

# A Tutorial on Sparse Variational Gaussian Processes

P.L.Green  
Institute for Risk and Uncertainty  
School of Engineering  
University of Liverpool  
Liverpool, L69 7ZF  
United Kingdom

p.l.green@liverpool.ac.uk  
<https://www.liverpool.ac.uk/engineering/staff/peter-green/>

May 22, 2019

## Acknowledgments

Thank you to Ryan Jackson (University of Liverpool, UK) and Kaikai Zhao (Indiana University, Bloomington, USA) for their feedback!

## 1 Introduction

This document is supposed to provide a tutorial-style introduction to Sparse Variational Gaussian Processes which are typically able to analyse larger data sets than standard Gaussian Processes. As with our other tutorials, the current document is accompanied by some example Python code. The code is heavily commented and is designed to help those new to the topic. It is recommended that this tutorial should only be studied once a basic understanding of standard Gaussian Processes has been established. We offer such a tutorial here: <https://github.com/plgreenLIRU/Gaussian-Process-regression-tutorial>.

## 2 Tutorial Structure

The tutorial is structured as follows. A general overview of the approach behind Sparse Variational Gaussian Processes (herein simply referred to as ‘Sparse GPs’) as well as a summary of key results is given in Section 3. Example code is then described immediately afterwards, in Section 4. At this point the reader may wish to stop, however, for those wanting more information, more detailed derivations are then given in Section 5 before a summary of some interesting references is given in Section 6.

## 3 Overview

In this tutorial, we specifically consider regression problems where we define  $\mathbf{y} = (y_1, \dots, y_N)^T$  as noisy observations of a function,  $f$ , at inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , such that our full set of training data is  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ . We assume that our observations are made according to the noise model

$$y_i = f(\mathbf{x}_i) + \epsilon, \quad \epsilon \sim \mathcal{N}(\epsilon; 0, \sigma^2) \quad (1)$$

Using the notation  $f_i \equiv f(\mathbf{x}_i)$ , the vector of true function values is written as  $\mathbf{f} = (f_1, \dots, f_N)^T$ . If  $N$  is large, computational cost often prevents us from using ‘standard’ GPs. With the Variational Sparse GP described in [1] and [2], we train our GP on a smaller set of data, consisting of observations at the set of inputs  $\mathbf{X}_m = \{\mathbf{x}_{m,1}, \dots, \mathbf{x}_{m,M}\}$  (where, for this tutorial, we will assume that  $\mathbf{X}_m$  is a subset of  $\mathbf{X}$ ). The aim is to make the Sparse GP ‘look like’ a standard GP that had been trained on the full set of training data as closely as possible. This approach therefore assumes that there is some redundancy in our data - it builds on the idea that we only really need a small subset of ‘important’ points to get a reasonable estimate of the true function,  $f$ .

Stated briefly, we proceed in 3 key steps. These are described in the boxes below:

### Step 1: Define the optimum inducing points.

In the following, we define the set of inputs  $\mathbf{X}_m$  as ‘inducing points’. As a first step, we have to define what we mean by the *optimum* set of inducing points - in other words, we ask the question: what do we use to help us choose the ‘best’  $\mathbf{X}_m$ ?

Say  $\mathbf{f}_m$  is the vector of true function values at the points in  $\mathbf{X}_m$  (such that  $f_{m,i} \equiv f(\mathbf{x}_{m,i})$ ). We begin by stating that the optimum choice of  $\mathbf{X}_m$  would lead to  $\mathbf{f}_m$  such that

$$p(\mathbf{y} | \mathbf{f}_m) = p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m) \quad (2)$$

(where dependence on  $\mathbf{X}$  and  $\mathbf{X}_m$  isn’t written for notational simplicity). Essentially, if we pick the best possible  $\mathbf{f}_m$ , then knowing  $\mathbf{f}$  wouldn’t change the probability of witnessing the observations  $\mathbf{y}$ . It turns out that enforcing this condition implies the following factorisation of our posterior probability distribution over  $\mathbf{f}$  and  $\mathbf{f}_m$ :

$$p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = p(\mathbf{f}_m | \mathbf{y})p(\mathbf{f} | \mathbf{f}_m) \quad (3)$$

From this we conclude that a good choice of  $\mathbf{X}_m$  is one that leads to equation (3) being true.

