

Imperial College  
London

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

---

# Identification of Gaussian Process State-Space Models

---

*Author:*  
William Sternberg

*Supervisor:*  
Marc Peter Deisenroth

Submitted in partial fulfilment of the requirements for the MSc degree in  
Computing (Machine Learning) of Imperial College London

September 2017



## Abstract

Gaussian process state-space models (GPSSMs) have been shown to be a competitive way of learning from time series due to their ability to identify complex systems; as a result, GPSSMs have the potential to be applied in areas such as Economics, Engineering and Physics. However, the current techniques for learning GPSSMs suffer from a variety of issues; for example, some methods have long training times or constraints on expressiveness while other methods have difficulties in dealing with high-dimensional latent spaces.

We provide an overview of all the key background material and focus on learning GPSSMs using a recent method that combines Hilbert reduced-rank Gaussian processes and sequential Monte Carlo. We also present several novel contributions, which add additional features to this recent method; for example, online learning, distributed learning, learning under Student-t noise and learning with Student-t processes. Additionally, we improve the Hilbert reduced-rank Gaussian process model using neural networks and state-of-the-art MCMC methods. This gives us a significant improvement in training times for high-dimensional inputs compared with the original Hilbert reduced-rank model. Finally, we use our improved reduced-rank Gaussian process model to create a novel GPSSM. We demonstrate all these contributions in a range of examples.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Background</b>	<b>5</b>
2.1	Gaussian Processes . . . . .	5
2.1.1	Gaussian Process Regression . . . . .	7
2.1.2	Reducing the Computational Cost . . . . .	10
2.1.3	Examples . . . . .	16
2.2	State-Space Models . . . . .	19
2.2.1	Inference: The General Framework . . . . .	20
2.2.2	Inference: Deterministic Methods . . . . .	21
2.2.3	Inference: Stochastic Methods . . . . .	24
2.2.4	Learning . . . . .	28
2.2.5	Examples . . . . .	29
2.3	Gaussian Process State-Space Models . . . . .	33
2.3.1	Inference in GPSSMs . . . . .	34
2.3.2	Learning in GPSSMs . . . . .	35
<b>3</b>	<b>Contributions</b>	<b>39</b>
3.1	Extending Hilbert reduced-rank Gaussian Processes . . . . .	39
3.1.1	The Adaptive Kernel . . . . .	39
3.1.2	The Manifold-Hilbert reduced-rank Gaussian Process . . . . .	41
3.1.3	Examples . . . . .	42
3.2	Extending Hilbert reduced-rank GPSSMs . . . . .	50
3.2.1	Online Learning . . . . .	50
3.2.2	Forgetful Online Learning . . . . .	51
3.2.3	Distributed Learning . . . . .	52
3.2.4	GPSSMs with Student-t noise . . . . .	56
3.2.5	STPSSMs: A State-Space Model with a Heavy Tailed Process . . . . .	60
3.2.6	Examples . . . . .	64
3.3	Learning DHGPMs . . . . .	71
3.3.1	The Deep Hilbert Gaussian Process Model (DHGPM) . . . . .	75
3.3.2	DHGPMs with non-Gaussian noise and Student-t processes . . . . .	79
3.3.3	Examples . . . . .	81
3.4	Learning DHGPSSMs . . . . .	92
3.4.1	The Deep Hilbert Gaussian Process State-Space Model . . . . .	92
3.4.2	Examples . . . . .	101
<b>4</b>	<b>Conclusion</b>	<b>109</b>
4.0.3	Conclusions . . . . .	109
4.0.4	Future Work . . . . .	110

# Chapter 1

## Introduction

Consider a sequence of time-ordered observations  $y_1, \dots, y_T$ . Our aim is to learn something about the system which has generated these observations and to predict future observations. Most interesting systems contain noise and this means that we cannot predict the future observations exactly even if we have perfect knowledge of the system; however, we can search for a predictive distribution. This distribution should ideally take into account all sources of uncertainty; for example, the system noise and the uncertainty in any estimated parameters. Although we are looking for a predictive distribution, we can compute statistics of this distribution such as the mean and variance in order to provide a point prediction and a confidence interval. The idea is that, though the point prediction might not be correct, the true value should have a high probability of falling inside the confidence interval.

As an example, consider the following autoregressive model known as an AR(1) model:

$$y_{t+1} = 0.5y_t + v_t \quad (1.1)$$

$$v_t \stackrel{iid}{\sim} \mathcal{N}(0, 1) \quad (1.2)$$

Despite the fact that we know all the parameters of this system, given  $y_{1:T} = y_1, \dots, y_T$  we cannot predict  $y_{T+1}$  with certainty; however, there does exist an optimal prediction. Let  $\hat{y}_{T+1}$  be our prediction of  $y_{T+1}$  given  $y_{1:T}$ , then the optimal  $\hat{y}_{T+1}$  is [Hamilton, 1994]:

$$\hat{y}_{T+1} = \mathbb{E}[y_{T+1}] \quad (1.3)$$

$$= 0.5y_T \quad (1.4)$$

Although (1.4) is the optimal prediction of  $y_{T+1}$ , the true value might be far away from this prediction and that is why we also wish to provide a distribution over the possible values of  $y_{T+1}$ . In this example, we can calculate the optimal predictive distribution as:

$$y_{T+1}|y_{1:T} \sim \mathcal{N}(0.5y_T, 1) \quad (1.5)$$

Now,  $y_{T+1}|y_{1:T} \sim \mathcal{N}(0.5y_T, 100)$  is also a *valid* predictive distribution but it is not the optimal one because it has more uncertainty than necessary.

Suppose next that we have observations  $y_1, \dots, y_T$  from the AR(1) model:

$$y_{t+1} = \alpha y_t + v_t \quad (1.6)$$

$$v_t \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \quad (1.7)$$

We could try and estimate the parameters  $\alpha$  and  $\sigma^2$  using the data  $y_1, \dots, y_T$  and then provide a predictive distribution of  $y_{T+1}$  by assuming that the estimated parameters have no uncertainty; for example:

$$y_{T+1}|y_{1:T}, \hat{\alpha}, \hat{\sigma}^2 \sim \mathcal{N}(\hat{\alpha}y_T, \hat{\sigma}^2) \quad (1.8)$$

where  $\hat{\alpha}$  and  $\hat{\sigma}^2$  is our estimate of the parameters  $\alpha$  and  $\sigma^2$  respectively. However, this predictive distribution is overconfident and we should try to take into account the uncertainty that we have in the model parameters in order to get a more accurate predictive distribution. One way to do this is to use Bayesian statistics. If we assume that we have some prior distributions for  $\alpha$  and  $\sigma^2$ :

$$\alpha \sim p(\alpha) \quad (1.9)$$

$$\sigma^2 \sim p(\sigma^2) \quad (1.10)$$

then, under the Bayesian framework, we can find the posterior distributions of the parameters given the data using Bayes' Theorem:

$$p(\alpha, \sigma^2|y_{1:T}) = \frac{p(y_{1:T}|\alpha, \sigma^2)p(\alpha)p(\sigma^2)}{p(y_{1:T})} \quad (1.11)$$

where we have assumed the prior distributions are independent. Using this, we can find a posterior predictive distribution of  $y_{T+1}|y_{1:T}$  which takes into account the uncertainty in the model parameters:

$$p(y_{T+1}|y_{1:T}) = \iint p(y_{T+1}|y_{1:T}, \alpha, \sigma^2)p(\alpha, \sigma^2|y_{1:T})d\alpha d\sigma^2 \quad (1.12)$$

$$= \iint p(y_{T+1}|y_T, \alpha, \sigma^2) \frac{p(y_{1:T}|\alpha, \sigma^2)p(\alpha)p(\sigma^2)}{p(y_{1:T})} d\alpha d\sigma^2 \quad (1.13)$$

$$= \frac{1}{p(y_{1:T})} \iint \mathcal{N}(\alpha y_T, \sigma^2)p(y_{1:T}|\alpha, \sigma^2)p(\alpha)p(\sigma^2)d\alpha d\sigma^2 \quad (1.14)$$

As you can see, it can get complicated quickly and with Bayesian statistics you often have intractable integrals: we will look at methods for dealing with these in Chapter 2.

### Why look at Gaussian process state-space models?

So far we have been focusing on predictive distributions rather than point estimates, and one of our many reasons for pursuing Gaussian process state-space models (GPSSMs) is that they naturally fit into this Bayesian framework. Furthermore, the more recent techniques for learning GPSSMs provide predictive distributions that take into account many different sources of uncertainty: from the observation

and latent noise, to the uncertainty in the latent states and kernel hyperparameters. Moreover, Gaussian process state-space models have been shown; for example, in Frigola [2015], to be a competitive method of learning from time series data. We believe there are two main reasons for this; firstly, time series data (especially if highly autocorrelated) does not contain as much information as the same amount of independent data but it is known that Gaussian processes are efficient learners in low data environments [Rasmussen and Williams, 2006]. Secondly, Gaussian processes are very flexible and this makes them ideal for putting inside a state-space model since often the interesting time series are rather complicated.

### **Other methods of learning time-ordered data:**

There are many techniques other than GPSSMs available in the machine learning and statistical literature for dealing with time series; here, we will briefly discuss a few of them. One of the simplest models is the autoregressive AR model, this is linear and it can learn simple systems via maximum likelihood or expectation-maximisation [Hamilton, 1994]. Other linear models include moving-average models (MA), autoregressive moving-average models (ARMA) and autoregressive integrated moving-average models (ARIMA). The key points for all the above linear models are that; firstly, they can be written as state-space models; secondly, they can only learn simple systems and thirdly, traditional forecasting methods for these models only take into account the system noise and do not include the uncertainty in the estimated parameters. More complex non-linear models include Generalized Autoregressive Conditional Heteroskedasticity (GARCH) models which allow for stochastic variance, and Holt-Winters/Exponential Smoothing models which are good at forecasting simple systems. However, none of these models are particularly expressive and they have difficulties dealing with high-dimensional data. A final example is recurrent neural networks, they have been shown to learn some complex systems but have long training times and only give point estimates. Overall, there seems to be a lack of fully Bayesian flexible models and GPSSMs aim to fill this gap.

The layout of this thesis is as follows, first we look at the key background material required to understand GPSSMs such as Gaussian processes, reduced-rank Gaussian processes, sequential Monte Carlo and state-space models. Then, we present our novel contributions with examples throughout to compare and contrast our new models with the previous ones. A summary of our contributions is stated below.

### **Contributions**

1. Created a new covariance function called the *Adaptive* covariance which is particularly suited to the Hilbert reduced-rank GP of Solin and Särkkä [2014] and useful for situations in which it is difficult to choose a good covariance function.
2. Combined the mGP of Calandra et al. [2016] with the Hilbert reduced-rank GP of Solin and Särkkä [2014] to create the manifold-Hilbert reduced-rank GP and we show how this model retains the properties of the original mGP but is much faster due to its use of reduced-rank GPs.

3. Added new features to the model of Svensson et al. [2016] such as learning with Student-t noise, online learning, forgetful online learning, distributed learning and learning with Student-t processes.
4. Improved the Hilbert reduced-rank Gaussian process model of Solin and Särkkä [2014] to allow it to scale much better as the input dimension increases. We call this model the deep Hilbert Gaussian process model (DHGPM).
5. Used the DHGPM to create a new GPSSM that can learn systems with high-dimensional latent spaces quicker than the model of Svensson et al. [2016].



# Chapter 2

## Background

### 2.1 Gaussian Processes

**Definition 2.1.** [Rasmussen and Williams, 2006, p. 13] A Gaussian process (GP) is a collection of random variables, any finite number of which have a joint Gaussian distribution.

As a result of these Gaussian *finite dimensional distributions* we can define a Gaussian process solely in terms of a mean function  $m(\cdot)$  and a covariance function  $k(\cdot, \cdot)$ .

**Definition 2.2.** [Rasmussen and Williams, 2006] A *scalar-valued covariance function* or *scalar-valued kernel* is a function  $k(\mathbf{x}, \mathbf{x}')$  with two inputs  $\mathbf{x} \in \mathbb{R}^{n_x}$  and  $\mathbf{x}' \in \mathbb{R}^{n_x}$  such that:

1.  $k(\mathbf{x}, \mathbf{x}')$  is real:  $k(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$ .
2.  $k(\mathbf{x}, \mathbf{x}')$  is symmetric:  $k(\mathbf{x}, \mathbf{x}') = k(\mathbf{x}', \mathbf{x}) \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^{n_x}$ .
3.  $k(\mathbf{x}, \mathbf{x}')$  is positive definite: given any  $n$  inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n$ , the matrix  $\mathbf{K} \in \mathbb{R}^{n_x \times n_x}$  with entries  $\mathbf{K}_{ij} = k(\mathbf{x}_i, \mathbf{x}_j) \in \mathbb{R}$  is positive definite. This matrix is known as the covariance or Gram matrix.

An example of a scalar valued covariance function is the squared exponential (SE) covariance function:

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l^2}\right) \in \mathbb{R} \quad (2.1)$$

where  $l$  is a hyperparameter and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{n_x}$ . See chapter 4 of Rasmussen and Williams [2006] for many more examples of covariance functions.

**Definition 2.3.** [Álvarez et al., 2012] A *multi-output covariance function* or *multi-output kernel* is a function  $k(\mathbf{x}, \mathbf{x}') : \mathbb{R}^{n_x \times n_x} \rightarrow \mathbb{R}^{D \times D}$  such that each  $(d, d')$  entry of the output matrix is a scalar-valued covariance function with inputs  $(\mathbf{x}, d)$  and  $(\mathbf{x}', d')$ . The  $(d, d')$  entry of the output matrix is written as  $k(\mathbf{x}, \mathbf{x}')_{d, d'} \in \mathbb{R}$  and we note that we must have  $k(\mathbf{x}, \mathbf{x}')_{d, d'} = k(\mathbf{x}', \mathbf{x})_{d', d}$  since  $k(\mathbf{x}, \mathbf{x}')_{d, d'}$  is a covariance function with inputs  $(\mathbf{x}, d)$  and  $(\mathbf{x}', d')$ .

An example of a multi-output covariance function is:

$$k(\mathbf{x}, \mathbf{x}') = \begin{pmatrix} k(\mathbf{x}, \mathbf{x}')_{1,1} & k(\mathbf{x}, \mathbf{x}')_{1,2} \\ k(\mathbf{x}, \mathbf{x}')_{2,1} & k(\mathbf{x}, \mathbf{x}')_{2,2} \end{pmatrix} \quad (2.2)$$

$$k(\mathbf{x}, \mathbf{x}')_{d,d'} = \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2l_1^2}\right) \exp\left(-\frac{(d - d')^2}{2l_2^2}\right) \in \mathbb{R} \quad (2.3)$$

where  $l_1, l_2$  are hyperparameters and  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{n_x}$ . See Álvarez et al. [2012] for many more examples of multi-output covariance functions.

Although multi-output covariance functions are not that popular, we are interested in them because of their application to multi-output Gaussian processes which are particularly important for Gaussian process state-space models. It is possible to avoid using multi-output kernels but we would like to investigate whether they could be useful.

**Notation:**

Let  $f$  be a Gaussian process (GP) with mean function  $m(\cdot)$  and covariance function  $k(\cdot, \cdot)$ , then  $f$  is written as:

$$f \sim \mathcal{GP}(m, k) \quad (2.4)$$

Given a single input  $\mathbf{x} \in \mathbb{R}^{n_x}$ , we have that  $f(\mathbf{x}) \in \mathbb{R}^D$  and we will write  $f_d(\mathbf{x}) \in \mathbb{R}$  for the  $d$ th component of  $f(\mathbf{x})$ . Also, given a single input  $\mathbf{x}$ , the mean function  $m(\mathbf{x})$  is related to  $f(\mathbf{x})$  by:

$$m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})] = \mathbb{E} \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_D(\mathbf{x}) \end{bmatrix} \in \mathbb{R}^D \quad (2.5)$$

Given two inputs  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^{n_x}$ , the covariance function  $k(\mathbf{x}, \mathbf{x}')$  is related to  $f(\mathbf{x})$  by:

$$k(\mathbf{x}, \mathbf{x}') = \text{cov}(f(\mathbf{x}), f(\mathbf{x}')) = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))^T] \in \mathbb{R}^{D \times D} \quad (2.6)$$

where the matrix  $k(\mathbf{x}, \mathbf{x}')$  has components [Álvarez et al., 2012]:

$$k(\mathbf{x}, \mathbf{x}')_{d,d'} = \text{cov}(f_d(\mathbf{x}), f_{d'}(\mathbf{x}')) \in \mathbb{R} \quad (2.7)$$

In the general case, there are no independence assumptions in (2.7) and so it is a dense matrix; however, as we shall see in section 2.1.2 making certain assumptions about this matrix can reduce the computational cost.

The notation in (2.4) is used to represent the joint distribution of the random variables  $f(\mathbf{x})$  at all possible inputs  $\mathbf{x} \in \mathbb{R}^{n_x}$ . However, if we take a finite number of inputs  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{n_x \times N}$ , then by the definition of Gaussian process we have that:

$$\text{vec}(f(\mathbf{X})) \sim \mathcal{N}(\text{vec}(m(\mathbf{X})), K(\mathbf{X}, \mathbf{X})) \quad (2.8)$$

where  $\text{vec}(\mathbf{Z})$  is the vectorisation of the matrix  $\mathbf{Z}$  (an operation which stacks the columns of  $\mathbf{Z}$  on top of each other to form a vector),  $f(\mathbf{X}) = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)] \in \mathbb{R}^{D \times N}$ ,  $m(\mathbf{X}) = [m(\mathbf{x}_1), \dots, m(\mathbf{x}_N)] \in \mathbb{R}^{D \times N}$  and

$$K(\mathbf{X}, \mathbf{X}) = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{ND \times ND} \quad (2.9)$$

In the case that  $D = 1$  we can drop the  $\text{vec}$  and then this matches the notation in Rasmussen and Williams [2006].

As an example, suppose we are interested in finding the joint distribution of  $f(\mathbf{x}_1) \in \mathbb{R}^2$  and  $f(\mathbf{x}_2) \in \mathbb{R}^2$  where  $f \sim \mathcal{GP}(m, k)$  and we use a multi-output kernel. In this case, we have that:

$$\text{vec}(f(\mathbf{X})) = \begin{pmatrix} f_1(\mathbf{x}_1) \\ f_2(\mathbf{x}_1) \\ f_1(\mathbf{x}_2) \\ f_2(\mathbf{x}_2) \end{pmatrix} \quad (2.10)$$

$$\sim \mathcal{N} \left( \begin{bmatrix} m_1(\mathbf{x}_1) \\ m_2(\mathbf{x}_1) \\ m_1(\mathbf{x}_2) \\ m_2(\mathbf{x}_2) \end{bmatrix}, \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1)_{1,1} & k(\mathbf{x}_1, \mathbf{x}_1)_{1,2} & k(\mathbf{x}_1, \mathbf{x}_2)_{1,1} & k(\mathbf{x}_1, \mathbf{x}_2)_{1,2} \\ k(\mathbf{x}_1, \mathbf{x}_1)_{2,1} & k(\mathbf{x}_1, \mathbf{x}_1)_{2,2} & k(\mathbf{x}_1, \mathbf{x}_2)_{2,1} & k(\mathbf{x}_1, \mathbf{x}_2)_{2,2} \\ k(\mathbf{x}_2, \mathbf{x}_1)_{1,1} & k(\mathbf{x}_2, \mathbf{x}_1)_{1,2} & k(\mathbf{x}_2, \mathbf{x}_2)_{1,1} & k(\mathbf{x}_2, \mathbf{x}_2)_{1,2} \\ k(\mathbf{x}_2, \mathbf{x}_1)_{2,1} & k(\mathbf{x}_2, \mathbf{x}_1)_{2,2} & k(\mathbf{x}_2, \mathbf{x}_2)_{2,1} & k(\mathbf{x}_2, \mathbf{x}_2)_{2,2} \end{bmatrix} \right) \quad (2.11)$$

### 2.1.1 Gaussian Process Regression

In this section, we will provide an overview of Gaussian process regression following the derivation in Rasmussen and Williams [2006]. However, we will also explicitly deal with the general case of  $f(\mathbf{x}) \in \mathbb{R}^D$ . In many ways, the general derivation is similar to the univariate case ( $D = 1$ ) and the results are well known but we aim to make clear the link between the different forms of the density functions which can arise in the general case. These different forms are a result of a relationship between the matrix-normal and multivariate-normal distributions.

Consider a set of data  $\mathcal{D} = \{(\mathbf{y}_i, \mathbf{x}_i) \mid i = 1, \dots, N\}$  where  $\mathbf{y}_i \in \mathbb{R}^{n_y}$  and  $\mathbf{x}_i \in \mathbb{R}^{n_x}$  for  $i = 1, \dots, N$ . In Gaussian process regression, the aim is to use this data  $\mathcal{D}$  to find the posterior distribution of a Gaussian process relating inputs  $\mathbf{x}_i$  to outputs  $\mathbf{y}_i$  and then to provide a predictive distribution for  $\mathbf{f}_* = f(\mathbf{x}_*)$  given a test input  $\mathbf{x}_*$ . Formally, this means that we have the following system ( $i = 1, \dots, N$ ):

$$\mathbf{y}_i = f(\mathbf{x}_i) + \boldsymbol{\epsilon}_i \quad (2.12a)$$

$$f \sim \mathcal{GP}(m, k) \quad (2.12b)$$

$$\boldsymbol{\epsilon}_i | \mathbf{Q} \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (2.12c)$$

$$p(\mathbf{y}_i | \mathbf{f}, \mathbf{x}_i, \boldsymbol{\theta}) \sim \mathcal{N}(f(\mathbf{x}_i), \mathbf{Q}) \quad (2.12d)$$

where  $\mathbf{Q} \in \mathbb{R}^{n_y \times n_y}$  (usually diagonal) and  $\boldsymbol{\theta}$  represents all the model hyperparameters which includes  $\mathbf{Q}$  and the hyperparameters of kernel  $k$ .

We wish to:

1. Learn the optimal hyperparameters  $\boldsymbol{\theta}_*$ .
2. Find the predictive distribution  $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}_*)$ .

where we have  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N] \in \mathbb{R}^{n_y \times N}$ ,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{n_x \times N}$ , test inputs  $\mathbf{X}_* = [\mathbf{x}_{*1}, \dots, \mathbf{x}_{*M}] \in \mathbb{R}^{n_x \times M}$  and function values  $\mathbf{f}_* = [f(\mathbf{x}_{*1}), \dots, f(\mathbf{x}_{*M})] \in \mathbb{R}^{n_y \times M}$ .

**Definition 2.4.** [Schäcke, 2013, p. 6]. Given matrices  $\mathbf{X} \in \mathbb{R}^{m \times n}$  and  $\mathbf{Y} \in \mathbb{R}^{p \times q}$  the Kronecker product is defined as:

$$\mathbf{X} \otimes \mathbf{Y} = \begin{pmatrix} x_{11}\mathbf{Y} & x_{12}\mathbf{Y} & \cdots & x_{1n}\mathbf{Y} \\ x_{21}\mathbf{Y} & x_{22}\mathbf{Y} & \cdots & x_{2n}\mathbf{Y} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1}\mathbf{Y} & x_{m2}\mathbf{Y} & \cdots & x_{mn}\mathbf{Y} \end{pmatrix} \in \mathbb{R}^{mp \times nq} \quad (2.13)$$

It has properties [Schäcke, 2013, pp. 7-9]:

1.  $(\mathbf{X} \otimes \mathbf{Y})(\mathbf{Z} \otimes \mathbf{T}) = \mathbf{XZ} \otimes \mathbf{YT}$   
for  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{Y} \in \mathbb{R}^{p \times q}$ ,  $\mathbf{Z} \in \mathbb{R}^{n \times r}$  and  $\mathbf{T} \in \mathbb{R}^{q \times s}$ .
2.  $(\mathbf{X} \otimes \mathbf{Y})^{-1} = \mathbf{X}^{-1} \otimes \mathbf{Y}^{-1}$  if  $\mathbf{X}$  and  $\mathbf{Y}$  are invertible.
3. Let  $\mathbf{X}$  and  $\mathbf{Y}$  have eigen-decompositions  $\mathbf{PAP}^{-1}$  and  $\mathbf{QBQ}^{-1}$  respectively where  $\mathbf{A}$  and  $\mathbf{B}$  are diagonal. Then the eigen-decomposition of  $\mathbf{X} \otimes \mathbf{Y}$  is  $(\mathbf{P} \otimes \mathbf{Q})(\mathbf{A} \otimes \mathbf{B})(\mathbf{P} \otimes \mathbf{Q})^{-1}$

**Definition 2.5.** [Gupta and Nagar, 2000] A random matrix  $\mathbf{X} \in \mathbb{R}^{n \times m}$  is said to have a matrix normal distribution  $\mathcal{MN}_{n,m}(\mathbf{M}, \mathbf{U}, \mathbf{V})$  if the density function satisfies:

$$p(\mathbf{X} | \mathbf{M}, \mathbf{U}, \mathbf{V}) = \frac{1}{(2\pi)^{\frac{nm}{2}} |\mathbf{V}|^{\frac{n}{2}} |\mathbf{U}|^{\frac{m}{2}}} \exp \left( -\frac{1}{2} \text{tr}[\mathbf{V}^{-1}(\mathbf{X} - \mathbf{M})^T \mathbf{U}^{-1}(\mathbf{X} - \mathbf{M})] \right) \quad (2.14)$$

where  $\mathbf{M} \in \mathbb{R}^{n \times m}$ ,  $\mathbf{U} \in \mathbb{R}^{n \times n}$ ,  $\mathbf{V} \in \mathbb{R}^{m \times m}$  and  $|\cdot|$  is the determinant.

An important result [Gupta and Nagar, 2000] is the relationship between the matrix normal and the multivariate normal distributions:

$$\mathbf{X} \sim \mathcal{MN}_{n,m}(\mathbf{M}, \mathbf{U}, \mathbf{V}) \iff \text{vec}(\mathbf{X}) \sim \mathcal{N}_{nm}(\text{vec}(\mathbf{M}), \mathbf{V} \otimes \mathbf{U}) \quad (2.15)$$

where  $\otimes$  is the Kronecker product and  $\text{vec}(\mathbf{X})$  is the vectorisation of  $\mathbf{X}$ .

Consider  $p(\mathbf{Y} | f(\mathbf{X}), \boldsymbol{\theta})$  where  $f(\mathbf{X}) = [f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)] \in \mathbb{R}^{n_y \times N}$ :

$$p(\mathbf{Y} | f(\mathbf{X}), \boldsymbol{\theta}) = \prod_{i=1}^N p(\mathbf{y}_i | f(\mathbf{x}_i), \boldsymbol{\theta}) \quad (2.16)$$

$$= \prod_{i=1}^N \mathcal{N}(\mathbf{y}_i | f(\mathbf{x}_i), \mathbf{Q}) \quad (2.17)$$

$$\propto \exp\left(-\frac{1}{2} \sum_{i=1}^N (\mathbf{y}_i - f(\mathbf{x}_i))^T \mathbf{Q}^{-1} (\mathbf{y}_i - f(\mathbf{x}_i))\right) \quad (2.18)$$

$$= \exp\left(-\frac{1}{2} [\text{vec}(\mathbf{Y}) - \text{vec}(f(\mathbf{X}))]^T (\mathbf{I}_N \otimes \mathbf{Q})^{-1} [\text{vec}(\mathbf{Y}) - \text{vec}(f(\mathbf{X}))]\right) \quad (2.19)$$

where  $\mathbf{I}_N$  is the  $N \times N$  identity matrix and to get (2.19) we used a standard result for the Kronecker product:  $\mathbf{X}^{-1} \otimes \mathbf{Y}^{-1} = (\mathbf{X} \otimes \mathbf{Y})^{-1}$  provided both  $\mathbf{X}$  and  $\mathbf{Y}$  are invertible.

From (2.19) we get that:

$$p(\text{vec}(\mathbf{Y})|f(\mathbf{X}), \boldsymbol{\theta}) = \mathcal{N}(\text{vec}(\mathbf{Y})| \text{vec}(f(\mathbf{X})), \mathbf{I}_N \otimes \mathbf{Q}) \quad (2.20)$$

And hence by (2.15) we get:

$$p(\mathbf{Y}|f(\mathbf{X}), \boldsymbol{\theta}) = \mathcal{MN}(\mathbf{Y}|f(\mathbf{X}), \mathbf{I}_N, \mathbf{Q}) \quad (2.21)$$

When we talk about the distribution of  $\mathbf{Y}$  what we are really talking about is the joint distribution of all the components  $Y_{ij}$  and so the distributions of  $\mathbf{Y}$  and  $\text{vec}(\mathbf{Y})$  must be the same: they contain the same random variables  $Y_{ij}$ . This means that at any particular test point  $\mathbf{y}$ , the densities  $p_{\mathbf{Y}}(\mathbf{y})$  and  $p_{\text{vec}(\mathbf{Y})}(\mathbf{y})$  are equal (note that the supports are also the same). However, the parametrisations will in general be different; for example  $\mathbf{Y}$  might follow a matrix-normal distribution and  $\text{vec}(\mathbf{Y})$  multivariate-normal distribution.

We can use the above results to get  $p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta})$ :

$$p(\mathbf{Y}|\mathbf{X}, \boldsymbol{\theta}) = \int p(\mathbf{Y}|f(\mathbf{X}), \boldsymbol{\theta}) p(f(\mathbf{X})|\boldsymbol{\theta}) df \quad (2.22)$$

$$= \int \mathcal{N}(\text{vec}(\mathbf{Y})| \text{vec}(f(\mathbf{X})), \mathbf{I}_N \otimes \mathbf{Q}) \mathcal{N}(\text{vec}(f(\mathbf{X}))| \text{vec}(m(\mathbf{X})), K(\mathbf{X}, \mathbf{X})) df \quad (2.23)$$

$$= \mathcal{N}(\text{vec}(\mathbf{Y})| \text{vec}(m(\mathbf{X})), K(\mathbf{X}, \mathbf{X}) + \mathbf{I}_N \otimes \mathbf{Q}) \quad (2.24)$$

To find the hyperparameters we maximise  $\log(p(\text{vec}(\mathbf{Y})|\mathbf{X}, \boldsymbol{\theta}))$  with respect to  $\boldsymbol{\theta}$ . In the standard marginal-likelihood way, this will find a balance between data fit and model complexity. See Rasmussen and Williams [2006, Ch. 5] for a discussion and detailed overview of how to find the optimal hyperparameters. Since  $\log(p(\text{vec}(\mathbf{Y})|\mathbf{X}, \boldsymbol{\theta}))$  is multivariate normal there are no additional technical difficulties over the univariate case; however, there may be more hyperparameters and bigger matrices involved so it can take more time for the optimisation algorithms to complete.

Finally using (2.8) and (2.24) we get:

$$\begin{pmatrix} \text{vec}(\mathbf{f}_*) \\ \text{vec}(\mathbf{Y}) \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \text{vec}(m(\mathbf{X}_*)) \\ \text{vec}(m(\mathbf{X})) \end{pmatrix}, \begin{pmatrix} K(\mathbf{X}_*, \mathbf{X}_*) & K(\mathbf{X}_*, \mathbf{X}) \\ K(\mathbf{X}, \mathbf{X}_*) & K(\mathbf{X}, \mathbf{X}) + \mathbf{I}_N \otimes \mathbf{Q} \end{pmatrix} \right) \quad (2.25)$$

where we note that:

$$\text{cov}(\text{vec}(\mathbf{f}_*), \text{vec}(\mathbf{Y})) \quad (2.26)$$

$$= \text{cov}(\text{vec}(f(\mathbf{X}_*)), \text{vec}(f(\mathbf{X}))) \quad (2.27)$$

$$= K(\mathbf{X}_*, \mathbf{X}) \in \mathbb{R}^{Mn_y \times Nn_y} \quad (2.28)$$

Hence, via Gaussian conditioning we get  $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{Y}, \mathbf{X}, \boldsymbol{\theta})$ :

$$p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{Y}, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\text{vec}(\mathbf{f}_*) | \mathbf{m}_*, \mathbf{V}_*) \quad (2.29)$$

with:

$$\mathbf{m}_* = \text{vec}(m(\mathbf{X}_*)) + K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \mathbf{I}_N \otimes \mathbf{Q})^{-1} \text{vec}(\mathbf{Y} - m(\mathbf{X})) \quad (2.30)$$

$$\mathbf{V}_* = K(\mathbf{X}_*, \mathbf{X}_*) - K(\mathbf{X}_*, \mathbf{X})(K(\mathbf{X}, \mathbf{X}) + \mathbf{I}_N \otimes \mathbf{Q})^{-1} K(\mathbf{X}, \mathbf{X}_*) \quad (2.31)$$

## 2.1.2 Reducing the Computational Cost

Both Gaussian process prediction (finding  $p(\mathbf{f}_* | \mathbf{X}_*, \mathbf{Y}, \mathbf{X}, \boldsymbol{\theta})$ ) and hyperparameter optimisation (training) involve inverting matrices of size  $Nn_y \times Nn_y$  and this requires  $\mathcal{O}(N^3n_y^3)$  operations in general. Therefore, in cases where we have a large amount of training data or in cases with high output dimension  $n_y$ , training and prediction can take a long time; furthermore, even storing the  $Nn_y \times Nn_y$  matrix in memory might not be possible [Rasmussen and Williams, 2006, p. 171]. In this section, we will briefly discuss various methods for dealing with these issues.

### Special multi-output kernels:

Consider the  $d, d'$  entry on the submatrix of  $K(\mathbf{X}, \mathbf{X})$  associated with  $\mathbf{x}, \mathbf{x}'$ :

$$k(\mathbf{x}, \mathbf{x}')_{d,d'} \in \mathbb{R} \quad (2.32)$$

This could be any scalar-valued covariance function with inputs:  $(\mathbf{x}, d)$  and  $(\mathbf{x}', d')$  [Álvarez et al., 2012, p. 8]. For example, a valid covariance function is the multi-output kernel from (2.3) which as we shall see is known as a separable kernel. Using a covariance function on  $(\mathbf{x}, d)$  and  $(\mathbf{x}', d')$  allows us to easily construct  $K(\mathbf{X}, \mathbf{X})$  but this leads to the covariance matrix being of size  $Nn_y \times Nn_y$ . To reduce the computational cost associated with the output dimension  $n_y$ , we can use special covariance functions [Álvarez et al., 2012]. A basic type of special covariance function is a separable covariance function; this is a covariance function of the form:

$$k(\mathbf{x}, \mathbf{x}')_{d,d'} = k_{\mathcal{X}}(\mathbf{x}, \mathbf{x}')k_D(d, d') \quad (2.33)$$

where  $k_{\mathcal{X}}(\mathbf{x}, \mathbf{x}') \in \mathbb{R}$  and  $k_D(d, d') \in \mathbb{R}$ . This leads to:

$$k(\mathbf{x}, \mathbf{x}') = k_{\mathcal{X}}(\mathbf{x}, \mathbf{x}')\mathbf{B} \quad (2.34)$$

$$K(\mathbf{X}, \mathbf{X}) = k_{\mathcal{X}}(\mathbf{X}, \mathbf{X}) \otimes \mathbf{B} \quad (2.35)$$

where  $\mathbf{B} \in \mathbb{R}^{n_y \times n_y}$  is a symmetric, positive definite matrix with entries  $B_{ij} = k_D(d_i, d_j)$  for  $i, j = 1, \dots, N$  and  $k_{\mathcal{X}}(\mathbf{X}, \mathbf{X})$  is a matrix with entries  $k_{\mathcal{X}}(\mathbf{X}, \mathbf{X})_{ij} =$

$k(\mathbf{x}_i, \mathbf{x}_j) \in \mathbb{R}$  for  $i, j = 1, \dots, N$  where  $k(\mathbf{x}_i, \mathbf{x}_j)$  is some scalar-valued covariance function. Immediately we can see that the storage costs have been reduced from  $\mathcal{O}(N^2 n_y^2)$  to  $\mathcal{O}(N^2 + n_y^2)$  and using standard properties of the Kronecker product we get:

$$K(\mathbf{X}, \mathbf{X})^{-1} = (k_{\mathcal{X}}(\mathbf{X}, \mathbf{X}) \otimes \mathbf{B})^{-1} = k_{\mathcal{X}}(\mathbf{X}, \mathbf{X})^{-1} \otimes \mathbf{B}^{-1} \quad (2.36)$$

and hence the number of operations for the inversion is reduced from  $\mathcal{O}(N^3 n_y^3)$  to  $\mathcal{O}(N^3 + n_y^3)$ . Using eigen-decompositions it is possible to invert  $K(\mathbf{X}, \mathbf{X}) + \mathbf{I}_N \otimes \mathbf{Q}$  when  $K(\mathbf{X}, \mathbf{X}) = k_{\mathcal{X}}(\mathbf{X}, \mathbf{X}) \otimes \mathbf{B}$  in  $\mathcal{O}(N^3 + n_y^3)$ . In the case where  $\mathbf{Q} = \sigma^2 \mathbf{I}$ , Álvarez et al. [2012, p. 26] provide a method to do just this by noting:

$$(k_{\mathcal{X}}(\mathbf{X}, \mathbf{X}) \otimes \mathbf{B} + \mathbf{I}_N \otimes \sigma^2 \mathbf{I}_{n_y})^{-1} \text{vec}(\mathbf{Y}) \quad (2.37)$$

$$= \sum_{d=1}^{n_y} \mathbf{u}_d [(\rho_d k_{\mathcal{X}}(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I}_N)^{-1} \bar{\mathbf{y}}^d]^T \quad (2.38)$$

where  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_N]$  and  $\rho_d, \mathbf{u}_d$  for  $d = 1, \dots, n_y$  are the eigenvalues/eigenvectors of  $\mathbf{B}$  and  $\bar{\mathbf{y}}^d = (\mathbf{y}_1^T \mathbf{u}_d, \dots, \mathbf{y}_N^T \mathbf{u}_d) \in \mathbb{R}^N$  for  $d = 1, \dots, n_y$ . This can also be applied in the case where  $\text{vec}(\mathbf{Y})$  is replaced by  $K(\mathbf{X}, \mathbf{X}_*)$  by using the same method on each of the columns of  $K(\mathbf{X}, \mathbf{X}_*)$ .

One thing we have not mentioned is how to find  $\mathbf{B}$ , GPy [2012] suggests setting  $\mathbf{B} = \mathbf{W}\mathbf{W}^T + \epsilon \mathbf{I}$  ( $\epsilon > 0$ ) and then viewing  $\mathbf{W}$  and  $\epsilon$  as covariance hyperparameters finding them in the usual way by maximising the marginal likelihood. This way is good because it ensures  $\mathbf{B}$  is positive definite and if  $\mathbf{W} \in \mathbb{R}^{n_y}$  then there are only  $n_y + 1$  additional hyperparameters to find.

Although we motivated these special covariance functions as a means to reduce computational cost, there are alternative constructions: the linear model of coregionalization (LMC) leads to the sums of separable covariances and a special case known as the intrinsic coregionalisation model (ICM) leads to the separable covariance case we have looked at above. Furthermore, it is possible to extend Gaussian process regression to the case where each output dimension has different training data and all the above is discussed in depth by Álvarez et al. [2012]. Sometimes the special covariance functions use techniques similar to the sparse GPs discussed below and introduce inducing point variables. It is worth pointing out that designing these special covariance functions is not only about computational cost (most sparse GP techniques can be used on top of the special covariance function if necessary) but also about creating covariance functions that can express the relationship between multiple output dimensions sufficiently; for example, see Álvarez and Lawrence [2009].

### Independent GPs for each output dimension:

Consider the dataset  $\mathcal{D} = \{(\mathbf{y}_i, \mathbf{x}_i) \mid i = 1, \dots, N\}$  with  $\mathbf{y}_i \in \mathbb{R}^{n_y}$  and  $\mathbf{x}_i \in \mathbb{R}^{n_x}$ . The independent GP for each output dimension approach splits  $\mathcal{D}$  up into datasets  $\mathcal{D}_j = \{(y_{ij}, \mathbf{x}_i) \mid i = 1, \dots, N\}$  for  $j = 1, \dots, n_y$  where  $y_{ij} \in \mathbb{R}$  is the  $j^{\text{th}}$  component of  $\mathbf{y}_i$ . This approach then trains independent univariate-output GPs on each of the

datasets  $\mathcal{D}_j$ . This removes the issue of multi-output kernels since all kernels will now be scalar-valued and this approach automatically reduces training times to  $\mathcal{O}(n_y N^3)$ . Moreover, by performing operations in parallel and reusing parts one can *effectively* reduce the training time to  $\mathcal{O}(N^3)$ .

A question that should be asked is whether the computational expense of multi-output covariance functions is necessary and worth it. The reason why we might wish to consider multi-output covariance functions is that the different dimensions can give information about each other: in other words the outputs are correlated.

Consider the following example:

$$\begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} f_1(x) \\ f_2(x) \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \end{pmatrix} = \begin{pmatrix} \tanh(x) \\ \tanh(x) \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \end{pmatrix} \quad (2.39)$$

where  $\epsilon_1, \epsilon_2 \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ . Now  $y_1$  and  $y_2$  are statistically independent, and  $f_1(x), f_2(x)$  are statistically independent (since  $x$  is a constant). That said knowing  $f_1(x)$  allows us to completely determine  $f_2(x)$  and vice versa. As a result, using a multi-output covariance function will give us a better predicted mean and variance for both  $f_1(x), f_2(x)$  compared with independent GPs since we can use information about  $f_1(x)$  to understand  $f_2(x)$  better and vice versa. To understand why this is the case, notice that using information about  $y_2$  for  $y_1$  provides two samples from  $\mathcal{N}(0, 1)$  at each  $x$  input instead of one and therefore for a single training point  $x$  we can use  $\frac{y_1 + y_2}{2}$  as a better (lower MSE) estimate of  $f_1$  than just  $y_1$ . Extending this to  $D$  output dimensions with a single training input  $x$  and  $y_i = \tanh(x) + \epsilon_i; i = 1, \dots, D$  where  $\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ , we see that as  $D \rightarrow \infty$ , we get that  $\frac{1}{D} \sum_{i=1}^D y_i \xrightarrow{a.s.} \mathbb{E}[y_i] = \tanh(x)$  by the strong law of large numbers.

In some cases, such as the above, we get better results using multi-output kernels; however, in general using independent GPs for each output dimension is a good choice which balances speed with model accuracy. Although it is an approximation, independent GPs can be a competitive alternative to multi-output kernels with the key advantage being a reduced computation cost; for example, see Turner et al. [2010]; Frigola [2015].

### Sparse Gaussian Processes:

Previously, we looked at reducing the cost associated with high output dimensions but now we will look at reducing the costs associated with a large amount of training data. One way to do this is to use sparse Gaussian process approximations and an extensive discussion can be found in Quiñonero-Candela and Rasmussen [2005]. The idea is to introduce a set  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_L]$  of latent (hidden) variables often called inducing variables where  $p(\mathbf{u}) \sim \mathcal{N}(\mathbf{0}, K(\mathbf{U}, \mathbf{U}))$  and then assume that  $\mathbf{f}_*$  and  $\mathbf{f}$  are independent given  $\mathbf{U}$ . This leads to [Quiñonero-Candela and Rasmussen, 2005]:

$$p(\mathbf{f}_*, \mathbf{f}) = \int p(\mathbf{f}_*, \mathbf{f} | \mathbf{u}) p(\mathbf{u}) d\mathbf{u} \quad (2.40)$$



$$\approx \int p(\mathbf{f}_*|\mathbf{u})p(\mathbf{f}|\mathbf{u})p(\mathbf{u})d\mathbf{u} \quad (2.41)$$

where  $p(\mathbf{f}_*, \mathbf{f})$  is the joint distribution of the training and test data. The key is to find a good set of inducing variables and the different sparse GP methods introduce different additional approximations on top of (2.41) as well as different methods for choosing the set  $\mathbf{U}$  of inducing variables [Quiñero-Candela and Rasmussen, 2005]. Examples include: *The Subset of Regressors*, *The Deterministic Training Conditional Approximation* and *The Fully Independent Conditional Approximation*. These are discussed in Quiñero-Candela and Rasmussen [2005] and a comparison of their performance is given in Rasmussen and Williams [2006]. Other methods include combining variational and sparse techniques; for example in Titsias [2009]. Furthermore, sometimes the computational cost associated with the amount of training data and the high output dimensions are dealt with in the same algorithm; for example, *Variational Inducing Kernels* [Álvarez et al., 2010].

### Hilbert reduced-rank Gaussian Processes:

The Hilbert reduced-rank Gaussian process approximation [Solin and Särkkä, 2014] (which we will also refer to as a reduced-rank Gaussian process) is an important Gaussian process approximation which does not fall under the above sparse GP framework. Instead of using inducing points, the Hilbert reduced-rank GP of Solin and Särkkä [2014] attempts to approximate the kernel. As we shall see, Hilbert reduced-rank GPs can be used to greatly reduce the computation cost and we are particularly interested in this approximation because of its application to GPSSMs.

