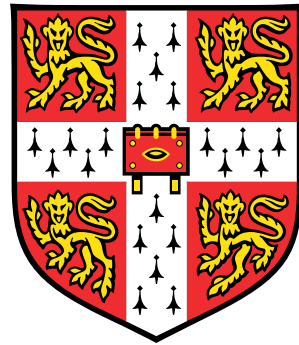


The Generalised Gaussian Process Convolution Model



Wessel Pieter Bruinsma

Department of Engineering

University of Cambridge

This dissertation is submitted for the degree of

Master of Philosophy

Sidney Sussex College

12 August 2016

Declaration

I, Wessel Pieter Bruinsma of Sidney Sussex College, being a candidate for the M.Phil in Machine Learning, Speech and Language Technology, hereby declare that this thesis and the work described in it are my own work, unaided except as may be specified below, and that the thesis does not contain material that has already been used to any substantial extent for a comparable purpose. This thesis, excluding appendices, contains approximately 7300 words.

Wessel Pieter Bruinsma

12 August 2016

Acknowledgements

First and foremost, I thank my supervisor Dr. Richard Turner. His guidance and support throughout the project have been indispensable.

I also thank Dr. Felipe Tobar, whose thoughts have been essential in the development of the project.

Abstract

This thesis formulates the Generalised Gaussian Process Convolution Model (GGPCM), which is a generalisation of the Gaussian Process Convolution Model presented by Tobar et al. [2015b]. The GGPCM provides a theoretical framework for nonparametric kernel models of multidimensional signals defined on multidimensional input spaces. We show that the GGPCM generalises and connects existing work; most notably, we derive a dual formulation of the cross-spectral mixture kernel presented by Ulrich et al. [2015]. Finally, we use the GGPCM to develop the Deep Kernel Model, which presents a new network structure for unsupervised learning.

Contents

List of Figures	xiii
List of Tables	xv
List of Models	xvii
Notation	xix
1 Introduction	1
1.1 Motivation	1
1.2 Contribution of Thesis and Outline	3
2 Modelling Dynamical Signals	5
2.1 Introduction	5
2.2 The Linear State-Space Model	5
2.3 A General-Purpose Model of Dynamical Signals	6
2.4 Conclusion	7
3 The Generalised Gaussian Process Convolution Model	9
3.1 Introduction	9
3.2 Gaussian Process Regression	9
3.3 The Generalised Gaussian Process Convolution Model	11
3.3.1 Interpretation and Choice of the Kernel in the Nonparametric Kernel Model	14
3.3.2 Illustrative Samples of the Nonparametric Kernel Model	16
3.3.3 Expressivity of White Noise Excitation	19
3.4 The Approximate Kernel Model	21
3.4.1 Interpretation of the Kernel Approximation	23

3.4.2	The Case of the Diagonal Multi-Output Decaying Exponentiated- Quadratic Kernel	23
3.5	Related Work	25
3.6	Inference in the Nonparametric Kernel Model	26
3.7	Conclusion	30
3.8	Discussion	30
4	Multi-Task Learning	31
4.1	Introduction	31
4.2	Mixing Models	31
4.3	The Mixing Model Hierarchy	32
4.4	Conclusion	33
5	The Deep Kernel Model	37
5.1	Introduction	37
5.2	The Deep Kernel Model	37
5.2.1	Network Interpretation	41
5.3	Illustrative Samples	43
5.4	Conclusion	44
5.5	Discussion	44
A	Solution of the Linear State-Space Model	47
A.1	Time-Variant Solution	47
A.2	Time-Invariant Solution	48
B	Properties of the Multivariate Gaussian Distribution	51
B.1	Marginal and Conditional Distribution	51
B.2	Kullback-Leibler Divergence	51
C	Multivariate Matrix-Valued Gaussian Processes	53
D	Gaussian Processes in Practice	55
D.1	Implementation of Gaussian Process Models	55
D.2	Nearest Symmetric Positive-Semidefinite Matrix	56
E	Circulant Approximation of Stationary Multi-Output Kernel Matrices	61
E.1	Introduction	61

Contents	xi
E.2 Circulant Approximation of Toeplitz Matrices	61
E.3 Circulant Approximation of Stationary Multi-Output Kernel Matrices	63
E.4 Approximating Determinants	65
E.5 Approximating Products Involving an Inverse	67
E.6 Conclusion	70
F Exponentiated Quadratic Forms	71
F.1 Introduction	71
F.2 General Form	71
F.3 Kronecker-Structured Form	75
F.4 Conclusion	78
G Roots of Kernels	79
H Approximate Kernel Model	81
I Variational Free Energy of the Nonparametric Kernel Model	83
I.1 Introduction	83
I.2 Predictive Mean	84
I.3 Predictive Autocovariance	85
I.4 Integrals $I_{\cdot}^{(\cdot)}$	87
I.4.1 Kernels	88
I.4.2 Integral $I_{\cdot}^{(L, \cdot)}$	89
I.4.3 Integral $I_{\cdot}^{(Q_1, \cdot)}$	90
I.4.4 Integral $I_{\cdot}^{(Q_2, \cdot)}$	91
I.4.5 Integral $I_{\cdot}^{(Q_3, \cdot)}$	92
I.4.6 Integral $I_{\cdot}^{(Q_4, \cdot)}$	94
I.5 Conclusion	94
References	95
Index	99

List of Figures

1.1	Posterior distribution over an unknown function as evidence is accumulated	2
1.2	Posterior distribution over an unknown function for different kernels . . .	2
3.1	Generation of the outputs of in Model 5	13
3.2	Generative process of Model 5	16
3.3	One-dimensional samples from Model 5	17
3.4	Two-dimensional samples from Model 5	17
3.5	Interpolation between two kernels sampled from Model 5	18
3.6	Relationship of Model 5 and Model 6 to current literature	26
4.1	Organisation of multi-output models from current literature	34
5.1	Graphical model of Model 13	42
5.2	Generative process of Model 13	45
5.3	Observations from Model 13	46
E.1	Circulant approximation of a stationary kernel	62

List of Tables

4.1	Identification of multi-output models from current literature	33
-----	---	----

List of Models

1	Time-Variant Model	6
2	Time-Invariant Model	6
3	Multidimensional Time-Invariant Model	9
4	Generalised Gaussian Process Convolution Model (GGPCM)	11
5	Nonparametric Kernel Model (NKM)	11
6	Approximate Kernel Model (AKM)	22
7	Basis Function Model	29
8	Instantaneous Mixing Model (IMM)	31
9	Convolutional Mixing Model (CMM)	32
10	Gaussian Process Convolution Model [Tobar et al., 2015b]	37
11	Gaussian Process Convolution Model (Explicit Decay)	38
12	Deep Gaussian Process Convolution Model	39
13	Deep Kernel Model	40
14	Approximate Deep Kernel Model	41

Notation

General

We use Lagrange's notation for differentiation.

$(1, 1), \dots, (N, M)$ Shorthand for

$$(1, 1), \dots, (N, 1), (1, 2), \dots, (N, 2), \dots, (1, M), \dots, (N, M)$$

i	Imaginary unit or, depending on the context, index
\cdot^*	Complex conjugate
$ \cdot $	Modulus or, depending on the context, determinant
$\mathbb{1}_X$	Indicator function
$\mathbb{1}$	Shorthand for $\mathbb{1}_{\{0\}}$
δ	Dirac delta function
R	Reversal function; for any function $f_1 : X \rightarrow Y$ the reversal function $R : Y^X \rightarrow Y^X$ yields the function $f_2 = R(f_1)$ such that $f_2(x) = f_1(-x)$ for all $x \in X$
\mathcal{O}	Bachmann-Landau notation for asymptotic behaviour

Probability Theory

p	Probability measure
$\stackrel{d}{=}$	Equality in distribution
\mathbb{E}	Expected value
$\mathcal{N}(\cdot; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Multivariate Gaussian distribution with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$
$D_{KL}(\cdot \parallel \cdot)$	Kullback-Leibler divergence

Linear Algebra

Unbolded symbols x denote scalars, bolded lower-case symbols \mathbf{x} vectors and bolded upper-case symbols \mathbf{X} matrices or arrays—by array we mean the multidimensional generalisation of the two-dimensional matrix. Vectors are column vectors unless specified otherwise.

\mathbf{I}	Identity matrix
\mathbf{I}_N	$N \times N$ identity matrix
\cdot^T	Transpose
\cdot^H	Conjugate transpose
$ \cdot $	Determinant or, depending on the context, modulus
$\ \cdot\ _p$	p -norm
$\ \cdot\ _F$	Frobenius norm
\otimes	Kronecker product
$\text{diag}(\mathbf{X}_1, \dots, \mathbf{X}_N)$	Block diagonal matrix; yields

$$\begin{bmatrix} \mathbf{X}_1 & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{X}_2 & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{X}_{N-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{X}_N \end{bmatrix}$$

$\text{circ}(\mathbf{X}_1, \dots, \mathbf{X}_N)$ Block circulant matrix; yields

$$\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_N & \cdots & \mathbf{X}_3 & \mathbf{X}_2 \\ \mathbf{X}_2 & \mathbf{X}_1 & \cdots & \mathbf{X}_4 & \mathbf{X}_3 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{X}_{N-1} & \mathbf{X}_{N-2} & \cdots & \mathbf{X}_1 & \mathbf{X}_N \\ \mathbf{X}_N & \mathbf{X}_{N-1} & \cdots & \mathbf{X}_2 & \mathbf{X}_1 \end{bmatrix}$$

$\text{stoepl}(\mathbf{X}_1, \dots, \mathbf{X}_N)$ Symmetric block Toeplitz matrix; yields

$$\begin{bmatrix} \mathbf{X}_1 & \mathbf{X}_2 & \cdots & \mathbf{X}_{N-1} & \mathbf{X}_N \\ \mathbf{X}_2 & \mathbf{X}_1 & \cdots & \mathbf{X}_{N-2} & \mathbf{X}_{N-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{X}_{N-1} & \mathbf{X}_{N-2} & \cdots & \mathbf{X}_1 & \mathbf{X}_2 \\ \mathbf{X}_N & \mathbf{X}_{N-1} & \cdots & \mathbf{X}_2 & \mathbf{X}_1 \end{bmatrix}$$

vec Matrix vectorisation function; if $\mathbf{X} \in \mathbb{R}^{N \times M}$, then

$$\text{vec } \mathbf{X} = [X_{1,1} \ \cdots \ X_{N,M}]^T$$

X_{i_1, \dots, i_N} Vector, matrix and array indexing; yields element (i_1, \dots, i_N) of the array \mathbf{X}

$\mathbf{X}_{i,:}$ Fixation of first dimension of a matrix and ranging over the second; yields vector containing the i 'th row of \mathbf{X}

$\mathbf{X}_{:,i}$ Fixation of second dimension of a matrix and ranging over the first; yields vector containing the i 'th column of \mathbf{X}

$\mathbf{X}_{\underbrace{i_1, \dots, i_n, \dots}_{N \text{ indices}}}$ Fixation of an arbitrary subset of dimensions and ranging over the remaining; yields the array whose (i_{n+1}, \dots, i_N) 'th element is given by X_{i_1, \dots, i_N}

$f(\mathbf{X})$ Element-wise function application; if \mathbf{X} is M -dimensional— $\mathbf{X} \in \mathbb{R}^{N_1 \times \cdots \times N_M}$ —but f takes a $(M - m)$ -dimensional object— $f : \mathbb{R}^{N_{m+1} \times \cdots \times N_M} \rightarrow C$ —then $f(\mathbf{X})$ is the m -dimensional array whose (i_1, \dots, i_m) 'th element is given by $f(\mathbf{X}_{i_1, \dots, i_m, :, \dots, :})$

$f(\mathbf{X}, \mathbf{Y})$ Element-wise function application; if $\mathbf{X} \in \mathbb{R}^{N_1 \times \cdots \times N_M}$ and $\mathbf{Y} \in \mathbb{R}^{N_1 \times \cdots \times N_K}$, but $f : \mathbb{R}^{N_{m+1} \times \cdots \times N_M} \times \mathbb{R}^{N_{k+1} \times \cdots \times N_K} \rightarrow C$, then $f(\mathbf{X}, \mathbf{Y})$ is the array whose $(i_1, \dots, i_m, j_1, \dots, j_k)$ 'th element is given by $f(\mathbf{X}_{i_1, \dots, i_m, :, \dots, :}, \mathbf{Y}_{i_1, \dots, i_k, :, \dots, :})$

Miscellaneous

* Convolution; for matrix-valued functions \mathbf{F}_1 and \mathbf{F}_2 defined as

$$(\mathbf{F}_1 * \mathbf{F}_2)(\mathbf{x}) = \int \mathbf{F}_1(\mathbf{x} - \mathbf{y}) \mathbf{F}_2(\mathbf{y}) \, d\mathbf{y}$$

\mathcal{F}_N $N \times N$ unitary discrete Fourier transform matrix

DFT Discrete Fourier transform

$\mathcal{F}_y\{f\}(\mathbf{x})$ Continuous Fourier transform of f ; defined as

$$\mathcal{F}_y\{f\}(\mathbf{x}) = \int f(\mathbf{y}) \exp(-2\pi i \mathbf{x}^T \mathbf{y}) \, d\mathbf{y}$$

$\mathcal{F}_f(\mathbf{x})$ Shorthand for $\mathcal{F}_y\{f\}(\mathbf{x})$

$\mathcal{S}_y\{f, g\}(\mathbf{x})$ Cross-spectral density between two wide-sense stationary processes f and g ; equal to

$$\mathcal{S}_y\{f, g\}(\mathbf{x}) = \mathcal{F}_z\{\mathbb{E}[f(\mathbf{x} + \mathbf{z})g^*(\mathbf{x})]\}(\mathbf{y})$$

$\mathcal{S}_{f,g}(\mathbf{y})$ Shorthand for $\mathcal{S}_y\{f, g\}(\mathbf{x})$

$(C, \mathbf{A}, \mathbf{b}, c)$ See Appendix F.2; shorthand for

$$C \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{b} + c\right)$$

(C, \mathbf{A}) See Appendix F.3; shorthand for $(C, \mathbf{A} \otimes \mathbf{I}, \mathbf{0}, 0)$

Acronyms

AKM Approximate Kernel Model (Model 6)

CGPM Collaborative Gaussian processes model [Nguyen and Bonilla, 2014]

CMOGPM Convolved multi-output Gaussian process model [Álvarez and Lawrence, 2011]

CSMK Cross-spectral mixture kernel [Ulrich et al., 2015]

GGPCM Generalised Gaussian Process Convolution Model (Model 4)

GPCM Gaussian Process Convolution Model [Tobar et al., 2015b]

GPRN Gaussian process regression network [Wilson et al., 2012]

ICM Intrinsic coregionalisation model [Goovaerts, 1997]

LCM	Linear coregionalisation model [Goovaerts, 1997]
LFM	Latent force model [Álvarez et al., 2009]
MTGPM	Multi-task Gaussian process [Bonilla et al., 2008]
NKM	Nonparametric Kernel Model (Model 5)
PSD	Power spectral density
SLFM	Semiparametric latent factor model [Teh and Seeger, 2005]
SMK	Spectral mixture kernel [Wilson and Adams, 2013]

1 | Introduction

1.1 Motivation

Gaussian processes elegantly provide means to model an unknown function. They give rise to Bayesian regression models in which one maintains a posterior distribution over the unknown function as evidence is accumulated; Figure 1.1 illustrates this process. Gaussian processes have been successfully applied in a wide variety of contexts. Rasmussen and Williams [2006] provide an excellent overview.

Gaussian processes are *nonparametric* models. As opposed to *parametric* models, there is no finite number of parameters that parametrises a Gaussian processes. Instead, the number of parameters grows with the amount of evidence that is accumulated. This property allows Gaussian processes to learn complex functions if plenty of evidence is available. Conversely, this property makes them robust against overfitting if only little evidence is available.

Their expressiveness and robustness however, come at a cost: Gaussian process models are often computationally expensive and one is forced to choose a *kernel*. The kernel of a Gaussian process reflects one's assumption on how the unknown function autocovaries.¹ The choice of the kernel is crucial; Figure 1.2 illustrates that the posterior distribution can wildly vary for different kernels. Unfortunately, determining which kernel to use is hard; the kernel thus poses a difficult design problem. We specifically refer to this problem as the *kernel design problem*.

A number of recent works address the kernel design problem: Duvenaud [2014] presents a way to search over a space of kernels through composition of existing kernels, and Wilson and Adams [2013] present a flexible kernel by modelling its power spectral

¹The autocovariance of an unknown function f specifies the covariance between any two function values $f(t_1)$ and $f(t_2)$.

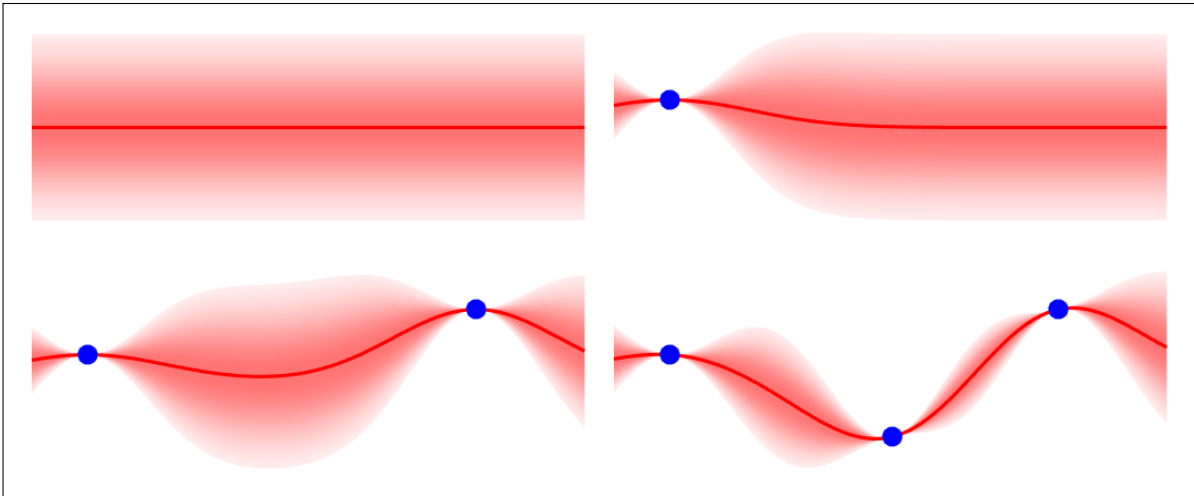


Figure 1.1: Posterior distribution over an unknown function as evidence is accumulated. The evidence is represented by blue dots. The red line represents the mean of the posterior distribution and the gradient indicates the marginal variance up to two standard deviations.

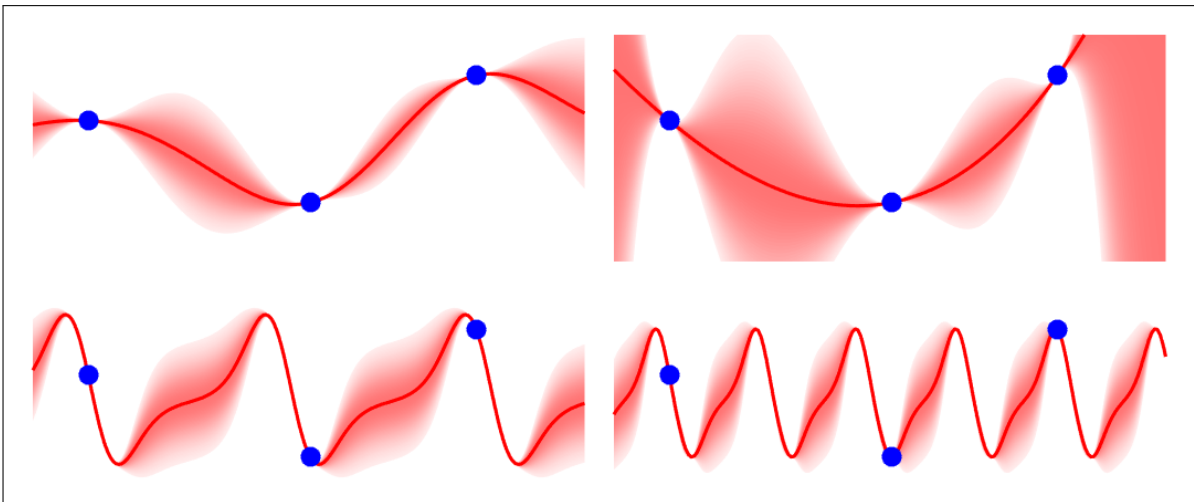


Figure 1.2: Posterior distribution over an unknown function for different kernels. The evidence is represented by the blue dots and is the same as in Figure 1.1. The red line represents the mean of the posterior distribution and the gradient indicates the marginal variance up to two standard deviations.

density with a sum of Gaussians. As noted by Tobar et al. [2015b], these approaches do not completely resolve the kernel design problem: computational expenses restrict the space of kernels that can be searched, and the large number of parameters introduced by the sum of Gaussians reintroduces the problem of overfitting. Instead, Tobar et al. [2015b] note that the kernel constitutes just another unknown function and thus propose to model it by another Gaussian; the results are encouraging. They call their model the Gaussian Process Convolution Model (GPCM).

The GPCM has only been formulated for one-dimensional signals defined on one-dimensional input spaces. Extending this model to multidimensional signals defined on multidimensional input spaces would enable numerous applications in econometrics, geostatistics and signal processing, among others.

1.2 Contribution of Thesis and Outline

The main contribution of this thesis is the formulation of the Generalised Gaussian Process Convolution Model (GGPCM). The GGPCM provides a theoretical framework for nonparametric kernel models of multidimensional signals defined on multidimensional input spaces.

Chapter 2 develops a general-purpose model of dynamical signals. This model will suggest a generalisation of the GPCM: the GGPCM. Thus, in hindsight, Chapter 2 motivates why this particular generalisation is appropriate.

Chapter 3 develops the GGPCM. We use the GGPCM to address the kernel design problem in multi-output Gaussian processes on multidimensional input spaces. Furthermore, we develop a dual formulation of the cross-spectral mixture kernel [Ulrich et al., 2015] and we show how the GGPCM connects to existing work.

Chapter 4 examines the GGPCM's connection to existing multi-output Gaussian process models. We provide an overview of the current literature.

Chapter 5 develops the Deep Kernel Model, which is a new network structure for unsupervised learning. We sample from the Deep Kernel Model and investigate its properties.

2 | Modelling Dynamical Signals

2.1 Introduction

For many years scientists have used the concept of a *system* to study dynamic processes of diverse nature. A system is a mathematical abstraction that is used to describe properties of the process we intend to study. *Dynamical systems* describe processes that evolve in time; the system then usually describes a transformation $T : (\mathbb{R}^M)^{\mathbb{R}} \rightarrow (\mathbb{R}^N)^{\mathbb{R}}$ from an input signal \mathbf{x} to an output signal $\mathbf{f} = T(\mathbf{x})$ —the latter is also called the *system response*.

This chapter develops a general-purpose model of dynamical signals. In doing so we take a systems-modelling perspective by postulating that every signal is the response of a particular system. Thus, equivalently, we develop a model of system responses. We therefore utilise a general description of dynamical *systems* to construct the desired model of dynamical *signals*.

2.2 The Linear State-Space Model

A widely used dynamical system model is one that describes mechanical, electrical, hydraulic, and thermal systems, among others. These systems are all compositions of resistive, inductive, capacitive and memristive elements that are connected through lossless transfer of energy; hence, they allow for a unified description. This description embodies the concept of *system state*, whose manifestation is a vector \mathbf{s} that at any time determines future output given future input, thus rendering past input irrelevant. More precisely, the description is a so-called *state-space model* whose general *linear* form is given by

$$\mathbf{s}'(t) = \mathbf{A}(t)\mathbf{s}(t) + \mathbf{B}(t)\mathbf{x}(t), \quad (2.1)$$

$$\mathbf{f}(t) = \mathbf{C}(t)\mathbf{s}(t) + \mathbf{D}(t)\mathbf{x}(t). \quad (2.2)$$

Examples of systems that are successfully described by linear state-space models are ubiquitous: linear state-space models describe electronics that are essential in everyday life; systems that control your car, airplanes and even rockets; and many phenomena in nature.

2.3 A General-Purpose Model of Dynamical Signals

Suppose that \mathbf{f} is a dynamical signal. By the postulate that every signal is the response of some system, we can assume that some system generated \mathbf{f} . Now, Section 2.2 suggests that we can safely assume this system to be of the form of Equations (2.1) and (2.2). In that case Appendix A.1 shows that \mathbf{f} admits the following parametrisation:

Model 1 (Time-Variant Model).

$$\mathbf{f}(t) = \int_{\mathbb{R}} \mathbf{H}(t, \tau) \mathbf{x}(\tau) \, d\tau.$$

Observe that \mathbf{f} is parametrised in terms of some matrix-valued function \mathbf{H} —the *impulse response*¹—and the postulated input \mathbf{x} .

In many applications of state-space models $\mathbf{A}(t)$, $\mathbf{B}(t)$, $\mathbf{C}(t)$ and $\mathbf{D}(t)$ are approximately constant.² In that case Appendix A.2 shows that $\mathbf{H}(t, \tau) = \mathbf{H}(t - \tau)$. Model 1 then reduces to the following model:

Model 2 (Time-Invariant Model).

$$\mathbf{f}(t) = \int_{\mathbb{R}} \mathbf{H}(t - \tau) \mathbf{x}(\tau) \, d\tau = (\mathbf{H} * \mathbf{x})(t).$$

Observe that Model 2 attains the form of a convolution.

¹Letting $\mathbf{x}(t) = \delta(t - t_0) \mathbf{x}_0$, $t \in \mathbb{R}$ yields the response $\mathbf{f}(t) = \mathbf{H}(t, t_0)$, $t \in \mathbb{R}$. That is, \mathbf{H} is the system response of an *impulse excitation*; in this sense \mathbf{H} is called the *impulse response*.

²For example, consider an electric circuit. In that case, approximating $\mathbf{A}(t)$, $\mathbf{B}(t)$, $\mathbf{C}(t)$ and $\mathbf{D}(t)$ as time invariant roughly corresponds to approximating resistances, capacitances and inductances as time invariant. This approximation is often reasonable.

2.4 Conclusion

We have developed a general-purpose model of dynamical signals: Model 2. Unfortunately, Model 2 leaves the impulse response \mathbf{H} and postulated input \mathbf{x} unspecified. This issue will be addressed in Chapter 3.

3 | The Generalised Gaussian Process Convolution Model

3.1 Introduction

Chapter 2 developed a general-purpose model of dynamical signals: the Time-Invariant Model (Model 2). By extending Model 2's input space to \mathbb{R}^K we derive the Multidimensional Time-Invariant Model:

Model 3 (Multidimensional Time-Invariant Model).

$$\mathbf{f}(\mathbf{t}) = \int_{\mathbb{R}^K} \mathbf{H}(\mathbf{t} - \boldsymbol{\tau}) \mathbf{x}(\boldsymbol{\tau}) \, \mathrm{d}\boldsymbol{\tau} = (\mathbf{H} * \mathbf{x})(\mathbf{t})$$

Model 3 is a deterministic model of multidimensional signals defined on multidimensional input spaces. This chapter shows that a stochastic version of Model 3 can be used to address the kernel design problem in multi-output Gaussian processes on multidimensional input spaces.

Recall that $\mathbf{f} : \mathbb{R}^K \rightarrow \mathbb{R}^N$, $\mathbf{H} : \mathbb{R}^K \rightarrow \mathbb{R}^{M \times N}$ and $\mathbf{x} : \mathbb{R}^K \rightarrow \mathbb{R}^M$.

3.2 Gaussian Process Regression

A Gaussian process defines a distribution over functions. More precisely, a stochastic process $f(t)$, $t \in \mathbb{R}$ is Gaussian if and only if for every $\mathbf{t} = [t_1 \ \dots \ t_T]^T$ it holds that $f(\mathbf{t}) = [f(t_1) \ \dots \ f(t_T)]^T$ is multivariate Gaussian distributed. We denote $f \sim \mathcal{GP}(m_f, \mathcal{K}_f)$ where m_f and \mathcal{K}_f denote respectively the mean function and kernel of f . This implies that $f(\mathbf{t}) \sim \mathcal{N}[m_f(\mathbf{t}), \mathcal{K}_f(\mathbf{t}, \mathbf{t})]$. Furthermore, we usually let $m_f = 0$ without loss of generality [Rasmussen and Williams, 2006] and denote $\mathbf{K}_{\mathbf{f}_1, \mathbf{f}_2} = \mathcal{K}_f(\mathbf{t}_1, \mathbf{t}_2)$ where $\mathbf{f}_1 = f(\mathbf{t}_1)$ and $\mathbf{f}_2 = f(\mathbf{t}_2)$.

Suppose that we observe data \mathbf{y} at \mathbf{t}_y generated by the process $y(t) = f(t) + \varepsilon(t)$, $t \in \mathbb{R}$ where $\varepsilon(t) \sim \mathcal{N}(0, \sigma^2)$, $t \in \mathbb{R}$. For unobserved function values \mathbf{f}^* at \mathbf{t}_{f^*} it then holds that

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}^* \end{bmatrix} \Bigg| \begin{bmatrix} \mathbf{t}_y \\ \mathbf{t}_{f^*} \end{bmatrix}, \boldsymbol{\theta} \sim \mathcal{N} \left(\mathbf{0}, \begin{bmatrix} \mathbf{K}_{\mathbf{y},\mathbf{y}} + \sigma^2 \mathbf{I} & \mathbf{K}_{\mathbf{y},\mathbf{f}^*} \\ \mathbf{K}_{\mathbf{f}^*,\mathbf{y}} & \mathbf{K}_{\mathbf{f}^*,\mathbf{f}^*} \end{bmatrix} \right)$$

where $\boldsymbol{\theta}$ denotes σ and the parameters of \mathcal{K}_f and m_f — $\boldsymbol{\theta}$ are also called the *hyperparameters*. We usually omit explicit conditioning \mathbf{t}_y and \mathbf{t}_{f^*} . According to Appendix B.1 we can perform prediction via

$$\mathbf{f}^* | \mathbf{y}, \boldsymbol{\theta} \sim \mathcal{N}[\mathbf{K}_{\mathbf{f}^*,\mathbf{y}}(\mathbf{K}_{\mathbf{y},\mathbf{y}} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, \mathbf{K}_{\mathbf{f}^*,\mathbf{f}^*} - \mathbf{K}_{\mathbf{f}^*,\mathbf{y}}(\mathbf{K}_{\mathbf{y},\mathbf{y}} + \sigma^2 \mathbf{I})^{-1} \mathbf{K}_{\mathbf{y},\mathbf{f}^*}].$$

This is the Gaussian process regression framework.

The values of the hyperparameters should be chosen such that the model best fits the true data-generating process. In this sense, as argued by Rasmussen and Williams [2006], their values can be determined by maximising the *marginal likelihood*—also called the *evidence*—given by

$$p(\mathbf{y} | \boldsymbol{\theta}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \propto \underbrace{|\boldsymbol{\Sigma}|^{-1/2}}_{\text{complexity penalty}} \underbrace{\exp \left[-\frac{1}{2} \|\boldsymbol{\Sigma}^{-1/2}(\mathbf{y} - \boldsymbol{\mu})\|_2^2 \right]}_{\text{compatibility with data}}.$$

for some $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$. Observe that the marginal likelihood is a trade-off between the model's compatibility with the data and $|\boldsymbol{\Sigma}|^{-1/2}$ being small; that is, the evidence is large if the model is compatible with the data, but only if the model is not too *complex* in the sense that it is compatible with *any* data—if $\boldsymbol{\Sigma}^{-1/2}$ is such that $\|\boldsymbol{\Sigma}^{-1/2}(\mathbf{y} - \boldsymbol{\mu})\|_2^2$ is small for any \mathbf{y} , then $|\boldsymbol{\Sigma}^{-1/2}| = |\boldsymbol{\Sigma}|^{-1/2}$ is small and thereby the evidence must be small. Thus maximising the evidence tends to explain the data in a *simple* way. This tendency is commonly recognised as a manifestation of *Occam's razor* [MacKay, 2002].

Finally, Appendix C shows how Gaussian processes can be used to define distributions over multivariate matrix-valued functions.

3.3 The Generalised Gaussian Process Convolution Model

Model 3 leaves the input signal \mathbf{x} and impulse response \mathbf{H} unspecified. Since they constitute unknown functions, it is sensible to model them by two Gaussian processes; we then obtain the Generalised Gaussian Process Convolution Model (GGPCM):

Model 4 (Generalised Gaussian Process Convolution Model (GGPCM)). *Draw*

$$\begin{aligned}\mathbf{H} &\sim \mathcal{GP}(\mathbf{0}, \mathcal{K}_{\mathbf{H}}), \\ \mathbf{x} &\sim \mathcal{GP}(\mathbf{0}, \mathcal{K}_{\mathbf{x}}), \\ \boldsymbol{\varepsilon} &\sim \mathcal{GP}(\mathbf{0}, \boldsymbol{\Lambda}^2)\end{aligned}$$

*independently for some kernel $\mathcal{K}_{\mathbf{x}}$, some kernel $\mathcal{K}_{\mathbf{H}}$, and some diagonal matrix $\boldsymbol{\Lambda}$. Then observations are generated by $\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon} = \mathbf{H} * \mathbf{x} + \boldsymbol{\varepsilon}$.*

In the remainder of this chapter we study the case that $\mathcal{K}_{\mathbf{x}}(\mathbf{t}_1, \mathbf{t}_2) = \delta(\mathbf{t}_1 - \mathbf{t}_2)\mathbf{I}$ —we let \mathbf{x} be white noise. The general case will be studied in Chapter 4.

