] show convergence proofs for function classes that include our objective, nonetheless the Markov noise of our problem still excludes it from this general theory. For this reason we leave a convergence analysis to future work.

#### 4.2 Computational complexity

In the following, we assume that the batch size  $B$  is larger than the number of inducing points  $M$ . For one mini-batch, the time complexity to update the posterior is dominated by matrix multiplications of size  $B$  and  $M$ , thus  $\mathcal{O}(B^2M)$ . In order to propagate the gradients of the posterior and to compute the derivative of the bound needs  $\mathcal{O}(BM^2)$  for a mini-batch and a parameter  $\theta \in \Theta$ . Thus, updating a mini-batch including all  $\mathcal{O}(MD)$  parameters costs  $\mathcal{O}(BM^3D + B^2M)$  for the SRGP method. Since SRGP stores the gradients of the posterior, it requires  $\mathcal{O}(M^3D + BM)$  storage.

On the other hand, SVGP needs  $\mathcal{O}(M^2 + BM)$  storage and  $\mathcal{O}(BM^3 + B^2M)$  time per mini-batch, where the latter can be broken down into once  $\mathcal{O}(B^2M)$  and  $\mathcal{O}(BM)$  for each of the  $\mathcal{O}(M^2)$  parameters. This means, for moderate dimensions, our algorithm has the same time complexity as state of the art method SVGP. However, due to the analytic updates of the posterior we achieve a higher accuracy and less epochs are needed as shown in Fig. 4 and in Sect. 5 empirically.

Fig. 4 shows the convergence of SRGP on a 1-D toy example with  $N = 1000$  data samples and  $M = 15$  inducing points. The parameters are sequentially optimized with our recursive approach (blue) and as comparison with SVGP (green) with a mini-batch size of  $B = 100$  over several epochs. The root-mean-squared-error (RMSE) computed on test points, the bound of the log marginal likelihood (LML) as well as the hyper-parameters converge in a few iterations to the corresponding batch values of VFE (red). Due to the analytic updates of the posterior, the accuracy is higher and SRGP needs much less epochs until convergence.

#### 4.3 Mini-batch size

The size of the mini-batches has an impact on the speed of convergence of the algorithm. Proposition 1 tells us that if we use a full batch, our algorithm requires the same number of gradient updates as a full batch method to converge. On the other hand smaller batches should require more updates and should lead to a higher variance in the results. Fig. 5 shows a comparison of different mini-batch sizes on a 1-D toy example with  $N = 10'000$  data samples generated with the same parameters as in Sect. 5.1. The convergence to the full batch value is slower as the batch size decreases. Moreover the variance of the error, over the repetitions, is much larger for smaller batch sizes: in the last 10 normalized gradient updates, the standard deviation of the error is on average  $3.8 \times 10^{-3}$  for  $B = 100$  and  $8.1 \times 10^{-4}$  for  $B = 5'000$ , denoting a more stable procedure for higher batch sizes. As the mini-batch size increases the computational cost for each gradient update also increases. In this example one

gradient update requires on average  $4.9 \times 10^{-3}$  sec and  $5.8 \times 10^{-2}$  sec with  $B = 100$  and  $B = 5'000$  respectively.<sup>2</sup> These considerations suggest that a reasonable choice is a large mini-batch size within the computational and time budgets.

## 5 Experiments

We first benchmark our method with  $N = 100'000$  synthetic data samples generated by a GP in several dimensions. Next, we apply our approach to the Airline data used in [15] with a million of data samples. Finally a more realistic setup is presented where we use up to a million data samples to train a nonlinear plant. We compare our SRGP method to full GP and sparse batch method VFE for a subset of data (using the implementation in GPy [13]) and to the state of the art stochastic parameter estimation method SVGP implemented in GPflow [22]. Our algorithm works also for many other sparse models, however, only large-scale implementations of standard SVGP are available (corresponding to the VFE model), thus we restrict the investigation to this model.

### 5.1 GP Simulation

In this section we test our proposed learning procedure on simulated GP data. We generate  $N = 100'000$  data samples from a zero-mean (sparse) GP with SE covariance kernel with hyper-parameters  $\sigma_0 = 1, \sigma_n = 0.1$  and  $l = \{0.1, 0.2, 0.5\}$  in  $D = \{1, 2, 5\}$  dimensions. The initial  $M = \{20, 50, 100\}$  inducing points are randomly selected points from the data and the hyper-parameters of a SE kernel with individual lengthscales for each dimension are initialized to the same values for both algorithms ( $\sigma_0 = 1, \sigma_n = 1, l_1, \dots, l_D = 1$ ). All parameters are sequentially optimized with our recursive approach and with SVGP with a mini-batch size of  $B = 5000$ . The stochastic gradient descent method ADAM [18] is employed for both methods with learning rates  $\{0.001, 0.005, 0.005\}$  for SVGP and  $\{0.0001, 0.001, 0.005\}$  for SRGP (based on some preliminary experiments). Each experiment is replicated 10 times.

Fig. 6 shows the bound to the log marginal likelihood, the RMSE and the coverage of 10'000 test points for the data dimensions  $D = \{1, 2, 5\}$  of both methods over 50 epochs. The shaded lines indicates the 10 repetitions and the thick line correspond to the mean. The recursive propagation of the gradients achieves faster convergence and more accurate performance regarding mean RMSE and smaller values for the log marginal likelihood. The higher accuracy and faster convergence can possibly be explained by the analytic updates of the posterior mean and covariance which leads to less parameters to be optimized numerically.

### 5.2 Airline Data

For the second example we apply our recursive method to the Airline Data used in [15]. It consists of flight arrival and departure times for more than 2 millions flights in the USA from January 2008 and April 2008. We preprocessed

<sup>2</sup> The times are measured on a laptop with a Intel i5-7300U CPU @ 2.6 GHz.