For now, suppose we have a scalar-valued (i.e.  $k(x, x') \in \mathbb{R}$ ) and *isotropic stationary* covariance function (kernel) (i.e.  $k(\mathbf{x}, \mathbf{x}')$  is a function of  $\|\mathbf{x} - \mathbf{x}'\|$ ) with input dimension  $n_x = 1$  (i.e.  $x, x' \in \mathbb{R}$ ). Furthermore, suppose that all the training and test data lies in the domain  $[-L, L]$  for some  $L \in \mathbb{R}$ . In this case, the  $m \in \mathbb{Z}_{>0}$  basis function approximation of the covariance function (kernel) is [Solin and Särkkä, 2014]:

$$k(x, x') \approx \sum_{j=1}^m S(\sqrt{\lambda_j})\phi_j(x)\phi_j(x') \quad (2.42)$$

$$\phi_j(x) = \frac{1}{\sqrt{L}} \sin\left(\frac{\pi j(x+L)}{2L}\right) \quad (2.43)$$

$$\lambda_j = \left(\frac{\pi j}{2L}\right)^2 \quad (2.44)$$

where  $x, x' \in \mathbb{R}$ , the  $\lambda_j$  and  $\phi_j(x)$  are eigenvalues and eigenfunctions of the Laplace operator in the domain  $[-L, L]$  and  $S$  is the spectral density of the covariance function. This spectral density always exists in the case of a stationary covariance functions (Bochner's Theorem) and is calculated as [Solin and Särkkä, 2014]:

$$S(\boldsymbol{\omega}) = \int k(\mathbf{r}) \exp(-i\boldsymbol{\omega}^T \mathbf{r}) d\mathbf{r} \quad (2.45)$$

where  $r = \|\mathbf{x} - \mathbf{x}'\|$  and we note that since we have assumed the covariance function is isotropic stationary by definition it is a function of  $\|\mathbf{x} - \mathbf{x}'\|$ . The actual derivation of (2.42)–(2.45) is complex and will not be discussed here but can be found in Solin and Särkkä [2014].

There are two sources of approximation in (2.42); firstly, we only take finitely ( $m$ ) many eigenvalues/eigenfunctions whereas they are infinity many and secondly, the domain is bounded by  $[-L, L]$ . The bounding of the domain means that even in the limit as  $m \rightarrow \infty$  this is still an approximation [Solín and Särkkä, 2014] and also means that we get edge effects (i.e. poor performance) near the domain boundary: we really do need the training and test data to be well within this domain boundary. As an example, consider the univariate RBF (Radial Basis Function) kernel; it turns out that we only need to take around the first 12 eigenfunctions (i.e.  $m = 12$ ) before adding any more has a negligible effect on the approximation accuracy (the eigenvalues decay relatively fast). However, more complex covariance functions do require more eigenvalues/eigenfunctions and performance depends on the decay of the eigenvalues [Svensson et al., 2016; Solin and Särkkä, 2014].

Solin and Särkkä [2014] also give the result for the common case where we have a stationary kernel with inputs  $\mathbf{x} \in \mathbb{R}^d$  and a domain  $\Omega = [-L_1, L_1] \times \dots \times [-L_d, L_d]$  with Dirichlet Boundary conditions (i.e. zero on the boundary). Suppose  $k(\mathbf{x}, \mathbf{x}')$  is a stationary kernel and let  $\mathbf{m} = (m_1, \dots, m_d) \in \mathbb{R}^d$  be a vector giving the chosen number of basis functions in each of the  $d$  input dimensions, then the  $\mathbf{m}$ -approximation of  $k(\mathbf{x}, \mathbf{x}')$  is [Solín and Särkkä, 2014]:

$$k_{\mathbf{m}}(\mathbf{x}, \mathbf{x}') \approx \sum_{(j_1, \dots, j_d)=\mathbf{1}}^{\mathbf{m}} S(\sqrt{\lambda_{j_1, \dots, j_d}}) \phi_{j_1, \dots, j_d}(\mathbf{x}) \phi_{j_1, \dots, j_d}(\mathbf{x}') \quad (2.46)$$

where  $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^d$ , each  $j_i$  can take any integer between 1 and  $\max(j_i) = m_i$  inclusive and the eigenvalues are:

$$\lambda_{j_1, \dots, j_d} = \sum_{k=1}^d \left( \frac{\pi j_k}{2L_k} \right)^2 \quad (2.47)$$

and the eigenfunctions are:

$$\phi_{j_1, \dots, j_d}(\mathbf{x}) = \prod_{k=1}^d \frac{1}{\sqrt{L_k}} \sin \left( \frac{\pi j_k (x_k + L_k)}{L_k} \right) \quad (2.48)$$

Firstly, note that the number of terms in the sum (2.46) is  $\prod_{i=1}^d m_i$  and secondly, note that this is a scalar-valued covariance function, so in the case we wish to turn it in a multi-output covariance function we can; for example, apply this approximation to a kernel on  $(\mathbf{x}, d)$  and  $(\mathbf{x}', d')$  for  $d, d' \in \{1, \dots, D\}$  where  $D$  is the output dimension.

Given training data  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , we can apply the above results to Gaussian

processes by noting that we can write the covariance matrix  $\mathbf{K} = K(\mathbf{X}, \mathbf{X})$  as [Solin and Särkkä, 2014]:

$$\mathbf{K} \approx \Phi \Lambda \Phi^T \quad (2.49)$$

where  $\Lambda$  is a diagonal matrix with  $\Lambda_{jj} = S(\sqrt{\lambda_j})$ ,  $\Phi_{ij} = \phi_j(\mathbf{x}_i)$  for  $\mathbf{j} = (j_1, \dots, j_d)$  and data  $\mathbf{x}_1, \dots, \mathbf{x}_N$  with  $\mathbf{x}_i \in \mathbb{R}^d$ .

Solin and Särkkä [2014] were not that clear on what should happen in the case  $d > 1$  since writing this out naively will result in complex interactions between high-dimensional objects. To deal with this issue, we introduce a simple method to map these high-dimensional arrays back to two dimensional arrays. The idea is to map each  $(j_1, \dots, j_d)$  to a unique integer between 1 and  $\prod_{i=1}^d \max(j_i)$  using the map:

$$\mathcal{M} : (j_1, \dots, j_d) \rightarrow j \in \left\{ 1, \dots, \prod_{i=1}^d \max(j_i) \right\} \quad (2.50)$$

$$\mathcal{M}(j_1, \dots, j_d) = 1 + \sum_{i=1}^d (j_i - 1) \prod_{k=i+1}^d \max(j_k) \quad (2.51)$$

where  $\max(j_k)$  is the maximum value each  $j_k$  can take and  $\prod_{k=d+1}^d \max(j_k) = 1$ . For example, in the case that we have an input dimension of 5 and each  $j_i$  takes a value between 1 and 12 (i.e. 12 eigenfunctions for each dimension) we would map:

$$\mathcal{M}(1, 1, 1, 1, 1) = 1 \quad (2.52)$$

$$\mathcal{M}(1, 1, 1, 1, 2) = 2 \quad (2.53)$$

$$\mathcal{M}(1, 1, 1, 2, 1) = 13 \quad (2.54)$$

$$\mathcal{M}(5, 1, 1, 1, 1) = 82945 \quad (2.55)$$

$$\mathcal{M}(6, 7, 4, 5, 3) = 114531 \quad (2.56)$$

$$\mathcal{M}(12, 1, 1, 1, 1) = 228097 \quad (2.57)$$

$$\mathcal{M}(\mathbf{m}) = \mathcal{M}(12, 12, 12, 12, 12) = 248832 = 12^5 \quad (2.58)$$

Going back to (2.49), we can now map each  $\mathbf{j}$  to an integer using  $\mathcal{M}(j_1, \dots, j_d)$  so we replace  $\mathbf{j}$  with  $\mathcal{M}(\mathbf{j}) = \mathcal{M}(j_1, \dots, j_d)$  and hence with  $\mathbf{j} = (j_1, \dots, j_d)$ ;  $\Lambda_{\mathbf{j}\mathbf{j}}$  becomes  $\Lambda_{\mathcal{M}(j_1, \dots, j_d)\mathcal{M}(j_1, \dots, j_d)}$  and  $\Phi_{i\mathbf{j}}$  become  $\Phi_{i\mathcal{M}(j_1, \dots, j_d)}$ . This allows us to map the original high dimension objects back to two dimensions and thus leaves  $\Lambda$  and  $\Phi$  as matrices. The fact that this scales poorly as the number of input dimensions  $d$  increases is not as a result of this mapping but inherent to the approximation and in section 3.3 we shall introduce a solution to this problem. Now we understand how to transform these high-dimensional arrays to two dimensions, we can write out the results needed for reduced rank GPs. Given the system (2.12) with  $n_y = 1$ ,  $\mathbf{Q} = \sigma_n^2$  and a set of test inputs  $\mathbf{x}_{*1}, \dots, \mathbf{x}_{*T}$  with  $\mathbf{x}_{*i} \in \mathbb{R}^d$  we have [Solin and Särkkä, 2014]:

$$\mathbb{E}[f_*] = \Phi^* (\Phi^T \Phi + \sigma_n^2 \Lambda^{-1})^{-1} \Phi^T \mathbf{y} \quad (2.59)$$

$$\mathbb{V}[f_*] = \sigma_n^2 \Phi^* (\Phi^T \Phi + \sigma_n^2 \Lambda^{-1})^{-1} \Phi^{*T} \quad (2.60)$$

where  $\Phi_{ij}^* = \phi_j(\mathbf{x}_{*i})$  and  $\mathbf{j}$  is mapped to an integer using  $\mathcal{M}$ . Something similar can be done for the marginal likelihood [Solin and Särkkä, 2014, p. 10]. As a result,

the computationally expensive matrix inversions are independent of the size of the training data and the Hilbert reduced-rank GP can train in  $\mathcal{O}(m^3)$  where  $m$  is the number of terms in the kernel approximation sum (2.46). However, this approximation scales very poorly as the input dimension  $d$  increases.

A final thing to note is that we can use the Karhunen-Loève expansion to write [Solin and Särkkä, 2014]:

$$f(\mathbf{X}) \approx \sum_{j=1}^m f_j \phi_j(\mathbf{X}) = \mathbf{F} \phi(\mathbf{X}) \quad (2.61)$$

where the  $\phi_j$  are the basis functions from (2.46) and the  $f_j$  are some constants. If  $d > 1$  we use the mapping  $j \rightarrow \mathcal{M}(j_1, \dots, j_d)$ .

### 2.1.3 Examples

In this section, we look at some examples of Gaussian process regression on synthetic data.

#### Example 1:

One of the most useful properties of covariance functions (kernels) is that summing (or multiplying) two different covariance functions returns another covariance function which has inherited the properties of both the original covariance functions and combined them in an additive (or multiplicative) way. For example, consider training data  $\mathbf{y}$ ,  $\mathbf{X}$  generated by the system:

$$y_i = f(x_i) + \epsilon_i \quad (2.62a)$$

$$f(x) = x^2 + 3 \operatorname{sgn}(x) \quad (2.62b)$$

$$\epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, 3^2) \quad (2.62c)$$

Using a kernel formed from summing a Rational Quadratic and MLP kernel, Gaussian process regression can learn something close to the true function: see Figure 2.1.

#### Example 2:

In this example, we compare some different forms of multi-output GP regression using the example we looked at in (2.39) but with four output variables:

$$\begin{pmatrix} y_{i1} \\ y_{i2} \\ y_{i3} \\ y_{i4} \end{pmatrix} = \begin{pmatrix} f_1(x_i) \\ f_2(x_i) \\ f_3(x_i) \\ f_4(x_i) \end{pmatrix} + \begin{pmatrix} \epsilon_{i1} \\ \epsilon_{i2} \\ \epsilon_{i3} \\ \epsilon_{i4} \end{pmatrix} \quad (2.63a)$$

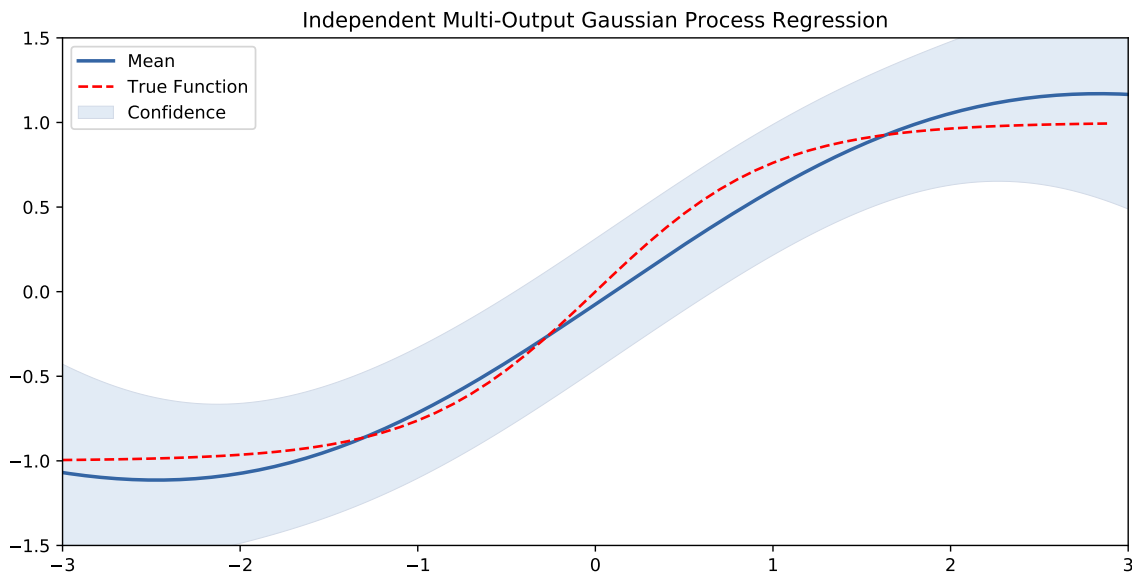
$$f_j(x) = \tanh(x) \quad (2.63b)$$

$$\epsilon_{ij} \stackrel{iid}{\sim} \mathcal{N}(0, 1^2) \quad (2.63c)$$

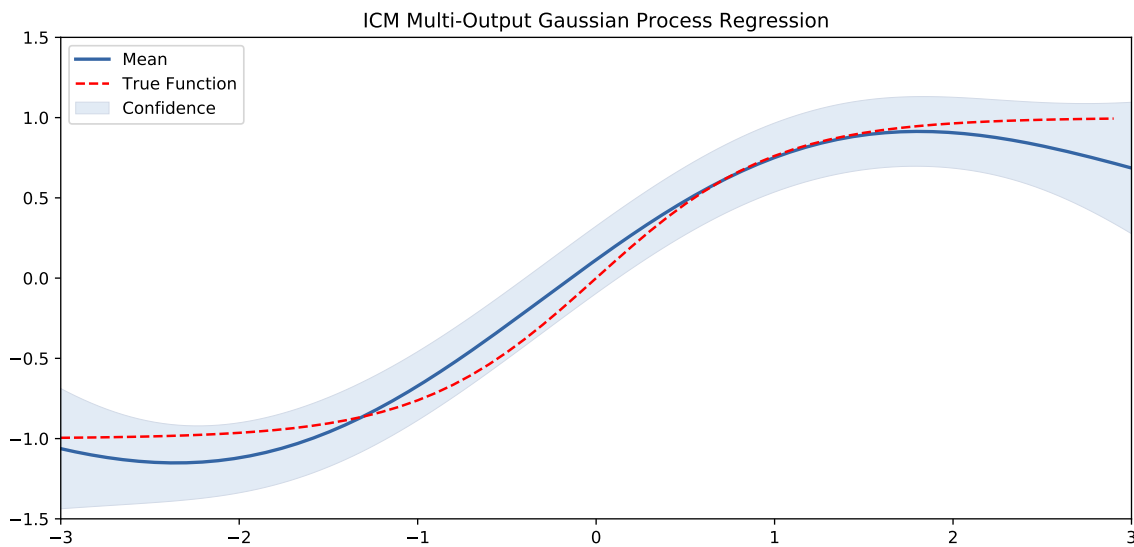
All the multi-output models were trained on 60 data points using 10 optimisation restarts and RBF covariance functions. The plots were produced using GPy [2012] and we write the training data as  $(\mathbf{X}, \mathbf{Y})$  where  $\mathbf{X}$  is the training inputs and  $\mathbf{Y} = (\mathbf{y}^1, \dots, \mathbf{y}^4)$  is the training outputs with  $\mathbf{y}^i$  being the training output for the  $i^{\text{th}}$  output dimension. The results are display in figures 2.2–2.4.



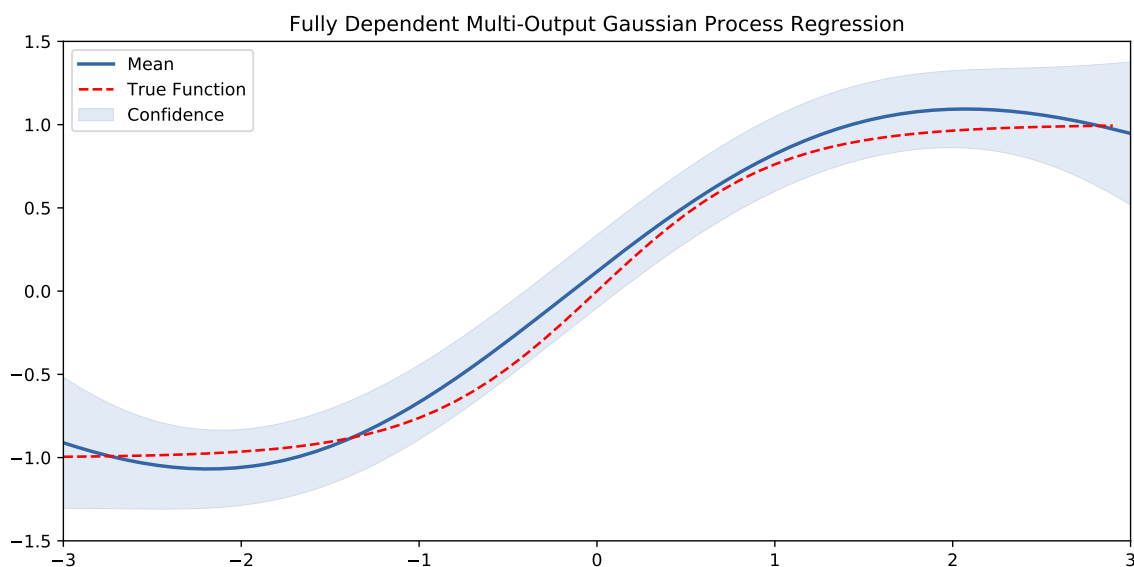
**Figure 2.1:** Univariate Gaussian process regression with a Rational Quadratic + MLP Kernel and 100 points  $(x_i, y_i)$  of training data generated by the system (2.62).



**Figure 2.2:** A plot of  $p(f_1(x_*)|x_*, \mathbf{X}, \mathbf{y}^1)$  against test inputs  $x_*$ . This is the posterior predictive distribution of the first dimension of  $f(x)$  in the case of each output dimension having a GP trained independently.



**Figure 2.3:** A plot of  $p(f_1(x_*)|x_*, \mathbf{X}, \mathbf{Y})$  against test inputs  $x_*$  in the case of a GP with a separable multi-output covariance:  $k(\mathbf{x}, \mathbf{x}')_{d,d'} = k_{\mathcal{X}}(\mathbf{x}, \mathbf{x}')k_D(d, d')$ . Notice how much more confident the model is compared to the independent GP case.



**Figure 2.4:** A plot of  $p(f_1(x_*)|x_*, \mathbf{X}, \mathbf{Y})$  against test inputs  $x_*$  in the case of a GP with a full multi-output covariance:  $k(\mathbf{x}, \mathbf{x}')_{d,d'} = \exp\left(-\frac{\|(\mathbf{x}, d) - (\mathbf{x}', d')\|^2}{2l^2}\right)$ . It is not that different to the separable multi-output covariance.

## 2.2 State-Space Models

A state-space model is a model of the form:

$$\mathbf{x}_1 \sim p(\mathbf{x}_1) \quad (2.64)$$

$$\mathbf{x}_t | \mathbf{x}_{t-1} \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}) \quad (2.65)$$

$$\mathbf{y}_t | \mathbf{x}_t \sim p(\mathbf{y}_t | \mathbf{x}_t) \quad (2.66)$$

where the  $\mathbf{x}_t \in \mathbb{R}^{n_x}$  for  $t = 1, \dots, T$  are latent (unobserved) variables and the  $\mathbf{y}_t \in \mathbb{R}^{n_y}$  for  $t = 1, \dots, T$  are observed variables. We will refer to (2.64) as the initial distribution, (2.65) as the transition distribution and (2.66) as the observation distribution. A graphical model is provided in figure 2.5.

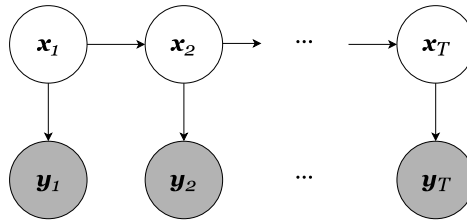


Figure 2.5: A state-space model

Two important special cases of the general state-space model are:

1. The Gaussian linear state-space model (also known as the Gaussian linear dynamical system):

$$\mathbf{x}_t = \mathbf{F}\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t \quad (2.67a)$$

$$\mathbf{y}_t = \mathbf{G}\mathbf{x}_t + \mathbf{w}_t \quad (2.67b)$$

$$\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \mathbf{P}_1) \quad (2.67c)$$

$$\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_t) \quad (2.67d)$$

$$\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{R}_t) \quad (2.67e)$$

where  $\mathbf{F}$ ,  $\mathbf{G}$  are matrices of suitable size and  $\mathbf{P}_1$ ,  $\mathbf{Q}_t$ ,  $\mathbf{R}_t$  are covariance matrices of suitable size. This model has a number of desirable properties such as a closed form recursive formula for filtering and smoothing called the Kalman filter and the Kalman smoother respectively.

2. The additive state-space model:

$$\mathbf{x}_t = f(\mathbf{x}_{t-1}) + \boldsymbol{\epsilon}_t \quad (2.68a)$$

$$\mathbf{y}_t = g(\mathbf{x}_t) + \mathbf{w}_t \quad (2.68b)$$

$$\mathbf{x}_1 \sim p(\mathbf{x}_1) \quad (2.68c)$$

$$\boldsymbol{\epsilon}_t \sim p(\boldsymbol{\epsilon}_t) \quad (2.68d)$$

$$\mathbf{w}_t \sim p(\mathbf{w}_t) \quad (2.68e)$$

where  $f$  and  $g$  are functions. As we shall see in section 2.3, a Gaussian process state-space model is an additive state-space model with Gaussian process priors on  $f$  and  $g$ .

Given an arbitrary state-space model and observations  $\mathbf{y}_1, \dots, \mathbf{y}_T$ , important problems include:

1. Finding the filtering distributions:  $p(\mathbf{x}_t | \mathbf{y}_1, \dots, \mathbf{y}_t)$  for  $t = 1, \dots, T$ .
2. Finding the smoothing distributions:  $p(\mathbf{x}_t | \mathbf{y}_1, \dots, \mathbf{y}_T)$  for  $t = 1, \dots, T$ .
3. Prediction of future  $\mathbf{y}$ :  $p(\mathbf{y}_{T+s} | \mathbf{y}_1, \dots, \mathbf{y}_T)$  for some  $s \in \mathbb{Z}_{\geq 1}$ .
4. Finding the unknown parts of the initial, transition and observation distributions.

The general name for the first three problems is inference while the general name for the fourth problem is learning. We will discuss learning briefly at the end of this section and cover it in more depth for the specific case of Gaussian process state-space models in section 2.3. For now, we will discuss inference.

### 2.2.1 Inference: The General Framework

Bayesian filtering and smoothing is a framework which encompasses all the inference (filtering, smoothing and prediction) that we will look at in this thesis. We will now state the key results in this general case and throughout this section suppress the dependence on possible parameters for clarity.

Given observations  $\mathbf{y}_{1:t}$  (where as usual  $\mathbf{y}_{1:t}$  means  $\mathbf{y}_1, \dots, \mathbf{y}_t$ ), the joint smoothing distribution is [Doucet et al., 2001]:

$$p(\mathbf{x}_{1:t} | \mathbf{y}_{1:t}) = \frac{p(\mathbf{y}_{1:t} | \mathbf{x}_{1:t}) p(\mathbf{x}_{1:t})}{\int p(\mathbf{y}_{1:t} | \mathbf{x}_{1:t}) p(\mathbf{x}_{1:t}) d\mathbf{x}_{1:t}} \quad (2.69)$$

this can be turned into a recursion version [Doucet and Johansen, 2012]:

$$p(\mathbf{x}_{1:t} | \mathbf{y}_{1:t}) = p(\mathbf{x}_{1:t-1} | \mathbf{y}_{1:t-1}) \frac{p(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{x}_{t-1})}{p(\mathbf{y}_t | \mathbf{y}_{1:t-1})} \quad (2.70)$$

where  $p(\mathbf{y}_t | \mathbf{y}_{1:t-1})$  can be found using the integral [Doucet and Johansen, 2012]:

$$p(\mathbf{y}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{y}_t | \mathbf{x}_t) d\mathbf{x}_{t-1} \quad (2.71)$$

The filtering distribution  $p(\mathbf{x}_t | \mathbf{y}_{1:t})$  can be found by a two step procedure [Doucet et al., 2001; Doucet and Johansen, 2012]:

$$p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \quad (2.72)$$

$$p(\mathbf{x}_t | \mathbf{y}_{1:t}) = \frac{p(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{y}_{1:t-1})}{p(\mathbf{y}_t | \mathbf{y}_{1:t-1})} \quad (2.73)$$

where (2.72) is called prediction and (2.73) is called updating. Finally given observations  $\mathbf{y}_{1:T}$ , the smoothing distribution  $p(\mathbf{x}_t | \mathbf{y}_{1:T})$  can be found via the Forward-Backward Recursions [Doucet and Johansen, 2012]:

$$p(\mathbf{x}_t | \mathbf{y}_{1:T}) = p(\mathbf{x}_t | \mathbf{y}_{1:t}) \int \frac{p(\mathbf{x}_{t+1} | \mathbf{x}_t)}{p(\mathbf{x}_{t+1} | \mathbf{y}_{1:t})} p(\mathbf{x}_{t+1} | \mathbf{y}_{1:T}) d\mathbf{x}_{t+1} \quad (2.74)$$



where the  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  for  $t = 1, \dots, T$  were found during filtering (forward pass) and after we have collected all the data we can perform the smoothing (backward pass) recursively via (2.74). An alternative to this is the Generalised Two-Filter Formula which is discussed in Doucet and Johansen [2012]. As noted in Doucet et al. [2001] these equations are ‘*deceptively simple*’ since they often require the solution of intractable integrals. To deal with this, some approximations are needed and we will now look at two different types of possible approximations: deterministic and stochastic.

### 2.2.2 Inference: Deterministic Methods

The two main deterministic inference frameworks are Gaussian filtering and variational inference. Gaussian filtering is in fact a form of the more general variational inference but it is very important because many methods traditionally used for inference in state-space models such as the Kalman filter, extended Kalman filter (EKF), unscented Kalman filter (UKF) and cubature Kalman filter (CKF) are all special cases of Gaussian filtering [Deisenroth, 2010].

Before discussing Gaussian filtering we will introduce the shorthand notation defined in Deisenroth [2010]:

$$\boldsymbol{\mu}_{t_1|t_2}^x = \mathbb{E}[\mathbf{x}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.75)$$

$$\boldsymbol{\Sigma}_{t_1|t_2}^x = \text{cov}[\mathbf{x}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.76)$$

$$\boldsymbol{\mu}_{t_1|t_2}^y = \mathbb{E}[\mathbf{y}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.77)$$

$$\boldsymbol{\Sigma}_{t_1|t_2}^y = \text{cov}[\mathbf{y}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.78)$$

$$\boldsymbol{\Sigma}_{t_1|t_2}^{xy} = \text{cov}[\mathbf{x}_{t_1}, \mathbf{y}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.79)$$

$$\boldsymbol{\Sigma}_{t_1|t_2}^{yx} = \text{cov}[\mathbf{y}_{t_1}, \mathbf{x}_{t_1} | \mathbf{y}_{1:t_2}] \quad (2.80)$$

The idea behind Gaussian filtering is to approximate most of the distributions which arise in the Bayesian filtering computations with Gaussian distributions by moment-matching; for example, we approximate  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  with a  $\mathcal{N}(\boldsymbol{\mu}_{t|t}^x, \boldsymbol{\Sigma}_{t|t}^x)$  distribution. Notice that in the case where  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  is already Gaussian then the above approximation will leave the distribution unchanged. The only distributions that are not (usually) approximated are the initial (2.64), transition (2.65) and observation (2.65) distributions because as long as we can find their first two moments we do not need to approximate them.

Assuming we know  $\boldsymbol{\mu}_{t-1|t-1}^x$  and  $\boldsymbol{\Sigma}_{t-1|t-1}^x$  and therefore have a Gaussian approximation of  $p(\mathbf{x}_{t-1}|\mathbf{y}_{1:t-1})$  as  $\mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{t-1|t-1}^x, \boldsymbol{\Sigma}_{t-1|t-1}^x)$ , it is possible to derive  $\boldsymbol{\mu}_{t|t}^x$  and  $\boldsymbol{\Sigma}_{t|t}^x$  and therefore a Gaussian approximation of  $p(\mathbf{x}_t|\mathbf{y}_{1:t})$  [Deisenroth, 2010]. Following this recursively will yield a sequence of Gaussian approximate filtering distributions. To derive the updates, we note that Gaussian filtering allows us to avoid direct computation of the prediction update integral in (2.72) because all we need are the first two moments (since we are going to approximate it with a Gaussian distribution)

and hence following Deisenroth [2010] we get:

$$p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) = \mathcal{N}(\boldsymbol{\mu}_{t|t-1}^x, \boldsymbol{\Sigma}_{t|t-1}^x) \quad (2.81)$$

where

$$\boldsymbol{\mu}_{t|t-1}^x = \mathbb{E}[p(\mathbf{x}_t | \mathbf{y}_{1:t-1})] \quad (2.82)$$

$$= \int \mathbf{x}_t p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) d\mathbf{x}_t \quad (2.83)$$

$$= \iint \mathbf{x}_t p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} d\mathbf{x}_t \quad (2.84)$$

$$= \iint \mathbf{x}_t p(\mathbf{x}_t | \mathbf{x}_{t-1}) d\mathbf{x}_t p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \quad (2.85)$$

$$= \int \mathbb{E}[p(\mathbf{x}_t | \mathbf{x}_{t-1})] p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \quad (2.86)$$

$$= \int \mathbb{E}[p(\mathbf{x}_t | \mathbf{x}_{t-1})] \mathcal{N}(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}^x, \boldsymbol{\Sigma}_{t-1|t-1}^x) d\mathbf{x}_{t-1} \quad (2.87)$$

where we note that if the expectation exists (i.e. is finite) swapping the integrals is allowed by Fubini's Theorem. Via a similar method,

$$\boldsymbol{\Sigma}_{t|t-1}^x = \int \mathbb{E}[p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_t | \mathbf{x}_{t-1})^T] \mathcal{N}(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}^x, \boldsymbol{\Sigma}_{t-1|t-1}^x) d\mathbf{x}_{t-1} - \boldsymbol{\mu}_{t|t-1}^x (\boldsymbol{\mu}_{t|t-1}^x)^T \quad (2.88)$$

Using the same approach it is possible to derive  $\boldsymbol{\Sigma}_{t|t-1}^{xy}$ ,  $\boldsymbol{\Sigma}_{t|t-1}^y$  and  $\boldsymbol{\Sigma}_{t|t-1}^{yx}$  [Deisenroth, 2010]. The filter updates are derived in Deisenroth [2010] and are reproduced here:

$$\boldsymbol{\mu}_{t|t}^x = \boldsymbol{\mu}_{t|t-1}^x + \boldsymbol{\Sigma}_{t|t-1}^{xy} (\boldsymbol{\Sigma}_{t|t-1}^y)^{-1} (\mathbf{y}_t - \boldsymbol{\mu}_{t|t-1}^y) \quad (2.89)$$

$$\boldsymbol{\Sigma}_{t|t}^x = \boldsymbol{\Sigma}_{t|t-1}^x + \boldsymbol{\Sigma}_{t|t-1}^{xy} (\boldsymbol{\Sigma}_{t|t-1}^y)^{-1} \boldsymbol{\Sigma}_{t|t-1}^{yx} \quad (2.90)$$

Although Gaussian filtering avoids computation of certain integrals such as (2.72), there is no guarantee that these new integrals (2.87), *etc* are tractable. As a result, the different Gaussian filtering methods make different additional approximations; for example, the extended Kalman filter performs filtering on the additive state-space model by linearisation of the non-linear functions (approximating  $f(\mathbf{x}_{t-1})$  with  $F_t \mathbf{x}_{t-1}$ ). This allows for the simple computation of the required integrals: see (2.91)–(2.96). Other Gaussian filtering methods; for example, the unscented Kalman filter and the cubature Kalman filter attempt to approximate the new integrals (2.87), *etc* directly. For more details see Deisenroth [2010]. In the case of the linear Gaussian dynamical system, there are no approximations made by Gaussian filtering: all the required equations can be computed exactly and the resulting update equations are the Kalman filter update equations; for example, (2.87) becomes:

$$\int \mathbb{E}[f(\mathbf{x}_t | \mathbf{x}_{t-1})] \mathcal{N}(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}^x, \boldsymbol{\Sigma}_{t-1|t-1}^x) d\mathbf{x}_{t-1} \quad (2.91)$$

$$= \int \mathbf{F} \mathbf{x}_{t-1} \mathcal{N}(\mathbf{x}_{t-1} | \boldsymbol{\mu}_{t-1|t-1}^x, \boldsymbol{\Sigma}_{t-1|t-1}^x) d\mathbf{x}_{t-1} \quad (2.92)$$

$$= \mathbf{F} \boldsymbol{\mu}_{t-1|t-1}^x \quad (2.93)$$

and the updates become (Deisenroth [2010]):

$$\boldsymbol{\mu}_{t|t}^x = \mathbf{F} \boldsymbol{\mu}_{t-1|t-1}^x + \boldsymbol{\Sigma}_{t|t-1}^x \mathbf{G}^T (\mathbf{G} \boldsymbol{\Sigma}_{t|t-1}^x \mathbf{G}^T + \mathbf{R}_t)^{-1} (\mathbf{y}_t - \mathbf{G} \boldsymbol{\mu}_{t-1|t-1}^x) \quad (2.94)$$

$$\boldsymbol{\Sigma}_{t|t}^x = \boldsymbol{\Sigma}_{t|t-1}^x + \boldsymbol{\Sigma}_{t|t-1}^x \mathbf{G}^T (\mathbf{G} \boldsymbol{\Sigma}_{t|t-1}^x \mathbf{G}^T + \mathbf{R}_t)^{-1} \mathbf{G} \boldsymbol{\Sigma}_{t|t-1}^x \quad (2.95)$$

where

$$\boldsymbol{\Sigma}_{t|t-1}^x = \mathbf{F} \boldsymbol{\Sigma}_{t-1|t-1}^x \mathbf{F}^T + \mathbf{Q}_t \quad (2.96)$$

Using the same ideas as Gaussian filtering, it is possible to derive a set of recursive equations to generate the Gaussian approximated smoothing distributions, these approximate distributions become the exact smoothing distributions in the linear Gaussian dynamical case and are known as the Kalman smoother: see section 4.2.2 in Deisenroth [2010].

A more general Bayesian filtering/smoothing method is variational inference which tries to approximate the distributions inside the Bayesian integrals (2.71)–(2.72) with distributions that make the computation of the integrals tractable.

**Definition 2.6.** [Bishop, 2009] The Kullback-Leibler divergence is defined as:

$$\mathcal{KL}(q||p) = - \int q(x) \log \left( \frac{p(x)}{q(x)} \right) dx \quad (2.97)$$

where  $p$  and  $q$  are functions. Its properties include:

1.  $\mathcal{KL}(q||p) \geq 0$ .
2.  $\mathcal{KL}(q||p) = 0 \iff q = p$ .

Consider (2.72) which is reproduced here for convenience:

$$p(\mathbf{x}_t | \mathbf{y}_{1:t-1}) = \int p(\mathbf{x}_t | \mathbf{x}_{t-1}) p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1}) d\mathbf{x}_{t-1} \quad (2.98)$$

A variational method would try to make this integral tractable by approximating  $p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})$  or  $p(\mathbf{x}_t | \mathbf{x}_{t-1})$ . One way to do this would be to consider a distribution  $q(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})$  which approximates  $p(\mathbf{x}_{t-1} | \mathbf{y}_{1:t-1})$ . To select the appropriate  $q$  one usually tries to find the  $q$  that minimises the Kullback-Leibler divergence  $\mathcal{KL}(q||p)$ . Now of course if no constraints are placed on  $q$  then the optimal  $q$  is  $p$  which is not very helpful. The key to variational inference is to come up with constraints on  $q$  that are not too tight and allow the required integrals to become tractable. These methods tend to be very problem specific and so we will not discuss them here: they will be introduced as needed.

### 2.2.3 Inference: Stochastic Methods

The main stochastic methods for inference in state-space models are sequential Monte Carlo (SMC) methods also known as particle filters and particle smoothers. In this section, we will look at a basic particle filter known as the bootstrap filter and then discuss the issues which can arise in SMC methods; finally, we will briefly look at smoothing and an important state-of-the-art particle smoother called Particle Gibbs with Ancestor Sampling (PGAS) [Lindsten et al., 2014].

In particle filters, the aim is (often) to devise a method to sample from the distributions  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  for  $t = 1, \dots, T$  and if we can sample from such distributions, it is trivial to sample from the filter distributions by just taking the last element of the sampled vectors. The assumption [Doucet et al., 2001] is that sampling from  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  directly is hard (or sampling takes a long time) and for state-space models this is almost always the case. As a result, sequential Monte Carlo (SMC) methods are based on importance sampling, which, instead of sampling directly from  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$ , samples from a *importance* or *proposal* distribution  $q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  that is similar to  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  but easier to sample from. These samples are then weighted; with samples most like samples from  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  given more weight. Formally, an importance distribution  $q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  can be any distribution that such that  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t}) > 0 \implies q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t}) > 0$  [Doucet et al., 2001]. In the case of importance sampling in state-space models we have [Doucet and Johansen, 2012]:

$$p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t}) = \frac{w_t(\mathbf{x}_{1:t})q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})}{\mathbf{Z}_t} \quad (2.99)$$

$$\mathbf{Z}_t = \int w_t(\mathbf{x}_{1:t})q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})d\mathbf{x}_{1:t} \quad (2.100)$$

$$w_t(\mathbf{x}_{1:t}) = \frac{p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t})}{q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})} \quad (2.101)$$

where  $q(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$  is the importance distribution and the  $w_t(\mathbf{x}_{1:t})$  are known as the unnormalised weights. This formulation relies on knowing  $p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t})$ , which, unlike  $p(\mathbf{x}_{1:t}|\mathbf{y}_{1:t})$ , is usually easy to calculate for state-space models. To calculate some quantity of interest say  $\mathbb{E}_p[h(\mathbf{x}_{1:t})]$  we can get independent samples  $\mathbf{x}_{1:t}^i$  from  $q$  and approximate it as [Doucet and Johansen, 2012]:

$$\mathbb{E}_p[h(\mathbf{x}_{1:t})] \approx \sum_i W_t^i h(\mathbf{x}_{1:t}^i) \quad (2.102)$$

where

$$W_t^i = \frac{w_t(\mathbf{x}_{1:t}^i)}{\sum_i w_t(\mathbf{x}_{1:t}^i)} \quad (2.103)$$

Now the main issue here is that it is very inefficient since as noted in Doucet et al. [2001] we are not really using the sequential nature of the problem and forgetting everything we learnt at time  $t - 1$  for time  $t$ . In Doucet and Johansen [2012] it is shown that one can rewrite the weights allowing for sequential updates:

$$w_t(\mathbf{x}_{1:t}) = w_{t-1}(\mathbf{x}_{1:t-1})\alpha_t(\mathbf{x}_{1:t}) \quad (2.104)$$

$$\alpha_t(\mathbf{x}_{1:t}) = \frac{p(\mathbf{x}_{1:t}, \mathbf{y}_{1:t})}{p(\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t-1})q(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t})} \quad (2.105)$$

A *particle* refers to a vector  $\mathbf{x}_{1:t}$  and so having  $N$  particles means we have  $N$  vectors  $\mathbf{x}_{1:t}^i$  for  $i = 1, \dots, N$ . Associated to each particle is a weight  $w_t(\mathbf{x}_{1:t}^i)$  which determines its importance in describing the true distribution. Before stating the bootstrap filter algorithm, there two more things we need to deal with: resampling and choosing a suitable importance/proposal distribution. One of the problems that can arise in SMC is *weight degeneracy* (also known as *path degeneracy*) [Doucet and Johansen, 2012]. This occurs when we have a continually increasing particle weight variance and eventually this leads to one particle having all the weight which is problematic since all the estimates will end up being based on only one particle. This *weight degeneracy* is measured by the effective sample size [Doucet and Johansen, 2012]:

$$ESS = \frac{1}{\sum_{i=1}^N (W_t^i)^2} \quad (2.106)$$

with smaller  $ESS$  meaning more weight degeneracy. This degeneracy can be reduced using resampling and Doucet and Johansen [2012, p. 13] gives a detailed overview of all the most common resampling techniques. Unfortunately, resampling causes another type of path degeneracy known as *sample degeneracy* or *sample impoverishment* whereby all particles eventually become the same [Doucet and Johansen, 2012]. One method to deal with this is to only resample when necessary (i.e. only resample if  $ESS$  is smaller than some threshold say  $\frac{N}{2}$ ). This is known as adaptive resampling. The last thing left to do is to find a suitable importance/proposal distribution, for SMC in state-space models all we need to find is a proposal distribution of the form [Doucet and Johansen, 2012]:

$$q(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t}) \quad (2.107)$$

It turns out [Doucet and Johansen, 2012, p. 20] that the optimal (in the sense of minimising weight variance which is important due to path degeneracy) proposal is:

$$q(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t}) = \frac{p(\mathbf{y}_t|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{x}_{t-1})}{p(\mathbf{y}_t|\mathbf{x}_{t-1})} \quad (2.108)$$

Although this is the optimal proposal distribution, choosing a different proposal distribution is often done and the bootstrap filter uses the bootstrap proposal distribution:  $q(\mathbf{x}_t|\mathbf{x}_{1:t-1}, \mathbf{y}_{1:t}) = p(\mathbf{x}_t|\mathbf{x}_{t-1})$ . This is non-optimal but it is still good because the optimal proposal requires computation of  $p(\mathbf{y}_t|\mathbf{x}_{t-1})$  and this might not be tractable. It is possible to use MCMC methods to approximate the integral in  $p(\mathbf{y}_t|\mathbf{x}_{t-1})$  but this will tend to slow everything down. Overall, the bootstrap proposal is a good and easy choice to make; however, in certain cases we will do better with different proposal distribution.

---

**Algorithm 1:** The Bootstrap Filter with Adaptive Resampling [Doucet et al., 2001; Doucet and Johansen, 2012]

---

**Input** : Number of Particles  $N$ , Observations  $\mathbf{y}_{1:T}$ , Function  $M(\mathbf{x}_{1:t})$   
**Output:**  $N$  samples from  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T})$ , Weights  $W_T^i$  for  $i = 1, \dots, N$  and  $\mathbb{E}[M(\mathbf{x}_t)]$  for  $t = 1, \dots, T$

- 1 set weights  $w_1^i = 1/N$  for  $i = 1, \dots, N$ ;
- 2 **for**  $t = 1, \dots, T$  **do**
- 3     **if**  $t = 1$  **then**
- 4         sample  $\mathbf{x}_1^i \sim p(\mathbf{x}_1)$  for  $i = 1, \dots, N$ ;
- 5     **else**
- 6         sample  $\mathbf{x}_t^i \sim p(\mathbf{x}_t|\mathbf{x}_{t-1}^i)$  for  $i = 1, \dots, N$ ;
- 7         set  $\mathbf{x}_{1:t}^i = (\mathbf{x}_{1:t-1}^i, \mathbf{x}_t^i)$  for  $i = 1, \dots, N$ ;
- 8     **end**
- 9     compute importance weights:  $w_t^i = w_{t-1}^i p(\mathbf{y}_t|\mathbf{x}_t^i)$ ;
- 10     normalise importance weights:  $W_t^i = \frac{w_t^i}{\sum_{i=1}^N w_t^i}$ ;
- 11     compute  $\mathbb{E}[M(\mathbf{x}_{1:t})] = \sum_{i=1}^N W_t^i M(\mathbf{x}_{1:t}^i)$ ;
- 12     compute effective sample size:  $ESS = \frac{1}{\sum_{i=1}^N (W_t^i)^2}$ ;
- 13     **if**  $ESS < \frac{N}{2}$  **then**
- 14         resample particles and set  $w_t^i = 1/N$  using techniques in Doucet and Johansen [2012, p. 13];
- 15     **end**
- 16 **end**
- 17 return  $\mathbf{x}_{1:T}^i$ ,  $W_T^i$  for  $i = 1, \dots, N$  and  $\mathbb{E}[M(\mathbf{x}_{1:t})]$  for  $t = 1, \dots, T$ ;

---

Due to sample impoverishment, most likely the  $N$  samples of  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T})$  from the bootstrap filter will all be exactly the same (if  $T$  is large) but by setting  $M(x) = x$ , we can get a good value for the mean of the sequence of filtering distributions over time and something similar can be done the other moments. See section 2.2.5 for some examples.