Observe that $\mathbf{f} | \mathbf{H}$ is a linear combination of Gaussian processes. Hence $\mathbf{f} | \mathbf{H}$ is another Gaussian process, which thus can be identified by its mean function and kernel:

$$\begin{aligned}\mathbb{E}[\mathbf{f}(\mathbf{t}) | \mathbf{H}] &= \int_{\mathbb{R}^K} \mathbf{H}(\mathbf{t} - \boldsymbol{\tau}) \mathbb{E}[\mathbf{x}(\boldsymbol{\tau})] d\boldsymbol{\tau} = \mathbf{0}, \\ \mathcal{K}_{\mathbf{f} | \mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2) &= \mathbb{E}[\mathbf{f}(\mathbf{t}_1) \mathbf{f}^T(\mathbf{t}_2) | \mathbf{H}] \\ &= \int_{\mathbb{R}^K} \int_{\mathbb{R}^K} \mathbf{H}(\mathbf{t}_1 - \boldsymbol{\tau}_1) \underbrace{\mathbb{E}[\mathbf{x}(\boldsymbol{\tau}_1) \mathbf{x}^T(\boldsymbol{\tau}_2)]}_{\delta(\boldsymbol{\tau}_1 - \boldsymbol{\tau}_2)\mathbf{I}} \mathbf{H}^T(\mathbf{t}_2 - \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\ &= \int_{\mathbb{R}^K} \mathbf{H}[\boldsymbol{\tau} - (\mathbf{t}_2 - \mathbf{t}_1)] \mathbf{H}^T(\boldsymbol{\tau}) d\boldsymbol{\tau} \\ &= [R(\mathbf{H}) * \mathbf{H}^T](\mathbf{t}_2 - \mathbf{t}_1) \\ &= \mathcal{K}_{\mathbf{f} | \mathbf{H}}(\mathbf{t}_1 - \mathbf{t}_2)\end{aligned}\tag{3.1}$$

where R is the reversal function. We have established the following equivalent model:

Model 5 (Nonparametric Kernel Model (NKM)). *Draw*

$$\begin{aligned}\mathbf{H} &\sim \mathcal{GP}(\mathbf{0}, \mathcal{K}_{\mathbf{H}}), \\ \varepsilon &\sim \mathcal{GP}[\mathbf{0}, \delta(\mathbf{t}_1 - \mathbf{t}_2)\mathbf{\Lambda}^2]\end{aligned}$$

independently for some kernel $\mathcal{K}_{\mathbf{H}}$ and some diagonal matrix $\mathbf{\Lambda}$. Afterwards draw

$$\mathbf{f} | \mathbf{H} \sim \mathcal{GP}\{\mathbf{0}, [R(\mathbf{H}) * \mathbf{H}^T](\mathbf{t}_2 - \mathbf{t}_1)\}.$$

Then observations are generated by $\mathbf{y} = \mathbf{f} + \varepsilon$.

This equivalent formulation reveals that white noise excitation in Model 4 yields ordinary Gaussian process regression in which the kernel is modelled nonparametrically. We cannot model the kernel of \mathbf{f} simply with a sample from a Gaussian process—a kernel has to be positive semidefinite, which in general a sample from a Gaussian process is not. Thus Model 5 shows a way to nonparametrically model a positive-semidefinite matrix-valued function. We verify that $\mathcal{K}_{\mathbf{f}|\mathbf{H}} = R(\mathbf{H}) * \mathbf{H}^T$ is indeed positive semidefinite: Let $\mathbf{g} : \mathbb{R}^K \rightarrow \mathbb{R}^N$ be square integrable. Then

$$\begin{aligned}\int_{\mathbb{R}^{2K}} \mathbf{g}^T(\mathbf{t}_1)\mathcal{K}(\mathbf{t}_1, \mathbf{t}_2)\mathbf{g}(\mathbf{t}_2) d\mathbf{t}_1 d\mathbf{t}_2 &= \int_{\mathbb{R}^{3K}} \mathbf{g}^T(\mathbf{t}_1)\mathbf{H}[\boldsymbol{\tau} - (\mathbf{t}_2 - \mathbf{t}_1)]\mathbf{H}^T(\boldsymbol{\tau})\mathbf{g}(\mathbf{t}_2) d\mathbf{t}_1 d\mathbf{t}_2 d\boldsymbol{\tau} \\ &= \int_{\mathbb{R}^{3K}} \mathbf{g}^T(\mathbf{t}_1)\mathbf{H}(\mathbf{t}_1 - \boldsymbol{\tau}) d\mathbf{t}_1 \mathbf{H}^T(\mathbf{t}_2 - \boldsymbol{\tau})\mathbf{g}(\mathbf{t}_2) d\mathbf{t}_2 d\boldsymbol{\tau} \\ &= \int_{\mathbb{R}^K} \left\| \int_{\mathbb{R}^K} \mathbf{H}^T(\mathbf{t} - \boldsymbol{\tau})\mathbf{g}(\mathbf{t}) d\mathbf{t} \right\|^2 d\boldsymbol{\tau} \\ &\geq 0.\end{aligned}$$

Alternatively, Model 5 can be interpreted in the frequency domain. It is clear that $\mathbf{f} | \mathbf{H}$

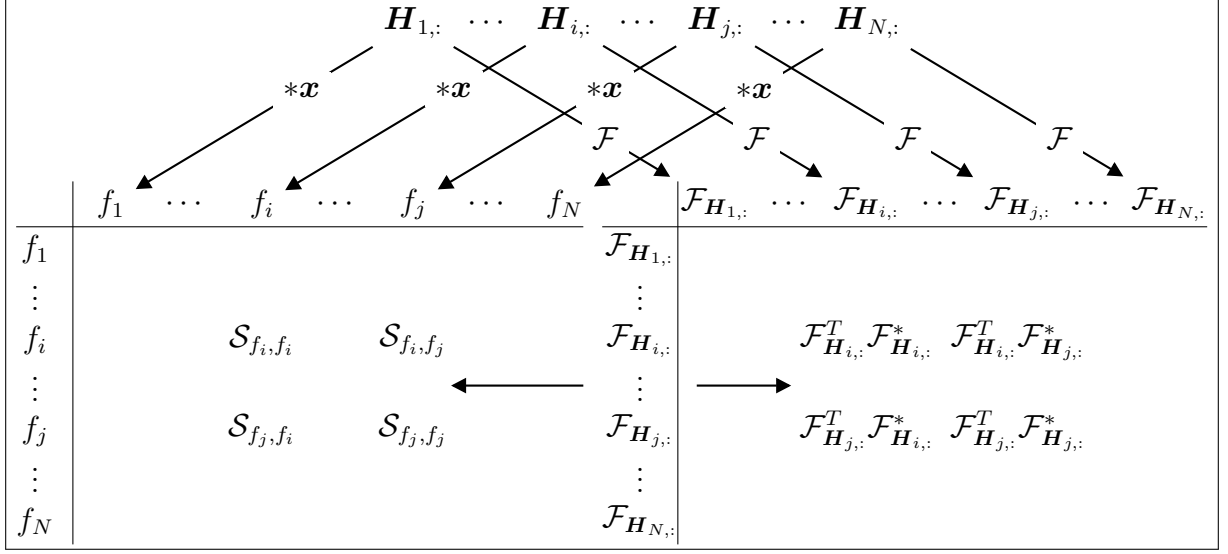


Figure 3.1: Generation of the outputs of in Model 5 and relationships between their cross-spectral densities

is wide-sense stationary; hence,

$$\begin{aligned}
 \mathcal{S}_{f_i, f_j | H}(\boldsymbol{\xi}) &= \mathcal{F}_{\boldsymbol{\tau}}\{\mathcal{K}_{f_i, f_j | H}(\mathbf{t} + \boldsymbol{\tau}, \mathbf{t})\}(\boldsymbol{\xi}) \\
 &= \sum_{m=1}^M \int_{\mathbb{R}^{2K}} H_{i,m}(\boldsymbol{\tau}' + \boldsymbol{\tau}) H_{j,m}(\boldsymbol{\tau}') \exp(-2\pi i \boldsymbol{\xi}^T \boldsymbol{\tau}) d\boldsymbol{\tau}' d\boldsymbol{\tau} \\
 &= \sum_{m=1}^M \int_{\mathbb{R}^{2K}} H_{i,m}(\boldsymbol{\tau}) \exp(-2\pi i \boldsymbol{\xi}^T \boldsymbol{\tau}) d\boldsymbol{\tau} [H_{j,m}(\boldsymbol{\tau}') \exp(-2\pi i \boldsymbol{\xi}^T \boldsymbol{\tau}')]^* d\boldsymbol{\tau}' \\
 &= \sum_{m=1}^M \mathcal{F}_{H_{i,m}}(\boldsymbol{\xi}) \mathcal{F}_{H_{j,m}}^*(\boldsymbol{\xi}) \tag{3.2} \\
 &= \mathcal{F}_{H_{i,:}}^T(\boldsymbol{\xi}) \mathcal{F}_{H_{j,:}}^*(\boldsymbol{\xi}). \tag{3.3}
 \end{aligned}$$

Since each $H_{i,j}$ is modelled nonparametrically, each $\mathcal{F}_{H_{i,j}}$ is modelled nonparametrically. Model 5 therefore models the cross-spectral densities \mathcal{S}_{f_i, f_j} nonparametrically by complex inner products between the stacked spectra $\mathcal{F}_{H_{i,:}}$ and $\mathcal{F}_{H_{j,:}}$ associated to f_i and f_j . This provides insight in how correlations between the outputs are induced; Figure 3.1 depicts in more detail how Model 5 generates outputs and how their cross-spectral densities relate.

Finally, Model 5 can be interpreted as a continuous-time moving average model. Equivalently, consider white noise excitation in Model 4. Then discretising the input space

results in a model of the form

$$\mathbf{f}_t = \sum_{\tau \in \mathbb{R}^K} \mathbf{H}_{t-\tau} \mathbf{x}_\tau$$

where \mathbf{x} is white noise. This is exactly a multi-output moving average model on a multidimensional input space. Hence Model 5 is a continuous-time Bayesian multi-output moving average model on a multidimensional input space.

3.3.1 Interpretation and Choice of the Kernel in the Nonparametric Kernel Model

In this section we resolve an apparent modelling issue concerning Model 5: it is not clear how to interpret and choose \mathcal{K}_H and thereby M —recall that $\mathcal{K}_H : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}^{NM \times NM}$.

First, we examine M . We already determined that $\mathbf{f} | \mathbf{H}$ is a Gaussian process with zero mean function. Thus Model 5 attains its full expressive power if it can generate any $\mathcal{K}_{f|H}$, or equivalently any $\mathcal{F}\{\mathcal{K}_{f|H}\} = \mathbf{S}_{f|H}$. To begin with,

$$\mathbf{S}_{f|H}^H = \mathcal{F}_\tau^* \{ \mathbb{E}^T [\mathbf{f}(t + \tau) \mathbf{f}^T(t)] \} = \mathcal{F}_\tau \{ \mathbb{E} [\mathbf{f}(t) \mathbf{f}^T(t - \tau)] \} = \mathbf{S}_{f|H}$$

shows that $\mathbf{S}_{f|H}$ is Hermitian. Therefore, by the Complex Spectral Theorem, generating any $\mathbf{S}_{f|H}$ is equivalent to generating any spectral decomposition. Now, utilising Equation (3.2) to express

$$\mathbf{S}_{f|H}(\boldsymbol{\xi}) = \sum_{m=1}^M \mathcal{F}_{H_{:,m}}(\boldsymbol{\xi}) \mathcal{F}_{H_{:,m}}^H(\boldsymbol{\xi})$$

shows that any spectral decomposition can be generated only if $M \geq N$. Thus $M = N$ is the minimal M for which Model 5 attains its full expressive power. Conversely, Model 5 has limited expressive power if $M < N$. More precisely, Equation (3.3) then shows that at each frequency $\boldsymbol{\xi}$ at most M outputs can be independent; f_i and f_j are independent if and only if $\mathcal{S}_{f_i, f_j | H}(\boldsymbol{\xi}) = \mathcal{F}_{H_{j,:}}^H \mathcal{F}_{H_{i,:}} = 0$, and at most M vectors $\mathcal{F}_{H_{i,:}}$ can simultaneously be orthogonal as they are elements of \mathbb{C}^M . We let $M = N$ in further development of Model 5.

Second, we examine \mathcal{K}_H . It turns out that we must be careful in choosing \mathcal{K}_H . Suppose the simplest case; that is, let the components of \mathbf{H} be independent and share an

exponentiated-quadratic kernel:

$$\mathcal{K}_{\mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2) = \sigma_h^2 \exp(-\gamma \|\mathbf{t}_1 - \mathbf{t}_2\|^2) \mathbf{I}.$$

Having already observed that Model 5 induces correlations between its outputs justifies modelling the components of \mathbf{H} independently. We call $\mathcal{K}_{\mathbf{H}}$ the diagonal multi-output exponentiated-quadratic kernel. Now, by equivalently considering white noise excitation in Model 4,

$$\begin{aligned} \mathbb{E}[f_i^2(\mathbf{t})] &= \int_{\mathbb{R}^{2K}} \mathbb{E}[\mathbf{H}_{i,:}^T(\mathbf{t} - \boldsymbol{\tau}_1) \mathbf{x}(\boldsymbol{\tau}_1) \mathbf{x}(\boldsymbol{\tau}_2)^T \mathbf{H}_{i,:}(\mathbf{t} - \boldsymbol{\tau}_2)] d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\ &= \int_{\mathbb{R}^{2K}} \text{tr} \mathbb{E}[\underbrace{\mathbf{H}_{i,:}(\mathbf{t} - \boldsymbol{\tau}_1) \mathbf{H}_{i,:}^T(\mathbf{t} - \boldsymbol{\tau}_2)}_{\mathcal{K}_{\mathbf{H}_{i,:}}[(\mathbf{t} - \boldsymbol{\tau}_1) - (\mathbf{t} - \boldsymbol{\tau}_2)]} \underbrace{\mathbf{x}(\boldsymbol{\tau}_1) \mathbf{x}^T(\boldsymbol{\tau}_2)}_{\delta(\boldsymbol{\tau}_1 - \boldsymbol{\tau}_2) \mathbf{I}}] d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\ &= \int_{\mathbb{R}^{2K}} \text{tr}(\sigma_h^2 \mathbf{I}) d\boldsymbol{\tau} \\ &= \infty. \end{aligned}$$

Intuitively, this means that the model prior assigns each f_i infinitely large error bars, or equivalently assigns each f_i infinite signal power. This is unlike real-world signals, which have finite power. We therefore let the components of \mathbf{H} have a decaying exponentiated-quadratic kernel [Tobar et al., 2015b] instead:

$$\mathcal{K}_{\mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2) = \sigma_h^2 \exp(-\alpha \|\mathbf{t}_1\|^2 - \alpha \|\mathbf{t}_2\|^2 - \gamma \|\mathbf{t}_1 - \mathbf{t}_2\|^2) \mathbf{I}.$$

We call $\mathcal{K}_{\mathbf{H}}$ the diagonal multi-output decaying exponentiated-quadratic kernel. Now, using the notation and identities from Appendix F with composite vector $[\mathbf{t}^T \quad \boldsymbol{\tau}^T]^T$,

$$\begin{aligned} \mathbb{E}[f_i^2(\mathbf{t})] &= \int_{\mathbb{R}^K} \text{tr}[\sigma_h^2 \exp(-2\alpha \|\mathbf{t} - \boldsymbol{\tau}\|^2) \mathbf{I}] d\boldsymbol{\tau} \\ &= I_{\boldsymbol{\tau}}[(K\sigma_h^2, 4 \begin{bmatrix} \alpha & -\alpha \\ -\alpha & \alpha \end{bmatrix})] \\ &= K\sigma_h^2 \frac{\pi^{K/2}}{(2\alpha)^{K/2}} \end{aligned}$$

and so our modelling issue is resolved. We decide to use the diagonal multi-output decaying exponentiated-quadratic kernel in further development of Model 5.

Note that the components of $\mathcal{K}_{\mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2)$ are small for $\|\mathbf{t}_1\|$ or $\|\mathbf{t}_2\|$ much larger than

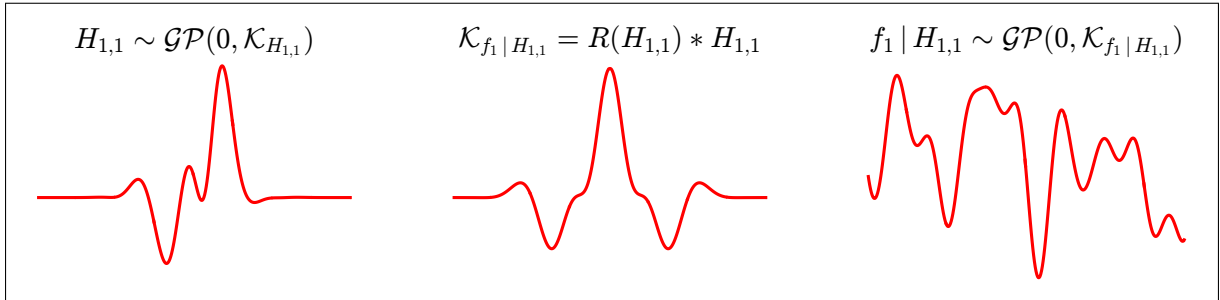


Figure 3.2: Generative process of Model 5 in the case that $N = M = K = 1$. Illustrates the interpretation that Model 5 corresponds ordinary Gaussian process regression in which the kernel is modelled nonparametrically.

$1/\sqrt{\alpha}$. Therefore for such \mathbf{t}_1 or \mathbf{t}_2 the posteriors of the components of \mathbf{H} are likely to be near zero; that is, $1/\sqrt{\alpha}$ is the effective extent of the components of \mathbf{H} . Similarly, the components of $\mathcal{K}_{\mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2)$ are small for $\|\mathbf{t}_1 - \mathbf{t}_2\|$ much larger than $1/\sqrt{\gamma}$; thus $1/\sqrt{\gamma}$ determines the effective length scale on which the components of \mathbf{H} vary.

We can let the effective extent and effective length scale be different for different dimensions of the input space. We then obtain

$$\mathcal{K}_{\mathbf{H}}(\mathbf{t}_1, \mathbf{t}_2) = \sigma_h^2 \exp[-\mathbf{t}_1^T \mathbf{A} \mathbf{t}_1 - \mathbf{t}_2^T \mathbf{A} \mathbf{t}_2 - (\mathbf{t}_1 - \mathbf{t}_2)^T \mathbf{\Gamma} (\mathbf{t}_1 - \mathbf{t}_2)] \mathbf{I}$$

for positive-semidefinite matrices \mathbf{A} and $\mathbf{\Gamma}$. Since \mathbf{A} and $\mathbf{\Gamma}$ reflect the relevance of each input variable, we call $\mathcal{K}_{\mathbf{H}}$ the diagonal multi-output decaying automatic relevance determination kernel [MacKay, 2002].

3.3.2 Illustrative Samples of the Nonparametric Kernel Model

Figure 3.2 shows the generative process of Model 5 in the simplest case that $M = N = 1$ and $K = 1$; that is, a sample is defined on a one-dimensional space—for example, time—and constitutes a single output. This figure illustrates the interpretation that Model 5 corresponds ordinary Gaussian process regression in which the kernel is modelled nonparametrically: First, $H_{1,1}$ is generated. Then, a kernel for $f_1 | H_{1,1}$ is constructed by computing $R(H_{1,1}) * H_{1,1}$. Finally, $f_1 | H_{1,1}$ is generated by drawing from a Gaussian process with $f_1 | H_{1,1}$'s generated kernel.

Furthermore, Figure 3.3 shows samples from Model 5 in the case that $M = N = 2$ and $K = 1$ —a sample now constitutes two outputs. First, observe that the sampled kernel

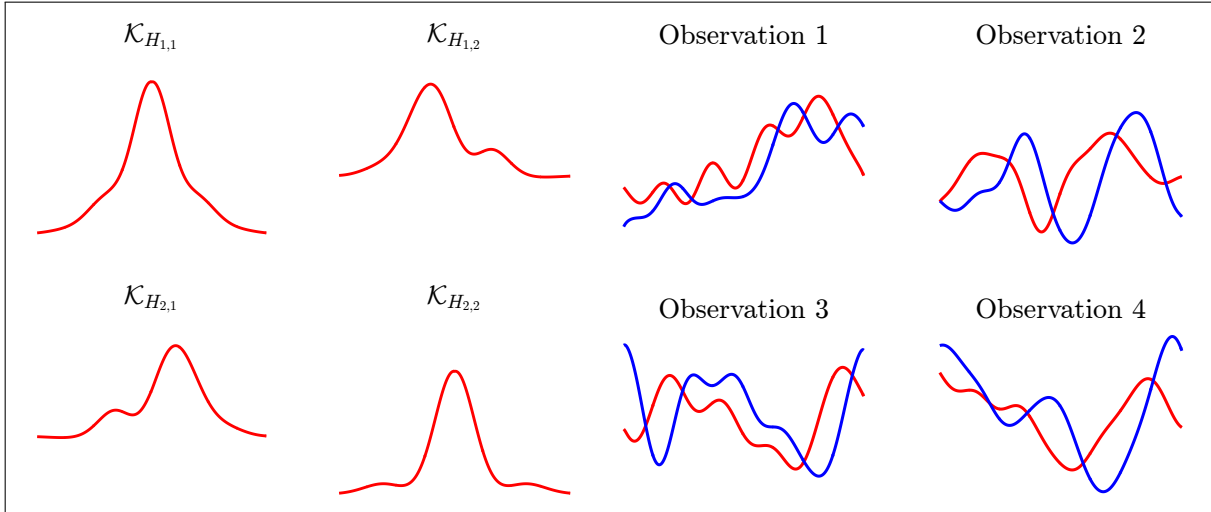


Figure 3.3: Samples from Model 5 in the case that $N = M = 2$ and $K = 1$. Shows the sampled kernel and four observations.

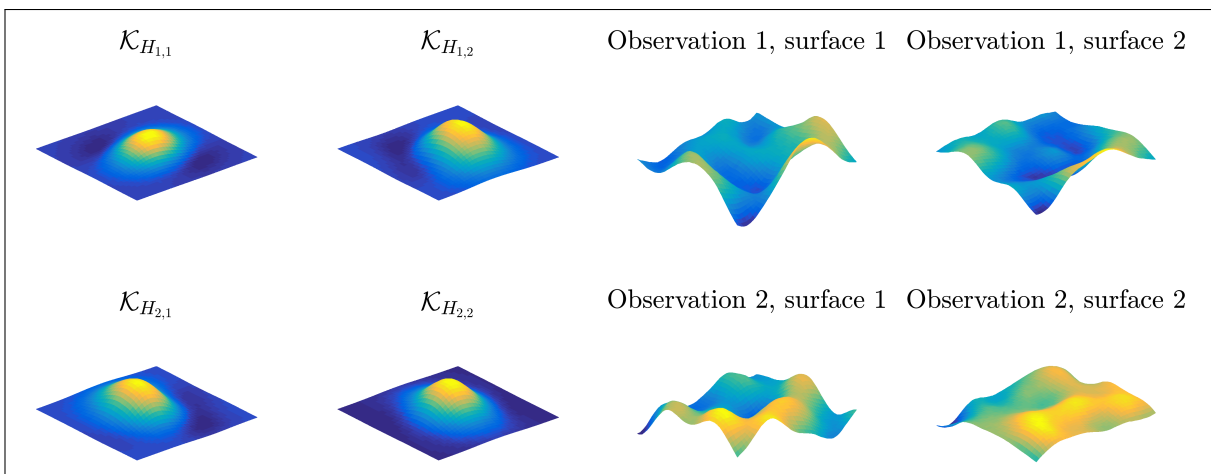


Figure 3.4: Samples from Model 5 in the case that $N = M = 2$ and $K = 2$. Shows the sampled kernel and two observations.

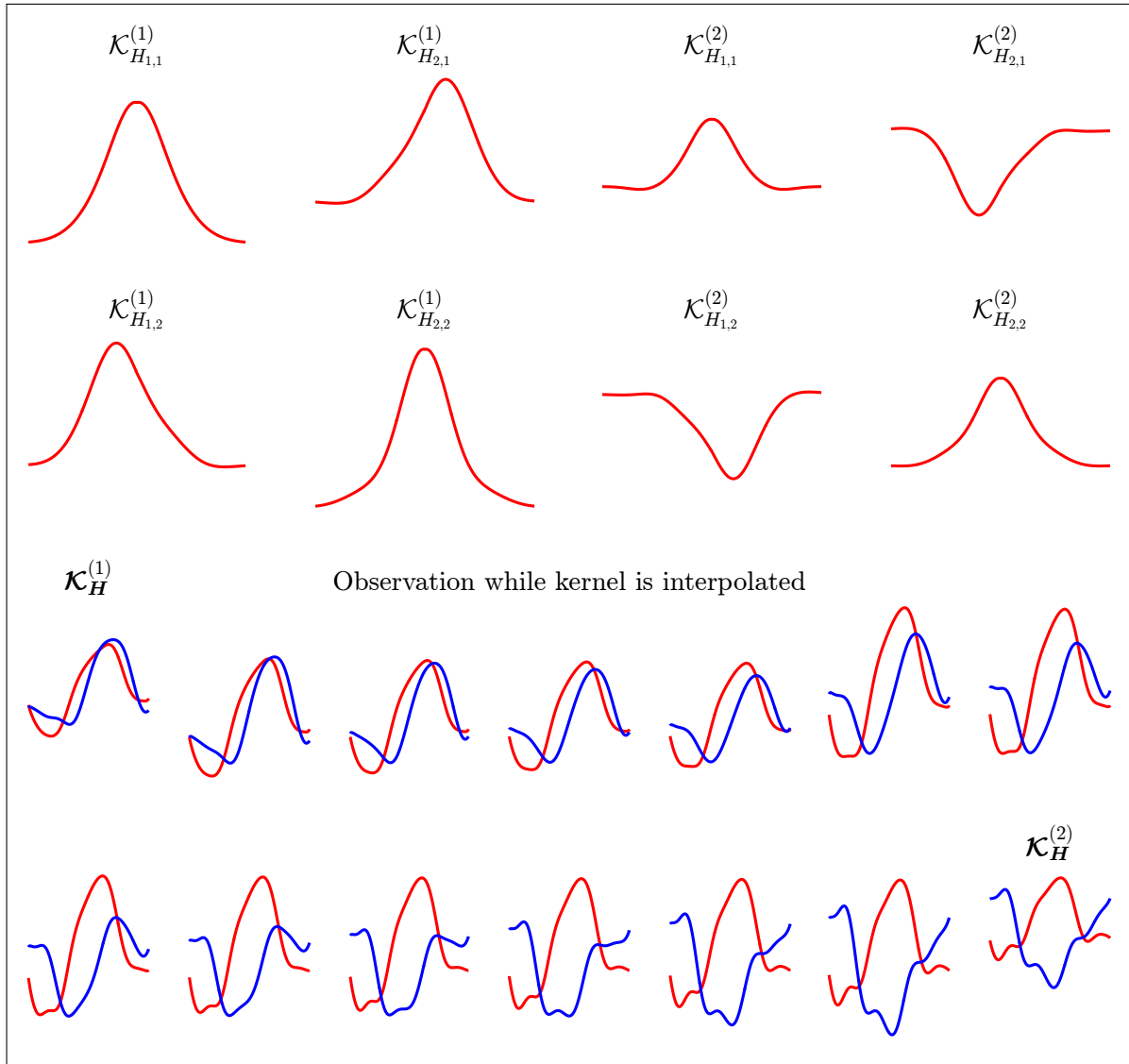


Figure 3.5: Interpolation between two kernels sampled from Model 5 in the case that $N = M = 2$ and $K = 1$. Shows the way an observation changes as its kernel is interpolated between the two kernels.

is indeed symmetric: $\mathcal{K}_{H_{1,1}}$ and $\mathcal{K}_{H_{2,2}}$ are symmetric and $\mathcal{K}_{H_{1,2}} = R(\mathcal{K}_{H_{2,1}})$. Second, observe that $\mathcal{K}_{H_{1,2}}$ and $\mathcal{K}_{H_{2,1}}$ attain reasonably peaked maxima at a nonzero lag. Hence we expect the outputs of a sample to be positively correlated at a nonzero lag; indeed, Figure 3.3 shows that the red lines essentially lead the blue lines. Thus Model 5 shows an certain ability to model multi-output time series. In the same way Model 5 is also able to model multi-output signals defined on multidimensional spaces; Figure 3.4 illustrates that the case $K = 2$ generates correlated surfaces.

Finally, Figure 3.5 illustrates that different samples from Model 5 yield different correlation structures between the outputs; in the case that $M = N = 2$ and $K = 1$, Figure 3.5 shows the way an observation changes as its kernel is interpolated between two kernels sampled from Model 5. Observe that kernel $\mathcal{K}_{\mathbf{H}}^{(1)}$ corresponds to positively correlated outputs, whereas kernel $\mathcal{K}_{\mathbf{H}}^{(2)}$ corresponds to negatively correlated outputs; indeed, as the kernel of the observation is interpolated, the outputs change from being positively correlated to being negatively correlated.

Appendix D was utilised in generating Figures 3.2 to 3.5.

3.3.3 Expressivity of White Noise Excitation

Thus far is not clear how Model 5's expressivity compares to that of Model 4. In this section we shed some light on this matter.

Model 4 shows that \mathbf{f} is generated by filtering \mathbf{x} with \mathbf{H} . This means that frequencies in \mathbf{x} are attenuated according to the spectrum of \mathbf{H} . In other words, \mathbf{f} consists only of frequencies that are also present in \mathbf{x} . Thus, if \mathbf{x} 's spectrum has effectively limited support, then \mathbf{f} is limited to only those frequencies. This makes choosing a white noise kernel for \mathbf{x} a great choice; a white noise process has a constant power spectrum and therefore contains all frequencies, meaning that \mathbf{f} can also contain all frequencies.

More concretely, we show that an instance $(\mathcal{K}_{\mathbf{x}^{(1)}}^{(1)}, \mathcal{K}_{\mathbf{H}^{(1)}}^{(1)}, \mathbf{\Lambda}^{(1)})$ of Model 4 can be expressed as an instance $(\mathcal{K}_{\mathbf{H}^{(2)}}^{(2)}, \mathbf{\Lambda}^{(2)})$ of Model 5 if $\mathcal{K}_{\mathbf{x}^{(1)}}^{(1)}(\mathbf{t}_1, \mathbf{t}_2) = [R(\mathbf{R}) * \mathbf{R}^T](\mathbf{t}_2 - \mathbf{t}_1)$ for some \mathbf{R} —we call \mathbf{R} the *root* of $\mathcal{K}_{\mathbf{x}^{(1)}}^{(1)}$. This is equivalent to showing that $(\mathcal{K}_{\mathbf{x}^{(1)}}^{(1)}, \mathcal{K}_{\mathbf{H}^{(1)}}^{(1)}, \mathbf{\Lambda}^{(1)})$ can be expressed as another instance $(\mathcal{K}_{\mathbf{x}^{(2)}}^{(2)}, \mathcal{K}_{\mathbf{H}^{(2)}}^{(2)}, \mathbf{\Lambda}^{(2)})$ of Model 4 where $\mathbf{x}^{(2)}$ is white noise.

Observe that $\mathbf{R} * \mathbf{x}^{(2)}$ and $\mathbf{H}^{(1)} * \mathbf{R}$ are linear combinations of Gaussian processes. Hence $\mathbf{R} * \mathbf{x}^{(2)}$ and $\mathbf{H}^{(1)} * \mathbf{R}$ are also Gaussian processes, which thus can be identified by their mean functions and kernels. It is readily verified that $\mathbf{R} * \mathbf{x}^{(2)}$ and $\mathbf{H}^{(1)} * \mathbf{R}$

have zero mean function. The kernel of $\mathbf{R} * \mathbf{x}^{(2)}$ is derived in a similar fashion as Equation (3.1), yielding $\mathcal{K}_{\mathbf{R} * \mathbf{x}^{(2)}}(\mathbf{t}_1, \mathbf{t}_2) = [R(\mathbf{R}) * \mathbf{R}^T](\mathbf{t}_2 - \mathbf{t}_1)$. Thus $\mathbf{x}^{(1)} \stackrel{d}{=} \mathbf{R} * \mathbf{x}^{(2)}$. Now let $\mathcal{K}_{\mathbf{H}^{(2)}} = \mathcal{K}_{\mathbf{H}^{(1)} * \mathbf{R}}$ and $\mathbf{\Lambda}^{(2)} = \mathbf{\Lambda}^{(1)}$. Then also $\mathbf{H}^{(1)} * \mathbf{R} \stackrel{d}{=} \mathbf{H}^{(2)}$ and $\boldsymbol{\varepsilon}^{(1)} \stackrel{d}{=} \boldsymbol{\varepsilon}^{(2)}$. Consequently, by associativity of the convolution operator,

$$\begin{aligned} \mathbf{f}^{(1)} &= \mathbf{H}^{(1)} * \mathbf{x}^{(1)} + \boldsymbol{\varepsilon}^{(1)} \\ &\stackrel{d}{=} \mathbf{H}^{(1)} * (\mathbf{R} * \mathbf{x}^{(2)}) + \boldsymbol{\varepsilon}^{(2)} \\ &= (\mathbf{H}^{(1)} * \mathbf{R}) * \mathbf{x}^{(2)} + \boldsymbol{\varepsilon}^{(2)} \\ &\stackrel{d}{=} \mathbf{H}^{(2)} * \mathbf{x}^{(2)} + \boldsymbol{\varepsilon}^{(2)} \\ &= \mathbf{f}^{(2)} \end{aligned}$$

so that $\mathbf{f}^{(1)} \stackrel{d}{=} \mathbf{f}^{(2)}$.

We have established that an instance of Model 4 can be expressed of an instance of Model 5 if \mathbf{x} 's kernel has a root and is thereby stationary. Appendix G shows that any kernel of exponentiated quadratic form has a root, that the diagonal multi-output exponentiated-quadratic kernel has a root and that the white noise kernel is its own root. This sheds some light on the class of kernels of \mathbf{x} for which Model 4 and Model 5 have equal expressivity.