### Step 2: Variational Bayes framework

Briefly stated, variational Bayesian approaches involve the construction of a *variational distribution*,  $q$ , that closely approximates some distribution of interest. A good choice of variational distribution is usually defined as that which minimises the *KL divergence* between  $q$  and the distribution of interest (see [3][4] for more information).

In our case, we choose to approximate the posterior with the PDF,  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y})$ , with

$$q(\mathbf{f}, \mathbf{f}_m) = \phi(\mathbf{f}_m)p(\mathbf{f} | \mathbf{f}_m) \quad (4)$$

where  $\phi(\mathbf{f}_m)$  is a PDF that is currently unknown. We hypothesise that, if we pick the optimum  $\mathbf{f}_m$ , then the posterior will factorise as shown in equation (3) and it would

then be possible to find the distribution,  $\phi(\mathbf{f}_m)$ , that sets the KL divergence between  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y})$  and  $q(\mathbf{f}, \mathbf{f}_m)$  equal to zero. Conveniently, we know that minimising the KL divergence is equivalent to maximising a lower bound on the log marginal likelihood, which we denote as  $L$  (again, see [3][4] for more details).

### Step 3: Maximise a lower bound

Rather than maximising  $L$  directly it turns out that, by instead maximising a lower bound on  $L$ , we can still select some ‘good’ inducing points while avoiding having to find the variational distribution  $\phi(\mathbf{f}_m)$ . This lower bound on  $L$  is denoted  $L_{lb}$ , and is given by

$$L_{lb} = \log [\mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}_n \sigma^2 + \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn})] - \frac{1}{2\sigma^2} \text{Tr}(\text{Cov}[\mathbf{f} | \mathbf{f}_m]) \quad (5)$$

where

$$\text{Cov}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nn} - \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \quad (6)$$

$$\mathbf{K}_{nn}(i, j) = k(\mathbf{x}_i, \mathbf{x}_j), \quad \mathbf{K}_{nm}(i, j) = k(\mathbf{x}_i, \mathbf{x}_{m,j}), \quad \mathbf{K}_{mm}(i, j) = k(\mathbf{x}_{m,i}, \mathbf{x}_{m,j}) \quad (7)$$

and  $k(\cdot, \cdot)$  represents the GP kernel function. The approach is effective as, if  $M$  is small, we can evaluate this lower bound relatively quickly, particularly as we only really need the diagonal terms of  $\mathbf{K}_{nn}$ . Note that the same objective function is also used to optimise the hyperparameters. The selection of inducing points and optimisation of the hyperparameters can be conducted sequentially, in a similar manner to the EM algorithm. Importantly, to evaluate our objective function, the only matrix we need to invert is  $\mathbf{K}_{mm}$ . When  $N$  is large and  $M$  is relatively small, the approach leads to significant computational savings compared to standard GPs.

## 4 Example Python Code

Code accompanying this tutorial is written in Python (version 2), is heavily commented and is designed to be readable. The module ‘LIRU\_SparseGP’ contains the functions needed to run a simple Sparse GP, while the script ‘1D\_example’ will run the Sparse GP on a test case that has univariate inputs. Specifically, ‘1D\_example’ will generate a full set of 500 training points, before attempting to fit a Sparse GP using 10 points only. To do so, it randomly selects 100 candidate sets of  $\mathbf{X}_m$ , before choosing the set that maximises the lower bound,  $L_{lb}$ . Once  $\mathbf{X}_m$  has been chosen, scipy’s ‘optimize.minimize’ module is used to minimise the (negative) lower bound with respect to the GP hyperparameters. Figure 1 shows an example of the results you should expect to see.

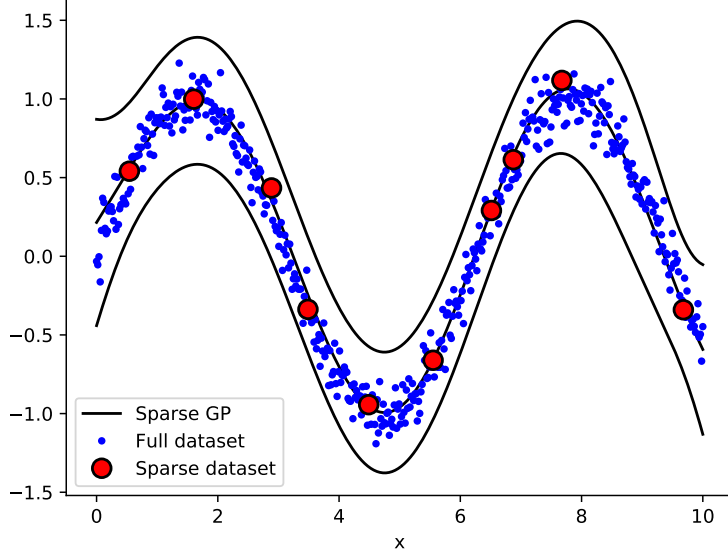


Figure 1: Example output of ‘1D\_example’ Python script.