SMC methods for smoothing follow the same ideas as filtering; however, they often have a higher time complexity. Although, in theory, we could get samples from the smoothing distribution  $p(\mathbf{x}_t|\mathbf{y}_{1:T})$  by just looking at the  $t$ th component of the vector  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T})$  for each of the  $N$  samples; in practice, due to sample impoverishment this will not give a good result [Doucet and Johansen, 2012]. Methods for smoothing include the *Fixed-Lag Approximation* and *Forward Filtering-Backward Smoothing* both of which are discussed in Doucet and Johansen [2012, pp. 34-36]. One of the most recent particle smoothing methods is Particle Gibbs with Ancestor Sampling (PGAS) [Lindsten et al., 2014], this requires an additional input compared with the bootstrap filter: a reference trajectory which the algorithm uses to guide itself to the invariant distribution of the latent states. This reference trajectory could be outputted from another method; for example, the mean trajectory in the bootstrap filter (i.e. the sequence  $M(\mathbf{x}_{1:t})$  when  $M(\mathbf{x}_{1:t}) = \mathbf{x}_t$ ) or it could be the output trajectory from the previous run of the PGAS algorithm. Unlike the bootstrap filter,

PGAS only produces one sample (trajectory) per run and although degeneracy still occurs in each run, only producing one sample per run mitigates this degeneracy and also means we do not need to use as many particles as in the bootstrap filter [Lindsten et al., 2014]. The PGAS algorithm is provided in Algorithm 2.

Although we have looked at SMC from the point of view of state-space models, it can be used in a wide variety of sequential and non-sequential problems [Doucet and Johansen, 2012, p. 35]; for example, sampling via Metropolis-Hastings can use SMC to improve mixing and this leads to particle Metropolis-Hastings and other particle MCMC methods [Dahlin and Schön, 2016]. We will compare a few SMC methods with some Gaussian filtering methods in section 2.2.5.

---

**Algorithm 2:** Particle Gibbs with Ancestor Sampling [Lindsten et al., 2014; Svensson et al., 2016]

---

**Input** : Number of Particles  $N$ , Observations  $\mathbf{y}_{1:T}$ , Reference Trajectory ( $\hat{\mathbf{x}}_{1:T}$ )  
**Output**: One sample from  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T})$

- 1 sample  $\mathbf{x}_1^i$  from the initial distribution  $p(\mathbf{x}_1)$  for  $i = 1, \dots, N - 1$ ;
- 2 set  $\mathbf{x}_1^N = \hat{\mathbf{x}}_1$ ;
- 3 set weights  $w^i = p(\mathbf{y}_1|\mathbf{x}_1^i)$  for  $i = 1, \dots, N$ ;
- 4 **for**  $t = 2, \dots, T$  **do**
  - 5 /\* resampling and ancestor sampling \*/  
sample with replacement  $N - 1$  particles  $a_1, \dots, a_{N-1}$  from the  $N$  particles  $\mathbf{x}_{1:t-1}^i$  with the probability of selecting particle  $\mathbf{x}_{1:t-1}^i$  proportional to its weight;
  - 6 compute  $p^i = w^i p(\hat{\mathbf{x}}_t|\mathbf{x}_{t-1}^i)$  for  $i = 1, \dots, N$  where  $p(\mathbf{x}_t|\mathbf{x}_{t-1})$  is the transition distribution;
  - 7 normalise:  $\hat{p}^i = \frac{p^i}{\sum_i p^i}$  for  $i = 1, \dots, N$ ;
  - 8 sample a single particle  $a_N$  from the  $N$  particles  $\mathbf{x}_{1:t-1}^i$  with the probability of selecting particle  $\mathbf{x}_{1:t-1}^i$  equal to  $\hat{p}^i$ ;
  - 9 /\* particle propagation \*/  
sample  $\mathbf{x}_t^i$  from the transition distribution  $p(\mathbf{x}_t|\mathbf{x}_{t-1} = a_i)$  for  $i = 1, \dots, N - 1$ ;
  - 10 set  $\mathbf{x}_t^N = \hat{\mathbf{x}}_t$ ;
  - 11 set  $\mathbf{x}_{1:t-1}^i = a_i$  for  $i = 1, \dots, N$ ;
  - 12 /\* weighting \*/  
set weights  $w^i = p(\mathbf{y}_t|\mathbf{x}_t^i)$  for  $i = 1, \dots, N$ ;
- 13 **end**
- 14 **return** one sample  $\mathbf{x}_{1:T}$  by sampling from the particles with the probability of selecting particle  $\mathbf{x}_{1:T}^i$  proportional to its weight  $w^i$ ;

---

### 2.2.4 Learning

There are two possible kinds of learning that can be done in state-space models. Firstly, there is parametric learning, this is where we assume that the functional forms of the initial, transition and observation distributions in (2.64)–(2.66) are known but there are some unknown parameters. In this case, the task in *learning* is to find these unknown parameters; for example, in the linear state-space model one way to find the optimal parameters is via a combination of Kalman filtering/smoothing and expectation-maximisation (see Hamilton [1994] for more details). Furthermore, parametric learning in non-linear non-Gaussian models is possible via PGAS. We can construct a Gibbs sampler and alternate between sampling from the distributions  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T}, \boldsymbol{\theta})$  and  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  where  $\boldsymbol{\theta}$  represents the unknown parameters. After an initial burn-in period this will generate (approximate) samples from the posterior distribution of the parameters: see Lindsten et al. [2014, p. 2161] for more details. Although these methods have good theoretical properties, there are still issues and one of the key problems is identifiability.

**Definition 2.7.** [Lehmann and Casella, 1998] Suppose the random variable  $X$  has a distribution  $p_{\boldsymbol{\theta}}$  where  $\boldsymbol{\theta}$  is a vector of parameters, then we say  $\boldsymbol{\theta}$  is identifiable if  $p_{\boldsymbol{\theta}_1} = p_{\boldsymbol{\theta}_2} \implies \boldsymbol{\theta}_1 = \boldsymbol{\theta}_2$ .

When optimising the parameters, having a lack of identifiability results in the model having many equivalent optima. This might not be a problem by itself, since by the very nature of identifiability they are all equally valid given the observations, and if we just want to perform forecasting this unidentifiability might not be a problem [Frigola, 2015, p. 19]. However, if we wish to be able to interpret the parameters, then identifiability is an issue since not all parameter values may even make sense. In this case, placing priors on the parameters or starting the optimisations in different places may help to alleviate this issue. For a general state-space model: (2.64)–(2.66) we cannot expect to have identifiability.

The second kind of learning in state-space models is non-parametric learning and this is where we do not assume any particular functional form for the initial, transition and observation distributions. Non-parametric learning is a hard problem and generally requires us to assume something about the final form; for example, an additive state-space model with a Gaussian likelihood. We will look at non-parametric learning in depth for the particular case of Gaussian process state-space models in section 2.3.



### 2.2.5 Examples

In this section, we will compare Gaussian filtering/smoothing with sequential Monte Carlo methods via a number of examples.

#### Example 1:

The first set of examples will be based on the model:

$$\mathbf{x}_t = 0.9\mathbf{x}_{t-1} + \boldsymbol{\epsilon}_t \quad (2.109)$$

$$\mathbf{y}_t = 3\mathbf{x}_t + \mathbf{w}_t \quad (2.110)$$

$$\mathbf{x}_1 \sim \mathcal{N}(0, 0.1^2) \quad (2.111)$$

$$\boldsymbol{\epsilon}_t \stackrel{iid}{\sim} \mathcal{N}(0, 0.3^2) \quad (2.112)$$

$$\mathbf{w}_t \stackrel{iid}{\sim} \mathcal{N}(0, 1^2) \quad (2.113)$$

The bootstrap filter will use 500 particles and PGAS will use 5 particles starting from a reference trajectory of all zeros and a burn-in of 30 samples. We generated a single set of 300 points of simulated observations/hidden states from the model above and each of the methods were given the same set of simulated observations. The aim was to find the hidden states and we scored each of the methods by calculating the MSE between the mean of filtering/smoothing distribution and the true hidden states. The MSEs for the first example are displayed in following table:

Method:	MSE:
Bootstrap Filter	0.066
Bootstrap Smoother	0.055
Kalman Filter	0.065
Kalman Smoother	0.053
PGAS	<b>0.052</b>

(2.114)

On this simple linear system, the above SMC methods perform similarly to their Kalman equivalents. This is as expected because the Kalman methods are meant to be the optimal (best unbiased estimator) filters/smoother on this sort of system [Hamilton, 1994]. The results are plotted in figures 2.6–2.8.

#### Example 2:

The second set of examples will compare the unscented Kalman filter (UKF), bootstrap filter, unscented Kalman smoother (UKS), bootstrap smoother and PGAS. This set of examples will be based on the model:

$$\mathbf{x}_t = \sin(\mathbf{x}_{t-1}) + \boldsymbol{\epsilon}_t \quad (2.115)$$

$$\mathbf{y}_t = \exp(\mathbf{x}_t) + \mathbf{w}_t \quad (2.116)$$

$$\mathbf{x}_1 \sim \mathcal{N}(0, 0.1^2) \quad (2.117)$$

$$\boldsymbol{\epsilon}_t \sim \mathcal{N}(0, 0.3^2) \quad (2.118)$$

$$\mathbf{w}_t \sim \mathcal{N}(0, 1^2) \quad (2.119)$$

with everything else the same as in the first set of examples. The results are plotted in figures 2.9–2.11 and the MSEs are displayed in the following table:

Method:	MSE:
Bootstrap Filter	0.17
Bootstrap Smoother	<b>0.14</b>
UKF	0.21
UKS	0.20
PGAS	<b>0.14</b>

(2.120)

The key takeaway from both examples is that SMC methods are a competitive inference method and can provide excellent estimates of the required distributions compared to the traditional methods such as the UKF/UKS.

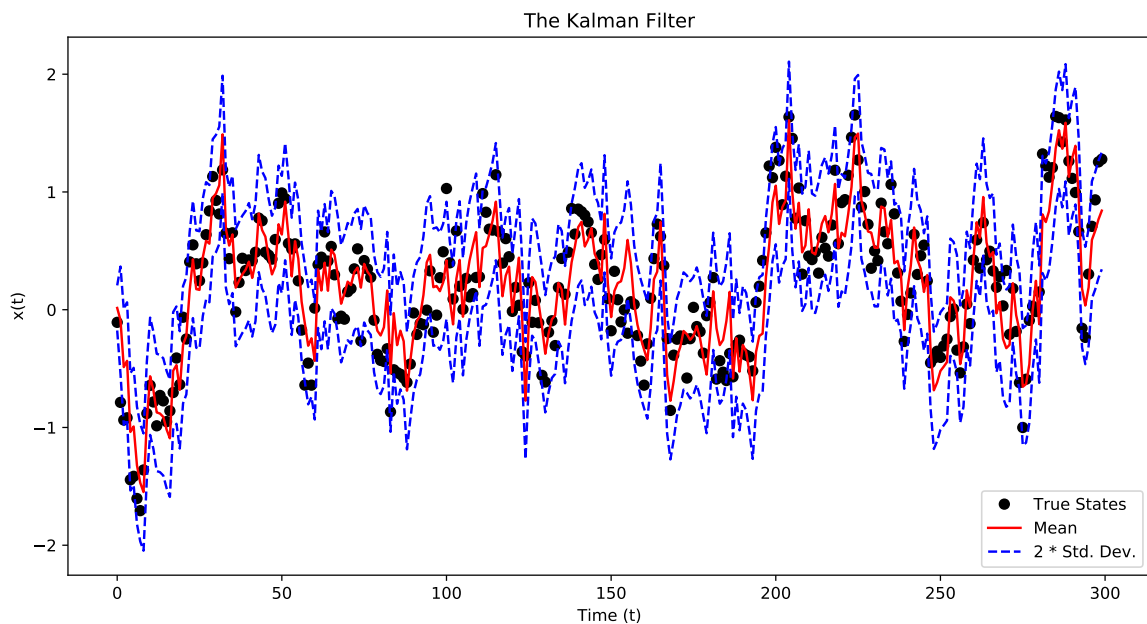


Figure 2.6: The Kalman Filter (Example 1)

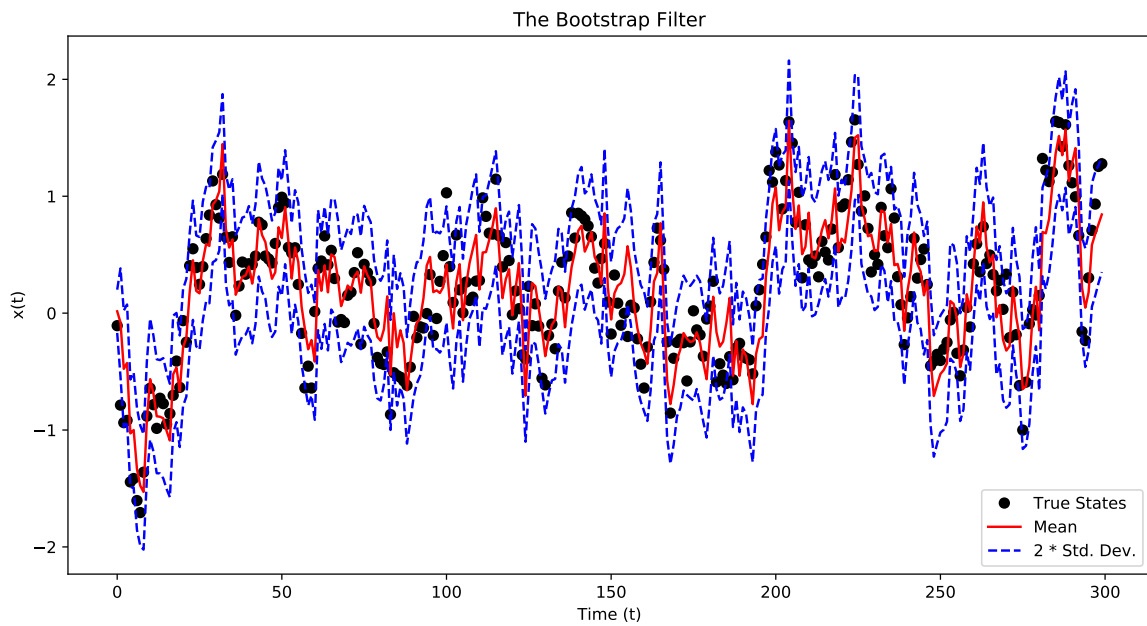


Figure 2.7: The Bootstrap Filter (Example 1)

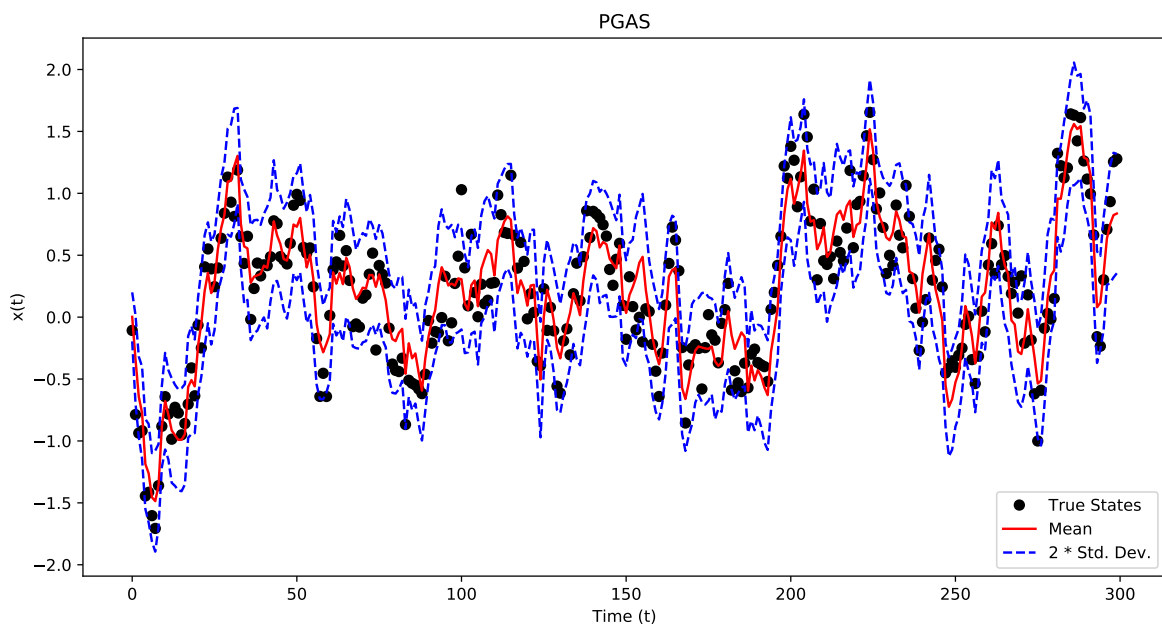


Figure 2.8: Particle Gibbs with Ancestor Sampling (PGAS) (Example 1)

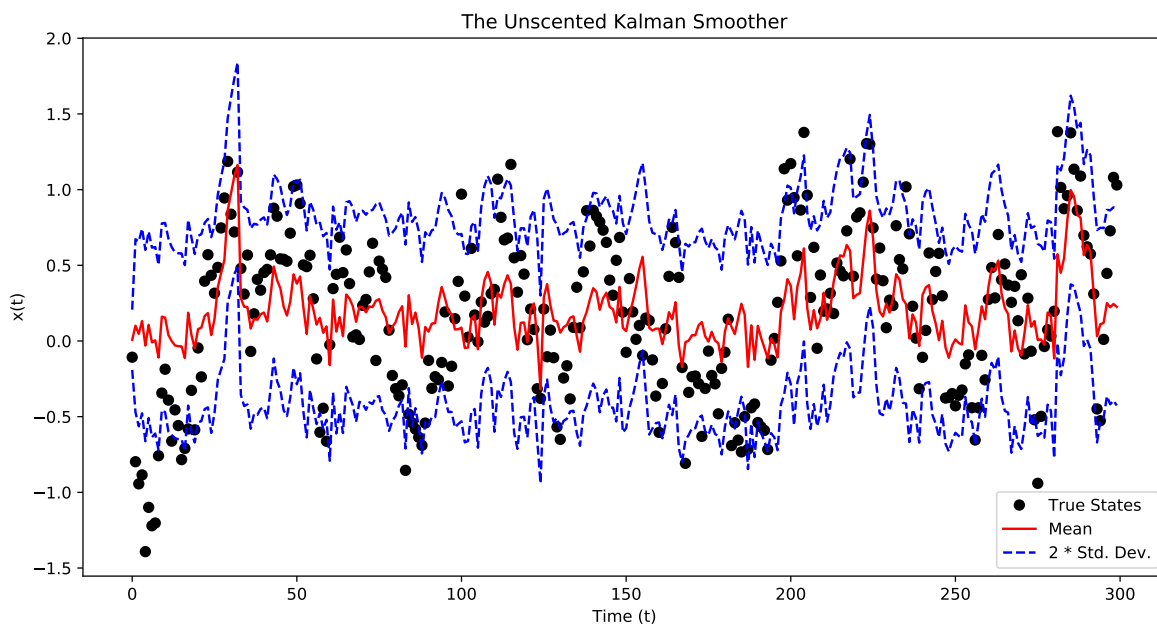


Figure 2.9: The Unscented Kalman Filter (Example 2)

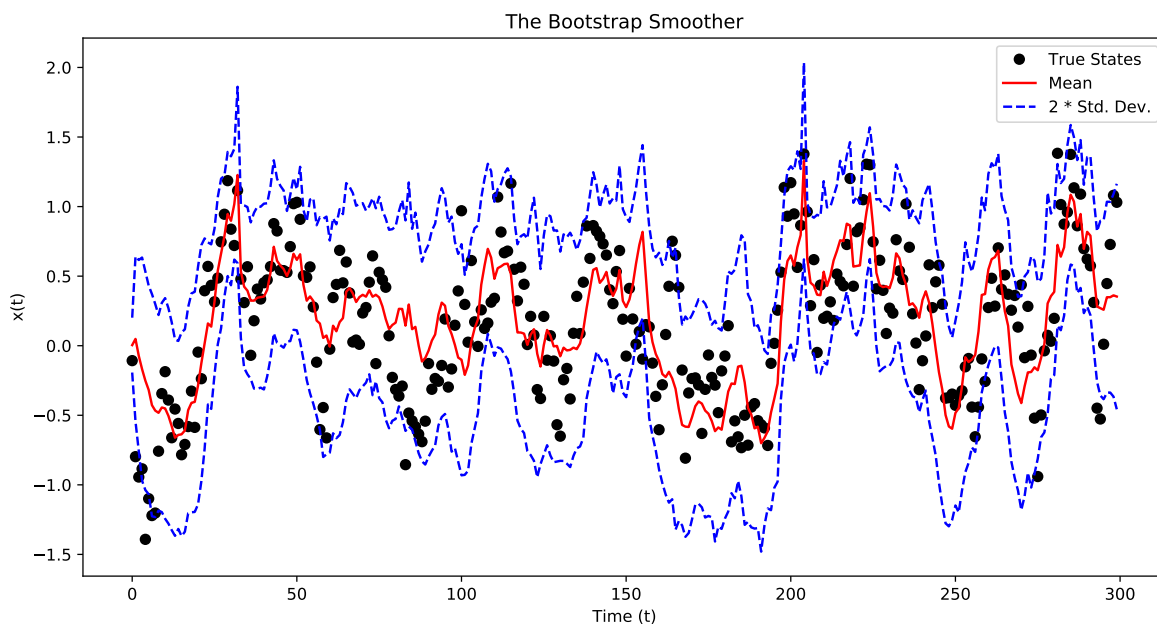


Figure 2.10: The Bootstrap Smoother (Example 2)

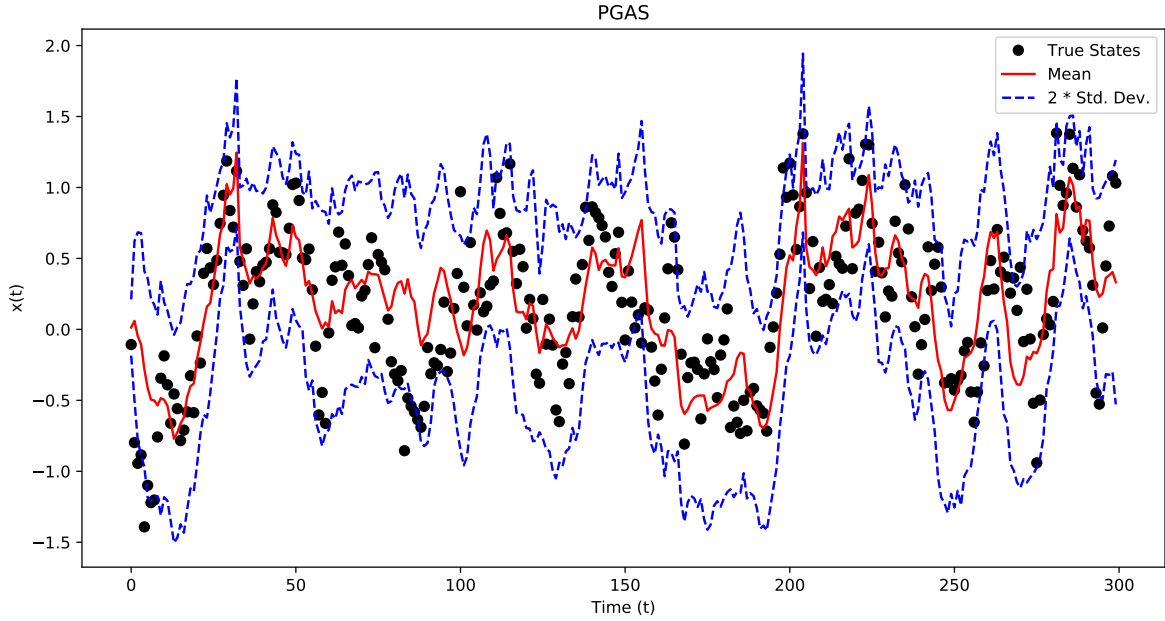


Figure 2.11: Particle Gibbs with Ancestor Sampling (PGAS) (Example 2)

## 2.3 Gaussian Process State-Space Models

A Gaussian process state-space model (GPSSM) is a model of the form:

$$\mathbf{x}_t = f(\mathbf{x}_{t-1}) + \boldsymbol{\epsilon}_t \quad (2.121a)$$

$$\mathbf{y}_t = g(\mathbf{x}_t) + \mathbf{w}_t \quad (2.121b)$$

$$\mathbf{x}_1 \sim p(\mathbf{x}_1) \quad (2.121c)$$

$$f \sim \mathcal{GP}(m_f, k_f) \quad (2.121d)$$

$$g \sim \mathcal{GP}(m_g, k_g) \quad (2.121e)$$

where  $\boldsymbol{\epsilon}_t$  and  $\mathbf{w}_t$  are random variables, the  $\mathbf{x}_t \in \mathbb{R}^{n_x}$  for  $t = 1, \dots, T$  are latent (unobserved) variables and the  $\mathbf{y}_t \in \mathbb{R}^{n_y}$  for  $t = 1, \dots, T$  are observed variables. In the case where we have:

$$\boldsymbol{\epsilon}_t \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (2.122)$$

$$\mathbf{w}_t \stackrel{iid}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{R}) \quad (2.123)$$

then this is called a GPSSM with (independent) Gaussian latent and observation noise. When we refer to a GPSSM we will be referring to a GPSSM with Gaussian latent and observation noise unless stated otherwise. Sometimes the functional form of  $f$  or  $g$  is known and so we only have one GP; in the case where the functional form of  $f$  is known and  $g$  is a GP we will call this a GPSSM with a known transition function and in the case where the functional form of  $g$  is known and  $f$  is a GP we will call this a GPSSM with a known observation function. As with all state-space

models there are two key problems we wish to solve: learning and inference. We will now look briefly at inference and then discuss the current status of research for learning GPSSMs.

### 2.3.1 Inference in GPSSMs

Once we have knowledge of the posterior predictive distributions for the transition function  $f$  in (2.121a) and observation function  $g$  in (2.121b), many of the inference techniques we discussed previously (unscented Kalman filter, bootstrap filter, etc) can be used directly for GPSSMs. However, this assumes that we can come up with some method for learning the covariance hyperparameters in  $f$  and  $g$  as well as finding their respective posterior predictive distributions. This is non-trivial because, in general, we will not have knowledge of the true hidden states  $\mathbf{x}_t$  even during training. Furthermore, even if somehow we could learn the functions  $f$  and  $g$ , as they are GPs there will always be some posterior uncertainty; but the inference techniques we have seen so far do not take into account this additional uncertainty and so will tend to give overconfident predictions [Deisenroth et al., 2009; Deisenroth, 2010].

As a result, several inference techniques have been devised which are tailored to GPSSMs. Gaussian filtering has been applied to GPSSMs [Deisenroth et al., 2009; Deisenroth, 2010] in a form known as the Gaussian process Assumed Density Filter (GP-ADF). Here, the Gaussian filtering framework is unchanged in the sense that we only need to compute to moments  $\Sigma_{t|t-1}^{xy}$ ,  $\Sigma_{t|t-1}^y$ , etc. However, in the GPSSM case the required moment calculations need to take into account the uncertainty in  $f$  and  $g$ . Moreover, the GPs  $f, g$  have uncertain inputs: when trying to calculate the next hidden state, the previous hidden state has a distribution of possible values and this uncertainty increases the GP posterior uncertainty. All this additional complexity leads to the required integrals (e.g. for computing  $\Sigma_{t|t-1}^{xy}$ ) only being tractable for a very small set of covariance functions such as the square exponential or rational quadratic covariances [Deisenroth et al., 2009]. Furthermore, during training, knowledge of the true hidden states is required [Deisenroth et al., 2009] which may not be possible depending on the problem. That said, in cases where these integrals are tractable the GP-ADF performs much better [Deisenroth et al., 2009, sec. 6.1] than the standard unscented Kalman filter (UKF) and it also performs better than a form of the UKF which takes GP uncertainty into account (GP-UFK) [Ko and Fox, 2008].

This has been extended to smoothing in Deisenroth [2010]; Deisenroth et al. [2012] and is known as GP-RTSS. This performs better than the standard unscented Kalman smoother (UKS) [Deisenroth et al., 2012, p. 5] but for the same reason as the GP-ADF, we have a very limited choice of covariance functions. For both the GP-ADF and GP-RTSS, the results are only derived for the case of the different output dimensions using independent GPs with scalar-valued covariances rather than using a single GP with a multi-output covariance. That said, it is probably possible to derive

something similar for the multi-output case.

Other inference techniques used successfully in GPSSMs are expectation propagation [Deisenroth and Mohamed, 2016] (although it was only derived and tested using an SE kernel) and particle filters [Frigola et al., 2013]. Although, the SE kernel is expressive, for some tasks other kernels are better suited (say for forecasting a periodic function) and as a result we are really looking for methods that work well across a wide range of covariance functions.

### 2.3.2 Learning in GPSSMs

Suppose we have a sequence of observations  $\mathbf{Y}_{1:T} = [\mathbf{y}_1, \dots, \mathbf{y}_T]$  and we wish to learn a GPSSM; this means that we wish to find the covariance hyperparameters for the functions  $f$  and  $g$  (or in the Bayesian case a posterior distribution for these hyperparameters). Wang et al. [2008] tried to combine learning the hyperparameters for the latent/observation GPs with learning the true hidden states, effectively trying to perform smoothing and hyperparameter learning at the same time. To do this, Wang et al. [2008] use a model based on the Gaussian process latent variable model (GPLVM) [Lawrence, 2004] and find the distributions of  $p(\mathbf{Y}|\mathbf{X})$  and  $p(\mathbf{X})$  by marginalising out  $g$  and  $f$  respectively. Although Wang et al. [2008] only states these results, a similar derivation can be found in Frigola [2015, pp. 29-30] and we can use the relationship between the matrix and multivariate normal distributions to get the results in Wang et al. [2008]. Wang et al. [2008] seem to get some good results but the issues are that training can take a long time, the model is very prone to overfitting since we have a lot of freedom when minimising over all the hidden states (especially if the hidden states have a higher dimension than the observations) and the model only provides point-estimates rather than distributions for the parameters. However, it is not limited in its choice of covariance function unlike many of the filtering methods we discussed before.

Turner et al. [2010] use the GP-Gaussian filtering methods from Deisenroth et al. [2009] to learn a GPSSM via expectation-maximisation. In the original GP filtering/smoothing methods knowledge of the true hidden states are required but Turner et al. [2010] side-step this by introducing pseudo training sets (for both the observation and transition equations) that are learnt along side the covariance hyperparameters. The idea is that in the  $E$ -step we can filter using a GP-ADF which has been constructed using the pseudo-training set rather than the actual true hidden states (which we do not know) and in the  $M$ -step as well as optimising the covariance hyperparameters we also optimise the pseudo-training set. Then, in the standard EM way we iterate between the two steps waiting for convergence; however, due to the approximations required in the  $M$ -step, we do not have any guarantees that it will converge. A limiting factor for learning using this method is the restriction on possible choices of covariance functions inherited from Gaussian filtering with Gaussian processes as well as only providing point-estimates of parameters; although, unlike Wang et al. [2008] we do get a distribution over the hidden states (using GP-ADF or GP-RTSS).

A popular set of methods for learning GPSSMs are variational methods. In all these methods, the aim is to derive an approximation to the marginal likelihood  $p(\mathbf{Y}_{1:T}|\boldsymbol{\theta})$  (i.e. marginalising out the hidden states) where  $\boldsymbol{\theta}$  are the covariance hyperparameters and then optimising it; finding (approximately) the covariance hyperparameters with a trade off between data fit and model complexity. This is done by defining the evidence lower bound (ELBO) as:

$$\mathcal{L}(q) = \int q(\mathbf{Z}|\boldsymbol{\theta}) \log \frac{p(\mathbf{Y}_{1:T}, \mathbf{Z}|\boldsymbol{\theta})}{q(\mathbf{Z}|\boldsymbol{\theta})} d\mathbf{Z} \quad (2.124)$$

where  $\mathbf{Z}$  is all the things we do not know (excluding the covariance hyperparameters  $\boldsymbol{\theta}$ ); for example, the hidden states  $\mathbf{X}_{1:T}$ , functions  $f$ ,  $g$  and any inducing variables for the sparse GP approximations. This is a lower bound to the marginal likelihood  $p(\mathbf{Y}_{1:T}|\boldsymbol{\theta})$  and so we use it to approximate the marginal likelihood and optimise it with respect to  $q$ , any other variational parameters and  $\boldsymbol{\theta}$ . Optimising with respect to  $q$  without any constraints would give the optimal  $q$  as  $p(\mathbf{Y}_{1:T}, \mathbf{Z}|\boldsymbol{\theta})$  resulting in an intractable integral. Hence, the different variational methods introduce different additional assumptions on  $q$ . Having decided on the additional assumptions about  $q$ , the ELBO is optimised using some numerical techniques; for example, gradient ascent. Damianou et al. [2011] use the GPLVMs applied to GPSSMs idea from Wang et al. [2008] but additionally marginalise out the hidden states  $\mathbf{X}$  using variational techniques and as a result overfitting is less of a problem compared to Wang et al. [2008]. Other uses of variational inference include Frigola et al. [2014]; Frigola [2015] which provide two variational methods (although certain calculations require a particle smoother) that use sparse GPs and learn from challenging data reasonable well and relatively fast (couple of minutes for 500 data points [Frigola et al., 2014]). Also, they allow for online learning [Frigola, 2015, pp. 66-67]. The main issue with all these variational methods is that they are not fully Bayesian and only provide point estimates of the parameters. Furthermore, they require that the variational assumptions are approximately correct. Other variational models include the model by Eleftheriadis et al. [2017] that combines variational methods and recurrent neural networks in order to efficiently learn a GPSSM.

Another set of methods for learning GPSSMs are based on sequential Monte Carlo and in particular Particle Gibbs with Ancestor Sampling (PGAS); for example, both Svensson et al. [2016] and Frigola et al. [2013] use PGAS to learn the hyperparameters of a GPSSM. The models suggested in these two papers both use a fully Bayesian approach (i.e. we get posterior distributions over all the hyperparameters) with Svensson et al. [2016] improving the work of Frigola et al. [2013] by using the Hilbert reduced-rank approximation of the GP to speed up learning quite considerably.



We will now focus on the model by Svensson et al. [2016]. Here, we will describe the model and in the next chapters we will attempt to improve upon this model. Suppose we have a set of observations  $\mathbf{y}_{1:T} = \mathbf{y}_1, \dots, \mathbf{y}_T$ , a reduced-rank GPSSM is a model of the form [Svensson et al., 2016]:

$$\mathbf{x}_t = \mathbf{A}\phi(\mathbf{x}_{t-1}) + \mathbf{w}_t \quad (2.125a)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t|\mathbf{x}_t) \quad (2.125b)$$

$$\mathbf{x}_1 \sim p(\mathbf{x}_1) \quad (2.125c)$$

$$\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (2.125d)$$

$$\mathbf{Q} \sim \mathcal{IW}(l_Q, \Lambda_Q) \quad (2.125e)$$

$$\mathbf{A}|\mathbf{Q} \sim \mathcal{MN}(\mathbf{A}|\mathbf{0}, \mathbf{Q}, \mathbf{V}) \quad (2.125f)$$

$$\mathbf{V} = \text{diag}(S^{-1}(\sqrt{\lambda_1}), \dots, S^{-1}(\sqrt{\lambda_m})) \quad (2.125g)$$

$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta}) \quad (2.125h)$$

#### Notes:

1.  $S$  is the spectral density (2.45) of the chosen covariance function for the transition GP. As in Hilbert reduced-rank GPs, this covariance function must be stationary.
2.  $\phi(\mathbf{x}_t) = [\phi_1(\mathbf{x}_t), \dots, \phi_m(\mathbf{x}_t)]^T$  with  $\phi_j, \lambda_j$  the eigenfunctions/eigenvalues used to approximate the covariance function. See (2.46)–(2.48) and (2.51) if required.
3. For simplicity,  $p(\mathbf{y}_t|\mathbf{x}_t)$  and  $p(\mathbf{x}_1)$  are assumed to be known; although, the model can be adapted to accommodate unknown parameters in these distributions as well.
4.  $\mathcal{IW}$  is the inverse Wishart distribution and  $\mathcal{MN}$  is the matrix-normal distribution described in (2.14).
5. The model comes from transforming the GPs to their reduced rank form using the Karhunen-Loève expansion: see (2.61).
6.  $\boldsymbol{\theta}$  represents all the covariance hyperparameters and  $p(\boldsymbol{\theta})$  is the prior distribution of  $\boldsymbol{\theta}$ .

Following Svensson et al. [2016], let

$$\boldsymbol{\Phi} = \sum_{t=1}^T \boldsymbol{\zeta}_t \boldsymbol{\zeta}_t^T \quad (2.126)$$

$$\boldsymbol{\Psi} = \sum_{t=1}^T \boldsymbol{\zeta}_t \mathbf{z}_t^T \quad (2.127)$$

$$\boldsymbol{\Sigma} = \sum_{t=1}^T \mathbf{z}_t \mathbf{z}_t^T \quad (2.128)$$

$$(2.129)$$

where  $\mathbf{z}_t = [\phi_1(\mathbf{x}_t), \dots, \phi_m(\mathbf{x}_t)]^T$ ,  $\boldsymbol{\zeta}_t = \mathbf{x}_{t+1}$ . Then some posterior distributions of  $\mathbf{A}$ ,  $\mathbf{Q}$  and  $\boldsymbol{\theta}$  given data  $\mathbf{y}_{1:T} = \mathbf{y}_1, \dots, \mathbf{y}_T$  are [Svensson et al., 2016]:

$$p(\mathbf{Q}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) = \mathcal{IW}(\mathbf{Q}|T + l_Q, \Lambda_Q + \boldsymbol{\Phi} - \boldsymbol{\Psi}(\boldsymbol{\Sigma} + \mathbf{V})^{-1}\boldsymbol{\Psi}^T) \quad (2.130)$$

$$p(\mathbf{A}|\mathbf{Q}, \mathbf{x}_{1:T}\mathbf{y}_{1:T}) = \mathcal{MN}(\mathbf{A}|\Psi(\Sigma + \mathbf{V})^{-1}, \mathbf{Q}, (\Sigma + \mathbf{V})^{-1}) \quad (2.131)$$

$$p(\boldsymbol{\theta}|\mathbf{Q}, \mathbf{A}, \mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \propto p(\boldsymbol{\theta})p(\mathbf{Q}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})p(\mathbf{A}|\mathbf{Q}, \mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \quad (2.132)$$

---

**Algorithm 3:** Learning of reduced-rank GPSSMs [Svensson et al., 2016]

---

**Input** : data  $\mathbf{y}_{1:T}$ , priors on  $\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}$   
**Output:** K samples with  $p(\mathbf{x}_{1:T}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}|\mathbf{y}_{1:T})$  as invariant distribution.

- 1 Sample initial  $\mathbf{x}_{1:T}[0], \mathbf{A}[0], \mathbf{Q}[0], \boldsymbol{\theta}[0]$  ;
- 2 **for**  $k = 0, \dots, K - 1$  **do**
- 3     sample  $\mathbf{x}_{1:T}[k + 1]|\mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]$  using PGAS [Lindsten et al., 2014, p. 2160];
- 4     sample  $\mathbf{Q}[k + 1]|\mathbf{A}[k], \boldsymbol{\theta}[k], \mathbf{x}_{1:T}[k + 1], \boldsymbol{\theta}[k]$  using (2.130);
- 5     sample  $\mathbf{A}[k + 1]|\mathbf{Q}[k + 1], \mathbf{x}_{1:T}[k + 1], \boldsymbol{\theta}[k]$  using (2.131);
- 6     sample  $\boldsymbol{\theta}[k + 1]|\mathbf{x}_{1:T}[k + 1], \mathbf{A}[k + 1], \mathbf{Q}[k + 1]$  using Metropolis-Hastings
- 7 **end**
- 8 return all samples of  $\mathbf{x}_{1:T}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}$

---

The training algorithm for reduced-rank GPSSMs is provided in Algorithm 3. From the output of this algorithm, we discard some of the early samples (burn-in) and then use the remaining samples to compute statistics of certain parameters such as the mean and variance. In particular, the mean and (marginal) variance of the predictive transition function given a test state  $\mathbf{x}_*$  is [Svensson et al., 2016]:

$$\mathbb{E}[f_*(\mathbf{x}_*)] \approx \text{mean}[\mathbf{A}]\phi(\mathbf{x}_*) \quad (2.133)$$

$$\mathbb{V}[f_*(\mathbf{x}_*)] \approx \text{diag}(\text{cov}[\mathbf{A}]\phi(\mathbf{x}_*) \text{cov}[\mathbf{A}]^T) \quad (2.134)$$

where  $\phi$  is as in (2.125a).

The state-of-the-art methods for learning GPSSMs all have a number of desirable properties; however, they vary in their expressiveness and how quickly they can learn. An ideal method for learning a GPSSM would have the following properties:

1. Allows for real-time online parameter learning.
2. Can efficiently learn from a wide range of possible datasets.
3. Can produce a predictive distribution of multi-step-ahead observations in real-time.
4. Robustness to outliers.
5. Allows for non-Gaussian latent and observation noise.
6. Can learn from large datasets in reasonable time.

The rest of this thesis will be focused on improving some of the state-of-the-art methods for learning GPSSMs by focusing on these ideal properties.

# Chapter 3

## Contributions

In this chapter, we introduce several novel contributions to the theory of Hilbert reduced-rank Gaussian processes and Hilbert reduced-rank Gaussian process state-space models.

### 3.1 Extending Hilbert reduced-rank Gaussian Processes

A key limitation of reduced-rank Gaussian processes is the difficulty in using anything more than the most basic covariance functions. In theory, the equations we described in (2.46)–(2.48) can be applied to any stationary covariance function but in practice it can be challenging to derive the spectral density. Moreover, although it is theoretically possible to extend the Hilbert reduced-rank GP to use any covariance function by changing the eigenfunctions and eigenvalues in (2.46), finding these eigenvalues and eigenfunctions is hard. To mitigate these issues we propose two solutions. Firstly, we introduce a new *adaptive* kernel that removes the need to find spectral densities while being well suited for use in reduced-rank GPs. Secondly, we combine the Hilbert reduced-rank GP model with the work of Calandra et al. [2016] to form the manifold-Hilbert reduced-rank Gaussian process model and this can introduce non-stationarity to stationary kernels.

#### 3.1.1 The Adaptive Kernel

In their derivation of Hilbert reduced-rank GPs, Solin and Särkkä [2014] define the following inner-product and operator:

$$\langle f, g \rangle = \int f(\mathbf{x})g(\mathbf{x})w(\mathbf{x})d\mathbf{x} \quad (3.1)$$

$$\mathcal{K}f(\mathbf{x}) = \int k(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')w(\mathbf{x}')d\mathbf{x}' \quad (3.2)$$

with  $w(\mathbf{x})$  some positive weight function such that  $\int w(\mathbf{x})$  exists (i.e. is finite). Then they note that since  $\mathcal{K}$  is self-adjoint with respect to the above inner-product (covariance functions are symmetric) via the spectral theorem there exists a decomposition

of  $\mathcal{K}$  into eigenvalues and orthogonal (with respect to (3.1)) eigenfunctions. Solin and Särkkä [2014] use this to write the covariance function as

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j_1, \dots, j_d} \gamma_{j_1, \dots, j_d} \phi_{j_1, \dots, j_d}(\mathbf{x}) \phi_{j_1, \dots, j_d}(\mathbf{x}') \quad (3.3)$$

for some positive eigenvalues  $\gamma_{j_1, \dots, j_d} \in \mathbb{R}_{>0}$ , orthogonal eigenfunctions  $\phi_{j_1, \dots, j_d}$  and inputs  $\mathbf{x} \in \mathbb{R}^d$ . Then they try to find these eigenvalues/eigenfunctions for particular covariance functions.

We use (3.3) to motivate the construction of a new covariance function: the *adaptive* covariance function. Instead of trying to find the eigenvalues/eigenfunctions for a particular covariance function, we fix a set of orthogonal eigenfunctions and let the eigenvalues become covariance function hyperparameters. This allows us to create some expressive covariance functions that naturally fit into the Hilbert reduced-rank GP framework while eliminating the need to find spectral densities, which can be challenging for certain covariance functions.

Solin and Särkkä [2014] required the covariance functions to be isotropic-stationary because it allowed them to find the eigenvalues/eigenfunctions in a simple manner. Furthermore, it turns out that the same set of eigenfunctions can be used for all *stationary* covariance functions [Solin and Särkkä, 2014]. However, we note that any valid covariance function will be self-adjoint with respect to (3.1) and hence can be decomposed in the form of (3.3). As a result, given the right set of eigenfunctions, we can construct any covariance function. In the adaptive covariance, we fix the eigenfunctions and this means that we are restricting ourselves to a particular *family* of covariance functions. The corresponding eigenvalues are covariance hyperparameters and so the adaptive covariance function can adapt itself into any covariance function of this chosen family.

**Definition 3.1** (The Adaptive Covariance Function). For inputs  $\mathbf{x} \in \mathbb{R}^d$  and given a fixed set of eigenfunctions  $\{\phi_{j_1, \dots, j_d}; j_i = 1, \dots, m_i; i = 1, \dots, d;\}$  such that  $\phi_{j_1, \dots, j_d} : \mathbb{R}^d \rightarrow \mathbb{R}$ , the adaptive covariance function is:

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j_1, \dots, j_d=1}^{\mathbf{m}} \gamma_{j_1, \dots, j_d} \phi_{j_1, \dots, j_d}(\mathbf{x}) \phi_{j_1, \dots, j_d}(\mathbf{x}') \quad (3.4)$$

where  $\gamma_{j_1, \dots, j_d} > 0$  are hyperparameters,  $\mathbf{m} = (m_1, \dots, m_d)$  and  $\mathbf{1} = (1, \dots, 1) \in \mathbb{R}^d$ . The  $\mathbf{m}$  determines the complexity with larger  $m_i$  creating a more complex and expressive kernel with increased computational cost. As in (2.46), if  $d > 1$  we can map  $(j_1, \dots, j_d) \rightarrow \mathbb{N}$  via (2.51). We also introduce a shorthand notation:

$$k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^{\mathbf{m}} \gamma_j \phi_j(\mathbf{x}) \phi_j(\mathbf{x}') \quad (3.5)$$

where  $\mathbf{j} = (j_1, \dots, j_d)$ . This equation is equivalent to (3.4).