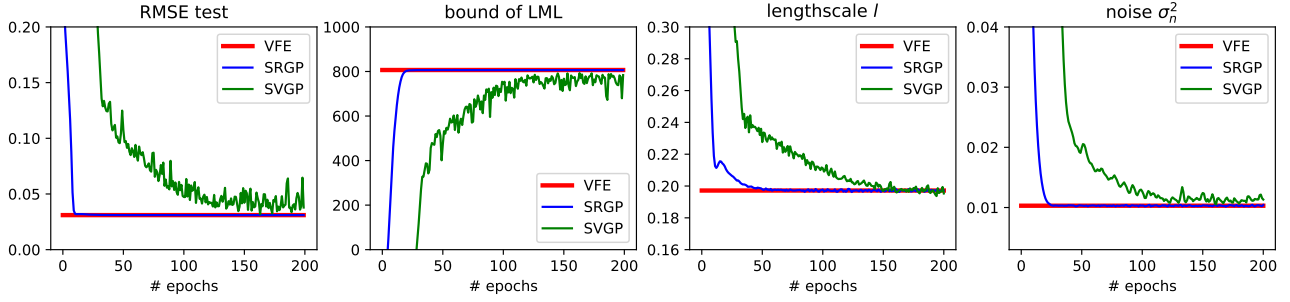


Fig. 4. Convergence of SRGP (blue) on a 1-D toy to batch version VFE (red). Compared to SVGP (green), the convergence of the root-mean-squared-error (RMSE) of test points, the bound of the log marginal likelihood (LML) as well as the hyper-parameters is faster and more accurate.

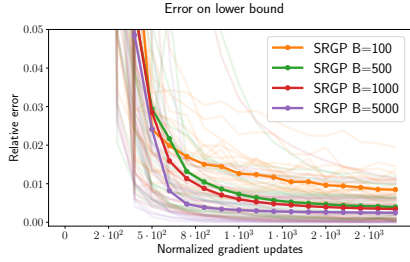


Fig. 5. Relative error in the bound of the log marginal likelihood between full batch and SRGP. Average over 20 repetitions.

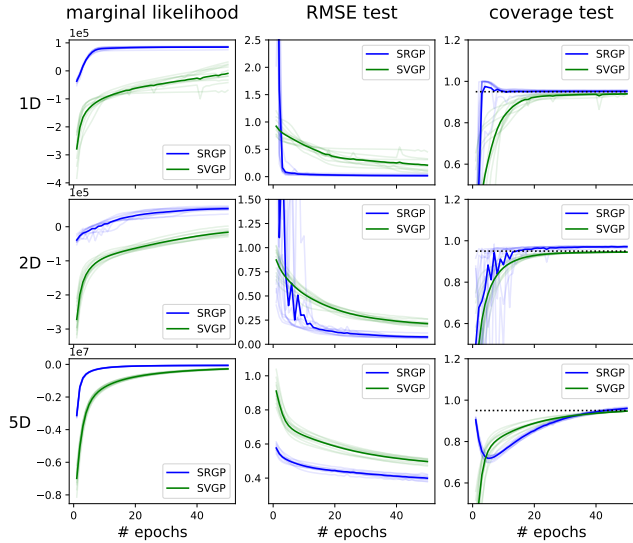


Fig. 6. Convergence over 50 epochs for  $N = 100'000$  synthetic GP data samples in several dimensions obtained by SVGP and our proposed method SRGP.

the data as similar as possible as described in [15] resulting in 8 variables: age of the aircraft, distance that needs to be covered, airtime, departure time, arrival time, day of the week, day of the month and month. We trained our recursive method as well as SVGP with a SE kernel on  $N = 1'000'000$  data samples with  $M = 500$  inducing points randomly selected from the data and a mini-batch size of  $B = 10'000$ . The ADAM learning rates are set to 0.005 for both methods and the size of the test set is 50'000. For 5 different repeti-

tions, the RMSE as a function of epochs is depicted in Fig. 7. The mean coverage on test data (at 95%) is comparable for both methods with values of 0.92 and 0.97 for SVGP and SRGP respectively. The overall performance of SRGP is superior to SVGP.

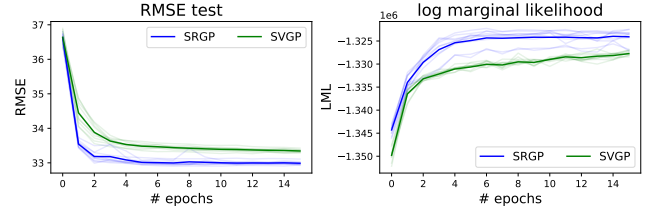


Fig. 7. Convergence over several epochs of RMSE and bound to log marginal likelihood for  $N = 1'000'000$  samples from the Airline data for SRGP and SVGP.

### 5.3 Non-Linear Plant

GPs are a powerful way to model complex functions in a non-parametric way, thus they are suitable to learn the complex input output behavior of a non-linear plant. However, with full or even sparse batch GP methods the use is restricted to a few thousands of samples. With our sequential learning method, we are able to exploit the huge amount of available data by training with up to a million of samples.

We consider a Continuous Stirred Tank Reactor (CSTR). The dynamic model of the plant is

$$\begin{aligned} \frac{d}{dt}h(t) &= w_1(t) + w_2(t) - 0.2\sqrt{h(t)} \\ \frac{d}{dt}C_b(t) &= (C_{b1} - C_b(t))\frac{w_1(t)}{h(t)} + (C_{b2} - C_b(t))\frac{w_2(t)}{h(t)} - \frac{k_1C_b(t)}{(1+k_2C_b(t))^2}, \end{aligned}$$

where  $C_b(t)$  is the product concentration at the output of the process,  $h(t)$  is the liquid level,  $w_1(t)$  is the flow rate of concentrated feed  $C_{b1}$ , and  $w_2(t)$  is the flow rate of the diluted feed  $C_{b2}$ . The input concentrations are  $C_{b1} = 24.9$  and  $C_{b2} = 0.1$ . The constants associated with the rate of consumption are  $k_1 = k_2 = 1$ . The objective of the controller is to maintain the product concentration by changing the flow  $w_1(t)$ . To simplify the example, we assume that  $w_2(t) = 0.1$  and that the level of the tank  $h(t)$  is not controlled. We denote the controlled outputs  $C_b(t), C_b(t-1), \dots, C_b(t-p)$  as  $f_t, f_{t-1}, \dots, f_{t-p}$



and the control variables as  $w_t, w_{t-1}, \dots, w_{t-p}$ . Therefore, the plant identification problem can be shaped into the problem of estimating the non-linear function  $f_t = g(f_{t-1}, \dots, f_{t-p}, w_t, w_{t-1}, \dots, w_{t-p})$  which depends on the  $p$  previous values as well as on the current and the  $p$  past control values  $w$ . However, we can only observe a