We can more generally identify which kernels have roots in the case that $M = N = 1$. Suppose that \mathcal{K}_x is stationary and Schwartz—that is, rapidly decreasing. Then, by Bochner's Theorem, its Fourier transform $\mathcal{F}_{\mathcal{K}_x}$ exists and is real. Therefore $\mathcal{F}_{\mathcal{K}_x}^{1/2}$ is well defined. As \mathcal{K}_x is Schwartz, $\mathcal{F}_{\mathcal{K}_x}$ is also Schwartz and thereby absolutely integrable. Thus, by

$$\int_{\mathbb{R}^K} |\mathcal{F}_{\mathcal{K}_x}(\boldsymbol{\xi})| \, d\boldsymbol{\xi} = \int_{\mathbb{R}^K} |\mathcal{F}_{\mathcal{K}_x}^{1/2}(\boldsymbol{\xi})|^2 \, d\boldsymbol{\xi} < \infty,$$

$\mathcal{F}_{\mathcal{K}_x}^{1/2}$ is square integrable and so $\mathcal{F}^{-1}\{\mathcal{F}_{\mathcal{K}_x}^{1/2}\}$ exists. Let $r = \mathcal{F}^{-1}\{\mathcal{F}_{\mathcal{K}_x}^{1/2}\}$. Then

$$\begin{aligned} [R(r) * r](\boldsymbol{\tau}) &= \int_{\mathbb{R}^{3K}} \mathcal{F}_{\mathcal{K}_x}^{1/2}(\boldsymbol{\xi}_1) \exp[2\pi i \boldsymbol{\xi}_1^T (\boldsymbol{\tau}' - \boldsymbol{\tau})] \mathcal{F}_{\mathcal{K}_x}^{1/2}(\boldsymbol{\xi}_2) \exp(2\pi i \boldsymbol{\xi}_2^T \boldsymbol{\tau}') \, d\boldsymbol{\xi}_1 \, d\boldsymbol{\xi}_2 \, d\boldsymbol{\tau}' \\ &= \int_{\mathbb{R}^{3K}} \mathcal{F}_{\mathcal{K}_x}^{1/2}(\boldsymbol{\xi}_1) \mathcal{F}_{\mathcal{K}_x}^{1/2}(\boldsymbol{\xi}_2) \exp(2\pi i \boldsymbol{\xi}_2^T \boldsymbol{\tau}) \underbrace{\exp[2\pi i (\boldsymbol{\xi}_1 + \boldsymbol{\xi}_2)^T \boldsymbol{\tau}']}_{\delta(-\boldsymbol{\xi}_1 - \boldsymbol{\xi}_2)} \, d\boldsymbol{\xi}_1 \, d\boldsymbol{\xi}_2 \, d\boldsymbol{\tau}' \\ &= \int_{\mathbb{R}^{3K}} \mathcal{F}_{\mathcal{K}_x}(\boldsymbol{\xi}) \exp(2\pi i \boldsymbol{\xi}^T \boldsymbol{\tau}) \, d\boldsymbol{\xi} \\ &= \mathcal{K}_x(\boldsymbol{\tau}) \end{aligned}$$

and so \mathcal{K}_x has a root.

In summary, white noise excitation in Model 4 yields great expressivity because the spectrum of white noise has infinite support. Concretely, an instance of Model 4 can be expressed as an instance of Model 5 if \boldsymbol{x} 's kernel has a root. In the case that $M = N = 1$, a kernel has a root if it is stationary and rapidly decreasing.

3.4 The Approximate Kernel Model

Model 5 is troublesome from a computational perspective. Namely, $\mathcal{K}_{f|\mathbf{H}} = R(\mathbf{H}) * \mathbf{H}^T$ is a convolution between two matrix-valued stochastic processes and it is not clear how such a quantity can be computed. We therefore develop a model that is approximate to Model 5, but whose computation is tractable.

Recall that \mathbf{H} is smooth and that its support is effectively limited due to use of the decaying exponentiated-quadratic kernel. In that case, as demonstrated by Minka [2000], we can numerically approximate the integral $R(\mathbf{H}) * \mathbf{H}^T$ by its expectation conditioned on some finite number of observations $\mathbf{H}(\mathbf{T})$ where $\mathbf{T} \in \mathbb{R}^{T \times K}$ —we call these observations *inducing points*. Symbolically,

$$\mathcal{K}_{f|\mathbf{H}} = R(\mathbf{H}) * \mathbf{H}^T \approx \mathbb{E}[R(\mathbf{H}) * \mathbf{H}^T | \mathbf{H}(\mathbf{T})] = \tilde{\mathcal{K}}_{f|\mathbf{H}(\mathbf{T})}.$$

Note that $\tilde{\mathcal{K}}_{f|\mathbf{H}(\mathbf{T})}$ numerically approximates $\mathcal{K}_{f|\mathbf{H}}$ and at the same time maintains knowledge of its own uncertainty; in this sense $\tilde{\mathcal{K}}_{f|\mathbf{H}(\mathbf{T})}$ is called a *Bayesian numerical approximation* of $\mathcal{K}_{f|\mathbf{H}}$ [Tobar et al., 2015b]. Appendix H shows that the analytical

form of $\tilde{\mathcal{K}}_{f|H(\mathbf{T})}$ is given by

$$\begin{aligned} \tilde{\mathcal{K}}_{f_i, f_j | H(\mathbf{T})} = \sum_{k=1}^M \left\{ \underbrace{\mathbb{1}(i-j) I^{(1, H_{i,k})}}_{\text{prior}} - \underbrace{\mathbb{1}(i-j) \text{tr}(\mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{i,k})})}_{\text{conditioning on inducing points}} \right. \\ \left. + \underbrace{H_{i,k}^T(\mathbf{T}) \mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{j,k})} \mathbf{K}_{H_{j,k}(\mathbf{T})}^{-1} H_{j,k}(\mathbf{T})}_{\text{learned through inducing points}} \right\}. \end{aligned}$$

We clearly distinguish the prior covariance structure, the covariance structure due to conditioning on the inducing points and the covariance structure learned through the inducing points.

We have derived the following approximation of Model 5:

Model 6 (Approximate Kernel Model (AKM)). *Draw*

$$\begin{aligned} \mathbf{H}(\mathbf{T}) &\sim \mathcal{N}[\mathbf{0}, \mathcal{K}_H(\mathbf{T}, \mathbf{T})], \\ \boldsymbol{\varepsilon} &\sim \mathcal{GP}[\mathbf{0}, \boldsymbol{\Lambda}^2 \delta(\mathbf{t}_1 - \mathbf{t}_2)] \end{aligned}$$

independently for some diagonal kernel \mathcal{K}_H and some diagonal matrix $\boldsymbol{\Lambda}$. Afterwards let $\tilde{\mathcal{K}}_{f|H(\mathbf{T})}$ be such that

$$\begin{aligned} \mathcal{K}_{f_i, f_j | H(\mathbf{T})} = \sum_{k=1}^M \left\{ \mathbb{1}(i-j) [I^{(1, H_{i,k})} - \text{tr}(\mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{i,k})})] \right. \\ \left. + H_{i,k}^T(\mathbf{T}) \mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{j,k})} \mathbf{K}_{H_{j,k}(\mathbf{T})}^{-1} H_{j,k}(\mathbf{T}) \right\}. \end{aligned}$$

Finally draw

$$\mathbf{f} | \mathbf{H}(\mathbf{T}) \sim \mathcal{GP}(\mathbf{0}, \tilde{\mathcal{K}}_{f|H(\mathbf{T})}).$$

Then observations are generated by $\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon}$.

Observe that Model 6 is ordinary Gaussian process regression with a specific kernel and a prior over the hyperparameters. Thus Model 6 can take advantage of existing literature on ordinary Gaussian process regression.

Furthermore, the fact that Model 6 approximates Model 5 implies that the prior over the

hyperparameters specified by Model 6 is one that is sensible to use. This is an advantage of Model 6 over most other multi-output kernels; priors over hyperparameters are usually chosen without principled motivation.

3.4.1 Interpretation of the Kernel Approximation

We investigate the kernel approximation in the case that $M = N = 1$. Denote $f = f_1$, $h = H_{1,1}$ and $h(\mathbf{T}) = \mathbf{h}$. Let $\mathcal{K}_{\tilde{h}|\mathbf{A}}$ be the posterior kernel of the sparse approximation of h according to [Titsias, 2009] where \mathbf{h} are the approximation's inducing points and \mathbf{A} is their posterior covariance matrix. We compare $\tilde{\mathcal{K}}_{f|\mathbf{h}}$ to $\mathcal{K}_{\tilde{h}|\mathbf{A}}$:

$$\begin{aligned}\tilde{\mathcal{K}}_{f|\mathbf{h}}(\mathbf{t}_1, \mathbf{t}_2) &= I^{(1,h)}(\mathbf{t}_1, \mathbf{t}_2) - \text{tr}[\mathbf{K}_{\mathbf{h}}^{-1} \mathbf{I}^{(2,h,h)}(\mathbf{t}_1, \mathbf{t}_2)] \\ &\quad + \text{tr}[\mathbf{K}_{\mathbf{h}}^{-1} \mathbf{I}^{(2,h,h)}(\mathbf{t}_1, \mathbf{t}_2) \mathbf{K}_{\mathbf{h}}^{-1} \mathbf{h} \mathbf{h}^T], \\ \mathcal{K}_{\tilde{h}|\mathbf{A}}(\mathbf{t}_1, \mathbf{t}_2) &= \mathcal{K}_h(\mathbf{t}_1, \mathbf{t}_2) - \text{tr}[\mathbf{K}_{\mathbf{h}}^{-1} \mathcal{K}_h(\mathbf{T}, \mathbf{t}_2) \mathcal{K}_h(\mathbf{t}_1, \mathbf{T})] \\ &\quad + \text{tr}[\mathbf{K}_{\mathbf{h}}^{-1} \mathcal{K}_h(\mathbf{T}, \mathbf{t}_2) \mathcal{K}_h(\mathbf{t}_1, \mathbf{T}) \mathbf{K}_{\mathbf{h}}^{-1} \mathbf{A}].\end{aligned}$$

Thus, by Equations (H.1) and (H.2), we establish that

$$\tilde{\mathcal{K}}_{f|\mathbf{h}}(\mathbf{t}_1, \mathbf{t}_2) = \int_{\mathbb{R}^K} \mathcal{K}_{\tilde{h}|\mathbf{A}=\mathbf{h}\mathbf{h}^T}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{t}_2 - \boldsymbol{\tau}) \, \text{d}\boldsymbol{\tau}.$$

This means that the approximate integration performed by $\tilde{\mathcal{K}}_{f|\mathbf{h}}$ corresponds to first sparsely approximating h and then exactly integrating its posterior kernel.

3.4.2 The Case of the Diagonal Multi-Output Decaying Exponentiated-Quadratic Kernel

Let \mathcal{K}_H be a diagonal multi-output decaying exponentiated-quadratic kernel. It then turns out that $I^{(1,H_{i,k})}$, $\mathbf{I}^{(2,H_{i,k},H_{j,k})}$ and their Fourier transforms admit simple forms. We now compute these simple forms, which in Section 3.5 will reveal Model 6's connection

to existing work. To begin with, Appendices I.4.3 and I.4.5 show that

$$\begin{aligned}
I^{(1, H_{i,k})}(\mathbf{t}_1, \mathbf{t}_2) &= \underbrace{\sigma_h^2 \frac{\pi^{K/2}}{(2\alpha)^{K/2}}}_{C^{(1)}} \exp \left[-\frac{1}{2} \underbrace{(\alpha + 2\gamma)}_{l^{-(1)}} \|\mathbf{t}_1 - \mathbf{t}_2\|^2 \right], \\
&= C^{(1)} \exp \left(-\frac{1}{2} l^{-(1)} \|\mathbf{t}_1 - \mathbf{t}_2\|^2 \right) \\
&= I^{(1)}(\mathbf{t}_1 - \mathbf{t}_2), \\
I_{m,n}^{(2, H_{i,k}, H_{j,k})}(\mathbf{t}_1, \mathbf{t}_2) &= \underbrace{\sigma_h^4 \frac{\pi^{K/2}}{(2\alpha + 2\gamma)^{K/2}} \exp \left[-\frac{1}{2} \frac{(\alpha + \gamma)^2 - \gamma^2}{\alpha + \gamma} (\|\mathbf{T}_{m,:} - \mathbf{T}_{n,:}\|^2 \right.}_{C_{m,n}^{(2)}} \\
&\quad \left. - \|\mathbf{T}_{m,:} + \mathbf{T}_{n,:}\|^2) \right]} \\
&\quad \exp \left[-\frac{1}{2} \underbrace{(\alpha + \gamma)}_{l^{-(2)}} \left\| (\mathbf{t}_1 - \mathbf{t}_2) - \underbrace{\frac{\gamma}{\alpha + \gamma} (\mathbf{T}_{m,:} - \mathbf{T}_{n,:})}_{\boldsymbol{\mu}_{m,n}^{(2)}} \right\|^2 \right] \\
&= C_{m,n}^{(2)} \exp \left[-\frac{1}{2} l^{-(2)} \|(\mathbf{t}_1 - \mathbf{t}_2) - \boldsymbol{\mu}_{m,n}^{(2)}\|^2 \right] \\
&= I_{m,n}^{(2)}(\mathbf{t}_1 - \mathbf{t}_2).
\end{aligned}$$

Therefore

$$\mathcal{F}_{\mathbf{t}_1 - \mathbf{t}_2} \{I^{(1)}\}(\mathbf{f}) = C^{(1)} \sqrt{2\pi l^{(1)}} \exp(-2\pi^2 l^{(1)} \|\mathbf{f}\|^2), \quad (3.4)$$

$$\mathcal{F}_{\mathbf{t}_1 - \mathbf{t}_2} \{I_{m,n}^{(2)}\}(\mathbf{f}) = C_{m,n}^{(2)} \sqrt{2\pi l^{(2)}} \exp(-2\pi^2 l^{(2)} \|\mathbf{f}\|^2 + 2\pi i \boldsymbol{\mu}_{m,n}^T \mathbf{f}). \quad (3.5)$$

Now, consider the case that $i = j$. First, all $\mathcal{K}_{H_{i,j}}$ are equal and so all $\mathbf{K}_{H_{i,j}(\mathbf{T})} = \mathbf{K}$; thus the kernel approximation simplifies to

$$\mathcal{K}_{f_i, f_i | \mathbf{H}(\mathbf{T})} = M [I^{(1)} - \text{tr}(\mathbf{K}^{-1} \mathbf{I}^{(2)})] + \sum_{k=1}^M H_{i,k}^T(\mathbf{T}) \mathbf{K}^{-1} \mathbf{I}^{(2)} \mathbf{K}^{-1} H_{i,k}(\mathbf{T}), \quad (3.6)$$

which shows that we can equivalently let $\mathbf{I}^{(2)} \leftarrow (\mathbf{I}^{(2)} + \mathbf{I}^{(2)T})/2$. Second, note that $C_{m,n}^{(2)} = C_{n,m}^{(2)}$ and $-\boldsymbol{\mu}_{n,m}^{(2)} = \boldsymbol{\mu}_{m,n}^{(2)}$; hence $\mathcal{F}\{I_{m,n}^{(2)}\} = \mathcal{F}^*\{I_{n,m}^{(2)}\}$. Therefore, if $i = j$, then

equivalently

$$\begin{aligned}
\mathcal{F}_{t_1-t_2}\{I_{m,n}^{(2)}\}(\mathbf{f}) &= \frac{1}{2}[\mathcal{F}_{t_1-t_2}\{I_{m,n}^{(2)}\}(\mathbf{f}) + \mathcal{F}_{t_1-t_2}\{I_{n,m}^{(2)}\}(\mathbf{f})] \\
&= \frac{1}{2}[\mathcal{F}_{t_1-t_2}\{I_{m,n}^{(2)}\}(\mathbf{f}) + \mathcal{F}_{t_1-t_2}^*\{I_{m,n}^{(2)}\}(\mathbf{f})] \\
&= C_{m,n}^{(2)}\sqrt{2\pi l^{(2)}}\exp(-2\pi^2 l^{(2)}\|\mathbf{f}\|^2)\cos(2\pi i\boldsymbol{\mu}_{m,n}^T\mathbf{f}). \tag{3.7}
\end{aligned}$$

3.5 Related Work

Model 5 and its approximation Model 6 are closely related to recent work.

First, Model 5 is a generalisation of the Gaussian Process Convolution Model (GPCM) [Tobar et al., 2015b]; the case $N = M = K = 1$ recovers their model exactly.

Second, Equations (3.4), (3.6) and (3.7) show that the power spectral densities of the kernel in Model 6 take the form of spectral mixture kernels (SMKs) [Wilson and Adams, 2013]. In other words, the diagonal entries of the kernel in Model 6 form *Fourier pairs* with spectral mixture kernels. Therefore the kernel in Model 6 is a *dual* multi-output generalisation of the spectral mixture kernel. Furthermore, recall that Model 6 approximates Model 5 more accurately as the number of inducing points increases. This reveals the kernel in Model 5 as a multi-output generalisation of the spectral mixture kernel with an infinite number of components.

Third, Equations (3.4) to (3.6) show that the kernel in Model 6 is a dual formulation of the cross-spectral mixture kernel (CSMK) [Ulrich et al., 2015]. Similar to the spectral mixture kernel, this reveals the kernel in Model 5 as a cross-spectral mixture kernel with an infinite number of components. Furthermore, the connection between Model 6 and the cross-spectral mixture kernel suggests that results in [Ulrich et al., 2015] might carry over to Model 6; indeed, inspired by [Ulrich et al., 2015], Appendix E presents an approximation of stationary multi-output kernel matrices that can be used to efficiently perform inference in Model 6 if $\mathbf{H}(\mathbf{T})$ are instead treated as hyperparameters.

Figure 3.6 summarises how Model 5 and Model 6 relate to current literature on flexible kernel models.

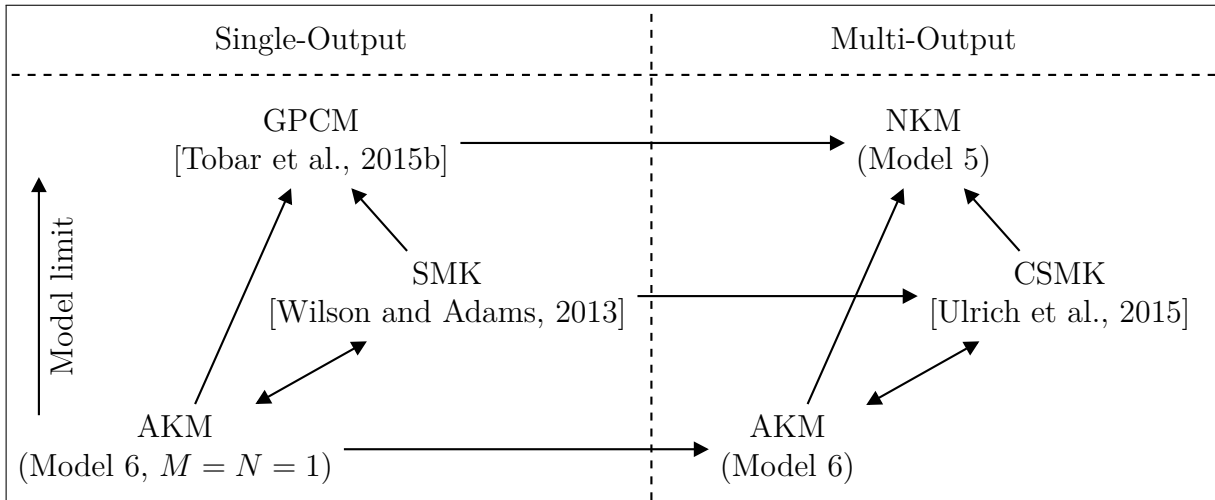


Figure 3.6: Relationship of Model 5 and Model 6 to current literature on flexible kernel models. Single-headed arrows indicate generalisation and double-headed arrows indicate duality. “Model limit” refers to taking the model’s number of components to infinity.

3.6 Inference in the Nonparametric Kernel

Model

Performing inference in Model 5 directly is troublesome, because \mathbf{H} parametrises the kernel of Model 5’s likelihood $p(\mathbf{y} | \mathbf{H})$ which thus is intractable. Instead, equivalently consider white noise excitation in Model 4. Model 4’s likelihood $p(\mathbf{y} | \mathbf{H}, \mathbf{x})$ is Gaussian with mean function $\mathbf{f} = \mathbf{H} * \mathbf{x}$. As we show shortly, this likelihood is manageable.

We perform inference in Model 5 by inferring \mathbf{H} and \mathbf{x} in Model 4 where \mathbf{x} is white noise. This immediately raises a concern: $\mathbf{f} = \mathbf{H} * \mathbf{x}$ requires knowledge of the whole processes \mathbf{H} and \mathbf{x} , but \mathbf{x} cannot wholly be learned because $\mathbf{x}(\mathbf{t}_1)$ and $\mathbf{x}(\mathbf{t}_2)$ are independent for $\mathbf{t}_1 \neq \mathbf{t}_2$; that is, we would have to learn $\mathbf{x}(\mathbf{t})$ for every \mathbf{t} separately, but a computer can only store $\mathbf{x}(\mathbf{t})$ for finitely many \mathbf{t} .

To resolve this issue, consider another instance of Model 4 whose excitation is denoted by $\tilde{\mathbf{x}}$. Let $\tilde{\mathbf{x}}$ have a diagonal multi-output exponentiated-quadratic kernel. In that case $\tilde{\mathbf{x}}$ is smooth and can therefore be learned by conditioning on $\tilde{\mathbf{x}}(\mathbf{t})$ for finitely many \mathbf{t} [Titsias, 2009]. Thus we should be able to perform inference in this new instance. Now, Appendix G shows that $\tilde{\mathbf{x}}$ ’s kernel has a root \mathbf{R} . Therefore, by Section 3.3.3, Model 5 can be expressed in terms of this new instance where $\tilde{\mathbf{x}} = \mathbf{R} * \mathbf{x}$. Crucially, the fact that we can learn this new instance then implies that we can learn Model 5 by learning

$\mathbf{R} * \mathbf{x}$ instead of \mathbf{x} .

Learning $\mathbf{R} * \mathbf{x}$ instead of \mathbf{x} has a number of equivalent interpretations. First, $\mathbf{R} * \mathbf{x}$ represents a smoothed version of \mathbf{x} where each point is a linear combination of the points of \mathbf{x} . Thus learning one point on $\mathbf{R} * \mathbf{x}$ induces knowledge about the whole process \mathbf{x} . Second, \mathbf{R} acts as a filter, which means that $\mathbf{R} * \mathbf{x}$ represents a band-limited version of \mathbf{x} . Hence we restrict ourselves to only learning frequencies in the band of \mathbf{R} ; this is an easier problem than learning all frequencies.

The latter interpretation raises another concern; if \mathbf{H} admits frequencies not in the band of \mathbf{R} , then those frequencies cannot be learned. Thus \mathbf{R} should be chosen carefully; in the case that $M = N = K = 1$, Tobar et al. [2015a] consider how \mathbf{R} can be designed to promote training performance. We decide to simply use

$$\mathbf{R}(\mathbf{t}) = \exp(-\omega \|\mathbf{t}\|^2) \mathbf{I}.$$

This choice should enable us to learn \mathbf{x} ; namely, Appendix G shows that this \mathbf{R} is a root of a diagonal multi-output exponentiated-quadratic kernel, and we previously determined that a process with a diagonal multi-output exponentiated-quadratic kernel can be learned.

Concretely, we perform approximate inference in Model 5 via a structured mean field approximation in which we sparsely approximate \mathbf{H} and $\mathbf{R} * \mathbf{x}$ through inducing points [Titsias, 2009]. Let \mathbf{Y} denote Y observations. Consider the family \mathcal{Q} of distributions of the form

$$q(\mathbf{H}, \mathbf{x}, \underbrace{\mathbf{u}_{H_{1,1}}, \dots, \mathbf{u}_{H_{N,M}}}_{\mathbf{U}_H}, \underbrace{\mathbf{u}_{\tilde{x}_1}, \dots, \mathbf{u}_{\tilde{x}_M}}_{\mathbf{U}_{\tilde{x}}}) = p(\mathbf{H}, \mathbf{x} \mid \mathbf{U}_H, \mathbf{U}_{\tilde{x}}) \prod_{i=1, j=1}^{N, M} q(\mathbf{u}_{H_{i,j}}) \prod_{j=1}^M q(\mathbf{u}_{\tilde{x}_j})$$

where $\mathbf{u}_{H_{i,j}} = H_{i,j}(\mathbf{T}_{H_{i,j}})$ and $\mathbf{u}_{\tilde{x}_j} = \tilde{x}_j(\mathbf{T}_{\tilde{x}_j})$ are inducing points for respectively the processes $H_{i,j}$ and $\tilde{x}_j = \exp(-\omega \|\mathbf{t}\|^2) * x_j$. Let

$$\begin{aligned} q(\mathbf{u}_{H_{i,j}}) &= \mathcal{N}(\mathbf{u}_{H_{i,j}}; \boldsymbol{\mu}_{H_{i,j}}, \boldsymbol{\Sigma}_{H_{i,j}}), \\ q(\mathbf{u}_{\tilde{x}_j}) &= \mathcal{N}(\mathbf{u}_{\tilde{x}_j}; \boldsymbol{\mu}_{\tilde{x}_j}, \boldsymbol{\Sigma}_{\tilde{x}_j}). \end{aligned}$$

We then perform approximate inference via

$$\begin{aligned}
& p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}} | \mathbf{Y}) \\
&= \operatorname{argmin}_q D_{KL}[q(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}}) \| p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}} | \mathbf{Y})] \\
&\approx \operatorname{argmin}_{q \in \mathcal{Q}} D_{KL}[q(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}}) \| p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}} | \mathbf{Y})] \\
&= \operatorname{argmax}_{q \in \mathcal{Q}} \{ \log p(\mathbf{Y}) - D_{KL}[q(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}}) \| p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}} | \mathbf{Y})] \} \\
&= \operatorname{argmax}_{q \in \mathcal{Q}} \mathcal{F}(q).
\end{aligned}$$

Since $D_{KL}(\cdot \| \cdot) \geq 0$, $\mathcal{F}(q)$ —also called the *variational free energy*, or *free energy* in short—is a lower bound on the marginal likelihood; hence performing inference is equivalent to maximising this lower bound. Now, recall that \mathbf{H} 's diagonal multi-output decaying exponentiated-quadratic kernel implies independence between different $H_{i,j}$'s; as a result,

$$p(\mathbf{H}, \mathbf{x} | \mathbf{U}_H, \mathbf{U}_{\tilde{x}}) = p(\mathbf{H} | \mathbf{U}_H) p(\mathbf{x} | \mathbf{U}_{\tilde{x}}) = \prod_{i=1, j=1}^{N, M} p(H_{i,j} | \mathbf{u}_{H_{i,j}}) \prod_{j=1}^M p(x_j | \mathbf{u}_{\tilde{x}_j}).$$

Then

$$\begin{aligned}
\mathcal{F}(q) &= \mathbb{E}_q[\log p(\mathbf{Y})] + \mathbb{E}_q \left[\log \frac{p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}} | \mathbf{Y})}{q(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}})} \right] \\
&= \mathbb{E}_q \left[\log \frac{p(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}}, \mathbf{Y})}{q(\mathbf{H}, \mathbf{x}, \mathbf{U}_H, \mathbf{U}_{\tilde{x}})} \right] \\
&= \mathbb{E}_q \left[\log \frac{p(\mathbf{Y} | \mathbf{H}, \mathbf{x}) p(\mathbf{H} | \mathbf{U}_H) p(\mathbf{x} | \mathbf{U}_{\tilde{x}}) \prod_{i=1, j=1}^{N, M} p(\mathbf{u}_{H_{i,j}}) \prod_{j=1}^M p(\mathbf{u}_{\tilde{x}_j})}{p(\mathbf{H} | \mathbf{U}_H) p(\mathbf{x} | \mathbf{U}_{\tilde{x}}) \prod_{i=1, j=1}^{N, M} q(\mathbf{u}_{H_{i,j}}) \prod_{j=1}^M q(\mathbf{u}_{\tilde{x}_j})} \right] \\
&= \mathbb{E}_q[\log p(\mathbf{Y} | \mathbf{H}, \mathbf{x})] - \sum_{i=1, j=1}^{N, M} D_{KL}[q(\mathbf{u}_{H_{i,j}}) \| p(\mathbf{u}_{H_{i,j}})] - \sum_{j=1}^M D_{KL}[q(\mathbf{u}_{\tilde{x}_j}) \| p(\mathbf{u}_{\tilde{x}_j})] \\
&= -\frac{1}{2} \log[(2\pi)^N |\mathbf{\Lambda}|^2] - \underbrace{\frac{1}{2} \sum_{i=1}^Y \mathbb{E}_q \{ \|\mathbf{\Lambda}^{-1}[\mathbf{Y}_{i,:} - \mathbf{y}(\mathbf{T}_{i,:})]\|^2 \}}_{\text{data reconstruction cost}} \\
&\quad - \underbrace{\sum_{i=1, j=1}^{N, M} D_{KL}[q(\mathbf{u}_{H_{i,j}}) \| p(\mathbf{u}_{H_{i,j}})] - \sum_{j=1}^M D_{KL}[q(\mathbf{u}_{\tilde{x}_j}) \| p(\mathbf{u}_{\tilde{x}_j})]}_{\text{divergence from model}}.
\end{aligned}$$

We observe that the variational free energy can be interpreted as negative a cost of reconstructing the data regularised by a divergence from the model prior; thus, maximisation of the free energy seeks an explanation of the data that is compatible with the model prior.

Appendix I derives an analytical expression for the variational free energy. Appendix D discusses some useful techniques concerning implementation of the free energy.

Furthermore, we consider the asymptotic time complexity of computing the free energy. Denote the number of inducing points for the processes $H_{i,j}$ and \tilde{x}_j by $T_{H_{i,j}}$ and $T_{\tilde{x}_j}$ respectively. Let $T_H = \max_{i,j} T_{H_{i,j}}$ and $T_{\tilde{x}} = \max_j T_{\tilde{x}_j}$. Then Appendix I.5 shows that the time complexity of computing the free energy is given by

$$\mathcal{O}[NMT_H^3 + MT_{\tilde{x}}^3 + YNM(T_H^2 T_{\tilde{x}} + KT_H^2 + KT_H T_{\tilde{x}} + KT_{\tilde{x}}^2 + T_H T_{\tilde{x}}^2) + YNM^2].$$

Recall that \mathbf{H} has effectively limited support due to use of the decaying exponentiated-quadratic kernel; thus likely $T_{\tilde{x}} \gg T_H$. Also, likely $T_{\tilde{x}} \gg M$. In that case the asymptotic time complexity of computing the free energy simplifies to $\mathcal{O}[MT_{\tilde{x}}^3 + YNM(K + T_H)T_{\tilde{x}}^2]$.

Finally, we discuss initialisation of the parameters of the variational free energy—that is, initialisation of $\boldsymbol{\mu}_{H_{i,j}}$, $\boldsymbol{\Sigma}_{H_{i,j}}$, $\boldsymbol{\mu}_{\tilde{x}_j}$ and $\boldsymbol{\Sigma}_{\tilde{x}_j}$. Optimisation of the free energy is usually performed via gradient-based methods. This means that the optimiser is likely to become trapped in a local minimum if the parameters of the free energy are initialised far from their optimal values. Hence proper initialisation of the parameters is of crucial importance.

We propose the following initialisation: To begin with, fit the predictive mean to the data using weighted least squares [Tobar and Turner, 2016]. Appendix I.2 shows that this is equivalent to performing inference in the following model:

Model 7 (Basis Function Model). *Let \mathbf{f} be such that*

$$f_i(\mathbf{t}) = \sum_{j=1}^M \boldsymbol{\mu}_{H_{i,j}}^T \mathbf{K}_{\mathbf{u}_{H_{i,j}}}^{-1} \mathbf{I}^{(L, H_{i,j}, x_j)}(\mathbf{t}) \mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}^{-1} \boldsymbol{\mu}_{\tilde{x}_j}$$

for some kernel \mathcal{K}_x , some kernel \mathcal{K}_H and some vectors $\boldsymbol{\mu}_{H_{i,j}}$ and $\boldsymbol{\mu}_{\tilde{x}_j}$. Afterwards draw $\boldsymbol{\varepsilon} \sim \mathcal{GP}(\mathbf{0}, \boldsymbol{\Lambda}^2)$ for some diagonal matrix $\boldsymbol{\Lambda}$. Then observations are generated by

$$\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon}.$$

Learning $\boldsymbol{\mu}_{H_{i,j}}$ and $\boldsymbol{\mu}_{\tilde{x}_j}$ through Model 7 yields arguably sensible initial values. Finally, we initialise $\boldsymbol{\Sigma}_{H_{i,j}}$ and $\boldsymbol{\Sigma}_{\tilde{x}_j}$ to fractions of their priors $\mathbf{K}_{\mathbf{u}_{H_{i,j}}}$ and $\mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}$.

Observe that Model 7 is a multi-output basis function model that is structured in its parametrisation. Model 7 is a generalisation of the basis function model presented by Tobar and Turner [2016]; the case $N = M = K = 1$ recovers their model exactly.

3.7 Conclusion

The Generalised Gaussian Process Convolution Model enabled us to formulate Model 5. Model 5 addresses the kernel design problem in multi-output Gaussian processes on multidimensional input spaces by parameterising the kernel with another Gaussian process. By modelling the kernel nonparametrically we avoid choosing a kernel of parametric form; instead, we infer its form from the data.

Finally, Model 5 is intimately connected to other models (Section 3.5). Most notably, Model 5 enabled us to formulate Model 6, which forms a dual formulation of the cross-spectral mixture kernel and reveals Model 5 as a cross-spectral mixture kernel with an infinite number of components.

3.8 Discussion

If $T_{\tilde{x}} \gg T_H$ and $T_{\tilde{x}} \gg M$, then Section 3.6 shows that the variational free energy of Model 5 can be computed in $\mathcal{O}[MT_{\tilde{x}}^3 + YNM(K + T_H)T_{\tilde{x}}^2]$ time. This complexity scales linearly in Y . However, in many settings—for example, time series—the size of the data is proportional to the size of the space the data occupies. Now, the inducing points for \boldsymbol{x} must cover the space the data occupies. Therefore, in such settings, $T_{\tilde{x}}$ effectively scales with Y , which means that cost of computing the free energy effectively scales with Y^3 .