Note that a more thorough approach to this problem would be to recursively add additional training points and update the GP hyperparameters. In [1][2] it is shown that the lower bound,  $L_{lb}$ , must increase every time an additional training point is included. This is convenient, as it allows the user to analyse convergence - new training points and hyperparameter updates can be run until  $L_{lb}$  is judged to have converged. For now though, for the sake of simplicity, our example code simply tries to identify a fixed number of training points (and optimum hyperparameters) for the Sparse GP implementation.

## 5 Detailed Methodology

Defining a prior  $p(\mathbf{f}, \mathbf{f}_m)$  and likelihood  $p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m)$  gives us the posterior

$$p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m)p(\mathbf{f}, \mathbf{f}_m)}{p(\mathbf{y})} \quad (8)$$

We decide that the optimum choice of  $\mathbf{f}_m$  would lead to

$$p(\mathbf{y} | \mathbf{f}_m) = p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m) \quad (9)$$

In other words, if we knew the optimum  $\mathbf{f}_m$  then, regarding the probability of witnessing  $\mathbf{y}$ ,  $\mathbf{f}$  doesn’t add any additional information<sup>1</sup>. It turns out that the optimum  $\mathbf{f}_m$  will factorise the posterior, so that we can write  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = p(\mathbf{f} | \mathbf{y})p(\mathbf{f}_m | \mathbf{y})$ . Proof:

$$p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{f}_m)p(\mathbf{f}, \mathbf{f}_m)}{p(\mathbf{y})} \quad (10)$$

$$= \frac{p(\mathbf{y} | \mathbf{f}_m)p(\mathbf{f}_m | \mathbf{f})p(\mathbf{f})}{p(\mathbf{y})} \times \frac{p(\mathbf{f}_m)}{p(\mathbf{f}_m)} \quad (11)$$

<sup>1</sup>Sometimes in the literature optimality is defined as  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = p(\mathbf{f} | \mathbf{f}_m)$ . In fact, expanding both sides using Bayes’ theorem you can show that  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = p(\mathbf{f} | \mathbf{f}_m) \implies p(\mathbf{y} | \mathbf{f}_m) = p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m)$  so it doesn’t make any difference. I prefer stating that  $p(\mathbf{y} | \mathbf{f}_m) = p(\mathbf{y} | \mathbf{f}, \mathbf{f}_m)$  is the optimal condition as it seems a bit more intuitive.

$$= \frac{p(\mathbf{f}_m | \mathbf{f})p(\mathbf{f})}{p(\mathbf{f}_m)} \frac{p(\mathbf{y} | \mathbf{f}_m)p(\mathbf{f}_m)}{p(\mathbf{y})} \quad (12)$$

$$= p(\mathbf{f} | \mathbf{f}_m)p(\mathbf{f}_m | \mathbf{y}) \quad (13)$$

So in the optimal case we can factorise the posterior pdf as

$$p(\mathbf{f}, \mathbf{f}_m | \mathbf{y}) = p(\mathbf{f}_m | \mathbf{y})p(\mathbf{f} | \mathbf{f}_m) \quad (14)$$

For this reason, we choose a variational distribution *which can be factorised in the same way* by defining

$$q(\mathbf{f}, \mathbf{f}_m) = \phi(\mathbf{f}_m)p(\mathbf{f} | \mathbf{f}_m) \quad (15)$$

and we try to minimise the KL divergence between  $q(\mathbf{f}, \mathbf{f}_m)$  and the posterior  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y})$ . This is based on the hypothesis that, if we find inputs  $\mathbf{X}_m$  that make  $q(\mathbf{f}, \mathbf{f}_m) = \phi(\mathbf{f}_m)p(\mathbf{f} | \mathbf{f}_m)$  and  $p(\mathbf{f}, \mathbf{f}_m | \mathbf{y})$  exactly the same, then it must be possible to factorise the posterior distribution into  $p(\mathbf{f}_m | \mathbf{y})p(\mathbf{f} | \mathbf{f}_m)$ , thus implying that our choice of  $\mathbf{X}_m$  is optimal. Noting that minimising