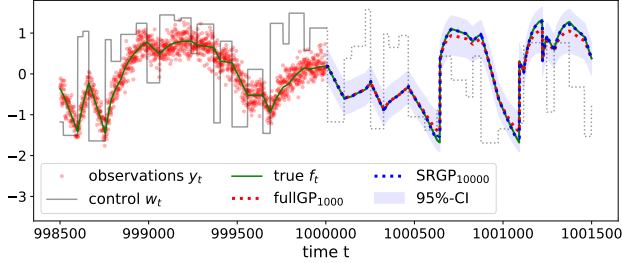


Fig. 8. Training and prediction phases for non-linear plant.

noisy version of the controlled response, that is  $y_t = f_t + \varepsilon_t$  with  $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$ . Using a sampling rate of  $0.2s$ , we have generated  $1'200'000$  observations (about 3 days of observations). The plant input is a series of steps, with random height (in the interval  $[0, 4]$ ), occurring at random intervals (in the interval  $[5, 20]s$ ). For different numbers  $N_{train}$ , we use the samples  $y_{10^6 - N_{train}}, \dots, y_{10^6}$  for training and the last  $200'000$  are used as a test set. The goal is to learn a model for the controlled response  $y_t$  given  $\mathbf{x}_t = [y_{t-1}, y_{t-2}, w_t, w_{t-1}, w_{t-2}]^T \in \mathbb{R}^5$  for the particular choice of  $p = 2$ . We model the non-linear function  $g$  with a GP with a SE kernel. For comparison, we train full GP and sparse batch GP (with 100 inducing points) on a time horizon  $N_{train}$  of up to  $10'000$  and  $50'000$  past values, respectively. With the sequential version SVGP and our recursive gradient propagation method SRGP (both with 100 inducing points and mini-batch size of  $1'000$ ), we use a time horizon of up to a million. This situation is depicted in Fig. 8, where for 1500 training samples  $y_t$  (red dots), the true (unknown) function  $f_t$  (green) and the control input  $w_t$  (grey) is shown together with the predicted values with full GP (red dotted) and recursive GP (blue dotted) trained on a time horizon of  $1'000$  and  $10'000$ , respectively. In Fig. 9, the RMSE and the median computed on the test set (with 10 repetitions) is depicted for full GP, sparse GP (VFE), SVGP and SRGP trained with varying time horizons. For small and medium training sizes, when the batch methods are applicable, our recursive method achieves the same performance as the batch counterpart (VFE) and is comparable to full GP. Due to the analytic updates of the posterior, SRGP outperforms SVGP regarding both RMSE and median for all training sizes. By exploiting more than several thousand past values, a significant increase in performance of SRGP can be still observed, thus it constitutes an approach to accurately scale GPs up to a million of past values.

## 6 Conclusion

In this paper we introduced a recursive inference and parameter estimation method SRGP for a general class of sparse GP approximations. Since the posterior updates are

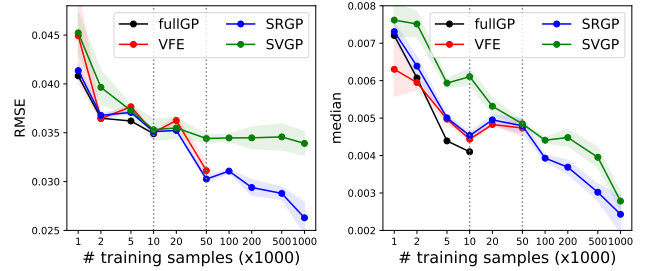


Fig. 9. RMSE and median for full GP, batch sparse GP (VFE), sequential SVGP, and our recursive method (SRGP) trained on varying time horizons (logarithmic scale). The grey dotted vertical lines at  $10'000$  and  $50'000$  indicate the maximal samples used for training full GP and sparse batch GP (VFE), respectively.

given analytically, one pass through the data is sufficient to compute the posterior for given parameters. For parameter estimation, we proposed a recursive collapsed bound to the log marginal likelihood that matches exactly the batch version but can be used for stochastic estimation. Due to the analytic updates of the posterior our method has much less parameters to be estimated numerically. As a consequence, the experimental section showed that our recursive method needs less epochs and has superior accuracy compared to state of the art, thus constitutes an efficient methodology for scaling GPs to big data problems.

Our approach could be enhanced in several directions. While the proposed method only exploits the update equations of the KF, an interesting direction would be to include a dynamic in a state space model that takes into account the varying hyper-parameters which makes it also applicable for the streaming setting as [6]. Moreover, we further plan to investigate distributed parameter estimation based on an information filter formulation of the problem. Finally, we aim to provide a convergence analysis of the proposed method to recursively learn the hyper-parameters using a similar approach recently employed for ADAM or AMSGrad [29,38].

## Acknowledgements

This work is supported by the Swiss National Research Programme 75 "Big Data" grant n. 407540\_167199 / 1.

## References

- [1] David Barber and Yali Wang. Gaussian processes for bayesian estimation in ordinary differential equations. In *ICML*, 2014.
- [2] Alessio Benavoli and Marco Zaffalon. State space representation of non-stationary gaussian processes. 2016.
- [3] Hildo Bijl, Thomas B. Schön, Jan-Willem van Wingerden, and Michel Verhaegen. Online sparse gaussian process training with input noise. *CoRR*, abs/1601.08068, 2016.
- [4] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, Berlin, Heidelberg, 2006.
- [5] Léon Bottou, Frank E. Curtis, and Jorge Nocedal. Optimization methods for large-scale machine learning. *SIAM Review*, 60(2):223–311, 2018.
- [6] Thang D Bui, Cuong Nguyen, and Richard E Turner. Streaming sparse gaussian process approximations. In *Advances in Neural Information Processing Systems*, pages 3301–3309, 2017.