**Theorem 1.** *The adaptive covariance function is a valid scalar-valued covariance function.*

*Proof.*

1. It is real:  $k(\mathbf{x}, \mathbf{x}') = \sum_{j=1}^m \gamma_j \phi_j(\mathbf{x}) \phi_j(\mathbf{x}') \in \mathbb{R}$
2. It is symmetric:

$$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= \sum_{j=1}^m \gamma_j \phi_j(\mathbf{x}) \phi_j(\mathbf{x}') \\ &= \sum_{j=1}^m \gamma_j \phi_j(\mathbf{x}') \phi_j(\mathbf{x}) \\ &= k(\mathbf{x}', \mathbf{x}) \end{aligned}$$

since  $\phi_j(\mathbf{x}) \in \mathbb{R}$ .

3. It is positive definite: given an arbitrary collection of inputs  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ , the matrix  $K = K(\mathbf{X}, \mathbf{X})$  with  $ij$  entry  $k(\mathbf{x}_i, \mathbf{x}_j)$  can be written as  $K = \Phi \Lambda \Phi^T$  where  $\Phi, \Lambda$  are defined in (2.49). Since the covariance hyperparameters (and hence eigenvalues) are constrained to be positive, we have that  $K$  is positive definite. □

Learning a reduced-rank Gaussian process with this adaptive kernel is simple, given an eigenbasis (i.e. given  $\phi_j$ ) replace (2.46) with (3.4) and let the  $\gamma_{j_1, \dots, j_d}$  in (3.4) become hyperparameters, then optimise the marginal likelihood [Solin and Särkkä, 2014, p. 10] with respect to the  $\gamma_{j_1, \dots, j_d}$  and the noise parameters.

One useful family of eigenfunctions (the  $\phi_j$ ) for the adaptive kernel is the Fourier eigenbasis from (2.48). This leads to every stationary covariance function being a special case of our adaptive covariance function because any (piecewise-continuous) periodic covariance function can be decomposed into these eigenfunctions (Fourier sine series) and non-periodic (piecewise-continuous) functions can be decomposed into this form as long as we restrict the domain (represented by lengths  $L_k$  in (2.48)). This means for periodic covariance functions the lengths  $L_k$  in (2.48) should become additional parameters so we can find the periodicity automatically and in this case we should be able to extrapolate to data outside of the domain  $\Omega$  (defined as in (2.48)). On the other hand, in the non-periodic case the  $L_k$  must be fixed in a manner that allows the domain  $\Omega$  to encompass all the training and test data and we cannot extrapolate outside of  $\Omega$ . This inability to extrapolate is usually not an issue, but as we shall see combining the *adaptive* covariance function with a manifold Gaussian process (mGP) [Calandra et al., 2016] can allow for extrapolation of non-periodic covariance functions.

### 3.1.2 The Manifold-Hilbert reduced-rank Gaussian Process

The manifold Gaussian process (mGP) uses a simple idea: given a valid covariance function  $k(\mathbf{x}, \mathbf{x}')$  where  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^d$  and an arbitrary function  $M : \mathbb{R}^d \rightarrow \mathbb{R}^f$  then

$k(M(\mathbf{x}), M(\mathbf{x}'))$  is also a valid covariance function [Mackay, 1998]. Calandra et al. [2016] introduced the mGP and they used a neural network to create the mapping function while finding the neural network parameters during hyperparameter optimisation. It turns out that a combination of the reduced-rank GP, adaptive kernel and the mGP can produce some very interesting results. From the perspective of the reduced-rank GP, using the mGP is equivalent to having an adaptive eigenbasis:

$$\hat{\phi}_{j_1, \dots, j_d} = \phi_{j_1, \dots, j_d} \circ M \quad (3.6)$$

$$k_{\mathbf{m}}(\mathbf{x}, \mathbf{x}') \approx \sum_{(j_1, \dots, j_d)=1}^{\mathbf{m}} S(\sqrt{\lambda_{j_1, \dots, j_d}}) \hat{\phi}_{j_1, \dots, j_d}(\mathbf{x}) \hat{\phi}_{j_1, \dots, j_d}(\mathbf{x}') \quad (3.7)$$

However, with this greatly increased flexibility comes an increased computational cost. We will refer to this combination of the mGP and reduced-rank GP as the manifold-Hilbert reduced-rank GP or the reduced-rank mGP.

### 3.1.3 Examples

In this section, we compare the Hilbert reduced-Rank GP and the manifold-Hilbert reduced-Rank GP to the full GP (the standard case from (2.12)) across a variety of kernels. Throughout this set of examples, we will only consider univariate input and output GPs. The performance of the various models will be tested using the root mean square error (on test data) and the mean log likelihood of the test data. Both of these statistics are defined below.

Given a trained model and a test dataset  $\{(x_i^*, y_i^*) \text{ for } i = 1, \dots, K\}$ , we can construct the posterior predictive distribution of  $f_* = f(x_i^*)$  given the training data  $\mathbf{X}, \mathbf{y}$  and a test input  $x_i^*$  using (2.29). Suppose that the posterior predictive distribution of  $f_* | \mathbf{X}, \mathbf{y}, x_i^*$  is  $\mathcal{N}(\mu_i^*, \sigma_i^{*2})$ , then we define:

1. The root mean square error (RMSE) (smaller is better):

$$RMSE = \sqrt{\frac{1}{K} \sum_{i=1}^K (y_i^* - \mu_i^*)^2} \quad (3.8)$$

2. The mean log likelihood (LL) (larger is better):

$$LL = \frac{1}{K} \sum_{i=1}^K \log \mathcal{N}(y_i^* | \mu_i^*, \sigma_i^{*2}) \quad (3.9)$$

Finally, all the models will be given 10 optimisation restarts, all reduced-rank models will use 12 basis functions and the training time will be defined as the mean time for an optimisation restart.

**Example 1:**

*Aim:* It is known from Calandra et al. [2016] that the mGP can produce better results than the full GP; however, the mGP is slower than using the full GP. The aim of this example is to demonstrate that the manifold-Hilbert reduced-rank GP retains the properties of the original mGP but is now faster than the full GP.

*Training Data:* 2000 pairs  $(x, y)$  with each  $x$  generated by sampling from  $\mathcal{N}(0, 1)$  and  $y$  generated using system:

$$y = \text{sgn}(x) + w_t \quad (3.10)$$

$$w_t \sim \mathcal{N}(0, 1^2) \quad (3.11)$$

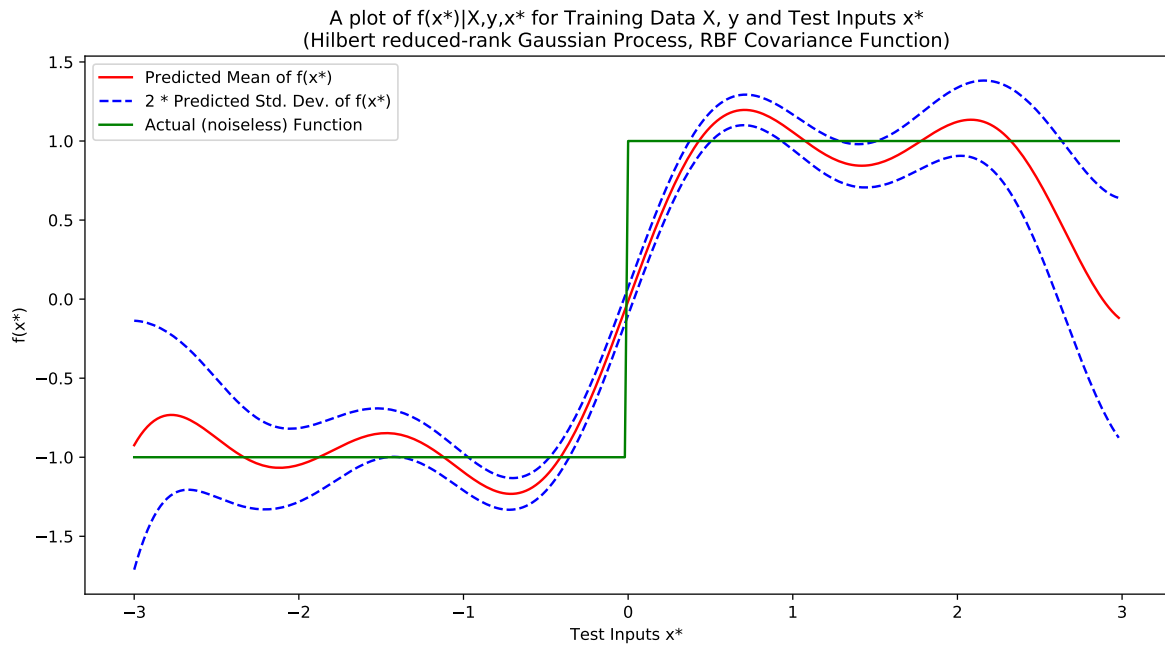
*Test Data:* 300 pairs  $(x, y)$  with  $x$  equally spaced between  $-3$  and  $3$  and  $y$  generated as above.

*Additional Information:* All the reduced-rank GP models will use a domain length  $L = 6$  and all the manifold-Hilbert reduced-rank models will use a neural network with 3 hidden layers of 1, 6 and 2 neurons respectively with log-sigmoid activations.

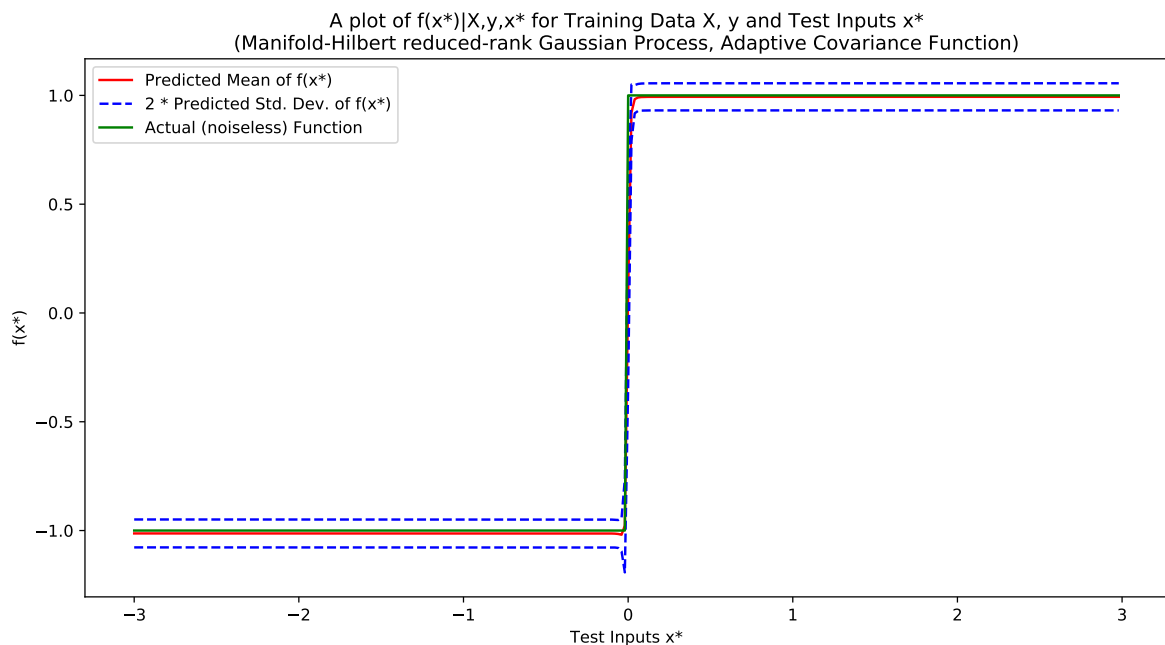
*Results:*

Model	Kernel	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP	RBF	1.10	-1.51	29.41
Full GP	Matern32	1.09	-1.51	34.64
Full GP	MLP	<b>1.04</b>	<b>-1.46</b>	160.00
Hilbert reduced-rank GP	RBF	1.10	-1.52	0.34
Hilbert reduced-rank GP	Matern32	1.10	-1.52	<b>0.31</b>
Hilbert reduced-rank GP	Adaptive	1.10	-1.52	0.95
Manifold-Hilbert reduced-rank GP	RBF	<b>1.04</b>	<b>-1.46</b>	5.49
Manifold-Hilbert reduced-rank GP	Matern32	<b>1.04</b>	<b>-1.46</b>	5.81
Manifold-Hilbert reduced-rank GP	Adaptive	<b>1.04</b>	<b>-1.46</b>	9.54

**Table 3.1:** Results for Example 1



**Figure 3.1:** A Hilbert reduced-rank GP with an RBF covariance function (table 3.1, row 4)



**Figure 3.2:** A manifold-Hilbert reduced-rank GP with an adaptive covariance function (table 3.1, row 9)



**Example 2:**

*Aim:* To demonstrate that the manifold-Hilbert reduced-rank GP can learn some challenging periodic functions.

*Training Data:* 1200 pairs  $(x, y)$  with  $x$  equally spaced between  $-15$  and  $15$  and  $y$  generated using system:

$$y = \text{sgn}(\cos(0.8x)) + w_t \quad (3.12)$$

$$w_t \sim \mathcal{N}(0, 0.5^2) \quad (3.13)$$

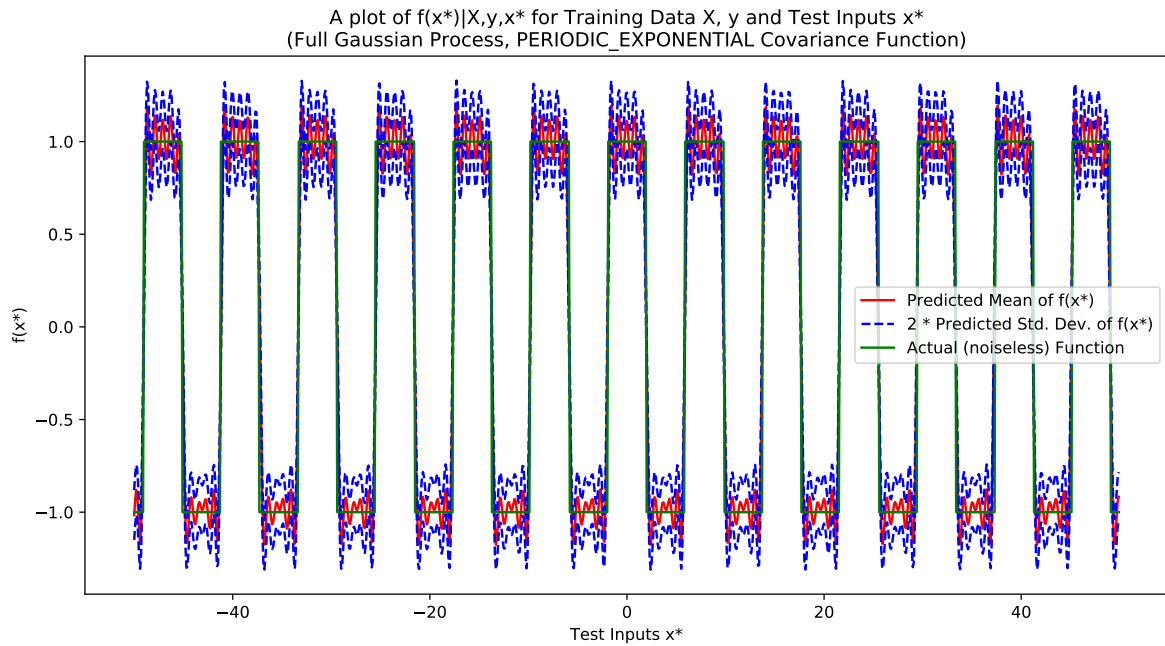
*Test Data:* 1000 pairs  $(x, y)$  with  $x$  equally spaced between  $-50$  and  $50$  and  $y$  generated as above.

*Additional Information:* All the reduced-rank GP models will use a variable domain length (i.e. we optimise the domain length) and all the manifold-Hilbert reduced-rank models will use a neural network with 3 hidden layers of 1, 6 and 2 neurons respectively with  $\sin(x)$  activation functions.

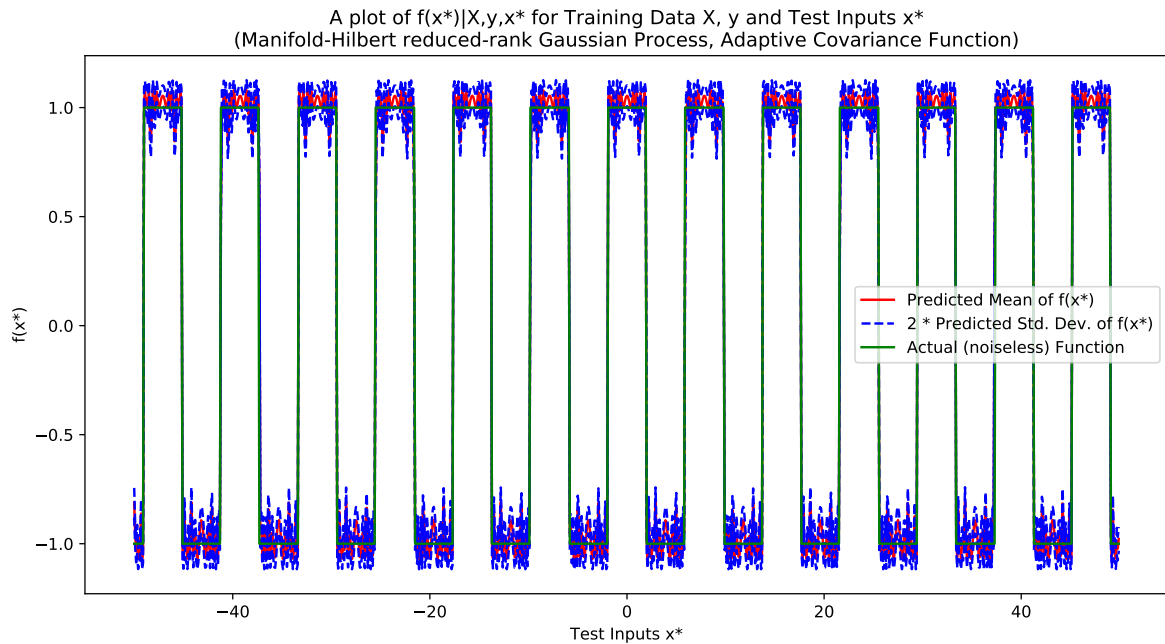
*Results:*

Model	Kernel	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP	RBF	0.99	-1.35	10.17
Full GP	Matern32	0.99	-1.33	6.71
Full GP	Periodic Exponential	0.56	-0.84	23.71
Hilbert reduced-rank GP	RBF	0.60	-0.92	0.37
Hilbert reduced-rank GP	Matern32	0.60	-0.92	<b>0.34</b>
Hilbert reduced-rank GP	Adaptive	0.60	-0.92	1.23
Manifold-Hilbert reduced-rank GP	RBF	<b>0.52</b>	<b>-0.76</b>	6.81
Manifold-Hilbert reduced-rank GP	Matern32	<b>0.52</b>	<b>-0.76</b>	7.48
Manifold-Hilbert reduced-rank GP	Adaptive	<b>0.52</b>	<b>-0.76</b>	8.12

**Table 3.2:** Results for Example 2



**Figure 3.3:** A full GP with a periodic exponential covariance function (table 3.2, row 3)



**Figure 3.4:** A manifold-Hilbert reduced-rank GP with an adaptive covariance function (table 3.2, row 9)

**Example 3:**

*Aim:* To demonstrate that the adaptive covariance function can sometimes give better results than the RBF or Matern32 covariance functions.

*Training Data:* 1000 pairs  $(x, y)$  with  $x$  equally spaced between  $-1$  and  $1$  and  $y$  generated using system:

$$y = \left| \frac{\sin(x)}{x} \right| + w_t \quad (3.14)$$

$$w_t \sim \mathcal{N}(0, 0.1^2) \quad (3.15)$$

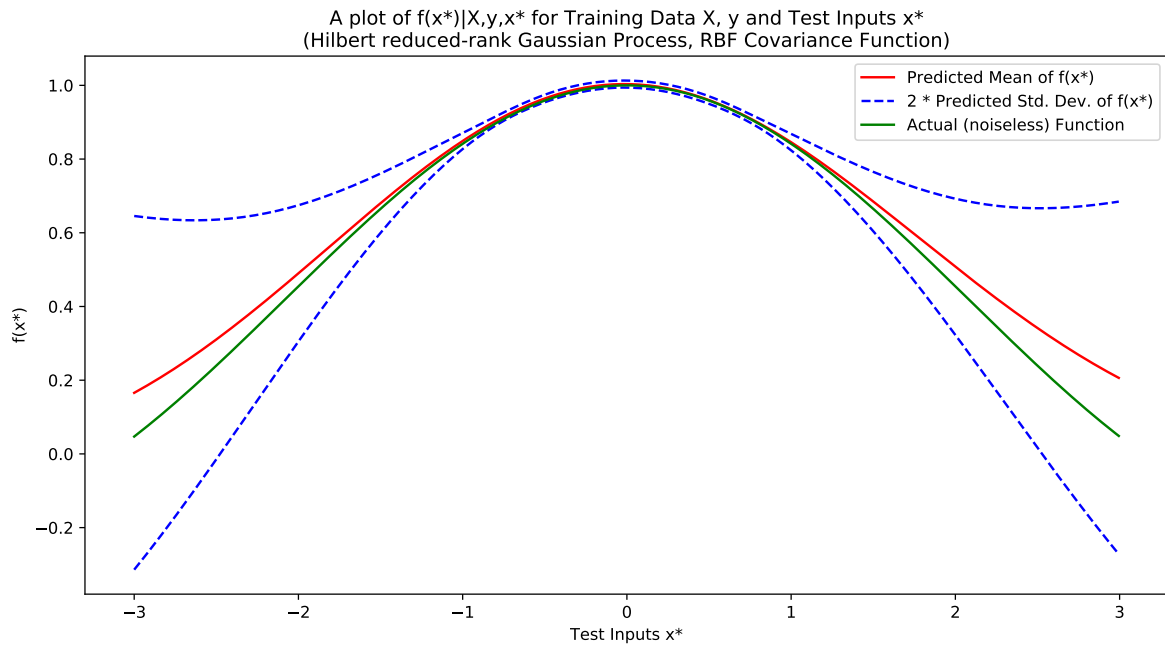
**Test Data:** 1200 pairs  $(x, y)$  with  $x$  equally spaced between  $-3$  and  $3$  and  $y$  generated as above.

*Additional Information:* All the reduced-rank GP models will use a variable domain length.

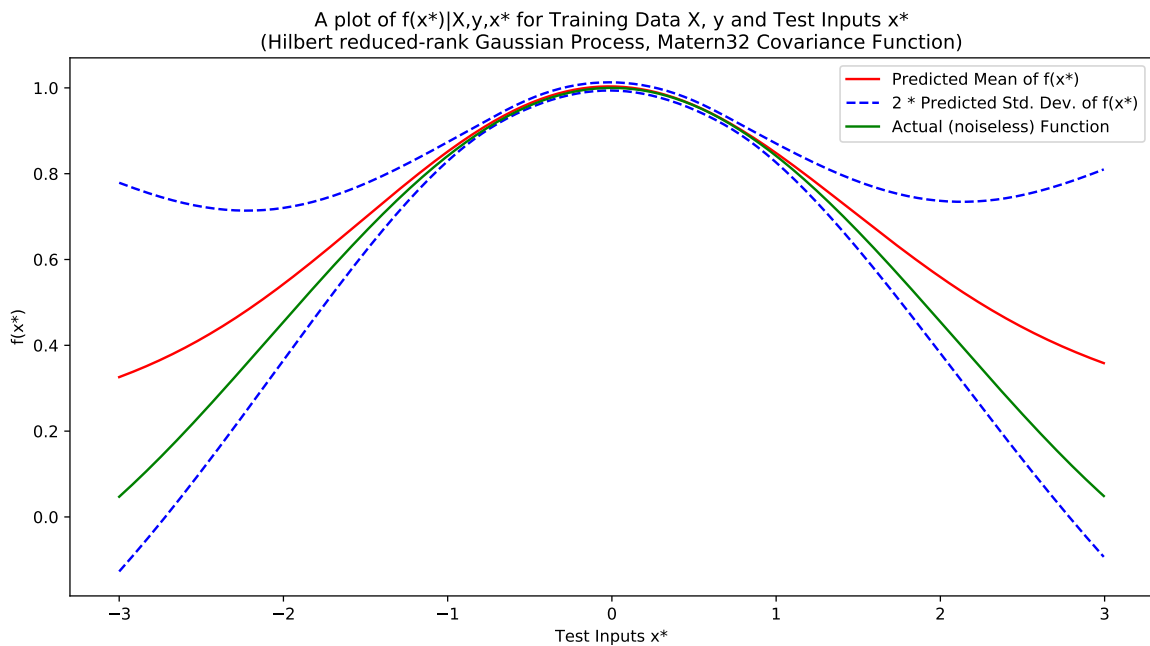
*Results:*

Model	Kernel	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP	RBF	0.03	1.16	7.00
Full GP	Matern32	0.17	0.77	7.89
Full GP	Periodic Exponential	0.18	0.52	17.25
Hilbert reduced-rank GP	RBF	0.06	1.10	<b>0.54</b>
Hilbert reduced-rank GP	Matern32	0.12	0.96	<b>0.54</b>
Hilbert reduced-rank GP	Adaptive	<b>0.02</b>	<b>1.23</b>	2.89

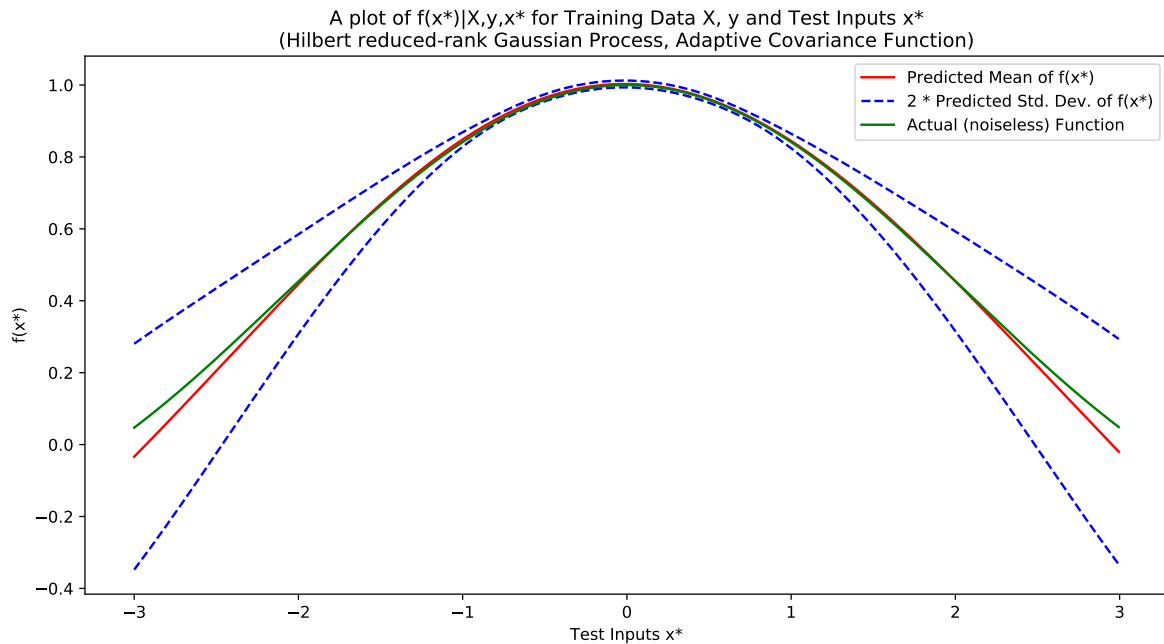
**Table 3.3:** Results for Example 3



**Figure 3.5:** A Hilbert reduced-rank GP with an RBF covariance function  
(table 3.3, row 4)



**Figure 3.6:** A Hilbert reduced-rank GP with a Matern32 covariance function  
(table 3.3, row 5)



**Figure 3.7:** A Hilbert reduced-rank GP with an adaptive covariance function (table 3.3, row 6)

#### Discussion:

These examples demonstrate that using the adaptive covariance function does incur a computation cost; however, its performance rivals picking the best covariance function for the job. This suggests that it could be very useful for the latent GP in a GPSSM: we cannot observe any of the latent data so trying to pick the best covariance function can be challenging. These examples also show that the manifold-Hilbert reduced-rank GP has retained the properties of the original mGP but is much faster.

## 3.2 Extending Hilbert reduced-rank GPSSMs

The work by Svensson et al. [2016] suggests that reduced-rank GPSSMs could be very useful for time series analysis; however, there are several desirable features that the model by Svensson et al. [2016] lacks. In this section, we introduce several novel learning methods for reduced-rank GPSSMs such as online learning, forgetful online learning, distributed learning, learning with student-t distributed transition noise and learning with Student-t process state-space models (STPSSMs). All these new methods augment the model and allow it to be used in many new situations; for example, in reinforcement learning. Throughout this section, we will assume that the basic GPSSM model is the same as in (2.125).

### 3.2.1 Online Learning

Currently, the reduced-rank GPSSM by Svensson et al. [2016] does not allow for online learning and so here we present a novel algorithm for online learning in Hilbert reduced-rank GPSSMs. Online learning means being able to update the model when we receive a new data point without having to retrain the whole model; in a sense, online learning can be viewed as efficient updates in the presence of new data. These online updates are useful for applications where we do not receive all the data in one go (e.g. reinforcement learning) because GPSSMs do take time to train: around a few minutes for 1000 single-dimension observations. Without online learning, every time we added new data, we would have to spend some time waiting for the model to update and it would not be realistic to continuously repeat the process of receiving new data and updating the model. However, we propose a new method that reduces update times to less than a second (depending on a variety of factors).

Before we introduce our new algorithm, we require a part of the original algorithm by Svensson et al. [2016], which is summarised in Algorithm 4. Svensson et al. [2016] gave a suggestion on how the model might be adapted for online learning: use the posterior distribution at time  $t$  for the prior when we receive the new data at time  $t + 1$ . For example, the update for (2.130) would be:

$$p(\mathbf{Q}|\mathbf{x}_{1:t+1}, \mathbf{y}_{1:t+1}) \propto p(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}|\mathbf{Q}, \mathbf{x}_{1:t}, \mathbf{y}_{1:t})p(\mathbf{Q}|\mathbf{x}_{1:t}, \mathbf{y}_{1:t}) \quad (3.16)$$

However, there is no point trying to make this distribution update efficient because we have to completely recompute (2.126)–(2.128) since we have a new set of states from the PGAS sample. A similar situation occurs for the other posterior distributions (2.131) and (2.132). Since we know (2.130)–(2.131) for arbitrary  $t$  and since there no need to produce efficient updates, we can just use the same distributions as before. However, the quality of the samples from the posterior distributions will increase as the amount of data increases and so we need to introduce sample weights. For our model, a given sample's weight is equal to the number of data points used for creating that sample and with online learning these weights will increase over time. When we compute various statistics based on the samples (e.g. mean or variance), we must take the weights into account i.e. calculate the mean/variance with samples weighted using these given weights.

### The Online Learning Algorithm

The algorithm is split into two parts; firstly, the inner loop (Algorithm 4), which samples from the latent space and parameter space given the observations and secondly, the main learning algorithm (Algorithm 5) which combines Algorithm 4 with our method for updating the GPSSM given new data.

### Tuning Parameters

In addition to the tuning parameters for the model by Svensson et al. [2016], there is one tuning parameter for online learning: the number of update samples. This is the number of times to run the sampling algorithm (Algorithm 4) when we receive a new data point. We recommend only using one sampling round per data update unless the data update provides a significant amount of extra data ( $\geq 5$  pieces of data). In a similar way to the model by Svensson et al. [2016], we assume that the initial distribution of the latent state  $\mathbf{x}_1$  is known but this has very little importance so if we do not know the initial distribution, we set it to  $\mathcal{N}(\mathbf{0}, 100\mathbf{I})$  for a suitably sized identity matrix  $\mathbf{I}$ .

### Using the Posterior Samples

As the online learning algorithm (Algorithm 5) progresses, we have some samples  $\mathbf{A}[1 : K]$ ,  $\mathbf{Q}[1 : K]$ ,  $\boldsymbol{\theta}[1 : K]$  and weights  $\omega[1 : K]$ . Let  $\omega[k]$  for  $k = 1, \dots, K$  be the normalised weights, then useful statistics such as the predictive mean and (marginal) variance of the transition function  $f$  given a test state  $\mathbf{x}_*$  can be computed as:

$$\text{mean}[A] = \sum_{k=\text{burn-in}}^K \mathbf{A}[k] \omega[k] \quad (3.17)$$

$$\mathbb{E}[f_*(\mathbf{x}_*)] \approx \text{mean}[A] \phi[\mathbf{x}_*] \quad (3.18)$$

$$\text{cov}[A] = \left( \sum_{k=\text{burn-in}}^K \omega[k] \mathbf{A}[k] \mathbf{A}[k]^T \right) - \text{mean}[A] \text{mean}[A]^T \quad (3.19)$$

$$\mathbb{V}[f_*(\mathbf{x}_*)] \approx \text{diag}(\text{cov}[A] \phi[\mathbf{x}_*] \text{cov}[A]^T) \quad (3.20)$$

where  $\phi$  is as in (2.125a). This can all be computed and updated during the algorithm letting us see how the model learns as it receives more data.

### 3.2.2 Forgetful Online Learning

Another type of online learning is *forgetful online learning* whereby while we add data, we also lose some data. This can be useful for systems that are evolving over time and is effective at keeping the amount of data in the system constant. This can be helpful because as the amount of data in the GPSSM increases, a slow down will be experienced. However, even if we start losing some past data, we do not lose all the past information because we retain all the latent state and parameter samples. The forgetful online learning algorithm uses a similar idea to the original online learning algorithm but it has to take into account that we lose information. As a result, there is an additional tuning parameter called the *Maximum Memory Length*:

this is the maximum amount of data to use for each call to InnerLoop (Algorithm 4). The forgetful online learning algorithm is summarised in Algorithm 6 and for performance statistics we can do something similar to the online learning case. For simplicity, we will keep number of samples per update (see tuning parameters for online learning) at one but the algorithm can be extended to multiple samples per update.

### 3.2.3 Distributed Learning

It is simple to extend all the above algorithms to make use of multicore/multiprocessor machines and this is summarised in Algorithm 8. Let  $\omega^i[k]$  for  $k = 1, \dots, K$  be the normalised (over both  $k$  and  $i$ ) weights for each thread  $i$ , then useful statistics such as the predictive mean and (marginal) variance of the transition function can be computed as:

$$\text{mean}[\mathbf{A}] = \sum_{i=1}^{\#threads} \left( \sum_{k=\text{burn-in}}^K \mathbf{A}^i[k] \omega^i[k] \right) \quad (3.21)$$

$$\mathbb{E}[f_*(\mathbf{x}_*)] \approx \text{mean}[\mathbf{A}] \phi[\mathbf{x}_*] \quad (3.22)$$

$$\text{cov}[\mathbf{A}] = \sum_{i=1}^{\#threads} \left( \sum_{k=\text{burn-in}}^K \mathbf{A}^i[k] (\mathbf{A}^i[k])^T \omega^i[k] \right) - \text{mean}[\mathbf{A}] \text{mean}[\mathbf{A}]^T \quad (3.23)$$

$$\mathbb{V}[f_*(\mathbf{x}_*)] \approx \text{diag}(\text{cov}[\mathbf{A}] \phi[\mathbf{x}_*] \text{cov}[\mathbf{A}]^T) \quad (3.24)$$

where  $\phi$  is as in (2.125a) and  $\mathbf{x}_*$  is a test point.

---

#### Algorithm 4: Inner Loop (based on Svensson et al. [2016])

---

**Input** :  $\mathbf{x}_{1:t}[k]$ ,  $\mathbf{A}[k]$ ,  $\mathbf{Q}[k]$ ,  $\boldsymbol{\theta}[k]$ ,  $\mathbf{y}_{1:t}$ , InitialDistribution  
**Output**:  $\mathbf{x}_{1:t}[k+1]$ ,  $\mathbf{A}[k+1]$ ,  $\mathbf{Q}[k+1]$ ,  $\boldsymbol{\theta}[k+1]$

- 1 sample  $\mathbf{x}_{1:t}[k+1] | \mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]$  using PGAS [Lindsten et al., 2014, p. 2160] with observations  $\mathbf{y}_{1:t}$ , reference trajectory  $\mathbf{x}_{1:t}[k]$  and initial distribution: InitialDistribution;  
 /\*  $\mathbf{A}[k], \mathbf{Q}[k]$  with (2.125a) gives us the transition distribution for PGAS and (2.125b) (which is assumed to be completely known) gives us the observation distribution. \*/
- 2 sample  $\mathbf{Q}[k+1] | \mathbf{A}[k], \boldsymbol{\theta}[k], \mathbf{x}_{1:t}[k+1], \boldsymbol{\theta}[k]$  using (2.130) with  $T = t$ ;
- 3 sample  $\mathbf{A}[k+1] | \mathbf{Q}[k+1], \mathbf{x}_{1:t}[k+1], \boldsymbol{\theta}[k]$  using (2.131) with  $T = t$ ;
- 4 sample  $\boldsymbol{\theta}[k+1] | \mathbf{x}_{1:t}[k+1], \mathbf{A}[k+1], \mathbf{Q}[k+1]$  using (2.132) with  $T = t$  and Metropolis-Hastings;
- 5 **return**  $\mathbf{x}_{1:t}[k+1], \mathbf{A}[k+1], \mathbf{Q}[k+1], \boldsymbol{\theta}[k+1]$ ;

---



**Algorithm 5:** Online learning of a reduced-rank GPSSM

---

**Input** : Initial data  $\mathbf{y}_{1:t}$ , Number of samples per update, InitialDistribution  
**Output**:  $K$  samples with the  $k^{\text{th}}$  sample having  $p(\mathbf{x}_{1:t_k}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}_A | \mathbf{y}_{1:t_k})$  as the invariant distribution.

- 1 initialise  $\mathbf{x}_{1:t}$ ,  $\mathbf{A}[0]$ ,  $\mathbf{Q}[0]$ ,  $\boldsymbol{\theta}[0]$  randomly;
- 2 set weights[0] = t;
- 3 set  $t_0 = t$ ;
- 4 set  $\mathbf{x}_{1:t_0}^0 = \mathbf{x}_{1:t}$ ;
- 5 set  $k = 1$ ;
- 6 **while**  $k \leq K$  **do**
- 7     newdata = getData();
- 8      $i$  = number of new data pieces;
- 9      $t_k = t_{k-1} + i$  ;
- 10    append the new data to the old data:  $\mathbf{y}_{1:t_k} = (\mathbf{y}_{1:t_{k-1}}, \text{new data})$ ;
- 11    **for**  $r = 1, \dots, i$  **do**
- 12     sample newstates[ $r$ ] from  $\mathcal{N}(\bar{\mathbf{x}}_{t_{k-1}}, \mathbf{I}_{d_{\text{Latent}}})$  where  $\bar{\mathbf{x}}_{t_{k-1}}$  is the state at time  $t_{k-1}$  (weighted average over all samples using weights[ $k-1$ ]) and  $d_{\text{Latent}}$  is the dimension of the latent space;  
     /\* Using  $\mathbf{I}_{d_{\text{Latent}}}$  rather than the sample covariance ensures that we always have some exploration and this helps PGAS. \*/
- 13    **end**
- 14    set the new reference trajectory:  $\hat{\mathbf{x}}_{1:t_k} = (\mathbf{x}_{1:t_{k-1}}^{k-1}, \text{newstates})$ ;
- 15     $\mathbf{x}_{1:t_k}^k, \mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k] = \text{InnerLoop}(\hat{\mathbf{x}}_{1:t_k}, \mathbf{A}[k-1], \mathbf{Q}[k-1], \boldsymbol{\theta}[k-1], \mathbf{y}_{1:t_k}, \text{InitialDistribution})$ ;
- 16    weights[ $k$ ] =  $t_k$ ;
- 17    **for**  $r = 1, \dots, \text{Number of samples per update} - 1$  **do**
- 18      $t_{k+r} = t_{k-1+r}$  ;
- 19     weights[ $k+r$ ] =  $t_{k+r}$ ;
- 20      $\mathbf{x}_{1:t_{k+r}}^{k+r}, \mathbf{A}[k+r], \mathbf{Q}[k+r], \boldsymbol{\theta}[k+r] = \text{InnerLoop}(\mathbf{x}_{1:t_{k+r-1}}^{k+r-1}, \mathbf{A}[k+r-1], \mathbf{Q}[k+r-1], \boldsymbol{\theta}[k+r-1], \mathbf{y}_{1:t_k}, \text{InitialDistribution})$ ;
- 21    **end**
- 22     $k = k + \text{Number of samples per update}$ ;
- 23 **end**
- 24 **return**  $\mathbf{x}_{1:t_1}^1, \dots, \mathbf{x}_{1:t_K}^K, \mathbf{A}[1 : K], \mathbf{Q}[1 : K], \boldsymbol{\theta}[1 : K]$  and weights[1 :  $K$ ].

---

**Algorithm 6:** Forgetful online learning of reduced-rank GPSSM

---

**Input** : Initial data  $\mathbf{y}_{1:t}$ , Maximum Memory Length  
**Output**:  $K$  samples with the  $k^{\text{th}}$  sample having  $p(\mathbf{x}_{t_k^{\text{start}}:t_k^{\text{end}}}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{y}_{t_k^{\text{start}}:t_k^{\text{end}}})$  as the invariant distribution.

- 1 Initialise  $\mathbf{x}_{1:t}$ ,  $\mathbf{A}[0]$ ,  $\mathbf{Q}[0]$ ,  $\boldsymbol{\theta}[0]$  randomly;
- 2 `weights[0] = t;`
- 3 `( $t_0^{\text{start}}$ ,  $t_0^{\text{end}}$ ) = (1, t);`
- 4  `$\mathbf{x}_{t_0^{\text{start}}:t_0^{\text{end}}}^0 = \mathbf{x}_{1:t}$ ;`
- 5 **for**  $k = 1, \dots, K$  **do**
- 6     `newdata = getData();`
- 7     `i = number of new data pieces;`
- 8      `$t_k^{\text{end}} = t_{k-1}^{\text{end}} + i$ ;`
- 9      `$t_k^{\text{start}} = \max(t_k^{\text{end}} - \text{Maximum Memory Length}, 0) + 1$ ;`
- 10    `weights[k] =  $t_k^{\text{end}} - t_k^{\text{start}}$ ;`
- 11    append the new data to the old data:  $\mathbf{y}_{1:t_k^{\text{end}}} = (\mathbf{y}_{1:t_{k-1}^{\text{end}}}, \text{new data})$ ;
- 12    **for**  $r = 1, \dots, i$  **do**
- 13     sample newstates[ $r$ ] from  $\mathcal{N}(\bar{\mathbf{x}}_{t_{k-1}}, \mathbf{I}_{d_{\text{Latent}}})$  where  $\bar{\mathbf{x}}_{t_{k-1}}$  is the state at time  $t_{k-1}^{\text{end}}$  (weighted average over all samples using weights[ $k-1$ ]) and  $d_{\text{Latent}}$  is the dimension of the latent space;  
    /\* Using  $\mathbf{I}_{d_{\text{Latent}}}$  rather than the sample covariance ensures that we always have some exploration and this helps PGAS.  
    \*/
- 14    **end**
- 15     `$\mathbf{y}^* = \mathbf{y}_{t_k^{\text{start}}:t_k^{\text{end}}}$ ;`
- 16     `$\mathbf{x}^* = (\mathbf{x}_{t_k^{\text{start}}:t_{k-1}^{\text{end}}}^{k-1}, \text{newstates})$ ;`
- 17     `$\mu = \text{mean}(\mathbf{x}_{t_k^{\text{start}}})$ ;`
- 18     `$\Sigma = \text{cov}(\mathbf{x}_{t_k^{\text{start}}})$ ;`  
    /\* We compute the mean and variance using the samples of the latent states at time  $t_k^{\text{start}}$ . Some of the latent state samples will not have a sample for time  $t_k^{\text{start}}$  and so they are ignored.     \*/
- 19    NewInitialDistribution =  $\mathcal{N}(\mu, \Sigma)$
- 20     `$\mathbf{x}_{t_k^{\text{start}}:t_k^{\text{end}}}^k | \mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k] = \text{Inner Loop}(\mathbf{x}^*, \mathbf{A}[k-1], \mathbf{Q}[k-1], \boldsymbol{\theta}[k-1], \mathbf{y}^*, \text{NewInitialDistribution})$ ;`
- 21 **end**
- 22 **return** `( $\mathbf{x}_{t_1^{\text{start}}:t_1^{\text{end}}}^1, \dots, \mathbf{x}_{t_K^{\text{start}}:t_K^{\text{end}}}^K, \mathbf{A}[1:K], \mathbf{Q}[1:K], \boldsymbol{\theta}[1:K]$  and weights[1:K].`

---

**Algorithm 7:** Distributed learning of a reduced-rank GPSSM

---

```

Input : Data  $\mathbf{y}_{1:T}$ , #threads
Output:  $K$  samples with  $p(\mathbf{x}_{1:T}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{y}_{1:T})$  as the invariant distribution.
1 split the observed data  $\mathbf{y}_{1:T}$  into #threads chunks of approximately equal size;
  /* One way to split is sequentially i.e. first chunk is  $\mathbf{y}_{1:t_1}$ ,
    second is  $\mathbf{y}_{t_1+1:t_2}$  with  $t_1 < t_2$ . An alternative way is ordered
    random splitting in which each chunk contains a random ordered
    set of  $\mathbf{y}_{1:T}$ ; for example, if the first chunk is  $\mathbf{y}_{t_1}, \mathbf{y}_{t_2}, \dots, \mathbf{y}_{t_m}$  then
    the times must satisfy  $t_1 < t_2 < \dots < t_m$ . Usually we would expect
    the chunks of data to be disjoint; however, this is not a
    requirement. */
2 let these data chunks be  $\mathbf{y}^1, \dots, \mathbf{y}^{\#threads}$ ;
3 for each  $\mathbf{y}^i$  in parallel do
4    $\mathbf{x}_{1:T_i}^i[1 : K], \mathbf{A}^i[1 : K], \mathbf{Q}^i[1 : K], \boldsymbol{\theta}^i[1 : K] = \text{LearningAlgorithm}(\mathbf{y}^i, \text{priors}$ 
     on  $\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}$ );
     /* LearningAlgorithm can either be the standard reduced rank
       learning of Svensson et al. [2016] or our online/forgetful
       online learning algorithm. For the chosen algorithm the
       observed data will be the  $\mathbf{y}^i$  only and the prior
       distributions should be the same for each thread. */
5   if not known already calculate the weight of each sample;
     /* For online/forgetful online learning we will already have a
       set of weights but for the standard algorithm there are no
       sample weights so we set the weights of each sample equal to
       the number of data points in the corresponding data chunk.
       */
6   discard burn in samples;
     /* This must be done for each thread. */
7 end
8 return weight-sample pairs of  $\mathbf{A}, \mathbf{Q}$  and  $\boldsymbol{\theta}$  for all threads.