$$\text{KL}(q(\mathbf{f}, \mathbf{f}_m) || p(\mathbf{f}, \mathbf{f}_m | \mathbf{y})) \quad (16)$$

is equivalent to maximising

$$L = \int \int q(\mathbf{f}, \mathbf{f}_m) \log \left[ \frac{p(\mathbf{f}, \mathbf{f}_m, \mathbf{y})}{q(\mathbf{f}, \mathbf{f}_m)} \right] d\mathbf{f} d\mathbf{f}_m \quad (17)$$

then, by enforcing the factorisation property of the variational distribution ( $q(\mathbf{f}, \mathbf{f}_m) = \phi(\mathbf{f}_m)p(\mathbf{f} | \mathbf{f}_m)$ ) and writing  $p(\mathbf{f}, \mathbf{f}_m, \mathbf{y}) = p(\mathbf{y} | \mathbf{f})p(\mathbf{f} | \mathbf{f}_m)p(\mathbf{f}_m)$  we have

$$L = \int \int p(\mathbf{f} | \mathbf{f}_m)\phi(\mathbf{f}_m) \log \left\{ \frac{p(\mathbf{y} | \mathbf{f})p(\mathbf{f}_m)}{\phi(\mathbf{f}_m)} \right\} d\mathbf{f} d\mathbf{f}_m \quad (18)$$

$$= \int \phi(\mathbf{f}_m) \left[ \underbrace{\int p(\mathbf{f} | \mathbf{f}_m) \log p(\mathbf{y} | \mathbf{f}) d\mathbf{f}}_{\log G(\mathbf{f}_m, \mathbf{y})} + \log \left( \frac{p(\mathbf{f}_m)}{\phi(\mathbf{f}_m)} \right) \right] d\mathbf{f}_m \quad (19)$$

To evaluate  $\log G = \int p(\mathbf{f} | \mathbf{f}_m) \log p(\mathbf{y} | \mathbf{f}) d\mathbf{f}$ , we first need to find an expression for  $p(\mathbf{f} | \mathbf{f}_m)$ . Noting that

$$p(\mathbf{f}, \mathbf{f}_m) = \mathcal{N} \left( \begin{pmatrix} \mathbf{f} \\ \mathbf{f}_m \end{pmatrix}; \mathbf{0}, \begin{bmatrix} \mathbf{K}_{nn} & \mathbf{K}_{nm} \\ \mathbf{K}_{mn} & \mathbf{K}_{mm} \end{bmatrix} \right) \quad (20)$$

where  $\mathbf{K}_{nn}(i, j) = k(\mathbf{x}_i, \mathbf{x}_j)$ ,  $\mathbf{K}_{nm}(i, j) = k(\mathbf{x}_i, \mathbf{x}_{m,j})$  and  $\mathbf{K}_{mm}(i, j) = k(\mathbf{x}_{m,i}, \mathbf{x}_{m,j})$ . Using standard properties of Gaussian distributions we can then show that

$$p(\mathbf{f} | \mathbf{f}_m) = \mathcal{N}(\mathbf{f}; \mathbf{0}, \text{E}[\mathbf{f} | \mathbf{f}_m], \text{Cov}[\mathbf{f} | \mathbf{f}_m]) \quad (21)$$

where

$$\text{E}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{f}_m \quad (22)$$

and

$$\text{Cov}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nn} - \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \quad (23)$$

We can now try to evaluate  $\log G$ . To simplify notation, we first define  $\boldsymbol{\mu} = \mathbb{E}[\mathbf{f} | \mathbf{f}_m]$  and  $\boldsymbol{\Sigma} = \text{Cov}[\mathbf{f} | \mathbf{f}_m]$  such that the aim is to evaluate

$$\log G = \int \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \log p(\mathbf{y} | \mathbf{f}) d\mathbf{f} \quad (24)$$

where

$$\log p(\mathbf{y} | \mathbf{f}) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_i)^2 \quad (25)$$

therefore

$$\log G = \int \mathcal{N}(\mathbf{f}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) \left[ -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_i)^2 \right] d\mathbf{f} \quad (26)$$

$$= -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N \underbrace{\left[ \int \mathcal{N}(f_i; \mu_i, \Sigma_{ii}) (y_i - f_i)^2 df_i \right]}_I \quad (27)$$