- [7] Thang D Bui, Josiah Yan, and Richard E Turner. A unifying framework for sparse gaussian process approximation using power expectation propagation. *Journal of Machine Learning Research*, 18:1–72, 2017.
- [8] Andrea Carron, Marco Todescato, Ruggero Carli, Luca Schenato, and Gianluigi Pillonetto. Machine learning meets kalman filtering. In *2016 IEEE 55th Conference on Decision and Control (CDC)*, pages 4594–4599. IEEE, 2016.
- [9] Tianshi Chen, Henrik Ohlsson, and Lennart Ljung. On the estimation of transfer functions, regularizations and gaussian processes—revisited. *Automatica*, 48(8):1525–1535, 2012.
- [10] Xiangyi Chen, Sijia Liu, Ruoyu Sun, and Mingyi Hong. On the convergence of a class of Adam-type algorithms for non-convex optimization. *7th International Conference on Learning Representations, ICLR 2019*, 2019.
- [11] Lehel Csató and Manfred Opper. Sparse online gaussian processes. *Neural computation*, 14(3):641–668, 2002.
- [12] Roger Frigola, Fredrik Lindsten, Thomas B Schön, and Carl Edward Rasmussen. Bayesian inference and learning in gaussian process state-space models with particle mcmc. In *Advances in Neural Information Processing Systems*, pages 3156–3164, 2013.
- [13] GPpy. GPpy: A gaussian process framework in python. <http://github.com/SheffieldML/GPy>, since 2012.
- [14] Jouni Hartikainen and Simo Särkkä. Kalman filtering and smoothing solutions to temporal gaussian process regression models. In *2010 IEEE International Workshop on Machine Learning for Signal Processing*, pages 379–384. IEEE, 2010.
- [15] James Hensman, Nicolo Fusi, and Neil D Lawrence. Gaussian processes for big data. In *Conference for Uncertainty in Artificial Intelligence*, 2013.
- [16] Trong Nghia Hoang, Quang Minh Hoang, and Bryan Kian Hsiang Low. A unifying framework of anytime sparse gaussian process regression models with stochastic variational inference for big data. In *ICML*, pages 569–578, 2015.
- [17] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [19] Juš Kocijan, Agathe Girard, Blaž Banko, and Roderick Murray-Smith. Dynamic systems identification with gaussian processes. *Mathematical and Computer Modelling of Dynamical Systems*, 11(4):411–424, 2005.
- [20] Haitao Liu, Yew-Soon Ong, Xiaobo Shen, and Jianfei Cai. When gaussian process meets big data: A review of scalable gps. *arXiv preprint arXiv:1807.01065*, 2018.
- [21] Benn Macdonald, Catherine Higham, and Dirk Husmeier. Controversy in mechanistic modelling with gaussian processes. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML’15*, pages 1539–1547. JMLR.org, 2015.
- [22] Alexander G. de G. Matthews, Mark van der Wilk, Tom Nickson, Keisuke Fujii, Alexis Boukouvalas, Pablo León-Villagrà, Zoubin Ghahramani, and James Hensman. Gpflow: A Gaussian process library using TensorFlow. *Journal of Machine Learning Research*, 18(40):1–6, apr 2017.
- [23] César Lincoln C. Mattos, Zhenwen Dai, Andreas Damianou, Jeremy Forth, Guilherme A. Barreto, and Neil D. Lawrence. Recurrent Gaussian processes. In Hugo Larochelle, Brian Kingsbury, and Samy Bengio, editors, *Proceedings of the International Conference on Learning Representations*, volume 3, Caribe Hotel, San Juan, PR, 00 2016.
- [24] Gianluigi Pillonetto. A new kernel-based approach to hybrid system identification. *Automatica*, 70:21–31, 2016.
- [25] Gianluigi Pillonetto and Giuseppe De Nicolao. A new kernel-based approach for linear system identification. *Automatica*, 46(1):81–93, 2010.
- [26] Gianluigi Pillonetto, Francesco Dinuzzo, Tianshi Chen, Giuseppe De Nicolao, and Lennart Ljung. Kernel methods in system identification, machine learning and function estimation: A survey. *Automatica*, 50(3):657–682, 2014.
- [27] Joaquin Quiñero-Candela and Carl Edward Rasmussen. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- [28] Carl Edward Rasmussen and Christopher KI Williams. *Gaussian processes for machine learning*, volume 1. MIT press Cambridge, 2006.
- [29] Sashank J. Reddi, Satyen Kale, and Sanjiv Kumar. On the convergence of adam and beyond. In *International Conference on Learning Representations*, 2018.
- [30] Simo Sarkka, Arno Solin, and Jouni Hartikainen. Spatiotemporal learning via infinite-dimensional bayesian filtering and smoothing: A look at gaussian process regression through kalman filtering. *IEEE Signal Processing Magazine*, 30(4):51–61, 2013.
- [31] Matthias Seeger, Christopher Williams, and Neil Lawrence. Fast forward selection to speed up sparse gaussian process regression. In *Artificial Intelligence and Statistics 9*, number EPFL-CONF-161318, 2003.
- [32] Bernhard W Silverman. Some aspects of the spline smoothing approach to non-parametric regression curve fitting. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–52, 1985.
- [33] Alex J Smola and Peter L Bartlett. Sparse greedy gaussian process regression. In *Advances in neural information processing systems*, pages 619–625, 2001.
- [34] Edward Snelson and Zoubin Ghahramani. Sparse gaussian processes using pseudo-inputs. In *Advances in Neural Information Processing Systems*, pages 1257–1264, 2006.
- [35] Andreas Svensson and Thomas B Schön. A flexible state–space model for learning nonlinear dynamical systems. *Automatica*, 80:189–199, 2017.
- [36] Michalis Titsias. Variational learning of inducing variables in sparse gaussian processes. In *Artificial Intelligence and Statistics*, pages 567–574, 2009.
- [37] Marco Todescato, Andrea Carron, Ruggero Carli, Gianluigi Pillonetto, and Luca Schenato. Efficient spatio-temporal gaussian regression via kalman filtering. *arXiv preprint arXiv:1705.01485*, 2017.
- [38] P. T. Tran and L. T. Phong. On the convergence proof of amsgrad and a new version. *IEEE Access*, 7:61706–61716, 2019.
- [39] Grace Wahba, Xiwu Lin, Fangyu Gao, Dong Xiang, Ronald Klein, and Barbara Klein. The bias-variance tradeoff and the randomized gacv. In *Advances in Neural Information Processing Systems*, pages 620–626, 1999.
- [40] Dongruo Zhou, Yiqi Tang, Ziyang Yang, Yuan Cao, and Quanquan Gu. On the Convergence of Adaptive Gradient Methods for Nonconvex Optimization. 2018.
- [41] Fangyu Zou, Li Shen, Zequn Jie, Weizhong Zhang, and Wei Liu. A sufficient condition for convergences of Adam and RMSProp. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2019-June(1):11119–11127, 2019.
- [42] M. A. Álvarez, D. Luengo, and N. D. Lawrence. Linear latent force models using gaussian processes. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(11):2693–2705, Nov 2013.