```

---

### 3.2.4 GPSSMs with Student-t noise

In this section, we extend the model of Svensson et al. [2016] to allow for Student-t noise in the latent states. Learning with Student-t noise is more challenging than with Gaussian noise because we no longer have closed form posterior distributions for  $p(\mathbf{Q}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  or  $p(\mathbf{A}|\mathbf{Q}, \mathbf{x}_{1:T}, \mathbf{y}_{1:T})$ . As a result, we either have to use Monte Carlo methods to sample from these posterior distributions or variational methods to find a closed form approximation to these posterior distributions that we can sample from. Since we wanted to keep the same structure as the model by Svensson et al. [2016], it seemed reasonable to try and find a Monte Carlo solution.

**Definition 3.2.** [Roth, 2013]

Let  $\mathbf{X} \in \mathbb{R}^d$  be a multivariate-t distributed random variable with mean  $\boldsymbol{\mu} \in \mathbb{R}^d$ , symmetric positive definite scale parameter  $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$  and  $\nu \in \mathbb{R}_{>0}$  degrees of freedom. The probability density of  $\mathbf{X}$  at the point  $\mathbf{x} \in \mathbb{R}^d$  is :

$$\mathcal{T}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \nu) = \frac{\Gamma(\frac{\nu+d}{2})}{\Gamma(\frac{\nu}{2})} \frac{1}{(\nu\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \left( 1 + \frac{1}{\nu} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)^{-\frac{\nu+d}{2}} \quad (3.25)$$

A key property of the multivariate-t distribution is that when the degrees of freedom  $\nu$  tends to infinity, the distribution tends to a multivariate normal distribution with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ .

A GPSSM with Student-t noise can be written as follows:

$$\mathbf{x}_t = \mathbf{A}\phi(\mathbf{x}_{t-1}) + \mathbf{w}_t \quad (3.26a)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t|\mathbf{x}_t) \quad (3.26b)$$

$$\mathbf{w}_t \sim \mathcal{T}(\mathbf{0}, \mathbf{Q}, \nu) \quad (3.26c)$$

$$\mathbf{A}|\mathbf{Q} \sim \mathcal{MN}(\mathbf{A}|\mathbf{0}, \mathbf{Q}, \mathbf{V}) \quad (3.26d)$$

$$\mathbf{Q} \sim \mathcal{IW}(l_Q, \boldsymbol{\Lambda}_Q) \quad (3.26e)$$

where  $\mathbf{V}$ ,  $\phi$  and the other terms are defined in (2.125a). We will now present an algorithm which uses the same Gibbs structure as the model by Svensson et al. [2016] but has different MCMC steps inside the Gibbs sampler to deal with the Student-t likelihood.

The first step is to find a suitable method of sampling from  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}|\mathbf{x}_{1:T})$  ( $\boldsymbol{\theta}$  is the kernel hyperparameters). To do this, the main options are either to try and sample from  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}|\mathbf{x}_{1:T})$  directly or to split it up into parts and use Gibbs sampling. However, unlike Svensson et al. [2016] who had closed form posterior distributions, splitting  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}|\mathbf{x}_{1:T})$  can cause us issues because traditional Monte Carlo methods, such as Metropolis-Hastings or hybrid Monte Carlo, and even newer methods, such as slice-sampling, can have poor performance when sampling from high-dimensional posteriors i.e. they can be slow at producing good samples. So either we split the  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}|\mathbf{x}_{1:T})$  into many low-dimensional parts and get very slow mixing in the Gibbs sampler or we split it into few high-dimensional parts and still get slow mixing because we are sampling from high-dimensional distributions. Note

that even if the latent dimension  $d = 1$ , then we still need to sample  $\mathbf{A}$  which will have  $m$  parameters if we use  $m$  basis functions. We tried Metropolis-Hastings, hybrid Monte Carlo and slice-sampling with none of them producing good results.

However, we were able to get good results by sampling from  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{x}_{1:T})$  directly using a recent MCMC method called the Affine-Invariant Ensemble Sampler which was originally developed by Goodman and Weare [2010] with a practical implementation (called EMCEE) produced by Foreman-Mackey et al. [2013]. EMCEE has several advantages compared with other methods; for example, it only has a few tuning parameters (two single dimensional parameters for most cases) unlike Metropolis-Hastings or its variants that have a performance highly dependant on the proposal distribution. Furthermore, EMCEE can give a measure of whether the drawn samples are ‘good’ samples from the distribution.

The Affine-Invariant Ensemble Sampler algorithm is relatively complex and is discussed in depth by Goodman and Weare [2010] and Foreman-Mackey et al. [2013]. Here, we will only give a brief overview of the EMCEE algorithm. Suppose we want to produce samples from the distribution  $p(\mathbf{x})$  but only have access to  $\hat{p}(\mathbf{x})$  where  $p(\mathbf{x}) = \frac{1}{Z}\hat{p}(\mathbf{x})$  and  $Z$  does not depend on  $\mathbf{x}$ . In this case, the EMCEE algorithm takes two (main) arguments:  $\log \hat{p}(\mathbf{x})$  and the initial positions of the ‘walkers’. The walkers can be seen as a similar idea to particles in SMC methods with more walkers leading to samples that more accurately represent the distribution. Each of the walkers perform their own semi-independent acceptance sampling and have their own Markov chain; however, more walkers allows EMCEE to better explore the sample space. One method suggested by Foreman-Mackey et al. [2013] to initialise the walkers is to start all the walkers very close to the point  $\mathbf{x}^*$  which maximises  $\log \hat{p}(\mathbf{x})$  (although the walkers must start in different places so random noise is added). After a brief burn-in period, each walker can start to generate samples from  $p(\mathbf{x})$ . Finally, the EMCEE algorithm provides a indicator of whether the samples are ‘good’ known as the acceptance fraction and ideally this should be between 0.2 and 0.5 (see Foreman-Mackey et al. [2013] for more details).

In order to use EMCEE for learning GPSSMs with Student-t latent noise we note:

$$p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{x}_{1:t}) \propto p(\mathbf{x}_{1:t} | \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}) p(\mathbf{A} | \mathbf{Q}, \boldsymbol{\theta}) p(\mathbf{Q} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (3.27)$$

$$\propto p(\mathbf{x}_1) \prod_{k=2}^t p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}) p(\mathbf{A} | \mathbf{Q}, \boldsymbol{\theta}) p(\mathbf{Q} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) \quad (3.28)$$

With Student-t noise the latent distribution is:

$$p(\mathbf{x}_k | \mathbf{x}_{k-1}, \mathbf{A}, \mathbf{Q}, \boldsymbol{\theta}) = \mathcal{T}(\mathbf{x}_k | \mathbf{A}\phi(\mathbf{x}_{k-1}), \mathbf{Q}, \nu) \quad (3.29)$$

Hence we can write that:

$$\begin{aligned} \log p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{x}_{1:t}) = Z + \sum_{k=2}^t \log \mathcal{T}(\mathbf{x}_k | \mathbf{A}\phi(\mathbf{x}_{k-1}), \mathbf{Q}, \nu) + \log \mathcal{MN}(\mathbf{A} | \mathbf{0}, \mathbf{Q}, \mathbf{V}) \\ + \log \mathcal{IW}(\mathbf{Q} | l_Q, \boldsymbol{\Lambda}_Q) + \log p(\boldsymbol{\theta}) \end{aligned} \quad (3.30)$$

Therefore, to produce samples from  $p(\mathbf{A}, \mathbf{Q}, \boldsymbol{\theta} | \mathbf{x}_{1:t})$  using EMCEE we use:

$$\log \hat{p} = \sum_{k=2}^t \log \mathcal{T}(\mathbf{x}_k | \mathbf{A}\phi(\mathbf{x}_{k-1}), \mathbf{Q}, \nu) + \log \mathcal{MN}(\mathbf{A} | \mathbf{0}, \mathbf{Q}, \mathbf{V}) \quad (3.31)$$

$$+ \log \mathcal{IW}(\mathbf{Q} | l_Q, \boldsymbol{\Lambda}_Q) + \log p(\boldsymbol{\theta}) \quad (3.32)$$

In many ways, the algorithm for learning GPSSMs with Student-t noise in the latent states is simple: all we do is switch between sampling the latent states using PGAS and sampling the parameters via EMCEE with occasional optimisation steps in order to initialise the EMCEE walkers. There is only one additional tuning parameter for this algorithm compared with that of Svensson et al. [2016]: the optimisation distance. This optimisation distance controls the number of sampling iterations before optimising  $\log \hat{p}$  again. Based on the examples we have tried, an optimisation distance of 5 is recommended since this balances speed (the optimisation step is slower than a sampling step) with performance.

Another thing to mention about our algorithm is that we only use two walkers in each call to EMCEE. Although this may seem at odds with what Foreman-Mackey et al. [2013] recommend (the number of walkers should be at least twice the dimension of the sample space), we are implicitly using more walkers since we initialise the start point at each call to EMCEE rather than using the last position of the previous walkers. Also, the number of samples we require at each stage is very small so there is little need for a large number of walkers (which would be able to produce a *large number* of quasi-independent samples that accurately represent the distribution).

Finally, we make use of automatic differentiation to find the gradient of (3.32) and hence we can use a gradient based optimisation method: this is much faster than a non-gradient based method such as Nelder-Mead. Our choice of optimisation method is L-BFGS-B because it allows for the use of bounds which is particularly important for several parameters. We can now write down the algorithm for learning GPSSMs with Student-t latent noise: see Algorithm 8. Furthermore, it is quite easy to adapt this algorithm to online learning: just use the same idea as in Algorithm 5.

**Algorithm 8:** Learning a GPSSM with Student-t latent noise

---

**Input** : Observed data  $\mathbf{y}_{1:T}$ , InitialDistribution  $p(\mathbf{x}_1)$ , *OptimisationDistance*.

**Output**:  $\mathbf{x}_{1:T}[k]$ ,  $\mathbf{A}[k]$ ,  $\mathbf{Q}[k]$ ,  $\boldsymbol{\theta}[k]$  for  $k = 1, \dots, K$ .

- 1 set  $\mathbf{x}_{1:T}[0]$ ,  $\mathbf{A}[0]$ ,  $\mathbf{Q}[0]$ ,  $\boldsymbol{\theta}[0]$  to arbitrary but valid (e.g.  $\mathbf{Q}[0]$  must be positive definite) values;
- 2 **for**  $k = 0, \dots, K-1$  **do**
- 3     sample  $\mathbf{x}_{1:T}[k+1] | \mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]$  using PGAS [Lindsten et al., 2014, p. 2160] with observations  $\mathbf{y}_{1:T}$ , reference trajectory  $\mathbf{x}_{1:T}[k]$  and initial distribution: InitialDistribution;  
      /\*  $\mathbf{A}[k], \mathbf{Q}[k]$  with (3.29) gives us the transition distribution for PGAS and (3.26c) (which is assumed to be completely known) gives us the observation distribution. \*/
- 4     **if**  $k \equiv 0 \pmod{\text{OptimisationDistance}}$  **then**
- 5         find the maximum  $(\mathbf{A}^*, \mathbf{Q}^*, \boldsymbol{\theta}^*)$  of (3.32) given  $\mathbf{x}_{1:T}[k+1]$  using L-BFGS-B with the maximum iterations not too large (50-100 is fine);  
       /\* The low maximum iterations ensure that we do not spend too much time in this step: for small  $k$  we will not have good samples of  $\mathbf{x}_{1:T}$  so there is no need to spend a lot of time in optimisation and for larger  $k$  we should not deviate too far from the maximum between each optimisation period. \*/
- 6         Set initial walker positions equal to  $(\mathbf{A}^*, \mathbf{Q}^*, \boldsymbol{\theta}^*) +$  vector of random noise;  
       /\* The random noise should be small ( $|\text{noise}_i| < 0.05$  for all  $i$ ) and independently chosen for each parameter of each walker (e.g. for each walker, each component of  $\mathbf{A}$  should have some independent noise). You must be careful that the noise does not lead to the initial walker positions being invalid; for example, degrees of freedom  $\nu$  must be positive and  $\mathbf{Q}$  must be positive definite. \*/
- 7     **end**
- 8     **else**
- 9         Set initial walker positions equal to  $(\mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]) +$  vector of random noise;  
       /\* above comments still apply \*/
- 10    **end**
- 11    sample  $\mathbf{A}[k+1], \mathbf{Q}[k+1], \boldsymbol{\theta}[k+1] | \mathbf{x}_{1:T}[k+1]$  using EMCEE and (3.32) with #walkers = 2, burn-in = 10;
- 12 **end**
- 13 **return**  $\mathbf{x}_{1:T}[1 : K]$ ,  $\mathbf{A}[1 : K]$ ,  $\mathbf{Q}[1 : K]$ ,  $\boldsymbol{\theta}[1 : K]$ ;

---

### 3.2.5 STPSSMs: A State-Space Model with a Heavy Tailed Process

Utilising the methods of the previous section allow us to present a new state-space model with a Student-t process rather than a Gaussian process for the prior on the latent state function  $f$ . Student-t processes (in the context of regression) were first introduced by Shah et al. [2014] and have several advantages over Gaussian processes; in particular, they are more flexible and are more robust against outliers because of the heavy-tailed nature of multivariate-t distribution [Shah et al., 2014]. Furthermore, Student-t processes also possess many of the important properties of Gaussian processes such as analytic marginal and predictive distributions [Shah et al., 2014]; however, to keep the analytic nature we have to incorporate the noise model into the kernel rather than having it separate (see Shah et al. [2014] for more details). If we do incorporate the noise model into the kernel then Student-t processes have the same computational cost as Gaussian processes [Shah et al., 2014].

**Definition 3.3.** Shah et al. [2014]

$f$  is a Student-t process with mean function  $m : \mathbb{R}^d \rightarrow \mathbb{R}$ , kernel function  $k : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$  and  $\nu \in \mathbb{R}_{>0}$  degrees of freedom, if any finite collection of function values satisfies:

$$(f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^T \sim \mathcal{T}(m(\mathbf{X}), K(\mathbf{X}, \mathbf{X}), \nu) \quad (3.33)$$

where  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T$ ,  $m(\mathbf{X}) = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_n))^T$  and  $K(\mathbf{X}, \mathbf{X})_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ . This is written as:  $f \sim \mathcal{TP}(m, k, \nu)$ .

On the whole, Student-t processes are very similar to Gaussian processes; for example, training still scales in  $\mathcal{O}(N^3)$  as the amount of data ( $N$ ) increases, and many of the approximation techniques used for Gaussian processes can be applied to Student-t processes. Also, for Student-t processes with multiple outputs, we can either train separate Student-t processes for each output dimension or use similar ideas (e.g. linear model of coregionalisation) as for Gaussian processes with multiple outputs. Furthermore, a key property of the Student-t process is that when the degrees of freedom  $\nu$  tend to infinity, the Student-t process tends towards a Gaussian process (proof in Shah et al. [2014]). One important difference is that previously we have used *kernel function* and *covariance function* interchangeably but in the case of the Student-t process technically  $k(\mathbf{x}, \mathbf{x}')$  is no longer the covariance. However, it is approximately the covariance and this approximation gets better as the degrees of freedom increases.

We will now look at reduced-rank Student-t processes. In the same way as reduced-rank Gaussian processes, the idea behind the reduced rank approximation is to find a decomposition of the kernel function  $\mathbf{K}$  into  $\mathbf{K} \approx \Phi \Lambda \Phi^T$  where we are assuming that the kernel function is stationary. The derivation of the basis function approximation for Student-t processes is exactly the same as for Gaussian processes since it relies on finding an approximation to the kernel function rather than any specific properties of Gaussian processes [Solin and Särkkä, 2014]. Hence, we can use the same decomposition of the kernel (2.46) and the same basis functions (3.70) combined with the matrix inversion lemma and the original posterior/predictive equations for  $f_*$  in Shah et al. [2014] to produce a reduced-rank Student-t process which



trains in  $\mathcal{O}(m^3)$  where  $m$  is the number of basis functions.

We will now define an important distribution required for reduced-rank Student-t processes: the matrixvariate-t distribution. As this distribution is not very common there are actually several competing probability density functions all of which have the multivariate-t distribution as a special case. Furthermore, not all the definitions of the matrixvariate-t distribution tend to the obvious matrix-normal distribution as the degrees of freedom tend to infinity. The definition we present below is slightly different to that in Gupta and Nagar [2000]; Zhu et al. [2008]; Iranmanesh et al. [2010] but tends to the obvious matrix-normal distribution as the degrees of freedom tends to infinity which is very important for Student-t processes.

**Definition 3.4.** Let  $\mathbf{X} \in \mathbb{R}^{m \times n}$  be a matrixvariate-t distributed random variable with mean  $\mathbf{M} \in \mathbb{R}^{m \times n}$ , symmetric positive definite row-scale parameter  $\Sigma \in \mathbb{R}^{m \times m}$ , symmetric positive definite column-scale parameter  $\Omega \in \mathbb{R}^{n \times n}$  and  $\nu \in \mathbb{R}_{>0}$  degrees of freedom. The probability density of  $\mathbf{X}$  at the point  $\mathbf{x} \in \mathbb{R}^{m \times n}$  is :

$$\mathcal{MT}(\mathbf{x}|\mathbf{M}, \Sigma, \Omega, \nu) = Z |\mathbf{I}_m + \frac{1}{\nu} \Sigma^{-1} (\mathbf{x} - \mathbf{M}) \Omega^{-1} (\mathbf{x} - \mathbf{M})^T|^{-\frac{\nu+m+n-1}{2}} \quad (3.34)$$

where

$$Z = \frac{\Gamma_n\left(\frac{\nu+m+n-1}{2}\right)}{\Gamma_n\left(\frac{\nu+n-1}{2}\right)} \frac{1}{(\nu\pi)^{\frac{mn}{2}} |\Sigma|^{\frac{n}{2}} |\Omega|^{\frac{m}{2}}} \quad (3.35)$$

and  $\Gamma_n(\cdot)$  is the  $n$ -multivariate Gamma function,  $\mathbf{I}_m$  is the  $m \times m$  identity matrix and  $|\cdot|$  is the determinant.

This definition has the property that:

$$\lim_{\nu \rightarrow \infty} \mathcal{MT}(\mathbf{x}|\mathbf{M}, \Sigma, \Omega, \nu) = \mathcal{MN}(\mathbf{x}|\mathbf{M}, \Sigma, \Omega) \quad (3.36)$$

For completeness, the definitions in Gupta and Nagar [2000]; Iranmanesh et al. [2010] have that a definition of the matrixvariate-t distribution that has the convergence property:

$$\lim_{\nu \rightarrow \infty} \mathcal{MT}(\mathbf{x}|\mathbf{M}, \Sigma, \nu\Omega, \nu) = \mathcal{MN}(\mathbf{x}|\mathbf{M}, \Sigma, \Omega) \quad (3.37)$$

so using our definition slightly simplifies matters.

Consider the regression problem:

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon_i \quad (3.38)$$

$$f \sim \mathcal{TP}(0, k, \nu) \quad (3.39)$$

where we have data  $(\mathbf{x}_i, \mathbf{y}_i)$  for  $i = 1, \dots, N$ ,  $\mathbf{y}_i \in \mathbb{R}^{d_y}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_x}$  and  $\epsilon_i$  is some noise distribution with location  $\mathbf{0}$  and scale  $\mathbf{Q}$ . Under the reduced rank approximation this can be written as [Solin and Särkkä, 2014]:

$$\mathbf{y}_i = \mathbf{A}\phi(\mathbf{x}_i) + \epsilon_i \quad (3.40)$$

$$\mathbf{A} \sim \mathcal{MT}(0, \mathbf{Q}, \mathbf{V}, \nu) \quad (3.41)$$

where  $\mathbf{A} \in \mathbb{R}^{d_y \times m}$ ,  $\phi(\mathbf{x}_i) = (\phi_1(\mathbf{x}_i), \dots, \phi_m(\mathbf{x}_i))^T$  for basis functions  $\phi_i$ ,  $\mathbf{Q}$  is the noise scale,  $\mathbf{V} = \text{diag}(S^{-1}(\sqrt{\lambda_1}), \dots, S^{-1}(\sqrt{\lambda_m}))$  where the  $\lambda_i$  are defined as in (2.47) and  $S$  is the spectral density of the kernel function (which is assumed to be stationary). However, unlike for Gaussian processes the posterior distribution of  $\mathbf{A}$  does not exist in closed form unless we incorporate the noise into the kernel function. If we wish to use PGAS to sample the latent states, then we cannot incorporate the noise into the kernel function. As a result, we are forced to use Monte Carlo methods to sample from the posterior of  $\mathbf{A}$  but fortunately this can be done by using the same method as in GPSSMs with Student-t latent noise: EMCEE.

We can now write down the equations of a Student-t process state-space model:

$$\mathbf{x}_t = \mathbf{A}\phi(\mathbf{x}_{t-1}) + \mathbf{w}_t \quad (3.42a)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t | \mathbf{x}_t) \quad (3.42b)$$

$$\mathbf{A} | \mathbf{Q} \sim \mathcal{MT}(\mathbf{A} | \mathbf{0}, \mathbf{Q}, \mathbf{V}, \nu_A) \quad (3.42c)$$

$$\mathbf{Q} \sim \mathcal{IW}(l_Q, \Lambda_Q) \quad (3.42d)$$

where  $\mathbf{w}_t$  can be any noise distribution with mean  $\mathbf{0}$  and scale  $\mathbf{Q}$ . For example, Gaussian noise:

$$\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (3.43)$$

or Student-t noise:

$$\mathbf{w}_t \sim \mathcal{T}(\mathbf{0}, \mathbf{Q}, \nu_w) \quad (3.44)$$

or even Laplace noise:

$$\mathbf{w}_t \sim \mathcal{L}(\mathbf{0}, \mathbf{Q}) \quad (3.45)$$

where  $\mathcal{L}$  is the multivariate Laplace distribution. Learning in all these models is very similar to learning in the GPSSM with Student-t noise: we use a Gibbs sampler with PGAS to sample the latent states and EMCEE to sample the parameter space. For EMCEE with latent states  $\mathbf{x}_{1:t}$  we would use:

$$\begin{aligned} \log \hat{p} = \sum_{k=2}^t \log p(\mathbf{x}_k | \mathbf{A}\phi(\mathbf{x}_{k-1}), \mathbf{Q}, \boldsymbol{\theta}) + \log \mathcal{MT}(\mathbf{A} | \mathbf{0}, \mathbf{Q}, \mathbf{V}, \nu_A) \\ + \log \mathcal{IW}(\mathbf{Q} | l_Q, \Lambda_Q) + \log p(\boldsymbol{\theta}) \end{aligned} \quad (3.46)$$

where  $p(\mathbf{x}_k | \mathbf{A}\phi(\mathbf{x}_{k-1}), \mathbf{Q}, \boldsymbol{\theta})$  is the transition pdf and it will depend on the chosen noise distribution.

We can now present an algorithm for online learning of Student-t process state-space models (STPSSMs). The algorithm uses the same idea as online learning for GPSSMs but is combined with the learning structure of GPSSM with latent Student-t noise. This can be easily extended for multiple data pieces per call to `getData()`, multiple PGAS/EMCEE samples per call to `getData()` or forgetful online learning by following the same ideas as in GPSSMs. Furthermore, this method can be extended to other function processes such as Dirichlet processes, Laplace processes and inverse Wishart processes by changing the prior distribution on  $\mathbf{A}$  to a matrixvariate-Dirichlet, matrixvariate-Laplace or an inverse Wishart distribution respectively.

**Algorithm 9:** Online learning of a STPSSM

---

```

Input : Initial observed data  $\mathbf{y}_{1:t_0}$ , InitialDistribution  $p(\mathbf{x}_1)$ ,
          OptimisationDistance.
Output:  $\mathbf{x}_{1:t_k}^k$ ,  $\mathbf{A}[k]$ ,  $\mathbf{Q}[k]$ ,  $\boldsymbol{\theta}[k]$  for  $k = 1, \dots, K$  and weights.
1 set  $\mathbf{x}_{1:t_0}^0$ ,  $\mathbf{A}[0]$ ,  $\mathbf{Q}[0]$ ,  $\boldsymbol{\theta}[0]$  to arbitrary but valid (e.g.  $\mathbf{Q}[0]$  must be positive
   definite) values;
2 set weights =  $[t_0]$ ;
3 for  $k = 0, \dots, K$  do
4   newData = getData();
5    $t_{k+1} = t_k + 1$ ;
6   append the new data to the old data:  $\mathbf{y}_{1:t_{k+1}} = (\mathbf{y}_{1:t_k}, \text{new data})$ ;
7   sample  $\hat{\mathbf{x}}$  from  $\mathcal{N}(\bar{\mathbf{x}}_{t_k}, \mathbf{I}_{d_{Latent}})$  where  $\bar{\mathbf{x}}_{t_k}$  is the state at time  $t_k$  (mean over
   all samples weighted using weights) and  $d_{Latent}$  is the dimension of the
   latent space;
8   append  $t_{k+1}$  to weights;
9   set  $\hat{\mathbf{x}}_{1:t_{k+1}} = (\mathbf{x}_{1:t_k}^k, \hat{\mathbf{x}})$ ;
10  sample  $\mathbf{x}_{1:t_{k+1}}^{k+1} | \mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]$  using PGAS [Lindsten et al., 2014, p. 2160]
   with observations  $\mathbf{y}_{1:t_{k+1}}$ , reference trajectory  $\hat{\mathbf{x}}_{1:t_{k+1}}$  and initial
   distribution: InitialDistribution;
   /*  $\mathbf{A}[k], \mathbf{Q}[k]$  with the noise distribution gives us the transition
      distribution for PGAS and (3.42c) (which is assumed to be
      completely known) gives us the observation distribution. */
11  if  $k \equiv 0 \pmod{\text{OptimisationDistance}}$  then
12    find the maximum  $(\mathbf{A}^*, \mathbf{Q}^*, \boldsymbol{\theta}^*)$  of (3.32) given  $\mathbf{x}_{1:t_{k+1}}^{k+1}$  using L-BFGS-B
      with the maximum iterations not too large (50-100 is fine);
13    Set initial walker positions equal to  $(\mathbf{A}^*, \mathbf{Q}^*, \boldsymbol{\theta}^*) +$  vector of random
      noise;
      /* see comments from learning GPSSMs with Student-t noise */
14  end
15  else
16    Set initial walker positions equal to  $(\mathbf{A}[k], \mathbf{Q}[k], \boldsymbol{\theta}[k]) +$  vector of
      random noise;
      /* see comments from learning GPSSMs with Student-t noise */
17  end
18  sample  $\mathbf{A}[k+1], \mathbf{Q}[k+1], \boldsymbol{\theta}[k+1] | \mathbf{x}_{1:t_{k+1}}^{k+1}$  using EMCEE and (3.32) with
      #walkers = 2, burn-in = 10;
19 end
20 return  $\mathbf{x}_{t_1^{start}:t_1^{end}}^1, \dots, \mathbf{x}_{t_K^{start}:t_K^{end}}^K$ ,  $\mathbf{A}[1:K]$ ,  $\mathbf{Q}[1:K]$ ,  $\boldsymbol{\theta}[1:K]$ , weights;
   /* In a similar manner to online Learning of GPSSMs, if we wish to
      compute values such as the mean or covariance of  $f(\mathbf{x}) = \mathbf{A}\phi(\mathbf{x})$ 
      then we should use a weighted mean and weighted covariance with
      the samples weighted by weights. */

```

---

### 3.2.6 Examples

In this section, we look at some examples of the various new features we have added to the model of Svensson et al. [2016]. We only consider models with fixed observation distributions in order to allow for an identifiable latent function and we only look at models with single-dimensional observation and latent spaces. Furthermore, we test the performance of the models by computing the root mean square error of the predicted latent states compared with the true latent states and by computing the mean log likelihood of the latent states given the learnt model. These statistics are defined below.

Given a trained model and a test set of latent states  $(x_1^*, \dots, x_T^*)$ , we can construct key statistics for the posterior predictive distribution of  $f(x_i^*)$  using the samples from the parameter posterior. This has already been looked at for online/forgetful online learning see: (3.18) and (3.20). Now, for learning with Student-t noise or Student-t processes, suppose that we have samples  $\theta_1, \dots, \theta_K$  from the parameter posterior (here we assume that  $\theta_i$  is a vector containing an  $\mathbf{A}$  sample, a  $\mathbf{Q}$  sample and a sample of each of the hyperparameters) and we write  $f_{\theta}(x)$  for the latent function  $f(x)$  with model parameters equal to  $\theta$ . Then we can define:

$$\mathbb{E}[f(x^*)] \approx \frac{1}{K} \sum_{i=1}^K f_{\theta_i}(x^*) \quad (3.47)$$

$$\mathbb{V}[f(x^*)] \approx \left( \frac{1}{K} \sum_{i=1}^K (f_{\theta_i}(x^*))^2 \right) - \mathbb{E}[f(x^*)]^2 \quad (3.48)$$

where these tend to the true results as  $K$  tends to infinity.

Now we can define:

1. The root mean square error (RMSE) (smaller is better):

$$RMSE = \sqrt{\frac{1}{T-1} \sum_{t=1}^{T-1} (x_{t+1}^* - \mathbb{E}[f(x_t^*)])^2} \quad (3.49)$$

2. For models which assume a Gaussian noise distribution for the latent states, the mean log likelihood (LL) is defined as (larger is better):

$$LL_{gauss} = \frac{1}{T-1} \sum_{t=1}^{T-1} \log \mathcal{N}(x_{t+1}^* | \mathbb{E}[f(x_t^*)], \mathbb{V}[f(x_t^*)] + \mathbb{E}[\mathbf{Q}]) \quad (3.50)$$

where  $\mathbb{E}[\mathbf{Q}] \approx \frac{1}{K} \sum_{i=1}^K \mathbf{Q}_{\theta_i}$  and  $\mathbf{Q}_{\theta_i}$  is the value of  $\mathbf{Q}$  when the model parameters are equal to  $\theta_i$ .

3. For models which assume a Student-t noise distribution for the latent states, calculating the mean log likelihood (LL) is more complex. The idea is still the

same: we wish to find the probability of the test data under the model so in this case, we define:

$$LL_{student-t} = \frac{1}{(T-1)K} \sum_{t=1}^{T-1} \sum_{i=1}^K \log \mathcal{T}(x_{t+1}^* | f_{\theta_i}(x_t^*), \mathbf{Q}_{\theta_i}, \nu_{\theta_i}) \quad (3.51)$$

where  $\mathcal{T}(x|\mu, \Sigma, \nu)$  is a multivariate-t density function (parametrised as in (3.25)) with mean  $\mu$ , scale  $\Sigma$  and  $\nu$  degrees of freedom. The notation  $\mathbf{Q}_{\theta_i}$  and  $\nu_{\theta_i}$  simply means the values of  $\mathbf{Q}_{\theta_i}$  and  $\nu_{\theta_i}$ , respectively, when the model parameters are equal to  $\theta_i$ .

### Example 1:

*Aim:* This is a ‘sanity-check’ example which allows us to compare the various models under a simple system. The idea is that we expect these models to produce similar results (RMSE and LL) not only to each other but also to the results stated in Svensson et al. [2016]. This should demonstrate that the different ways we compare particular performance statistics lead to the same results when the models are equivalent. Note that we do expect the training times to be different.

*Training Data:* A sequence of 500 observations generated by the system:

$$x_{t+1} = f(x_t) + w_t \quad (3.52)$$

$$y_t = x_t + e_t \quad (3.53)$$

$$w_t \sim \mathcal{N}(0, 1) \quad (3.54)$$

$$e_t \sim \mathcal{N}(0, 1) \quad (3.55)$$

where

$$f(x) = \begin{cases} x + 1 & x < 4 \\ -4x + 21 & x \geq 4 \end{cases} \quad (3.56)$$

The models only receive the observed states  $y_1, \dots, y_{500}$ . The latent states  $x_1, \dots, x_{500}$  are discarded.

*Test Data:* A sequence of 10000 observations and latent states generated by the above system.

*Additional Information:* The initial point will be zero (known), the observation distribution will be known beforehand (i.e. not learnt), all the reduced-rank GPSSMs will use 12 basis functions, 20 PGAS particles, domain length  $L = 12$ , the  $\mathbf{Q}$  prior hyperparameters will both be equal to 1 and we will use a burn-in of 50 samples. The offline models will have  $K = 180$  sampling rounds and the online models will have  $K = 200$  sampling rounds; furthermore, all the online/forgetful online models will start with two data points and receive an additional 5 data points every other sampling round (i.e. number of update samples = 2). Finally, all forgetful models will have a maximum memory length of 250 data points (see section 3.2.2).

Results:

Method	Kernel	Assumed Latent Noise	Process	RMSE	LL	Training Time (seconds)
Offline	RBF	Gaussian	Gaussian	1.13	-1.50	83.36
Online	RBF	Gaussian	Gaussian	1.12	-1.52	58.06
Forgetful	RBF	Gaussian	Gaussian	<b>1.11</b>	<b>-1.50</b>	<b>45.66</b>
Offline	Adaptive	Student-t	Gaussian	1.13	-1.53	126.01
Online	RBF	Student-t	Student-t	1.12	-1.59	77.50

Table 3.4: Results for Example 1

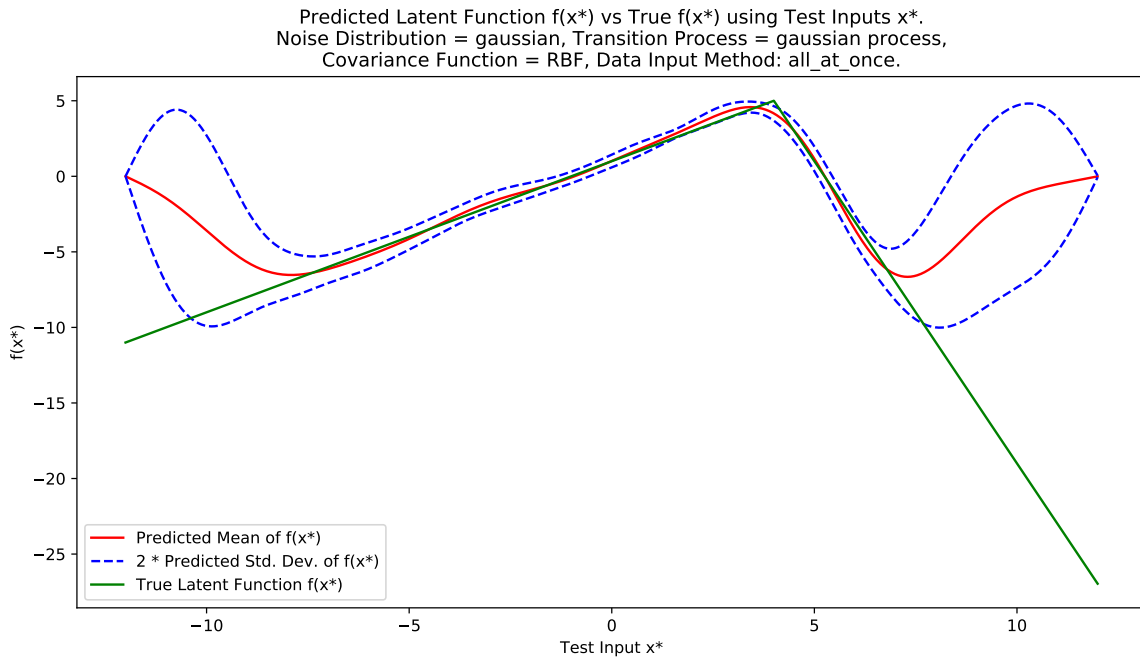
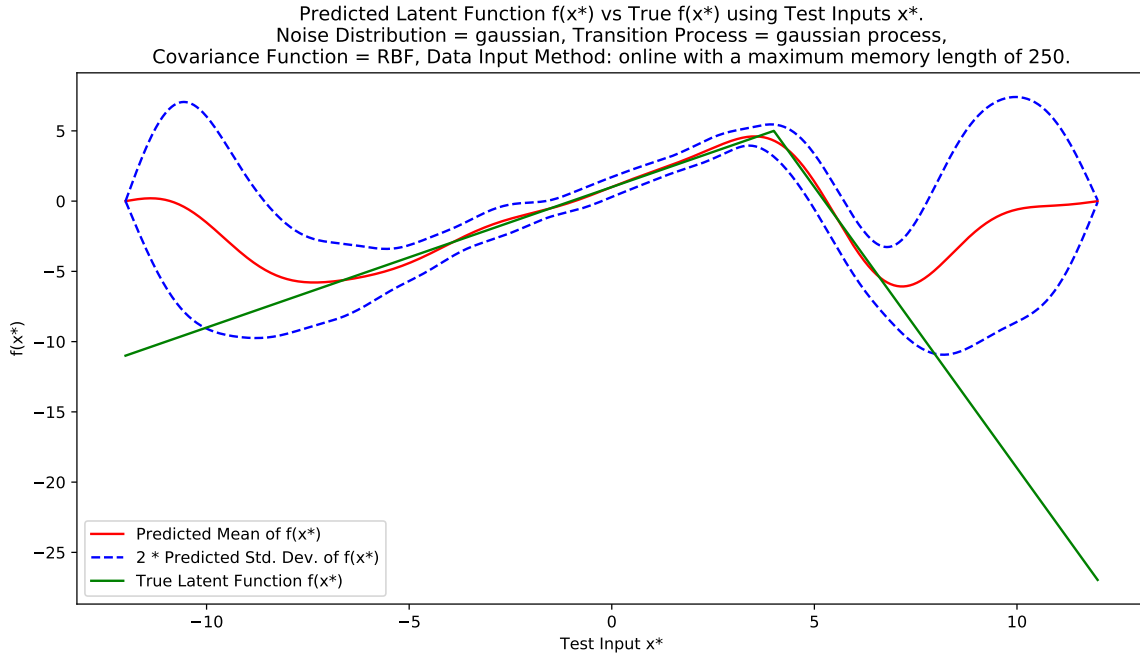


Figure 3.8: A GPSSM learnt using the algorithm of Svensson et al. [2016] (Table 3.4, Row 1)



**Figure 3.9:** A GPSSM learnt using our forgetful online learning algorithm (Table 3.4, Row 3)

### Example 2:

*Aim:* To compare the different models under a system with Student-t noise in the latent states.

*Training Data:* A sequence of 1000 observations generated by the system:

$$x_{t+1} = 0.5x_t + 3 \sin(x_t) + w_t \quad (3.57)$$

$$y_t = x_t + e_t \quad (3.58)$$

$$w_t \sim \mathcal{T}(0, 1, \nu = 1) \quad (3.59)$$

$$e_t \sim \mathcal{N}(0, 1) \quad (3.60)$$

The models only receive the observed states  $y_1, \dots, y_{1000}$ . The latent states  $x_1, \dots, x_{1000}$  are discarded.

*Test Data:* A sequence of 10000 observations and latent states generated by the above system.

*Additional Information:* The initial point will be zero (known), the observation distribution will be known beforehand (i.e. not learnt), all the reduced-rank GPSSMs will use 8 basis functions, 20 PGAS particles, domain length  $L = 10$ , the  $Q$  prior hyperparameters will both be equal to 1 and we will use a burn-in of 100 samples. The offline models will have  $K = 180$  sampling rounds and the online models will have  $K = 200$  sampling rounds; furthermore, all the online/forgetful online models will start with two data points and receive an additional 10 data points every other sampling round (i.e. number of update samples = 2). Finally, all forgetful models

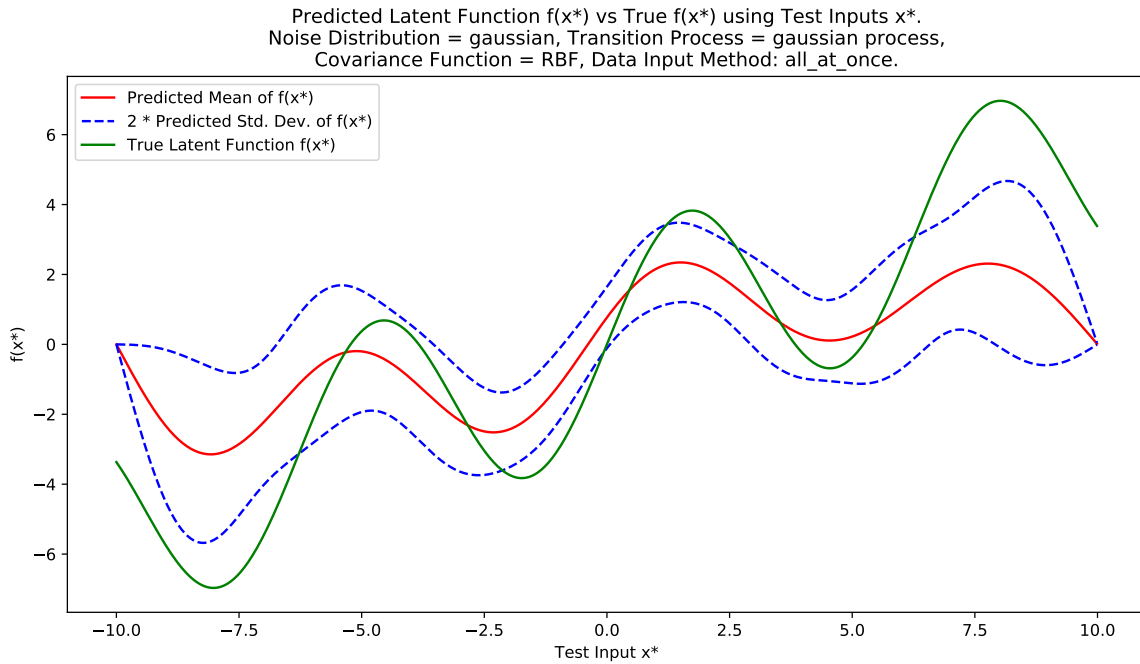
will have a maximum memory length of 250 data points (see section 3.2.2).

*Results:*

Since the true noise distribution is Student-t, we do not provide a RMSE because it is unhelpful under Student-t noise due to the presence of large outliers.

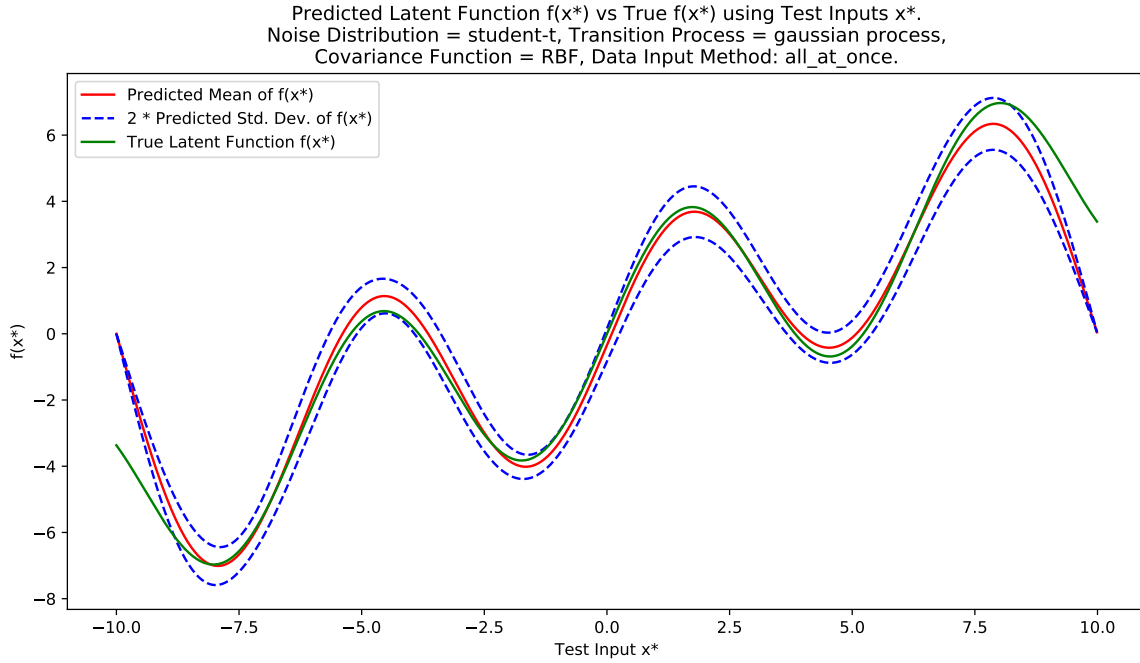
Method	Kernel	Assumed Latent Noise	Process	LL	Training Time (seconds)
Offline	RBF	Gaussian	Gaussian	-573.47	138.72
Online	RBF	Student-t	Gaussian	-17.30	121.85
Forgetful	RBF	Student-t	Gaussian	-3.23	<b>59.83</b>
Offline	RBF	Student-t	Gaussian	<b>-3.06</b>	175.64

**Table 3.5:** Results for Example 2.



**Figure 3.10:** A GPSSM with Student-t noise in the latent states learnt using the algorithm of Svensson et al. [2016] (Table 3.5, Row 1)





**Figure 3.11:** A GPSSM with Student-t noise in the latent states learnt using algorithm 8 (Table 3.5, Row 4)

### Example 3:

*Aim:* To compare some combinations of features we have not looked at previously.

