



Universiteit
Leiden
The Netherlands

Bayesian inference for Gaussian models: Inverse problems and evolution equations

Yan, D.

Citation

Yan, D. (2020, March 3). *Bayesian inference for Gaussian models: Inverse problems and evolution equations*. Retrieved from <https://hdl.handle.net/1887/86070>

Version: Publisher's Version

License: [Licence agreement concerning inclusion of doctoral thesis in the Institutional Repository of the University of Leiden](#)

Downloaded from: <https://hdl.handle.net/1887/86070>

Note: To cite this publication please use the final published version (if applicable).

Cover Page



Universiteit Leiden



The handle <http://hdl.handle.net/1887/86070> holds various files of this Leiden University dissertation.

Author: Yan, D.

Title: Bayesian inference for Gaussian models: Inverse problems and evolution equations

Issue Date: 2020-03-03

BAYESIAN INFERENCE
FOR
GAUSSIAN MODELS

INVERSE PROBLEMS AND EVOLUTION EQUATIONS

DONG YAN

Bayesian Inference for Gaussian Models

Inverse Problems and Evolution Equations

Dong Yan

ISBN 978-94-028-1963-2

This publication is typeset in L^AT_EX using the Memoir class

Printed by Ipskamp Printing B.V., Enschede

Cover by Dong Yan

Copyright © 2020 by Dong Yan

All rights reserved

The research of this doctoral thesis received financial assistance from the European Research Council (ERC; 320637)

Bayesian Inference for Gaussian Models

PROEFSCHRIFT

TER VERKRIJGING VAN

DE GRAAD VAN DOCTOR AAN DE UNIVERSITEIT LEIDEN,

OP GEZAG VAN RECTOR MAGNIFICUS PROF. MR. C. J. J. M. STOLKER,

VOLGENS BESLUIT VAN HET COLLEGE VOOR PROMOTIES

TE VERDEDIGEN OP DINSDAG 3 MAART 2020

KLOKKE 13:45 UUR

DOOR

DONG YAN

GEBOREN TE GUANGDONG, CHINA

IN 1989

Samenstelling van de promotiecommissie:

Promotor:	Prof. dr. A.W. van der Vaart	(Universiteit Leiden)
Copromotor:	Dr. S. Gugushvili	(Wageningen University & Research)
Overige leden:	Prof. dr. P. Steenhagen	(Universiteit Leiden, voorzitter)
	Prof. dr. J.J. Meulman	(Universiteit Leiden, secretaris)
	Prof. dr. M. Reiss	(Humboldt University of Berlin)
	Prof. dr. A.J. Schmidt-Hieber	(Universiteit Twente)
	Prof. dr. J.H. van Zanten	(Vrije Universiteit Amsterdam)
	Dr. K. Ray	(Imperial College London)

*To my grandparents, who enlightened me
To my mother, for everything*

Contents

1	Introduction	1
1.1	Infinite-Dimensional Parameter Spaces	2
1.2	Gaussian Linear Models	4
1.2.1	Forward Mapping \mathcal{A}	5
1.2.2	Gaussian noise	6
1.2.3	Concrete Models	7
1.3	Bayesian Methodology	8
1.4	Overview	8
1.5	Notations	9
1.6	Notes	10
I	Foundations	13
2	Smoothness Class	17
2.1	Smoothness scales	17
2.2	Hilbert Scale	20
2.2.1	Relation to Boundary Conditions	25
2.3	Smoothness in Higher Dimensions	26
2.3.1	Multi-dimensional Smoothness	27
2.4	Approximation Number and Metric Entropy	31
2.5	Notes	33
3	Gaussian Analysis	35
3.1	Probability Measures on Banach Spaces	35
3.2	Gaussian Measures	38
3.2.1	Covariance Structure	39
3.2.2	Examples	43
3.3	Radonification of Cylindrical Measures	46
3.4	Notes	48
4	Bayesian Nonparametrics for Gaussian Linear Models	49
4.1	Bayesian Nonparametrics	49
4.1.1	Bayes' Rule	49
4.1.2	Bayesian Asymptotics	50
4.2	Gaussian Linear Models	52
4.2.1	Continuous Observation	52
4.2.2	Discrete Observation	55
4.3	Bayesian Contraction for Gaussian Linear Models	56

II Inverse Problems	61
5 Linear Inverse Problems	65
5.1 Introduction	65
5.2 Inverse Nature	66
5.3 Galerkin Projection	70
5.4 Notes	72
6 Inverse Problems with Continuous Observations in Smoothness Scales	75
6.1 Introduction	75
6.2 General Result	75
6.3 Random Series Priors	80
6.4 Gaussian Priors	82
6.5 Gaussian Mixtures	85
6.6 Proofs	85
6.6.1 Proof of Theorem 6.5	85
6.6.2 Proof of Theorem 6.7	87
6.6.3 Proof of Theorem 6.11	90
6.7 Discussion and Comments	92
7 Inverse Problems with Discrete Observations: Gaussian Conjugacy	95
7.1 Introduction	95
7.1.1 Notation	96
7.2 Sequence model	97
7.2.1 Singular value decomposition	97
7.2.2 Equivalent formulation	98
7.3 Main results	101
7.3.1 Contraction rates	101
7.3.2 Credible sets	103
7.4 Simulation examples	104
7.5 Proofs	106
7.5.1 Proof of Lemma 7.9	106
7.5.2 Proof of Theorem 7.12	108
7.5.3 Proof of Theorem 7.14	110
7.5.4 Proof of Theorem 7.15	110
7.5.5 Proof of Theorem 7.16	112
7.6 Auxiliary lemmas	113
8 Inverse Problems with Discrete Observations in Smoothness Scales	115
8.1 Signal Reconstruction	116
8.2 General Contraction Rates	120
8.3 Random Series Priors	124
8.4 Gaussian Priors	125
8.5 Gaussian Mixtures	126
8.6 Proofs	127
8.6.1 Proof of Theorem 8.10	127
8.6.2 Proof of Theorem 8.11	129
8.6.3 Proof of Theorem 8.13	133

III Evolution Equations	137
9 Linear Evolution Equations	141
9.1 \mathcal{Q} -Wiener Processes	142
9.2 Stochastic Integrals in Hilbert Spaces	143
9.3 Extension of Stochastic Integrals	146
9.3.1 Cylindrical Wiener Process	146
9.3.2 Stochastic Integral with Cylindrical Wiener Process	147
9.4 Deterministic Evolution Equations	147
9.5 Solutions of SPDEs	148
9.6 Notes	152
10 Bayesian Inference for Linear Evolution Equations	153
10.1 Recovery of the Initial condition	154
10.1.1 Spatial Gaussian Priors	154
10.1.2 Contraction Rate for Initial Condition	155
10.2 Recovery of the Drift	156
10.2.1 Spatial-Temporal Gaussian Priors	156
10.2.2 Contraction Rate for Drift Recovery	157
10.3 Proofs	158
10.3.1 Proofs in Section 10.1	158
10.3.2 Proofs in Section 10.2	162
10.3.2.1 Whitening Ornstein-Uhlenbeck processes	163
10.3.2.2 Complete Sequence Model	165
10.3.2.3 Gaussian Posterior Contraction for Multi-dimensional White Noise Model	165
10.3.2.4 Proof of Theorem 10.5	167
10.4 Entropy Number with Non-Polynomial Rates	167
 Appendix	 169
A Mathematical Tools	171
A.1 Miscellaneous Lemmas	172
A.2 Pseudo-Inverse	174
A.3 Compact Operators	175
 References	 179
 Index	 187
 Summary	 189
 Samenvatting	 193
 Acknowledgement	 197
 Curriculum Vitae	 199

Chapter 1

Introduction

The objective of statistical inference is to infer a quantity from observations. A ‘useful’ inference procedure should disclose a certain truth. Intuitively, it sounds reasonable to wish that, given the existence of a (fixed) underlying truth, the inference procedure produces approximations closer to the truth, when more information is available. Following this intuition, a rigorous theory has been well developed in the twentieth century, which is now known as the field of *asymptotic theory* (alternatively, *large sample theory*). To reflect realistic observational processes, it is customary to number observations by the natural number $\mathbb{N} = \{1, 2, 3, \dots\}$. In the asymptotic framework, the performance of inferences is studied using various criteria measuring the ‘accuracy’ of an estimation, when the sample size n tends to infinity. The behaviour of an inferential procedure by taking the limit as $n \rightarrow \infty$ is called *asymptotics*.

The idea in the previous paragraph can be formulated in mathematical language as follows. Given a measurable space $(\mathbb{X}, \mathcal{X})$, i.e. the *sample space*, an *experiment* \mathcal{E} is a set $\{\mathbb{P}_\theta : \theta \in \Theta\}$ of probability distributions indexed by a parameter family Θ . In practice, the indexation is guided by a model that generates a probability distribution $\{\mathbb{P}_\theta\}$ for each $\theta \in \Theta$. In this situation, the experiment is also called a *statistical model*. For each $n \in \mathbb{N}$, the n th observation $X^{(n)}$ is a random element whose distribution $\mathbb{P}_\theta^{(n)}$ is from an experiment $\mathcal{E}_n = \{\mathbb{P}_\theta^{(n)}\}_{\theta \in \Theta}$. The term *sample* and *observation* are used interchangeably for $\{X^{(n)}\}$. Notice that the experiments are not necessarily identical (as they depend on n), but they are indexed with the same parameter space, and additionally, each sample $X^{(n)}$ encodes some information on the same parameter θ . Statistical inference is to propose an estimate¹ $\hat{\theta}$, which only depends on the observations, that is a sensible approximation to the true parameter θ .

In the general asymptotic framework, the following components are decisive to the formulation of asymptotics: the parameter space Θ , the experiment \mathcal{E}_n , and the methodology guiding the inferential process. In this thesis, we will study the *asymptotics of Bayesian nonparametric inference for Gaussian linear models*

¹In literature, statistical inference is often categorised into estimation and hypothesis testing. However, we use inference and estimation interchangeably, as hypothesis testing is not touched in this thesis.

(GLMs). Admittedly confusing, the term ‘Nonparametric’ actually refers to the fact that the dimension of the parameter space Θ is infinite. In order to handle statistical analysis in infinite-dimensional space, a set of mathematical tools, which differs from the tools used in parametric statistics, is desired. We will provide a short survey in Chapter 2. The statistical model considered is the Gaussian linear model (GLM). To be precise, we will study several models, varying in the level of abstraction and generality, all of which fall into the class of GLMs. The phrase ‘Bayesian’ refers to the methodology that is based on Bayes’ rule. In Chapter 3, the general framework of GLMs will be outlined and a Bayesian asymptotic result coping with GLM will also be given.

The rest of the chapter is arranged as follows. First we sketch the characteristics of parameter spaces used in nonparametric statistics in Section 1.1. Then, the statistical model considered in this thesis is introduced in Section 1.2. The Bayesian approach to statistical inference is briefly reviewed in Section 1.3. After specifying all the necessary components to formulate a feasible asymptotic study, we provide in Section 1.4 an overview of the main topics covered in this thesis. We conclude this chapter with a summary of notation in Section 1.5, and supplementary notes with references in Section 1.6.

1.1 Infinite-Dimensional Parameter Spaces

To accommodate complex stochastic phenomena, many advanced probabilistic models have been naturally developed in infinite-dimensional spaces. Meanwhile, due to the growing capacity of computation and the increasing volume of data storage, it is also practical to consider inference for models of infinite-dimensional nature, which leads to the study of nonparametric estimation. The parameter spaces in parametric inference are essentially finite-dimensional vector spaces, i.e. Euclidean spaces. Due to its plain structure, it does not raise any concerns that the chosen space may give rise to complications in the estimation. On the contrary, in nonparametric statistics, the basic assumption is that the true parameter belongs to an infinite-dimensional space Θ , whose structure is more involved. For example, Lebesgue measure does not extend to infinite-dimensional spaces; closed and bounded sets are not necessarily compact, etc. Therefore, infinite-dimensional spaces require more careful examination. In this section, we will introduce the fundamental concept of smoothness of Hölder and Sobolev types. Other related notions separability, approximation property and compactness are also important criteria of nonparametric parameter spaces. A detailed discussion with the emphasis on Hilbert spaces can be found in Chapter 2.

In nonparametric statistics, the infinite-dimensional parameter space Θ is often postulated to be a complete space, i.e. a Banach space. Sometimes, it is also hypothesised that the space Θ possesses an inner product structure, and hence Θ is a Hilbert space. The Hilbertian assumption is often driven by underlying (physical) models, such as the postulates of quantum mechanics. We consider the following two canonical examples of parameter spaces: on $\mathfrak{X} = [0, 1] \subset \mathbb{R}$, the space $C[0, 1]$ of continuous functions, and the space $L^2[0, 1]$ of square integrable functions, which has an inner product structure.

With the index n indicating the growth of information, for a sequence $\{\widehat{\theta}_n\}$ of estimates for the parameter θ , it is desirable that $d(\widehat{\theta}_n, \theta)$ converges to zero in probability as n tends to infinity. If the convergence is true for all possible parameters, then the statistical procedure is (asymptotically) *consistent*. Without further assumptions, the convergence rate may be arbitrarily slow, which is not so useful in practice.

To obtain a reasonable convergence rate, e.g. the polynomial rate $n^{-\beta}$ with $\beta \in \mathbb{R}^+$, the logarithmic rate $\log n$, etc., a *smoothness* (or *regularity*) condition is required. Briefly, a set $\{\Theta_s : s \in \mathcal{S}\}$ of subspaces of a parameter space Θ is called a *smoothness class*, if a reasonable convergence rate can be obtained for certain statistical models, when the parameter is known to belong to Θ_s . The index s in Θ_s symbolises the ‘smoothness’. To illustrate this concept, we use the two previously mentioned examples $C[0, 1]$ and $L^2[0, 1]$.