4 | Multi-Task Learning

4.1 Introduction

Many inference problems involve dealing with multiple signals. These problems are best solved not only by learning signals individually, but also by simultaneously exploiting dependencies between signals. The signals can be interpreted as the output of multiple *tasks*; a problem involving multiple signals is therefore commonly referred to as a *multi-task* problem. We will use the terms “multi-output” and “multi-task” interchangeably.

A prominent example of a multi-task problem is the prediction of concentration levels of pollutants in geostatistics [Álvarez and Lawrence, 2011; Álvarez et al., 2009; Wilson et al., 2012]. Pollutants are often expensive to sample, which hinders accurately predicting their concentration levels. Fortunately, pollutants often strongly correlate with other substances, which can be cheap to sample. The multi-task predictor aims to improve predicting concentration levels of scarcely sampled pollutants by exploiting their correlation with densely sampled cheap substances.

Numerous multi-task models have been developed, many of which are Gaussian process models. Now, the Generalised Gaussian Process Convolution Model (Model 4) is also a multi-output model. This chapter provides an overview of multi-output models from the geostatistics and machine learning literature and shows how Model 4 fits in.

4.2 Mixing Models

A model for a multi-output signal \mathbf{f} can be constructed by assuming that at any \mathbf{t} the output $\mathbf{f}(\mathbf{t})$ is explained by the value $\mathbf{x}(\mathbf{t})$ of some latent process \mathbf{x} with independent components. If the relationship between $\mathbf{f}(\mathbf{t})$ and $\mathbf{x}(\mathbf{t})$ is linear, then we arrive at the Instantaneous Mixing Model (IMM):

Model 8 (Instantaneous Mixing Model (IMM)). *Let \mathbf{H} be a matrix. Draw*

$$\begin{aligned}\mathbf{x} &\sim \mathcal{GP}(\mathbf{0}, \mathcal{K}_x), \\ \boldsymbol{\varepsilon} &\sim \mathcal{GP}(\mathbf{0}, \boldsymbol{\Lambda}^2)\end{aligned}$$

independently for some diagonal kernel \mathcal{K}_x and some diagonal matrix $\boldsymbol{\Lambda}$. Then observations are generated by $\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon} = \mathbf{H}\mathbf{x} + \boldsymbol{\varepsilon}$.

Observe that Model 8 corresponds to ordinary Gaussian process regression with a kernel of the form $\mathcal{K}_f = \mathbf{H}\mathcal{K}_x\mathbf{H}^T$ for some matrix \mathbf{H} and some diagonal kernel \mathcal{K}_x . Alternatively, Model 8 can be interpreted as a factor analysis model at all t .

Model 8 assumes that $\mathbf{f}(t)$ depends only on $\mathbf{x}(t)$, meaning that $\mathbf{f}(t)$ is independent of $\mathbf{x}(t')$ for $t' \neq t$. This assumption limits the expressivity of \mathbf{f} significantly; for example, \mathbf{f} cannot be a smoothed version of \mathbf{x} . We therefore relax Model 8 by instead letting $\mathbf{f}(t)$ depend on $\mathbf{x}(t)$ for all t . If their relationship is again linear, then we arrive at the Convolutional Mixing Model (CMM):

Model 9 (Convolutional Mixing Model (CMM)). *Let \mathbf{H} be a matrix-valued function. Draw*

$$\begin{aligned}\mathbf{x} &\sim \mathcal{GP}(\mathbf{0}, \mathcal{K}_x), \\ \boldsymbol{\varepsilon} &\sim \mathcal{GP}(\mathbf{0}, \boldsymbol{\Lambda}^2)\end{aligned}$$

*independently for some diagonal kernel \mathcal{K}_x and some diagonal matrix $\boldsymbol{\Lambda}$. Then observations are generated by $\mathbf{y} = \mathbf{f} + \boldsymbol{\varepsilon} = \mathbf{H} * \mathbf{x} + \boldsymbol{\varepsilon}$.*

Observe that Model 9 is a generalisation of Model 8, since letting $\mathbf{H} = \delta\mathbf{H}'$ for some constant matrix \mathbf{H}' yields that

$$\mathbf{f}(t) = (\delta\mathbf{H}' * \mathbf{x})(t) = \int_{\mathbb{R}^K} \delta(t - \tau)\mathbf{H}'\mathbf{x}(\tau) d\tau = \mathbf{H}'\mathbf{x}(t).$$

4.3 The Mixing Model Hierarchy

Many multi-output Gaussian process models from the geostatistics and machine learning literature can be identified as specialisations of Models 8 and 9; Table 4.1 shows the identification of the intrinsic coregionalisation model (ICM) [Goovaerts, 1997], the lin-

Model	Form of \mathbf{H}	Form of \mathcal{K}_x	Mixing type
ICM [Goovaerts, 1997]	\mathbf{H}	$(\sum_{q=1}^Q k_x^{(q)})\mathbf{I}$	Instantaneous
LCM [Goovaerts, 1997]	$[\mathbf{H}_1 \ \cdots \ \mathbf{H}_Q]$	$\text{diag}(k_x^{(1)}\mathbf{I}, \dots, k_x^{(Q)}\mathbf{I})$	Instantaneous
SLFM [Teh and Seeger, 2005]	\mathbf{H}	\mathcal{K}_x	Instantaneous
MTGPM [Bonilla et al., 2008]	\mathbf{H}	$k_x\mathbf{I}$	Instantaneous
LFM [Álvarez et al., 2009]	Green’s function	\mathcal{K}_x	Convolutional
CMOGPM [Álvarez and Lawrence, 2011]	$[\mathbf{H}_1 \ \cdots \ \mathbf{H}_Q]$	$\text{diag}(k_x^{(1)}\mathbf{I}, \dots, k_x^{(Q)}\mathbf{I})$	Convolutional
CGPM [Nguyen and Bonilla, 2014]	$[\mathbf{H} \ \mathbf{I}]$	\mathcal{K}_x	Instantaneous

Table 4.1: Identification of multi-output models from the geostatistics and machine learning literature as specialisations of Models 8 and 9

ear coregionalisation model (LCM) [Goovaerts, 1997], the semiparametric latent factor model (SLFM) [Teh and Seeger, 2005], the multi-task Gaussian process model (MT-GPM) [Bonilla et al., 2008], the latent force model (LFM) [Álvarez et al., 2009], the convolved multi-output Gaussian process model (CMOGPM) [Álvarez and Lawrence, 2011] and the collaborative Gaussian processes model (CGPM) [Nguyen and Bonilla, 2014].

Model 9 is also closely related to Model 4. Specifically, Model 4 is a generalisation of Model 9 where \mathbf{H} is modelled stochastically. This shows that Models 4, 8 and 9 form a hierarchy—the *mixing model hierarchy*—in which Model 4 generalises Model 9 and Model 9 generalises Model 8.

Figure 4.1 organises the main models presented in this thesis, the cross-spectral mixture kernel (CSMK) [Ulrich et al., 2015], and the models in Table 4.1 according to the mixing model hierarchy. This yields an overview of many multi-output Gaussian process models from the geostatistics and machine learning literature that emphasises their distinctive modelling assumptions. An immediate result is that the cross-spectral mixture kernel can now be connected to other multi-output models.

An important model that does not fit in the mixing model hierarchy is the Gaussian process regression network (GPRN) [Wilson et al., 2012].

4.4 Conclusion

We have presented the mixing model model hierarchy, which organises many multi-output Gaussian process from the geostatistics and machine learning literature according

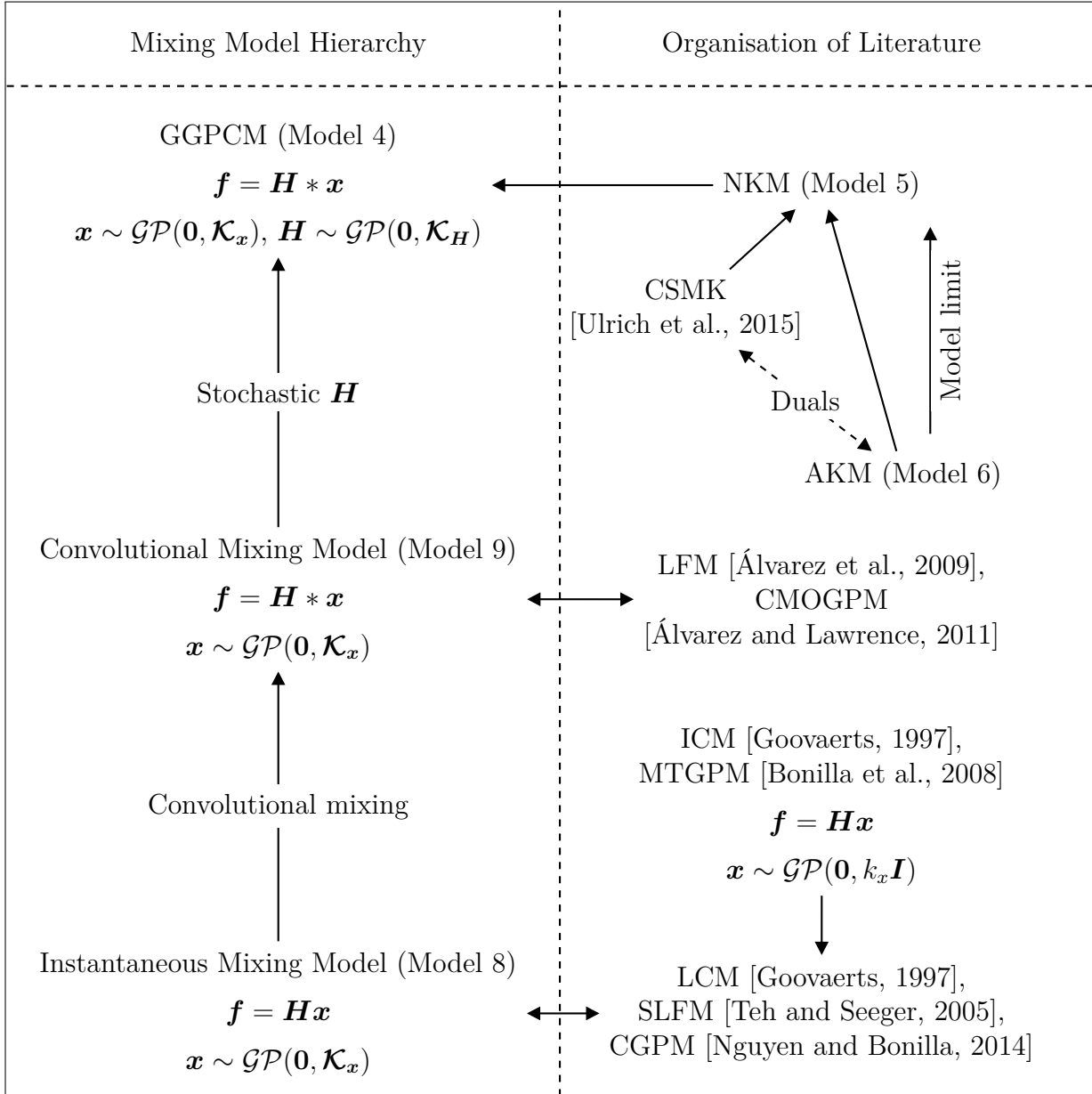


Figure 4.1: Organisation of multi-output models from the geostatistics and machine learning literature according to the mixing model hierarchy. Single-header arrows indicate generalisation and double-headed arrows indicate equivalence. “Model limit” refers to taking the model’s number of components to infinity.

to their distinctive modelling assumptions.

Model 4 fits in as a generalisation of Model 9 and Model 8. This shows that many multi-output Gaussian process models from the current literature can be derived from Model 4, with the exception being the Gaussian process regression network.

5 | The Deep Kernel Model

5.1 Introduction

Chapter 3 developed the idea that the kernel design problem can be addressed by modelling the kernel nonparametrically. In retrospect, this approach actually replaces the kernel design problem with another one: we now have to design the *kernel of the kernel*. Fortunately, we know already how to solve this kernel design problem: we can parametrise the *kernel of the kernel* with another Gaussian process. But then we have to design the *kernel of the kernel of the kernel...*

This chapter investigates a model that not only models the kernel nonparametrically, but also the *kernel of the kernel*, the *kernel of the kernel of the kernel*, and further “*deeper*” kernels. We confine the presentation to one-dimensional signals f on one-dimensional input spaces.

5.2 The Deep Kernel Model

Consider the Nonparametric Kernel Model (Model 5). In the case of one-dimensional signals on one-dimensional input spaces Model 5 recovers the Gaussian Process Convolution Model (GPCM) [Tobar et al., 2015b] (Section 3.5):

Model 10 (Gaussian Process Convolution Model [Tobar et al., 2015b]). *Draw*

$$\begin{aligned}h &\sim \mathcal{GP}(0, \mathcal{K}_h), \\x &\sim \mathcal{GP}[0, \delta(t_1 - t_2)], \\ \varepsilon &\sim \mathcal{GP}[0, \sigma^2 \delta(t_1 - t_2)]\end{aligned}$$

*independently for some kernel \mathcal{K}_h and some constant σ . Then observations are generated by $y = f + \varepsilon = h * x + \varepsilon$.*

Recall that Model 10 corresponds to ordinary Gaussian process regression in which the kernel is modelled nonparametrically (Section 3.3). We immediately recognise that the kernel design problem is replaced with another one: we now have to choose \mathcal{K}_h —the *kernel of the kernel*. We investigate whether \mathcal{K}_h can be modelled nonparametrically.

Modelling \mathcal{K}_h nonparametrically immediately presents a problem. In Section 3.3.1 we ensured that f has finite variance by letting \mathcal{K}_h be a *decaying* exponentiated-quadratic kernel. However, if we model \mathcal{K}_h nonparametrically, then it is not clear how we must restrict the prior on \mathcal{K}_h to ensure that f has finite variance. We investigate by examining the case that \mathcal{K}_h is a *non-decaying* exponentiated-quadratic kernel. In that case f has infinite variance (Section 3.3.1). Now, let $w(t) = \exp(-\alpha t^2)$, $t \in \mathbb{R}$ and consider the process wh . It holds that

$$\begin{aligned} \mathcal{K}_{wh}(t_1, t_2) &= \mathbb{E}[w(t_1)h(t_1)w(t_2)h(t_2)] \\ &= w(t_1)w(t_2)\mathbb{E}[h(t_1)h(t_2)] \\ &= \exp(-\alpha t_1^2) \exp(-\alpha t_2^2) \sigma_h^2 \exp[-\gamma(t_1 - t_2)^2] \\ &= \sigma_h^2 \exp[-\alpha t_1^2 - \alpha t_2^2 - \gamma(t_1 - t_2)^2], \end{aligned}$$

which is a *decaying* exponentiated-quadratic kernel. In other words, despite h having a kernel for which f has infinite variance, multiplication by w yields a kernel for which f has finite variance. This suggests the following, more robust formulation of the GPCM:

Model 11 (Gaussian Process Convolution Model (Explicit Decay)). *Draw*

$$\begin{aligned} h &\sim \mathcal{GP}(0, \mathcal{K}_h), \\ x &\sim \mathcal{GP}[0, \delta(t_1 - t_2)], \\ \varepsilon &\sim \mathcal{GP}[0, \sigma^2 \delta(t_1 - t_2)] \end{aligned}$$

*independently for some kernel \mathcal{K}_h and some constant σ . Let $w(t) = \exp(-\alpha t^2)$, $w \in \mathbb{R}$. Then observations are generated by $y = f + \varepsilon = wh * x + \varepsilon$.*

It now holds that

$$\begin{aligned}\mathbb{E}[f^2(t)] &= \int_{\mathbb{R}^2} w(t - \tau_1)w(t - \tau_2)\mathbb{E}[h(t - \tau_1)h(t - \tau_2)] \underbrace{\mathbb{E}[x(\tau_1)x(\tau_2)]}_{\delta(\tau_1 - \tau_2)} d\tau_1 d\tau_2 \\ &= \int_{\mathbb{R}} w^2(\tau)\mathcal{K}_h(\tau, \tau) d\tau,\end{aligned}$$

which clearly is finite if \mathcal{K}_h is bounded. In other words, if h has finite variance, then f has finite variance. Therefore, we can safely model \mathcal{K}_h nonparametrically if we ensure that h has finite variance.

To model \mathcal{K}_h nonparametrically, consider the following extension of Model 11:

Model 12 (Deep Gaussian Process Convolution Model). *Let N be the model's order. Draw*

$$\begin{aligned}h_1 &\sim \mathcal{GP}(0, \mathcal{K}_h), \\ x_1 &\sim \mathcal{GP}[0, \delta(t_1 - t_2)], \\ &\vdots \\ x_N &\sim \mathcal{GP}[0, \delta(t_1 - t_2)], \\ \varepsilon &\sim \mathcal{GP}[0, \sigma^2\delta(t_1 - t_2)]\end{aligned}$$

independently for some kernel \mathcal{K}_h and some constant σ . Let $w(t) = \exp(-\alpha t^2)$, $t \in \mathbb{R}$. Then observations are generated by y where

$$\begin{aligned}h_1 &= h, \\ h_2 &= wh_1 * x_1, \\ &\vdots \\ h_{N+1} &= wh_N * x_N, \\ f &= h_{N+1}, \\ y &= f + \varepsilon.\end{aligned}$$

Observe that each $h_{i+1} | h_i$ is a linear combination of Gaussian processes. Hence $h_{i+1} | h_i$ is another Gaussian process, which thus can be identified by its mean function and

kernel:

$$\begin{aligned}
\mathbb{E}[h_{i+1}(t) | h_i] &= \int_{\mathbb{R}} w(t - \tau) h_i(t - \tau) \mathbb{E}[x_i(\tau)] d\tau = 0, \\
\mathcal{K}_{h_{i+1} | h_i}(t_1, t_2) &= \mathbb{E}[h_{i+1}(t_1) h_{i+1}(t_2) | h_i] \\
&= \int_{\mathbb{R}^2} w(t_1 - \tau_1) h_i(t_1 - \tau_1) w(t_2 - \tau_2) h_i(t_2 - \tau_2) \underbrace{\mathbb{E}[x_i(\tau_1) x_i(\tau_2)]}_{\delta(\tau_1 - \tau_2)} d\tau_1 d\tau_2 \\
&= \int_{\mathbb{R}} w[\tau - (t_2 - t_1)] h_i[\tau - (t_2 - t_1)] w(\tau) h_i(\tau) d\tau \\
&= [R(wh_i) * wh_i](t_2 - t_1) \\
&= \mathcal{K}_{h_{i+1}}(t_1 - t_2).
\end{aligned}$$

We have established the following equivalent model:

Model 13 (Deep Kernel Model). *Let N be the model's order. Draw*

$$\begin{aligned}
h &\sim \mathcal{GP}(0, \mathcal{K}_h), \\
\varepsilon &\sim \mathcal{GP}[0, \sigma^2 \delta(t_1 - t_2)]
\end{aligned}$$

independently for some kernel \mathcal{K}_h and some constant σ . Let $h = h_1$ and $w(t) = \exp(-\alpha t^2)$, $t \in \mathbb{R}$. Afterwards draw

$$h_{i+1} | h_i \sim \mathcal{GP}\{0, [R(wh_i) * wh_i](t_2 - t_1)\}$$

in the order $i = 1, \dots, N$. Then observations are generated by $y = f + \varepsilon = h_{N+1} + \varepsilon$.

Model 13 shows that the kernel of $f = h_{N+1}$ is now parametrised by h_N , whose kernel is parametrised by h_{N-1} , whose kernel is parametrised by h_{N-2} , et cetera. That is, we have successfully formulated a nonparametric model of the kernel, the *kernel of the kernel*, the *kernel of the kernel of the kernel*, and further “*deeper*” kernels. Furthermore, it holds

that

$$\begin{aligned}
\mathbb{E}[f^2(t)] &= \mathbb{E}[(w\{\cdots [w(wh_1 * x_1) * x_2] \cdots\} * x_N)^2(t)] \\
&= \int \left[\prod_{i=1}^N w\left(t - \sum_{j=i}^N \tau_j^{(1)}\right) w\left(t - \sum_{j=i}^N \tau_j^{(2)}\right) \right] \\
&\quad \mathbb{E}\left[h_1\left(t - \sum_{i=1}^N \tau_i^{(1)}\right) h_1\left(t - \sum_{i=1}^N \tau_i^{(2)}\right)\right] \prod_{i=1}^N \underbrace{\mathbb{E}[x_i(\tau_i^{(1)}) x_i(\tau_i^{(2)})]}_{\delta(\tau_i^{(1)} - \tau_i^{(2)})} d\tau_i^{(1)} d\tau_i^{(2)} \\
&= \int \left[\prod_{i=1}^N w^2\left(t - \sum_{j=i}^N \tau_j\right) \right] \mathcal{K}_{h_1}\left(t - \sum_{i=1}^N \tau_i, t - \sum_{i=1}^N \tau_i\right) \prod_{i=1}^N d\tau_i \\
&= \int_{\mathbb{R}} w^2(\tau_1) \int_{\mathbb{R}} w^2(\tau_1 + \tau_2) \cdots \int_{\mathbb{R}} w^2\left(\sum_{i=1}^N \tau_i\right) \mathcal{K}_{h_1}\left(\sum_{i=1}^N \tau_i, \sum_{i=1}^N \tau_i\right) d\tau_N \cdots d\tau_2 d\tau_1.
\end{aligned}$$

Observe that the most inner integral is finite if \mathcal{K}_{h_1} is bounded, in which case all integrals are finite. In other words, if h_1 has finite variance, then f has finite variance. Therefore Models 12 and 13 are well defined.

5.2.1 Network Interpretation

The left side of Figure 5.1 depicts the graphical model of Model 13. We recognise the “deep” kernel structure by the chain of latent h_i ’s.

The right side of Figure 5.1 depicts the graphical model of Model 13 if the latent processes h_i are expanded into their function values; that is, any h_i is instead represented by $h_i(t)$ for all t —uncountably many t . Now, since the kernel of h_{i+1} is a convolution of h_i , it holds that every $h_{i+1}(t)$ depends on $h_i(t')$ for all t' . Hence the resulting graphical model is a fully connected network with layers of infinite size. This is very reminiscent of neural networks. The key difference between a neural network and Model 13 is that in a neural network layer i forms the input of layer $i + 1$, whereas in Model 13 layer i *parametrises* layer $i + 1$.

We can truncate these layers of infinite size to finite size by approximating the kernels in Model 13 by their expectations conditioned on some finite number of observations. We then obtain the Approximate Deep Kernel Model:

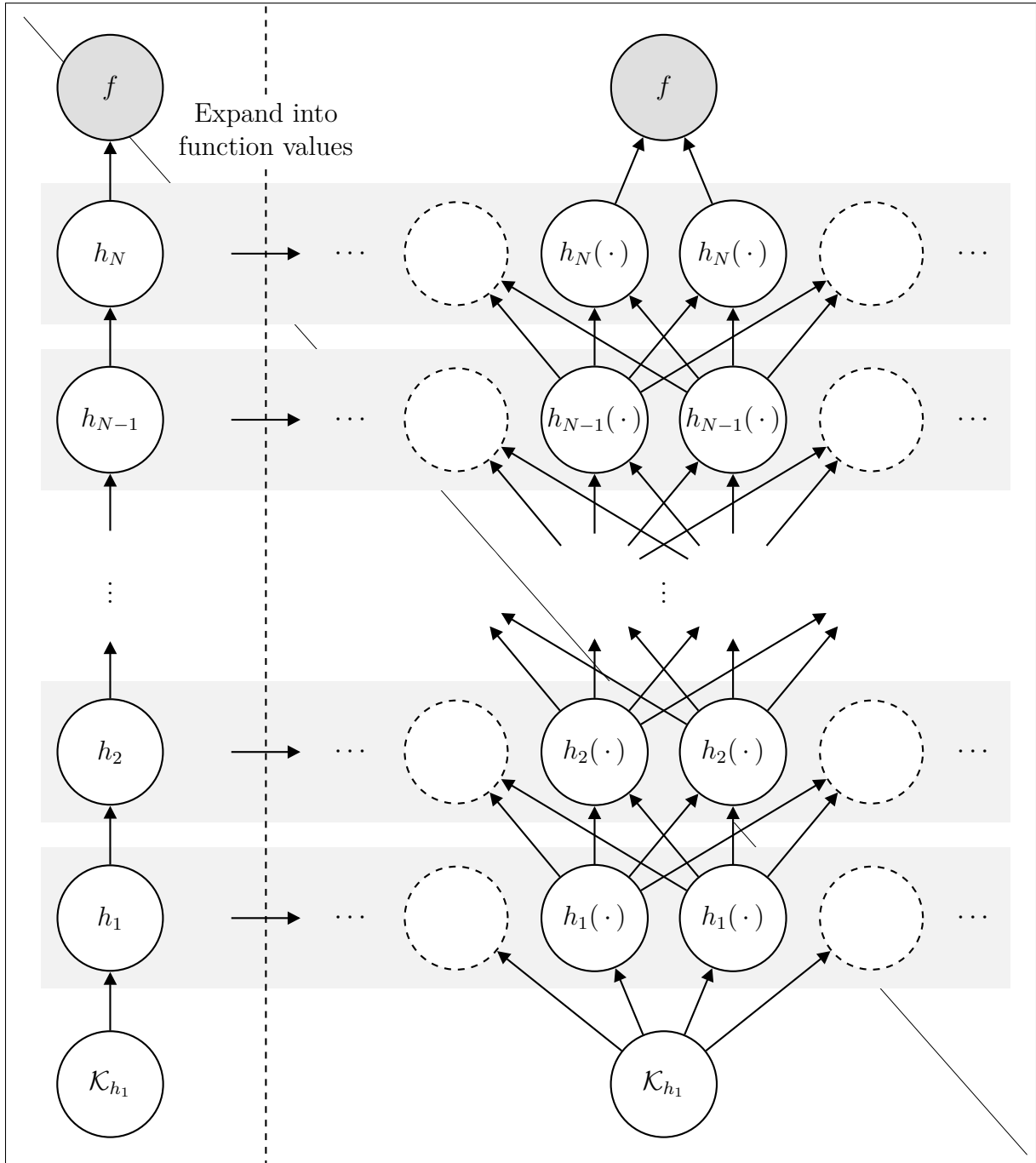


Figure 5.1: Graphical model of Model 13. Shows the resulting network structure if the latent h_i are expanded into their function values.

Model 14 (Approximate Deep Kernel Model). *Let N be the model's order. Draw*

$$\begin{aligned} h &\sim \mathcal{GP}(0, \mathcal{K}_h), \\ \varepsilon &\sim \mathcal{GP}[0, \sigma^2 \delta(t_1 - t_2)] \end{aligned}$$

independently for some kernel \mathcal{K}_h and some constant σ . Let $h = h_1$ and $w(t) = \exp(-\alpha t^2)$, $w \in \mathbb{R}$. Afterwards draw

$$\begin{aligned} \mathcal{K}_{h_{i+1}|h_i}(t_1, t_2) &= \mathbb{E}\{[R(wh_i) * wh_i](t_2 - t_1) | h_i(\mathbf{t})\}, \\ h_{i+1} | h_i &\sim \mathcal{GP}(0, \mathcal{K}_{h_{i+1}|h_i}) \end{aligned}$$

in the order $i = 1, \dots, N$. Then observations are generated by $y = f + \varepsilon = h_{N+1} + \varepsilon$.

Note that Model 14 approximates Model 13 in same way that Model 6 approximates Model 5.

5.3 Illustrative Samples

In the following we let h_1 have an exponentiated-quadratic kernel.

Figure 5.2 illustrates the generative process of Model 13 in the case that $N = 10$. Observe that the produced $h_i | h_{i-1}$'s and $\mathcal{K}_{h_i|h_{i-1}}$'s change in complexity as the generation progresses; that is, if by chance $h_i | h_{i-1}$ is slightly more complicated, then $\mathcal{K}_{h_{i+1}|h_i}$ and thereby $h_{i+1} | h_i$ tend to be more complicated as well. Thus, in each generation step Model 13 can either simplify or complicate the kernel, which means that $\mathcal{K}_{f|h_{10}}$ can be of a vastly different form than the kernel produced by h_1 ; Model 13 therefore seems to be less restricted by h_1 's kernel than Model 10 is by h 's kernel.

Figure 5.3 shows observations from Model 13 in the cases that $N = 1$ and $N = 20$ —a *shallow* and *deep* model respectively. Consider the shallow model. Observe that its observations exhibit barely any variation in smoothness and degree of periodicity; that is, their kernels are all of a similar form and their power spectral densities all show a similar low-pass structure. This is due h_1 's exponentiated-quadratic kernel, which imposes smoothness on h_1 and thereby directly on the observation's kernel. On the other hand, consider the deep model. Observe that its observations are much richer: they can be periodic, aperiodic, smoothly varying, or anywhere inbetween. This is reflected in the generated kernels, which now show great diversity. Furthermore, instead of the low-pass structure observed for the shallow model, the power spectral densities

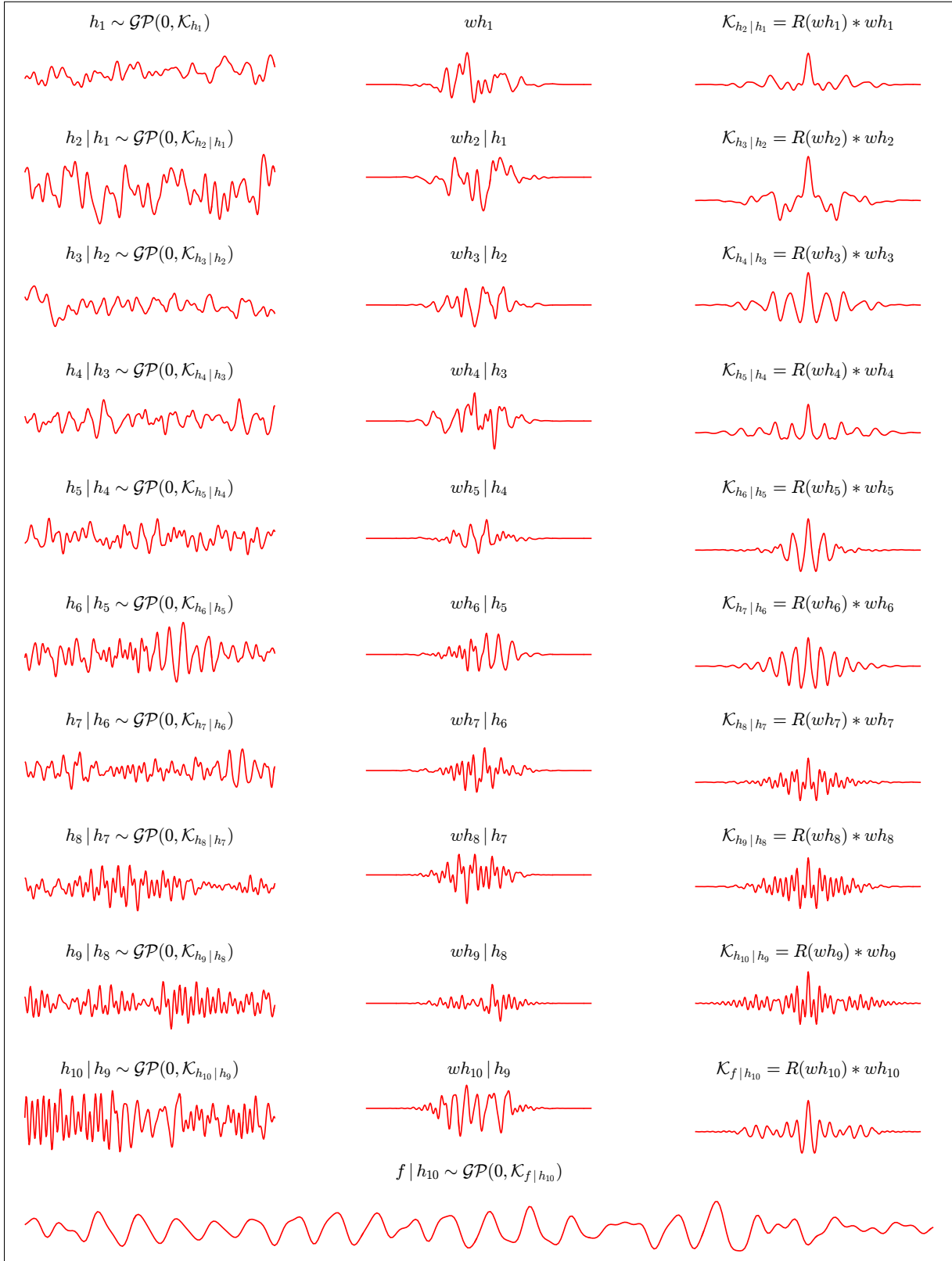
now show low-pass structures, high-pass structures, band-pass structures, and structures inbetween.

5.4 Conclusion

We have presented Model 13, which not only models the kernel nonparametrically, but also the *kernel of the kernel*, the *kernel of the kernel of the kernel*, and further “*deeper*” kernels. Experiments showed that Model 13 exhibits greatly increased expressivity compared to its shallow counterpart Model 10, or equivalently Model 5 where $M = N = K = 1$.

5.5 Discussion

We can perform approximate inference in Model 13 along the lines of Section 3.6. Although the free energy allows to be solved for analytically, its resulting complexity is exponential in the number of layers. Therefore, to perform inference in Model 13, additional approximations are probably necessary.