*Training Data:* A sequence of 1000 observations generated by the system:

$$x_{t+1} = 3 \tanh(0.5x_t) + w_t \quad (3.61)$$

$$y_t = x_t + e_t \quad (3.62)$$

$$w_t \sim \mathcal{N}(0, 1.2) \quad (3.63)$$

$$e_t \sim \mathcal{N}(0, 0.8) \quad (3.64)$$

The models only receive the observed states  $y_1, \dots, y_{1000}$ . The latent states  $x_1, \dots, x_{1000}$  are discarded.

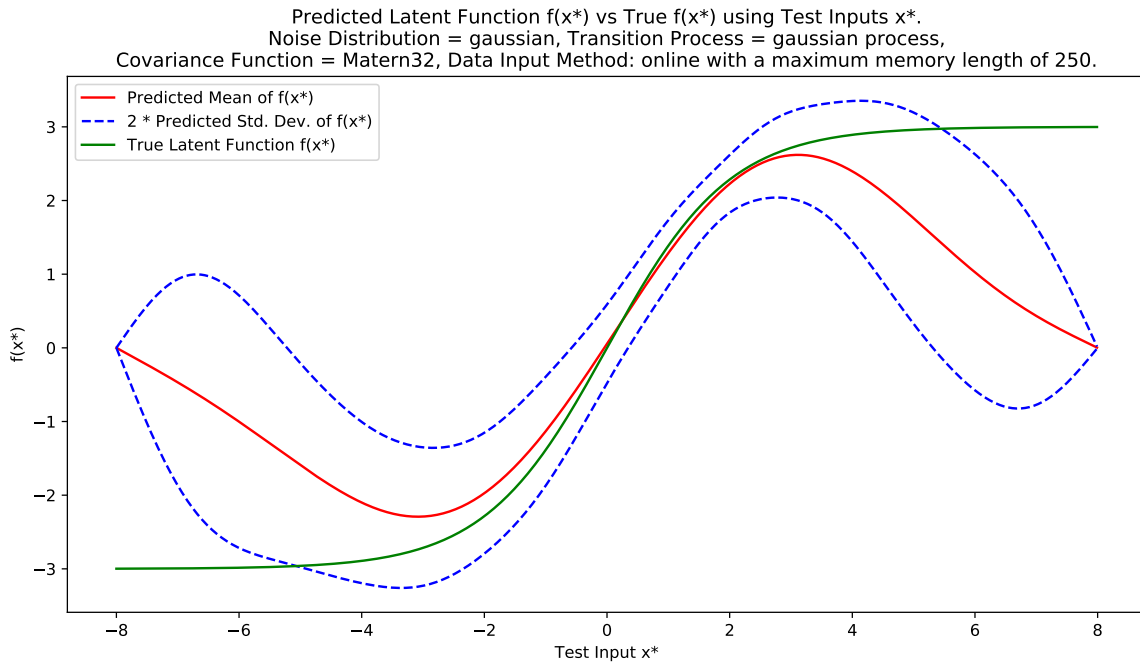
*Test Data:* A sequence of 10000 observations and latent states generated by the above system.

*Additional Information:* The initial point will be zero (known), the observation distribution will be known beforehand (i.e. not learnt), all the reduced-rank GPSSMs will use 8 basis functions, 20 PGAS particles, domain length  $L = 8$ , the  $Q$  prior hyperparameters will both be equal to 1 and we will use a burn-in of 50 samples. The offline models will have  $K = 180$  sampling rounds and the online models will have  $K = 200$  sampling rounds; furthermore, all the online/forgetful online models will start with two data points and receive an additional 10 data points every other sampling round (i.e. number of update samples = 2). Finally, all forgetful models will have a maximum memory length of 250 data points (see section 3.2.2).

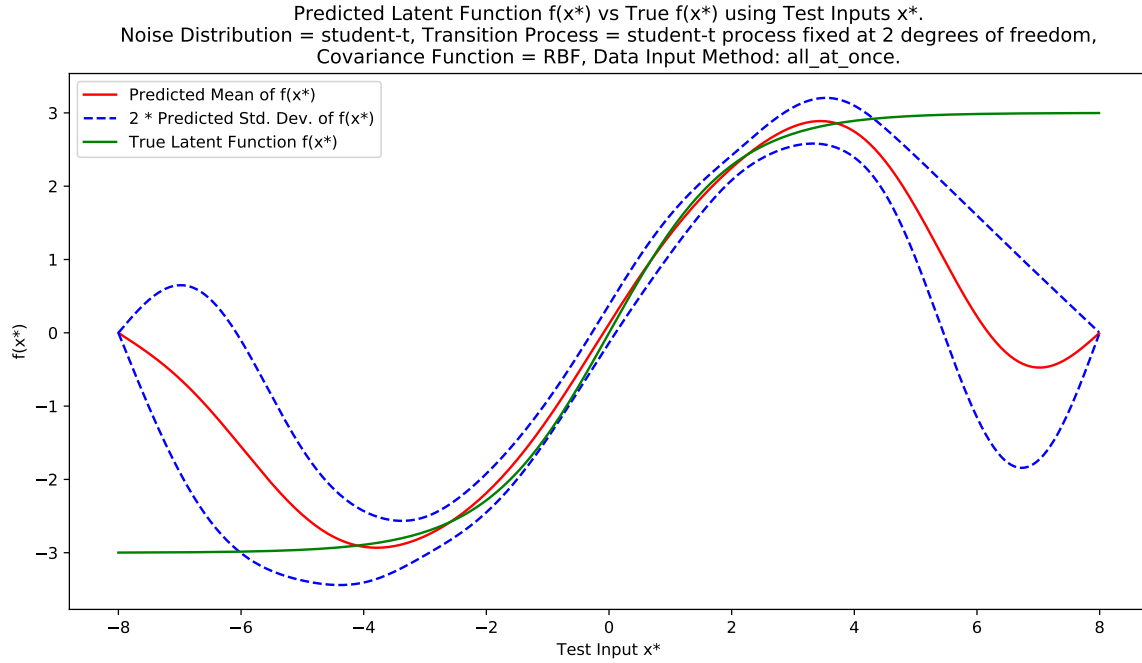
Results:

Method	Kernel	Assumed Latent Noise	Process	RMSE	LL	Training Time (seconds)
Offline	Adaptive	Gaussian	Gaussian	1.16	-1.59	143.38
Online	Matern52	Gaussian	Gaussian	1.25	-1.70	92.62
Forgetful	Matern32	Gaussian	Gaussian	1.16	<b>-1.57</b>	<b>43.80</b>
Offline	RBF	Student-t	Gaussian	<b>1.11</b>	-1.59	177.96
Offline	RBF	Student-t	Student-t	1.12	-1.59	178.55

**Table 3.6:** Results for Example 3.



**Figure 3.12:** Forgetful online learning with a Matern32 covariance function (Table 3.6, Row 3)



**Figure 3.13:** An example of learning with a Student-t process  
 (Table 3.6, Row 5)

#### Discussion:

From example 1, we can see that all the different models will produce similar results in a simple system: this is as expected. Furthermore, in example 2, we see that our novel algorithm for learning in the presence of latent Student-t noise outperforms the algorithm of Svensson et al. [2016]. This is not surprising because Svensson et al. [2016] assume the latent noise is Gaussian. Overall, we see that as the amount of data increases, online learning and forgetful online learning tend to similar results as the offline learning case. Sometimes these methods perform better than offline learning and we believe this might be because the weighting system can lead to bad samples (which are often the early samples even after burn-in) having a small weight. Moreover, forgetful online learning, in particular, can ‘forget’ outliers which would usually have a detrimental effect on the sample quality.

### 3.3 Learning Hilbert reduced-rank Gaussian Processes in High Dimensions

Consider the Gaussian process regression problem:

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon_i \quad (3.65)$$

$$f \sim \mathcal{GP}(\mathbf{0}, k) \quad (3.66)$$

$$\epsilon_i \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}) \quad (3.67)$$

where  $\mathbf{y}_i \in \mathbb{R}^{d_{out}}$ ,  $\mathbf{x}_i \in \mathbb{R}^{d_{in}}$ ,  $\mathbf{Q} \in \mathbb{R}^{d_{out} \times d_{out}}$  is the noise covariance and we have training data  $\{(\mathbf{y}_i, \mathbf{x}_i) \text{ for } i = 1, \dots, N\}$ .

When we talk about dimensionality in the context of Gaussian processes, we are either referring to the input dimension ( $\mathbb{R}^{d_{in}}$ ) or the output dimension ( $\mathbb{R}^{d_{out}}$ ). Usually, high input dimensions are not mentioned that much because the training time of a full Gaussian process largely depends on the amount of training data ( $N$ ) and not the dimension of the input data ( $\mathbb{R}^{d_{in}}$ ). However, it could be said that as the input dimension increases, more training data is required to produce a good model. When we use the Hilbert reduced-rank approximation of the Gaussian process, the picture is different. As we have noted previously: Hilbert reduced-rank Gaussian processes do not scale well when the input dimension increases. For completeness, neither the full Gaussian process nor the Hilbert reduced-rank Gaussian process have much issue with high output dimensions ( $\mathbb{R}^{d_{out}}$ ) because we can just train independent GPs for each output dimension and although this is an approximation compared with a full multi-output GP, it is generally a very good approximation. We will now discuss in depth why the Hilbert reduced-rank Gaussian process model of Solin and Särkkä [2014] does not scale well as the input dimension increases and the approach we take to deal with this.

Consider the above regression problem but take  $d_{out} = 1$  and  $d_{in} = d$ . In this case, the Hilbert reduced-rank model of Solin and Särkkä [2014] with domain  $\Omega = [-L_1, L_1] \times \dots \times [-L_d, L_d]$  plus Dirichlet boundary conditions (zero on the boundary) would approximate the kernel as:

$$k_{\mathbf{m}}(\mathbf{x}, \mathbf{x}') \approx \sum_{j_1, \dots, j_d=1}^m S(\sqrt{\lambda_{j_1, \dots, j_d}}) \phi_{j_1, \dots, j_d}(\mathbf{x}) \phi_{j_1, \dots, j_d}(\mathbf{x}') \quad (3.68)$$

where  $\mathbf{1} \in \mathbb{R}^d$  is a vector of all ones and  $\mathbf{m} = (m_1, \dots, m_d) \in \mathbb{R}^d$  is a vector giving the chosen number of basis functions in each of the  $d$  input dimensions. As before we have eigenvalues:

$$\lambda_{j_1, \dots, j_d} = \sum_{k=1}^d \left( \frac{\pi j_k}{2L_k} \right)^2 \quad (3.69)$$

and the eigenfunctions/basis functions:

$$\phi_{j_1, \dots, j_d}(\mathbf{x}) = \prod_{k=1}^d \frac{1}{\sqrt{L_k}} \sin \left( \frac{\pi j_k (x_k + L_k)}{L_k} \right) \quad (3.70)$$

For clarification, when we refer to the *number of basis functions in dimension  $i$* , we mean  $m_i$ ; furthermore, when we refer to the *total number of basis functions*, we mean the total number of  $\phi_{j_1, \dots, j_d}(\mathbf{x})$  in the sum (3.68) which is  $\prod_{i=1}^d m_i$ .

As the dimension of the input space  $d$  increases, the number of terms in the kernel approximation sum (3.68) is  $\prod_{i=1}^d m_i$  and if the number of basis functions is the same (say  $m$ ) in all dimensions then this simplifies to  $m^d$ . The chosen number of basis functions per dimension (the  $m_i$ ) can be small but not too small: we do need at least 8 basis functions per dimension (i.e.  $m_i \geq 8$  for  $i = 1, \dots, d$ ) for a good performance. To emphasise the scaling issue, if we have 5 input dimensions and 8 basis

functions for each of the input dimensions (i.e.  $m_i = 8$  for  $i = 1, \dots, 5$ ), then the total number of basis functions is  $8^5 = 32768$ . This means not only do we have to sum up  $8^5 = 32768$  terms every time we need to compute the kernel but we also have to invert a  $\prod_{i=1}^d m_i \times \prod_{i=1}^d m_i$  (here equal  $32768 \times 32768$ ) basis function matrix during training/prediction. Clearly this becomes unreasonable to use for input dimensions higher than 2 or 3.

We will now present a novel reduced-rank Gaussian process model that scales considerably better with high input dimensions than the model by Solin and Särkkä [2014]. This new model is based on a combination of neural networks and Gaussian processes inspired by the mGP of Calandra et al. [2016] and the deep Gaussian processes of Damianou and Lawrence [2013]. Previous methods involving deep Gaussian process, such as Damianou and Lawrence [2013], use GPs for the activation functions but the method here is different and, in fact, we fix the activation functions and refer to them as *kernel augmentation functions*.

Consider the regression problem in (3.65)–(3.67) with  $d_{out} = 1$  and  $d_{in} = d$ . As discussed previously in (2.61), under the Hilbert reduced-rank Gaussian process model, this can be written as Svensson et al. [2016]:

$$y_i = \mathbf{A}\phi(\mathbf{x}_i) + \epsilon_i \quad (3.71)$$

$$\mathbf{A} \sim \mathcal{MN}(\mathbf{0}, \mathbf{Q}, \mathbf{V}) \quad (3.72)$$

$$\epsilon_i \sim \mathcal{N}(0, \mathbf{Q}) \quad (3.73)$$

where  $\mathbf{A} \in \mathbb{R}^{1 \times d}$ ,  $\mathbf{Q} \in \mathbb{R}^{1 \times 1}$  is the noise covariance,  $\phi(\mathbf{x}_i) = (\phi_1(\mathbf{x}_i), \dots, \phi_{\mathcal{M}(\mathbf{m})}(\mathbf{x}_i))^T$  for basis functions  $\phi_i$  and  $\mathbf{m} = (m_1, \dots, m_d)$ ,  $\mathbf{V} = \text{diag}(S^{-1}(\sqrt{\lambda_1}), \dots, S^{-1}(\sqrt{\lambda_{\mathcal{M}(\mathbf{m})}}))$  where the  $\lambda_i$  are defined as in (3.69),  $S$  is the spectral density of the chosen kernel function (which is assumed to be stationary) and we have used the mapping  $\mathcal{M}$  from (2.51). Of course, this could be simplified since  $y_i \in \mathbb{R}$  but we keep it in this format just to keep continuity between the previous discussion and the model we are about to present. The Hilbert reduced-rank GP demonstrates nicely that GP regression is ultimately just Bayesian linear regression with infinitely many features: the Hilbert reduced-rank GP just uses finitely many of those features.

In the same way as the neural networks extend linear regression, our new model which we call the deep Hilbert Gaussian process model will extend the Hilbert reduced-rank Gaussian process model. Consider a basic multilayer perceptron (MLP) with activation function  $\sigma$  and  $n \in \mathbb{Z}_{\geq 0}$  hidden layers. In the case of regression, for a given input  $\mathbf{x}$  the output  $y$  is constructed as:

$$z_1 = \sigma(\mathbf{W}_0 \mathbf{x} + \mathbf{b}_0) \quad (3.74)$$

$$z_{i+1} = \sigma(\mathbf{W}_i z_i + \mathbf{b}_i) \text{ for } i = 1, \dots, n-1 \quad (3.75)$$

$$y = \mathbf{W}_n z_n + \mathbf{b}_n \quad (3.76)$$

Important things to note are that we have a linear final activation and that the activation function  $\sigma$  is applied to each element of the input vector (i.e. is vectorised).

Also, the backpropagation algorithm is used to find the weights  $\mathbf{W}_i$  and biases  $\mathbf{b}_i$  for  $i = 0, \dots, n$ . As with all standard neural networks, the MLP is prone to overfitting; although, there are a variety of methods (e.g.  $l_1, l_2$  regularisation and dropout) designed to reduce it. Furthermore, we only get point estimates of the predicted outputs but having an understanding of the confidence in our estimate can be very useful so this is a big drawback of the MLP. The deep Hilbert Gaussian process model (DHGPM) combines the MLP with reduced-rank GPs, and this reduces the risk of overfitting while providing both point estimates and confidence regions (which take into account all sources of uncertainty) for the model's predictions.

The DHGPM arises from breaking the neural net into more pieces:

$$v_i = \mathbf{W}_i \mathbf{z}_i + \mathbf{b}_i \quad (3.77)$$

$$z_{i+1} = \sigma(v_i) \quad (3.78)$$

hence we can write:

$$v_i = \mathbf{W}_i \sigma(v_{i-1}) + \mathbf{b}_i \quad (3.79)$$

Therefore, the MLP with activation  $\phi$ ,  $n \in \mathbb{Z}_{\geq 0}$  hidden layers and given input  $\mathbf{x}$  could have the output  $y$  constructed as:

$$\mathbf{v}_0 = \mathbf{A}_0 \mathbf{x} + \mathbf{b}_0 \quad (3.80)$$

$$\mathbf{v}_i = \mathbf{A}_i \phi(\mathbf{v}_{i-1}) + \mathbf{b}_i \text{ for } i = 1, \dots, n-1 \quad (3.81)$$

$$y = \mathbf{A}_n \phi(\mathbf{v}_{n-1}) + \mathbf{b}_n \quad (3.82)$$

for weights  $\mathbf{A}_i$  and biases  $\mathbf{b}_i$ . If we imagine that the  $\phi$  in (3.81)–(3.82) is the same the  $\phi$  in (3.70), then this is starting to look very much like layers of Hilbert reduced-rank Gaussian processes but without any prior distributions. Hence, to construct the deep Hilbert Gaussian process model (DHGPM) we add these priors.

### 3.3.1 The Deep Hilbert Gaussian Process Model (DHGPM)

Suppose we have training data  $(y_k, \mathbf{x}_k)$  for  $k = 1, \dots, N$  where  $y_k \in \mathbb{R}$  and  $\mathbf{x}_k \in \mathbb{R}^d$ , then for a particular input  $\mathbf{x}_k \in \mathbb{R}^d$  and kernel augmentation function  $\phi$ , the DHGPM with  $n \in \mathbb{Z}_{\geq 0}$  hidden layers is constructed as follows:

$$\mathbf{v}_0^k = \mathbf{A}_0 \mathbf{x}_k + \mathbf{b}_0 \quad (3.83)$$

$$\mathbf{v}_j^k = \mathbf{A}_j \phi(\mathbf{v}_{j-1}^k) + \mathbf{b}_j \text{ for } j = 1, \dots, n-1 \quad (3.84)$$

$$f(\mathbf{x}_k) = \mathbf{A}_n \phi(\mathbf{v}_{n-1}^k) \quad (3.85)$$

$$y_k = f(\mathbf{x}_k) + \epsilon_k \quad (3.86)$$

$$\epsilon_k \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_n) \quad (3.87)$$

$$\mathbf{A}_j | \mathbf{Q}_j \sim \mathcal{MN}(\mathbf{0}, \mathbf{Q}_j, \mathbf{V}_j) \text{ for } j = 0, \dots, n \quad (3.88)$$

$$\mathbf{Q}_j \sim \mathcal{IW}(r_j + l, \mathbf{\Lambda}) \text{ for } j = 0, \dots, n \quad (3.89)$$

$$\mathbf{V}_j = \text{diag}(S_j^{-1}(\sqrt{\lambda_1^j}), \dots, S_j^{-1}(\sqrt{\lambda_{c_j}^j})) \text{ for } j = 0, \dots, n \quad (3.90)$$

$$\lambda_s^j = \left( \frac{\pi s}{2L_j} \right)^2 \text{ for } j = 0, \dots, n \quad (3.91)$$

As  $k$  goes from 1 to  $N$ , the above set-up gives us the input-output pairs  $(y_k, \mathbf{x}_k)$  for  $k = 1, \dots, N$ . The various variables and constants are discussed below.

#### Notes:

1. Inputs are assumed to be column vectors. This is important as it means the given matrix normal distribution is the correct one.
2. The kernel augmentation function  $\phi$  is vectorised; for example, if we have a vector  $\mathbf{v} = (v_1, \dots, v_m)^T \in \mathbb{R}^m$  then  $\phi(\mathbf{v}) = (\phi(v_1), \dots, \phi(v_m))^T \in \mathbb{R}^m$ .
3. Since  $y_k \in \mathbb{R}$ , (3.87) is just a univariate normal distribution with mean 0 and variance  $\mathbf{Q}_n \in \mathbb{R}$ .
4.  $\mathcal{MN}(\mathbf{M}, \mathbf{Q}, \mathbf{V})$  is the matrixvariate normal distribution with mean  $\mathbf{M}$ , row covariance  $\mathbf{Q}$  and column covariance  $\mathbf{V}$ .
5.  $\mathcal{IW}(\nu, \mathbf{\Lambda})$  is the inverse Wishart distribution with degrees of freedom  $\nu$  and scale  $\mathbf{\Lambda}$ . Although  $\mathbf{Q}_n \in \mathbb{R}$  and so the inverse Wishart distribution simplifies to an inverse Gamma distribution, the remaining  $\mathbf{Q}_i$  for  $i = 0, \dots, n-1$  will be matrices in general.
6. For simplicity, the same kernel should be used for all layers. This chosen kernel appears in its spectral density form  $S_j$  in each layer  $j$ . Although the form of the kernel should be the same for all layers; for example, all layers could have an RBF kernel, these kernels are not constrained to have the same hyperparameters. In other words, for each layer  $j$ , if the chosen kernel  $S_j$  has hyperparameters  $\theta_j$ , then  $\theta_j$  is not constrained to equal  $\theta_i$  for different layers  $i$  and  $j$  (of course within a layer the  $\theta_j$  must be the same for all occurrences of  $S_j$ ).

7. The outputs  $y_k$  for  $k = 1, \dots, N$  should be single-dimensional i.e.  $y_k \in \mathbb{R}$ . For multiple outputs, fit separate DHGPMS for each output dimension.
8. A layout  $[d_0, \dots, d_{n-1}]$  gives the dimensions of  $\mathbf{v}_0, \dots, \mathbf{v}_{n-1}$ . So when we refer to DHGPM with layout  $[d_0, \dots, d_{n-1}]$  we mean a model with  $n$  hidden layers with input  $\mathbf{x}_k \in \mathbb{R}^d$  mapped to  $v_0 \in \mathbb{R}^{d_0}$ ,  $v_0$  mapped to  $v_1 \in \mathbb{R}^{d_1}$ , ...,  $v_{n-2}$  mapped to  $v_{n-1} \in \mathbb{R}^{d_{n-1}}$  and then  $v_{n-1}$  mapped to  $y_k \in \mathbb{R}$ . This gives us the dimensions of the  $\mathbf{A}_j$  and  $\mathbf{b}_j$  for  $j = 0, \dots, n$ ; for example, suppose with have a DHGPM with layout  $[3, 12]$  with input dimension 4, then this is a DHGPM with 2 hidden layers and  $\mathbf{A}_0 \in \mathbb{R}^{3 \times 4}$ ,  $\mathbf{A}_1 \in \mathbb{R}^{12 \times 3}$ ,  $\mathbf{A}_2 \in \mathbb{R}^{1 \times 12}$ ,  $\mathbf{b}_0 \in \mathbb{R}^3$  and  $\mathbf{b}_1 \in \mathbb{R}^{12}$  (assuming the input is a column vector).
9.  $\mathbf{Q}_j \in \mathbb{R}^{r_j \times r_j}$  where  $r_j$  is the number of rows in  $\mathbf{A}_j$ . In the inverse Wishart prior (3.89) the degrees of freedom depends on this parameter  $r_j$ , which equals the number of columns (or rows) of  $\mathbf{Q}_j$ .
10. For simplicity, either we take the  $\mathbf{Q}_j$  to be diagonal matrices or we take  $\mathbf{Q}_j = \sigma_j \mathbf{I}$  where  $\sigma_j \in \mathbb{R}_{>0}$  and  $\mathbf{I}$  is the identity matrix or we take the  $\mathbf{Q}_i$  to be of the form  $\mathbf{Q}_j = \mathbf{p}_j \mathbf{p}_j^T + \sigma_j \mathbf{I}$  where  $\sigma_j \in \mathbb{R}_{>0}$  and  $\mathbf{p} \in \mathbb{R}^{r_j}$ .
11.  $\mathbf{V}_j \in \mathbb{R}^{c_j \times c_j}$  where  $c_j$  is the number of columns in  $\mathbf{A}_j$ . In the definition of  $\mathbf{V}_j$  (3.90) there is a parameter  $c_j$  and this  $c_j$  equals the number of columns (or rows) of  $\mathbf{V}_j$ .
12. The prior parameters  $l$  and  $\Lambda$  for the Inverse Wishart prior (3.89) should be the same for all layers and fixed before training. The degrees of freedom ( $r_j + l$ ) might not be the same for all layers (depending on the  $r_j$ ) but the  $l$  should be the same.
13. The last layer has no bias term: this helps prevent overfitting.
14. The biases have no explicit priors; however, they are somewhat constrained by the priors placed on the  $\mathbf{A}_j$ .
15. There is no explicit noise terms in each of the hidden or inputs layers (unlike the output layer) but the  $\mathbf{Q}_j$  inside the matrix-normal distribution do take potential noise into account which is why we need to use the inverse Wishart prior: it stops the model assuming everything is noise.
16. Unlike in the Hilbert reduced-rank Gaussian process model, the  $L_j$  in  $\lambda_s^j$  (3.91) are not fixed: they are variables for optimisation and sampling.

### Training:

Let  $\boldsymbol{\theta}$  be a vector<sup>1</sup> of all the unknown parameters in the model, this includes all the weights  $\mathbf{A}_j$ , the biases  $\mathbf{b}_j$ , the  $\mathbf{Q}_j$ , the  $L_j$  from (3.91) and the unknown parameters  $\boldsymbol{\theta}_j$  of the spectral densities  $S_j$ . It does not include various fixed parameters like

<sup>1</sup>When we refer to  $\boldsymbol{\theta}$ , we assume that it is an element of  $\mathbb{R}^D$  for some  $D > 0$  i.e. all the matrices have been flattened.



the  $l$  and  $\Lambda$  in (3.89). Given training data  $(\mathbf{y}, \mathbf{X}) = \{(y_k, \mathbf{x}_k) \text{ for } k = 1, \dots, N\}$ , the parameters  $\theta$  are learnt by maximising the posterior  $p(\theta|\mathbf{y}, \mathbf{X})$ . To do this, we note that

$$p(\theta|\mathbf{y}, \mathbf{X}) \propto p(\mathbf{y}|\theta, \mathbf{X})p(\theta|\mathbf{X}) \quad (3.92)$$

$$\propto \left( \prod_{i=1}^N \mathcal{N}(y_i|f(\mathbf{x}_i), \mathbf{Q}_n) \right) \left( \prod_{j=0}^n \mathcal{MN}(\mathbf{A}_j|\mathbf{0}, \mathbf{Q}_j, \mathbf{V}_j) \mathcal{IW}(\mathbf{Q}_j|r_j + l, \Lambda) \right) p(\hat{\theta}) \quad (3.93)$$

where  $n$  is the number of hidden layers and  $p(\hat{\theta})$  is the joint prior of all the model parameters in  $\theta$  except for the  $\mathbf{A}_j$  and the  $\mathbf{Q}_j$ . Finding an analytic formula for the posterior distribution of the parameters is impossible so we learn the posterior via a two step process. Firstly, we find the maximum of the posterior log likelihood via L-BFGS-B and then we use EMCEE to sample from the posterior (recall that we get good results starting the sampler at the maximum). We use L-BFGS-B since gradient based optimisers are faster than non-gradient based optimisers like Nelder-Mead and it is trivial to compute the derivatives of (3.93) (with respect to  $\theta$ ) using automatic differentiation. Furthermore, L-BFGS-B allows for the use of bounds which is particularly important for several parameters. The above is summarised in algorithm 10.

The model we presented above aims to balance the number of parameters to be optimised with the expressiveness of the model while also limiting the number of user chosen tuning parameters. However, there are still some tuning parameters which are set before training and we will discuss them below.

### Tuning Parameters

1. **Covariance (kernel) function:** Due to the requirement to have a spectral density, we are limited to choosing a stationary kernel. However, the neural network mitigates this limitation somewhat and extrapolation outside the training data is better than a standard Gaussian process (with a stationary kernel) particularly when using the adaptive kernel we introduced previously.
2. **Prior hyperparameters for  $Q_i$ :**  $l = 1$  and  $\Lambda = \mathbf{I}$  are reasonably good choices but if it is believed that the data is very noisy, then they could be changed.
3. **Model Layout  $[d_0, \dots, d_{n-1}]$ :** This is the most important tuning parameter because the performance of the DHGPM is heavily dependant on the size of the neural network. The model of Solin and Särkkä [2014] depends most heavily on the input dimension and the number of basis functions; in a sense, the DHGPM depends on similar things both of which manifest themselves in the model layout  $[d_0, \dots, d_{n-1}]$ : the larger the input dimension and the more complex the function, the more the DHGPM will require a bigger number of hidden layers and larger  $d_i$  in  $[d_0, \dots, d_{n-1}]$  in order to produce good results. We recommend starting with the layout as small as possible (i.e. one hidden layer and small  $d_0$ ) and if the results are not good (i.e. underfitting) expand the size of the layout. Also, it is important to remember that for one dimensional inputs a

DHGPM with kernel augmentation function  $\sin(x)$  and layout  $[m]$  has a similar flexibility to the model of Solin and Särkkä [2014] with  $m$  basis functions. We have found good results for DHGPMS with kernel augmentation function  $\sin(x)$  and layouts of the form  $[d, 3d]$  where  $d$  is the input dimension.

4. **Kernel augmentation function (activation function):** This function is similar to the activation function in a neural network and it is fixed before training. We call it a kernel augmentation function because it is able to extend the power of the kernel beyond those functions which are expected to be drawn from it. For example, an RBF kernel augmented with  $\sin(x)$  (i.e.  $\sin(x)$  augmentation function) allows us to find smooth periodic functions. In fact,  $\sin(x)$  is almost always a good choice of augmentation function due to its ability to extrapolate periodic functions and learn non-periodic functions inside a fixed domain (similar to Fourier analysis). In the context of DHGPMS,  $\sin(x)$  is able to learn more complex functions in a much smaller layout than other kernel augmentation functions like  $\tanh(x)$ . However, for extrapolating non-periodic functions outside the data regime a  $\tanh(x)$  or log-sigmoid kernel augmentation function can also work well.

---

**Algorithm 10:** Learning of a DHGPM
 

---

**Input** : Training data  $\mathbf{X}, \mathbf{y}$ .  
**Output:** Samples from parameter posterior  $p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y})$ .

- 1 initialise unknown parameters  $\boldsymbol{\theta}$  randomly;  
 /\* Make sure the random initial parameters are valid; for example, some parameters must be positive. See the *Notes* section. \*/
- 2 maximise (3.93) via L-BFGS-B to give the optimal parameters  $\boldsymbol{\theta}_*$ ;  
 /\* Take note of parameter bounds. Also, using several restarts (with different random initial parameters) can alleviate the issue of getting stuck in a poor local maximum. \*/
- 3 sample from posterior  $p(\boldsymbol{\theta}|\mathbf{X}, \mathbf{y})$  using EMCEE with each walker having initial parameters  $\boldsymbol{\theta}_* +$  small random noise;  
 /\* See comments about noise in algorithm 8. For DHGPMS use at least  $2p + 2$  walkers where  $p$  is the length of the vector  $\boldsymbol{\theta}_*$ . A total of 500 posterior samples with each walker having a burn-in of 60 samples usually gives a good characterisation of the posterior. \*/
- 4 **return** samples from parameter posterior;

---

Let  $\Theta = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_m\}$  be our collection of samples from parameter posterior and define  $f_{\hat{\boldsymbol{\theta}}}(\mathbf{x})$  to be  $f(\mathbf{x})$  from (3.85) but with the parameters of the DHGPM equal to  $\hat{\boldsymbol{\theta}}$ . Given a test input  $\mathbf{x}_*$ , we can define (implicitly conditioned on the training data):

$$\mathbb{E}[y_*|\mathbf{x}_*] = \mathbb{E}[f(\mathbf{x}_*)|\mathbf{x}_*] \approx \frac{1}{m} \sum_{i=1}^m f_{\boldsymbol{\theta}_i}(\mathbf{x}_*) \quad (3.94)$$

$$\mathbb{V}[f(\mathbf{x}_*)|\mathbf{x}_*] \approx \left( \frac{1}{m} \sum_{i=1}^m (f_{\boldsymbol{\theta}_i}(\mathbf{x}_*))^2 \right) - (\mathbb{E}[f(\mathbf{x}_*)|\mathbf{x}_*])^2 \quad (3.95)$$

$$\mathbb{V}[y_* | \mathbf{x}_*] \approx \mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*] + \frac{1}{m} \sum_{i=1}^m \mathbf{Q}_n^{\theta_i} \quad (3.96)$$

where  $\mathbf{Q}_n^{\theta_i}$  is the  $\mathbf{Q}_n$  from (3.87) with DHGPM parameters  $\boldsymbol{\theta} = \boldsymbol{\theta}_i$  and  $y_*$  is the test observation corresponding to test input  $\mathbf{x}_*$ .

### Issues associated with DHGPMS

As with all Gaussian process models, there are situations in which the DHGPM performs poorly. The potential issues which can occur are discussed below as well as potential remedies.

1. **Local Optima:** In the same way as with standard Gaussian process regression and with the Hilbert reduced-rank model of Solin and Särkkä [2014], it is possible for the optimisation procedure to return a local optimum, and this can lead to poor performance. The solution is to restart the optimisation from several random starting points and this usually solves the issue but has the draw back of increasing training time.
2. **Amount of Data:** Generally we need more data for the DHGPM than with the standard GP. In fact, in the low data regime (less than 200 data points) we would recommend using the standard GP rather than the DHGPM.
3. **Overfitting:** Fortunately, overfitting is rare in this model but there are two situations in which overfitting can occur. Firstly, when the number of data points is small compared with the dimensions of  $\mathbf{v}_0, \dots, \mathbf{v}_{n-1}$  and secondly when the training data contains outliers. In the first case, we simply have to reduce the dimensions of  $\mathbf{v}_0, \dots, \mathbf{v}_{n-1}$  and in the second case, we should replace (3.87) with a Student-t distribution (see next section). That said, if it is thought that overfitting has occurred in the DHGPM, it is worth testing whether we are in a local optima since the nature of the model (especially with a  $\sin(x)$  augmentation function) means that local optima have many of the same characteristics as overfitting. Finally, it is worth pointing out that overfitting is not unique to the DHGPM; in fact, overfitting affects all Gaussian process models and even standard Gaussian process regression often overfits in the presence of outliers.
4. **Underfitting:** We recommended having the layout as small as possible (i.e. the dimensions of  $\mathbf{v}_0, \dots, \mathbf{v}_{n-1}$  as small as possible) because this can reduce training times significantly. However, this can lead to underfitting and the solution here is to increase the dimensions of  $\mathbf{v}_0, \dots, \mathbf{v}_{n-1}$ .

### 3.3.2 DHGPMS with non-Gaussian noise and Student-t processes

The DHGPM can be adapted to use different noise distributions or to use Student-t processes: all we have to do is change some prior distributions; for example, if we wish to change the noise distribution, we should change (3.87) to the desired noise

distribution and if we wish to use Student-t processes we should change (3.88) to a matrixvariate-t distribution (3.34). However, changing the model may introduce additional parameters that will need to be included in  $\theta$  if we want them to be optimised/sampled from. The changes to the prior distributions lead to changes in (3.93); for example, using a Student-t process with Gaussian noise leads to (3.93) becoming:

$$\left( \prod_{i=1}^N \mathcal{N}(y_i | f(\mathbf{x}_i), \mathbf{Q}_n) \right) \left( \prod_{j=0}^n \mathcal{MT}_\nu(\mathbf{A}_j | \mathbf{0}, \mathbf{Q}_j, \mathbf{V}_j) \mathcal{IW}(\mathbf{Q}_j | r_j + l, \mathbf{\Lambda}) \right) p(\hat{\theta}) \quad (3.97)$$

where  $\mathcal{MT}_\nu$  is a matrixvariate-t distribution with  $\nu$  degrees of freedom and  $\nu$  is usually fixed before training and not optimised. Another example is a Student-t process with Student-t noise which leads to (3.93) becoming:

$$\left( \prod_{i=1}^T \mathcal{T}(y_i | f(\mathbf{x}_i), \mathbf{Q}_n, \nu_1) \right) \left( \prod_{j=0}^n \mathcal{MT}_{\nu_2}(\mathbf{A}_j | \mathbf{0}, \mathbf{Q}_j, \mathbf{V}_j) \mathcal{IW}(\mathbf{Q}_j | r_j + l, \mathbf{\Lambda}) \right) p(\hat{\theta}) \quad (3.98)$$

where  $\mathcal{T}$  is a multivariate-t distribution. In this case, the  $\nu_1$  in  $\mathcal{T}$  is optimised and so must be added to  $\theta$ . Since we use EMCEE and automatic differentiation these changes to (3.93) do not introduce any additional difficulties; we just follow the same learning algorithm as before while using a suitable replacement to (3.93) and potentially optimising some additional parameters. However, note that (3.96) assumes that the noise distribution is Gaussian and for a general noise distribution (with mean zero) we use that (via iterated variances):

$$\mathbb{V}[y_* | \mathbf{x}_*] = \mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*] + \frac{1}{m} \sum_{i=1}^m \mathbb{V}[p(y_* | \mathbf{x}_*, \theta_i)] \quad (3.99)$$

where as before  $\Theta = \{\theta_1, \dots, \theta_m\}$  is our collection of samples from the parameter posterior and  $\mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*]$  is calculated using (3.94) and (3.95). In the case of Student-t noise, this becomes:

$$\mathbb{V}[y_* | \mathbf{x}_*] = \mathbb{V}[f(\mathbf{x}_*) | \mathbf{x}_*] + \frac{1}{m} \sum_{i=1}^m \frac{\nu^{\theta_i}}{\nu^{\theta_i} - 2} \mathbf{Q}_n^{\theta_i} \quad (3.100)$$

if  $\min_i(\nu^{\theta_i}) > 2$  otherwise it does not exist (is infinite). Here,  $\nu^{\theta_i}$  and  $\mathbf{Q}_n^{\theta_i}$  are respectively the noise degrees of freedom  $\nu$  and the noise scale  $\mathbf{Q}_n$  under parameters  $\theta = \theta_i$ .

### 3.3.3 Examples

In this section, we look at several examples of DHGPMS on a variety of synthetic and non-synthetic datasets. When comparing models, we look at three performance statistics: the root mean square error (RMSE), the mean log likelihood (LL) and the mean training time. Suppose we have a test dataset  $\mathcal{D}^* = \{(\mathbf{x}_i^*, y_i^*) \text{ for } i = 1, \dots, N^*\}$ , then for the full GP and Hilbert reduced-rank GP we define the RMSE and LL as in (3.8) and (3.9) respectively. Further suppose that we have samples  $\theta_1, \dots, \theta_K$  from the parameter posterior, then for the DHGPM we define the RMSE as:

$$RMSE = \sqrt{\frac{1}{K} \sum_{i=1}^{N^*} (y_i^* - \mathbb{E}[f(\mathbf{x}_i^*)])^2} \quad (3.101)$$

where  $\mathbb{E}[f(\mathbf{x}_i^*)]$  is defined in (3.94). If the DHGPM assumes a Gaussian noise distribution, then the mean log likelihood (LL) is defined as:

$$LL_{gauss} = \frac{1}{N^*} \sum_{i=1}^{N^*} \log \mathcal{N}(y_i^* | \mathbb{E}[f(\mathbf{x}_i^*)], \mathbb{E}[f(\mathbf{x}_i^*)] + \mathbb{E}[\mathbf{Q}_n]) \quad (3.102)$$

where  $\mathbb{E}[\mathbf{Q}_n] \approx \frac{1}{K} \sum_{j=1}^K \mathbf{Q}_n^{\theta_j}$  and  $\mathbf{Q}_n^{\theta_j}$  is the value of  $\mathbf{Q}_n$  when the model parameters are equal to  $\theta_j$ . Finally, for models which assume a Student-t noise distribution, we define the mean log likelihood (LL) as:

$$LL_{Student-t} = \frac{1}{KN^*} \sum_{i=1}^{N^*} \sum_{j=1}^K \log \mathcal{T}(y_i^* | f_{\theta_j}(\mathbf{x}_i^*), \mathbf{Q}_n^{\theta_j}, \nu^{\theta_j}) \quad (3.103)$$

where  $\mathcal{T}(x | \mu, \Sigma, \nu)$  is a Student-t density function with mean  $\mu$ , scale  $\Sigma$  and  $\nu$  degrees of freedom (using the parametrisation of (3.25)). The notation  $\mathbf{Q}_n^{\theta_j}$  and  $\nu^{\theta_j}$  simply means the values of  $\mathbf{Q}_n$  and  $\nu$  (noise degrees of freedom), respectively, when the model parameters are equal to  $\theta_j$ .

All models will be given 10 optimisation restarts and the training time will be defined as the mean time for an optimisation restart. In the DHGPM diagrams, we plot both the model mean and the model optimum. The *model optimum* is the value of the model under the parameters  $\theta_*$  found during the optimisation (Algorithm 10, line 2) and the *model mean* is the mean value of the model using the parameter samples and (3.94). Ideally, these two values should be very close and if they are far apart, then this suggests that the sampling was not successful.

#### Example 1:

*Aim:* This is a ‘sanity-check’ example. It compares the performance of the full GP, Hilbert reduced-rank GP and the DHGPM on a simple system. We expect all the models to produce similar results.

*Training Data:* 250 pairs  $(x, y)$  with each  $x$  generated by sampling from  $\mathcal{U}(-3, 3)$  and  $y$  generated using system:

$$y = 2.5e^{-x^2} + w_t \quad (3.104)$$

$$w_t \sim \mathcal{N}(0, 1^2) \quad (3.105)$$

*Test Data:* 500 pairs  $(x, y)$  with  $x$  equally spaced between  $-3$  and  $3$  and  $y$  generated as above.

*Additional Information:*

- The manifold-Hilbert reduced-rank models will use a neural network with tanh activations and 2 hidden layers of 1 and 6 neurons respectively.
- The DHGPMS will use  $\sin(x)$  activation functions.

More information is summarised in table 3.7.