Hölder Smoothness

Recall that $C[0, 1]$ is the space of continuous functions. For $k \in \mathbb{N}$, let $C^k[0, 1]$ be the space of k -times differentiable (in the classical sense) continuous functions and $C^0[0, 1] = C[0, 1]$. For a non-integer s , define

$$C^s[0, 1] := \left\{ f \in C^{\lfloor s \rfloor} : \|f\|_s := \sup_{\substack{x, y \in [0, 1] \\ x \neq y}} \frac{|f^{\lfloor s \rfloor}(x) - f^{\lfloor s \rfloor}(y)|}{|x - y|^{s - \lfloor s \rfloor}} < \infty \right\},$$

where $\lfloor s \rfloor$ is the largest integer strictly smaller than s and $f^{\lfloor s \rfloor}$ is the $\lfloor s \rfloor$ -th derivative. The spaces $C^s, s \in \mathbb{R}_0^+$ are *Hölder spaces*, which is the canonical smoothness class for continuous functions.

Continuity plays an important role in the study of inference. The standard estimator for a Hölder class, the *kernel estimator*, is constructed by utilising the continuity property, and consequently the convergence rate also depends on the Hölder smoothness.

Sobolev Smoothness

Sobolev smoothness is directly related to the concept of (weak) *differentiability*. Let \mathfrak{D} be a bounded domain in Euclidean space \mathbb{R}^d . With $k \in \mathbb{N}$, the Sobolev spaces $H_k(\mathfrak{D})$, containing the $L^2(\mathfrak{D})$ functions whose L^2 -weak derivatives exist up to the order k , were first introduced to study partial differential equations (PDEs). Soon it evolved into an important (sub)field in the theory of function spaces, while maintaining an intimate connection to PDEs. One discovery of great importance is that the integer indexed Sobolev spaces, in many situations, e.g. certain boundary conditions being satisfied, can be identified with the domains of the integer powers of a differential operator that is densely defined, strictly positive, and self-adjoint. Subsequently, *fractional* Sobolev spaces can be defined using spectral theory.

With the standard result that the Fourier basis is the eigenbasis of the second order differential operator on $L^2[0, 1]$ with the periodic boundary condition, a spectral argument leads to the well known Statistician’s Sobolev spaces as follows. The *projection* mapping from a function f in $L^2[0, 1]$ to its Fourier coefficients

$\{f_j\}_{j \in \mathbb{N}}$ is an isometric isomorphism from L^2 to the space ℓ^2 of square summable sequences. In addition, with the singular values $\{j : j \in \mathbb{N}\}$ of the second order differential operator, for $s \in \mathbb{R}_0^+$, the Sobolev norms admits a form of weighted ℓ^2 norms, i.e.

$$\|f\|_s^2 = \sum_{j \in \mathbb{N}} j^{2s} f_j^2,$$

in particular, $\|f\|_0 = \|f\|_{\ell^2}$.

In the example $L^2[0, 1]$, the Sobolev smoothness is directly linked to an *unbounded operator* (the second differential operator), which induces a scale of Hilbert spaces with duality relations (see related sections in Chapter 2). Due to the Hilbert space structure, the typical estimators constructed for L^2 are commonly based on projections. Similar to Hölder class, the achieved convergence rate depends on the Sobolev smoothness.

In Hölder smoothness and the example of Sobolev smoothness, the index set is $\mathfrak{T} = [0, 1]$ for simplicity. As seen in the general description on Sobolev smoothness, it is often possible to be replaced by a more general setting, e.g. *compact metric spaces*, and in particular, closed and bounded domains in Euclidean spaces. When the domain $\mathfrak{T} \subset \mathbb{R}^d$ but not bounded, e.g. $\mathfrak{T} = [0, \infty)^d$, the spaces defined on \mathfrak{T} become too ‘big’ to recover the parameter. This difficulty is usually removed by introducing additional properties to the parameter, for example, periodicity or tail conditions.

Although having been mentioned separately in the related paragraphs, we want to stress a fundamental difference between Hölder and Sobolev smoothness. In the Hölder case, the smoothness is characterised by the *local properties of paths* of a function, and the index set may directly influence the properties of functions. In statistics, Hölder smoothness is usually chosen for the recovery of stochastic processes whose index sets have straightforward interpretations, such as time, spatial domains, etc. On the other hand, Sobolev smoothness directly relates to the *duality structure*, as in the previous example via weak differentiability, or in general via an (unbounded) operator. Since probability measures on infinite-dimensional spaces are usually (if not always) characterised via the topological duals, Sobolev smoothness serves as a natural stage for studying statistical problems in infinite-dimensional spaces. In particular, if the underlying space is a Hilbert space, the full scale of the induced dual spaces admits explicit representations, which significantly remedies the technicality involved in the development of the theory. A well known case is the L^2 theory in PDEs. In this thesis, we exclusively focus on the estimation in the smoothness of *Hilbertian Sobolev* type, i.e. Sobolev spaces that are also Hilbert spaces.

1.2 Gaussian Linear Models

Arguably, the Gaussian linear model has occupied a central role since the birth of statistics. Its importance cannot be exaggerated. In this section, we heuristically introduce a general (nonparametric) Gaussian linear model in Hilbert space setting. This model provides a common ground for the (more concrete) models

studied in the subsequent chapters. We will systematically examine the model in Chapter 3.

Formally, the model can be represented as a very simple additive model,

$$X^{(n)} = \mathcal{A}^{(n)}\theta + \xi^{(n)}. \quad (1.1)$$

The parameter of interest is the *signal* θ . In practice, however, it is often the case that only a transform of the original signal is observable, because of e.g. the experimental set-up, etc. A bounded operator $\mathcal{A}^{(n)}$, from the parameter space Θ to another Hilbert space \mathbb{X} , is used to characterise the *transform* (alternatively, *forward mapping*). The randomness arising in the observational process is modelled by a stochastic noise $\xi^{(n)}$ with a structure of Gaussian type.

The model consists of three components: the signal θ , the *forward mappings* $\mathcal{A}^{(n)}$, and the noise $\xi^{(n)}$. Throughout this thesis, the parameter θ is understood as an infinitely-dimensional object, while the features of $\mathcal{A}^{(n)}$ and ξ depend on observation schemes.

Continuous observation

Continuous observation is often an idealization of experiments. The observation $X^{(n)}$ is ‘complete’: for example, the entire trajectory of a process, all functionals on the parameter space, etc. Although in most cases the observation scheme is unrealistic, the model is fruitful for gaining insight. In this situation, the forward mapping is set to $\mathcal{A}^{(n)} = \mathcal{A}$, see Section 1.2.1 for more information. The noise structure is more involved, see the subsequent section Section 1.2.2.

Discrete observation

Discrete observation reflects a more realistic situation. We consider the concrete situation that the image $\mathcal{A}(\Theta)$ is contained in a function space defined on a domain \mathfrak{D} , and noisy samples of the unknown function $\mathcal{A}f$ at a finite number of locations, called *design points*, in its domain are recorded. The observation $Y^{(n)}$ is a random vector $(Y_i)_{i \leq n}$ in \mathbb{R}^n , with the coordinates given by

$$Y_i = \mathcal{A}f(x_i) + z_i, \quad i = 1, \dots, n,$$

where $(z_i)_{i \leq n}$ is often assumed to be standard Gaussian in \mathbb{R}^n . As a consequence, the forward mappings are $\mathcal{A}^{(n)} = \mathcal{E}^{(n)}\mathcal{A}$, where $\mathcal{E}^{(n)}$ is the evaluation operator at (deterministic) design points. Alternatively, discrete observation may also refer to noisy finite-dimensional projections of the signal $\mathcal{A}f$. We do not pursue this direction. In fact, under mild conditions, the discrete observation on design points can be translated to finite-dimensional projections, see the relevant chapter Chapter 8.

We are going to describe the forward mapping and the noise below. After that, we also briefly mention the models that will be investigated in later chapters.

1.2.1 Forward Mapping \mathcal{A}

Since our goal is to recover the parameter θ , a procedure to reconstruct θ from the image $\mathcal{A}\theta$ is desired, i.e. to ‘invert’ the forward operator $\mathcal{A}^{(n)}$. The methods

working with noiseless observations (i.e. $\xi^{(n)} = 0$ almost surely) do not necessarily remain valid for noisy observations. In fact, most of them break down and modifications are needed. The recovery of θ depends on the following components: the transform $\mathcal{A}^{(n)}$ and the structure of the noise $\xi^{(n)}$. Heuristically, the idea for the recovery is described as follows. For clarity, the superscript (n) is omitted. The structure of ξ will determine whether the observations realise in the space \mathbb{X} or are actually certain ‘generalised’ processes over \mathbb{X} . An approximate ‘inverse’² \mathcal{A}^\dagger of \mathcal{A} would help to recover θ . Assuming that \mathcal{A}^\dagger is available, formally applying \mathcal{A}^\dagger to X , we obtain the sum of $\mathcal{A}^\dagger \mathcal{A} \theta$ and $\mathcal{A}^\dagger \xi$. \mathcal{A}^\dagger would be acceptable if $\mathcal{A}^\dagger \mathcal{A} \theta$ is a reasonable approximation of θ , while the spread of $\mathcal{A}^\dagger \xi$ is under control. It is worth noticing that \mathcal{A}^\dagger interplays with ξ . Precisely, if \mathcal{A}^\dagger is linear, the $\mathcal{A}^\dagger \xi$ has the covariance³ $\mathcal{A}^\dagger \Sigma (\mathcal{A}^\dagger)^*$. A desirable ‘inverse’ \mathcal{A}^\dagger should not amplify the noise too much. From the heuristics above, we conclude that the combination of the transform and the noise to a large extent determines the inferential procedure.

1.2.2 Gaussian noise

As already mentioned at the beginning of this section, the noise ξ is directly related to the observation scheme adopted. For discrete observation, the standard Gaussian on Euclidean space is well-known and it does not need explanation. On the other hand, for continuous observation, the noise deserves a closer look.

Consider the case that ξ is a Gaussian random element that lives ‘around’⁴ \mathbb{X} . There are two interpretations to characterise a Gaussian element. When the underlying space \mathbb{X} contains the functions on a domain \mathfrak{D} , the noise ξ can be treated as a *Gaussian process* on the same domain. Alternatively, ξ can also be viewed as a random element whose law is given by a *Gaussian measure* on \mathbb{X} . In many situations, they are equivalent and can be translated from one to the other. Since a large class of the functionals of Wiener process are continuous, it is customary to consider the process version when Hölder smoothness is considered. On the other hand, due to the intrinsic Hilbert structure (i.e. *reproducing kernel Hilbert space*, alias RKHS), Gaussian measure version is often adopted for Sobolev smoothness. More details regarding the features of Gaussian elements are given in Chapter 3.

A Gaussian element ξ in a Hilbert space \mathbb{X} is fully characterised by a mean vector and a covariance operator Σ on \mathbb{X} . Similar to real-valued Gaussian variables, the mean and the covariance operator specify, respectively, the location and the ‘spread’ of the distribution. The noise ξ will always be zero mean in this study. If Σ is an operator of trace class, the noise ξ is a *proper* random element in \mathbb{X} , i.e. $\xi \in \mathbb{X}$ almost surely. Otherwise, ξ is a *generalised* random element *over* \mathbb{X} , that means, ξ takes values in an extension of \mathbb{X} . One noteworthy example of a generalised random element is *white noise*, i.e. $\Sigma = \text{id}$, where id is the identity operator on an infinite-dimensional Hilbert space \mathbb{X} .

²Precisely, it is the *psudoinverse*. This concept is only used in this section for heuristics, and hence we refer to [29] for the detailed treatment.

³ \mathcal{A}^* is the adjoint of \mathcal{A} , see Definition A.1.

⁴As we do not specify if ξ takes values in \mathbb{X} almost surely.

1.2.3 Concrete Models

We will examine two types of Gaussian linear models in this thesis. Although the two types are not completely disjoint from each other and have some overlap, each of them has its own flavour and interest.

Inverse Problems

We only consider the linear case, i.e. \mathcal{A} being a linear operator. When \mathcal{A} has no bounded inverse, which frequently occurs when \mathcal{A} has a ‘smoothing’ property, the (unbounded) \mathcal{A}^\dagger amplifies the noise ξ . The situation is *ill-posed*, as small perturbations (modelled by stochastic noise) in the observation result in a poor reconstruction, dominated by the noise $\mathcal{A}^\dagger \xi$ magnified by the unbounded \mathcal{A}^\dagger . As a consequence, *regularization* techniques need to be developed to handle the ill-posedness of \mathcal{A}^\dagger . The models satisfying the description above are customarily considered *inverse problems*, and the ill-posedness is also known as *inverse nature*. To focus on examining the ramifications of ill-posedness, the noise structure is selected to be a simple form, white noise, which has been widely accepted as a reasonable choice for nonparametric models, see e.g. [36].

The inverse problem is the main theme of Part II.

Evolution Equations

Stochastic partial differential equations (SPDEs) are broadly applied to model stochastic (differential) dynamical systems, time dependent processes whose dynamics possess differential structures. Due to the involvement of time, they are also known as *evolution equations*. A typical form of linear evolution equations is

$$\begin{cases} dX(t) = [\mathcal{L}X(t) + f(t)] dt + B dW(t), \\ X(0) = u \end{cases} \quad (1.2)$$

where \mathcal{L} is a differential operator, f is a *source term* representing the *drift*, B is a linear operator, and dW is the differential of a Wiener process. We highlight the following facts. The state space of $X(t)$ is an infinite-dimensional space; the process W is as well an infinite-dimensional object, i.e. a vector-valued process; the drift f is a deterministic driving force of the dynamics; and $B dW$ shapes a stochastic driving force, which is a well-defined object shown with a vector-valued stochastic integration theory.