Noting that we can write

$$(y_i - f_i)^2 = (y_i^2 - 2f_i y_i + (f_i - \mu_i)^2) + 2f_i \mu_i - \mu_i^2 \quad (28)$$

then the intergral  $I$  becomes

$$y_i^2 - 2\mu_i y_i + \mu_i^2 + \Sigma_{ii} = (y_i - \mu_i)^2 + \Sigma_{ii} \quad (29)$$

therefore

$$\log G = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N [(y_i - \mu_i)^2 + \Sigma_{ii}] \quad (30)$$

$$= \log(\mathcal{N}(\mathbf{y}; \boldsymbol{\mu}, \mathbf{I}_N \sigma^2)) - \frac{1}{2\sigma^2} \text{Tr}(\boldsymbol{\Sigma}) \quad (31)$$

This is a good point to summarise what we've done so far! We want to find the inputs,  $\mathbf{X}_m$ , that maximise

$$L = \int \phi(\mathbf{f}_m) \left[ \log G(\mathbf{f}_m, \mathbf{y}) + \log \left( \frac{p(\mathbf{f}_m)}{\phi(\mathbf{f}_m)} \right) \right] d\mathbf{f}_m \quad (32)$$

where

$$\log G = \log(\mathcal{N}(\mathbf{y}; \mathbb{E}[\mathbf{f} | \mathbf{f}_m], \mathbf{I}_N \sigma^2)) - \frac{1}{2\sigma^2} \text{Tr}(\text{Cov}[\mathbf{f} | \mathbf{f}_m]) \quad (33)$$

$$\mathbb{E}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{f}_m \quad (34)$$

and

$$\text{Cov}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nn} - \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn} \quad (35)$$

Noting that we can write

$$L = \int \phi(\mathbf{f}_m) \left[ \log \left( G(\mathbf{f}_m, \mathbf{y}) \frac{p(\mathbf{f}_m)}{\phi(\mathbf{f}_m)} \right) \right] d\mathbf{f}_m \quad (36)$$

then, using Jensen's inequality, we can now calculate a lower bound on  $L$ :

$$L \geq L_{lb} = \log \left[ \int \phi(\mathbf{f}_m) \frac{G(\mathbf{f}_m, \mathbf{y}) p(\mathbf{f}_m)}{\phi(\mathbf{f}_m)} d\mathbf{f}_m \right] \quad (37)$$

$$= \log \int G(\mathbf{f}_m, \mathbf{y}) p(\mathbf{f}_m) d\mathbf{f}_m \quad (38)$$

Note that, to maximise  $L_{lb}$ , we don't actually have to find the variational distribution  $\phi(\mathbf{f}_m)$  (a nice trick!). To evaluate the lower bound:

$$L_{lb} = \log \underbrace{\left[ \int \mathcal{N}(\mathbf{y}; \mathbf{E}[\mathbf{f} | \mathbf{f}_m], \mathbf{I}_N \sigma^2) p(\mathbf{f}_m) d\mathbf{f}_m \right]}_{I_2} - \frac{1}{2\sigma^2} \text{Tr}(\text{Cov}[\mathbf{f} | \mathbf{f}_m]) \quad (39)$$

The integral  $I_2$  can be evaluated using standard expressions for Gaussian distributions (see Appendix A) such that, finally, we have our objective function:

$$L_{lb} = \log [\mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}_n \sigma^2 + \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn})] - \frac{1}{2\sigma^2} \text{Tr}(\text{Cov}[\mathbf{f} | \mathbf{f}_m]) \quad (40)$$

where  $\text{Cov}[\mathbf{f} | \mathbf{f}_m] = \mathbf{K}_{nn} - \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}$ .

## 6 Some Interesting References

[5] produced a Variational Sparse GP that could also treat heteroscedastic problems. [6] specifically considered classification but, in the process, looked into how a Variational Sparse GP could be applied when there are non-conjugate likelihoods. [7] introduced ‘Deep Gaussian Processes’ (A deep belief network based on GP mappings), where variational GPs was central to the approach. They showed that Deep GPs are applicable to large data sets (like other deep methods but also small data sets (unlike other deep methods). This was illustrated on a case study with more dimensions than data. [8] produced a variant of Variational Sparse GPs that is suitable for stochastic variational inference. This was applied to datasets with around 700,000 points. [9] looked at Deep GPs again and proposed a ‘nested’ approach to variational inference that makes them more suitable for parallel computer architectures. Notes the danger of overfitting when using deep methods. [10] proposed a Variational Sparse GP that is suitable for a Map-Reduce setting and applied to datasets with around 2 million data points (implemented on open-source MapReduce software in Python).