Figure 5.2: Generative process of Model 13 in the case that $N = 10$

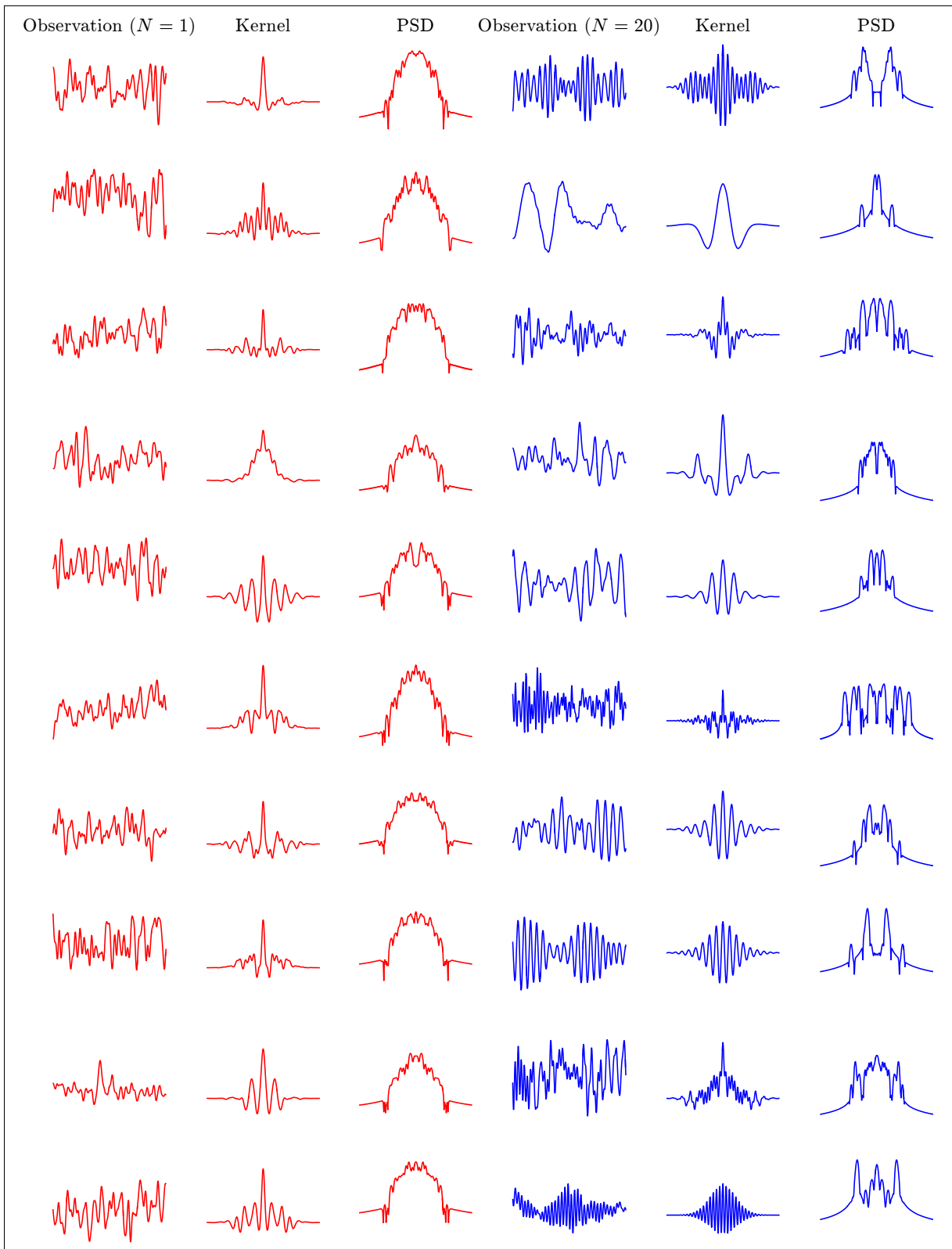


Figure 5.3: Observations from Model 13 in the cases that $N = 1$ and $N = 20$. Also shows the kernels and power spectral densities (PSDs) of the observations.

A | Solution of the Linear State-Space Model

A.1 Time-Variant Solution

Consider

$$\begin{aligned} \mathbf{s}'(t) &= \mathbf{A}(t)\mathbf{s}(t) + \mathbf{B}(t)\mathbf{x}(t), \\ \mathbf{f}(t) &= \mathbf{C}(t)\mathbf{s}(t) + \mathbf{D}(t)\mathbf{x}(t). \end{aligned} \tag{A.1}$$

Let $\Psi(t)$ be a fundamental matrix of the homogeneous system—that is, $\Psi'(t) = \mathbf{A}(t)\Psi(t)$ where $\Psi(t)$'s columns are independent. We proceed to find the general solution via variation of parameters. Assume a solution of the form $\mathbf{s}(t) = \Psi(t)\mathbf{c}(t)$. Then application of the product rule yields that

$$\mathbf{s}'(t) = \Psi'(t)\mathbf{c}(t) + \Psi(t)\mathbf{c}'(t) = \mathbf{A}(t)\Psi(t)\mathbf{c}(t) + \Psi(t)\mathbf{c}'(t)$$

while substitution into Equation (A.1) yields that

$$\mathbf{s}'(t) = \mathbf{A}(t)\Psi(t)\mathbf{c}(t) + \mathbf{B}(t)\mathbf{x}(t).$$

Hence $\Psi(t)\mathbf{c}'(t) = \mathbf{B}(t)\mathbf{x}(t)$. Since $\Psi(t)$ is invertible we have that $\mathbf{c}'(t) = \Psi^{-1}(t)\mathbf{B}(t)\mathbf{x}(t)$. Therefore

$$\begin{aligned}
 \mathbf{f}(t) &= \mathbf{C}(t)\Psi(t)\mathbf{c}(t) + \mathbf{D}(t)\mathbf{x}(t) \\
 &= \mathbf{C}(t)\Psi(t) \int_{-\infty}^t \Psi^{-1}(\tau)\mathbf{x}(\tau)\mathbf{B}(t) \, d\tau + \mathbf{D}(t)\mathbf{x}(t) \\
 &= \int_{-\infty}^t \mathbf{C}(t)\Psi(t)\Psi^{-1}(\tau)\mathbf{B}(t)\mathbf{x}(\tau) \, d\tau + \mathbf{D}(t)\mathbf{x}(t) \\
 &= \int_{\mathbb{R}} \underbrace{[\mathbb{1}(\tau \leq t)\mathbf{C}(t)\Psi(t)\Psi^{-1}(\tau)\mathbf{B}(t) + \delta(t - \tau)\mathbf{D}(t)]}_{\mathbf{H}(t,\tau)} \mathbf{x}(\tau) \, d\tau \\
 &= \int_{\mathbb{R}} \mathbf{H}(t,\tau)\mathbf{x}(\tau) \, d\tau.
 \end{aligned}$$

A.2 Time-Invariant Solution

Let $\mathbf{A}(t)$, $\mathbf{B}(t)$, $\mathbf{C}(t)$ and $\mathbf{D}(t)$ be independent of t ; that is, let $\mathbf{A}(t) = \mathbf{A}$, $\mathbf{B}(t) = \mathbf{B}$, $\mathbf{C}(t) = \mathbf{C}$ and $\mathbf{D}(t) = \mathbf{D}$. The fundamental matrix is then given by the matrix exponential

$$\Psi(t) = \sum_{i=0}^{\infty} \frac{t^i \mathbf{A}^i}{i!},$$

denoted as $\Psi(t) = \exp t\mathbf{A}$. To show this, note that by the binomial theorem and a change of variables

$$\begin{aligned}
 \exp[(s+t)\mathbf{A}] &= \sum_{i=0}^{\infty} \frac{(s+t)^i \mathbf{A}^i}{i!} \\
 &= \sum_{i=0, j=0}^{\infty, \infty} \frac{s^{i-j} t^j \mathbf{A}^i}{(i-j)! j!} \\
 &= \sum_{i=0, j=0}^{\infty, \infty} \frac{s^i t^j \mathbf{A}^{i+j}}{i! j!} \\
 &= (\exp s\mathbf{A})(\exp t\mathbf{A})
 \end{aligned}$$

so that

$$\begin{aligned}\Psi'(t) &= \lim_{h \rightarrow 0} \frac{\exp[(s+h)\mathbf{A}] - \exp s\mathbf{A}}{h} \\ &= \left(\lim_{h \rightarrow 0} \frac{\exp h\mathbf{A} - \mathbf{I}}{h} \right) \exp s\mathbf{A} \\ &= \mathbf{A}\Psi(t).\end{aligned}$$

Now

$$\mathbf{I} = \exp \mathbf{0} = \exp(t\mathbf{A} - t\mathbf{A}) = (\exp t\mathbf{A})[\exp(-t\mathbf{A})]$$

implies that $(\exp t\mathbf{A})^{-1} = \exp(-t\mathbf{A})$. Consequently

$$\begin{aligned}\mathbf{H}(t, \tau) &= \mathbf{1}(t - \tau \geq 0)\mathbf{C}\Psi(t)\Psi^{-1}(\tau)\mathbf{B} + \delta(t - \tau)\mathbf{D} \\ &= \mathbf{1}(t - \tau \geq 0)\mathbf{C}(\exp t\mathbf{A})(\exp \tau\mathbf{A})^{-1}\mathbf{B} + \delta(t - \tau)\mathbf{D} \\ &= \mathbf{1}(t - \tau \geq 0)\mathbf{C}(\exp t\mathbf{A})[\exp(-\tau\mathbf{A})]\mathbf{B} + \delta(t - \tau)\mathbf{D} \\ &= \mathbf{1}(t - \tau \geq 0)\mathbf{C} \exp[(t - \tau)\mathbf{A}]\mathbf{B} + \delta(t - \tau)\mathbf{D} \\ &= \mathbf{H}(t - \tau).\end{aligned}$$

B | Properties of the Multivariate Gaussian Distribution

B.1 Marginal and Conditional Distribution

Let

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{1,1} & \boldsymbol{\Sigma}_{1,2} \\ \boldsymbol{\Sigma}_{2,1} & \boldsymbol{\Sigma}_{2,2} \end{bmatrix} \right).$$

Then Murphy [2012] shows that $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{1,1})$ and

$$\mathbf{x}_1 | \mathbf{x}_2 \sim \mathcal{N}[\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{1,2}\boldsymbol{\Sigma}_{2,2}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{1,1} - \boldsymbol{\Sigma}_{1,2}\boldsymbol{\Sigma}_{2,2}^{-1}\boldsymbol{\Sigma}_{2,1}].$$

B.2 Kullback-Leibler Divergence

Let $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $\mathbf{x}_2 \sim \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ where \mathbf{x}_1 and \mathbf{x}_2 attain values in \mathbb{R}^N . Then Murphy [2012] shows that

$$D_{KL}[p(\mathbf{x}_1) || p(\mathbf{x}_2)] = \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1}\boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) + \log \frac{|\boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|} - N \right].$$

C | Multivariate Matrix-Valued Gaussian Processes

A stochastic process $\mathbf{F}(\mathbf{t}) \in \mathbb{R}^{N \times M}$, $\mathbf{t} \in \mathbb{R}^K$ is Gaussian if and only if for every $\mathbf{T} \in \mathbb{R}^{T \times K}$ it holds that $(\text{vec } \mathbf{F})(\mathbf{T})$ is multivariate Gaussian distributed where

$$\begin{aligned} (\text{vec } \mathbf{F})(\mathbf{T}) &= \begin{bmatrix} \text{vec } \mathbf{F}(\mathbf{T}_{1,:}) \\ \vdots \\ \text{vec } \mathbf{F}(\mathbf{T}_{T,:}) \end{bmatrix} \\ &= \left[F_{1,1}(\mathbf{T}_{1,:}) \quad \cdots \quad F_{N,M}(\mathbf{T}_{1,:}) \quad \cdots \quad F_{1,1}(\mathbf{T}_{T,:}) \quad \cdots \quad F_{N,M}(\mathbf{T}_{T,:}) \right]^T. \end{aligned}$$

The mean function is then defined by

$$\begin{aligned} \mathbf{m}_{\mathbf{F}}(\mathbf{T}) &= \left[\mathbf{m}_{\mathbf{F}}^T(\mathbf{T}_{1,:}) \quad \cdots \quad \mathbf{m}_{\mathbf{F}}^T(\mathbf{T}_{T,:}) \right]^T \\ &= \left[m_{F_{1,1}}(\mathbf{T}_{1,:}) \quad \cdots \quad m_{F_{N,M}}(\mathbf{T}_{1,:}) \quad \cdots \quad m_{F_{1,1}}(\mathbf{T}_{T,:}) \quad \cdots \quad m_{F_{N,M}}(\mathbf{T}_{T,:}) \right]^T \end{aligned}$$

and the kernel by

$$\mathcal{K}_{\mathbf{F}}(\mathbf{T}^{(1)}, \mathbf{T}^{(2)}) = \begin{bmatrix} \mathcal{K}_{\mathbf{F}}(\mathbf{T}_{1,:}^{(1)}, \mathbf{T}_{1,:}^{(2)}) & \cdots & \mathcal{K}_{\mathbf{F}}(\mathbf{T}_{1,:}^{(1)}, \mathbf{T}_{T,:}^{(2)}) \\ \vdots & \ddots & \vdots \\ \mathcal{K}_{\mathbf{F}}(\mathbf{T}_{T,:}^{(1)}, \mathbf{T}_{1,:}^{(2)}) & \cdots & \mathcal{K}_{\mathbf{F}}(\mathbf{T}_{T,:}^{(1)}, \mathbf{T}_{T,:}^{(2)}) \end{bmatrix}$$

where

$$\mathcal{K}_{\mathbf{F}}(\mathbf{T}_{i,:}^{(1)}, \mathbf{T}_{j,:}^{(2)}) = \begin{bmatrix} \mathcal{K}_{F_{1,1}, F_{1,1}}(\mathbf{T}_{i,:}^{(1)}, \mathbf{T}_{j,:}^{(2)}) & \cdots & \mathcal{K}_{F_{1,1}, F_{N,M}}(\mathbf{T}_{i,:}^{(1)}, \mathbf{T}_{j,:}^{(2)}) \\ \vdots & \ddots & \vdots \\ \mathcal{K}_{F_{N,M}, F_{1,1}}(\mathbf{T}_{i,:}^{(1)}, \mathbf{T}_{j,:}^{(2)}) & \cdots & \mathcal{K}_{F_{N,M}, F_{N,M}}(\mathbf{T}_{i,:}^{(1)}, \mathbf{T}_{j,:}^{(2)}) \end{bmatrix}.$$

D | Gaussian Processes in Practice

D.1 Implementation of Gaussian Process Models

We briefly discuss some useful techniques concerning the implementation of Gaussian process models. Assume that all matrices have shape $N \times N$.

First, carefully considering the memory layout of large matrices can be beneficial. Namely, contiguous memory is accessed more cheaply than non-contiguous memory. This means that the cost of operations on large matrices greatly depends on the memory layout of the matrices.

Second, a covariance matrix Σ can be made positive semidefinite by adding a small diagonal. In theory, Σ is always positive semidefinite. However, when it is evaluated numerically, it can be indefinite or negative definite due to round-off errors introduced by the floating-point representation of numbers. In that case Σ 's positive semidefiniteness can be ensured by adding a diagonal comparable to $\varepsilon \|\Sigma\|_\infty$ where ε is the machine epsilon. Usually the noise introduced by the diagonal is negligible.

In the case that Σ is indefinite or negative definite due to noise other than round-off errors, the diagonal required to make Σ positive semidefinite might be so large that the introduced noise becomes significant. In that case Theorem 1 can be used to compute the symmetric positive semidefinite matrix nearest to Σ in Frobenius norm.

Third, quantities of the form $\mathbf{A}\Sigma^{-1}\mathbf{B}$ where Σ is a covariance matrix should not be computed by first computing $\Sigma^{-1} = \mathbf{X}$ explicitly. Namely, the error in $\mathbf{A}\mathbf{X}\mathbf{B}$ due to round-off errors can be large, especially when Σ is ill-conditioned [Trefethen and Bau, 1997]. A numerically more stable approach is to first compute the Cholesky decomposi-

tion $\Sigma = \mathbf{L}\mathbf{L}^T$ where \mathbf{L} is lower triangular. Then

$$\mathbf{A}\Sigma^{-1}\mathbf{B} = \mathbf{A}(\mathbf{L}\mathbf{L}^T)^{-1}\mathbf{B} = (\mathbf{A}\mathbf{L}^{-T})(\mathbf{L}^{-1}\mathbf{B})$$

where the systems $\mathbf{A}\mathbf{L}^{-T}$ and $\mathbf{L}^{-1}\mathbf{B}$ can be solved using respectively back substitution and forward substitution. The Cholesky decomposition has time complexity $\mathcal{O}(N^3)$ and back substitution and forward substitution have time complexity $\mathcal{O}(N^2)$ [Trefethen and Bau, 1997]. Thus $\mathbf{A}\Sigma^{-1}\mathbf{B}$ can be computed in $\mathcal{O}(N^3)$ time.

Fourth, the determinant Σ of a covariance matrix can efficiently be computed via its Cholesky decomposition $\Sigma = \mathbf{L}\mathbf{L}^T$. Specifically, as \mathbf{L} is lower triangular,

$$|\Sigma| = |\mathbf{L}|^2 = \prod_{i=1}^N L_{i,i}^2.$$

Thus the determinant can be computed in $\mathcal{O}(N^3)$ time.

D.2 Nearest Symmetric Positive-Semidefinite Matrix

Lemma 1 (Polar Decomposition of Symmetric Matrix). *Let \mathbf{A} be symmetric and let $\mathbf{B} = \mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^T$ be its spectral decomposition. Then there exists an orthogonal matrix \mathbf{U} such that*

$$\mathbf{A} = \mathbf{U}\mathbf{Z}\mathbf{S}\mathbf{Z}^T$$

where $\mathbf{S} = \text{diag}(|\Lambda_{1,1}|, \dots, |\Lambda_{N,N}|)$.

Proof. Let \mathbf{L} be diagonal such that $L_{i,i} = \text{sign } \Lambda_{i,i}$. Then $\mathbf{L}\mathbf{L} = \mathbf{I}$. Therefore

$$\mathbf{A} = \mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^T = \mathbf{Z}\mathbf{L}\mathbf{L}\mathbf{\Lambda}\mathbf{Z}^T = \underbrace{\mathbf{Z}\mathbf{L}\mathbf{Z}^T}_{\mathbf{U}} \underbrace{\mathbf{Z}\mathbf{L}\mathbf{\Lambda}\mathbf{Z}^T}_{\mathbf{S}}$$

where $\mathbf{L}\mathbf{\Lambda} = \text{diag}(|\Lambda_{1,1}|, \dots, |\Lambda_{N,N}|)$ and

$$\mathbf{U}\mathbf{U}^T = \mathbf{Z}\mathbf{L}\mathbf{Z}^T \mathbf{Z}\mathbf{L}\mathbf{Z}^T = \mathbf{Z}\mathbf{L}\mathbf{L}\mathbf{Z}^T = \mathbf{Z}\mathbf{Z}^T = \mathbf{I}.$$

□

Lemma 2. *Let \mathbf{A} be symmetric and let \mathbf{B} be antisymmetric. Then $\|\mathbf{A} + \mathbf{B}\|_F^2 = \|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2$.*

Proof. To begin with, $\mathbf{A}^T \mathbf{B} = -\mathbf{A} \mathbf{B}^T$; thus $\text{tr}(\mathbf{A}^T \mathbf{B}) + \text{tr}(\mathbf{A} \mathbf{B}^T) = 0$. Therefore

$$\begin{aligned} \|\mathbf{A} + \mathbf{B}\|_F^2 &= \text{tr}[(\mathbf{A} + \mathbf{B})^T (\mathbf{A} + \mathbf{B})] \\ &= \text{tr}(\mathbf{A}^T \mathbf{A}) + \text{tr}(\mathbf{A}^T \mathbf{B}) + \text{tr}(\mathbf{B}^T \mathbf{A}) + \text{tr}(\mathbf{B}^T \mathbf{B}) \\ &= \text{tr}(\mathbf{A}^T \mathbf{A}) + \text{tr}(\mathbf{B}^T \mathbf{B}) \\ &= \|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2. \end{aligned}$$

□

Lemma 3. *Let \mathbf{A} be a matrix and let \mathbf{U} be orthogonal. Then $\|\mathbf{U} \mathbf{A} \mathbf{U}^T\|_F^2 = \|\mathbf{A}\|_F^2$.*

Proof. Direct computation yields that

$$\begin{aligned} \|\mathbf{U} \mathbf{A} \mathbf{U}^T\|_F^2 &= \text{tr}(\mathbf{U} \mathbf{A} \mathbf{U}^T \mathbf{U} \mathbf{A}^T \mathbf{U}^T) \\ &= \text{tr}(\mathbf{U} \mathbf{A} \mathbf{A}^T \mathbf{U}^T) \\ &= \text{tr}(\mathbf{U}^T \mathbf{U} \mathbf{A} \mathbf{A}^T) \\ &= \text{tr}(\mathbf{A} \mathbf{A}^T) \\ &= \|\mathbf{A}\|_F^2. \end{aligned}$$

□

The following proof is based on the proof by Higham [1988].

Theorem 1 (Nearest Symmetric Positive-Semidefinite Matrix). *Let \mathbf{A} be a matrix and $\mathbf{B} = (\mathbf{A} + \mathbf{A}^T)/2$. Furthermore, let $\mathbf{B} = \mathbf{Z} \mathbf{\Lambda} \mathbf{Z}^T$ be \mathbf{B} 's spectral decomposition and let $\mathbf{B} = \mathbf{U} \mathbf{Z} \mathbf{S} \mathbf{Z}^T$ be \mathbf{B} 's polar decomposition (Lemma 1). Finally, let $\mathbf{A}_F = (\mathbf{B} + \mathbf{Z} \mathbf{S} \mathbf{Z}^T)/2$. Then*

1. \mathbf{A}_F is symmetric positive semidefinite and
2. for any other symmetric positive semidefinite matrix \mathbf{X} it holds that $\|\mathbf{A} - \mathbf{X}\|_F^2 \geq \|\mathbf{A} - \mathbf{A}_F\|_F^2$.

Proof. First, we verify that \mathbf{A}_F is symmetric positive semidefinite. Trivially, \mathbf{A}_F is symmetric. To see that \mathbf{A}_F is positive semidefinite, consider

$$\mathbf{A}_F = \frac{1}{2}(\mathbf{B} + \mathbf{Z}\mathbf{S}\mathbf{Z}^T) = \frac{1}{2}(\mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^T + \mathbf{Z}\mathbf{S}\mathbf{Z}^T) = \mathbf{Z}\frac{\mathbf{\Lambda} + \mathbf{S}}{2}\mathbf{Z}^T. \quad (\text{D.1})$$

By Lemma 1,

$$\frac{\mathbf{\Lambda} + \mathbf{S}}{2} = \text{diag}\left(\frac{\Lambda_{1,1} + |\Lambda_{1,1}|}{2}, \dots, \frac{\Lambda_{N,N} + |\Lambda_{N,N}|}{2}\right), \quad (\text{D.2})$$

whose diagonal elements are all nonnegative. Therefore, for any \mathbf{x} , it holds that

$$\mathbf{x}^T \mathbf{A}_F \mathbf{x} = \mathbf{x}^T \mathbf{Z} \left(\frac{\mathbf{\Lambda} + \mathbf{S}}{2}\right)^{1/2} \left(\frac{\mathbf{\Lambda} + \mathbf{S}}{2}\right)^{1/2} \mathbf{Z}^T \mathbf{x} = \left\| \left(\frac{\mathbf{\Lambda} + \mathbf{S}}{2}\right)^{1/2} \mathbf{Z}^T \mathbf{x} \right\|_2^2 \geq 0.$$

Second, let \mathbf{X} be symmetric positive semidefinite. We then show that $\|\mathbf{A} - \mathbf{X}\|_F^2 \geq \|\mathbf{A} - \mathbf{A}_F\|_F^2$. To begin with, let $\mathbf{C} = (\mathbf{A} - \mathbf{A}_F)/2$. Then $\mathbf{A} = \mathbf{B} + \mathbf{C}$ where \mathbf{B} is symmetric and \mathbf{C} is antisymmetric. Thus, by Lemma 2,

$$\|\mathbf{A} - \mathbf{X}\|_F^2 = \|(\mathbf{B} - \mathbf{X}) + \mathbf{C}\|_F^2 = \|\mathbf{B} - \mathbf{X}\|_F^2 + \|\mathbf{C}\|_F^2 \geq \|\mathbf{B} - \mathbf{X}\|_F^2.$$

Let $\mathbf{Y} = \mathbf{Z}^T \mathbf{X} \mathbf{Z}$. As \mathbf{X} is positive semidefinite, \mathbf{Y} is also positive semidefinite and so $Y_{i,i} \geq 0$. Therefore, by Lemma 3,

$$\begin{aligned} \|\mathbf{A} - \mathbf{X}\|_F^2 &\geq \|\mathbf{B} - \mathbf{X}\|_F^2 \\ &= \|\mathbf{Z}\mathbf{\Lambda}\mathbf{Z}^T - \mathbf{X}\|_F^2 \\ &= \|\mathbf{Z}(\mathbf{\Lambda} - \mathbf{Z}^T \mathbf{X} \mathbf{Z})\mathbf{Z}^T\|_F^2 \\ &= \|\mathbf{\Lambda} - \mathbf{Y}\|_F^2 \\ &\geq \sum_{i:\Lambda_{i,i} < 0} (\Lambda_{i,i} - Y_{i,i})^2 \\ &\geq \sum_{i:\Lambda_{i,i} < 0} \Lambda_{i,i}^2. \end{aligned}$$

This lower bound holds for any positive semidefinite \mathbf{X} . Therefore, if \mathbf{A}_F achieves this lower bound, then the result follows.

To this end, let $\mathbf{X} = \mathbf{A}_F$. As \mathbf{A}_F is symmetric, $\mathbf{C} = \mathbf{0}$ and so

$$\|\mathbf{A} - \mathbf{A}_F\|_F^2 = \|\mathbf{\Lambda} - \mathbf{Y}\|_F^2$$

where $\mathbf{Y} = \mathbf{Z}\mathbf{A}_F\mathbf{Z}^T$. Now, Equation (D.1) and Equation (D.2) show that

$$\mathbf{Y} = \text{diag}\left(\frac{\Lambda_{1,1} + |\Lambda_{1,1}|}{2}, \dots, \frac{\Lambda_{N,N} + |\Lambda_{N,N}|}{2}\right).$$

Hence

$$\begin{aligned} \|\mathbf{A} - \mathbf{A}_F\|_F^2 &= \|\mathbf{\Lambda} - \mathbf{Y}\|_F^2 \\ &= \sum_{i=1}^N \left(\Lambda_{i,i} - \frac{\Lambda_{i,i} + |\Lambda_{i,i}|}{2} \right)^2 \\ &= \sum_{i:\Lambda_{i,i} < 0} \left(\Lambda_{i,i} - \frac{\Lambda_{i,i} + |\Lambda_{i,i}|}{2} \right)^2 \\ &= \sum_{i:\Lambda_{i,i} < 0} \Lambda_{i,i}^2. \end{aligned}$$

□

E | Circulant Approximation of Stationary Multi-Output Kernel Matrices

E.1 Introduction

Let $\mathbf{K} : \mathbb{R}^K \times \mathbb{R}^K \rightarrow \mathbb{R}^{N \times N}$ be a stationary multi-output kernel and let \mathbf{K}_F be \mathbf{K} evaluated for points that lie on an evenly-spaced grid $\{t_1^{(1)}, \dots, t_{Y_1}^{(1)}\} \times \dots \times \{t_1^{(K)}, \dots, t_{Y_K}^{(K)}\}$. The total number of points is then given by $Y_1 \cdots Y_K = Y$. Assume that all Y_i are even.

Consider computation of $|\mathbf{K}_F|$ and a product of the form $\mathbf{P} = \mathbf{U}^T \mathbf{K}_F^{-1} \mathbf{V}$ for some $NY \times P$ matrix \mathbf{U} and $NY \times Q$ matrix \mathbf{V} . Computing $|\mathbf{K}_F|$ directly costs $\mathcal{O}(N^3 Y^3)$ time (Appendix D). This does not scale for large Y . Furthermore, inverting \mathbf{K}_F costs $\mathcal{O}(N^3 Y^3)$ time and then computing $\mathbf{U}^T \mathbf{K}_F^{-1} \mathbf{V}$ costs $\mathcal{O}(PQN^2 Y^2)$ time, resulting in time complexity $\mathcal{O}(PQN^2 Y^2 + N^3 Y^3)$. This also does not scale for large Y .

This chapter develops an approximation of \mathbf{K}_F based on the facts that \mathbf{K}_F is stationary and is evaluated for points on a grid [Ulrich et al., 2015]. We show that this approximation can be leveraged to compute $|\mathbf{K}_F|$ and $\mathbf{U}^T \mathbf{K}_F^{-1} \mathbf{V}$ with complexities that scale favourably in Y .

E.2 Circulant Approximation of Toeplitz Matrices

Let $k : \mathbb{R} \rightarrow \mathbb{R}$ be a stationary kernel and let \mathbf{K}_f be k evaluated for points that lie on an evenly-spaced grid $\{t_1, \dots, t_Y\}$. Then, by k 's stationary and the fact that the points

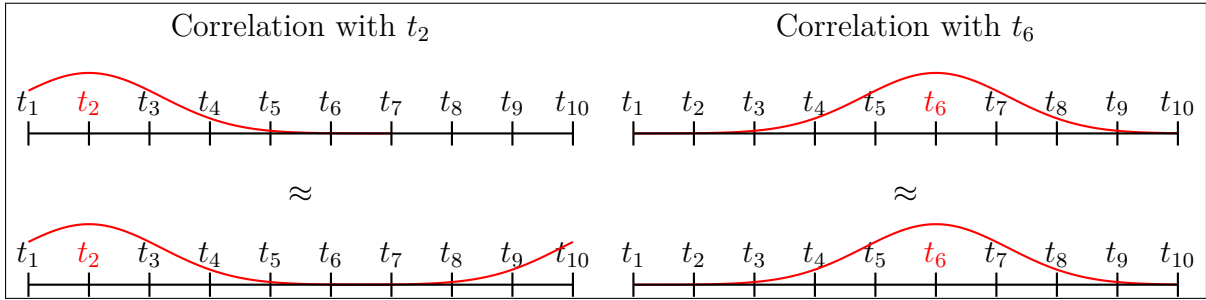


Figure E.1: Circulant approximation of a stationary kernel. Shows how the space “wraps around” at the edge of the grid. Also shows that the approximation is excellent for points not in the vicinity of the edge.

lie on a grid,

$$\mathbf{K}_f = \text{stoep}[k(t_1 - t_1), \dots, k(t_1 - t_Y)].$$

A banded Toeplitz matrix can be approximated by a particular circulant matrix where the approximation becomes more accurate as the matrix becomes large [Gray, 2006]. More specifically, if k has effectively limited support, then we can approximate

$$\begin{aligned} \mathbf{K}_f &= \text{stoep}[k(t_1 - t_1), \dots, k(t_1 - t_{Y/2}), k(t_1 - t_{Y/2+1}), k(t_1 - t_{Y/2+2}), \dots, k(t_1 - t_Y)] \\ &\approx \text{stoep}[k(t_1 - t_1), \dots, k(t_1 - t_{Y/2}), k(t_1 - t_{Y/2+1}), k(t_1 - t_{Y/2}), \dots, k(t_1 - t_2)] \\ &= \text{circ}[k(t_1 - t_1), \dots, k(t_1 - t_{Y/2}), k(t_1 - t_{Y/2+1}), k(t_1 - t_{Y/2}), \dots, k(t_1 - t_2)]. \quad (\text{E.1}) \end{aligned}$$

We call this approximation the *circulant approximation* of \mathbf{K}_f . Essentially, the approximation lets the space in which distance is measured now “wrap around” at the edge of the grid. As a consequence, points near the edge of the grid correlate. However, if the kernel is local, then correlations between points not in the vicinity of the edge are preserved; in that case \mathbf{K}_f ’s circulant approximation is an excellent approximation for most points. Figure E.1 illustrates the circulant approximation of a stationary kernel.

Finally, the $Y \times Y$ unitary discrete Fourier transform matrix \mathcal{F}_Y diagonalises any $Y \times Y$ circulant matrix [Gray, 2006]. Thus, by $\mathbf{K}_f = \mathbf{K}_f^T$, it holds that $\mathbf{K}_f \approx \mathcal{F}_Y \mathbf{\Lambda} \mathcal{F}_Y^H = \mathcal{F}_Y^H \mathbf{\Lambda} \mathcal{F}_Y$ where $\mathbf{\Lambda}$ is diagonal.

E.3 Circulant Approximation of Stationary Multi-Output Kernel Matrices

In the same way that a Toeplitz matrix has a circulant approximation, a block Toeplitz matrix has a block circulant approximation. We use this observation to approximate \mathbf{K}_F .