The solution of (1.2) is given by a stochastic version of the variation of parameters formula,

$$X(t) = S(t)u + \int_0^t S(t-s)f(s) ds + \int_0^t S(t-s)B dW(s), \quad (1.3)$$

where $S(t)$ is the semigroup generated by the differential operator \mathcal{L} .

It is noteworthy that in (1.3), the stochastic noise is still Gaussian (from the theorem that integration of a deterministic process with respect to a Wiener process is a Gaussian process), but no longer ‘white’. The parameters of interest

in the model (1.2) are the initial condition u and the drift f . The structures of these two parameters have different levels of similarity to the noise pattern. Hence, depending on the parameter to estimate, different approaches are more preferable.

The inference for evolution equations is the main subject of Part III.

1.3 Bayesian Methodology

In short, Bayesian methods utilize Bayes' rule to achieve the goal of inference. The term *Bayesian nonparametrics* refers to the Bayesian methods for infinite-dimensional models. In the Bayesian framework, a *prior distribution* (a probability measure) Π is assigned to the parameter θ . As a consequence, the prior induces a probability measure on the statistical model $\{\mathbb{P}_\theta : \theta \in \Theta\}$, where \mathbb{P}_θ is the law of X given θ . The *posterior distribution* is the conditional distribution of θ given X . The Bayesian procedure is identical for both parametric and nonparametric cases. However, the subtlety of Bayesian nonparametrics is higher, because of measurability concerns, and it will be considered in a subsequent chapter.

Our perspective on the Bayesian framework is that it offers a universal approach for nonparametric inference, while there are several issues necessary to address in order to obtain reasonable asymptotic results. Since we treat the Bayesian framework as a methodological device, and we are interested in the asymptotics of Posterior distributions, our standing is still in the frequentist regime. In other words, we want to understand the asymptotic performance of Bayesian methods from a frequentist perspective. In this thesis we solely focus on the goal just mentioned, and we have no intention to go further towards the long lasting disputation about the two regimes of frequentist and Bayesian at the philosophical level.

A more detailed review on Bayesian nonparametrics is given in Chapter 4.

1.4 Overview

This thesis is organized as follows.

In Part I, we prepare the basic elements that will serve as the building blocks for the later study on Bayesian inference for Gaussian linear models. Chapter 2 deals with the smoothness classes that will be used as the parameter spaces for the statistical study in the later chapters. Chapter 3 treats Gaussian measures on Banach spaces, which is the noise structure in all the models considered in this thesis. Chapter 4 presents the Bayesian nonparametric framework. We explore the noise structure of Gaussian linear models, with continuous and discrete observations, using the results from Chapter 3. In particular, we develop a posterior contraction theorem suitable for Gaussian linear models.

In Part II, we study Bayesian inference for linear inverse problems with Gaussian noise. We consider two types of problems, distinguished by their observation schemes: continuous and discrete observations, corresponding to the white noise model and regression model with transformed signals. In Chapter 5, we formally formulate the smoothing property of the forward operator \mathcal{A} , or equivalently, the *ill-posedness* of its inverse \mathcal{A}^\dagger , in the framework introduced in Chapter 2. In

Chapter 6, we systematically investigate the inverse problem with continuous observations. For the inverse problem with discrete observations, it is studied with two approaches. First, in Chapter 7, the inverse problem is studied in a concrete setting, by leveraging the Gaussian conjugacy in linear models. Then, in Chapter 8, we generalise the methodology developed in Chapter 6 to study the regression model.

In Part III, we study the inference for evolution equations. In Chapter 9, we present the semigroup approach to SPDEs, and as well introduce the additional structure suitable for statistical study. Subsequently, in Chapter 10 the Bayesian approach for the recovery of the parameters of evolution equations is examined.

1.5 Notations

We outline our conventions on the notations used in this work.

- The sets of Natural numbers, real numbers and complex numbers are $\mathbb{N} = \{1, 2, 3, \dots\}$, \mathbb{R} and \mathbb{C} , respectively. The *imaginary* unit is denoted by i . Special subsets are $\mathbb{N}_0 = \{0\} \cup \mathbb{N}$, $\mathbb{R}^+ = (0, \infty)$ and $\mathbb{R}_0^+ = [0, \infty)$. Other similar notations are defined accordingly.
- \mathbb{R}^d vectors and \mathbb{N}^d multi-indices are denoted by $k = (k_1, \dots, k_d)$. When a multi-index has identical entries, the following convention is adopted: $\beta = (\beta, \dots, \beta) \in \mathbb{R}_+^d$. For $p \in (0, \infty]$, the canonical p -norm (quasi-norm when $p < 1$) is defined as $|k|_p := (\sum_{i=1}^d k_i^p)^{1/p}$ with the usual modification for the case $p = \infty$. When $p = 2$, the norm is the standard Euclidean norm on \mathbb{R}^d . In this situation, the subscript is often omitted and we simply use $|\cdot|$. The following convention⁵ is also used $k^\beta = (k^{\beta_1}, \dots, k^{\beta_d})$. Partial orders are denoted by \leq, \geq, \dots . For example, $j \leq k$ is understood as $j_i \leq k_i$, for all $i = 1, \dots, d$.
- Constants are usually designated by capital letters I, J, M, N , etc. Index sets, domains in \mathbb{R}^d , are denoted by Fraktur letters $\mathfrak{J}, \mathfrak{J}, \mathfrak{D}$, etc.
- Real vector spaces are normally denoted by capital letters H, G, X, Y , etc., and their elements by small letters h, g, x, y , etc. Less often, blackboard bold is also used to designate spaces in the following situations. First, the elements in the space are customarily denoted by capital letters. One example is the sample space \mathbb{X} , which contains the observations customarily denoted by X . The other situation is when the letter indicates a particular space, for example, the reproducing kernel Hilbert space \mathbb{H} of a Gaussian measure.
- The symbols $\lesssim, \gtrsim, \simeq$ mean $\leq, \geq, =$ up to a positive multiple independent of n (or another asymptotic parameter). The constant may be stated explicitly in subscripts, and e.g. \lesssim_f means that it depends on f .
- For a normed space $(E, \|\cdot\|)$, the closed unit ball is denoted by $U(E) = \{x \in E : \|x\| \leq 1\}$. Given the topological dual E^* of E , the duality pair

⁵Notice its difference with the notation $k^\beta = \prod_{i=1}^d k_i^{\beta_i}$ sometimes appeared in literature.

is denoted by $\langle \cdot, \cdot \rangle : E^* \times E \rightarrow \mathbb{R}$. The same notation is as well used for inner products on Hilbert spaces. The notation is consistent, since for real spaces sesquilinear form reduces to bilinear form. If there is a danger of confusion, subscripts are used.

Consider function spaces consisting of the functions with domain \mathfrak{D} and codomain a vector space H . The function spaces are denoted similarly as for the real-valued case, e.g. $L^2(\mathfrak{D}; H)$ the space of square integrable H -valued functions, i.e.

$$\int_{\mathfrak{D}} \|f\|_H^2 dx < \infty.$$

For \mathbb{R}^d -valued function spaces, the codomain is often omitted.

- General operators are designated by capital calligraphic letters $\mathcal{A}, \mathcal{T}, \dots$ and Greek letters Λ, Φ, \dots , while there are exceptional cases, e.g. the expectation operator \mathbb{E} , identity mapping id , embedding ι , etc.
- Let μ and ν be two measures. If μ is dominated by (i.e. absolutely continuous to) ν , then it is denoted as $\mu \ll \nu$. If $\mu \ll \nu$ and $\nu \ll \mu$, i.e. they are equivalent measures, then we write $\mu \sim \nu$. Mutual singularity is denoted by $\mu \perp \nu$.
- The following abbreviations are used in this thesis.

almost everywhere		a.e.
almost sure		a.s.

1.6 Notes

Asymptotic statistics, function spaces, are well established research fields and there exist numerous outstanding references. We only list a very small collection here, which by no means intends to be complete or exclusive. It is merely based on the author's familiarity.

Asymptotic Statistics

While written in 1940s, Cramér's book [22] still in large captures the essence of asymptotic theory, and additionally it well reflects the development of asymptotic theory at the early stage. Le Cam's noted treatise [66] largely presents the whole picture of asymptotic methods up to 1980s. Van der Vaart's textbook [97] is another standard reference in asymptotic theory, which provides a comprehensive introduction on the subject and also includes the new developments in 1990s. The more recent monograph [36] systematically depicts the asymptotic theory in nonparametric statistics. For the more detailed literature review, we refer to [66] for results up to 1980s and [36] for nonparametric asymptotics.

Function Spaces

Functions and function spaces serve as the fundamental element for many mathematical studies, and its own study has become an independent field long time ago. In this thesis, we do not use any advanced function spaces. Many classical results of the modern theory of function spaces can be found in [90]. The more recent contributions are largely collected in the sequel [91, 92, 93, 94] from the same author.

Separability

One fact which has been less addressed is the separability of parameter space Θ . A topological space is called *separable* if it contains a countable, dense subset. The previous examples of function spaces $C[0, 1]$ and $L^2[0, 1]$ are both separable. To begin with, separability is important in terms of approximation. Since separability imposes the existence of a countable dense subset, many approximation results can be stated and proved using induction, without invoking axiom of choice (or equivalently, Zorn's lemma). In other words, the induction arguments can be translated into implementable numerical algorithms. If the underlying space is nonseparable, it is no longer necessarily true. Additionally, since separability is one of the fundamental assumptions for the development of probability theory in Banach spaces (see [67]), it is appropriate to adopt the same notion.

For a separable Hilbert space, there always exists a countable orthonormal basis. It serves as the a cornerstone for the development of L^2 estimation theory. In general, for an arbitrary separable Hilbert space H , an isometric isomorphism can be established between H and ℓ^2 using projection (see Sobolev smoothness for example). The separability in Banach spaces is more involved. A separable Banach space does not necessarily have a (Schauder) basis, while a Banach space with a Schauder basis is necessarily separable. More importantly, Banach spaces with Schauder bases also possess *approximation property*, which roughly means that elements in the space can be approximated by finite-dimensional subjects.

Bayesian Methodology

The monograph [35] published in 2017 provides an extensive and thoroughly survey, covering almost all aspects, on the development of Bayesian nonparametrics up to the publication date. In addition, very comprehensive literature reviews are given in many places in the book. Therefore, here we simply refer to it for the reference on Bayesian nonparametrics.

Part I

Foundations

In Part I, we prepare the material essential to the statistical inference for Gaussian linear models, which can be represented in the following form

$$X^{(n)} = \mathcal{A}^{(n)}\theta + \xi^{(n)}. \tag{I.1}$$

Most of the content are established subjects or relatively recent results but known by researchers working in the relevant fields. The readers familiar with the content may skip this part.

In Chapter 2, we survey a Sobolev type of regularity scales that characterises the smoothness used in this thesis.

In Chapter 3, Gaussian measures on Banach spaces are introduced. Especially, we collect the properties of Gaussian measures relevant to the statistical study of Gaussian linear models.

In Chapter 4, we present the framework of Bayesian nonparametric inference. In particular, we demonstrate a general contraction theorem tailored to Gaussian linear models with transformed signals.

Chapter 2

Smoothness Class

Smoothness, also known as *regularity*, characterises how ‘well’ a function behaves. A set of function spaces, each of which consists of functions with the same smoothness property, is called a *smoothness class*. We introduce a scale of smoothness classes, i.e. smoothness classes that are indexed by an ordered set, that is particularly suitable for studying Gaussian linear models in Hilbert space framework.

While most of the statements in this chapter are given in a general Hilbert space setting, they can always be translated to a concrete L^2 function space $L^2(\mathfrak{D}, \mu)$, square integrable functions on a domain \mathfrak{D} with respect to measure μ . When the underlying domain \mathfrak{D} is one-dimensional, the smoothness is naturally considered to be a scalar. In the higher dimensional case, i.e. the dimension of \mathfrak{D} being $d > 1$, a function $f \in L^2(\mathfrak{D}, \mu)$ does not necessarily behave identically along each coordinate direction. Hence, the scalar smoothness should as well adapt to the multi-dimensional nature. In this chapter, we first present the *isotropic* scale, that is a class of spaces parametrised by scalar-valued indices. Subsequently, the concept is generalised to the *anisotropic* scale, with \mathbb{R}^d -valued indices. Studied later in Part II, the inverse problems are placed in isotropic scales. In contrast, for evolution equations examined in Part III, it is more adequate to consider the spatial and temporal regularity of functions separately, and as a consequence, anisotropic scales are used.

In this chapter, we define smoothness scales, introduce important examples that will be used as underlying parameter spaces for the later statistical study, and examine their properties, especially those important for establishing approximation error estimate.

2.1 Smoothness scales

The parameter θ of interest in the Gaussian linear model (I.1) is assumed to be an element of a Hilbert space H . We embed this space as the space $H = H_0$ in a ‘scale of smoothness classes’, defined as follows.

Definition 2.1 (Smoothness scale). For every $s \in \mathbb{R}$ the space H_s is an infinite-dimensional, separable Hilbert space, with inner product $\langle \cdot, \cdot \rangle_s$ and induced norm $\|\cdot\|_s$. The spaces $(H_s)_{s \in \mathbb{R}}$ satisfy the following conditions:

- (i) For $s < t$ the space H_t is a dense subspace of H_s and $\|f\|_s \lesssim \|f\|_t$, for $f \in H_t$.
- (ii) For $s \geq 0$ and $f \in H_0$ viewed as element of $H_{-s} \supset H_0$,

$$\|f\|_{-s} = \sup_{\|g\|_s \leq 1} \langle f, g \rangle_0, \quad f \in H_0. \quad (2.1)$$

The notion of scales of smoothness classes is standard in the literature. In the preceding definition we have stripped it to the bare essentials needed in our general result on posterior contraction. Concrete examples, as well as more involved structures such as Hilbert scales, are introduced in the subsequent section.