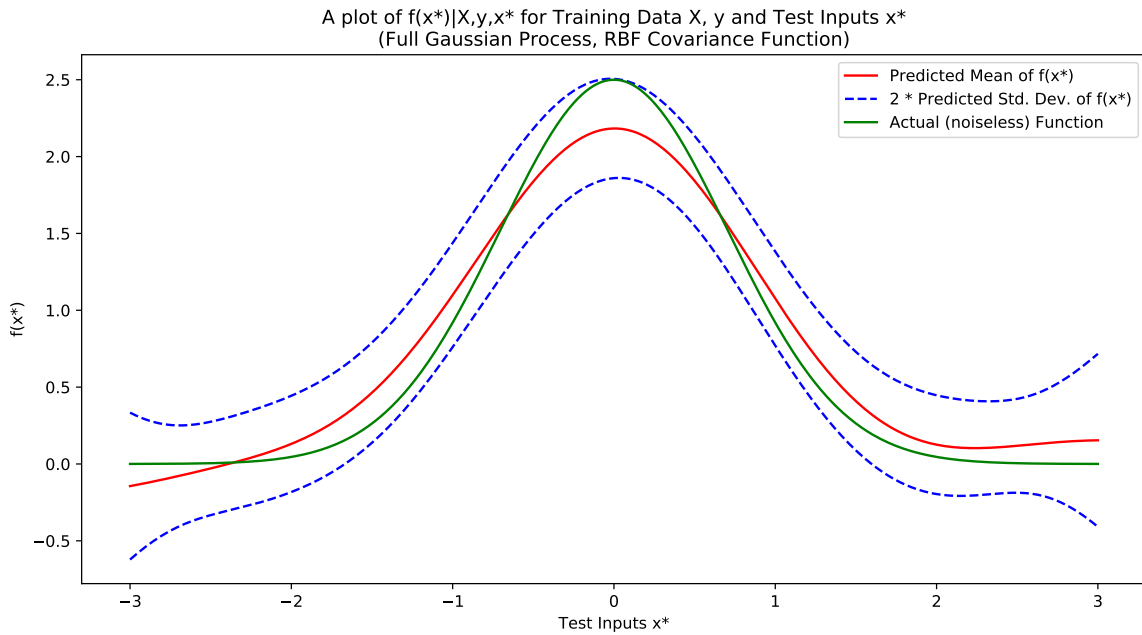
Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP 1	RBF	Gaussian	Gaussian	N/A
Full GP 2	Matern32 + MLP	Gaussian	Gaussian	N/A
Hilbert reduced-rank GP	RBF	Gaussian	Gaussian	12 basis functions, $L = 5$
Manifold-Hilbert reduced-rank GP	Matern32	Gaussian	Gaussian	12 basis functions, $L = 5$
DHGPM 1	Matern52	Gaussian	Student-t	layout = [6]
DHGPM 2	Adaptive	Student-t	Gaussian	layout = [2, 2]

**Table 3.7:** Additional Information for Example 1

*Results:*

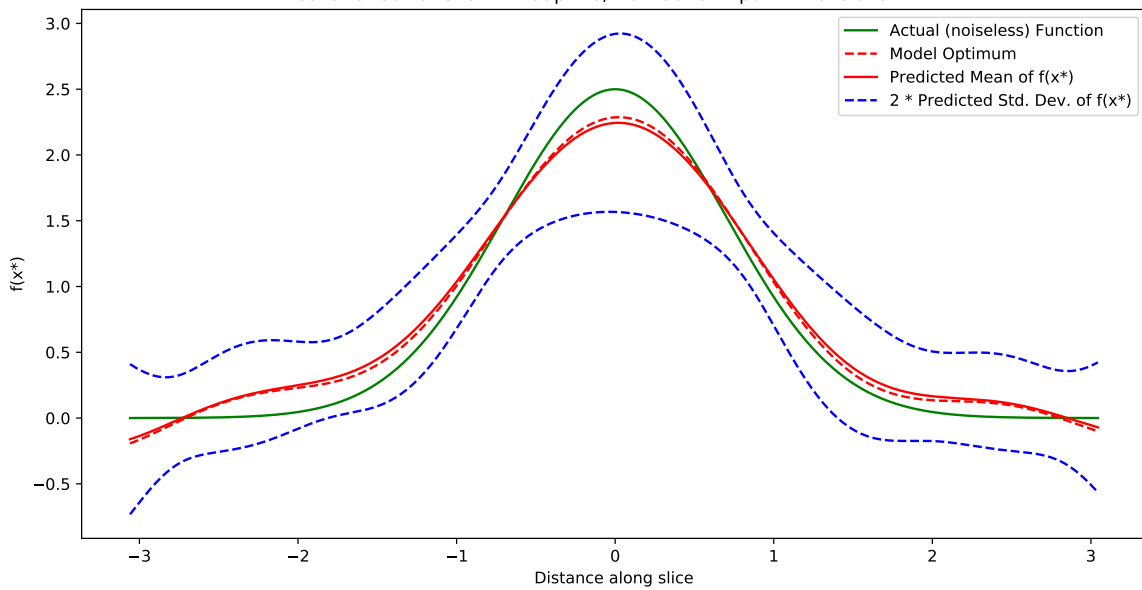
Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	<b>0.99</b>	-1.42	0.14
Full GP 2	<b>0.99</b>	<b>-1.41</b>	0.90
Hilbert reduced-rank GP	<b>0.99</b>	<b>-1.41</b>	<b>0.19</b>
Manifold-Hilbert reduced-rank GP	<b>0.99</b>	-1.42	3.20
DHGPM 1	<b>0.99</b>	-1.44	2.95
DHGPM 2	<b>0.99</b>	<b>-1.41</b>	3.10

**Table 3.8:** Results for Example 1



**Figure 3.14:** A full Gaussian process  
(Table 3.7, Row 1)

Predicted  $f(x^*)$  vs True  $f(x^*)$  using a Random Slice through the Test Inputs  $x^*$ .  
Noise Distribution = gaussian, Function Process = student-t process with 3 degrees of freedom,  
Covariance Function = Adaptive, Number of Input Dimensions = 1.



**Figure 3.15:** A DHGPM with a Student-t process  
(Table 3.7, Row 6)

**Example 2:**

*Aim:* To compare the performance of the full GP, Hilbert reduced-rank GP and the DHGPM on dataset with Student-t noise.

*Training Data:* 250 pairs  $(x, y)$  with each  $x$  generated by sampling from  $\mathcal{U}(-3, 3)$  and  $y$  generated using system:

$$y = 4|\tanh(0.5x + 0.5 \cos(0.8x))| + w_t \quad (3.106)$$

$$w_t \sim \mathcal{N}(0, 1^2, \nu = 1) \quad (3.107)$$

*Test Data:* 2000 pairs  $(x, y)$  with  $x$  equally spaced between  $-3.3$  and  $3.3$  and  $y$  generated as above.

*Additional Information:*

- The manifold-Hilbert reduced-rank models will use a neural network with log-sigmoid activations and 2 hidden layers of 3 and 2 neurons respectively.
- The DHGPMS will use  $\sin(x)$  activation functions.

More information is summarised in table 3.9.

Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP 1	RBF	Gaussian	Gaussian	N/A
Full GP 2	Matern32 + RatQuad	Gaussian	Gaussian	N/A
Hilbert reduced-rank GP	RBF	Gaussian	Gaussian	12 basis functions, $L = 5$
Manifold-Hilbert reduced-rank GP	Matern32	Gaussian	Gaussian	12 basis functions, $L = 5$
DHGPM 1	RBF	Gaussian	Student-t	layout = $[4, 1]$
DHGPM 2	Adaptive	Student-t	Student-t	layout = $[4, 1]$

**Table 3.9:** Additional Information for Example 2

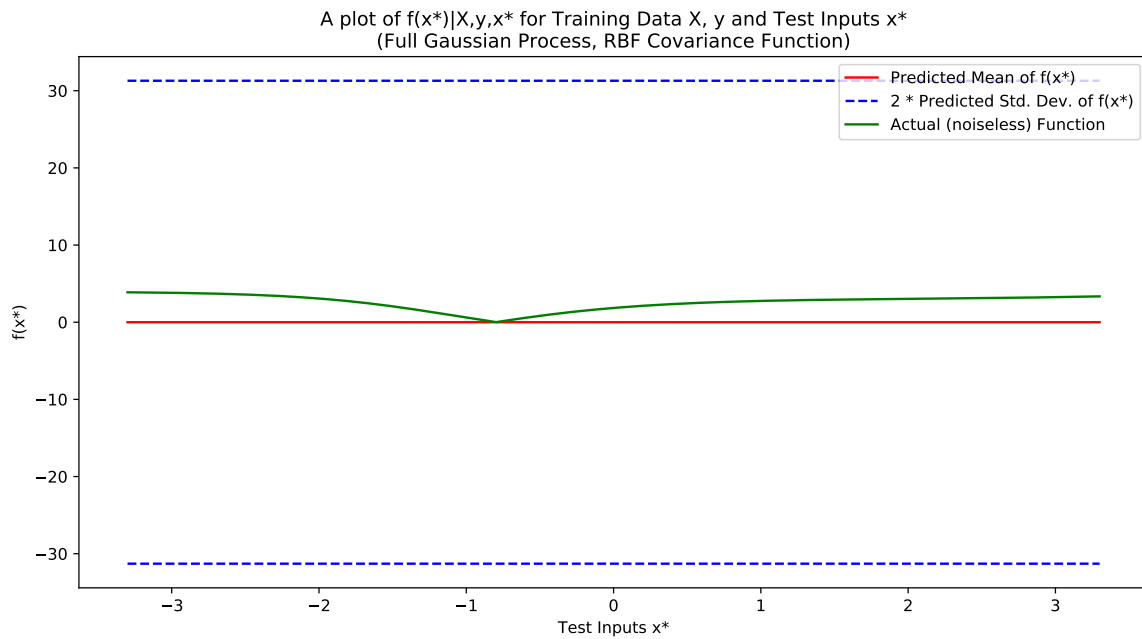


*Results:*

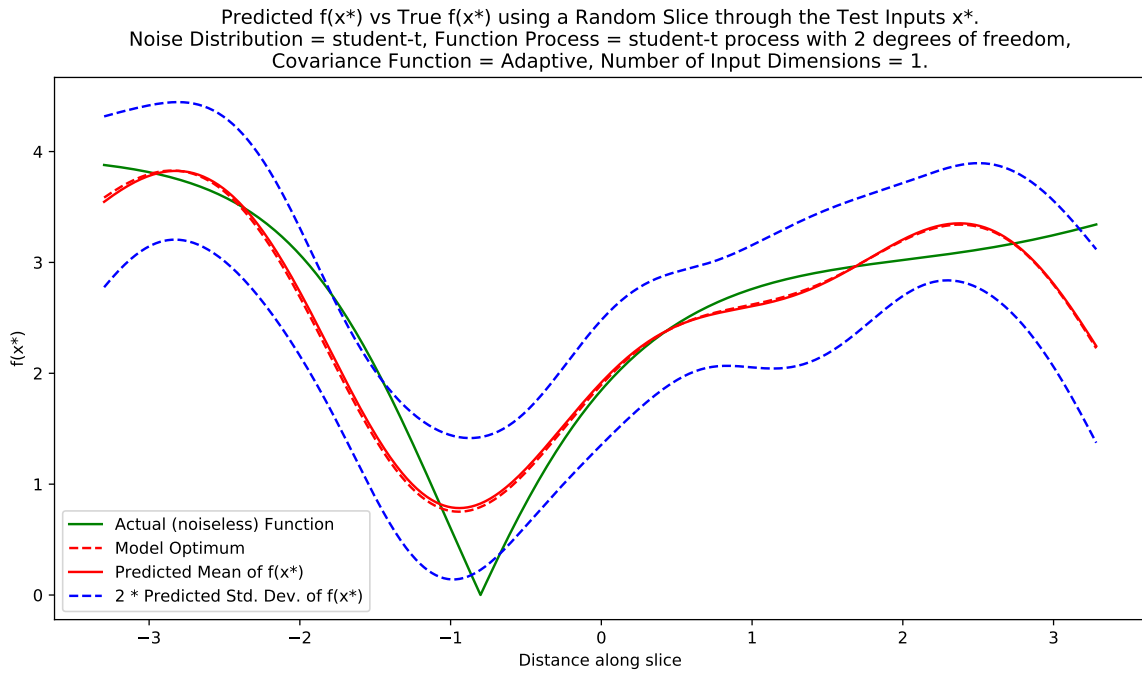
Since the true noise distribution is Student-t, we do not provide a RMSE because it is unhelpful under Student-t noise due to the presence of large outliers.

Model	LL (on test data)	Mean Time (seconds)
Full GP 1	-8.91	0.40
Full GP 2	-9.04	0.92
Hilbert reduced-rank GP	-26.57	<b>0.11</b>
Manifold-Hilbert reduced-rank GP	-27.10	0.72
DHGPM 1	-2.63	3.52
DHGPM 2	<b>-2.62</b>	2.76

**Table 3.10:** Results for Example 2



**Figure 3.16:** A full Gaussian process with the noise assumed to be Gaussian  
(Table 3.9, Row 1)



**Figure 3.17:** A DHGPM with the noise assumed to be Student-t (Table 3.9, Row 6)

Now consider the same example but with 2500 data points for training:

Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP	RBF	Gaussian	Gaussian	N/A
DHGPM	RBF	Gaussian	Student-t	layout = [4, 2]

**Table 3.11:** Additional Information for Example 2 (part 2)

Model	LL (on test data)	Mean Time (seconds)
Full GP	-6.08	101.19
DHGPM	-2.55	10.45

**Table 3.12:** Results for Example 2 (part 2)

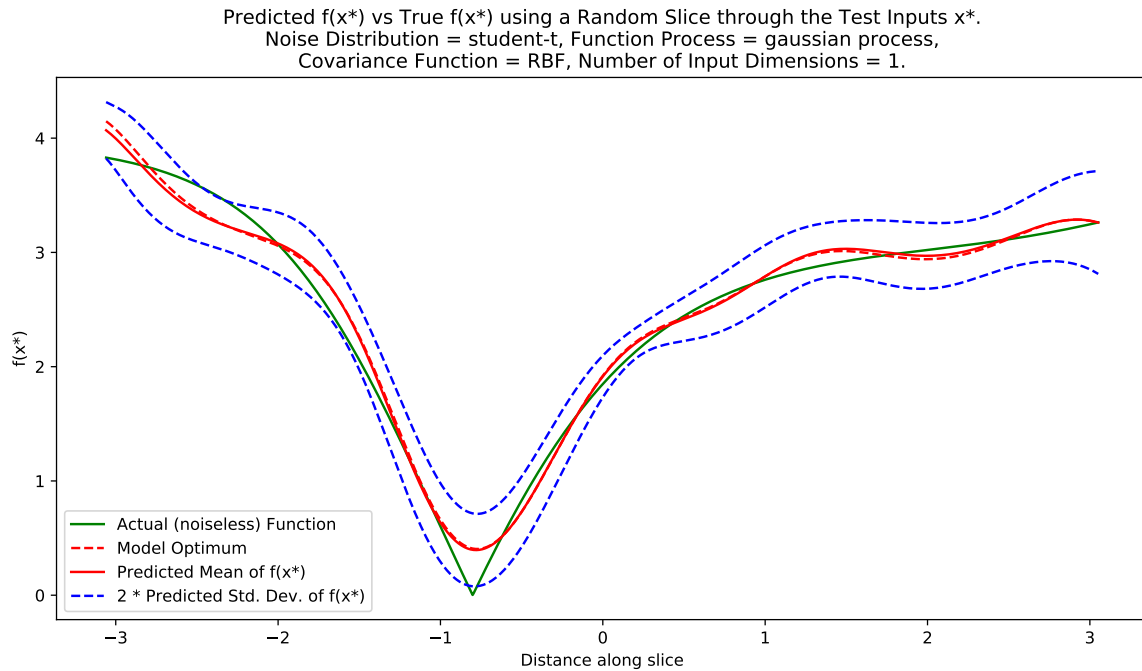


Figure 3.18: A DHGPM (Table 3.11, Row 2)

### Example 3:

*Aim:* To demonstrate that the DHGPM can work well with multidimensional input data and also to show that the DHGPM does not overfit in the presence of very noisy data.

*Training Data:* 1000 pairs  $(x, y)$  with each  $x = (x_1, x_2, x_3)$  generated by sampling  $x_i$  from  $\mathcal{U}(-3, 3)$  for  $i = 1, \dots, 3$  and  $y$  generated using system:

$$y = x_1 + x_2 + x_3 + w_t \quad (3.108)$$

$$w_t \sim \mathcal{N}(0, 5^2) \quad (3.109)$$

*Test Data:* 10000 pairs  $(x, y)$  with each  $x = (x_1, x_2, x_3)$  generated by sampling  $x_i$  from  $\mathcal{U}(-3.6, 3.6)$  for  $i = 1, \dots, 3$  and  $y$  generated as above.

*Additional Information:*

- The DHGPMS will use  $\sin(x)$  activation functions.

More information is summarised in table 3.13.

Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP 1	RBF	Gaussian	Gaussian	N/A
Full GP 2	Matern32 + RatQuad	Gaussian	Gaussian	N/A
DHGPM 1	Adaptive	Gaussian	Gaussian	layout = [2, 1]
DHGPM 2	Adaptive	Student-t	Gaussian	layout = [2, 1]

Table 3.13: Additional Information for Example 3

Results:

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	5.03	-3.03	10.24
Full GP 2	5.03	-3.03	31.96
DHGPM 1	5.06	-3.04	<b>2.22</b>
DHGPM 2	<b>5.02</b>	<b>-3.03</b>	2.91

Table 3.14: Results for Example 3

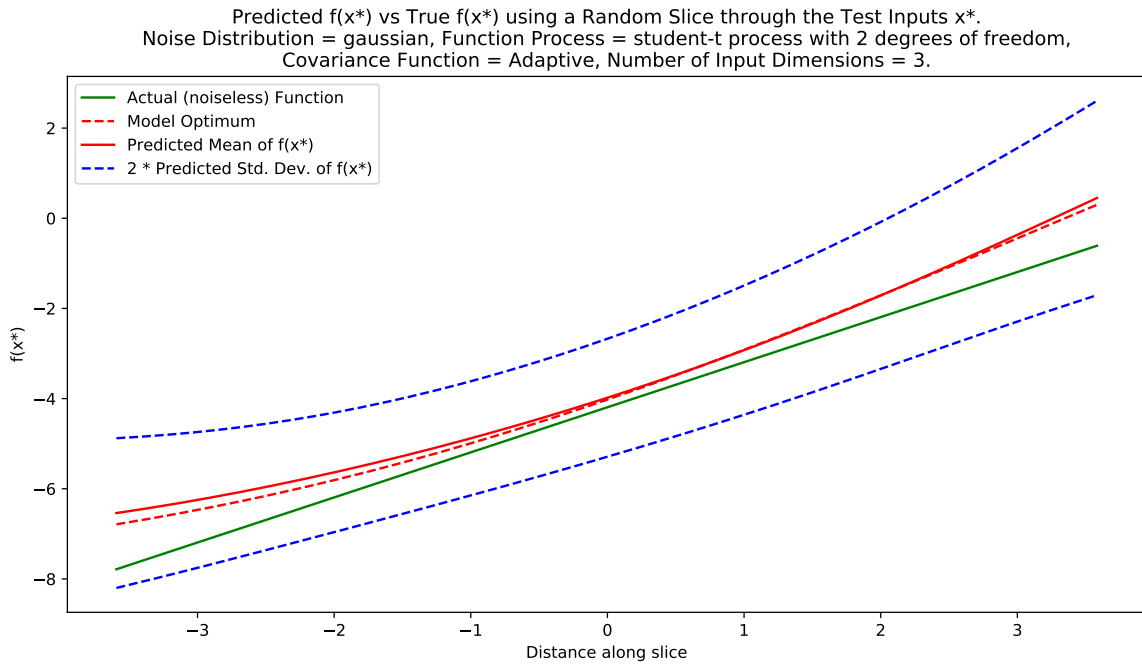


Figure 3.19: A DHGPM with a Student-t process (Table 3.13, Row 4)

**Example 4:**

*Aim:* To demonstrate that the DHGPM can work well with multidimensional input data in low noise environments.

*Training Data:* 1000 pairs  $(\mathbf{x}, y)$  with each  $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5)$  generated by sampling  $x_i$  from  $\mathcal{U}(-3, 3)$  for  $i = 1, \dots, 5$  and  $y$  generated using system:

$$y = |0.1x_1 - 0.2x_2 + 0.3x_3 - 0.4x_4 + 0.5x_5| + w_t \quad (3.110)$$

$$w_t \sim \mathcal{N}(0, 0.1^2) \quad (3.111)$$

*Test Data:* 10000 pairs  $(x, y)$  with each  $\mathbf{x} = (x_1, x_2, x_3, x_4, x_5)$  generated by sampling  $x_i$  from  $\mathcal{U}(-3.6, 3.6)$  for  $i = 1, \dots, 5$  and  $y$  generated as above.

*Additional Information:*

- The DHGPMS will use  $\sin(x)$  activation functions.

More information is summarised in table 3.15.

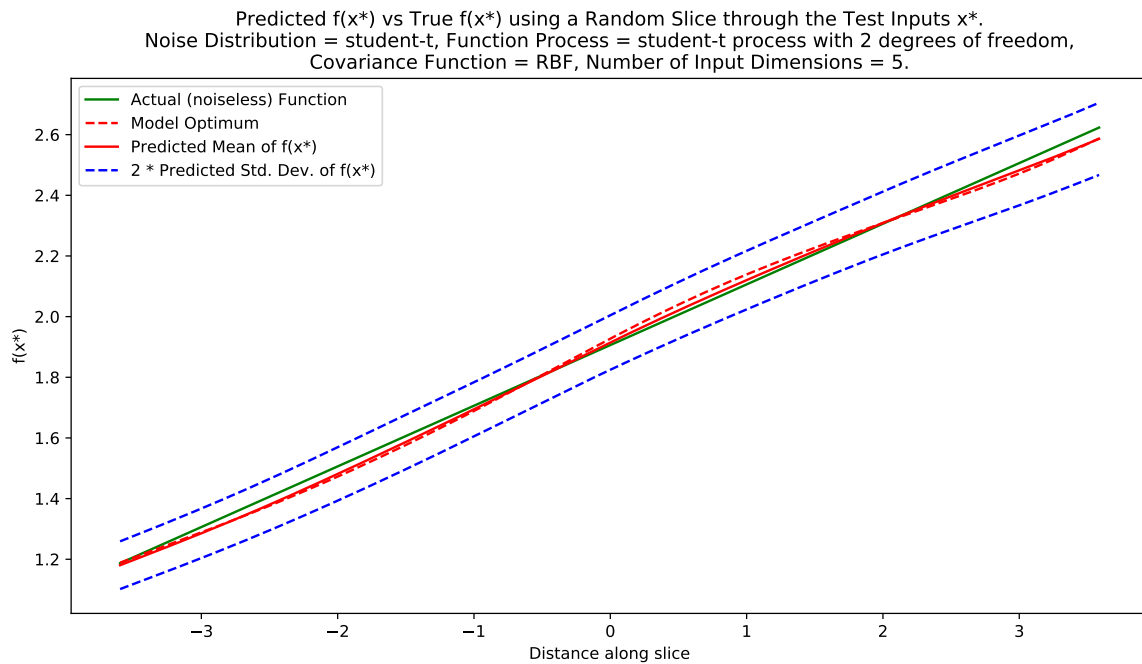
Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP 1	RBF	Gaussian	Gaussian	N/A
Full GP 2	Matern32 + RatQuad	Gaussian	Gaussian	N/A
DHGPM 1	RBF	Gaussian	Gaussian	layout = [10]
DHGPM 2	RBF	Student-t	Student-t	layout = [5, 3]

**Table 3.15:** Additional Information for Example 4

*Results:*

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	0.22	0.20	<b>7.85</b>
Full GP 2	0.17	0.37	34.70
DHGPM 1	0.16	0.31	12.35
DHGPM 2	<b>0.11</b>	<b>0.72</b>	15.18

**Table 3.16:** Results for Example 4



**Figure 3.20:** A DHGPM (Table 3.15, Row 4)

### Example 5:

*Aim:* To demonstrate that the DHGPM can work well on non-synthetic datasets.

#### *The Datasets:*

The first three datasets are from Lichman [2013] (see footnotes for the links to the datasets) and the final dataset was created using data from Chicago Board Options Exchange [2017] and Yahoo Finance [2017]. Associated with each of the datasets is a corresponding task, for the first three datasets this task is described in the provided links and for the final dataset the task is described below. We compare the performance of the DHGPM and the full GP for each of these tasks.

1. The Airfoil Noise Dataset<sup>2</sup>: this dataset contains 1503 datapoints of which we use a randomly selected 1000 points for training and the rest as a test set. The inputs have 5 dimensions and we try to learn a 1 dimensional output.
2. The Combined Cycle Power Plant Dataset<sup>3</sup>: this dataset contains 9568 datapoints of which we use a randomly selected 5000 points for training and the rest as a test set. The inputs have 4 dimensions and we try to learn a 1 dimensional output.
3. The Concrete Compressive Strength Dataset<sup>4</sup>: this dataset contains 1039 datapoints of which we use a randomly selected 800 points for training and the rest as a test set. The inputs have 8 dimensions and we try to learn a 1 dimensional output.

<sup>2</sup><https://archive.ics.uci.edu/ml/datasets/Airfoil+Self-Noise>

<sup>3</sup><https://archive.ics.uci.edu/ml/datasets/Combined+Cycle+Power+Plant>

<sup>4</sup><https://archive.ics.uci.edu/ml/datasets/Concrete+Compressive+Strength>

4. The final dataset consists of a collection of Open, High, Low, Close and Volume data points for the S&P 500<sup>5</sup> and the Open Price of the VIX Index<sup>6</sup>. The aim is to calculate the VIX Open Price from the S&P 500 data (on the same day). We note that VIX is the implied volatility of the S&P 500 and so there is indeed a relation between the two. This dataset contains 3436 datapoints of which we use first 3000 for training and the rest as a test set. The inputs have 5 dimensions and we try to learn a 1 dimensional output.

*Additional Information:*

- The DHGPMS will use  $\sin(x)$  activation functions.

More information is summarised in table 3.17.

Model	Kernel	Process	Assumed Noise Distribution	Additional Information
Full GP 1	RBF	Gaussian	Gaussian	N/A
DHGPM 1	RBF	Gaussian	Gaussian	layout = [5, 8]
DHGPM 2	RBF	Gaussian	Student-t	layout = [8, 24]
DHGPM 3	RBF	Gaussian	Gaussian	layout = [16]
DHGPM 4	Adaptive	Student-t	Student-t	layout = [3, 2]

**Table 3.17:** Additional Information for Example 5

*Results:*

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	2.35	-0.41	5.71
DHGPM 1	<b>2.11</b>	<b>-0.26</b>	35.18

**Table 3.18:** Results for the Airfoil Noise Dataset

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	3.94	0.04	414.94
DHGPM 2	<b>3.88</b>	<b>0.07</b>	<b>257.07</b>

**Table 3.19:** Results for the Combined Cycle Power Plant Dataset

<sup>5</sup><https://finance.yahoo.com/quote/%5EGSPC/>

<sup>6</sup><http://www.cboe.com/products/vix-index-volatility/vix-options-and-futures/vix-index/vix-historical-data>

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	5.31	-0.27	<b>4.92</b>
DHGPM 3	<b>5.01</b>	<b>-0.22</b>	28.78

**Table 3.20:** Results for the Concrete Compressive Strength Dataset

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP 1	10.15	-1.25	145.58
DHGPM 4	<b>4.92</b>	<b>-0.91</b>	<b>24.62</b>

**Table 3.21:** Results for the S&P 500 – VIX Dataset

#### Discussion:

Overall these examples demonstrate that the DHGPM is a very competitive model that can often outperform the full GP. While on small datasets ( $< 1000$  data points) the DHGPM is slower than the full GP, the training time scales well as the dataset size increases and on larger datasets, the DHGPM will be much faster than the full GP. In comparison to the Hilbert reduced-rank GP of Solin and Särkkä [2014], the DHGPM scales significantly better as the input dimension increases.

## 3.4 Learning GPSSMs with High-Dimensional Latent Spaces

### 3.4.1 The Deep Hilbert Gaussian Process State-Space Model

In a similar way to how Svensson et al. [2016] uses the Hilbert reduced-rank GP to create a reduced-rank GPSSM, we will use the DHGPM to construct a novel GPSSM that is particularly useful when we wish to have high-dimensional latent spaces. We will call our new GPSSM the deep Hilbert Gaussian process state-space model (DHGPSSM) and in this section, we will explain the model and provide both a training and prediction algorithm.

The main idea for learning DHGPSSMs is similar to the learning algorithm we created for reduced-rank GPSSMs with Student-t noise: we alternate between sampling the latent states given the model parameters and optimising/sampling the model parameters given the latent states. We will focus on DHGPSSMs with a latent DHGPM but a known observation distribution because there is little reason to have both an observation and transition DHGPM since this increases the computational cost while not providing any additional flexibility. The reason is due to identifiability and the same issue affects all GPSSMs; for more details see Frigola [2015, pp. 32-33].



**Set-up:**

Suppose we have observations  $\mathbf{y}_1, \dots, \mathbf{y}_T$  where  $\mathbf{y}_t \in \mathbb{R}^{n_y}$  for  $t = 1, \dots, T$  and that the latent state dimension is  $n_x$  i.e.  $\mathbf{x}_t \in \mathbb{R}^{n_x} \forall t$ . A deep Hilbert Gaussian process State Space model (DHGPSSM) with  $n \in \mathbb{Z}_{\geq 0}$  hidden layers and kernel augmentation function  $\phi$  is constructed as follows:

$$\mathbf{x}_{t+1}^i = f_i(\mathbf{x}_t) + \epsilon_t^i \text{ for } i = 1, \dots, n_x \quad (3.112)$$

$$\mathbf{x}_{t+1} = (x_{t+1}^1, \dots, x_{t+1}^{n_x}) \in \mathbb{R}^{n_x} \quad (3.113)$$

$$\epsilon_t^i \sim \mathcal{N}(\mathbf{0}, \mathbf{Q}_{n,i}) \text{ for } i = 1, \dots, n_x \quad (3.114)$$

$$\mathbf{v}_{0,i}^t = \mathbf{A}_{0,i} \mathbf{x}_t + \mathbf{b}_{0,i} \text{ for } i = 1, \dots, n_x \quad (3.115)$$

$$\mathbf{v}_{j,i}^t = \mathbf{A}_{j,i} \phi(\mathbf{v}_{j-1,i}^t) + \mathbf{b}_{j,i} \text{ for } j = 1, \dots, n-1 \text{ for } i = 1, \dots, n_x \quad (3.116)$$

$$f_i(\mathbf{x}_t) = \mathbf{A}_{n,i} \phi(\mathbf{v}_{n-1,i}^t) \text{ for } i = 1, \dots, n_x \quad (3.117)$$

$$f(\mathbf{x}_t) = (f_1(\mathbf{x}_t), \dots, f_{n_x}(\mathbf{x}_t)) \in \mathbb{R}^{n_x} \quad (3.118)$$

$$\mathbf{A}_{j,i} | \mathbf{Q}_{j,i} \sim \mathcal{MN}(\mathbf{0}, \mathbf{Q}_{j,i}, \mathbf{V}_{j,i}) \text{ for } j = 0, \dots, n \text{ for } i = 1, \dots, n_x \quad (3.119)$$

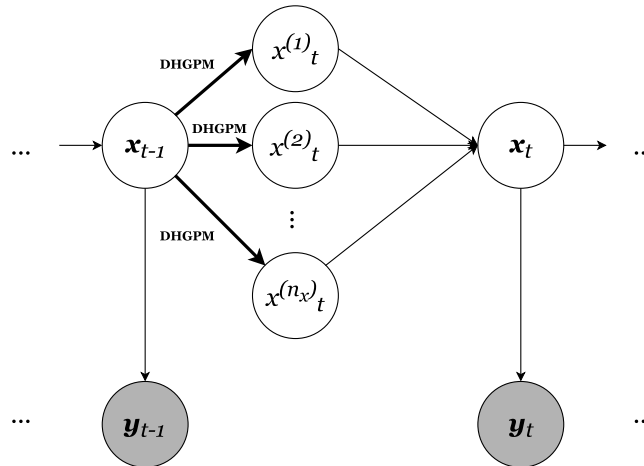
$$\mathbf{Q}_{j,i} \sim \mathcal{IW}(r_{j,i} + l, \mathbf{\Lambda}) \text{ for } j = 0, \dots, n \text{ for } i = 1, \dots, n_x \quad (3.120)$$

$$\mathbf{V}_{j,i} = \text{diag}(S_{j,i}^{-1}(\sqrt{\lambda_1^{j,i}}), \dots, S_{j,i}^{-1}(\sqrt{\lambda_{c_{j,i}}^{j,i}})) \text{ for } j = 0, \dots, n \text{ for } i = 1, \dots, n_x \quad (3.121)$$

$$\lambda_s^{j,i} = \left( \frac{\pi s}{2L_{ji}} \right)^2 \text{ for } j = 0, \dots, n \text{ for } i = 1, \dots, n_x \quad (3.122)$$

$$\mathbf{y}_t \sim p(\mathbf{y}_t | \mathbf{x}_t, \boldsymbol{\theta}_y) \quad (3.123)$$

where equations (3.112) to (3.122) form a latent multi-output DHGPM on  $\mathbf{x}_1, \dots, \mathbf{x}_T$  and we have an observation equation (3.123) linking the latent and observed states. The various variables are discussed in the *Notes* section and a simple graphical model is provided in figure 3.21.



**Figure 3.21:** A DHGPSSM: the model trains independent DHGPMS that map  $\mathbf{x}_{t-1}$  to  $x_t^i$  for  $i = 1, \dots, n_x$  where  $x_t^i$  (written as  $x_t^{(i)}$  in the diagram) is the  $i^{\text{th}}$  dimension of  $\mathbf{x}_t$ .

Although this model looks complex, the idea is conceptually straightforward: learn a latent DHGPM  $f$  that maps  $\mathbf{x}_t$  to  $\mathbf{x}_{t+1}$ . This means that given the latent states

$\mathbf{x}_1, \dots, \mathbf{x}_T$ , we need to find the parameters of the DHGPM using training data  $\mathcal{D} = \{(\mathbf{x}_{t+1}, \mathbf{x}_t) \text{ for } t = 1, \dots, T - 1\}$  where for each  $t$ ,  $\mathbf{x}_t$  is the input of the DHGPM and  $\mathbf{x}_{t+1}$  is the output. The challenge here is that the output  $\mathbf{x}_{t+1}$  is multidimensional in general, so we could have multidimensional inputs and outputs whereas the DHGPM assumes we have a single-dimensional output. However, as we mentioned in the section on DHGPMs, to deal with multiple outputs, we fit independent DHGPMs to each of the output dimensions. For example, suppose that the latent dimension is  $n_x$  and hence  $\mathbf{x}_t = (x_t^1, \dots, x_t^{n_x}) \forall t$ , then all we need to do is to split the training data  $\mathcal{D}$  into  $\mathcal{D}_i = \{(x_{t+1}^i, x_t) \text{ for } t = 1, \dots, T - 1\}$  for  $i = 1, \dots, n_x$  and for each  $\mathcal{D}_i$  we train an independent univariate output DHGPM i.e. we learn a map between  $\mathbf{x}_t$  and  $x_{t+1}^i$  for  $i = 1, \dots, n_x$ .

#### Notes:

1. We will use  $\rho_i$  to refer to all the parameters in (3.112)–(3.123) with an  $i$  subscript; for example, the  $A_{j,i} \forall j$ ,  $\mathbf{b}_{j,i} \forall j$ ,  $\mathbf{Q}_{j,i} \forall j$ ,  $L_{j,i} \forall j$  and the hyperparameters of the spectral density  $S_{j,i}$  for all  $j$  are included in  $\rho_i$ . Therefore, all the parameters  $\theta$  can be written as  $(\rho_1, \dots, \rho_{n_x}, \theta_y)$ .
2. We use independent DHGPMs when dealing with multiple outputs i.e. a multidimensional latent space. This means that  $\rho_{i_1}$  and  $\rho_{i_2}$  are independent for all  $i_1 \neq i_2$  and hence can be optimised/sampled separately.
3. For fixed  $i$ , equations (3.112)–(3.122) form a DHGPM which maps  $\mathbf{x}_t \in \mathbb{R}^{n_x}$  to  $x_{t+1}^i \in \mathbb{R}$ . All the points from the *Notes* part in the section on DHGPMs still apply to each of these independent DHGPMs and explain all the parameters not defined here.
4. For simplicity, we generally use the same tuning parameters for each of the independent latent DHGPMs; for example, for each  $i$  the form covariance function of each DHGPM should be the same; however, there is no such constraint on the covariance hyperparameters since the  $\rho_i$  are independent. Moreover, the layout should be the same for all DHGPMs and so the layout of the  $i^{\text{th}}$  DHGPM  $[d_0^i, \dots, d_{n-1}^i] = [d_0, \dots, d_{n-1}] \forall i$  which leads to  $c_{j,i} = c_j \forall i$  and  $r_{j,i} = r_j \forall i$ . Hence, when we refer to the layout of a DHGPSSM we are referring to  $[d_0, \dots, d_{n-1}]$ , which is the layout for all the independent latent DHGPMs.
5. It is simple to adapt the model to have both an observation and latent DHGPM but we assume the form of the observation distribution is known (3.123) with some unknown parameters  $\theta_y$  which are optimised/sampled.
6. Adapting the model for latent Student-t noise or latent Student-t processes is simple. In a similar way to in DHGPMs, for Student-t noise just replace (3.114) with a Student-t distribution and for a Student-t process replace (3.119) with a matrixvariate-t distribution. Of course, this may introduce some additional parameters which need to be optimised.

7. Due to the independent DHGPMs, the noise (3.114) must be independent for different  $i$ . In particular, this means that we cannot have a multivariate-normal or multivariate-t distribution for  $\epsilon_t = (\epsilon_t^1, \dots, \epsilon_t^{n_x})$ : we require independent univariate-normal or univariate-t distributions for each of the  $\epsilon_t^i$ .
8. Adapting the model for online or forgetful online learning is simple using exactly the same techniques we explained in section 3.2.1.
9. The model is also required to have an initial latent distribution  $p(\mathbf{x}_1)$ ; however, this distribution is of little importance and so we usually take  $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{0}, 100\mathbf{I})$  where  $\mathbf{I}$  is an identity matrix of suitable size.

### Training:

The training process for this model is slightly more complex than in DHGPMs. The idea behind the training algorithm for DHGPSSMs is to balance optimisation/sampling of parameters with learning the latent states: we alternate between these two tasks with our knowledge of the latent states determining how much we can optimise.

One of the key parts of the training process is to optimise the model parameters given the observations and latent states i.e. to maximise  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  with respect to  $\boldsymbol{\theta}$ . Now given observations  $\mathbf{y}_{1:T}$  and latent states  $\mathbf{x}_{1:T}$  we get:

$$p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T}) \quad (3.124)$$

$$\propto p(\mathbf{x}_{1:T}, \mathbf{y}_{1:T}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (3.125)$$

$$\propto p(\mathbf{y}_{1:T}|\mathbf{x}_{1:T}, \boldsymbol{\theta}_y)p(\mathbf{x}_{1:T}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (3.126)$$

since we have independent DHGPMs for each output dimension in the latent states we get:

$$p(\mathbf{y}_{1:T}|\mathbf{x}_{1:T}, \boldsymbol{\theta}_y)p(\mathbf{x}_{1:T}|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (3.127)$$

$$\propto p(\mathbf{y}_{1:T}|\mathbf{x}_{1:T}, \boldsymbol{\theta}_y)p(\mathbf{x}_{1:T}|\boldsymbol{\rho}_1, \dots, \boldsymbol{\rho}_{n_x})p(\boldsymbol{\rho}_1, \dots, \boldsymbol{\rho}_{n_x}, \boldsymbol{\theta}_y) \quad (3.128)$$

$$\propto p(\mathbf{y}_{1:T}|\mathbf{x}_{1:T}, \boldsymbol{\theta}_y)p(\boldsymbol{\theta}_y) \prod_{i=1}^{n_x} p(\mathbf{x}_{1:T}|\boldsymbol{\rho}_i)p(\boldsymbol{\rho}_i) \quad (3.129)$$

It can be seen that we have split the parameters into independent positive parts that are multiplied and hence they can be optimised separately. Now consider  $p(\mathbf{x}_{1:T}|\boldsymbol{\rho}_i)p(\boldsymbol{\rho}_i)$ :

$$p(\mathbf{x}_{1:T}|\boldsymbol{\rho}_i)p(\boldsymbol{\rho}_i) \quad (3.130)$$

$$\propto \left( \prod_{t=1}^{T-1} \mathcal{N}(x_{t+1}^i | f_i(\mathbf{x}_t), \mathbf{Q}_{n,i}) \right) \left( \prod_{j=0}^n \mathcal{MN}(\mathbf{A}_{j,i} | \mathbf{0}, \mathbf{Q}_{j,i}, \mathbf{V}_{j,i}) \mathcal{IW}(\mathbf{Q}_{j,i} | r_{j,i} + l, \boldsymbol{\Lambda}) \right) p(\hat{\boldsymbol{\rho}}_i) \quad (3.131)$$

where  $n$  is the number of hidden layers for the  $i^{\text{th}}$  latent DHGPM (which is the same for all  $i$ ) and  $\hat{\boldsymbol{\rho}}_i$  is all the parts of  $\boldsymbol{\rho}_i$  except for the  $\mathbf{A}_{j,i}$  and the  $\mathbf{Q}_{j,i}$ . As expected, this is the likelihood of a DHGPM with parameters  $\boldsymbol{\rho}_i$  and training data  $\mathcal{D}_i = \{(x_{t+1}^i, \mathbf{x}_t)$

for  $t = 1, \dots, T - 1$ }. We can use this to write the optimisation part of the DHGPSSM training algorithm: see Algorithm 11.

Another key part of the training process is to sample new latent states  $\mathbf{x}_{1:T}^{new}$  given the observations  $\mathbf{y}_{1:T}$ , the previous latent states  $\mathbf{x}_{1:T}^{old}$  and the current model parameters  $\theta$ . We do this using the PGAS sampler of Lindsten et al. [2014, p. 2160] with observations  $\mathbf{y}_{1:T}$ , reference trajectory  $\mathbf{x}_{1:T}^{old}$ , initial samples drawn using a  $\mathcal{N}(\mathbf{0}, 100\mathbf{I})$  distribution and both the transition distribution and observation distribution constructed using the current model parameters  $\theta$  with (3.112) and (3.123) respectively. We will now present the training algorithms and discuss some tuning parameters.

#### List of Algorithms:

1. **Optimisation:** This algorithm maximises (3.124) with respect to  $\theta$  using the independence of DHGPMS to split up the optimisation as we have discussed in (3.129).
2. **OptimisationLoop:** This uses Algorithm 11 to quickly get both the latent states  $\mathbf{x}_{1:T}$  and parameters  $\theta$  near to their respective invariant distributions. The idea is that we alternative between PGAS and Algorithm 11 until (3.124) decreases for the first time and then we stop. Since, at this point (3.124) is no longer just strictly increasing so we must be relatively near to some sort of invariant distribution.
3. **BurnerLoop:** In a similar way to the OptimisationLoop Algorithm, the aim is to get to the invariant distribution but here we add sampling of  $\theta$  which allows us to burn-in EMCEE while getting even closer to the invariant distribution compared with OptimisationLoop.
4. **Learning of a DHGPSSM:** This combines the above algorithms to go from a set of observations to a learnt DHGPSSM.
5. **Prediction:** Given a learnt DHGPSSM and test observations  $\mathbf{y}_1^*, \dots, \mathbf{y}_t^*$ , this algorithm samples from  $\mathbf{y}_{t+k}^*$  for some  $k \in \mathbb{Z}_{>0}$ .

#### Tuning Parameters:

1. **Latent state dimension:** Unless we know the dimension of the latent state, we will need to choose it. In general, the dimension of the latent states depends on the complexity of the process which generates the observations as well as the level of dependency of the future of the process on its history. In other words, if we have a complex system or the system depends on observations which happened many time steps ago, then we will need a correspondingly bigger latent state dimension.
2. **Form of the observational distribution:** This might be known depending on the use case, otherwise set (3.123) to:

$$\mathbf{y}_t = \mathbf{W}\mathbf{x}_t + \mathbf{b} + \mathbf{v}_t \quad (3.132)$$

for suitably sized matrices  $\mathbf{W}$  and  $\mathbf{b}$  and a chosen noise distribution  $\mathbf{v}_t$ . The unknown parameters form  $\theta_y$ .

3. **Number of PGAS particles:** 20 particles is almost always a good choice but lowering it might improve PGAS speed (not by much though).
4. **Tuning parameters for the latent DHGPMS:** See section on tuning parameters for DHGPMS.

### Prediction:

Having looked at how to train the model, the only remaining question is how to use a trained model i.e. prediction. The aim of prediction is: given a learnt model and a collection of ordered test observations  $\mathbf{y}_1^*, \dots, \mathbf{y}_t^*$ , find the predictive distribution of  $\mathbf{y}_{t+k}^*$  for some  $k \in \mathbb{Z}_{>0}$ . Unfortunately for DHGPSSMs, it is not possible to find the exact predictive distribution; however, we can construct a procedure for sampling from the predictive distribution. Furthermore, given samples  $\mathbf{y}_{t+k}^i$  for  $i = 1, \dots, N$ , we can approximate statistics such as the mean and variance of the predictive distribution; for example,

$$\mathbb{E}[\mathbf{y}_{t+k}] \approx \frac{1}{N} \sum_{i=1}^N \mathbf{y}_{t+k}^i \quad (3.133)$$

$$\mathbb{V}[\mathbf{y}_{t+k}] \approx \left( \frac{1}{N} \sum_{i=1}^N (\mathbf{y}_{t+k}^i)^2 \right) - \left( \frac{1}{N} \sum_{i=1}^N \mathbf{y}_{t+k}^i \right)^2 \quad (3.134)$$

The prediction algorithm for DHGPSSMs is summarised in Algorithm 15.