We construct \mathbf{K}_F in a particular way. Specifically, let $\mathbf{K}^{(0)}(\mathbf{t}, \mathbf{t}') = \mathbf{K}(\mathbf{t}, \mathbf{t}')$ and let

$$\begin{aligned} \mathbf{t}^{(>i)} &= [t_{i+1} \ \cdots \ t_K]^T, \\ \mathbf{t}^{(>i)'} &= [t'_{i+1} \ \cdots \ t'_K]^T, \\ \mathbf{K}^{(i)}(\mathbf{t}^{(>i)}, \mathbf{t}^{(>i)'}) &= \begin{bmatrix} \mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_1^{(i)}, \mathbf{t}^{(>i)'}) & \cdots & \mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_{Y_i}^{(i)}, \mathbf{t}^{(>i)'}) \\ \vdots & \ddots & \vdots \\ \mathbf{K}^{(i-1)}(t_{Y_i}^{(i)}, \mathbf{t}^{(>i)}, t_1^{(i)}, \mathbf{t}^{(>i)'}) & \cdots & \mathbf{K}^{(i-1)}(t_{Y_i}^{(i)}, \mathbf{t}^{(>i)}, t_{Y_i}^{(i)}, \mathbf{t}^{(>i)'}) \end{bmatrix} \end{aligned}$$

for $i = 1, \dots, K$. Then $\mathbf{K}_F = \mathbf{K}^{(K)}$. Now, \mathbf{K} is stationary and the points for which we evaluate \mathbf{K} lie on a grid; thus we can use Equation (E.1) to approximate

$$\begin{aligned} &\mathbf{K}^{(i)}(\mathbf{t}^{(>i)}, \mathbf{t}^{(>i)'}) \\ &= \text{stoep}[\mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_1^{(i)}, \mathbf{t}^{(>i)'}), \dots, \mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_{Y_i}^{(i)}, \mathbf{t}^{(>i)'})] \\ &\approx \text{circ}[\mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_{w^{(i)}(1)}^{(i)}, \mathbf{t}^{(>i)'}), \dots, \mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_{w^{(i)}(Y_i)}^{(i)}, \mathbf{t}^{(>i)'})] \end{aligned}$$

where $w^{(i)}(j_i) = \min\{j, Y_i - j_i + 2\}$. Then, by decomposing each $\mathbf{K}^{(i)}(\mathbf{t}^{(>i)}, \mathbf{t}^{(>i)'})$ as

$$\mathbf{K}^{(i)}(\mathbf{t}^{(>i)}, \mathbf{t}^{(>i)'}) \approx \sum_{j_i=1}^{Y_i} \overbrace{\text{circ}(\underbrace{0, \dots, 0}_{j_i-1 \text{ times}}, 1, \underbrace{0, \dots, 0}_{Y_i-j_i \text{ times}})}_{\mathcal{C}^{(i, j_i)}} \otimes \mathbf{K}^{(i-1)}(t_1^{(i)}, \mathbf{t}^{(>i)}, t_{w^{(i)}(j_i)}^{(i)}, \mathbf{t}^{(>i)'}),$$

we obtain that

$$\begin{aligned}
\mathbf{K}^{(K)} &\approx \sum_{j_K=1}^{Y_K} \mathbf{C}^{(K,j_K)} \otimes \mathbf{K}^{(K-1)}(t_1^{(K)}, t_{w^{(K)}(j_K)}^{(K)}) \\
&\approx \sum_{j_K=1, j_{K-1}=1}^{Y_K, Y_{K-1}} \mathbf{C}^{(K,j_K)} \otimes \mathbf{C}^{(K-1,j_{K-1})} \otimes \mathbf{K}^{(K-1)}(t_1^{(K-1)}, t_1^{(K)}, t_{w^{(K-1)}(j_{K-1})}^{(K-1)}, t_{w^{(K)}(j_K)}^{(K)}) \\
&\vdots \\
&\approx \sum_{j_K=1, \dots, j_1=1}^{Y_K, \dots, Y_1} \mathbf{C}^{(K,j_K)} \otimes \dots \otimes \mathbf{C}^{(1,j_1)} \otimes \mathbf{K}^{(0)}(\underbrace{t_1^{(1)}, \dots, t_1^{(K)}}_{\mathbf{t}_1}, \underbrace{t_{w^{(1)}(j_1)}^{(1)}, \dots, t_{w^{(K)}(j_K)}^{(K)}}_{\mathbf{t}_{w(j)}}) \\
&= \sum_{j_K=1, \dots, j_1=1}^{Y_K, \dots, Y_1} \mathbf{C}^{(K,j_K)} \otimes \dots \otimes \mathbf{C}^{(1,j_1)} \otimes \mathbf{K}(\mathbf{t}_1, \mathbf{t}_{w(j)}).
\end{aligned}$$

Let $\mathcal{F}^{(i)} = \mathcal{F}_{Y_i} \otimes \dots \otimes \mathcal{F}_{Y_1} \otimes \mathbf{I}$. Then

$$\begin{aligned}
&\mathcal{F}^{(K)} \mathbf{K}^{(K)} \mathcal{F}^{(K)H} \\
&\approx \sum_{j_K=1, \dots, j_1=1}^{Y_K, \dots, Y_1} \underbrace{(\mathcal{F}_{Y_K} \mathbf{C}^{(K,j_K)} \mathcal{F}_{Y_K}^H)}_{\Lambda^{(K,j_K)}} \otimes \dots \otimes \underbrace{(\mathcal{F}_{Y_1} \mathbf{C}^{(1,j_1)} \mathcal{F}_{Y_1}^H)}_{\Lambda^{(1,j_k)}} \otimes \mathbf{K}(\mathbf{t}_1, \mathbf{t}_{w(j)}) \quad (\text{E.2})
\end{aligned}$$

where all $\Lambda^{(i,j)}$ are diagonal because all $\mathbf{C}^{(i,j)}$ are circulant by definition. Thus $\mathcal{F}^{(K)} \mathbf{K}^{(K)} \mathcal{F}^{(K)H}$ is approximately block diagonal. Let

$$\mathbf{B}_{m_1, \dots, m_K} = \sum_{j_K=1, \dots, j_1=1}^{Y_K, \dots, Y_1} \Lambda_{m_K, m_K}^{(K,j_K)} \dots \Lambda_{m_1, m_1}^{(1,j_1)} \mathbf{K}(\mathbf{t}_1, \mathbf{t}_{w(j)}). \quad (\text{E.3})$$

Then Equation (E.2) tells us that $\mathbf{B}_{m_1, \dots, m_K}$ is the $d(m_1, \dots, m_K)$ 'th block on the diagonal of $\mathcal{F}^{(K)} \mathbf{K}^{(K)} \mathcal{F}^{(K)H}$ where

$$d(m_1, \dots, m_K) = (m_K - 1)(Y_{K-1} \dots Y_1) + \dots + (m_2 - 1)Y_1 + (m_1 - 1) + 1.$$

Finally, let $m^{(i)}(d)$ be such that $m^{(i)}[d(m_1, \dots, m_K)] = m_i$.¹ Then \mathbf{K}_F is approximated by

$$\mathbf{K}_F \approx \mathcal{F}^{(K)H} \text{diag}(\mathbf{B}_{m^{(1)}(1), \dots, m^{(K)}(1)}, \dots, \mathbf{B}_{m^{(1)}(Y), \dots, m^{(K)}(Y)}) \mathcal{F}^{(K)}. \quad (\text{E.4})$$

E.4 Approximating Determinants

We can use Equation (E.4) to efficiently approximate $|\mathbf{K}_F|$. We first derive an approximation of $|\mathbf{K}_F|$ and then show that this approximation can be computed efficiently.

Note that $\mathcal{F}^{(K)H} \mathcal{F}^{(K)} = \mathbf{I}$. Thus, by Equation (E.4),

$$\begin{aligned} |\mathbf{K}_F| &\approx |\mathcal{F}^{(K)H}| \left| \text{diag}(\mathbf{B}_{m^{(1)}(1), \dots, m^{(K)}(1)}, \dots, \mathbf{B}_{m^{(1)}(Y), \dots, m^{(K)}(Y)}) \right| |\mathcal{F}^{(K)}| \\ &= \prod_{d=1}^Y |\mathbf{B}_{m^{(1)}(d), \dots, m^{(K)}(d)}|. \end{aligned} \quad (\text{E.5})$$

We now show that Equation (E.5) can be computed efficiently. To begin with, we show that all $\mathbf{B}_{m_1, \dots, m_K}$ can be computed efficiently. Let $\mathbf{S}^{(0)}(\mathbf{t}) = \mathbf{K}(\mathbf{t}_1, \mathbf{t})$ and let

$$\begin{aligned} \mathbf{t}_{w(j)}^{(>i)} &= \left[t_{w^{(i+1)}(j_{i+1})}^{(i+1)} \quad \dots \quad t_{w^{(K)}(j_K)}^{(K)} \right]^T, \\ \mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)}) &= \sum_{j_i=1}^{Y_i} \Lambda_{m_i, m_i}^{(i, j_i)} \mathbf{S}_{m_1, \dots, m_{i-1}}^{(i-1)}(t_{w^{(i)}(j_i)}^{(i)}, \mathbf{t}_{w(j)}^{(>i)}) \end{aligned} \quad (\text{E.6})$$

for $i = 1, \dots, K$. Then one verifies that $\mathbf{S}_{m_1, \dots, m_K}^{(K)} = \mathbf{B}_{m_1, \dots, m_K}$. Now, fix i and fix

¹That is, let $m^{(i)}(d) = 1 + [(d-1)/(Y_{i-1} \cdots Y_i)] \bmod Y_i$. Then

$$\begin{aligned} m^{(i)}(d) \Big|_{d=d(m_1, \dots, m_K)} &= 1 + \left(\left\lfloor \frac{d-1}{Y_{i-1} \cdots Y_1} \right\rfloor \bmod Y_i \right) \Big|_{d=d(m_1, \dots, m_K)} \\ &= 1 + \left((m_K - 1)(Y_K \cdots Y_i) + \dots + (m_{i+1} - 1)Y_i \right. \\ &\quad \left. + (m_i - 1) + \frac{m_{i-1} - 1}{Y_{i-1}} + \dots + \frac{m_1 - 1}{Y_{i-1} \cdots Y_1} \right) \bmod Y_i \\ &= 1 + \left[(m_K - 1)(Y_K \cdots Y_i) + \dots + (m_{i+1} - 1)Y_i + m_i - 1 \right] \bmod Y_i \\ &= m_i. \end{aligned}$$

m_1, \dots, m_{i-1} and j_{i+1}, \dots, j_K . It holds that [Gray, 2006]

$$\begin{aligned} \Lambda_{m_i, m_i}^{(i, j_i)} &= (\mathcal{F}_{Y_i} \mathbf{C}^{(i, j_i)} \mathcal{F}_{Y_i}^H)_{m_i, m_i} \\ &= \sum_{n=1}^{Y_i} \mathbf{1}(j_i - n) \exp \left[-2\pi\sqrt{-1} \frac{(m_i - 1)(n - 1)}{Y_i} \right] \\ &= \exp \left[-2\pi\sqrt{-1} \frac{(m_i - 1)(j_i - 1)}{Y_i} \right]. \end{aligned}$$

Hence, by Equation (E.6),

$$\underbrace{\mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)})}_{\mathbf{A}_{m_i}} = \sum_{j_i=1}^{Y_i} \underbrace{\mathbf{S}_{m_1, \dots, m_{i-1}}^{(i-1)}(\mathbf{t}_{w^{(i)}(j_i)}^{(i)}, \mathbf{t}_{w(j)}^{(>i)})}_{\mathbf{B}_{j_i}} \exp \left[-2\pi\sqrt{-1} \frac{(m_i - 1)(j_i - 1)}{Y_i} \right].$$

Note that all \mathbf{A}_{m_i} and \mathbf{B}_{j_i} are $N \times N$ matrices. Importantly, observe that $\mathbf{A}_{:,n,n'}$ = DFT $\mathbf{B}_{:,n,n'}$. Thus, the fast Fourier transform algorithm can be used to compute a single $\mathbf{A}_{:,n,n'}$ in $\mathcal{O}(Y_i \log Y_i)$ time. By doing this for every (n, n') we can compute $\mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)})$ for all m_i in $\mathcal{O}(N^2 Y_i \log Y_i)$ time. Hence, we can compute $\mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)})$ for all m_1, \dots, m_i and j_{i+1}, \dots, j_K in

$$\mathcal{O}(Y_1 \cdots Y_{i-1} Y_{i+1} \cdots Y_K N^2 Y_i \log Y_i) = \mathcal{O}(N^2 Y \log Y_i)$$

time, which shows that recursively applying Equation (E.6) to compute $\mathbf{S}_{m_1, \dots, m_K}^{(K)} = \mathbf{B}_{m_1, \dots, m_K}$ for all m_1, \dots, m_K costs

$$\mathcal{O}[N^2 Y (\log Y_1 + \dots + \log Y_K)]$$

time. This time complexity is bounded by $\mathcal{O}(N^2 Y \log Y)$.

Finally, Equation (E.5) shows that after all $\mathbf{B}_{m_1, \dots, m_K}$ have been computed, $|\mathbf{K}_F|$ can be computed in $\mathcal{O}(N^3 Y)$ time (Appendix D). Therefore, the resulting complexity of computing $|\mathbf{K}_F|$ is $\mathcal{O}(N^2 Y \log Y + N^3 Y)$. The complete procedure to compute $|\mathbf{K}_F|$ is outlined in Algorithm 1.

Algorithm 1 Efficient approximation of the determinant of a multi-output stationary kernel matrix. Runs in $\mathcal{O}(N^2Y \log Y + N^3Y)$ time.

```

1: function DETERMINANT( $\mathbf{K}$ ,  $\{t_1^{(1)}, \dots, t_{Y_1}^{(1)}\} \times \dots \times \{t_1^{(K)}, \dots, t_{Y_K}^{(K)}\}$ )
2:   for  $i = 1, \dots, K$  do
3:     for  $m_1 = 1, \dots, Y_1$  to  $m_{i-1} = 1, \dots, Y_{i-1}$  and
        $j_{i+1} = 1, \dots, Y_{i+1}$  to  $j_K = 1, \dots, Y_K$  do
4:       for  $n = 1, \dots, N$  and  $n' = 1, \dots, N$  do
5:          $\mathbf{A}_{:,n,n'} \leftarrow \text{DFT } \mathbf{B}_{:,n,n'}$ 
                                      $\triangleright$  Via fast Fourier transform algorithm
                                      $\triangleright \mathbf{S}^{(0)}(\mathbf{t}) = \mathbf{K}(\mathbf{t}_1, \mathbf{t})$ ,  $\mathbf{A}_{m_i} = \mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)})$  and  $\mathbf{B}_{j_i} =$ 
                                      $\mathbf{S}_{m_1, \dots, m_{i-1}}^{(i-1)}(\mathbf{t}_{w^{(i)}(j_i)}^{(i)}, \mathbf{t}_{w(j)}^{(>i)})$  where  $w^{(i)}(j_i) = \min\{j, Y_i - j_i + 2\}$ 
6:   for  $m_1 = 1, \dots, Y_1$  to  $m_K = 1, \dots, Y_K$  do
7:     Compute  $|\mathbf{S}_{m_1, \dots, m_K}^{(K)}|$ 
8:   return  $\prod_{d=1}^Y |\mathbf{S}_{m^{(1)}(d), \dots, m^{(K)}(d)}^{(K)}|$     $\triangleright m^{(i)}(d) = 1 + \lfloor [(d-1)/(Y_{i-1} \dots Y_i)] \bmod Y_i \rfloor$ 

```

E.5 Approximating Products Involving an Inverse

We can also use Equation (E.4) to efficiently approximate a product of the form $\mathbf{P} = \mathbf{U}^T \mathbf{K}_F^{-1} \mathbf{V}$ for some $NY \times P$ matrix \mathbf{U} and $NY \times Q$ matrix \mathbf{V} . We first derive an approximation of \mathbf{P} and then show that this approximation can be computed efficiently.

Consider $P_{p,q}$. Recall that by Equation (E.4), $\mathcal{F}^{(K)} \mathbf{K}^{(K)} \mathcal{F}^{(K)H}$ is approximately block diagonal with blocks $\mathbf{B}_{m_1, \dots, m_K}$; hence $\mathcal{F}^{(K)} \mathbf{K}^{-(K)} \mathcal{F}^{(K)H}$ is approximately block diagonal with blocks $\mathbf{B}_{m_1, \dots, m_K}^{-1}$. Thus

$$\begin{aligned}
P_{p,q} &= \underbrace{(\mathbf{U}_{:,p})^T}_{\mathbf{u}^T} \mathbf{K}^{-(K)} \underbrace{\mathbf{V}_{:,q}}_{\mathbf{v}} \\
&= (\mathcal{F}^{(K)} \mathbf{u})^H (\mathcal{F}^{(K)} \mathbf{K}^{-(K)} \mathcal{F}^{(K)H}) (\mathcal{F}^{(K)} \mathbf{v}) \\
&\approx \sum_{d=1}^Y \left[u_{(d-1)N+1} \ \dots \ u_{dN} \right] \mathbf{B}_{m^{(1)}(d), \dots, m^{(K)}(d)}^{-1} \left[v_{(d-1)N+1} \ \dots \ v_{dN} \right]^T. \quad (\text{E.7})
\end{aligned}$$

We now show that Equation (E.7) can be computed efficiently. First, we show that all $\mathbf{B}_{m_1, \dots, m_K}^{-1}$ can be computed efficiently. Appendix E.4 showed that all $\mathbf{B}_{m_1, \dots, m_K}$ can be computed in $\mathcal{O}(N^2Y \log Y)$ time. Therefore all $\mathbf{B}_{m_1, \dots, m_K}^{-1}$ can be computed in $\mathcal{O}(N^2Y \log Y + N^3Y)$ time.

Second, we show that $\mathcal{F}^{(K)}\mathbf{u}$ and thereby $\mathcal{F}^{(K)}\mathbf{v}$ from Equation (E.7) can be computed efficiently. It holds that

$$\begin{aligned}\mathcal{F}^{(K)}\mathbf{u} &= (\mathcal{F}_{Y_K} \otimes \cdots \otimes \mathcal{F}_{Y_1} \otimes \mathbf{I})\mathbf{u} \\ &= \underbrace{(\mathcal{F}_{Y_K} \otimes \mathbf{I}_{Y_{K-1}} \otimes \cdots \otimes \mathbf{I}_{Y_1} \otimes \mathbf{I}_N)}_{\mathbf{T}_{Y_K}} \cdots \underbrace{(\mathbf{I}_{Y_K} \otimes \cdots \otimes \mathbf{I}_{Y_2} \otimes \mathcal{F}_{Y_1} \otimes \mathbf{I}_N)}_{\mathbf{T}_{Y_1}} \mathbf{u}.\end{aligned}$$

Observe that each \mathbf{T}_{Y_i} is of the form $\mathbf{I}_{Z_i^{(1)}} \otimes \mathcal{F}_{Y_i} \otimes \mathbf{I}_{Z_i^{(2)}}$ where $Z_i^{(1)} = Y_K \cdots Y_{i+1}$ and $Z_i^{(2)} = Y_{i-1} \cdots Y_1 N$. Inspection of

$$(\mathcal{F}_{Y_i} \otimes \mathbf{I}_{Z_i^{(2)}})\mathbf{x} = \begin{bmatrix} \mathcal{F}_{Y_i,1,1}\mathbf{I}_{Z_i^{(2)}} & \cdots & \mathcal{F}_{Y_i,1,Y_i}\mathbf{I}_{Z_i^{(2)}} \\ \vdots & \ddots & \vdots \\ \mathcal{F}_{Y_i,Y_i,1}\mathbf{I}_{Z_i^{(2)}} & \cdots & \mathcal{F}_{Y_i,Y_i,Y_i}\mathbf{I}_{Z_i^{(2)}} \end{bmatrix} \mathbf{x} = \mathbf{y}$$

shows that

$$\mathcal{F}_{Y_i} \underbrace{\begin{bmatrix} x_z & x_{z+Z_i^{(2)}} & \cdots & x_{z+(Y_i-1)Z_i^{(2)}} \end{bmatrix}^T}_{\text{every } Z_i^{(2)}\text{'th element of } \mathbf{x}, \text{ starting at } z} = \underbrace{\begin{bmatrix} y_z & y_{z+Z_i^{(2)}} & \cdots & y_{z+(Y_i-1)Z_i^{(2)}} \end{bmatrix}^T}_{\text{every } Z_i^{(2)}\text{'th element of } \mathbf{y}, \text{ starting at } z} \quad (\text{E.8})$$

for $z = 1, \dots, Z_i^{(2)}$, which completely specifies \mathbf{y} . The fast Fourier transform algorithm can be used to compute Equation (E.8) in $\mathcal{O}(Y_i \log Y_i)$ time, which means that \mathbf{y} can be computed in $\mathcal{O}(Z_i^{(2)} Y_i \log Y_i)$ time. Similarly, by directly applying the definition of the Kronecker product we see that $(\mathbf{I}_{Z_i^{(1)}} \otimes \mathcal{F}_{Y_i} \otimes \mathbf{I}_{Z_i^{(2)}})\mathbf{x} = \mathbf{y}$ yields the \mathbf{y} such that

$$\begin{aligned}(\mathcal{F}_{Y_i} \otimes \mathbf{I}_{Z_i^{(2)}}) &\underbrace{\begin{bmatrix} x_{(z-1)Z_i^{(2)}Y_i+1} & x_{(z-1)Z_i^{(2)}Y_i+2} & \cdots & x_{zZ_i^{(2)}Y_i} \end{bmatrix}^T}_{z\text{'th consecutive group of } Z_i^{(2)}Y_i \text{ consecutive elements of } \mathbf{x}} \\ &= \underbrace{\begin{bmatrix} y_{(z-1)Z_i^{(2)}Y_i+1} & y_{(z-1)Z_i^{(2)}Y_i+2} & \cdots & y_{zZ_i^{(2)}Y_i} \end{bmatrix}^T}_{z\text{'th consecutive group of } Z_i^{(2)}Y_i \text{ consecutive elements of } \mathbf{y}}\end{aligned}$$

for $z = 1, \dots, Z_{Y_i}^{(1)}$, which again completely specifies \mathbf{y} . Hence multiplication by any \mathbf{T}_{Y_i} can be done in $\mathcal{O}(Z_{Y_i}^{(1)} Z_{Y_i}^{(2)} Y_i \log Y_i) = \mathcal{O}(NY \log Y_i)$ time. Thus $\mathcal{F}^{(K)}\mathbf{u}$ and $\mathcal{F}^{(K)}\mathbf{v}$ can be computed in

$$\mathcal{O}[NY(\log Y_1 + \dots + \log Y_K)].$$

time. This complexity is bounded by $\mathcal{O}(NY \log Y)$.

Finally, Equation (E.7) shows that after all $\mathbf{B}_{m_1, \dots, m_K}^{-1}$, $\mathcal{F}^{(K)} \mathbf{u}$ and $\mathcal{F}^{(K)} \mathbf{v}$ have been computed, $P_{p,q}$ can be computed in $\mathcal{O}(N^2 Y)$ time. Therefore, since all $\mathbf{B}_{m_1, \dots, m_K}^{-1}$ have to be computed only once, the resulting complexity of computing \mathbf{P} is $\mathcal{O}[N^2 Y \log Y + N^3 Y + PQ(NY \log Y + N^2 Y)]$. The complete procedure to compute \mathbf{P} is outlined in Algorithm 2.

Algorithm 2 Efficient approximation of a product involving an inverse stationary multi-output kernel matrix. Runs in $\mathcal{O}[N^2 Y \log Y + N^3 Y + PQ(NY \log Y + N^2 Y)]$ time.

```

1: function PRODUCT( $\mathbf{K}$ ,  $\{t_1^{(1)}, \dots, t_{Y_1}^{(1)}\} \times \dots \times \{t_1^{(K)}, \dots, t_{Y_K}^{(K)}\}$ ,  $\mathbf{U}$ ,  $\mathbf{V}$ )
2:   for  $i = 1, \dots, K$  do
3:     for  $m_1 = 1, \dots, Y_1$  to  $m_{i-1} = 1, \dots, Y_{i-1}$  and
4:        $j_{i+1} = 1, \dots, Y_{i+1}$  to  $j_K = 1, \dots, Y_K$  do
5:         for  $n = 1, \dots, N$  and  $n' = 1, \dots, N$  do
6:            $\mathbf{A}_{:,n,n'} \leftarrow$  DFT  $\mathbf{B}_{:,n,n'}$ 
7:            $\triangleright$  Via fast Fourier transform algorithm
8:            $\triangleright \mathbf{S}^{(0)}(\mathbf{t}) = \mathbf{K}(\mathbf{t}_1, \mathbf{t})$ ,  $\mathbf{A}_{m_i} = \mathbf{S}_{m_1, \dots, m_i}^{(i)}(\mathbf{t}_{w(j)}^{(>i)})$  and  $\mathbf{B}_{j_i} =$ 
9:              $\mathbf{S}_{m_1, \dots, m_{i-1}}^{(i-1)}(\mathbf{t}_{w^{(i)}(j_i)}^{(i)}, \mathbf{t}_{w(j)}^{(>i)})$  where  $w^{(i)}(j_i) = \min\{j, Y_i - j_i + 2\}$ 
10:        for  $m_1 = 1, \dots, Y_1$  to  $m_K = 1, \dots, Y_K$  do
11:          Compute  $\mathbf{S}_{m_1, \dots, m_K}^{-K}$ 
12:        for  $p = 1, \dots, P$  and  $q = 1, \dots, Q$  do
13:           $\mathbf{U}_{:,p} \leftarrow$  TRANSFORM  $\mathbf{U}_{:,p}$ 
14:           $\mathbf{U}_{:,q} \leftarrow$  TRANSFORM  $\mathbf{U}_{:,q}$ 
15:           $P_{p,q} \leftarrow \sum_{d=1}^Y [U_{(d-1)N+1,p} \cdots U_{dN,p}] \mathbf{S}_{m^{(1)}(d), \dots, m^{(K)}(d)}^{-K} [V_{(d-1)N+1,q} \cdots V_{dN,q}]^T$ 
16:           $\triangleright m^{(i)}(d) = 1 + [(d-1)/(Y_{i-1} \cdots Y_i)] \bmod Y_i$ 
17:        return  $\mathbf{P}$ 
18: function TRANSFORM( $\mathbf{x}$ )
19:   for  $i = 1, \dots, K$  do
20:     for  $z_1 = 1, \dots, Z_1$  do  $\triangleright Z_1 = Y_K \cdots Y_{i+1}$ 
21:       for  $z_2 = 1, \dots, Z_2$  do  $\triangleright Z_2 = Y_{i-1} \cdots Y_1 N$ 
22:          $\begin{bmatrix} y_{z_2} & y_{z_2+Z_2} & \cdots & y_{z_2+(Y_i-1)Z_2} \end{bmatrix}^T$ 
23:          $\leftarrow Y_i^{-1/2}$  DFT  $\begin{bmatrix} y_{z_2} & y_{z_2+Z_2} & \cdots & y_{z_2+(Y_i-1)Z_2} \end{bmatrix}^T$ 
24:          $\triangleright$  Via fast Fourier transform algorithm
25:          $\triangleright \mathbf{y} = \begin{bmatrix} x_{(z_1-1)Z_2 Y_i + 1} & x_{(z_1-1)Z_2 Y_i + 2} & \cdots & x_{z_1 Z_2 Y_i} \end{bmatrix}^T$ 
26:   return  $\mathbf{x}$ 

```

E.6 Conclusion

We have derived procedures to efficiently approximate the determinant of a stationary multi-output kernel matrix and a product involving an inverse stationary multi-output kernel matrix. These procedures have respectively time complexities $\mathcal{O}(N^2Y \log Y + N^3Y)$ and $\mathcal{O}[N^2Y \log Y + N^3Y + PQ(NY \log Y + N^2Y)]$. Importantly, these complexities scale favourably in Y .

F | Exponentiated Quadratic Forms

F.1 Introduction

This chapter develops notation that makes working with exponentiated quadratic forms more convenient.

F.2 General Form

Notation 1 (Exponentiated Quadratic Form). *Let $\mathbf{x} \in \mathbb{R}^N$. Then denote an exponentiated quadratic form*

$$p(\mathbf{x}) = C \exp\left(-\frac{1}{2}\mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{x}^T \mathbf{b} + c\right)$$

by $(C, \mathbf{A}, \mathbf{b}, c)$.

We call \mathbf{x} the *composite vector of variables*, or *composite vector* in short.

Lemma 4 (Product Identity (General Form)). *Let p_1 and p_2 be two exponentiated quadratic forms. Then*

$$p_1 p_2 = (C_1, \mathbf{A}_1, \mathbf{b}_1, c_1)(C_2, \mathbf{A}_2, \mathbf{b}_2, c_2) = (C_1 C_2, \mathbf{A}_1 + \mathbf{A}_2, \mathbf{b}_1 + \mathbf{b}_2, c_1 + c_2).$$

Proof. Direct computation yields that

$$\begin{aligned} p_1 p_2 &= C_1 C_2 \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{A}_1 \mathbf{x} + \mathbf{x}^T \mathbf{b}_1 - \frac{1}{2} \mathbf{x}^T \mathbf{A}_2 \mathbf{x} + \mathbf{x}^T \mathbf{b}_2 + c\right) \\ &= C_1 C_2 \exp\left[-\frac{1}{2} \mathbf{x}^T (\mathbf{A}_1 + \mathbf{A}_2) \mathbf{x} + \mathbf{x}^T (\mathbf{b}_1 + \mathbf{b}_2) + (c_1 + c_2)\right] \\ &= (C_1 C_2, \mathbf{A}_1 + \mathbf{A}_2, \mathbf{b}_1 + \mathbf{b}_2, c_1 + c_2). \end{aligned}$$

□

Lemma 5 (Integration Identity (General Form)). *Let the composite vector be $[\mathbf{x}_1^T \quad \mathbf{x}_2^T]^T$ and consider*

$$p = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, c).$$

Denote integration over the subset of variables \mathbf{x}_1 by $I_{\mathbf{x}_1}(p)$:

$$I_{\mathbf{x}_1}(p)(\mathbf{x}_2) = \int (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, c) ([\mathbf{x}_1^T \quad \mathbf{x}_2^T]^T) d\mathbf{x}_1.$$

Then

$$I_{\mathbf{x}_1}(p) = \left(C \frac{(2\pi)^{N/2}}{|\mathbf{A}_{1,1}|^{1/2}}, \mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2}, \mathbf{b}_2 - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{b}_1, c + \frac{1}{2} \mathbf{b}_1^T \mathbf{A}_{1,1}^{-1} \mathbf{b}_1 \right).$$

Proof. An exponentiated quadratic form can be integrated as follows:

$$\begin{aligned} & \int_{\mathbb{R}^N} \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c\right) d\mathbf{x} \\ &= \exp\left(\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + c\right) \int_{\mathbb{R}^N} \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} - \frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}\right) d\mathbf{x} \\ &= \frac{(2\pi)^{N/2}}{|\mathbf{A}|^{1/2}} \exp\left(\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + c\right) \int_{\mathbb{R}^N} \frac{(2\pi)^{-N/2}}{|\mathbf{A}|^{-1/2}} \exp\left(-\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} - \frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}\right) d\mathbf{x} \\ &= \frac{(2\pi)^{N/2}}{|\mathbf{A}|^{1/2}} \exp\left(\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + c\right) \int_{\mathbb{R}^N} \mathcal{N}(\mathbf{x}; \mathbf{A}^{-1} \mathbf{b}, \mathbf{A}) d\mathbf{x} \\ &= \frac{(2\pi)^{N/2}}{|\mathbf{A}|^{1/2}} \exp\left(\frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b} + c\right). \end{aligned}$$

Now,

$$\begin{aligned} & -\frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + c \\ & = -\frac{1}{2} \mathbf{x}_1^T \mathbf{A}_{1,1} \mathbf{x}_1 + \mathbf{x}_1^T \left[\mathbf{b}_1 - \underbrace{\frac{1}{2} (\mathbf{A}_{1,2} + \mathbf{A}_{2,1}^T)}_{\mathbf{A}_{1,2}} \mathbf{x}_2 \right] + \left(\mathbf{x}_2^T \mathbf{b}_2 - \frac{1}{2} \mathbf{x}_2^T \mathbf{A}_{2,2} \mathbf{x}_2 + c \right). \end{aligned}$$

Thus

$$\begin{aligned} I_{\mathbf{x}_1}(p)(\mathbf{x}_2) &= C \frac{(2\pi)^{N/2}}{|\mathbf{A}_{1,1}|^{1/2}} \exp \left[\frac{1}{2} (\mathbf{b}_1 - \mathbf{A}_{1,2} \mathbf{x}_2)^T \mathbf{A}_{1,1}^{-1} (\mathbf{b}_1 - \mathbf{A}_{1,2} \mathbf{x}_2) \right. \\ &\quad \left. + \left(\mathbf{x}_2^T \mathbf{b}_2 - \frac{1}{2} \mathbf{x}_2^T \mathbf{A}_{2,2} \mathbf{x}_2 + c \right) \right] \\ &= C \frac{(2\pi)^{N/2}}{|\mathbf{A}_{1,1}|^{1/2}} \exp \left[-\frac{1}{2} \mathbf{x}_2^T (\mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2}) \mathbf{x}_2 \right. \\ &\quad \left. + \mathbf{x}_2^T (\mathbf{b}_2 - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{b}_1) + \left(c + \frac{1}{2} \mathbf{b}_1^T \mathbf{A}_{1,1}^{-1} \mathbf{b}_1 \right) \right] \end{aligned}$$

and so the result follows. \square

Lemma 6 (Linear Expansion Property (General Form)). *Let the composite vector be $\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T \end{bmatrix}^T$ and consider*

$$p_0 = \left(C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, c \right).$$

Then the exponentiated quadratic forms

$$\begin{aligned} p_1 \left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_3^T \end{bmatrix}^T \right) &= p_0 \left(\begin{bmatrix} (\mathbf{x}_1 + \mathbf{B} \mathbf{x}_3)^T & \mathbf{x}_2^T \end{bmatrix}^T \right), \\ p_2 \left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_3^T \end{bmatrix}^T \right) &= p_0 \left(\begin{bmatrix} \mathbf{x}_1^T & (\mathbf{x}_2 + \mathbf{B} \mathbf{x}_3)^T \end{bmatrix}^T \right) \end{aligned}$$

are given by

$$p_i = \left(C, \begin{bmatrix} \mathbf{A} & \mathbf{A}_{1,i} \mathbf{B} \\ \mathbf{B}^T \mathbf{A}_{i,1} & \mathbf{B}^T \mathbf{A}_{i,2} & \mathbf{B}^T \mathbf{A}_{i,i} \mathbf{B} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{B}^T \mathbf{b}_i \end{bmatrix}, c \right)$$

for respectively $i = 1, 2$.