Remark 2.2 (Norm duality). The *norm duality* (2.1) is implied if, for $s > 0$, the space H_{-s} can be identified with the dual space H_s^* of H_s and the embedding $\iota : H_0 \rightarrow H_{-s}$ is the adjoint of the embedding $\iota : H_s \rightarrow H_0$, after the usual identification of H_0 and its dual space H_0^* . (The three nested spaces $H_{-s} \supset H_0 \supset H_s$ then form a ‘Gelfand triple’.) Indeed, by definition the image $\iota^* f$ of $f \in H_0 = H_0^*$ under the adjoint $\iota^* : H_0^* \rightarrow H_s^*$ is the map $g \mapsto (\iota^* f)(g) = \langle \iota g, f \rangle_0 = \langle g, f \rangle_0$ from $H_s \rightarrow \mathbb{R}$. The norm of this map as an element of H_s^* is $\sup_{\|g\|_s \leq 1} (\iota^* f)(g)$. The norm duality follows if $\iota^* f$ is identified with the element $f \in H_0 \subset H_{-s}$.

Since every H_s is a Hilbert space, one can also identify H_s^* with itself in the usual way, but this involves the inner product in H_s , and is different from the identification of H_s^* with the ‘bigger space’ H_{-s} .

We assume that the smoothness scale allows good finite-dimensional approximations, as in the following condition.

Assumption 2.3 (Approximation property). For every $j \in \mathbb{N}$ and $s \in (0, S)$, for some $S > 0$, there exists a $(j - 1)$ -dimensional linear subspace $V_j \subset H_0$ and a number $\delta(j, s)$ such that $\delta(j, s) \rightarrow 0$ as $j \rightarrow \infty$, and such that

$$\inf_{g \in V_j} \|f - g\|_0 \lesssim \delta(j, s) \|f\|_s, \quad (2.2)$$

$$\|g\|_s \lesssim \frac{1}{\delta(j, s)} \|g\|_0, \quad \forall g \in V_j. \quad (2.3)$$

This assumption is also common in the literature on numerical analysis, approximation theory, and inverse problems, etc. The two inequalities (2.2) and (2.3) are known as of Jackson and Bernstein type, respectively, see, e.g., [12]. The approximation property (2.2) shows that ‘smooth elements’ $f \in H_s$ are well approximated in $\|\cdot\|_0$ by their projection onto a finite-dimensional space V_j , with approximation error tending to zero as the dimension of V_j tends to infinity. Naturally one expects the numbers $\delta(j, s)$ that control the approximation to be decreasing in both j and s . In our examples we shall mostly have polynomial dependence $\delta(j, s) = j^{-s/d}$, in the case that H_0 consists of functions on a d -dimensional domain. The stability property (2.3) quantifies the smoothness norm of the projections in terms of the approximation numbers. Both conditions are assumed up to a maximal order of smoothness $S > 0$, and it follows from (2.3) that V_j must be contained in the space H_S .

The following estimates derived from the approximation property are convenient for the later study.

Lemma 2.4. *If V_j is a finite-dimensional space as in Assumption 2.3 such that (2.2) and (2.3) hold, then, for $P_j : H_0 \rightarrow V_j$ the orthogonal projection onto V_j , and $0 \leq s, t < S$,*

$$\|f - P_j f\|_{-t} \lesssim \delta(j, t) \delta(j, s) \|f\|_s, \quad f \in H_0, \quad (2.4)$$

$$\|g\|_s \lesssim \frac{1}{\delta(j, s) \delta(j, t)} \|g\|_{-t}, \quad g \in V_j. \quad (2.5)$$

Proof. By the dual norm relation in (ii) of Definition 2.1, and the orthogonality of $f - P_j f$ to V_j ,

$$\begin{aligned} \|f - P_j f\|_{-t} &= \sup_{\|g\|_t \leq 1} \langle f - P_j f, g \rangle_0 = \sup_{\|g\|_t \leq 1} \langle f - P_j f, g - P_j g \rangle_0 \\ &\leq \|f - P_j f\|_0 \sup_{\|g\|_t \leq 1} \|g - P_j g\|_0, \end{aligned}$$

by the Cauchy-Schwarz inequality. Here $\|f - P_j f\|_0 \lesssim \delta(j, s) \|f\|_s$ and $\|g - P_j g\|_0 \lesssim \delta(j, t) \|g\|_t$, both by (2.2). Inequality (2.4) follows.

For the second inequality we have, for $g \in V_j$,

$$\|g\|_0 = \sup_{f \in V_j: \|f\|_0 \leq 1} \langle g, f \rangle_0 \lesssim \sup_{f \in V_j: \|f\|_0 \leq 1} \|g\|_{-t} \|f\|_t,$$

again by the dual norm relation. Here we can bound $\|f\|_t$ by $\|f\|_0 / \delta(j, t)$, with the help of (2.3). We obtain (2.5) by first bounding $\|g\|_s$ with the help of (2.3) and next using the preceding display. \square

The approximation property (2.2) can also be stated in terms of the ‘approximation numbers’ of the canonical embedding $\iota : H_s \rightarrow H_0$. The j th *approximation number* of a general bounded linear operator $T : G \rightarrow H$ between normed spaces is defined as

$$a_j(T : G \rightarrow H) = \inf_{U: \text{Rank } U < j} \sup_{f: \|f\|_G \leq 1} \|(T - U)f\|_H, \quad (2.6)$$

where the infimum is taken over all linear operators $U : G \rightarrow H$ of rank less than j . It is immediate from the definitions that the numbers $\delta(j, s)$ in (2.2) can be taken equal to the approximation numbers $a_j(\iota : H_s \rightarrow H_0)$. The set of approximation numbers $a_j(\iota : H_{s+t} \rightarrow H_t)$ of the canonical embedding describes many characteristics of the smoothness scale $(H_s)_{s \in \mathbb{R}}$. In particular, Assumption 2.3 implies that the canonical embedding $\iota : H_s \rightarrow H_0$ is a limit of a sequence of finite rank operators, and hence is compact. Later we give a brief discussion in Section 2.4.

Example 2.5 (Sobolev classes). The most important examples of smoothness classes satisfying Definition 2.1 are fractional Sobolev spaces on a bounded domain $\mathfrak{D} \subset \mathbb{R}^d$. For a natural number $s \in \mathbb{N}$ the Sobolev space of order s can be defined by

$$H_s(\mathfrak{D}) = W^{s,2}(\mathfrak{D}) := \left\{ f \in \mathcal{D}'(\mathfrak{D}) : \|f\|_s := \sum_{|\alpha| \leq s} \|D^\alpha f\|_{L^2(\mathfrak{D})} < \infty \right\}.$$

Here $\mathcal{D}'(\mathfrak{D})$ is the space of generalized functions on \mathfrak{D} (distributions), i.e. the topological dual space of the space $C_c^\infty(\mathfrak{D})$ of infinitely differentiable functions with compact support in \mathfrak{D} ; the sum ranges over the multi-indices $\alpha = (\alpha_1, \dots, \alpha_d) \in (\{0\} \cup \mathbb{N})^d$ with $|\alpha| := \sum_{i=1}^s \alpha_i \leq s$; and D^α is the differential operator

$$D^\alpha := \frac{\partial^{\alpha_1} \partial^{\alpha_2} \dots \partial^{\alpha_d}}{\partial x_1^{\alpha_1} \partial x_2^{\alpha_2} \dots \partial x_d^{\alpha_d}}.$$

The definition can be extended to $s \in \mathbb{R} \setminus \mathbb{N}$ in several ways. All constructions are equivalent to the Besov space $B_{2,2}^s(\mathfrak{D})$, see [92, 93].

It is well known that the approximation numbers of the scale of Sobolev spaces satisfy Assumption 2.3 with $\delta(j, t) = j^{-t/d}$, see [45].

Example 2.6 (Sequence spaces). Suppose $(\phi_i)_{i \in \mathbb{N}}$ is a given orthonormal sequence in a given Hilbert space H , and $1 \leq b_i \uparrow \infty$ is a given sequence of numbers. For $s \geq 0$, define H_s as the set of all elements $f = \sum_{i \in \mathbb{N}} f_i \phi_i \in H$ with $\sum_{i \in \mathbb{N}} b_i^{2s} f_i^2 < \infty$, equipped with the norm

$$\|f\|_s = \left(\sum_{i \in \mathbb{N}} b_i^{2s} f_i^2 \right)^{1/2}.$$

Then $H_0 = H$ is embedded in H_s , for every $s > 0$, and the norms $\|f\|_s$ are increasing in s . Every space H_s is a Hilbert space; in fact H_s is isometric to H_0 under the map $(f_i) \rightarrow (f_i b_i^s)$, where we have identified the series with their coefficients for simplicity of notation.

For $s < 0$, we equip the elements $f = \sum_{i \in \mathbb{N}} f_i \phi_i$ of H , where $(f_i) \in \ell_2$, with the norm as in the display, which is now automatically finite, and next define H_s as the metric completion of H under this norm. The space H_s is isometric to the set of all sequences $(f_i)_{i \in \mathbb{N}}$ with $\sum_{i \in \mathbb{N}} f_i^2 b_i^{2s} < \infty$ equipped with the norm given on the right hand side of the preceding display, but the series $\sum_{i \in \mathbb{N}} f_i \phi_i$ may not possess a concrete meaning, for instance as a function if H is a function space.

By Parseval's identity the inner product on $H = H_0$ is given by $\langle f, g \rangle_0 = \sum_{i \in \mathbb{N}} f_i g_i$, and the norm duality (2.1) follows with the help of the Cauchy-Schwarz inequality.

The natural approximation spaces for use in Assumption 2.3 are $V_j = \text{Span}(\phi_i : i < j)$. Inequalities are satisfied with the approximation numbers taken equal to $\delta(j, t) = b_j^{-t}$.

2.2 Hilbert Scale

As we are going to show in this section, two Hilbert spaces such that one is dense in the other, naturally define a smoothness scale with additional structure, which is called Hilbert scale. In many applications, the smoothness scales turn out to be Hilbert scales.

A Hilbert scale is generated by a *densely defined unbounded self-adjoint strictly positive operator* $\Lambda : \text{Dom}(\Lambda) \subset H_0 \rightarrow H_0$, with domain $\text{Dom}(\Lambda)$ such that

- (a) $\text{Dom}(\Lambda)$ is dense in H_0 , (i.e. ' Λ is densely defined'),

- (b) $\text{Dom}(\Lambda) = \text{Dom}(\Lambda^*)$,
- (c) $\langle \Lambda x, y \rangle = \langle x, \Lambda y \rangle$ for all $x, y \in \text{Dom}(\Lambda)$, (i.e. ‘ Λ is symmetric’),
- (d) $\langle \Lambda x, x \rangle \geq \kappa \|x\|^2$, for all $x \in \text{Dom}(\Lambda)$, and some $\kappa > 0$.

The set $\text{Dom}(\Lambda^*)$ in (b) is the domain of the adjoint Λ^* of Λ , which is defined as the set of all $y \in H$ such that the map $x \mapsto \langle \Lambda x, y \rangle$ from $\text{Dom}(\Lambda)$ to \mathbb{R} is continuous (see Appendix A). Note that this depends on the domain $\text{Dom}(L)$, which is considered part of the definition of Λ and is restricted by (a) only. Together, requirements (b) and (c) are equivalent to the requirement that Λ be *self-adjoint*.

The domain of the k -th power of the operator Λ is defined, by induction for $k = 2, 3, \dots$, as (with $\Lambda^1 = \Lambda$)

$$\text{Dom}(\Lambda^k) = \{f \in \text{Dom}(\Lambda^{k-1}) : \Lambda f \in \text{Dom}(\Lambda)\}, \quad k > 1.$$

All powers Λ^k , for $k \in \mathbb{N}$, are defined on

$$H_\infty := \bigcap_{k \in \mathbb{N}} \text{Dom}(\Lambda^k). \quad (2.7)$$

It can be shown that H_∞ is dense in H_0 (Lemma 8.17 in [29]). Next, using spectral theory, fractional powers Λ^s can be defined as well on the domain H_∞ , for every $s \in \mathbb{R}$, through integration with respect to the spectral family (E_λ) of Λ , i.e.

$$\Lambda^s := \int_{\mathbb{R}} \lambda^s dE_\lambda = \int_{\kappa}^{\infty} \lambda^s dE_\lambda.$$

This allows to define an inner product on H_∞ by, for $h, g \in H_\infty$ and $s \in \mathbb{R}$,

$$\langle h, g \rangle_s := \langle \Lambda^s h, \Lambda^s g \rangle. \quad (2.8)$$

Definition 2.7 (Hilbert scale). The Hilbert space H_s is the completion of H_∞ with respect to the norm induced by the inner product $\langle \cdot, \cdot \rangle_s$ defined in (2.8). The family $(H_s)_{s \in \mathbb{R}}$ is called the *Hilbert scale generated by Λ* .

The following proposition, adapted from Proposition 8.19 in [29], lists basic properties of Hilbert scales.

Proposition 2.8. *Let Λ be a densely defined unbounded operator satisfying (a)–(d). Then the Hilbert scale $(H_s)_{s \in \mathbb{R}}$ is a smoothness scale in the sense of Definition 2.1, with $\|f\|_s \leq \kappa^{s-t} \|f\|_t$, for $f \in H_t$, and $s < t$.*

In addition, a Hilbert scale possesses the following properties.

- (i) *If $s \geq 0$, then $H_s = \text{Dom}(\Lambda^s)$, and H_{-s} is the dual space of H_s , i.e.*

$$H_{-s} = (H_s)^*.$$

- (ii) *The following interpolation inequality holds. If $f \in H_t$,*

$$\|f\|_s \leq \|f\|_r^\lambda \|f\|_t^{1-\lambda}, \quad \text{with } \lambda = (t-s)/(t-r), \quad (2.9)$$

for $-\infty < r < s < t < \infty$.