## A Details for Recursive Gradient Propagation

We show here the computation for the recursive gradient propagation from Sect. 4.1 for the PEP model. For other models,  $a_k$  and  $\bar{\mathbf{V}}$  from the Table 2 could be used correspondingly. We will use the following notation:  $\text{diag}[\mathbf{A}] = \mathbf{d}, d_i = a_{ii}$ ,  $\text{Diag}[\mathbf{d}] = \mathbf{A}, a_{ii} = d_i, a_{ij} = 0$ ,  $\mathbf{A} \odot \mathbf{B} = \mathbf{C}, c_{ij} = a_{ij}b_{ij}$ ,  $\mathbf{A} \div \mathbf{B} = \mathbf{C}, c_{ij} = \frac{a_{ij}}{b_{ij}}$ ,  $\mathbf{A}^{\odot 2} = \mathbf{C}, c_{ij} = a_{ij}^2$ ,  $\dot{\mathbf{A}} = \frac{\partial \mathbf{A}(\theta)}{\partial \theta}, \forall \theta \in \Theta$ ,  $\text{sum}[\mathbf{A}] = \sum_{i,j} a_{ij}$ ,  $1_{[z]} = 1$  if  $z = \text{true}$ , 0 other.

### Initialization

$$\begin{aligned} \eta_0 &= \mathbf{0}; & \eta_0 &= \mathbf{0}; \\ \Lambda_0 &= \mathbf{K}_{RR}^{-1}; & \dot{\Lambda}_0 &= -\mathbf{K}_{RR}^{-1} \dot{\mathbf{K}}_{RR} \mathbf{K}_{RR}^{-1}; \\ \psi_0 &= -\frac{N}{2} \log 2\pi; & \psi_0 &= 0; \\ \Sigma_0 &= \mathbf{K}_{RR}; & \log \text{Det}_0 &= \log |\Lambda_0|; \end{aligned}$$

### Natural Mean and Precision Updates

$$\begin{aligned} \mathbf{H}_k &= \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1}; \\ \mathbf{d}_k &= \text{diag}[\mathbf{K}_{X_k X_k} - \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1} \mathbf{K}_{R X_k}]; \\ \mathbf{v}_k &= \alpha \mathbf{d}_k + \sigma_n^2 \mathbb{1}; \quad \mathbf{V}_k^{-1} = \text{Diag}[\mathbb{1} \div \mathbf{v}_k]; \\ a_k &= \frac{1-\alpha}{\alpha} \left( \sum_{i=1}^B \log([\mathbf{v}_k]_i) - B \log \sigma_n^2 \right); \end{aligned}$$

$$\begin{aligned} \mathbf{r}_k &= \mathbf{y}_k - \mathbf{H}_k \Sigma_{k-1} \eta_{k-1}; \\ \eta_k &= \eta_{k-1} + \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{y}_k; \\ \Lambda_k &= \Lambda_{k-1} + \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{H}_k; \end{aligned}$$

$$\begin{aligned} \Sigma_k, \log |\Lambda_k| &= \Lambda_k^{-1}, \log |\Lambda_k|; \\ \mathbf{S}_k^{-1} &= \mathbf{V}_k^{-1} - \mathbf{V}_k^{-1} \mathbf{H}_k \Sigma_k \mathbf{H}_k^T \mathbf{V}_k^{-1}; \\ \psi_k &= \psi_{k-1} - \frac{1}{2} (\log |\Lambda_k| - \log |\Lambda_{k-1}| \\ &\quad - \log |\mathbf{V}_k^{-1}| + \mathbf{r}_k^T \mathbf{S}_k^{-1} \mathbf{r}_k + a_k) \end{aligned}$$