Proof. The exponent of p_1 can be simplified as follows:

$$\begin{aligned}
& -\frac{1}{2} \begin{bmatrix} \mathbf{x}_1 + \mathbf{B}\mathbf{x}_3 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 + \mathbf{B}\mathbf{x}_3 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 + \mathbf{B}\mathbf{x}_3 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + c \\
&= -\frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} \mathbf{B}\mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \\
&\quad - \frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{B}\mathbf{x}_3 \\ \mathbf{0} \end{bmatrix} - \frac{1}{2} \begin{bmatrix} \mathbf{B}\mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{B}\mathbf{x}_3 \\ \mathbf{0} \end{bmatrix} \\
&\quad + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{B}\mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + c \\
&= -\frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} - \frac{1}{2} \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{0} & \mathbf{B}^T \mathbf{A}_{1,2} \\ \mathbf{0} & \mathbf{B}^T \mathbf{A}_{2,2} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \\
&\quad - \frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} \mathbf{B} & \mathbf{0} \\ \mathbf{A}_{2,1} \mathbf{B} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{0} \end{bmatrix} - \frac{1}{2} \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{B}^T \mathbf{A}_{1,1} \mathbf{B} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{0} \end{bmatrix} \\
&\quad + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_3 \\ \mathbf{0} \end{bmatrix}^T \begin{bmatrix} \mathbf{B}^T \mathbf{b}_1 \\ \mathbf{0} \end{bmatrix} + c \\
&= -\frac{1}{2} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,1} \mathbf{B} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,1} \mathbf{B} \\ \mathbf{B}^T \mathbf{A}_{1,1} & \mathbf{B}^T \mathbf{A}_{1,2} & \mathbf{B}^T \mathbf{A}_{1,1} \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{B}^T \mathbf{b}_1 \end{bmatrix} + c.
\end{aligned}$$

Thus

$$p_1 = (C, \begin{bmatrix} \mathbf{A} & \mathbf{A}_{1,1} \mathbf{B} \\ \mathbf{A}_{2,1} \mathbf{B} & \mathbf{B}^T \mathbf{A}_{1,1} \mathbf{B} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{B}^T \mathbf{b}_1 \end{bmatrix}, c).$$

The case p_2 follows similarly. □

Lemma 7 (Compression Property (General Form)). *Let the composite vector be $[\mathbf{x}_1^T \ \mathbf{x}_2^T \ \mathbf{x}_3^T]^T$ and consider*

$$p_1 = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} \\ \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{b}_3 \end{bmatrix}, c).$$

Then the exponentiated quadratic form

$$p_2\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T \end{bmatrix}^T\right) = p_1\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_2^T \end{bmatrix}^T\right)$$

is given by

$$p_2 = \left(C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} + \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} + \mathbf{A}_{3,1} & \mathbf{A}_{2,2} + \mathbf{A}_{2,3} + \mathbf{A}_{3,2} + \mathbf{A}_{3,3} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 + \mathbf{b}_3 \end{bmatrix}, c\right).$$

Proof. The exponent of p_1 can be simplified as follows:

$$\begin{aligned} & \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} \\ \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \mathbf{b}_3 \end{bmatrix} + c \\ &= \mathbf{x}_1^T \mathbf{A}_{1,1} \mathbf{x}_1 + \mathbf{x}_1^T (\mathbf{A}_{1,2} + \mathbf{A}_{1,3}) \mathbf{x}_2 + \mathbf{x}_2^T (\mathbf{A}_{2,1} + \mathbf{A}_{2,3}) \mathbf{x}_1 \\ &\quad + \mathbf{x}_2^T (\mathbf{A}_{2,2} + \mathbf{A}_{2,3} + \mathbf{A}_{3,2} + \mathbf{A}_{3,3}) \mathbf{x}_2 + \mathbf{x}_1^T \mathbf{b}_1 + \mathbf{x}_2^T (\mathbf{b}_2 + \mathbf{b}_3) + c \\ &= \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} + \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} + \mathbf{A}_{3,1} & \mathbf{A}_{2,2} + \mathbf{A}_{2,3} + \mathbf{A}_{3,2} + \mathbf{A}_{3,3} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} + \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}^T \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 + \mathbf{b}_3 \end{bmatrix} + c. \end{aligned}$$

Hence the result follows. \square

F.3 Kronecker-Structured Form

Notation 2 (Kronecker-Structured Exponentiated Quadratic Form). *Let the composite vector be $\begin{bmatrix} \mathbf{x}_1^T & \dots & \mathbf{x}_M^T \end{bmatrix}^T \in \mathbb{R}^{MN}$. Then denote a Kronecker-structured exponentiated quadratic form*

$$(C, \mathbf{A} \otimes \mathbf{I}_N, \mathbf{0}, 0)$$

by (C, \mathbf{A}) .

Lemma 8 (Product Identity (Kronecker-Structured Form)). *Let p_1 and p_2 be two Kronecker-structured exponentiated quadratic forms. Then*

$$p_1 p_2 = (C_1, \mathbf{A}_1)(C_2, \mathbf{A}_2) = (C_1 C_2, \mathbf{A}_1 + \mathbf{A}_2).$$

Proof. Follows from Lemma 4 and bilinearity of the Kronecker product. \square

Lemma 9 (Integration Identity (Kronecker-Structured Form)). *Let the composite vector be $[\mathbf{x}_1^T \ \mathbf{x}_2^T]^T$ where $\mathbf{x}_1 \in \mathbb{R}^{KN}$ and $\mathbf{x}_2 \in \mathbb{R}^{MN}$. Consider*

$$p = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}).$$

Denote integration over the subset of variables \mathbf{x}_1 by $I_{\mathbf{x}_1}(p)$:

$$I_{\mathbf{x}_1}(p)(\mathbf{x}_2) = \int (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}, \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix}, c) \left([\mathbf{x}_1^T \ \mathbf{x}_2^T]^T \right) d\mathbf{x}_1.$$

Then

$$I_{\mathbf{x}_1}(p) = \left(C \frac{(2\pi)^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}}, \mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2} \right).$$

Proof. Direct computation yields that

$$\begin{aligned} I_{\mathbf{x}_1} \left[(C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}) \right] &= I_{\mathbf{x}_1} \left[(C, \begin{bmatrix} \mathbf{A}_{1,1} \otimes \mathbf{I} & \mathbf{A}_{1,2} \otimes \mathbf{I} \\ \mathbf{A}_{2,1} \otimes \mathbf{I} & \mathbf{A}_{2,2} \otimes \mathbf{I} \end{bmatrix}, \mathbf{0}, 0) \right] \\ &= \left(C \frac{(2\pi)^{KN/2}}{|\mathbf{A}_{1,1} \otimes \mathbf{I}|^{1/2}}, \right. \\ &\quad \left. \mathbf{A}_{2,2} \otimes \mathbf{I} - (\mathbf{A}_{2,1} \otimes \mathbf{I})(\mathbf{A}_{1,1} \otimes \mathbf{I})^{-1}(\mathbf{A}_{1,2} \otimes \mathbf{I}), \mathbf{0}, 0 \right) \\ &= \left(C \frac{(2\pi)^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}}, (\mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2}) \otimes \mathbf{I}, \mathbf{0}, 0 \right) \\ &= \left(C \frac{(2\pi)^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}}, \mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2} \right). \end{aligned}$$

□

Lemma 10 (Linear Expansion Property (Kronecker-Structured Form)). *Let the composite vector be $[\mathbf{x}_1^T \ \mathbf{x}_2^T]^T$ and consider*

$$p_0 = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix}).$$

Then the exponentiated quadratic forms

$$\begin{aligned} p_1\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_3^T \end{bmatrix}^T\right) &= p_0\left(\begin{bmatrix} (\mathbf{x}_1 + \mathbf{B}\mathbf{x}_3)^T & \mathbf{x}_2^T \end{bmatrix}^T\right), \\ p_2\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_3^T \end{bmatrix}^T\right) &= p_0\left(\begin{bmatrix} \mathbf{x}_1^T & (\mathbf{x}_2 + \mathbf{B}\mathbf{x}_3)^T \end{bmatrix}^T\right) \end{aligned}$$

are given by

$$p_i = (C, \begin{bmatrix} \mathbf{A} & \mathbf{A}_{1,i}\mathbf{B} \\ \mathbf{B}^T\mathbf{A}_{i,1} & \mathbf{B}^T\mathbf{A}_{i,2} & \mathbf{B}^T\mathbf{A}_{i,i}\mathbf{B} \end{bmatrix})$$

for respectively $i = 1, 2$.

Proof. Follows from Lemma 6 and bilinearity of the Kronecker product. \square

Lemma 11 (Compression Property (Kronecker-Structured Form)). *Let the composite vector be $\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_2^T \end{bmatrix}^T$ and consider*

$$p_1 = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} & \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \mathbf{A}_{2,3} \\ \mathbf{A}_{3,1} & \mathbf{A}_{3,2} & \mathbf{A}_{3,3} \end{bmatrix}).$$

Then the exponentiated quadratic form

$$p_2\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T \end{bmatrix}^T\right) = p_1\left(\begin{bmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \mathbf{x}_2^T \end{bmatrix}^T\right)$$

is given by

$$p_2 = (C, \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} + \mathbf{A}_{1,3} \\ \mathbf{A}_{2,1} + \mathbf{A}_{3,1} & \mathbf{A}_{2,2} + \mathbf{A}_{2,3} + \mathbf{A}_{3,2} + \mathbf{A}_{3,3} \end{bmatrix}).$$

Proof. Follows from Lemma 7 and bilinearity of the Kronecker product. \square

Finally, we prove one additional identity.

Lemma 12 (Direct Integration Identity (Kronecker-Structured Form)). *Let the com-*

posite vector be $[\mathbf{x}_1^T \ \mathbf{x}_2^T]^T$ where $\mathbf{x}_1 \in \mathbb{R}^{KN}$ and $\mathbf{x}_2 = [\mathbf{y}_1^T \ \cdots \ \mathbf{y}_M^T]^T \in \mathbb{R}^{MN}$. Then

$$I_{\mathbf{x}_1}[(C, 2 \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix})](\mathbf{x}_2) = C \frac{\pi^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}} \exp \left\{ - \sum_{m_1=1, m_2=1}^{M, M} A_{2,2, m_1, m_2} \mathbf{y}_{m_1}^T \mathbf{y}_{m_2} \right. \\ \left. + \sum_{k_1=1, k_2=1}^{K, K} A_{1,1, k_1, k_2}^{-1} \left(\sum_{m=1}^M A_{2,1, k_1, m} \mathbf{y}_m \right)^T \right. \\ \left. \left(\sum_{m=1}^M A_{2,1, k_2, m} \mathbf{y}_m \right) \right\}.$$

Proof. Direct computation yields that

$$I_{\mathbf{x}_1}[(C, 2 \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{A}_{1,2} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} \end{bmatrix})] = \left(C \frac{\pi^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}}, 2(\mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2}) \right) \\ = C \frac{\pi^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}} \exp \{ -\mathbf{x}_2^T [(\mathbf{A}_{2,2} - \mathbf{A}_{2,1} \mathbf{A}_{1,1}^{-1} \mathbf{A}_{1,2}) \otimes \mathbf{I}_N] \mathbf{x}_2 \} \\ = C \frac{\pi^{KN/2}}{|\mathbf{A}_{1,1}|^{N/2}} \exp \{ -\mathbf{x}_2^T (\mathbf{A}_{2,2} \otimes \mathbf{I}_N) \mathbf{x}_2 \\ + \mathbf{x}_2^T (\mathbf{A}_{2,1} \otimes \mathbf{I}_N) (\mathbf{A}_{1,1}^{-1} \otimes \mathbf{I}_N) \\ [(\mathbf{A}_{1,2} \otimes \mathbf{I}_N) \mathbf{x}_2] \}.$$

Hence the result follows. \square

F.4 Conclusion

We have developed notation that makes working with exponentiated quadratic forms more convenient. Importantly, the notation reduces manipulation of exponentiated quadratic forms to relatively simple matrix operations.

G | Roots of Kernels

First, using the notation and identities from Appendix F with composite vector $[\mathbf{t}^T \quad \boldsymbol{\tau}^T]^T$, any kernel of exponentiated quadratic form has a root:

$$\begin{aligned}
 & R \left[\left(\sigma \frac{|2\mathbf{A}|^{1/4}}{\pi^{K/4}}, 2\mathbf{A}, \mathbf{0}, 0 \right) \right] * \left(\sigma \frac{|2\mathbf{A}|^{1/4}}{\pi^{K/4}}, 2\mathbf{A}, \mathbf{0}, 0 \right) \\
 &= I_{\boldsymbol{\tau}} \left[\left(\sigma \frac{|2\mathbf{A}|^{1/4}}{\pi^{K/4}}, 2 \begin{bmatrix} \mathbf{A} & -\mathbf{A} \\ -\mathbf{A} & \mathbf{A} \end{bmatrix}, \mathbf{0}, 0 \right) \left(\sigma \frac{|2\mathbf{A}|^{1/4}}{\pi^{K/4}}, 2 \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{A} \end{bmatrix}, \mathbf{0}, 0 \right) \right] \\
 &= I_{\boldsymbol{\tau}} \left[\left(\sigma^2 \frac{|2\mathbf{A}|^{1/2}}{\pi^{K/2}}, 2 \begin{bmatrix} \mathbf{A} & -\mathbf{A} \\ -\mathbf{A} & 2\mathbf{A} \end{bmatrix}, \mathbf{0}, 0 \right) \right] \\
 &= \left(\sigma^2 \frac{|2\mathbf{A}|^{1/2}}{\pi^{K/2}} \frac{\pi^{K/2}}{|2\mathbf{A}|^{1/2}}, 2[\mathbf{A} - \mathbf{A}(2\mathbf{A})^{-1}\mathbf{A}], \mathbf{0}, 0 \right) \\
 &= (\sigma^2, \mathbf{A}, \mathbf{0}, 0).
 \end{aligned}$$

Second, consider the diagonal multi-output exponentiated-quadratic kernel. Its diagonal entries are of exponentiated quadratic form and thus have roots. Now, for any diagonal \mathbf{H} it holds that

$$R(\mathbf{H}) * \mathbf{H}^T = \text{diag}[R(H_{1,1}) * H_{1,1}, \dots, R(H_{N,N})].$$

Thus the diagonal entries having roots implies that the diagonal multi-output exponentiated-quadratic kernel has a root. Specifically, if

$$\mathcal{K}(\mathbf{t}_1, \mathbf{t}_2) = \sigma^2 \exp(-\gamma \|\mathbf{t}_1 - \mathbf{t}_2\|^2) \mathbf{I},$$

then $\mathcal{K}(\mathbf{t}_1, \mathbf{t}_2) = [R(\mathbf{R}) * \mathbf{R}](\mathbf{t}_2 - \mathbf{t}_1)$ where

$$\mathbf{R}(\mathbf{t}) = \sigma \frac{(4\gamma)^{K/4}}{\pi^{K/4}} \exp(-2\gamma\|\mathbf{t}\|^2) \mathbf{I}.$$

Finally, the white noise kernel is its own root:

$$[R(\delta\mathbf{I}) * \delta\mathbf{I}](\mathbf{t}) = \int_{\mathbb{R}^K} \delta(\mathbf{t}' - \mathbf{t}) \mathbf{I} \delta(\mathbf{t}') \mathbf{I} d\mathbf{t}' = \delta(\mathbf{t}) \mathbf{I}.$$

H | Approximate Kernel Model

We derive an analytical expression for $\tilde{\mathcal{K}}_{f|H(\mathbf{T})}$. Recall that \mathbf{H} 's components are modelled independently. Then

$$\tilde{\mathcal{K}}_{f_i, f_j | \mathbf{H}(\mathbf{T})} = \sum_{k=1}^M \int_{\mathbb{R}^K} \mathbb{E}[H_{i,k}(\mathbf{t}_1 - \boldsymbol{\tau}) H_{j,k}(\mathbf{t}_2 - \boldsymbol{\tau}) | \mathbf{H}(\mathbf{T})] d\boldsymbol{\tau}$$

where, denoting $h_i = H_{i,k}$, $h_j = H_{j,k}$, $\mathbf{u}_{h_i} = h_i(\mathbf{T})$ and $\mathbf{u}_{h_j} = h_j(\mathbf{T})$,

$$\begin{aligned} & \int_{\mathbb{R}^K} \mathbb{E}[H_{i,k}(\mathbf{t}_1 - \boldsymbol{\tau}) H_{j,k}(\mathbf{t}_2 - \boldsymbol{\tau}) | \mathbf{H}(\mathbf{T})] d\boldsymbol{\tau} \\ &= \int_{\mathbb{R}^K} \{ \mathcal{K}_{h_i, h_j | \mathbf{u}_{h_i}, \mathbf{u}_{h_j}}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{t}_2 - \boldsymbol{\tau}) \\ & \quad - \mathbb{E}[h_i(\mathbf{t}_1 - \boldsymbol{\tau}) | \mathbf{u}_{h_i}] \mathbb{E}[h_j(\mathbf{t}_2 - \boldsymbol{\tau}) | \mathbf{u}_{h_j}] \} d\boldsymbol{\tau} \\ &= \int_{\mathbb{R}^{2K}} [\mathbb{1}(i - j) \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{t}_2 - \boldsymbol{\tau}) \\ & \quad + \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{T}) \mathbf{M}^{(h_i, h_j)} \mathcal{K}_{h_j}(\mathbf{T}, \mathbf{t}_2 - \boldsymbol{\tau})] d\boldsymbol{\tau} \\ &= \mathbb{1}(i - j) \underbrace{\int_{\mathbb{R}^K} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{t}_2 - \boldsymbol{\tau}) d\boldsymbol{\tau}}_{I^{(1, h_i)}(\mathbf{t}_1, \mathbf{t}_2)} \end{aligned} \tag{H.1}$$

$$+ \sum_{m=1, n=1}^{T, T} M_{m, n}^{(h_i, h_j)} \underbrace{\int_{\mathbb{R}^K} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{T}_{m, :}) \mathcal{K}_{h_j}(\mathbf{T}_{n, :}, \mathbf{t}_2 - \boldsymbol{\tau}) d\boldsymbol{\tau}}_{I_{m, n}^{(2, h_i, h_j)}(\mathbf{t}_1, \mathbf{t}_2)} \tag{H.2}$$

and

$$\mathbf{M}^{(h_i, h_j)} = \mathbf{K}_{\mathbf{u}_{h_i}}^{-1} \mathbf{u}_{h_i} \mathbf{u}_{h_j}^T \mathbf{K}_{\mathbf{u}_{h_j}}^{-1} - \mathbb{1}(i - j) \mathbf{K}_{\mathbf{u}_{h_i}}^{-1}.$$

Hence,

$$\begin{aligned}
\tilde{\mathcal{K}}_{f_i, f_j | \mathbf{H}(\mathbf{T})} &= \sum_{k=1}^M \left[\mathbb{1}(i - j) I^{(1, H_{i,k})} + \text{tr}(\mathbf{M}^{(H_{i,k}, H_{j,k})T} \mathbf{I}^{(2, H_{i,k}, H_{j,k})}) \right] \\
&= \sum_{k=1}^M \left\{ \mathbb{1}(i - j) [I^{(1, H_{i,k})} - \text{tr}(\mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{i,k})})] \right. \\
&\quad \left. + H_{i,k}^T(\mathbf{T}) \mathbf{K}_{H_{i,k}(\mathbf{T})}^{-1} \mathbf{I}^{(2, H_{i,k}, H_{j,k})} \mathbf{K}_{H_{j,k}(\mathbf{T})}^{-1} H_{j,k}(\mathbf{T}) \right\}.
\end{aligned}$$

I | Variational Free Energy of the Nonparametric Kernel Model

I.1 Introduction

This chapter derives an analytical expression for the variational free energy of Model 5. We also determine the asymptotic time complexity of computing the free energy.

Recall that the free energy is given by

$$\begin{aligned} \mathcal{F}(q) = & \underbrace{-\frac{1}{2} \log[(2\pi)^N |\mathbf{\Lambda}|^2] - \frac{1}{2} \sum_{i=1}^Y \mathbb{E}_q \{ \|\mathbf{\Lambda}^{-1} [\mathbf{Y}_{i,:} - \mathbf{y}(\mathbf{T}_{i,:})]\|^2 \}}_{\mathbb{E}_q[\log p(\mathbf{Y} | \mathbf{H}, \mathbf{x})]} \\ & - \sum_{i=1, j=1}^{N, M} D_{KL}[q(\mathbf{u}_{H_{i,j}}) \| p(\mathbf{u}_{H_{i,j}})] - \sum_{j=1}^M D_{KL}[q(\mathbf{u}_{\tilde{x}_j}) \| p(\mathbf{u}_{\tilde{x}_j})]. \end{aligned}$$

In more detail, the Kullback-Leibler divergences are given by (Appendix B.2)

$$D_{KL}[q(\mathbf{u}_{H_{i,j}}) \| p(\mathbf{u}_{H_{i,j}})] = \frac{1}{2} \left[\text{tr}(\mathbf{K}_{\mathbf{u}_{H_{i,j}}}^{-1} \mathbf{\Sigma}_{H_{i,j}}) + \boldsymbol{\mu}_{H_{i,j}}^T \mathbf{K}_{\mathbf{u}_{H_{i,j}}}^{-1} \boldsymbol{\mu}_{H_{i,j}} + \log \frac{|\mathbf{K}_{\mathbf{u}_{H_{i,j}}}|}{|\mathbf{\Sigma}_{H_{i,j}}|} - N \right], \quad (\text{I.1})$$

$$D_{KL}[q(\mathbf{u}_{\tilde{x}_j}) \| p(\mathbf{u}_{\tilde{x}_j})] = \frac{1}{2} \left[\text{tr}(\mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}^{-1} \mathbf{\Sigma}_{\tilde{x}_j}) + \boldsymbol{\mu}_{\tilde{x}_j}^T \mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}^{-1} \boldsymbol{\mu}_{\tilde{x}_j} + \log \frac{|\mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}|}{|\mathbf{\Sigma}_{\tilde{x}_j}|} - N \right] \quad (\text{I.2})$$

and the likelihood is given by

$$\begin{aligned} \mathbb{E}_q[\log p(\mathbf{Y} | \mathbf{H}, \mathbf{x})] = & -\frac{1}{2} \log[(2\pi)^N |\mathbf{\Lambda}|^2] - \frac{1}{2} \sum_{i=1}^Y \left\{ \mathbf{Y}_{i,:}^T \mathbf{\Lambda}^{-2} \mathbf{Y}_{i,:} - 2 \mathbf{Y}_{i,:}^T \mathbf{\Lambda}^{-2} \mathbb{E}_q[\mathbf{y}(\mathbf{T}_{i,:})] \right. \\ & \left. + \text{tr} \{ \mathbf{\Lambda}^{-2} \mathbb{E}_q[\mathbf{y}(\mathbf{T}_{i,:}) \mathbf{y}^T(\mathbf{T}_{i,:})] \} \right\}. \end{aligned}$$

It thus remains to determine the predictive mean $\mathbb{E}_q[\mathbf{y}(\mathbf{t})]$ and predictive autocovariance $\mathbb{E}_q[\mathbf{y}(\mathbf{t}_1)\mathbf{y}^T(\mathbf{t}_2)]$. Note that the latter is only evaluated on its diagonal.

Recall that $H_{i,j}$ and \tilde{x}_j denote the number of inducing points for the processes $T_{H_{i,j}}$ and $T_{\tilde{x}_j}$ respectively. Also recall that $T_H = \max_{i,j} T_{H_{i,j}}$ and $T_{\tilde{x}} = \max_j T_{\tilde{x}_j}$. Finally, recall that Y denotes the number of observed data points. Then Equations (I.1) and (I.2) have respectively time complexities $\mathcal{O}(T_H^3)$ and $\mathcal{O}(T_{\tilde{x}}^3)$ (Appendix D). Thus $\mathcal{F}(q)$ has time complexity $\mathcal{O}[NMT_H^3 + MT_{\tilde{x}}^3 + YN(C_\mu + C_\Sigma)]$ where C_μ and C_Σ represent the number of operations required to compute respectively $\mathbb{E}_q[y_i(\mathbf{t})]$ and $\mathbb{E}_q[y_i(\mathbf{t})y_i(\mathbf{t})]$.

In the remainder of this chapter we utilise the notation and identities from Appendix F without further notion.

I.2 Predictive Mean

We have that

$$\mathbb{E}_q[y_i(\mathbf{t})] = \mathbb{E}_q\{[(\mathbf{H} * \mathbf{x})(\mathbf{t})]_i\} = \sum_{j=1}^M \underbrace{\int_{\mathbb{R}^K} \mathbb{E}_q[H_{i,j}(\mathbf{t} - \boldsymbol{\tau})x_j(\boldsymbol{\tau})] d\boldsymbol{\tau}}_{T_{i,j}^{(1)}}$$

where, denoting $h = H_{i,j}$ and $x = x_j$,

$$\begin{aligned} T_{i,j}^{(1)} &= \int h(\mathbf{t} - \boldsymbol{\tau})p(h | \mathbf{u}_h) dh q(\mathbf{u}_h) d\mathbf{u}_h \\ &\quad x(\boldsymbol{\tau})p(x | \mathbf{u}_{\tilde{x}}) dx q(\mathbf{u}_{\tilde{x}}) d\mathbf{u}_{\tilde{x}} d\boldsymbol{\tau} \\ &= \int_{\mathbb{R}^2} \mathcal{K}_h(\mathbf{t} - \boldsymbol{\tau}, T_h) \mathbf{K}_{\mathbf{u}_h}^{-1} \boldsymbol{\mu}_h \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, T_{\tilde{x}}) \mathbf{K}_{\mathbf{u}_{\tilde{x}}}^{-1} \boldsymbol{\mu}_{\tilde{x}} d\boldsymbol{\tau} \\ &= \boldsymbol{\mu}_h^T \mathbf{K}_{\mathbf{u}_h}^{-1} \underbrace{\int_{\mathbb{R}^K} \mathcal{K}_h(T_h, \mathbf{t} - \boldsymbol{\tau}) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, T_{\tilde{x}}) d\boldsymbol{\tau}}_{\mathbf{I}^{(L,h,x)}(\mathbf{t})} \mathbf{K}_{\mathbf{u}_{\tilde{x}}}^{-1} \boldsymbol{\mu}_{\tilde{x}}. \end{aligned}$$

In summary,

$$\mathbb{E}_q[y_i(\mathbf{t})] = \sum_{j=1}^M \boldsymbol{\mu}_{H_{i,j}}^T \mathbf{K}_{\mathbf{u}_{H_{i,j}}}^{-1} \mathbf{I}^{(L,H_{i,j},x_j)}(\mathbf{t}) \mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}^{-1} \boldsymbol{\mu}_{\tilde{x}_j}.$$

The integral $\mathbf{I}^{(L,\cdot,\cdot)}$ is computed in Appendix I.4.

To begin with, there is a one-time cost of constructing and inverting all $\mathbf{K}_{\mathbf{u}_{H_{i,j}}}$ and $\mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}$, which has time complexity $\mathcal{O}[M(T_H^3 + KT_H^2 + KT_{\tilde{x}}^2 + T_{\tilde{x}}^3)]$. Afterwards, all $\mathbf{I}^{(L,H_{i,j},\tilde{x}_j)}$ are

constructed in $\mathcal{O}(MKT_H T_{\tilde{x}})$ time and the predictive mean is computed in $\mathcal{O}[M(T_H^2 + T_H T_{\tilde{x}} + T_{\tilde{x}}^2)]$. Hence $\mathcal{O}(C_\mu) = \mathcal{O}[M(T_H^2 + KT_H T_{\tilde{x}} + T_{\tilde{x}}^2)]$.

I.3 Predictive Autocovariance

We have that

$$\begin{aligned}
\mathbb{E}_q[y_i(\mathbf{t}_1)y_j(\mathbf{t}_2)] &= \sum_{k=1, l=1}^{M, M} \int_{\mathbb{R}^{2K}} \mathbb{E}_q[H_{i,k}(\mathbf{t}_1 - \boldsymbol{\tau}_1)x_k(\boldsymbol{\tau}_1)H_{j,l}(\mathbf{t}_2 - \boldsymbol{\tau}_2)x_l(\boldsymbol{\tau}_2)] d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\
&= \sum_{k=1}^M \int_{\mathbb{R}^K} \mathbb{E}_q[H_{i,k}(\mathbf{t}_1 - \boldsymbol{\tau}_1)x_k(\boldsymbol{\tau})] d\boldsymbol{\tau} \sum_{l=1, l \neq k}^M \int_{\mathbb{R}^K} \mathbb{E}_q[H_{j,l}(\mathbf{t}_2 - \boldsymbol{\tau})x_l(\boldsymbol{\tau})] d\boldsymbol{\tau} \\
&\quad + \underbrace{\sum_{k=1}^M \int_{\mathbb{R}^{2K}} \mathbb{E}_q[H_{i,k}(\mathbf{t}_1 - \boldsymbol{\tau}_1)x_k(\boldsymbol{\tau}_1)H_{j,k}(\mathbf{t}_2 - \boldsymbol{\tau}_2)x_k(\boldsymbol{\tau}_2)] d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2}_{T_{i,j,k}^{(2)}} \\
&= \sum_{k=1}^M T_{i,k}^{(1)} \sum_{l=1, l \neq k}^M T_{j,l}^{(1)} + \sum_{k=1}^M T_{i,j,k}^{(2)}.
\end{aligned}$$

Then, denoting $h_i = H_{i,k}$, $h_j = H_{j,k}$ and $x = x_k$,

$$\begin{aligned}
T_{i,j,k}^{(2)} &= \int h_i(\mathbf{t}_1 - \boldsymbol{\tau}_1)h_j(\mathbf{t}_2 - \boldsymbol{\tau}_2)p(h_i, h_j | \mathbf{u}_{h_i}, \mathbf{u}_{h_j}) d(h_i, h_j)q(\mathbf{u}_{h_i}, \mathbf{u}_{h_j}) d(\mathbf{u}_{h_i}, \mathbf{u}_{h_j}) \\
&\quad x(\boldsymbol{\tau}_1)x(\boldsymbol{\tau}_2)p(x | \mathbf{u}_{\tilde{x}}) dxq(\mathbf{u}_{\tilde{x}}) d\mathbf{u}_{\tilde{x}} d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\
&= \int_{\mathbb{R}^{2K}} [\mathbb{1}(i - j)\mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{t}_2 - \boldsymbol{\tau}_2) + \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{T}_{h_i})\mathbf{M}^{(h_i, h_j)}\mathcal{K}_{h_j}(\mathbf{T}_{h_j}, \mathbf{t}_2 - \boldsymbol{\tau}_2)] \\
&\quad [\mathcal{K}_x(\boldsymbol{\tau}_1, \boldsymbol{\tau}_2) + \mathcal{K}_{x, \tilde{x}}(\boldsymbol{\tau}_1, \mathbf{T}_{\tilde{x}})\mathbf{M}^{(\tilde{x})}\mathcal{K}_{\tilde{x}, x}(\mathbf{T}_{\tilde{x}}, \boldsymbol{\tau}_2)] d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2
\end{aligned}$$

where

$$\begin{aligned}
\mathbf{M}^{(h_i, h_j)} &= \mathbf{K}_{\mathbf{u}_{h_i}}^{-1} [\mathbb{1}(i - j)\boldsymbol{\Sigma}_{h_i} + \boldsymbol{\mu}_{h_i}\boldsymbol{\mu}_{h_j}^T] \mathbf{K}_{\mathbf{u}_{h_j}}^{-1} - \mathbb{1}(i - j)\mathbf{K}_{\mathbf{u}_{h_i}}^{-1}, \\
\mathbf{M}^{(\tilde{x})} &= \mathbf{K}_{\mathbf{u}_{\tilde{x}}}^{-1} [\boldsymbol{\Sigma}_{\tilde{x}} + \boldsymbol{\mu}_{\tilde{x}}\boldsymbol{\mu}_{\tilde{x}}^T] \mathbf{K}_{\mathbf{u}_{\tilde{x}}}^{-1} - \mathbf{K}_{\mathbf{u}_{\tilde{x}}}^{-1}.
\end{aligned}$$

Thus $T_{i,j,k}^{(2)} = Q_1 + Q_2 + Q_3 + Q_4$ where

$$Q_1 = \mathbb{1}(i - j) \underbrace{\int_{\mathbb{R}^{2K}} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{t}_2 - \boldsymbol{\tau}_2) \mathcal{K}_x(\boldsymbol{\tau}_1, \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2}_{I^{(Q_1, h_i, x)}(\mathbf{t}_1, \mathbf{t}_2)},$$

$$Q_2 = \mathbb{1}(i - j) \sum_{m,n} M_{m,n}^{(\tilde{x})} \underbrace{\int_{\mathbb{R}^{2K}} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{t}_2 - \boldsymbol{\tau}_2) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}_1, \mathbf{T}_{\tilde{x},m,:}) \mathcal{K}_{\tilde{x},x}(\mathbf{T}_{\tilde{x},n,:}, \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2}_{I_{m,\tilde{n}}^{(Q_2, h_i, x)}(\mathbf{t}_1, \mathbf{t}_2)},$$

$$Q_3 = \sum_{m,n} M_{m,n}^{(h_i, h_j)} \underbrace{\int_{\mathbb{R}^{2K}} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{T}_{h_i,m,:}) \mathcal{K}_{h_j}(\mathbf{T}_{h_j,n,:}, \mathbf{t}_2 - \boldsymbol{\tau}_2) \mathcal{K}_x(\boldsymbol{\tau}_1, \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2}_{I_{m,n}^{(Q_3, h_i, h_j, x)}(\mathbf{t}_1, \mathbf{t}_2)},$$

$$\begin{aligned} Q_4 &= \sum_{m,n,o,p} M_{m,n}^{(h_i, h_j)} M_{o,p}^{(\tilde{x})} \int_{\mathbb{R}^{2K}} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{T}_{h_i,m,:}) \mathcal{K}_{h_j}(\mathbf{T}_{h_j,n,:}, \mathbf{t}_2 - \boldsymbol{\tau}_2) \\ &\quad \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}_1, \mathbf{T}_{\tilde{x},o,:}) \mathcal{K}_{\tilde{x},x}(\mathbf{T}_{\tilde{x},p,:}, \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\ &= \sum_{m,n,o,p} M_{m,n}^{(h_i, h_j)} M_{o,p}^{(\tilde{x})} \int_{\mathbb{R}^K} \mathcal{K}_{h_i}(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{T}_{h_i,m,:}) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, \mathbf{T}_{\tilde{x},o,:}) d\boldsymbol{\tau} \\ &\quad \underbrace{\int_{\mathbb{R}^K} \mathcal{K}_{h_j}(\mathbf{t}_2 - \boldsymbol{\tau}, \mathbf{T}_{h_j,n,:}) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, \mathbf{T}_{\tilde{x},p,:}) d\boldsymbol{\tau}}_{I_{n,p}^{(Q_4, h_j, x)}(\mathbf{t}_2)} \\ &= \sum_{m,n,o,p} M_{m,n}^{(h_i, h_j)} M_{o,p}^{(\tilde{x})} I_{m,o}^{(Q_4, h_i, x)}(\mathbf{t}_1) I_{n,p}^{(Q_4, h_j, x)}(\mathbf{t}_2). \end{aligned}$$

In summary,

$$\begin{aligned}
\mathbb{E}_q[y_i(\mathbf{t}_1)y_j(\mathbf{t}_2)] &= \sum_{k=1}^M \underbrace{\boldsymbol{\mu}_{H_{i,k}}^T \mathbf{K}_{\mathbf{u}_{H_{i,k}}}^{-1} \mathbf{I}^{(L,H_{i,k},x_k)}(\mathbf{t}_1) \mathbf{K}_{\mathbf{u}_{\tilde{x}_k}}^{-1} \boldsymbol{\mu}_{\tilde{x}_k}}_{T^{(1)}} \\
&\quad + \underbrace{\sum_{l=1, l \neq k}^M \boldsymbol{\mu}_{H_{j,l}}^T \mathbf{K}_{\mathbf{u}_{H_{j,l}}}^{-1} \mathbf{I}^{(L,H_{j,l},x_l)}(\mathbf{t}_2) \mathbf{K}_{\mathbf{u}_{\tilde{x}_l}}^{-1} \boldsymbol{\mu}_{\tilde{x}_l}}_{T^{(1)}} \\
&\quad + \mathbb{1}(i-j) \sum_{k=1}^M I^{(Q_1, H_{i,k}, x_k)}(\mathbf{t}_1, \mathbf{t}_2) \\
&\quad + \mathbb{1}(i-j) \sum_{k=1}^M \text{tr}[\mathbf{M}^{(\tilde{x}_k)T} \mathbf{I}^{(Q_2, H_{i,k}, x_k)}(\mathbf{t}_1, \mathbf{t}_2)] \\
&\quad + \sum_{k=1}^M \text{tr}[\mathbf{M}^{(H_{i,k})T} \mathbf{I}^{(Q_3, H_{i,k}, H_{j,k}, x_k)}(\mathbf{t}_1, \mathbf{t}_2)] \\
&\quad + \underbrace{\sum_{k=1}^M \text{tr}[\mathbf{M}^{(H_{i,k}, H_{j,k})T} \mathbf{I}^{(Q_4, H_{i,k}, x_k)}(\mathbf{t}_1) \mathbf{M}^{(\tilde{x}_k)} \mathbf{I}^{(Q_4, H_{j,k}, x_k)T}(\mathbf{t}_2)]}_{T^{(2)}}.
\end{aligned}$$

The integrals $I^{(\cdot)}$ are computed in Appendix I.4.