Furthermore, for any $s, t \in \mathbb{R}$ the operator Λ^{t-s} has a unique extension from H_∞ to a bounded, self-adjoint operator $\Lambda^{t-s} : H_t \rightarrow H_s$, satisfying

$$(iii) \quad \|\Lambda^{t-s} f\|_s \simeq \|f\|_t, \text{ for } f \in H_t.$$

$$(iv) \quad \Lambda^{t-s} = \Lambda^t \Lambda^{-s}.$$

$$(v) \quad (\Lambda^s)^{-1} = \Lambda^{-s}.$$

Somewhat abusing notation, we have denoted the extension of Λ^{t-s} in the proposition using the same symbol Λ^{s-t} . Taking $s = 0$ or $t = 0$, we see that $\Lambda^s : H_s \rightarrow H_0$ and $\Lambda^s : H_0 \rightarrow H_{-s}$ are norm isomorphisms, for every $s \in \mathbb{R}$. In particular, the unbounded densely defined operator $\Lambda : D(\Lambda) \subset H_0 \rightarrow H_0$ that generates the scale can be extended to a bounded operator $\Lambda : H_1 \rightarrow H_0$, by strengthening the norm on its domain, and also to a bounded operator $\Lambda : H_0 \rightarrow H_{-1}$, by extending its range space and weakening the norm of its range space. Moreover, the inverse map is a norm isomorphism $\Lambda^{-1} : H_0 \rightarrow H_1$, and hence is certainly bounded as an operator $\Lambda^{-1} : H_0 \rightarrow H_0$.

The eigenvalues of Λ^{-1} are closely connected to the approximation property in Assumption 2.3.

Proposition 2.9. *If $\Lambda^{-1} : H_0 \rightarrow H_0$ is compact with eigenvalues $\lambda_j \downarrow 0$, then Assumption 2.3 is satisfied in the Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by Λ , with $\delta(j, t) \simeq \lambda_j^t$ and $S = \infty$. In fact, there exist linear spaces V_j of dimension $j - 1$ such that, for $s \geq 0$ and $t \in \mathbb{R}$,*

$$\inf_{g \in V_j} \|f - g\|_t \lesssim \delta(j, s) \|f\|_{s+t}, \quad (2.10)$$

$$\|g\|_{s+t} \lesssim \frac{1}{\delta(j, s)} \|g\|_t, \quad \forall g \in V_j. \quad (2.11)$$

Proof. Because $\Lambda^{-1} : H_0 \rightarrow H_0$ is compact, there exists an orthonormal basis $(\phi_i)_{i \in \mathbb{N}}$ of eigenfunctions in H_0 . It may be checked that $f = \sum_{i \in \mathbb{N}} f_i \phi_i$ has $L^s f = \sum_{i \in \mathbb{N}} f_i \lambda_i^{-s} \phi_i$, and square norm $\|f\|_s^2 = \sum_{i \in \mathbb{N}} f_i^2 \lambda_i^{-2s}$, provided the latter series converges. Take V_j equal to the linear span of the first $j - 1$ eigenfunctions. Then $f - P_j f = \sum_{i \geq j} f_i \phi_i$ and hence $\|f - P_j f\|_t^2 = \sum_{i \geq j} f_i^2 \lambda_i^{-2t} \leq \lambda_j^{2s} \sum_{i \geq j} f_i^2 \lambda_i^{-2t-2s} \leq \lambda_j^{2s} \|f\|_{s+t}^2$, for $s, t \geq 0$, and for $f \in V_j$ we have $\|f\|_{s+t}^2 = \sum_{i < j} f_i^2 \lambda_i^{-2s-2t} \leq \lambda_j^{-2s} \sum_{i \leq j} f_i^2 \lambda_i^{-2t} = \lambda_j^{-2s} \|f\|_t^2$. \square

The sequence spaces of Example 2.6 are one class of examples of Hilbert scales, generated by the operator $L : (f_i) \mapsto (f_i b_i)$. More intricate Hilbert scales arise from (elliptic) differential operators. These are useful in that they can incorporate boundary conditions, which are then automatically inherited by a Gaussian prior attached to such a scale. The following one-dimensional example is simplistic, but illustrative.

Example 2.10 (Sobolev scales). Consider the one-dimensional negative Laplacian

$$-\Delta = -\frac{d^2}{dx^2}$$

as an operator on the space $C_c^\infty(0, 1)$ of infinitely often differentiable functions with compact support in $(0, 1)$, viewed as subset of $L^2(0, 1)$, with range space $L^2(0, 1)$. On this domain this operator is not self-adjoint, but it has a self-adjoint extension (with differentiation interpreted in the sense of distributions) to the space of all functions $f \in W^{2,2}(0, 1)$ satisfying the *Dirichlet boundary condition*

$$f(0) = 0 = f(1). \quad (2.12)$$

(See Theorem 4.23 in [41].) The eigenfunctions of the Laplacian under the Dirichlet boundary condition are the functions $x \mapsto \sin(j\pi x)$, for $j \in \mathbb{N}$, with eigenvalues of the order $b_j \asymp j^{-1}$. The corresponding Hilbert scale can also be described as the sequence space generated by this orthogonal basis.

Because the Laplacian is a second derivative it is natural to half the scale parameter, or equivalently use the root negative Laplacian $\Lambda := \sqrt{-\Delta}$ as the generator of the scale (where the root is defined through the spectral decomposition).

The boundary conditions play an important role in defining the scale. Technically they are needed to create a domain on which the operator is self-adjoint. An alternative choice to the Dirichlet is the *Cauchy boundary condition*

$$f'(0) = 0 = f(1).$$

This leads to the sequence scale generated by the eigenfunctions $x \mapsto \cos((j - 1/2)\pi x)$, for $j \in \mathbb{N}$, and is different from the Dirichlet scale. Again the eigenvalues of Λ^{-1} are of the order j^{-1} .

Incidentally, it is shown in [73] that the full Sobolev scale ($s \in \mathbb{R}$) of Example 2.5 is not a Hilbert scale for any generating operator Λ . Also in that sense the boundary conditions are essential.

Until now, the Hilbert scales have been constructed from H_0 with the help of a generating operator $(\Lambda, \text{Dom}(\Lambda))$. Alternatively, given a dense subset G of H , the existence of a generating operator Λ with $\text{Dom} \Lambda = G$ is guaranteed, as shown in the following theorem.

Theorem 2.11. *Let G and H be two Hilbert spaces such that G is densely and continuously embedded into H such that $\|g\|_H \leq \|g\|_G$ for all $g \in G$. Then there exists a unique operator Λ , which is positive-definite and self-adjoint, and generates a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ such that $H_1 = G$ and $H_0 = H$.*

Proof. The Λ is constructed using Lemma A.8. The rest is to check that Λ satisfies the properties of a generating operator. \square

In Proposition 2.8, a natural duality structure is possessed by Hilbert scales, i.e. $H_{-s} = (H_s)^*$, given $s \geq 0$. Now we are going to have a closer look at the norm duality mentioned in Remark 2.2 and its connection to Hilbert scales. First let us recall a noted duality structure.

Definition 2.12. Let G and H be two Hilbert spaces such that G is a dense subset of H and the canonical embedding $G \hookrightarrow H$ is continuous. The triplet

$$G \subset H = H^* \subset G^*$$

is a *Gelfand triple*, where the subsets are dense in all bigger spaces.

Remark 2.13. The Gelfand triple is a well-defined object, see Lemma A.7. Its general version is that the *pivot space* H is a Hilbert space and the dense space G is a reflexive Banach space, but for this thesis it is sufficient to only work with Hilbert spaces.

We further elaborate on how to identify dual spaces with the spaces with negative indices in Hilbert scales.

Lemma 2.14. *Let $\{H_s\}_{s \in \mathbb{R}}$ be a Hilbert scale. For any fixed $t > 0$ and $f \in H_{-t}$, the mapping,*

$$\begin{aligned} \mathcal{J} : H_{-t} &\rightarrow (H_t)^*, \\ (\mathcal{J}f)(g) &:= \langle \Lambda^{-2t}f, g \rangle_{H_0}, \end{aligned} \tag{2.13}$$

with the generating operator Λ of $\{H_s\}_{s \in \mathbb{R}}$, is an isometric isomorphism between H_{-t} and $(H_t)^*$.

Proof. It is convenient to recall that $\Lambda^t : H_{s+t} \rightarrow H_s$ is an isometric isomorphism.

For $f \in H_{-t}$, the map $g \mapsto \langle f, g \rangle_{H_{-t}}$ is a bounded linear functional on H_t , and therefore by the Riesz representation theorem, there is a unique $g_f \in H_t$ such that for all $g \in H_t$,

$$\langle f, g \rangle_{H_{-t}} = \langle \Lambda^{-t}f, \Lambda^{-t}g \rangle_{H_0} = \langle \Lambda^{-2t}f, g \rangle_{H_0} = \langle \Lambda^t \Lambda^{-4t}f, \Lambda^t g \rangle_{H_0} = \langle g_f, g \rangle_{H_t},$$

where $g_f = \Lambda^{-4t}f$. Consequently, for any $f \in H_{-t}$ and $g \in H_t$,

$$(\mathcal{J}f)(g) = \langle \Lambda^{-4t}f, g \rangle_{H_t} = \langle \Lambda^{-2t}f, g \rangle_{H_0}.$$

Conversely, if $\ell \in H_t^*$ is a bounded linear functional on H_t , then again by the Riesz representation theorem, there is a unique $h_\ell \in H_t$ such that

$$\ell(h) = \langle h_\ell, h \rangle_{H_t} = \langle \Lambda^{2t}f_\ell, h \rangle_{H_0},$$

where $f_\ell = \Lambda^{2t}h_\ell \in H_{-t}$. Furthermore, since $\|f_\ell\|_{H_{-t}} = \|h_\ell\|_{H_t}$,

$$|\mathcal{J}(f_\ell)(h)| = |\ell(h)| \leq \|h_\ell\|_{H_t} \|h\|_{H_t} = \|f_\ell\|_{H_{-t}} \|h\|_{H_t},$$

which implies $\|\ell\|_{H_t^*} \leq \|f_\ell\|_{H_{-t}}$. On the other hand,

$$\ell(\Lambda^{-2t}f_\ell) = \langle h_\ell, \Lambda^{-2t}f_\ell \rangle_{H_t} = \langle \Lambda^{-t}f_\ell, \Lambda^{-t}f_\ell \rangle_{H_0} = \|f_\ell\|_{H_{-t}}^2.$$

Because $\|f_\ell\|_{H_{-t}} = \|h_\ell\|_{H_t}$, the preceding equation implies $\|\ell\| \geq \|f_\ell\|_{H_{-t}}$.

Combining the results above, we conclude that $\mathcal{J} : H_{-t} \rightarrow (H_t)^*$ is an isometric isomorphism. \square

The connection between Gelfand triples and Hilbert scales is stated in the following result.

Theorem 2.15. *Let $(H_s)_{s \in \mathbb{R}}$ be a Hilbert scale. Then, (H_{s+t}, H_s, H_{s-t}) is a Gelfand triple. Conversely, for any Gelfand triple (G, H, G^*) , there exists a unique Hilbert scale $(H_s)_{s \in \mathbb{R}}$ such that $H_1 = G, H_0 = H$, and $H_{-1} = G^*$.*

Proof. Identify the dual space $(H_s)^*$ of H_s with itself. the first statement follows the same argument from Lemma 2.14 with

$$\begin{aligned} \mathcal{J} &: H_{s-t} \rightarrow (H_{s+t})^*, \\ (\mathcal{J}f)(g) &:= \langle \Lambda^{-2t}f, g \rangle_{H_s}. \end{aligned}$$

The second part is a corollary of Theorem 2.11. \square

2.2.1 Relation to Boundary Conditions

Boundary conditions play an important role in the formulation of multi-dimensional problems, and Hilbert scales naturally cope with this issue. While for functions with domain an interval of the real line a boundary condition just concerns the values at the two endpoints of the interval, on multi-dimensional domains boundary conditions are a subtle issue. In the latter case the boundary is itself a continuous, possibly multi-dimensional, domain, and the boundary condition will involve a space of functions defined on the boundary, an infinite-dimensional space. Generally speaking, Hilbert scales are useful in the sense that the functions in the Hilbert scale automatically satisfy the boundary condition if L is chosen properly. It is tightly connected to the Hilbert (L^2) theory of elliptic equations.

To see this in more detail, consider the case of a partial differential equation

$$\mathcal{L}u = f, \tag{2.14}$$

where \mathcal{L} is a second order elliptic differential operator, and $u, f : \mathfrak{D} \rightarrow \mathbb{R}$ are functions on a bounded domain $\mathfrak{D} \subset \mathbb{R}^d$ with sufficiently regular boundary (so that trace operators are well-defined). Given $f \in L^2(\mathfrak{D})$, different choices of the domain $\text{Dom}(\mathcal{L})$ of the operator L lead to different realizations of solution spaces. The closure of $\mathcal{L}|_{C_c^\infty(\mathfrak{D})}$ as an operator in $L^2(\mathfrak{D})$ is known as the *minimal realization* associated with \mathcal{L} , denoted L_{min} (c.f., Definition 4.2 in [41]). On the other hand, the *maximal realization*, denoted L_{max} , has domain of definition $\{u \in L^2(\mathfrak{D}) : \exists f \in L^2(\mathfrak{D}) \text{ such that } \mathcal{L}u = f \text{ weakly}\}$ (see Definition 4.1 in [41]). In the one-dimensional case when \mathfrak{D} is a bounded interval in \mathbb{R} , the domains $\text{Dom}(\mathcal{L}_{min})$ and $\text{Dom}(\mathcal{L}_{max})$ of the two operators differ only by a two-dimensional space, but when $d \geq 2$ the difference is an infinite-dimensional space. Moreover, the domain $\text{Dom}(\mathcal{L}_{max})$ of the maximal realization can be larger than the canonical Sobolev space $W^{2,2}(\mathfrak{D})$, see Example 2.5 for the definition. In this section we also use

$$W_0^{k,2}(\mathfrak{D}) = \{f \in W^{k,2}(\mathfrak{D}) : D^\alpha f|_{\partial\mathfrak{D}} = 0, |\alpha| \leq k-1\}, \quad k = 1, 2, \dots,$$

which is the closure of $C_c^\infty(\mathfrak{D})$ under $W^{k,2}$ -norm.