### Intermediate Derivatives

$$\begin{aligned} \dot{L}_{d\mathbf{H}_k} &= 2(\mathbf{V}_k^{-1} \mathbf{H}_k \Sigma_k - \mathbf{S}_k^{-1} \mathbf{r}_k (\Sigma_{k-1} \eta_{k-1} \\ &\quad + \Sigma_k \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{r}_k)^T) \\ \dot{L}_{d\mathbf{v}_k} &= - \left( \text{diag}[\mathbf{H}_k \Sigma_k \mathbf{H}_k^T] - \frac{1}{\alpha} \mathbf{v}_k \right. \\ &\quad \left. + (\mathbf{r}_k - \mathbf{H}_k \Sigma_k \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{r}_k)^{\odot 2} \right) \div \mathbf{v}_k^{\odot 2} \\ \dot{L}_{d\mathbf{K}_{X_k R}} &= \dot{L}_{d\mathbf{H}_k} \mathbf{K}_{RR}^{-1} - 2\alpha \text{Diag}[\dot{L}_{d\mathbf{v}_k}] \mathbf{H}_k \\ \dot{L}_{d\mathbf{K}_{RR}} &= -\mathbf{H}_k^T (\dot{L}_{d\mathbf{H}_k} \mathbf{K}_{RR}^{-1} - \alpha \text{Diag}[\dot{L}_{d\mathbf{v}_k}] \mathbf{H}_k) \\ \dot{L}_{d\mathbf{K}_{X_k X_k}} &= \alpha \dot{L}_{d\mathbf{v}_k} \\ \dot{L}_{d\Lambda_k} &= \Sigma_k - \Sigma_{k-1} + 2\Sigma_{k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1} \mathbf{r}_k \eta_{k-1}^T \Sigma_{k-1} \\ &\quad + \Sigma_k \mathbf{H}_k^T \mathbf{V}_k^{-1} \mathbf{r}_k \mathbf{r}_k^T \mathbf{V}_k^{-1} \mathbf{H}_k \Sigma_k \\ \dot{L}_{d\eta_k} &= -2\Sigma_{k-1} \mathbf{H}_k^T \mathbf{S}_k^{-1} \mathbf{r}_k \\ \dot{L}_{d.d_n} &= 2\sigma_n^2 \text{sum}[\dot{L}_{d\mathbf{v}_k}] - 2B \frac{1-\alpha}{\alpha} \end{aligned}$$

## Derivative Updates

### Loop over $\theta_i \in \Theta$ :

$$\begin{aligned} \dot{\psi}_k &= \dot{\psi}_{k-1} - \frac{1}{2} (\text{sum}[\dot{L}_{d\eta_k} \odot \dot{\eta}_{k-1}] \\ &\quad + \text{sum}[\dot{L}_{d\Lambda_k} \odot \dot{\Lambda}_{k-1}] + \text{sum}[\dot{L}_{d\mathbf{K}_{RR}} \odot \dot{\mathbf{K}}_{RR}]) \\ &\quad + \text{sum}[\dot{L}_{d\mathbf{K}_{X_k R}} \odot \dot{\mathbf{K}}_{X_k R}] \\ &\quad + \text{sum}[\dot{L}_{d\mathbf{K}_{X_k X_k}} \odot \dot{\mathbf{K}}_{X_k X_k}] + 1_{[\theta_k = \sigma_n]} \dot{L}_{d.d_n}) \\ \dot{\mathbf{H}}_k &= \dot{\mathbf{K}}_{X_k R} \mathbf{K}_{RR}^{-1} - \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1} \dot{\mathbf{K}}_{RR} \mathbf{K}_{RR}^{-1}; \\ \dot{\mathbf{d}}_k &= \text{diag}[\dot{\mathbf{K}}_{X_k X_k} - \dot{\mathbf{K}}_{X_k R} \mathbf{K}_{RR}^{-1} \mathbf{K}_{R X_k} \\ &\quad + \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1} \dot{\mathbf{K}}_{RR} \mathbf{K}_{RR}^{-1} \mathbf{K}_{R X_k} \\ &\quad - \mathbf{K}_{X_k R} \mathbf{K}_{RR}^{-1} \dot{\mathbf{K}}_{R X_k}]; \\ \dot{\mathbf{V}}_k^{-1} &= -1_{[\theta_i \neq \sigma_n]} \alpha \mathbf{V}_k^{-1} \text{Diag}[\dot{\mathbf{d}}_k] \mathbf{V}_k^{-1} \\ &\quad - 1_{[\theta_i = \sigma_n]} 2\sigma_n^2 \dot{\mathbf{V}}_k^{-1} \dot{\mathbf{V}}_k^{-1}; \\ \dot{\eta}_k &= \dot{\eta}_{k-1} + \dot{\mathbf{H}}_k^T \mathbf{V}_k^{-1} \mathbf{y}_k + \mathbf{H}_k^T \dot{\mathbf{V}}_k^{-1} \mathbf{y}_k; \\ \dot{\Lambda}_k &= \dot{\Lambda}_{k-1} + \dot{\mathbf{H}}_k^T \mathbf{V}_k^{-1} \mathbf{H}_k \\ &\quad + \mathbf{H}_k^T \dot{\mathbf{V}}_k^{-1} \mathbf{H}_k + \mathbf{H}_k^T \mathbf{V}_k^{-1} \dot{\mathbf{H}}_k \end{aligned}$$

For the noise  $\sigma_n$ , all kernel derivatives are zero, therefore the calculations simplify significantly.

## B Derivation of Recursive Collapsed Bound

We provide more details and a detailed derivation for the recursive collapsed lower bound (25) for the VFE model from Sect. 4.1. Instead of lower bounding directly the batch log marginal likelihood as done by Titsias [36] in the batch case, our approach relies on the recursive factorization of the joint log marginal likelihood  $\log p(\mathbf{y}|\theta) = \log \prod_{k=1}^K p(\mathbf{y}_k | \mathbf{y}_{1:k-1}, \theta) = \sum_{k=1}^K \log p(\mathbf{y}_k | \mathbf{y}_{1:k-1}, \theta)$ . The properties induced by the sparse augmented inducing point model yield  $\log p(\mathbf{y}|\theta) = \sum_{k=1}^K \log \int p(\mathbf{y}_k | \mathbf{f}_k, \theta) p(\mathbf{f}_k | \mathbf{u}, \theta) p(\mathbf{u} | \mathbf{y}_{1:k-1}, \theta) d\mathbf{f}_k d\mathbf{u}$ . We first introduce the variational distributions  $q_k(\mathbf{f}_k, \mathbf{u}) = p(\mathbf{f}_k | \mathbf{u}, \theta) q_k(\mathbf{u}) \approx p(\mathbf{f}_k, \mathbf{u} | \mathbf{y}_{1:k}, \theta)$ , then by applying Jensen's inequality to each individual term in the true log marginal likelihood, we obtain the lower bound

$$\begin{aligned} \log p(\mathbf{y}|\theta) &\geq \sum_{k=1}^K \int p(\mathbf{f}_k | \mathbf{u}, \theta) q_k(\mathbf{u}) \dots \\ &\dots \log \frac{p(\mathbf{y}_k | \mathbf{f}_k, \theta) p(\mathbf{f}_k | \mathbf{u}, \theta) p(\mathbf{u} | \mathbf{y}_{1:k-1}, \theta)}{p(\mathbf{f}_k | \mathbf{u}, \theta) q_k(\mathbf{u})} d\mathbf{f}_k d\mathbf{u}. \end{aligned}$$