To begin with, there is a one-time cost of constructing and inverting all $\mathbf{K}_{\mathbf{u}_{H_{i,j}}}$ and $\mathbf{K}_{\mathbf{u}_{\tilde{x}_j}}$, which has time complexity $\mathcal{O}[M(T_H^3 + KT_H^2 + KT_{\tilde{x}}^2 + T_{\tilde{x}}^3)]$, and a one-time cost of computing all $\mathbf{M}^{(H_{i,j}, H_{i,j})}$ and $\mathbf{M}^{(\tilde{x}_j)}$, which has time complexity $\mathcal{O}[M(T_H^3 + T_{\tilde{x}}^3)]$. Thus, the resulting one-time cost has time complexity $\mathcal{O}[M(T_H^3 + KT_H^2 + KT_{\tilde{x}}^2 + T_{\tilde{x}}^3)]$. Afterwards, all $\mathbf{I}^{(Q_2, H_{i,j}, x_j)}$, $\mathbf{I}^{(Q_3, H_{i,j}, H_{i,j}, x_j)}$ and $\mathbf{I}^{(Q_4, H_{i,j}, x_j)}$ are constructed in $\mathcal{O}[M(KT_H^2 + KT_H T_{\tilde{x}} + KT_{\tilde{x}}^2)]$ time and $T^{(1)}$ and $T^{(2)}$ are computed in respectively $\mathcal{O}[M(T_H^2 + T_H T_{\tilde{x}} + T_{\tilde{x}}^2) + M^2]$ and $\mathcal{O}(M[T_H^2 T_{\tilde{x}} + T_H T_{\tilde{x}}^2])$ time. Hence

$$\mathcal{O}(C_\Sigma) = \mathcal{O}[M(T_H^2 T_{\tilde{x}} + KT_H^2 + KT_H T_{\tilde{x}} + KT_{\tilde{x}}^2 + T_H T_{\tilde{x}}^2) + M^2].$$

I.4 Integrals $I^{(\cdot)}$

In this section we compute the integrals $I^{(\cdot)}$ from Appendix I.3.

I.4.1 Kernels

First, let the composite vector be $[\mathbf{t}_1^T \ \mathbf{t}_2^T \ \boldsymbol{\tau}]^T$. Then

$$\begin{aligned}
\mathcal{K}_{\tilde{x}_j}(\mathbf{t}_1, \mathbf{t}_2) &= \int_{\mathbb{R}^{2K}} \exp(-\omega \|\mathbf{t}_1 - \boldsymbol{\tau}_1\|^2) \underbrace{\mathbb{E}[x_j(\boldsymbol{\tau}_1)x_j(\boldsymbol{\tau}_2)]}_{\delta(\boldsymbol{\tau}_1 - \boldsymbol{\tau}_2)} \exp(-\omega \|\mathbf{t}_2 - \boldsymbol{\tau}_2\|^2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\
&= \int_{\mathbb{R}^K} \exp(-\omega \|\mathbf{t}_1 - \boldsymbol{\tau}\|^2 - \omega \|\mathbf{t}_2 - \boldsymbol{\tau}\|^2) d\boldsymbol{\tau} \\
&= I_{\boldsymbol{\tau}} \left[(1, 2 \left[\begin{array}{cc|c} \omega & 0 & -\omega \\ 0 & \omega & -\omega \\ \hline -\omega & -\omega & 2\omega \end{array} \right]) \right] \\
&= \left(\frac{\pi^{K/2}}{(2\omega)^{K/2}}, 2 \left[\begin{array}{cc} \omega & 0 \\ 0 & \omega \end{array} \right] - \frac{1}{\omega} \begin{bmatrix} \omega \\ \omega \end{bmatrix} \begin{bmatrix} \omega \\ \omega \end{bmatrix}^T \right) \\
&= \left(\frac{\pi^{K/2}}{(2\omega)^{K/2}}, \begin{bmatrix} \omega & -\omega \\ -\omega & \omega \end{bmatrix} \right) \\
&= \frac{\pi^{K/2}}{(2\omega)^{K/2}} \exp\left(-\frac{1}{2}\omega \|\mathbf{t}_1 - \mathbf{t}_2\|^2\right) \\
&= \mathcal{K}_{\tilde{x}_j}(\mathbf{t}_1 - \mathbf{t}_2).
\end{aligned}$$

Second, let the composite vector be $[\mathbf{t}_1^T \ \mathbf{t}_2^T]^T$. Then

$$\begin{aligned}
\mathcal{K}_{x_j, \tilde{x}_j}(\mathbf{t}_1, \mathbf{t}_2) &= \int_{\mathbb{R}^K} \underbrace{\mathbb{E}[x_j(\mathbf{t}_1)x_j(\mathbf{t}_2 - \boldsymbol{\tau})]}_{\delta(\mathbf{t}_1 - \mathbf{t}_2 + \boldsymbol{\tau})} \exp(-\omega \|\boldsymbol{\tau}\|^2) d\boldsymbol{\tau} \\
&= \exp(-\omega \|\mathbf{t}_1 - \mathbf{t}_2\|^2) \\
&= (1, 2 \left[\begin{array}{cc} \omega & -\omega \\ -\omega & \omega \end{array} \right]) \\
&= \mathcal{K}_{x_j, \tilde{x}_j}(\mathbf{t}_1 - \mathbf{t}_2)
\end{aligned}$$

and $\mathcal{K}_{\tilde{x}_j, x_j}(\mathbf{t}_1, \mathbf{t}_2) = \mathcal{K}_{x_j, \tilde{x}_j}(\mathbf{t}_1, \mathbf{t}_2)$.

Finally, let the composite vector be $[\mathbf{t}_1^T \quad \mathbf{t}_2^T]^T$. Then

$$\begin{aligned} \mathcal{K}_{H_{i,j}}(\mathbf{t}_1, \mathbf{t}_2) &= \sigma_h^2 \exp(-\alpha \|\mathbf{t}_1\|^2 - \alpha \|\mathbf{t}_2\|^2 - \gamma \|\mathbf{t}_2 - \mathbf{t}_1\|^2) \\ &= (\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma \\ -\gamma & \alpha + \gamma \end{bmatrix}). \end{aligned}$$

I.4.2 Integral $I_{\cdot, \cdot}^{(L, \cdot, \cdot)}$

Let the composite vector be $[\mathbf{t}^T \quad \mathbf{T}_{h,i,:}^T \quad \mathbf{T}_{\tilde{x},j,:}^T \quad \boldsymbol{\tau}^T]^T$ and let $L = \alpha + \gamma + \omega$. Then

$$\begin{aligned} I_{i,j}^{(L,h,x)}(\mathbf{t}) &= \int_{\mathbb{R}^K} \mathcal{K}_h(\mathbf{T}_{h,i,:}, \mathbf{t} - \boldsymbol{\tau}) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, \mathbf{T}_{\tilde{x},j,:}) d\boldsymbol{\tau} \\ &= I_{\boldsymbol{\tau}}[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & -\alpha - \gamma \\ -\gamma & \alpha + \gamma & 0 & \gamma \\ 0 & 0 & 0 & 0 \\ -\alpha - \gamma & \gamma & 0 & \alpha + \gamma \end{bmatrix}) (1, 2 \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \omega & -\omega \\ 0 & 0 & -\omega & \omega \end{bmatrix})], \\ &= I_{\boldsymbol{\tau}}[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & -\alpha - \gamma \\ -\gamma & \alpha + \gamma & 0 & \gamma \\ 0 & 0 & \omega & -\omega \\ \hline -\alpha - \gamma & \gamma & -\omega & L \end{bmatrix})] \\ &= \sigma_h^2 \frac{\pi^{K/2}}{L^{K/2}} \exp(-\alpha \|\mathbf{t}\|^2 - \alpha \|\mathbf{T}_{h,i,:}\|^2 - \gamma \|\mathbf{t} - \mathbf{T}_{h,i,:}\|^2 - \omega \|\mathbf{T}_{\tilde{x},j,:}\|^2 \\ &\quad + L^{-1} \|(\alpha + \gamma)\mathbf{t} - \gamma \mathbf{T}_{h,i,:} + \omega \mathbf{T}_{\tilde{x},j,:}\|^2). \end{aligned}$$

I.4.3 Integral $I^{(Q_1, \cdot, \cdot)}$

Let the composite vector be $[\mathbf{t}_1^T \quad \mathbf{t}_2^T \quad \boldsymbol{\tau}^T]^T$. Then

$$\begin{aligned}
I^{(Q_1, \cdot, \cdot)}(\mathbf{t}_1, \mathbf{t}_2) &= \int_{\mathbb{R}^K} \mathcal{K}_h(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{t}_2 - \boldsymbol{\tau}) \, d\boldsymbol{\tau} \\
&= (\sigma_h^2, 2 \begin{bmatrix} \gamma & -\gamma \\ -\gamma & \gamma \end{bmatrix}) I_{\boldsymbol{\tau}} \left[\left(1, 2 \begin{array}{cc|c} \alpha & 0 & -\alpha \\ 0 & \alpha & -\alpha \\ \hline -\alpha & -\alpha & 2\alpha \end{array} \right) \right] \\
&= (\sigma_h^2, 2 \begin{bmatrix} \gamma & -\gamma \\ -\gamma & \gamma \end{bmatrix}) \left(\frac{\pi^{K/2}}{(2\alpha)^{K/2}}, 2 \begin{bmatrix} \alpha & 0 \\ 0 & \alpha \end{bmatrix} - \frac{1}{\alpha} \begin{bmatrix} -\alpha \\ -\alpha \end{bmatrix} \begin{bmatrix} -\alpha \\ -\alpha \end{bmatrix}^T \right) \\
&= (\sigma_h^2, 2 \begin{bmatrix} \gamma & -\gamma \\ -\gamma & \gamma \end{bmatrix}) \left(\frac{\pi^{K/2}}{(2\alpha)^{K/2}}, \begin{bmatrix} \alpha & -\alpha \\ -\alpha & \alpha \end{bmatrix} \right) \\
&= \left(\sigma_h^2 \frac{\pi^{K/2}}{(2\alpha)^{K/2}}, \begin{bmatrix} \alpha + 2\gamma & -\alpha - 2\gamma \\ -\alpha - 2\gamma & \alpha + 2\gamma \end{bmatrix} \right) \\
&= \sigma_h^2 \frac{\pi^{K/2}}{(2\alpha)^{K/2}} \exp \left[-\frac{1}{2} (\alpha + 2\gamma) \|\mathbf{t}_1 - \mathbf{t}_2\|^2 \right].
\end{aligned}$$

I.4.4 Integral $I_{\cdot, \cdot}^{(Q_2, \cdot, \cdot)}$

Let the composite vector be $= [\mathbf{t}_1^T \quad \mathbf{t}_2^T \quad \mathbf{T}_{\tilde{x},i,:}^T \quad \mathbf{T}_{\tilde{x},j,:}^T \quad \boldsymbol{\tau}_1^T \quad \boldsymbol{\tau}_2^T]^T$. Then

$$\begin{aligned}
& I_{i,j}^{(Q_2, h, x)}(\mathbf{t}_1, \mathbf{t}_2) \\
&= \int_{\mathbb{R}^{2K}} \mathcal{K}_h(\mathbf{t}_1 - \boldsymbol{\tau}_1, \mathbf{t}_2 - \boldsymbol{\tau}_2) \mathcal{K}_{x, \tilde{x}}(\boldsymbol{\tau}_1, \mathbf{T}_{\tilde{x},i,:}) \mathcal{K}_{\tilde{x}, x}(\mathbf{T}_{\tilde{x},j,:}, \boldsymbol{\tau}_2) d\boldsymbol{\tau}_1 d\boldsymbol{\tau}_2 \\
&= I_{\boldsymbol{\tau}_1, \boldsymbol{\tau}_2}[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & 0 & -\alpha - \gamma & \gamma \\ -\gamma & \alpha + \gamma & 0 & 0 & \gamma & -\alpha - \gamma \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ -\alpha - \gamma & \gamma & 0 & 0 & \alpha + \gamma & \gamma \\ \gamma & -\alpha - \gamma & 0 & 0 & \gamma & \alpha + \gamma \end{bmatrix} \\
&\quad (1, 2 \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \omega & 0 & -\omega & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -\omega & 0 & \omega & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \quad (1, 2 \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \omega & 0 & -\omega \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & -\omega & 0 & \omega \end{bmatrix})], \\
&= I_{\boldsymbol{\tau}_1, \boldsymbol{\tau}_2}[(\sigma_h^2, 2 \begin{array}{c|c} \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & 0 \\ -\gamma & \alpha + \gamma & 0 & 0 \\ 0 & 0 & \omega & 0 \\ 0 & 0 & 0 & \omega \end{bmatrix} & \begin{bmatrix} -\alpha - \gamma & \gamma \\ \gamma & -\alpha - \gamma \\ -\omega & 0 \\ 0 & -\omega \end{bmatrix} \\ \hline \begin{bmatrix} -\alpha - \gamma & \gamma & -\omega & 0 \\ \gamma & -\alpha - \gamma & 0 & -\omega \end{bmatrix} & \begin{bmatrix} L & \gamma \\ \gamma & L \end{bmatrix} \end{array}) \\
&= \sigma_h^2 \frac{\pi^K}{(L^2 - \gamma^2)^{K/2}} \\
&\quad \exp \left[-\alpha \|\mathbf{t}_1\|^2 - \alpha \|\mathbf{t}_2\|^2 - \gamma \|\mathbf{t}_1 - \mathbf{t}_2\|^2 - \omega \|\mathbf{T}_{\tilde{x},i,:}\|^2 - \omega \|\mathbf{T}_{\tilde{x},j,:}\|^2 \right. \\
&\quad \left. + (L^2 - \gamma^2)^{-1} \left\{ L \|(\alpha + \gamma)\mathbf{t}_1 - \gamma\mathbf{t}_2 + \omega\mathbf{T}_{\tilde{x},i,:}\|^2 \right. \right. \\
&\quad \left. \left. + L \|-\gamma\mathbf{t}_1 + (\alpha + \gamma)\mathbf{t}_2 + \omega\mathbf{T}_{\tilde{x},j,:}\|^2 \right. \right. \\
&\quad \left. \left. + 2\gamma [(\alpha + \gamma)\mathbf{t}_1 - \gamma\mathbf{t}_2 + \omega\mathbf{T}_{\tilde{x},i,:}]^T [-\gamma\mathbf{t}_1 + (\alpha + \gamma)\mathbf{t}_2 + \omega\mathbf{T}_{\tilde{x},j,:}] \right\} \right].
\end{aligned}$$

I.4.5 Integral $I_{i,j}^{(Q_3, \cdot, \cdot, \cdot)}$

Let the composite vector be $\mathbf{t} = [\mathbf{t}_1^T \quad \mathbf{t}_2^T \quad \mathbf{T}_{h_k, i, \cdot}^T \quad \mathbf{T}_{h_l, j, \cdot}^T \quad \boldsymbol{\tau}^T]^T$. Then

$$\begin{aligned}
& I_{i,j}^{(Q_3, h_k, h_l, x)}(\mathbf{t}_1, \mathbf{t}_2) \\
&= \int_{\mathbb{R}^K} \mathcal{K}_h(\mathbf{t}_1 - \boldsymbol{\tau}, \mathbf{T}_{h_k, i, \cdot}) \mathcal{K}_h(\mathbf{T}_{h_l, j, \cdot}, \mathbf{t}_2 - \boldsymbol{\tau}) \, d\boldsymbol{\tau} \\
&= I_{\boldsymbol{\tau}} \left[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & 0 & -\gamma & 0 & -\alpha - \gamma \\ 0 & 0 & 0 & 0 & 0 \\ -\gamma & 0 & \alpha + \gamma & 0 & \gamma \\ 0 & 0 & 0 & 0 & 0 \\ -\alpha - \gamma & 0 & \gamma & 0 & \alpha + \gamma \end{bmatrix}) (\sigma_h^2, 2 \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & \alpha + \gamma & 0 & -\gamma & -\alpha - \gamma \\ 0 & 0 & 0 & 0 & 0 \\ 0 & -\gamma & 0 & \alpha + \gamma & \gamma \\ 0 & -\alpha - \gamma & 0 & \gamma & \alpha + \gamma \end{bmatrix}) \right] \\
&= I_{\boldsymbol{\tau}} \left[(\sigma_h^4, 2 \left[\begin{array}{cccc|c} \alpha + \gamma & 0 & -\gamma & 0 & -\alpha - \gamma \\ 0 & \alpha + \gamma & 0 & -\gamma & -\alpha - \gamma \\ -\gamma & 0 & \alpha + \gamma & 0 & \gamma \\ 0 & -\gamma & 0 & \alpha + \gamma & \gamma \\ \hline -\alpha - \gamma & -\alpha - \gamma & \gamma & \gamma & 2\alpha + 2\gamma \end{array} \right]) \right] \\
&= \left(\sigma_h^4 \frac{\pi^{K/2}}{(2\alpha + 2\gamma)^{K/2}}, 2 \begin{bmatrix} \alpha + \gamma & 0 & -\gamma & 0 \\ 0 & \alpha + \gamma & 0 & -\gamma \\ -\gamma & 0 & \alpha + \gamma & 0 \\ 0 & -\gamma & 0 & \alpha + \gamma \end{bmatrix} - \frac{1}{\alpha + \gamma} \begin{bmatrix} \alpha + \gamma \\ \alpha + \gamma \\ -\gamma \\ -\gamma \end{bmatrix} \begin{bmatrix} \alpha + \gamma \\ \alpha + \gamma \\ -\gamma \\ -\gamma \end{bmatrix}^T \right).
\end{aligned}$$

Now

$$\begin{aligned}
& 2 \begin{bmatrix} \alpha + \gamma & 0 & -\gamma & 0 \\ 0 & \alpha + \gamma & 0 & -\gamma \\ -\gamma & 0 & \alpha + \gamma & 0 \\ 0 & -\gamma & 0 & \alpha + \gamma \end{bmatrix} - \frac{1}{\alpha + \gamma} \begin{bmatrix} \alpha + \gamma \\ \alpha + \gamma \\ -\gamma \\ -\gamma \end{bmatrix} \begin{bmatrix} \alpha + \gamma \\ \alpha + \gamma \\ -\gamma \\ -\gamma \end{bmatrix}^T \\
&= 2 \begin{bmatrix} \alpha + \gamma & 0 & -\gamma & 0 \\ 0 & \alpha + \gamma & 0 & -\gamma \\ -\gamma & 0 & \alpha + \gamma & 0 \\ 0 & -\gamma & 0 & \alpha + \gamma \end{bmatrix} - \begin{bmatrix} \alpha + \gamma & \alpha + \gamma & -\gamma & -\gamma \\ \alpha + \gamma & \alpha + \gamma & -\gamma & -\gamma \\ -\gamma & -\gamma & \alpha + \gamma & \alpha + \gamma \\ -\gamma & -\gamma & \alpha + \gamma & \alpha + \gamma \end{bmatrix} \\
&\quad - \frac{\gamma^2 - (\alpha + \gamma)^2}{\alpha + \gamma} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \\
&= \begin{bmatrix} \alpha + \gamma & -\alpha - \gamma & -\gamma & \gamma \\ -\alpha - \gamma & \alpha + \gamma & \gamma & -\gamma \\ -\gamma & \gamma & \alpha + \gamma & -\alpha - \gamma \\ \gamma & -\gamma & -\alpha - \gamma & \alpha + \gamma \end{bmatrix} - \frac{\gamma^2 - (\alpha + \gamma)^2}{\alpha + \gamma} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix}
\end{aligned}$$

and so

$$\begin{aligned}
I_{i,j}^{(Q_3, h_k, h_l, x)}(\mathbf{t}_1, \mathbf{t}_2) &= \sigma_h^4 \frac{\pi^{K/2}}{(2\alpha + 2\gamma)^{K/2}} \exp \left[-\frac{1}{2} \alpha \|\mathbf{t}_1 - \mathbf{t}_2\|^2 - \frac{1}{2} \alpha \|\mathbf{T}_{h_k, i, :} - \mathbf{T}_{h_l, j, :}\|^2 \right. \\
&\quad \left. - \frac{1}{2} \gamma \|(\mathbf{t}_1 - \mathbf{t}_2) - (\mathbf{T}_{h_k, i, :} - \mathbf{T}_{h_l, j, :})\|^2 \right. \\
&\quad \left. + \frac{1}{2} \frac{\gamma^2 - (\alpha + \gamma)^2}{\alpha + \gamma} \|\mathbf{T}_{h_k, i, :} + \mathbf{T}_{h_l, j, :}\|^2 \right].
\end{aligned}$$

Furthermore, denote $\mathbf{t}_1 - \mathbf{t}_2 = \mathbf{t}$ and $\mathbf{T}_{h_k, i, :} - \mathbf{T}_{h_l, j, :} = \mathbf{T}$. Then

$$\begin{aligned}
\alpha \|\mathbf{t}\|^2 + \alpha \|\mathbf{T}\|^2 + \gamma \|\mathbf{t} - \mathbf{T}\|^2 &= (\alpha + \gamma) \|\mathbf{t}\|^2 - 2\gamma \mathbf{t}^T \mathbf{T} + (\alpha + \gamma) \|\mathbf{T}\|^2 \\
&= (\alpha + \gamma) \left\| \mathbf{t} - \frac{\gamma}{\alpha + \gamma} \mathbf{T} \right\|^2 + \frac{(\alpha + \gamma)^2 - \gamma^2}{\alpha + \gamma} \|\mathbf{T}\|^2
\end{aligned}$$

so that

$$I_{i,j}^{(Q_3, h_k, h_l, x)}(\mathbf{t}_1, \mathbf{t}_2) = \sigma_h^4 \frac{\pi^{K/2}}{(2\alpha + 2\gamma)^{K/2}} \exp \left[-\frac{1}{2}(\alpha + \gamma) \left\| (\mathbf{t}_1 - \mathbf{t}_2) - \frac{\gamma}{\alpha + \gamma} (\mathbf{T}_{h_k, i, :} - \mathbf{T}_{h_l, j, :}) \right\|^2 + \frac{\gamma^2 - (\alpha + \gamma)^2}{\alpha + \gamma} (\|\mathbf{T}_{h_k, i, :}\|^2 + \|\mathbf{T}_{h_l, j, :}\|^2) \right].$$

I.4.6 Integral $I_{i,j}^{(Q_4, \cdot, \cdot)}$

Let the composite vector be $[\mathbf{t}^T \quad \mathbf{T}_{h,i,:}^T \quad \mathbf{T}_{\tilde{x},j,:}^T \quad \boldsymbol{\tau}^T]^T$. Then

$$\begin{aligned} I_{i,j}^{(Q_4, h, x)}(\mathbf{t}) &= \int_{\mathbb{R}^K} \mathcal{K}_h(\mathbf{t} - \boldsymbol{\tau}, \mathbf{T}_{h,i,:}) \mathcal{K}_{x,\tilde{x}}(\boldsymbol{\tau}, \mathbf{T}_{\tilde{x},j,:}) d\boldsymbol{\tau} \\ &= I_{\boldsymbol{\tau}} \left[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & -\alpha - \gamma \\ -\gamma & \alpha + \gamma & 0 & \gamma \\ 0 & 0 & 0 & 0 \\ -\alpha - \gamma & \gamma & 0 & \alpha + \gamma \end{bmatrix}) (1, 2 \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & \omega & -\omega \\ 0 & 0 & -\omega & \omega \end{bmatrix}) \right] \\ &= I_{\boldsymbol{\tau}} \left[(\sigma_h^2, 2 \begin{bmatrix} \alpha + \gamma & -\gamma & 0 & -\alpha - \gamma \\ -\gamma & \alpha + \gamma & 0 & \gamma \\ 0 & 0 & \omega & -\omega \\ -\alpha - \gamma & \gamma & -\omega & L \end{bmatrix}) \right] \\ &= \sigma_h^2 \frac{\pi^{K/2}}{L^{K/2}} \exp[-\alpha \|\mathbf{t}\|^2 - \alpha \|\mathbf{T}_{h,i,:}\|^2 - \gamma \|\mathbf{t} - \mathbf{T}_{h,i,:}\|^2 - \omega \|\mathbf{T}_{\tilde{x},j,:}\|^2 \\ &\quad + L^{-1} \|(\alpha + \gamma)\mathbf{t} - \gamma \mathbf{T}_{h,i,:} + \omega \mathbf{T}_{\tilde{x},j,:}\|^2]. \end{aligned}$$

I.5 Conclusion

We have derived an analytical expression for the variational free energy of Model 5. The asymptotic time complexity of computing the free energy including one-time costs is given by

$$\begin{aligned} &\mathcal{O}[NMT_H^3 + MT_{\tilde{x}}^3 + YN(C_\mu + C_\Sigma) + \text{one-time costs}] \\ &= \mathcal{O}[NMT_H^3 + MT_{\tilde{x}}^3 + YNM(T_H^2 T_{\tilde{x}} + KT_H^2 + KT_H T_{\tilde{x}} + KT_{\tilde{x}}^2 + T_H T_{\tilde{x}}^2) + YNM^2 \\ &\quad + M(T_H^3 + KT_H^2 + KT_{\tilde{x}}^2 + T_{\tilde{x}}^3)] \\ &= \mathcal{O}[NMT_H^3 + MT_{\tilde{x}}^3 + YNM(T_H^2 T_{\tilde{x}} + KT_H^2 + KT_H T_{\tilde{x}} + KT_{\tilde{x}}^2 + T_H T_{\tilde{x}}^2) + YNM^2]. \end{aligned}$$

References

- Álvarez, M. A. and Lawrence, N. D. (2011). Computationally Efficient Convolved Multiple Output Gaussian Processes. *Journal of Machine Learning Research*, (12):1459–1500. (Pages xxii, 31, 33, and 34)
- Álvarez, M. A., Luengo, D., and Lawrence, N. D. (2009). Latent Force Models. *Artificial Intelligence and Statistics*, 5:9–16. (Pages xxiii, 31, 33, and 34)
- Bonilla, E. V., Chai, K. M., and Williams, C. K. I. (2008). Multi-Task Gaussian Process Prediction. *Advances in Neural Information Processing Systems*, 20:153–160. (Pages xxiii, 33, and 34)
- Duvenaud, D. (2014). *Automatic Model Construction with Gaussian Processes*. PhD thesis, Computational and Biological Learning Laboratory, University of Cambridge. (Page 1)
- Goovaerts, P. (1997). *Geostatistics for Natural Resources Evaluation*. Oxford University Press. (Pages xxii, xxiii, 32, 33, and 34)
- Gray, R. M. (2006). Toeplitz and Circulant Matrices: A Review. *Foundations and Trends in Communications and Information Theory*, 2(3):155–239. (Pages 62 and 66)
- Higham, N. J. (1988). Computing a Nearest Symmetric Positive Semidefinite Matrix. *Linear Algebra and its Applications*, 103:103–118. (Page 57)
- MacKay, D. J. C. (2002). *Information Theory, Inference & Learning Algorithms*. Cambridge University Press. (Pages 10 and 16)
- Minka, T. (2000). Deriving Quadrature Rules from Gaussian Processes. Technical report, Statistics Department, Carnegie Mellon University. (Page 21)

- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press. (Page 51)
- Nguyen, T. V. and Bonilla, E. V. (2014). Collaborative Multi-Output Gaussian Processes. *Conference on Uncertainty in Artificial Intelligence*, 30. (Pages xxii, 33, and 34)
- Rasmussen, C. E. and Williams, C. K. I. (2006). *Gaussian Processes for Machine Learning*. MIT Press. (Pages 1, 9, and 10)
- Teh, Y. W. and Seeger, M. (2005). Semiparametric Latent Factor Models. *International Workshop on Artificial Intelligence and Statistics*, 10. (Pages xxiii, 33, and 34)
- Titsias, M. K. (2009). Variational Learning of Inducing Variables in Sparse Gaussian Processes. *Artificial Intelligence and Statistics*, 12:567–574. (Pages 23, 26, and 27)
- Tobar, F., Bui, T. D., and Turner, R. E. (2015a). Design of Covariance Functions Using Inter-Domain Inducing Variables. *Time Series Workshop on Advances in Neural Information Processing Systems*. (Page 27)
- Tobar, F., Bui, T. D., and Turner, R. E. (2015b). Learning Stationary Time Series using Gaussian Processes with Nonparametric Kernels. *Advances in Neural Information Processing Systems*, 29:3501–3509. (Pages vii, xvii, xxii, 3, 15, 21, 25, 26, and 37)
- Tobar, F. and Turner, R. E. (2016). Modelling Time Series via Automatic Learning of Basis Functions. *Ninth IEEE Sensor Array and Multichannel Signal Processing Workshop*. (Pages 29 and 30)
- Trefethen, L. N. and Bau, D. (1997). *Numerical Linear Algebra*. Society for Industrial and Applied Mathematics. (Pages 55 and 56)
- Ulrich, K., Carlson, D. E., Dzirasa, K., and Carin, L. (2015). GP Kernels for Cross-Spectrum Analysis. *Advances in Neural Information Processing Systems*, 28:1999–2007. (Pages vii, xxii, 3, 25, 26, 33, 34, and 61)
- Wilson, A. G. and Adams, R. P. (2013). Gaussian Process Kernels for Pattern Discovery and Extrapolation. *International Conference on Machine Learning*, 3:1067–1075. (Pages xxiii, 1, 25, and 26)

Wilson, A. G., Knowles, D. A., and Ghahramani, Z. (2012). Gaussian process Regression Networks. *International Conference on Machine Learning*, 29.

(Pages xxii, 31, and 33)

Index

- Approximate Deep Kernel Model, 43
- approximate inference, 27
- Approximate Kernel Model, 22
- automatic relevance determination, 16
- Basis Function Model, 29
- Bayesian numerical approximation, 21
- Cholesky decomposition, 56
- circulant approximation, 62
- Convolutional Mixing Model, 32
- cross-spectral mixture kernel, 25
- Deep Gaussian Process Convolution Model, 39
- Deep Kernel Model, 40
- evidence, 10
- free energy, 28
- Gaussian process, 9
- Gaussian Process Convolution Model, 25, 38
 - explicit decay, 38
- Gaussian process regression, 10
- Generalised Gaussian Process Convolution Model, 11
- Instantaneous Mixing Model, 32
- inter-domain transformation, 27
- kernel, 1
 - decaying exponentiated-quadratic kernel, 15
 - diagonal multi-output decaying automatic relevance determination kernel, 16
 - diagonal multi-output decaying exponentiated-quadratic kernel, 15
 - diagonal multi-output exponentiated-quadratic kernel, 15
 - exponentiated-quadratic kernel, 15
 - root, 19
- kernel design problem, 1
- Kullback-Leibler divergence, 51
- linear state-space model, 5
- marginal likelihood, 10
- mixing model hierarchy, 33

multi-task learning, 31

Multidimensional Time-Invariant
Model, 9

nearest symmetric positive-semidefinite
matrix, 57

neural network, 41

Nonparametric Kernel Model, 12

nonparametric model, 1

Occam's razor, 10

spectral mixture kernel, 25

Time-Invariant Model, 6

Time-Variant Model, 6

variational free energy, 28