In the definition of a Hilbert scale it is assumed that \mathcal{L} is self-adjoint, which requires both the structural property $\langle \mathcal{L}x, y \rangle = \langle x, \mathcal{L}y \rangle$ and that the domains of \mathcal{L} and its adjoint \mathcal{L}^* be identical. The domain of \mathcal{L}^* is determined by the domain of \mathcal{L} (see Appendix A) and hence the latter must be chosen carefully.

As an example, consider the operator $\mathcal{L} = -\Delta$, where Δ is the d -dimensional Laplacian. Given $\text{Dom}(\mathcal{L}_{min}) = W_0^{2,2}(\mathfrak{D})$ (see Theorem 10.19 in [84]), \mathcal{L}_{min} is too

small to be self-adjoint. Indeed, by the Green's identity (only real functions for simplicity)

$$\int_{\mathfrak{D}} (-\Delta u)v dx = \int_{\mathfrak{D}} u((-\Delta)^\dagger v) dx + \int_{\partial\mathfrak{D}} u \partial_\nu v - (\partial_\nu u)v d\sigma,$$

where $(-\Delta)^\dagger$ is the formal adjoint and ∂_ν is the directional derivative in the direction of outward pointing normal ν to the surface element $d\sigma$. This implies that $D(-\Delta)^\dagger \supseteq W^{2,2}(\mathfrak{D}) \cap W_0^{1,2}(\mathfrak{D}) \supsetneq W_0^{2,2}(\mathfrak{D})$. On the other hand, $\text{Dom}(\mathcal{L}_{max}) \supset W^{2,2}(\mathfrak{D})$ is too big (see Exercise 11.10 in [84] for example). There are several self-adjoint extensions \mathcal{L} of the minimal operator \mathcal{L}_{min} that represent boundary conditions. For example, $\mathcal{L} = -\Delta$ with $\text{Dom}(\mathcal{L}_{min}) = W_0^{2,2}(\mathfrak{D})$ has a self-adjoint extension $-\Delta_D$ with $\text{Dom}(-\Delta_D) = W^{2,2}(\mathfrak{D}) \cap W_0^{1,2}(\mathfrak{D})$, corresponding to the Dirichlet condition $f|_{\partial\mathfrak{D}} = 0$ (Theorem 10.19, [84]). For the Neumann condition $\frac{\partial f}{\partial \nu}|_{\partial\mathfrak{D}} = 0$, one may use the variational form of $-\Delta u = f$ (see Section 10.6.2, [84]), to show that there exists a self-adjoint extension $-\Delta_N$ with domain $\text{Dom}(-\Delta_N) = W^{1,2}(\mathfrak{D})$ (Theorem 10.20, [84]).

In the situations given above, $\sqrt{-\Delta}$ can be defined using the spectral measure of the self-adjoint extension of $-\Delta$. Consequently, the Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by $\sqrt{-\Delta}$ (see Definition 2.7) is the scale of Sobolev spaces $\widetilde{W}^{s,2}(\mathfrak{D})$ of (generalized) functions $u \in W^{s,2}(\mathfrak{D})$ that satisfy the corresponding boundary condition. More advanced techniques, such as the pseudo-differential method, are necessary for more sophisticated boundary conditions, see [41, 55].

It is reasonable to consider Sobolev scales of functions that satisfy boundary conditions, since the existence and uniqueness of the solution of the forward problem (2.14) is proved by establishing the fact that \mathcal{L} is isomorphism between Sobolev spaces satisfying boundary conditions, see [55]. As an immediate consequence, the isomorphism of the forward operator $\mathcal{A} = \mathcal{L}^{-1}$ is clear in the context of the corresponding inverse problem in form (I.1).

2.3 Smoothness in Higher Dimensions

In this section, we briefly review how to construct Hilbert spaces of functions on higher dimensional domains. We adopt the convention of notations for multidimensional vectors and multi-indices from Section 1.5.

First we recall some facts about tensor products. Let H_1, H_2 be two real separable Hilbert spaces. For each $f_1 \in H_1, f_2 \in H_2$, define a bilinear form $f_1 \otimes f_2$ acting on $H_1 \times H_2$ such as

$$f_1 \otimes f_2(h_1, h_2) := \langle f_1, h_1 \rangle_{H_1} \langle f_2, h_2 \rangle_{H_2}.$$

On the set \mathcal{T} of finite linear combinations of the bilinear forms defined above, one can define an inner product

$$\langle f_1 \otimes f_2, g_1 \otimes g_2 \rangle_{\mathcal{T}} = \langle f_1, g_1 \rangle_{H_1} \langle f_2, g_2 \rangle_{H_2}. \quad (2.15)$$

The *tensor product* of H_1 and H_2 is the completion of \mathcal{T} under the inner product defined in (2.15), denoted by $H_1 \otimes H_2$. A tensor product space inherits a base

from the original spaces. That is, if $\{\varphi_i\}$ and $\{\psi_j\}$ are orthonormal bases for H_1 and H_2 , then $\{\varphi_i \otimes \psi_j\}_{i,j}$ is an orthonormal basis for $H_1 \otimes H_2$.

We are only interested in the case that H_1 and H_2 are L^2 spaces on Euclidean domains. The results needed are summarised in the following theorem.

Theorem 2.16 (Theorem II.10, [81]). *For $i = 1, 2$, let (\mathfrak{D}_i, μ_i) be measure spaces such that $L^2(\mathfrak{D}_i, \mu_i)$ are separable. Then, the following statements hold.*

(i) *There is a unique isomorphism from $L^2(\mathfrak{D}_1, \mu_1) \otimes L^2(\mathfrak{D}_2, \mu_2)$ to $L^2(\mathfrak{D}_1 \otimes \mathfrak{D}_2, \mu_1 \otimes \mu_2)$ such that*

$$f \otimes g \mapsto fg.$$

(ii) *Let \tilde{H} be a separable Hilbert space. Then there is a unique isomorphism from $L^2(\mathfrak{D}_1, \mu_1) \otimes \tilde{H}$ to $L^2(\mathfrak{D}_1, \mu_1; \tilde{H})$ such that*

$$f(x) \otimes h \mapsto f(x)h.$$

(iii) *In particular, there is a unique isomorphism from $L^2(\mathfrak{D}_1 \otimes \mathfrak{D}_2, \mu_1 \otimes \mu_2)$ to $L^2(\mathfrak{D}_1; L^2(\mathfrak{D}_2, \mu_2))$ such that*

$$f(x, y) \mapsto (x \mapsto f(x, \cdot)),$$

and

$$\int_{\mathfrak{D}_1 \times \mathfrak{D}_2} |f(x, y)|^2 dx dy = \int_{\mathfrak{D}_1} \|f(x, \cdot)\|_{L^2(\mathfrak{D}_2, \mu_2)}^2 dx.$$

2.3.1 Multi-dimensional Smoothness

Sobolev spaces on multi-dimensional domains are defined using the general statement from the previous subsection. Throughout this subsection, we assume that \mathfrak{D} is a bounded domain in \mathbb{R}^d with sufficiently regular boundary, e.g C^k with $k \in \mathbb{N}$ larger than the order of \mathcal{L} as in Section 2.2.1.

Let H_1 be $L^2([0, T]; \mathbb{R})$ and H_2 be $L^2(\mathfrak{D}; \mathbb{R})$. We will omit the codomain if it is the real space \mathbb{R} . Following from Theorem 2.16,

$$L^2([0, T]; L^2(\mathfrak{D})) \cong L^2([0, T] \times \mathfrak{D}) \cong L^2(\mathfrak{D}) \otimes L^2([0, T]) \cong L^2([0, T]) \otimes L^2(\mathfrak{D}),$$

Moreover, by Theorem 2.16, if there exist orthonormal bases $\{\varphi_i\}_{i \in \mathbb{N}^d}$ for $L^2(\mathfrak{D})$ and $\{\psi_i\}_{i \in \mathbb{N}}$ for $L^2([0, T])$, the tensor orthonormal basis $\{\varphi_i \otimes \psi_j\}_{(i,j) \in \mathbb{N}^{d+1}}$ is an orthonormal basis for $L^2(\mathfrak{D}) \otimes L^2([0, T])$. In particular, there exists a unique isomorphism from $H \otimes L^2([0, T])$ to $L^2([0, T]; L^2(\mathfrak{D}))$ so that $\varphi \otimes f(t) \mapsto \varphi f(t)$. As a consequence, for any element $f \in L^2([0, T]; L^2(\mathfrak{D}))$, it admits the following representation

$$f(x, t) = \sum_{(i,j) \in \mathbb{N}^{d+1}} f_{(i,j)} \varphi_i \otimes \psi_j(x, t) = \sum_{(i,j) \in \mathbb{N}^{d+1}} f_{(i,j)} \varphi_i(x) \psi_j(t),$$

with

$$f_{(i,j)} = \langle f(x, t), \varphi_i \otimes \psi_j \rangle_{L^2(dx \times dt)} = \int_{[0, T]} \int_{\mathfrak{D}} f(x, t) \varphi_i(x) dx \psi_j(t) dt.$$

Introduce $\mathfrak{D}_T := \mathfrak{D} \times [0, T]$. The tensor space is isomorphic to the ordinary space $L^2(\mathfrak{D}_T)$, by the isomorphism statement above. On the other hand, the series representation above sheds light on how to obtain concrete smoothness scales centred at $L^2(\mathfrak{D}_T)$.

Consider an arbitrary L^2 space $L^2(\tilde{\mathfrak{D}})$ with a bounded domain $\tilde{\mathfrak{D}} \subset \mathbb{R}^m, m \in \mathbb{N}$. Since it is a separable Hilbert space, there exists an orthonormal basis $\{\varphi_k\}_{k \in \mathbb{N}^m}$. By abstract Parseval's identity, any function f in $L^2(\tilde{\mathfrak{D}})$ admits a series expansion such that $f = \sum_k f_k \varphi_k$ in L^2 sense. With a multi-index $\beta = (\beta_1, \dots, \beta_m) \in \mathbb{R}_+^m$, an anisotropic smoothness class is defined to be the completion of

$$\left\{ f = \sum_{k \in \mathbb{N}^m} f_k \varphi_k \mid \|f\|_{H_\beta} := \left(\sum_{k \in \mathbb{N}^m} |\lambda_k^\beta|^2 f_k^2 \right)^{1/2} < \infty \right\}, \quad (2.16)$$

where $\lambda_k = (\lambda_{k_1}, \dots, \lambda_{k_m})$ for $k \in \mathbb{N}^m$, and

$$|\lambda_k^\beta| = |\lambda_k^\beta|_2 = \left(\sum_{i \leq m} \lambda_{k_i}^{2\beta_i} \right)^{1/2}.$$

Similarly, the type of smoothness classes above can be introduced to sequence spaces. Let $\ell^2(\mathbb{N}^m)$ with $m \in \mathbb{N}$ be the m -dimensional square integrable sequence space, i.e. for any $f = \{f_k\}_{k \in \mathbb{N}^m} \in \ell^2(\mathbb{N}^m)$,

$$\|f\|_{\ell^2} = \left(\sum_{k \in \mathbb{N}^m} f_k^2 \right)^{1/2} < \infty.$$

The anisotropic ellipsoid h_β with $\beta \in \mathbb{R}_+^m$ is the completion of the set $f \in \ell^2(\mathbb{N}^m)$ such that, with same $\{\lambda_k\}_k$,

$$\|f\|_{h_\beta} = \left(\sum_{k \in \mathbb{N}^m} |\lambda_k^\beta|^2 f_k^2 \right)^{1/2} < \infty \quad (2.17)$$

under the norm $\|\cdot\|_{h_\beta}$.

The Sequence space is convenient as we often deal with the coefficients. Once the basis $\{\varphi_k\}_{k \in \mathbb{N}^m}$ is fixed, given $f = \sum_k f_k \varphi_k$ and $\tilde{f} = \{f_k\}_{k \in \mathbb{N}^m}$, we have the isometry,

$$\|f\|_{L^2} = \|\tilde{f}\|_{\ell^2}, \quad \|f\|_{H_\beta} = \|\tilde{f}\|_{h_\beta}.$$

Remark 2.17. It is worth noting that so far we only assume that $\{\varphi_k\}$ is an orthonormal basis of $L^2(\mathfrak{D})$ and the smoothness is characterised by weighted ℓ^2 norms of the coefficients. In order to establish the connection to canonical Sobolev spaces, additional requirements are necessary, see the upcoming subsections.

Different types of smoothness might be more suitable for various problems. As we will see in the study of the inference for evolution equations in Part III, for the

recovery of an initial condition, a type of spatial smoothness classes is expected, while for the recovery of a drift term, we need to introduce a proper smoothness class to describe the space-time regularity. To distinguish the different types of smoothness, we call a smoothness class *isotropic* when the smoothness index β satisfies $\beta_i = \beta_j$ for all $i, j \leq m$, or otherwise *anisotropic*.

It is also convenient to introduce the harmonic mean, which will be used to describe the ‘balanced’ smoothness of anisotropic Sobolev spaces (see Section 10.4). For a multi-index $\beta \in \mathbb{R}_+^m$, the harmonic mean is defined as

$$\mathcal{H}(\beta) := \frac{m}{\sum_{i=1}^m (1/\beta_i)}. \quad (2.18)$$

Below we collect some elementary but useful lemmas related to the harmonic mean, which will be used mainly in Part III. We restrict to the case that $\lambda_{k_i} \simeq k_i, i = 1, \dots, m$, in (2.16).