It is worth noting that, while these datasets are quite large, [11] presented the ‘robust Bayesian Committee Machine’ (extension of a product-of-experts model) that was applicable to data sets of the order  $10^7$ . They also have some interesting criticisms of Sparse GPs, particularly with regard to the number of parameters that need to be trained.

Finally, Sparse GP code is available for free through the University of Sheffield’s Python module GPpy. Other interesting variants of GPs for large data sets can also be found in [12][13].

## A Evaluating equation (39)

Generally speaking, for the situation where

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a; \boldsymbol{\mu}, \mathbf{A}^{-1}), \quad p(\mathbf{x}_b | \mathbf{x}_a) = \mathcal{N}(\mathbf{x}_b; \mathbf{H} \mathbf{x}_a + \mathbf{b}, \mathbf{L}^{-1}) \quad (41)$$

it can be shown that

$$p(\mathbf{x}_b) = \int p(\mathbf{x}_b | \mathbf{x}_a) p(\mathbf{x}_a) d\mathbf{x}_a = \mathcal{N}(\mathbf{x}_b; \mathbf{H} \boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} \mathbf{H} \mathbf{A}^{-1} \mathbf{H}^T) \quad (42)$$

This is proved, for example, in Bishop's book [14] and also forms part of the derivation of the well-known Kalman filter equations. Applying this generalised result to the integral in equation (39), we find the following:

$$\int \mathcal{N}(\mathbf{y}; \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{f}_m, \mathbf{I}_n \sigma^2) \mathcal{N}(\mathbf{f}_m; \mathbf{0}, \mathbf{K}_{mm}) d\mathbf{f}_m \quad (43)$$

$$= \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}_n \sigma^2 + \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mm} \mathbf{K}_{mm}^{-T} \mathbf{K}_{mn}) \quad (44)$$

$$= \mathcal{N}(\mathbf{y}; \mathbf{0}, \mathbf{I}_n \sigma^2 + \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}) \quad (45)$$

## References

- [1] Michalis Titsias. Variational Learning of Inducing Variables in Sparse Gaussian Processes. *Artificial Intelligence and Statistics*, 5:567–574, 2009.
- [2] Michalis K Titsias. Variational Model Selection for Sparse Gaussian Process Regression. (2006):1–20, 2009.
- [3] David JC MacKay and David JC Mac Kay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
- [4] Christopher M Bishop. *Pattern Recognition and Machine Learning*, volume 53. 2013.
- [5] Miguel Lazaro-Gredilla and Michalis Titsias. Variational Heteroscedastic Gaussian Process Regression. *Proceedings of International Conference on Machine Learning (ICML 2011)*, Bellevue, Washington, USA, page 8, 2011.
- [6] Kian Ming A. Chai. Variational Multinomial Logit Gaussian Process. *J. Mach. Learn. Res.*, 98888:1745–1808, 2012.
- [7] Andreas C. Damianou and Neil D. Lawrence. Deep Gaussian Processes. *Artificial Intelligence and Statistics*, 31, 2012.
- [8] James Hensman, N Fusi, and Neil D. Lawrence. Gaussian Processes for Big Data. *Uncertainty in Artificial Intelligence*, pages 282–290, 2013.
- [9] James Hensman and Neil D. Lawrence. Nested Variational Compression in Deep Gaussian Processes. pages 1–21, 2014.
- [10] Yarin Gal, Mark van der Wilk, and Carl Rasmussen. Distributed Variational Inference in Sparse Gaussian Process Regression and Latent Variable Models. *Proc. Neural Information Processing Systems (NIPS)*, pages 1–9, 2014.
- [11] Marc Peter Deisenroth and Jun Wei Ng. Distributed Gaussian Processes. *International Conference on Machine Learning*, 37, 2015.



- [12] Andrew Gordon Wilson, Christoph Dann, and Hannes Nickisch. Thoughts on Massively Scalable Gaussian Processes. pages 1–25, 2015.
- [13] Andrew Gordon Wilson and Hannes Nickisch. Kernel Interpolation for Scalable Structured Gaussian Processes (KISS-GP). 37, 2015.
- [14] Christopher M Bishop. *Pattern recognition and machine learning*. springer, 2006.