---

#### Algorithm 11: Optimisation

---

**Input** : Observations  $\mathbf{y}_{1:T}$ , Latent States  $\mathbf{x}_{1:T}$ , Initial Point  $\theta_0$ , maxIterations

**Output:** Optimised Parameters  $\theta_*$

/\* All the following optimisations should be done using L-BFGS-B with the maximum iterations equal to maxIterations. Also start the optimisation at the corresponding part of  $\theta_0$ . \*/

1 optimise  $\log p(\mathbf{y}_{1:T} | \mathbf{x}_{1:T}, \theta_y)$  with respect to  $\theta_y$  to get  $\theta_y^*$ ;

2 **for**  $i = 1, \dots, n_x$  **do**

3     optimise  $\log p(\mathbf{x}_{1:T} | \rho_i) + \log p(\rho_i)$  with respect to  $\rho_i$  to get  $\rho_i^*$ ;

   /\* In other words, create a DHGPM with initial parameters

$\rho_{i_0} \subset \theta_0$  and training data  $\mathcal{D}_i = \{(x_{t+1}^i, \mathbf{x}_t) \text{ for } t = 1, \dots, T-1\}$

   then optimise and return  $\rho_i^*$ . \*/

4 **end**

5 **return**  $\theta_* = (\rho_1^*, \dots, \rho_{n_x}^*, \theta_y^*)$ ;

---

---

**Algorithm 12:** OptimisationLoop

---

**Input** : Observations  $\mathbf{y}_{1:T}$ , Current Latent States  $\mathbf{x}_{1:T}$ , Current Parameters  $\theta$   
**Output**: New Latent States  $\mathbf{x}_{1:T}^{new}$ , New Parameters  $\theta_{new}$ , #Loops  $maxJ$

```

1  $bestObjective = \log p(\theta | \mathbf{x}_{1:T}, \mathbf{y}_{1:T});$ 
2 for  $j = 1, 2, \dots$  do
   /* Usually this reaches the break point but it is sensible to
   set a maximum number of iterations. */
3 sample new latent states  $\mathbf{x}_{1:T}^{new}$  given observations  $\mathbf{y}_{1:T}$ , current latent states
 $\mathbf{x}_{1:T}$  and current parameters  $\theta$  via PGAS;
4  $\theta_{new} = \text{Optimisation}(\text{Observations}=\mathbf{y}_{1:T}, \text{Latent States}=\mathbf{x}_{1:T}^{new}, \text{Initial Point}
= \theta, \text{maxIterations}=2j + 5);$ 
5 if  $\log p(\theta_{new} | \mathbf{x}_{1:T}^{new}, \mathbf{y}_{1:T}) < bestObjective$  then
6   |  $maxJ = j;$ 
7   | break;
8 end
9 else
10  |  $bestObjective = \log p(\theta_{new} | \mathbf{x}_{1:T}^{new}, \mathbf{y}_{1:T});$ 
11  |  $\mathbf{x}_{1:T} = \mathbf{x}_{1:T}^{new};$ 
12  |  $\theta = \theta_{new};$ 
13 end
14 end
15 return  $\mathbf{x}_{1:T}^{new}, \theta_{new}, maxJ.$ 

```

---

**Algorithm 13: BurnerLoop**


---

```

Input : Observations  $\mathbf{y}_{1:T}$ , Current Latent States  $\mathbf{x}_{1:T}$ , Current Parameters  $\boldsymbol{\theta}$ ,
          Base Iterations  $maxJ$ 
Output: New Latent States  $\mathbf{x}_{1:T}^{new}$ , New Parameters  $\boldsymbol{\theta}_{new}$ 
1  $bestObjective = \log p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$ ;
2 for  $i = 1, 2, \dots$  do
   | /* Usually this reaches the break point but it is sensible to
   |   set a maximum number of iterations. */
3   sample new latent states  $\mathbf{x}_{1:T}^{new}$  given observations  $\mathbf{y}_{1:T}$ , current latent states
   |  $\mathbf{x}_{1:T}$  and current parameters  $\boldsymbol{\theta}$  via PGAS;
4   if  $i \equiv 0 \pmod{5}$  and  $i > 0$  then
5     |  $\boldsymbol{\theta}_* = \text{Optimisation}(\text{Observations}=\mathbf{y}_{1:T}, \text{Latent States}=\mathbf{x}_{1:T}^{new}, \text{Initial Point}$ 
6     |  $=\boldsymbol{\theta}, \text{maxIterations}=2(maxJ + i) + 5)$ ;
7     | if  $\log p(\boldsymbol{\theta}_*|\mathbf{x}_{1:T}^{new}, \mathbf{y}_{1:T}) < bestObjective$  then
8     |   | break;
9     | end
10    | else
11    |   |  $bestObjective = \log p(\boldsymbol{\theta}_*|\mathbf{x}_{1:T}^{new}, \mathbf{y}_{1:T})$ ;
12    |   |  $\mathbf{x}_{1:T} = \mathbf{x}_{1:T}^{new}$ ;
13    |   |  $\boldsymbol{\theta} = \boldsymbol{\theta}_*$ ;
14    | end
15    | sample parameters  $\boldsymbol{\theta}_{new}$  from  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  using EMCEE with initial
16    | positions  $\boldsymbol{\theta} +$  vector of random noise for each walker and some EMCEE
17    | burn-in steps. We recommend using 10 burn-in steps per walker and 16
18    | walkers;
19    | /* See previous discussions about adding noise to  $\boldsymbol{\theta}$ ; for
20    | example, in Algorithm 8. Also, we can use a small amount
21    | of walkers for the same reasons as in Algorithm 8. */
22  end
23 else
24   |  $\mathbf{x}_{1:T} = \mathbf{x}_{1:T}^{new}$ ;
25   | sample parameters  $\boldsymbol{\theta}$  from  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  using EMCEE, with initial
26   | walker positions equal to the walker positions from the previous loop
27   | and 1 burn-in step per walker and 16 walkers;
28   | /* For the first loop, set the initial walker positions to  $\boldsymbol{\theta}$ 
29   |   + vector of random noise for each walker. */
30 end
31 return  $\mathbf{x}_{1:T}^{new}, \boldsymbol{\theta}_*$ 

```

---

**Algorithm 14:** Learning of a DHGPSSM

---

```

Input : Observations  $\mathbf{y}_{1:T}$ , Number of Samples  $K$ 
Output: Samples from the parameter posterior  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}, \mathbf{y}_{1:T})$  and samples
          from the latent state distribution  $p(\mathbf{x}_{1:T}|\mathbf{y}_{1:T}, \boldsymbol{\theta})$ .
1 initialise unknown parameters  $\boldsymbol{\theta}$  randomly;
  /* Make sure the random initial parameters are valid. */
  /* Optimisation Stage: */
2 initialise latent states  $\mathbf{x}_{1:T}$  by setting each dimension of  $\mathbf{x}_{1:T}$  to a random
  dimension of  $\mathbf{y}_{1:T}$ . Repeat this random allocation of  $\mathbf{x}_{1:T}$  around 10-100 times
  and finally set the initial latent states equal to the random allocation that
  maximised (3.123);
3  $\mathbf{x}_{1:T}, \boldsymbol{\theta}, \text{max}J = \text{OptimisationLoop}(\text{Observations}=\mathbf{y}_{1:T}, \text{Latent States}=\mathbf{x}_{1:T},$ 
   $\text{Current Parameters}=\boldsymbol{\theta})$ ;
  /* Burn-in Stage: */
4  $\mathbf{x}_{1:T}[0], \boldsymbol{\theta}[0] = \text{BurnerLoop}(\text{Observations}=\mathbf{y}_{1:T}, \text{Latent States}=\mathbf{x}_{1:T}, \text{Current}$ 
   $\text{Parameters}=\boldsymbol{\theta}, \text{Base Iterations} = \text{max}J)$ ;
  /* Sampling Stage: */
5 for  $k = 1, \dots, K$  do
6   sample new latent states  $\mathbf{x}_{1:T}[k]$  given observations  $\mathbf{y}_{1:T}$ , latent states
    $\mathbf{x}_{1:T}[k-1]$  and parameters  $\boldsymbol{\theta}[k-1]$  via PGAS;
7   if  $i \equiv 0 \pmod{5}$  then
8      $\boldsymbol{\theta}_* = \text{Optimisation}(\text{Observations}=\mathbf{y}_{1:T}, \text{Latent States}=\mathbf{x}_{1:T}[k], \text{Initial}$ 
      $\text{Point}=\boldsymbol{\theta}[k-1], \text{maxIterations}=25)$ ;
9     sample parameters  $\boldsymbol{\theta}[k]$  from  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}[k], \mathbf{y}_{1:T})$  using EMCEE with initial
     positions  $\boldsymbol{\theta}_* +$  vector of random noise for each walker and some
     EMCEE burn-in steps. We recommend using 10 burn in steps per walker
     and 16 walkers;
     /* See previous discussions about adding noise to  $\boldsymbol{\theta}_*$ ; for
     example, in Algorithm 8. Also, we can use a small amount
     of walkers for the same reasons as in Algorithm 8. */
10  end
11  else
12    sample parameters  $\boldsymbol{\theta}[k]$  from  $p(\boldsymbol{\theta}|\mathbf{x}_{1:T}[k], \mathbf{y}_{1:T})$  using EMCEE with initial
    walker positions equal to the walker positions from the previous loop
    with 1 burn in step per walker and 16 walkers;
13  end
14 end
15 return  $\mathbf{x}_{1:T}[1:K], \boldsymbol{\theta}[1:K]$ 

```

---



**Algorithm 15:** DHGPSSM Prediction

---

**Input** : Samples  $\theta^1, \dots, \theta^K$  from the posterior parameter distribution of a DHGPSSM, test observations  $\mathbf{y}_1^*, \dots, \mathbf{y}_t^*$

**Output:** Samples from  $\mathbf{y}_{t+k}^*$  for some  $k \in \mathbb{Z}_{>0}$

- 1 **for**  $i = 1, \dots, nSamples$  **do**
- 2     select  $\theta^i$  uniformly from  $\theta^1, \dots, \theta^K$ ;
- 3     set parameters of DHGPSSM equal to  $\theta^i$ ;
- 4     find the latent states  $\mathbf{x}_1^*, \dots, \mathbf{x}_t^*$  corresponding to  $\mathbf{y}_1^*, \dots, \mathbf{y}_t^*$  via PGAS;
- 5     **for**  $j = 1, \dots, k$  **do**
- 6         sample  $\mathbf{x}_{t+j}$  from  $p(\mathbf{x}_{t+j} | \mathbf{x}_{t+j-1}, \theta^i)$  (equation (3.112));
- 7     **end**
- 8     sample  $\mathbf{y}_{t+k}$  from  $p(\mathbf{y}_{t+k} | \mathbf{x}_{t+k}, \theta^i)$  (equation (3.123));
- 9 **end**
- 10 **return** all samples of  $\mathbf{y}_{t+k}$ .

---

### 3.4.2 Examples

In this section, we look at several examples of DHGPSSMs and compare them with an autoregressive full GP and the GPSSM by Svensson et al. [2016]. In the examples that look at comparing the latent states, we use the formulas of (3.49), (3.50) and (3.51) to calculate the RMSE, LL (Gaussian) and LL (Student-t) respectively. Due to unidentifiability, in models with unknown observation distributions or multi-dimensional latent states, we are unable to compare the latent states of different models; however, we can look at comparing predictions from the model and we will look at two types of predictions:

1. Prediction of fixed  $k$ : Given a sequence of test observations  $y_1^*, \dots, y_{T^*}^*$ , we use the models to predict  $y_{t+k}^*$  given  $y_{1:t}^*$  for  $t = 1, \dots, T^* - k$ . Note that we use a particle filter (bootstrap filter) to get the states at time  $t$  rather than a smoother because we wish to compare predictions given only the history of the process. We use the formulas of (3.101), (3.102) and (3.103) to compute the statistics RMSE, LL (Gaussian) and LL (Student-t), respectively, with  $y_i^*$  replaced with  $y_{t+k}^*$  and  $f(x_i^*)$  replaced with the model's prediction for  $y_{t+k}^*$ . Note that the assumed model distribution now refers to the assumed observation distribution.
2. Multi-step forecast: Suppose that the model training data is  $y_1, \dots, y_{T_{train}}$  and the model test data is  $y_{T_{train}+1}, \dots, y_{T^*}$ , then given the training data we predict  $y_{T_{train}+k}$  for  $k = 1, \dots, \min(30, T^*)$  using Algorithm 15. We compare the predictions to the true values using (3.101), (3.102) and (3.103).

**Example 1:**

*Aim:* This is a ‘sanity-check’ example comparing the performance of the DHGPSSM with the model of Svensson et al. [2016] on a simple system. We expect all models to produce similar results. The set-up is the same as in section 3.2.6, example 1 and for clarification, we have the same training and test data as in the original example.

*Additional Information:*

- The model called GPSSM refers to the model of Svensson et al. [2016] with none of our new features added.
- The model called DHGPSSM is a deep Hilbert Gaussian process state-space model.
- The latent DHGPMs in the DHGPSSM will use  $\sin(x)$  activation functions.

More information is summarised in table 3.22.

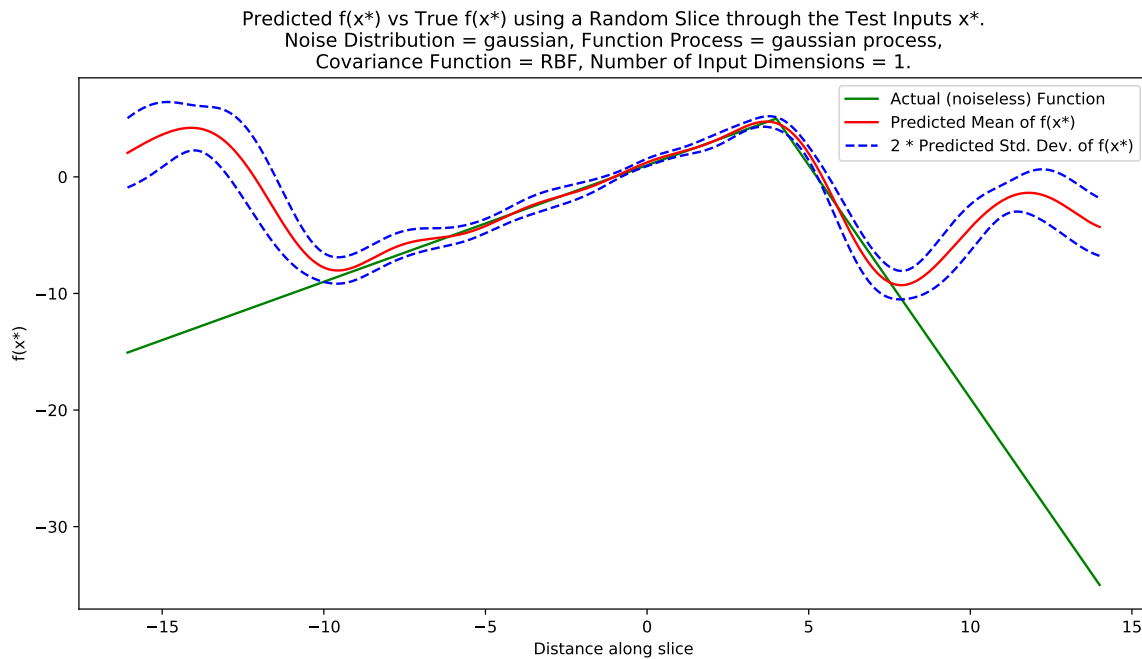
Model	Latent Kernel	Latent Process	Assumed Latent Noise	Observational Distribution	Additional Information
GPSSM	RBF	Gaussian	Gaussian	Fixed	N/A
DHGPSSM	RBF	Gaussian	Gaussian	Fixed	layout = [12]

**Table 3.22:** Additional Information for Example 1

*Results:*

Model	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
GPSSM	1.13	-1.50	83.36
DHGPSSM	<b>1.06</b>	<b>-1.48</b>	<b>45.53</b>

**Table 3.23:** Results for Example 1



**Figure 3.22:** A DHGPSSM  
 (Table 3.22, Row 2)

### Example 2:

*Aim:* This is another ‘sanity-check’ example comparing the performance of the DHGPSSM with Student-t noise and our Student-t noise version of the model of Svensson et al. [2016]. We expect all models to produce similar results. The set-up is the same as in section 3.2.6, Example 2 and for clarification, we have the same training and test data as in the original example.

#### Additional Information:

- The model called GPSSM refers to our Student-t noise version of the model of Svensson et al. [2016].
- The model called DHGPSSM is a deep Hilbert Gaussian process state-space model with latent Student-t noise assumed.
- The latent DHGPMs in the DHGPSSM will use  $\sin(x)$  activation functions.

More information is summarised in table 3.24.

Model	Latent Kernel	Latent Process	Assumed Latent Noise	Observational Distribution	Additional Information
GPSSM	RBF	Gaussian	Student-t	Fixed	See 3.2.6, Ex. 2
DHGPSSM	RBF	Gaussian	Student-t	Fixed	layout = [4]

**Table 3.24:** Additional Information for Example 2

Results:

Model	LL (on test data)	Mean Time (seconds)
GPSSM	-3.06	175.64
DHGPSSM	<b>-2.85</b>	<b>69.97</b>

Table 3.25: Results for Example 2

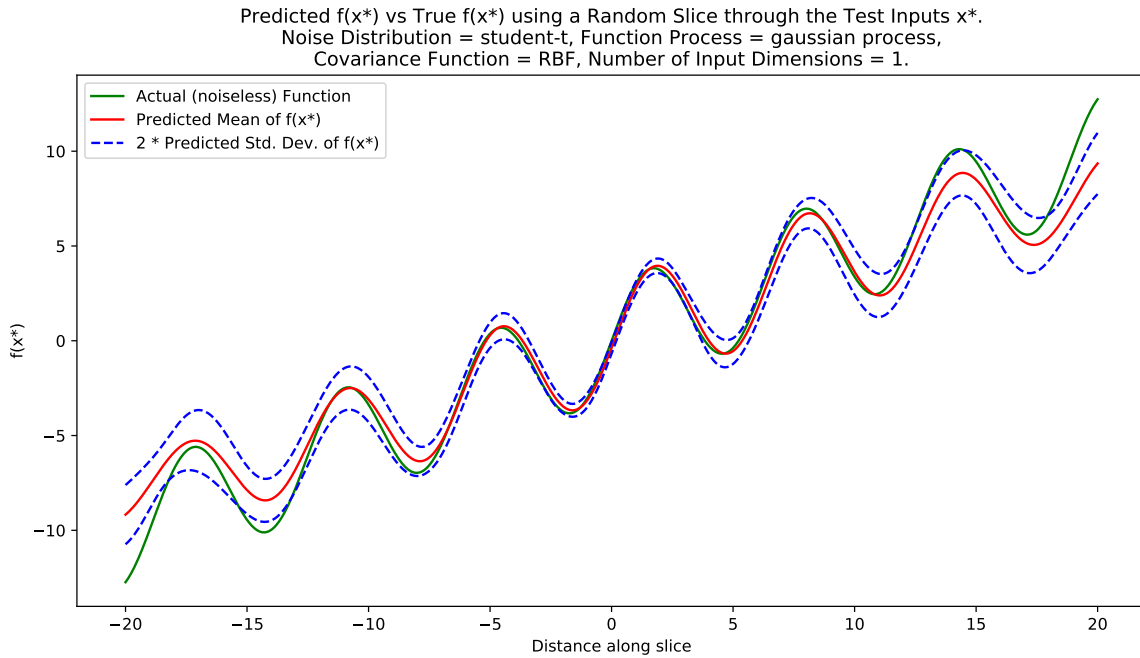


Figure 3.23: A DHGPSSM  
(Table 3.24, Row 2)

### Example 3:

*Aim:* To demonstrate that the DHGPSSM can learn with multi-dimensional latent spaces in reasonable time.

*Training/Test Data:* Given 500 points  $x_1, \dots, x_{500}$  evenly spread between -50 and 50, we generate the corresponding observations  $y_1, \dots, y_{500}$  using the system:

$$y = \sin(\pi x) + w_t \quad (3.135)$$

$$w_t \sim \mathcal{N}(0, 0.1^2) \quad (3.136)$$

Then, we discard the  $x_i$ . The training observations are  $y_1, \dots, y_{300}$  and the test observations are  $y_{301}, \dots, y_{500}$ .

Additional Information:

- We test the ability of the models to predict 4 steps ahead.
- The model called full GP is a full Gaussian process and it learns a mapping between  $y_t$  and  $y_{t+4}$  using training observations  $y_1, \dots, y_{300}$ .
- The model called DHGPSSM is a deep Hilbert Gaussian process state-space model with a three dimensional latent space and layout = [6, 2]. We will also have an observation distribution:

$$y_t = \boldsymbol{\alpha}^T \mathbf{x}_t + v_t \quad (3.137)$$

$$v_t \stackrel{iid}{\sim} \mathcal{N}(0, R) \quad (3.138)$$

where  $\boldsymbol{\alpha}^T = (a, b, c)$  and the parameters  $a, b, c$  and  $R$  are unknown and found during training.

- The latent DHGPMs in the DHGPSSM will use  $\tanh(x)$  activation functions (only because using  $\sin(x)$  might be seen as an unfair test).

Results:

Model	Kernel	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP	RBF	0.51	-0.80	<b>2.93</b>
DHGPSSM	RBF	<b>0.14</b>	<b>0.40</b>	174.22

Table 3.26: Results for Example 3

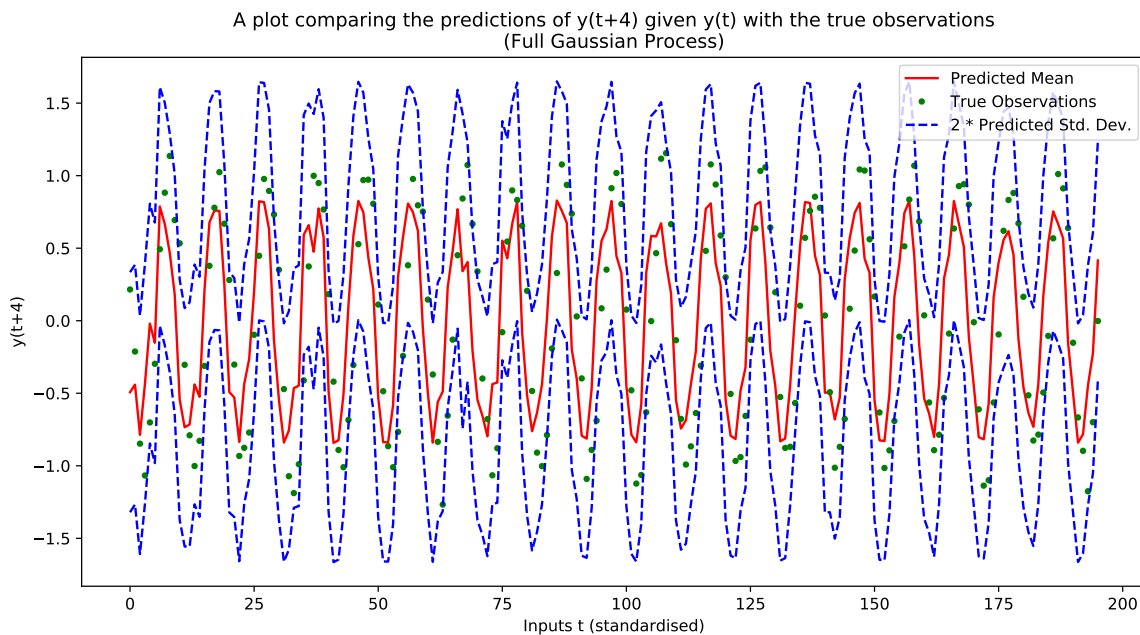
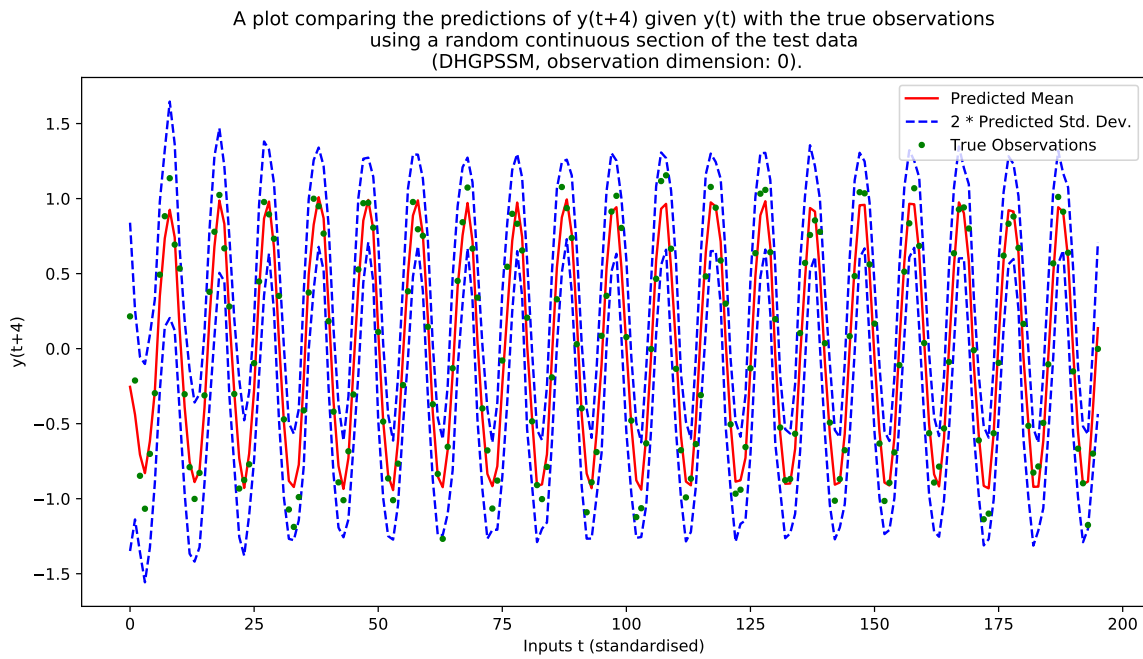
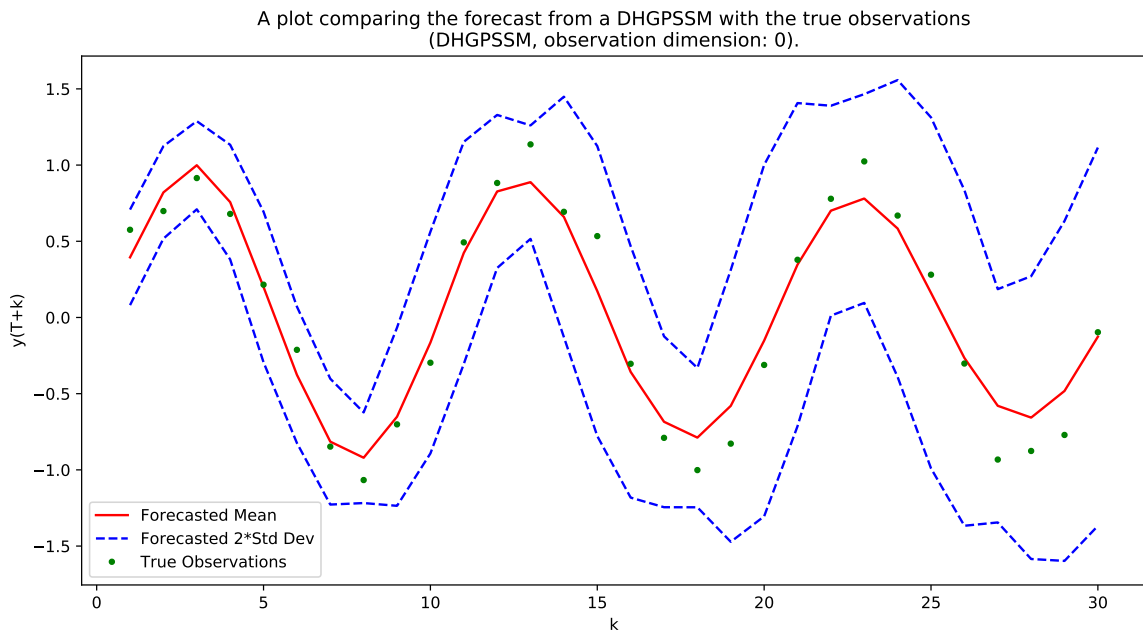


Figure 3.24: Predicting  $y_{t+4}$  given  $y_t$  using the full GP



**Figure 3.25:** Predicting  $y_{t+4}$  given  $y_t$  using the DHGPSSM. The poor performance with small  $t$  is because we use a particle filter and do not have a lot of information which can help to calculate the latent states at this point.



**Figure 3.26:** Forecasting using a DHGPSSM (RMSE: 0.17, LL = 0.05)

**Example 4:**

*Aim:* To demonstrate that the DHGPSSM can learn with multi-dimensional latent spaces and with non-synthetic data.

*The dataset and task:* Predicting yearly sunspot numbers using the SILSO sunspot dataset<sup>7</sup>.

*Training Data:* The first 200 observations.

*Test Data:* The remaining 115 observations.

*Additional Information:*

- We test the ability of the models to predict 2 steps ahead.
- The model called full GP is a full Gaussian process and it learns a mapping between  $y_t$  and  $y_{t+2}$  using training observations  $y_1, \dots, y_{200}$ .
- The model called DHGPSSM is a deep Hilbert Gaussian process state-space model with a three dimensional latent space and layout = [2]. We will also have an observation distribution:

$$y_t = \boldsymbol{\alpha}^T \mathbf{x}_t + v_t \quad (3.139)$$

$$v_t \stackrel{iid}{\sim} \mathcal{N}(0, R) \quad (3.140)$$

where  $\boldsymbol{\alpha}^T = (a, b, c)$  and the parameters  $a, b, c$  and  $R$  are unknown and found during learning.

- The latent DHGPMs in the DHGPSSM will use  $\sin(x)$  activation functions.

*Results:*

Model	Kernel	RMSE (on test data)	LL (on test data)	Mean Time (seconds)
Full GP	RBF	43.57	-1.77	1.47
DHGPSSM	RBF	<b>34.94</b>	<b>-1.68</b>	96.39

**Table 3.27:** Results for Example 4

<sup>7</sup><https://datamarket.com/data/set/4apm/yearly-mean-total-sunspot-number>

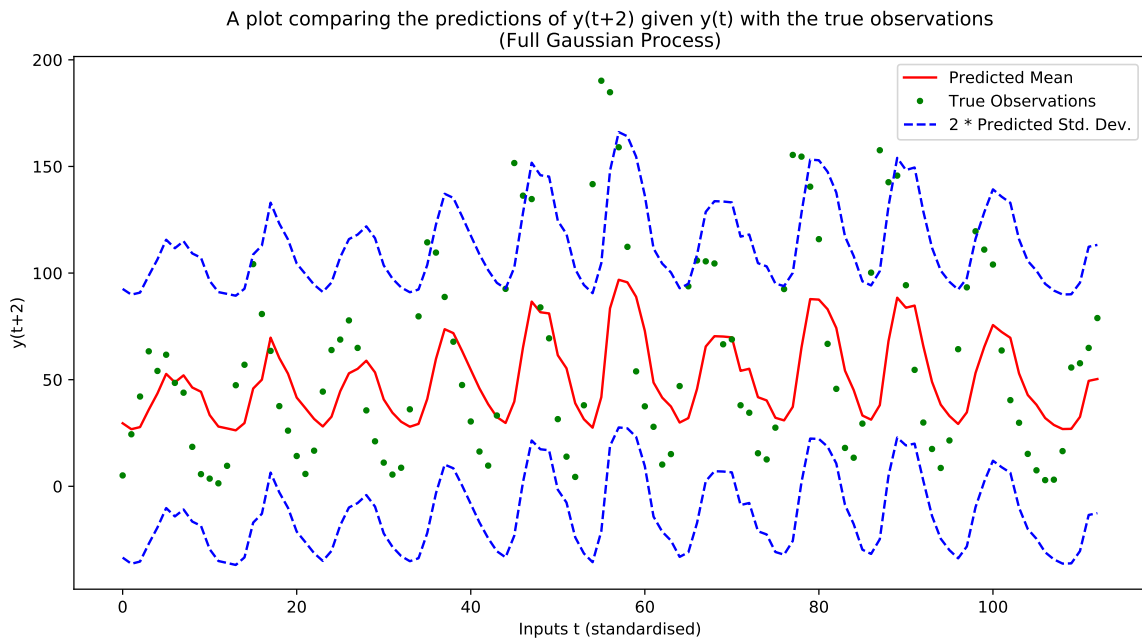


Figure 3.27: Predicting  $y_{t+2}$  given  $y_t$  using the full GP

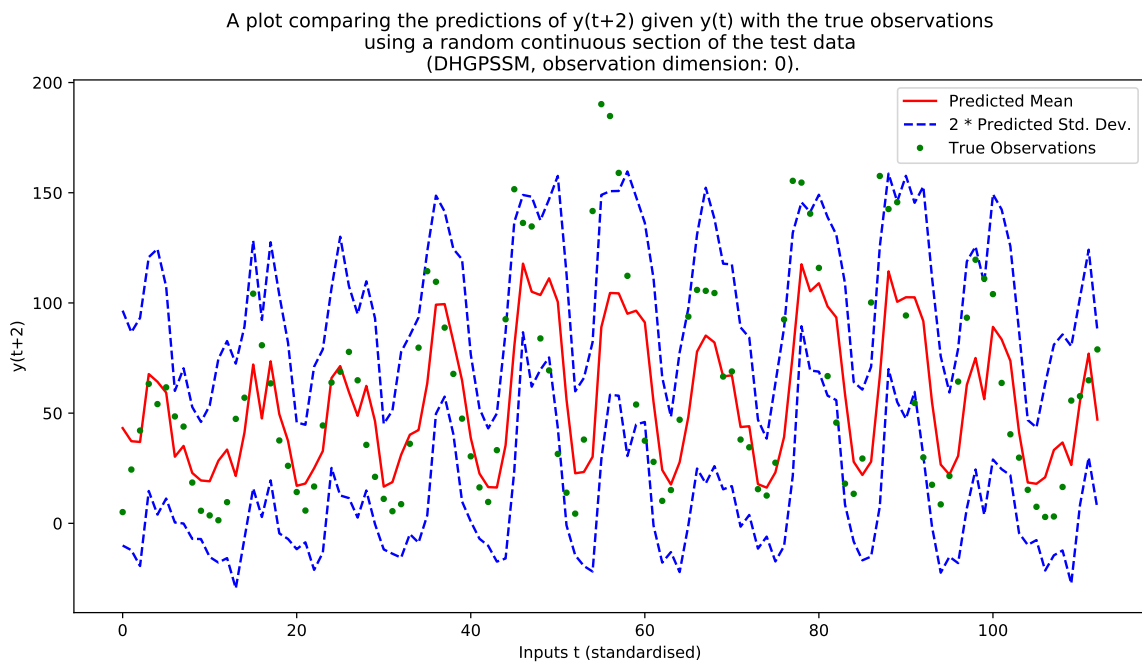


Figure 3.28: Predicting  $y_{t+2}$  given  $y_t$  using the DHGPSSM.



# Chapter 4

## Conclusion

### 4.0.3 Conclusions

In this thesis, we have added features to the GPSSM of Svensson et al. [2016], improved the Hilbert reduced-rank GP of Solin and Särkkä [2014] and created a novel GPSSM called the deep Hilbert Gaussian process state-space model. Here, we look at each contribution in turn and provide some concluding remarks.

1. The *Adaptive* covariance seems to be promising and although there is a computation cost compared with other covariance functions such as the RBF or Matern32 covariance, its performance is good and it seems to be particularly useful for DHGPMs and when combined with the manifold-Hilbert reduced-rank Gaussian process.
2. The manifold-Hilbert reduced-rank GP retains the properties of the original mGP; moreover, it is faster than the full GP even though the original mGP was slower than the full GP. From the examples, we can see that it is particularly good at learning functions with sharp edges. Furthermore, when using a  $\sin(x)$  activation, the manifold-Hilbert reduced-rank GP can learn some challenging periodic functions.
3. Our methods for learning with Student-t noise work even under distributions with small degrees of freedom and it is not surprising that it outperforms methods that do not assume a Student-t noise distribution. From the examples, we can see that it is almost always a good choice to use a Student-t noise distribution even when the true noise distribution is Gaussian. This is because it does not incur a significant computational cost (especially in the DHGPM) but with high degrees of freedom, the Student-t distribution tends to a Gaussian distribution: the models can *learn* the noise distribution.
4. Online learning and forgetful online learning both tend to a good solution as the amount of samples and data in the system increases. It is worth pointing out that sometimes these methods perform better than offline learning and we believe this might be because the weighting system can lead to bad samples (which are often the early samples even after burn-in) having a small weight.

Moreover, forgetful online learning, in particular, can ‘forget’ outliers which would usually have a detrimental effect on the sample quality.

5. Learning with Student-t process state-space models was somewhat disappointing because it often had slightly worse (or similar) performance compared with using Gaussian processes. This is contrary to the results stated in Shah et al. [2014]; however, we believe this might be because EMCEE has a challenging time sampling from models which use Student-t processes.
6. The deep Hilbert Gaussian process model trains much faster than the Hilbert reduced-rank GP when we have high-dimensional inputs. However, for small datasets it often takes longer to train than the full GP but it does seem to consistently produce better results than the full GP. With large datasets, the DHGPM produces good results while also having a training time much less than the full GP. As a result, we only recommend using the DHGPM for datasets larger than 1000 data points.
7. The deep Hilbert Gaussian process state-space model (DHGPSSM) can learn with high-dimensional latent spaces significantly quicker than the model of Svensson et al. [2016]. However, the limiting factor is the speed of latent state samples via PGAS, which can be very with large datasets. That said, it has shown competitive performance compared with the autoregressive full GP.

#### 4.0.4 Future Work

Having studied GPSSMs in depth, we have identified three areas for future work.

1. High-Performance implementation of GPSSMs: We have implemented our models in pure Python; however, rewriting some of the key bottlenecks (optimisation steps, PGAS samples) in a low-level language or using GPUs should give a significant performance benefit.
2. Combine DHGPMs and recurrent neural networks: Instead of using a DHGPM inside a state-space model, it might be worth trying to combine DHGPMs with recurrent neural networks by placing matrix-normal priors on all the weight matrices in a similar way to how we constructed the DHGPMs. It would be interesting to see if this improves the performance of the recurrent neural network, and it could provide predictive distributions rather than point estimates.
3. Understand relationship between machine learning and functional analysis: it seems that there is a link between statistical machine learning and functional analysis. Considering that machine learning is a relatively young topic compared with function analysis, it is worth investigating whether any of the techniques of function analysis can be applied to machine learning problems.

Overall, we believe that GPSSMs could become one of the best methods for identification of time-dependant complex systems; however, scalability needs to be improved.

# Bibliography

- Álvarez, M. A. and Lawrence, N. D. (2009). Sparse Convolved Gaussian Processes for Multi-output Regression. In *Advances in Neural Information Processing Systems 21*, pages 57–64. Curran Associates, Inc. [Cited on page 11.]
- Álvarez, M. A., Luengo, D., Titsias, M. K., and Lawrence, N. D. (2010). Efficient Multioutput Gaussian Processes through Variational Inducing Kernels. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics (AISTATS)*. [Cited on page 13.]
- Álvarez, M. A., Rosasco, L., and Lawrence, N. D. (2012). Kernels for Vector-Valued Functions: a Review. <https://arxiv.org/pdf/1106.6251.pdf>. [Cited on pages 5, 6, 10, and 11.]
- Bishop, C. M. (2009). *Pattern Recognition and Machine Learning*. Springer, New York, USA. [Cited on page 23.]
- Calandra, R., Peters, J., Rasmussen, C. E., and Deisenroth, M. P. (2016). Manifold Gaussian processes for regression. In *Proceedings of the IEEE International Joint Conference on Neural Networks*. [Cited on pages 3, 39, 41, 42, 43, and 73.]
- Chicago Board Options Exchange (2017). VIX Historical Data. <http://www.cboe.com/products/vix-index-volatility/vix-options-and-futures/vix-index/vix-historical-data>. [Cited on page 90.]
- Dahlin, J. and Schön, T. (2016). Getting started with particle Metropolis-Hastings for inference in nonlinear dynamical models. *Journal of Statistical Software*. [Cited on page 27.]
- Damianou, A., Titsias, M. K., and Lawrence, N. D. (2011). Variational Gaussian process dynamical systems. In *Advances in Neural Information Processing Systems*. [Cited on page 36.]
- Damianou, A. C. and Lawrence, N. D. (2013). Deep Gaussian Processes. In *Proceedings of the 16th International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 31. [Cited on page 73.]
- Deisenroth, M. P. (2010). *Efficient Reinforcement Learning using Gaussian Processes*. Karlsruhe Institut für Technologie Scientific Publishing, Straße am Forum 2, D-76131 Karlsruhe. [Cited on pages 21, 22, 23, and 34.]

- Deisenroth, M. P., Huber, M. F., and Hanebeck, U. D. (2009). Analytic Moment-based Gaussian process filtering. In *Proceedings of the 26th International Conference on Machine Learning (ICML)*, pages 225–232. [Cited on pages 34 and 35.]
- Deisenroth, M. P., Huber, M. F., and Hanebeck, U. D. (2012). Robust Filtering and Smoothing with Gaussian Processes. In *IEEE Transactions on Automatic Control*, volume 57(7), pages 1865–1871. [Cited on page 34.]
- Deisenroth, M. P. and Mohamed, S. (2016). Expectation propagation in Gaussian process dynamical systems. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2609–2617. [Cited on page 35.]
- Doucet, A., Freitas, N., and Gordon, N. (2001). *An Introduction to Sequential Monte Carlo Methods*. Springer, New York, USA. [Cited on pages 20, 21, 24, and 26.]
- Doucet, A. and Johansen, A. M. (2012). A Tutorial On Particle Filtering and Smoothing: Fifteen years later. Oxford, UK. [Cited on pages 20, 21, 24, 25, 26, and 27.]
- Eleftheriadis, S., Nicholson, T. F. W., Deisenroth, M. P., and Hensman, J. (2017). Identification of Gaussian Process State Space Models. <https://arxiv.org/pdf/1705.10888.pdf>. [Cited on page 36.]
- Foreman-Mackey, D., Hogg, D. W., Lang, D., and Goodman, J. (2013). EMCEE: The MCMC Hammer. *PASP*, 125:306–312. [Cited on pages 57 and 58.]
- Frigola, R. (2015). *Bayesian Time Series Learning with Gaussian Processes*. Cambridge, UK. [Cited on pages 3, 12, 28, 35, 36, and 92.]
- Frigola, R., Lindsten, F., Schön, T., and Rasmussen, C. E. (2013). Bayesian inference and learning in Gaussian process state-space models with particle MCMC. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3156–3164. [Cited on pages 35 and 36.]
- Frigola, R., Lindsten, F., Schön, T., and Rasmussen, C. E. (2014). Variational Gaussian process state-space models. In *Advances in Neural Information Processing Systems (NIPS)*, pages 3680–3688. [Cited on page 36.]
- Goodman, J. and Weare, J. (2010). Ensemble samplers with affine invariance. In *Communications in Applied Mathematics and Computational Science*. [Cited on page 57.]
- GPy (since 2012). GPy: A Gaussian process framework in python. <http://github.com/SheffieldML/GPy>. [Cited on pages 11 and 17.]
- Gupta, A. K. and Nagar, D. K. (2000). *Matrix Variate Distributions*. CRC Press. [Cited on pages 8 and 61.]
- Hamilton, J. D. (1994). *Time Series Analysis*. Princeton University Press, Princeton, New Jersey. [Cited on pages 1, 3, 28, and 29.]

- Iranmanesh, A., Arashi, M., and Tabatabaey, S. M. M. (2010). On Conditional Applications of Matrix Variate Normal Distributions. In *Iranian Journal of Mathematical Sciences and Informatics*, volume Vol. 5, No. 2, pages 33–44. [Cited on page 61.]
- Ko, J. and Fox, D. (2008). GP-BayesFilters: Bayesian Filtering Using Gaussian Process Prediction and Observation Models. In *Proceedings of the 2008 IEEE/RSJ International Conference on Intelligent Robots and Systems*. [Cited on page 34.]
- Lawrence, N. D. (2004). Gaussian Process Latent Variable Models for Visualisation of High Dimensional Data. In *Advances in Neural Information Processing Systems (NIPS)*. [Cited on page 35.]
- Lehmann, E. L. and Casella, G. (1998). *Theory of Point Estimation (Second Edition)*. Springer, New York, USA. [Cited on page 28.]
- Lichman, M. (2013). UCI machine learning repository. <http://archive.ics.uci.edu/ml>. [Cited on page 90.]
- Lindsten, F., Jordan, M. I., and Schön, T. (2014). Particle Gibbs with Ancestor Sampling. In *Journal of Machine Learning Research*, volume 15, pages 2145–2184. [Cited on pages 24, 26, 27, 28, 38, 52, 59, 63, and 96.]
- Mackay, D. J. C. (1998). *Introduction to Gaussian Processes*. Springer-Verlag. [Cited on page 42.]
- Quiñonero-Candela, J. and Rasmussen, C. E. (2005). A Unifying View of Sparse Approximate Gaussian Process Regression. In *Journal of Machine Learning Research*, volume 6, pages 1939–1959. [Cited on pages 12 and 13.]
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press, Cambridge, Massachusetts. [Cited on pages 3, 5, 7, 9, 10, and 13.]
- Roth, M. (2013). On the Multivariate- $t$  Distribution. <http://users.isy.liu.se/en/rt/roth/student.pdf>. [Cited on page 56.]
- Schäcke, K. (2013). On the Kronecker Product. <https://www.math.uwaterloo.ca/~hwoikowi/henry/reports/kronthesisschaecke04.pdf>, Department of Mathematics, University of Waterloo, Canada. [Cited on page 8.]
- Shah, A., Wilson, A., and Ghahramani, Z. (2014). Student- $t$  Processes as Alternatives to Gaussian Processes. In *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics (AISTATS)*, Reykjavik, Iceland. [Cited on pages 60 and 110.]
- Solin, A. and Särkkä (2014). Hilbert Space Methods for Reduced-Rank Gaussian Process Regression. <https://arxiv.org/pdf/1401.5508.pdf>. [Cited on pages 3, 4, 13, 14, 15, 16, 39, 40, 41, 60, 61, 72, 73, 77, 78, 79, 92, and 109.]

Svensson, A., Solin, A., Särkkä, S., and Schön, T. B. (2016). Computationally efficient Bayesian learning of Gaussian process state space models. In *Proceedings of 19<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS)*, Cadiz, Spain. [Cited on pages 4, 14, 27, 36, 37, 38, 50, 51, 52, 55, 56, 58, 64, 65, 66, 68, 71, 73, 92, 101, 102, 103, 109, and 110.]

Titsias, M. K. (2009). Variational Learning of Inducing Variables in Sparse Gaussian Processes. In *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics (AISTAT)*. [Cited on page 13.]

Turner, R., Deisenroth, M. P., and Rasmussen, C. E. (2010). State-Space Inference and Learning with Gaussian Processes. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 9, pages 868–875. [Cited on pages 12 and 35.]

Wang, J. M., Fleet, D. J., and Hertzmann, A. (2008). Gaussian process dynamical models. In *IEEE Transactions on Pattern Analysis and Machine Intelligence*, volume 30, pages 283–298. [Cited on pages 35 and 36.]

Yahoo Finance (2017). S&P 500 Index. <https://finance.yahoo.com/quote/>. [Cited on page 90.]

Zhu, S., Yu, K., and Gong, Y. (2008). Predictive Matrix-Variate  $t$  Models. In *Advances in Neural Information Processing Systems*. [Cited on page 61.]