Lemma 2.18. *Given a multi-index $\alpha = (\alpha_1, \dots, \alpha_m) \in \mathbb{R}_+^m$, if $\mathcal{H}(\alpha) > m/2$, then $\sum_{k \in \mathbb{N}^m} |k^\alpha|^{-2}$ is finite.*

Proof. Since $|k^\alpha|^2 = \sum_{i \leq m} k_i^{2\alpha_i} \geq \frac{1}{m} \prod_{i \leq m} k_i^{2\alpha_i/m}$,

$$\sum_{k \in \mathbb{N}^m} |k^\alpha|^{-2} \leq m \sum_{k \in \mathbb{N}^m} \prod_{i \leq m} k_i^{-2\alpha_i/m}.$$

Introduce the hypercubes

$$C_n = \{k \in \mathbb{N}^m : k_i \lesssim n^{\mathcal{H}(\alpha)/\alpha_i}, i = 1, \dots, m\},$$

where the constants are independent from m, i and n . The number of points in \mathbb{N}^m covered by C_n is $\#C_n \simeq n^m$, and consequently $\#[C_n \setminus C_{n-1}] \simeq n^{m-1}$. Hence, the summation can be estimated as,

$$\sum_{n \in \mathbb{N}} \sum_{k \in [C_n \setminus C_{n-1}]} \prod_{i \leq m} k_i^{-2\alpha_i/m} \simeq \sum_{n \in \mathbb{N}} n^{m-1} n^{-2\mathcal{H}(\alpha)}.$$

The result immediately follows. \square

The following elementary lemma entitles us the freedom to choose between equivalent norms.

Lemma 2.19. *Let $\beta = (\beta_1, \dots, \beta_m) \in \mathbb{R}_+^m$ and $k = (k_1, \dots, k_m)$ be a multi-index. Then, with a constant only dependent on m ,*

$$|k^\beta|^2 = \sum_{i=1}^m k_i^{2\beta_i} \simeq_m \left(\sum_{i=1}^m k_i^{\beta_i} \right)^2.$$

In particular, if $\beta_1 = \dots = \beta_m = \beta > 0$, for any $q \in [0, \infty]$,

$$|k^\beta|_2 \simeq_m |k|_q^\beta,$$

where the constant is between 1 and m .

Proof. The first equivalence is elementary. The second one is argued as follows. Recall $|\cdot|_p$ is the p -norm on \mathbb{R}^d . From power mean inequality one can derive $|v|_p \leq |v|_r \leq d^{1/r-1/p} |v|_p$ for $0 < r < p$. The proof is concluded by applying the inequality to k^β with $\beta = (\beta, \dots, \beta) \in \mathbb{R}_+^d$. \square

The following lemma provides an error estimate on the truncation of a sequence in an anisotropic Sobolev ellipsoid, which can be easily translated to the projection in a Sobolev space.

Lemma 2.20. *Let P_N be the projection of a sequence to the coordinates*

$$\{k < N : k_i < N_i, i \leq m\}$$

with $N = (N_1, \dots, N_m)$. For $f \in h_\beta$ with $\lambda_{k_i} \simeq k_i, i = 1, \dots, m$,

$$\|P_N f - f\|_{\ell^2}^2 \leq \left(\sum_{i \leq m} N_i^{-2\beta_i} \right) \|f\|_{h_\beta}^2.$$

In particular, when $N_i \simeq n^{\mathcal{H}(\beta)/(\beta_i m)}$, we have $\prod_{i \leq m} N_i \simeq n$ and

$$\|P_N f - f\|_{\ell^2} \lesssim_m n^{-\mathcal{H}(\beta)/m} \|f\|_{h_\beta}.$$

Proof.

$$\|P_N f - f\|_{\ell^2}^2 \leq \sum_{i \leq m} \left[\sum_{k_i \geq N_i} \sum_{\substack{k_j \in \mathbb{N}: \\ j \leq m \\ j \neq i}} f_k^2 \right] = \sum_{i \leq m} S_i.$$

The items S_i can be bounded from above by,

$$\sum_{k_i \geq N_i} \sum_{\substack{k_j \in \mathbb{N}: \\ j \leq m \\ j \neq i}} \frac{|k^\beta|^2}{N_i^{-2\beta_i}} f_k^2 \leq N_i^{-2\beta_i} \|f\|_{h_\beta}^2.$$

\square

Using the lemmas above, one can show that the anisotropic smoothness class is a multi-dimensional version of a smoothness class.

Corollary 2.21. *Consider an anisotropic smoothness class defined in (2.16) with $\lambda_{k_i} \simeq k_i, i = 1, \dots, m$, then it is a (multi-dimensional) smoothness scale in the following sense.*

- (i) For $\mathbf{s} < \mathbf{t}$, i.e. $s_i \leq t_i$ for all $i \leq m$ and $s_i < t_i$ for at least one $i \leq m$, the space $H_{\mathbf{t}}$ is a dense subspace of $H_{\mathbf{s}}$ and $\|f\|_{\mathbf{s}} \lesssim \|f\|_{\mathbf{t}}$, for $f \in H_{\mathbf{t}}$.
- (ii) For $\mathbf{s} \geq 0$, i.e. $s_i \geq 0$ for all $i \leq m$, $f \in H_0$ can be viewed as element of $H_{-\mathbf{s}} \supset H_0$,

$$\|f\|_{-\mathbf{s}} = \sup_{\|g\|_{\mathbf{s}} \leq 1} \langle f, g \rangle_0, \quad f \in H_0. \quad (2.19)$$

Furthermore, Assumption 2.3 is satisfied with $\delta(j, \mathbf{s}) \simeq j^{-\mathcal{H}(\mathbf{s})/m}$.

2.4 Approximation Number and Metric Entropy

The j th *approximation number* of a bounded linear operator $T : G \rightarrow H$ between normed spaces is defined as

$$a_j(T : G \rightarrow H) = \inf_{U: \text{Rank } U < j} \|T - U\|_{G \rightarrow H}, \quad (2.20)$$

where the infimum is taken over all linear operators $U : G \rightarrow H$ of rank (i.e., dimension of the range space) strictly less than j , and the norm on the right is the operator norm $\|T - U\|_{G \rightarrow H} = \sup_{f: \|f\|_G \leq 1} \|(T - U)f\|_H$. The approximation numbers measure the possibility of approximating an operator by simpler operators of finite-dimensional rank. There is a rich literature on approximation numbers. The main purpose of the present section is to note their relationship to singular values and to metric entropy. Metric entropy plays an important role in the characterization of contraction rates of Bayesian posterior distributions.

If $G \subset H$, we can take T equal to the embedding $\iota : G \rightarrow H$, and then by linearity we see that there exists an operator U of rank smaller than j such that

$$\|f - Uf\|_H \lesssim a_j(\iota : G \rightarrow H) \|f\|_G, \quad \forall f \in G.$$

If H is a Hilbert space, then the minimizing finite-rank operator U is of course the orthogonal projection P_j on V_j . However, the approximation numbers also ‘search’ an optimal projection space. If we take $G = H_s$ and $H = H_0$, then the range space V_j of U satisfies the approximation property (2.2), with the numbers $\delta(j, s)$ taken equal to the approximation numbers $a_j(\iota : H_s \rightarrow H_0)$.

The approximation number is an example of an *s-number*, as introduced in [76]. In general *s-numbers* are defined as maps $T \mapsto (s_j(T))_{j \in \mathbb{N}}$, attaching to every operator T a sequence of nonnegative numbers $s_j(T)$, satisfying certain axiomatic properties. In general, approximation numbers attached to operators $T : H \rightarrow H$ are the ‘largest’ possible *s-numbers*, but on Hilbert spaces there is only one *s-number*: all *s-numbers* are the same (see 2.11.9 in [77]). Because the singular values are also *s-numbers*, the latter unicity yields the important relation that the approximation numbers of operators on Hilbert spaces are equal to their singular values. Recall here that the singular values of a compact operator $T : G \rightarrow H$ are the roots of the eigenvalues of the self-adjoint operator $T^*T : G \rightarrow G$.

The finite-rank approximations U that (nearly) achieve the infimum in the definition of the approximation numbers for different j are not a-priori ordered. However, in many cases there exists a basis $(\phi_i)_{i \in \mathbb{N}}$ such that the projections on the linear span of the first $j - 1$ basis elements achieve the infimum. For Sobolev spaces e.g. spline bases, the Fourier basis, or wavelet bases are all ‘optimal’ in this sense (see [18, 79]).

Approximation numbers are strongly connected to metric entropy. In the literature the connection is usually made through the notion of ‘entropy numbers’, which are defined as follows. The j -th *entropy number* $e_j(T)$ of an operator $T : G \rightarrow H$ is defined as the infimum of the numbers $\varepsilon > 0$ so that the image $T(U_G) \subset H$ of the unit ball U_G in G can be covered by 2^{j-1} balls of radius ε in

H ; or more formally, with U_H the unit ball in H ,

$$e_j(T) = \inf \left\{ \varepsilon > 0 : T(U_G) \subset \bigcup_{i=1}^{2^{j-1}} (h_i + \varepsilon U_H), \text{ for some } h_1, \dots, h_{2^{j-1}} \in H \right\}.$$

The function $j \mapsto e_j(T)$ is roughly the inverse function of the metric entropy of $T(U_G)$ relative to the metric induced by $\|\cdot\|_H$. Recall that the *metric entropy* of a metric space (U, d) is the logarithm of the covering number $N(\varepsilon, U, d)$, which is the minimal number of d -balls of radius $\varepsilon > 0$ needed to cover the space U . Presently we consider the metric entropy $H(\varepsilon, T) = \log N(\varepsilon, T(U_G), \|\cdot\|_H)$ of $T(U_G)$ under the metric of H . Roughly we have that

$$N(\varepsilon, T(U_G), \|\cdot\|_H) \simeq 2^{j-1}, \quad \text{if} \quad e_j(T) \simeq \varepsilon.$$

If we use the logarithm at base 2, then the map $\varepsilon \mapsto H(\varepsilon, T)$ is approximately inverse to the map $j \mapsto e_j(T)$.

Now it is proved in [26] that for any operator $T : G \rightarrow H$ between Hilbert spaces with infinite-dimensional ranges:

$$e_{j+1}(T) \leq 2a_{J+1}(T) \leq 2\sqrt{2}e_{J+2}(T),$$

for any natural numbers j, J satisfying:

$$j \log 2 \geq 2 \sum_{i=1}^J \log \frac{3a_i(T)}{a_{J+1}(T)}.$$

As shown in [26] this relationship between entropy numbers and approximation numbers may be solved to derive the entropy number from the approximation numbers in many cases.

The following lemma gives one example, important to the present thesis.

Lemma 2.22 (Metric entropy). *For a smoothness scale $(H_s)_{s \in \mathbb{R}}$ satisfying (2.2) with $\delta(j, s) = j^{-s/d}$, and $s > 0$ and $t \geq 0$,*

$$\log N(\varepsilon, \{f \in H_s : \|f\|_s \leq 1\}, \|\cdot\|_{-t}) \sim \varepsilon^{-d/(s+t)}. \quad (2.21)$$

Proof. By (2.5) the approximation number $a_j(\iota : H_s \rightarrow H_{-t})$ is of the order $\delta(j, s)\delta(j, t) = j^{-(s+t)/d}$. It is shown in [26] that the entropy numbers $e_j(\iota : H_s \rightarrow H_{-t})$ are of the order $j^{-(s+t)/d}$. By the preceding reasoning this can be inverted to obtain the order of the metric entropy of the image of the unit ball in H_{-t} . \square

Similarly, the results attained above can be extended to the multi-dimensional scale as follows.

Lemma 2.23. *Let $h_\beta(\mathbb{N}^m)$ be a Sobolev ellipsoid given in (2.17). When $\beta > 0$, i.e. $\beta_i > 0$, for all $1 \leq i \leq m$, then both the approximation number and entropy number of the canonical embedding $\iota : h_\beta(\mathbb{N}^m) \rightarrow \ell^2(\mathbb{N}^m)$ are of the order*

$$a_j(\iota : h_\beta \rightarrow \ell^2) \simeq e_j(\iota : h_\beta \rightarrow \ell^2) \simeq j^{-\mathcal{H}(\beta)/m}.$$

Proof. Let P_N be the truncation of a sequence to $\{k < N : k_i < N_i, i \leq m\}$ with $N = (N_1, \dots, N_m)$. From Lemma 2.20, for the truncation at n coefficients, i.e. $\prod_{i \leq m} N_i \simeq n$, the minimal projection error

$$\|P_N f_0 - f_0\|_{\ell^2} \lesssim_d n^{-\mathcal{H}(\beta)/m} \|f\|_{h_\beta}$$

is achieved when N_i are balanced, i.e. $N_i \simeq n^{\mathcal{H}(\beta)/(\beta_i m)}$ for all $i \leq m$. This estimate leads to the upper bound of the approximation number

$$a_n(\iota : h_\beta \rightarrow \ell^2) \lesssim_d n^{-\mathcal{H}(\beta)/m}.$$

Taking a sequence whose only nonzero entry is 1 and its multi-index satisfies $k_i \simeq n^{\mathcal{H}(\beta)/(\beta_i m)}$ for $1 \leq i \leq d$, it is also straightforward to show that

$$a_n(\iota : h_\beta \rightarrow \ell^2) \gtrsim n^{-\mathcal{H}(\beta)/m}.$$

Then by the uniqueness of s -number on Hilbert spaces, we conclude that the entropy number is of the same order as the approximation number. \square

Corollary 2.24 (Metric entropy). *Under the same assumption in Lemma 2.23, the metric entropy is given by, as $\varepsilon \downarrow 0$,*

$$H(\varepsilon, \iota) := \log N\left(\varepsilon, \{f \in h_\beta : \|f\|_{h_\beta} \leq 1\}, \|\cdot\|_{\ell^2}\right) \sim \varepsilon^{-m/\mathcal{H}(\beta)}.$$

Proof. The proof is same as Lemma 2.22 and hence is omitted. \square

In a similar way it is possible to invert approximation numbers that are not of the polynomial form $j^{-s/d}$. There are many examples of this type, for instance, involving additional logarithmic terms, or exponentially decreasing rates. We defer the related discussion until Chapter 10 in Part III.