The quantity  $p(\mathbf{u} | \mathbf{y}_{1:k-1}, \theta)$  is unknown, however, we can replace it with  $q_{k-1}(\mathbf{u})$  leading to  $\mathcal{L}(q_1, \dots, q_K, \theta)$

$$\sum_{k=1}^K \int p(\mathbf{f}_k | \mathbf{u}, \theta) q_k(\mathbf{u}) \log \frac{p(\mathbf{y}_k | \mathbf{f}_k, \theta) q_{k-1}(\mathbf{u})}{q_k(\mathbf{u})} d\mathbf{f}_k d\mathbf{u}. \quad (\text{B.1})$$

Maximizing this lower bound recursively with respect to the distributions  $q_k(\mathbf{u})$  leads to a sequence of optimal variational

distributions  $q_k^*(\mathbf{u})$  for the inducing outputs

$$\mathcal{N}\left(\mathbf{u}|\Sigma_k \left\{ \frac{1}{\sigma_n^2} \mathbf{H}_k^T \mathbf{y}_k + \Sigma_{k-1}^{-1} \boldsymbol{\mu}_{k-1} \right\}, \Sigma_k\right), \quad (\text{B.2})$$

where  $\Sigma_k = \left(\Sigma_{k-1}^{-1} + \frac{1}{\sigma_n^2} \mathbf{H}_k^T \mathbf{H}_k\right)^{-1}$  and  $\mathbf{H}_k = \mathbf{K}_{\mathbf{X}_k \mathbf{R}} \mathbf{K}_{\mathbf{R} \mathbf{R}}^{-1}$ . Plugging  $q_k^*(\mathbf{u})$  into (B.1) yields

$$\begin{aligned} \mathcal{L}_{REC}(\boldsymbol{\theta}) = & \sum_{k=1}^K [\log \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \boldsymbol{\mu}_{k-1}, \dots \\ & \dots \mathbf{H}_k \Sigma_{k-1} \mathbf{H}_k^T + \sigma_n^2 \mathbb{I}) - \frac{\text{Tr}[\mathbf{D}_{\mathbf{X}_k \mathbf{X}_k}]}{2\sigma_n^2}]. \end{aligned} \quad (\text{B.3})$$

The recursive bound  $\mathcal{L}_{REC}(\boldsymbol{\theta})$  is equivalent to the batch collapsed bound  $\mathcal{L}_{VFE}(\boldsymbol{\theta})$  in (5) and it holds  $\mathcal{L}_{SVGP}(\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\theta}) \leq \mathcal{L}_{REC}(\boldsymbol{\theta})$  with equality when inserting the optimal variational posterior. The variational posterior update in (B.2) has the same form as the recursive updates in (15). Similarly, the recursive collapsed lower bound in (B.3) is equal to (25).

We provide below a detailed derivation that follows closely the proof in [36] for the recursive collapsed bound (B.3) as well as the sequence of optimal distributions (B.2). We assume mini-batches of size  $B$ , that is, we have training data  $\mathcal{D} = \{\mathbf{y}_k, \mathbf{X}_k\}_{k=1}^K$  and the corresponding latent function values  $\{\mathbf{f}_k\}_{k=1}^K$ . We briefly recap the involved quantities and introduce abbreviations:

$$\begin{aligned} p(\mathbf{y}_k | \mathbf{f}_k, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{y}_k | \mathbf{f}_k, \sigma_n^2 \mathbb{I}); \\ p(\mathbf{f}_k | \mathbf{u}, \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{f}_k | \mathbf{H}_k \mathbf{u}, \mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}); \\ p(\mathbf{u} | \boldsymbol{\theta}) &= \mathcal{N}(\mathbf{u} | \mathbf{0}, \mathbf{K}_{\mathbf{R} \mathbf{R}}) = q_0(\mathbf{u}); \\ q_{k-1}(\mathbf{u}) &= \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}_{k-1}, \Sigma_{k-1}) \approx p(\mathbf{u} | \mathbf{y}_{1:k-1}, \boldsymbol{\theta}) \\ \mathbf{H}_k &= \mathbf{K}_{\mathbf{X}_k \mathbf{R}} \mathbf{K}_{\mathbf{R} \mathbf{R}}^{-1}; \\ \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k} &= \mathbf{K}_{\mathbf{X}_k \mathbf{R}} \mathbf{K}_{\mathbf{R} \mathbf{R}}^{-1} \mathbf{K}_{\mathbf{R} \mathbf{X}_k}. \end{aligned}$$

Starting from (B.1), we have the bound

$$\sum_{k=1}^K \int p(\mathbf{f}_k | \mathbf{u}, \boldsymbol{\theta}) q_k(\mathbf{u}) \log \frac{p(\mathbf{y}_k | \mathbf{f}_k, \boldsymbol{\theta}) q_{k-1}(\mathbf{u})}{q_k(\mathbf{u})} d\mathbf{f}_k d\mathbf{u}$$

which can be rearranged to

$$\sum_{k=1}^K \int q_k(\mathbf{u}) \left\{ \log G(\mathbf{u}, \mathbf{y}_k) + \log \frac{q_{k-1}(\mathbf{u})}{q_k(\mathbf{u})} \right\} d\mathbf{u},$$

where  $\log G(\mathbf{u}, \mathbf{y}_k) = \int p(\mathbf{f}_k | \mathbf{u}, \boldsymbol{\theta}) \log p(\mathbf{y}_k | \mathbf{f}_k, \boldsymbol{\theta}) d\mathbf{f}_k$ . The integral involving  $\mathbf{f}_k$  is computed as  $\log G(\mathbf{u}, \mathbf{y}_k) = \int p(\mathbf{f}_k | \mathbf{u}, \boldsymbol{\theta}) \log p(\mathbf{y}_k | \mathbf{f}_k, \boldsymbol{\theta}) d\mathbf{f}_k$ , which equals

$$\begin{aligned} &= \mathbb{E}_{\mathbf{f}_k | \mathbf{u}} \left[ -\frac{B}{2} \log(2\pi\sigma_n^2) - \frac{1}{2} [\mathbf{y}_k - \mathbf{f}_k]^T \frac{1}{\sigma_n^2} [\mathbf{y}_k - \mathbf{f}_k] \right] \\ &= -\frac{B}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (\mathbf{y}_k^T \mathbf{y}_k + \mathbb{E}_{\mathbf{f}_k | \mathbf{u}} [\mathbf{f}_k^T \mathbf{f}_k - 2\mathbf{y}_k^T \mathbf{f}_k]). \end{aligned}$$

Using  $\mathbb{E}[\mathbf{x}^T \mathbf{A} \mathbf{x}] = \text{Tr}[\mathbf{A} \boldsymbol{\Sigma}] + \boldsymbol{\mu}^T \mathbf{A} \boldsymbol{\mu}$  with  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma})$  yields

$$\begin{aligned} \log G(\mathbf{u}, \mathbf{y}_k) &= -\frac{K}{2} \log(2\pi\sigma_n^2) - \frac{1}{2\sigma_n^2} (\mathbf{y}_k^T \mathbf{y}_k \\ &+ \text{Tr}[\mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}] + \mathbf{u}^T \mathbf{H}_k^T \mathbf{H}_k \mathbf{u} - 2\mathbf{y}_k^T \mathbf{H}_k \mathbf{u}) \\ &= \log[\mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \sigma_n^2 \mathbb{I})] - \frac{1}{2\sigma_n^2} \text{Tr}[\mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}]. \end{aligned}$$

Substitute this expression back, the lower bound becomes

$$\begin{aligned} \sum_{k=1}^K \left[ \int q_k(\mathbf{u}) \log \frac{\mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \sigma_n^2 \mathbb{I}) q_{k-1}(\mathbf{u})}{q_k(\mathbf{u})} d\mathbf{u} \right. \\ \left. - \frac{1}{2\sigma_n^2} \text{Tr}[\mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}] \right]. \end{aligned}$$

We can now maximize this bound with respect to  $q_k(\mathbf{u})$ . Here since we have not constrained  $q_b$  to belong to any fixed family of distributions, we can compute the optimal bound by reversing the Jensen's inequality leading to

$$\begin{aligned} \mathcal{L}_{REC}(\boldsymbol{\theta}) &= \sum_{k=1}^K \left[ \log \int \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \sigma_n^2 \mathbb{I}) q_{k-1}(\mathbf{u}) d\mathbf{u} \right. \\ &\quad \left. - \frac{1}{2\sigma_n^2} \text{Tr}[\mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}] \right] \\ &= \sum_{k=1}^K \left[ \log \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \boldsymbol{\mu}_{k-1}, \mathbf{H}_k \Sigma_{k-1} \mathbf{H}_k^T + \sigma_n^2 \mathbb{I}) \right. \\ &\quad \left. - \frac{1}{2\sigma_n^2} \text{Tr}[\mathbf{K}_{\mathbf{X}_k \mathbf{X}_k} - \mathbf{Q}_{\mathbf{X}_k \mathbf{X}_k}] \right] \end{aligned}$$

where we used a linear Gaussian identity (see, e.g., [4] Ch. 2.3) in the last step. This is equal to Eq. (25). The optimal distribution  $q_k^*$  that gives rise to this bound is proportional to  $\mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \sigma_n^2 \mathbb{I}) q_{k-1}(\mathbf{u}) = \mathcal{N}(\mathbf{y}_k | \mathbf{H}_k \mathbf{u}, \sigma_n^2 \mathbb{I}) \mathcal{N}(\mathbf{u} | \boldsymbol{\mu}_{k-1}, \Sigma_{k-1})$  and can be analytically computed leading to

$$q_k^*(\mathbf{u}) = \mathcal{N}\left(\mathbf{u} | \frac{1}{\sigma_n^2} \Sigma_k \mathbf{H}_k^T \mathbf{y}_k + \Sigma_{k-1}^{-1} \boldsymbol{\mu}_{k-1}, \Sigma_k\right)$$

where  $\Sigma_k = \left(\Sigma_{k-1}^{-1} + \frac{1}{\sigma_n^2} \mathbf{H}_k^T \mathbf{H}_k\right)^{-1}$ . This matches the result in Eq. (15) and (16) and completes the proof.

## C Proof of proposition 1

We showed (App. B) how the batch approximate log-likelihood  $\mathcal{L}_{PEP}$  in (4) can be recursively computed. In Sect. 4.1, we showed an equivalent way to compute it, that is, we have the cumulative bound  $\Psi^{(K)}(\boldsymbol{\theta}) = \sum_k^K \psi_k(\boldsymbol{\theta})$  and it satisfies  $\mathcal{L}_{PEP}(\boldsymbol{\theta}) = \sum_k^K \psi_k(\boldsymbol{\theta})$ . By induction and linearity of derivatives, we therefore get  $\frac{\partial \Psi^{(K)}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{\partial \psi_k(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} + \frac{\partial \Psi^{(K-1)}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$  and  $\frac{\partial \mathcal{L}_{PEP}}{\partial \boldsymbol{\theta}} = \frac{\partial \Psi^{(K)}}{\partial \boldsymbol{\theta}}$ .