2.5 Notes

Hilbert scales

Hilbert scales and the relate concepts such as Gelfand triples have been studied be many authors from various fields with different motivations. The Gelfand triples were introduced in [32] to study the theory of generalised functions¹. The idea of scales of function spaces can be traced back to Krein's work, [62, 63], whose main focus was on the interpolation theory of linear operators. In [28], rigged Hilbert space was studied under another name 'Sobolev towers' in the context of operator semigroups. In the field of stochastic analysis, Hilbert scales can be found in Hida's analysis of white noise functionals (as known as generalised stochastic processes) [47], and it is also used to construct the solution spaces in some studies of stochastic partial differential equations [48, 85]. Its application in regularization theory can be found e.g. in Chapter 8 of [29]. For the application in physics, we refer to the survey paper [31].

¹We would like to point out that [32] contains a flaw as addressed by the translator. On the bottom of page 122, the translator raised a concern on the proof of spectral theorem of normal operators in rigged Hilbert spaces. The correct proof is given in [40].

Compactness

Compared to separable spaces mentioned in the Section 1.6, the elements in a compact space can be estimated by a finite dimensional approximation but yet retaining reasonable accuracy. This property is important for any estimation procedure. [50] systematically describes the connections between compactness and statistical estimation. Applications can be found in [51, 52], which are also related to Part III of this thesis.

Chapter 3

Gaussian Analysis

The study of Gaussian distributions has been an active field since the first invention of the normal distribution by Gauss in the nineteenth century. By generalizing the elementary concepts of mean and covariance of a Gaussian distribution to vector spaces, it leads to the study of Gaussian measure on vector spaces. The importance of Gaussian measure cannot be overstated. One classical example is Brownian motion, investigated by Wiener in the first half of the twentieth century, which is the central achievement in the development of stochastic processes. In 1940s and 1950s, several attempts to develop stochastic integration theory with Brownian motion as the integrator resulted in the well known Ito and Stratonovich calculus, which were later generalized to more general integrators such as Lévy process. In seventies, Malliavin in his celebrated seminal work established the stochastic version of the calculus of variation, now known as Malliavin calculus.

Gaussian analysis is also of great importance in applications. In 1973, Black and Scholes published the famous Black-Scholes formula, which revived the mathematical study of financial market using Gaussian models, which can be traced back to Bachelier's work in 1900. Since then, the stochastic differential equation driven by Wiener process has been a fundamental tool to construct many market models, see e.g. [86].

First, in Section 3.1 we give an overview on measures on Banach spaces. Gaussian measures are introduced in Section 3.2. We also briefly discuss how to radonify cylindrical Gaussian measures in Section 3.3.

In this chapter, we always consider the following situation, unless explicitly stating otherwise. Let $(E, \|\cdot\|)$ be a *separable* Banach space. Denote its (topological) dual space by E^* with the (canonical) duality pair by $\langle \cdot, \cdot \rangle : E^* \times E \rightarrow \mathbb{R}$, i.e. $\langle x^*, x \rangle := x^*(x)$, and its Borel σ -algebra by $\mathcal{B}(E)$ (or simply \mathcal{B}).

3.1 Probability Measures on Banach Spaces

A Gaussian measure can be defined at different levels of generality. While it is standard to consider locally convex topological space as the working space when studying Gaussian measures, it is sufficient for us to only consider Gaussian measure on separable Banach spaces. Actually, we even only use the results of Gaussian

measures on Hilbert spaces, but it does not require many extra efforts to state the theory in Banach spaces. Since completeness is standard, we explain in short the reason for imposing separability. First, since the parameter space Θ is separable and the transform operator \mathcal{A} is bounded, $\mathcal{A}(\Theta)$ is necessarily separable, because separability is topologically invariant. Therefore, it is reasonable to only consider the case that $\mathcal{A}(\Theta)$ embeds in another separable space. Consequently, it is sufficient to define the Gaussian noise in a separable space. Second, when introducing Gaussian priors, the parameter space Θ is also assumed to be separable. Hence, separable Banach spaces are sufficient to serve our purposes.

We first consider some properties of general probability measures on vector spaces. The following type of probability measures is of special interest to us, as many theorems below use it as a premise.

Definition 3.1 (Radon measure). A finite Borel measure μ on a Hausdorff topological space E is a *Radon measure* if

$$\mu(B) = \sup\{\mu(K) : K \subset B, K \text{ compact}\},$$

for each $B \in \mathcal{B}(E)$.

Due to the following result, we will always deal with Radon measures.

Theorem 3.2 (Theorem 3.1, Chapter II, [96]). *Every Borel probability measure on a Polish space, i.e. a complete separable metric space, is a Radon measure.*

When considering a probability measure, the σ -algebra has great influence on the properties of the probability. The following remark explains why it is always sufficient to consider Borel σ -algebra.

Remark 3.3 (Choice of σ -algebra). Strictly speaking, probability measures on a vector space E are defined with the *cylindrical* σ -algebra $\mathcal{C}(E)$ generated by the cylindrical sets of E (see Section 3.3), the smallest σ -algebra such that all continuous functionals in E^* are measurable. However, for a F chet space F , $\mathcal{C}(F) = \mathcal{B}(F)$. Since all separable Banach spaces are F chet, there is no need to distinguish \mathcal{C} and \mathcal{B} in our setting. See Appendix A.3 in [8].

Fourier transform is a useful tool for showing the uniqueness of measures.

Lemma 3.4 (Uniqueness of measures by Fourier transform (Lemma 7.13.5, [9])). *For a measure μ on (E, \mathcal{C}) , its Fourier transform, also known as characteristic functional, $\hat{\mu}$ is defined by*

$$\hat{\mu} : E^* \rightarrow \mathbb{C}, \quad \hat{\mu}(f) = \int_E e^{if(x)} \mu(dx).$$

Two measures on \mathcal{C} are identical, if they have same Fourier transform. In particular, it is true for Radon measures.

For a random element in an infinite-dimensional Banach space E , the most fundamental concepts, mean and covariance, are defined as linear actions on the dual space E^* . While, topological support is also important for measures on vector spaces, on the contrary to measures on Euclidean spaces.

Definition 3.5. Let μ be a probability measure on (E, \mathcal{B}) such that $E^* \subset L^2(\mu)$. The *mean* is an element a_μ in the algebraic dual $(E^*)'$ of E^* satisfying

$$a_\mu f = \int_E f(x) \mu(dx).$$

The *covariance* operator $\mathcal{C}_\mu : E^* \rightarrow (E^*)'$ is given by

$$\mathcal{C}_\mu(f)(g) := \int_E [f(x) - a_\mu(f)][g(x) - a_\mu(g)] \mu(dx). \quad (3.1)$$

The *topological support* of μ is the smallest closed set $S_\mu \in E$ such that

$$\mu(E \setminus S_\mu) = 0.$$

Remark 3.6. The general definition of mean and covariance is given by replacing E by a locally convex topological space and \mathcal{B} by the σ -algebra generated by cylindrical sets, see Section II.3 and III.2 in [96].

Similar to the finite-dimensional case, the existence of mean and covariance relies on certain moment conditions. Let the law of X be μ on E . The *p th-weak moment* of X is given by

$$m_p(X) := \sup_{f \in U(E^*)} \int_E |\langle f, x \rangle|^p \mu(dx), \quad (3.2)$$

where $U(E^*)$ is the unit ball of E^* , and $\mathbb{E}\|X\|^p$ is the *p th-strong moment*. The criteria for the existence of mean and covariance are given as follows.

Lemma 3.7. *Let X be a random element on E with distribution μ .*

- (i) *If μ is Borel measurable and $\mathbb{E}\|X\| < \infty$, then a_μ is an element of X .*
- (ii) *If μ is Radon and $m_2(X) < \infty$, then $\mathcal{C}_\mu(E^*) \subset E$.*
- (iii) *If μ is Radon and $\mathbb{E}\|X\|^2 < \infty$, then in addition to (ii), $\mathcal{C}_\mu : E^* \rightarrow E$ is nuclear, i.e. it admits the representation,*

$$\mathcal{C}_\mu f = \sum_{j \in \mathbb{N}} \kappa_j \langle f, a_j \rangle b_j,$$

where $\{a_j\}, \{b_j\}$ are sequences in E such that $\|a_j\| = \|b_j\| = 1$ for $j \in \mathbb{N}$, and $\{\kappa_j\}$ is a real sequence such that $\sum_{j \in \mathbb{N}} |\kappa_j| < \infty$.

Proof. The statement (i) is Lemma 2.1 in [100], and the rest are Theorem 2.1 and 2.3 in chapter III, [96]. \square

Some properties of the covariance operator \mathcal{C}_μ are summarised in the following lemma.

Lemma 3.8 (Factorization of covariance operators.). *Let μ be a Radon probability measure with finite second order strong moment. The covariance operator $\mathcal{C}_\mu : E^* \rightarrow E$, given by*

$$\text{Cov}(f(X), g(X)) = \langle f, \mathcal{C}_\mu g \rangle, \quad f, g \in E^*,$$

is symmetric, i.e.

$$\langle f, \mathcal{C}_\mu g \rangle = \langle g, \mathcal{C}_\mu f \rangle, \quad f, g \in E^*,$$

positive¹, i.e. $\text{Var } f = \langle f, \mathcal{C}_\mu f \rangle \geq 0$ for all $f \in E^*$.

In addition, there exists a Hilbert space H and a bounded linear operator $\mathcal{S} : E^* \rightarrow H$ with its dual $\mathcal{S}^* : H \rightarrow (E^*)^*$, whose range is $\mathcal{S}^*(H) \subset E$, such that the following factorization holds,

$$\mathcal{C}_\mu = \mathcal{S}^* \mathcal{S}.$$

In particular, $\mathcal{C}_\mu(E^*) = \mathcal{S}^*(H)$.

Proof. The factorization is Lemma 1.1 and 1.2 in chapter III, [96]. The rest are immediate consequences from the definition. \square

In Lemma 3.8, the factorization is not unique in the following sense. If there exists another Hilbert space \tilde{H} and an isometry $\mathcal{U} : H \rightarrow \tilde{H}$, then \mathcal{C}_μ can as well be factorized as $\mathcal{C}_\mu = \tilde{\mathcal{S}}^* \tilde{\mathcal{S}}$ with

$$\tilde{\mathcal{S}} = \mathcal{U} \mathcal{S} : E^* \rightarrow \tilde{H} \quad \text{and} \quad \tilde{\mathcal{S}}^* = \mathcal{S}^* \mathcal{U}^{-1} : \tilde{H} \rightarrow E^*.$$

We will see in the next Section 3.2, that for Gaussian measures there is a natural choice of H . However, the factorization offers additional flexibility in terms of calculation, see also Section 3.3.

3.2 Gaussian Measures

The Gaussianity of a random element on Banach space E is characterised by the continuous functionals in E^* . Due to Remark 3.6, it is sufficient to consider measures on \mathcal{B} .

Definition 3.9 (Gaussian measure). A probability measure γ on (E, \mathcal{B}) is called *Gaussian* if, for any $f \in E^*$, the induced measure $\gamma \circ f^{-1}$ is a Gaussian distribution on \mathbb{R} . A random element X in E is *Gaussian* if its law is Gaussian.

Directly from the definition, the quantity $\mathbb{E} \|X\|^p$ is finite for $p = 1, 2$. Consequently $m_2(X)$ also exists. Because of Lemma 3.7, for X , its mean a_γ of X is an element of E and the mapping $a_\gamma : E^* \rightarrow E$ is bounded. If $a_\gamma = 0$, we say it is a *centred* (i.e. *zero mean*) Gaussian measure. In addition, the covariance operator is given by a symmetric positive nuclear operator $\mathcal{C}_\gamma : E^* \rightarrow E$.

As an analogy to finite-dimensional Gaussian distributions, one may conjecture that the Gaussian measures on vector spaces are as well uniquely determined by its mean and covariance. It turns out to be true due to the lemma below, taken from Proposition 2.8, Chapter IV, [96].

¹To be precise, it is positive semi-definite or nonnegative definite. However, we use the term ‘positive’ for brevity, and strictly positive if $\langle f, \mathcal{C}_\mu f \rangle > 0$ for any $f \neq 0$.

Lemma 3.10 (Fourier transform and uniqueness of Gaussian measures). *For a Gaussian measure γ on E with mean a_γ and covariance operator \mathcal{C}_γ , its Fourier transform is given by*

$$\widehat{\gamma}(f) = \exp(i f(a_\gamma) - \frac{1}{2} \langle f, \mathcal{C}_\gamma f \rangle).$$

As shown in Lemma 3.10, the Fourier transform of Gaussian measures only involves the pair $(a_\gamma, \mathcal{C}_\gamma)$. Then by Lemma 3.4, Gaussian measures are completely determined by its mean and covariance. Henceforth, we will use the notation $\mathcal{N}_E(a_\gamma, \mathcal{C}_\gamma)$ to denote a Gaussian measure on E with mean a_γ and covariance operator \mathcal{C}_γ . In addition, the following corollary is an immediate consequence of Lemma 3.10.

Before starting investigating fine properties of the covariance structure, we collect two celebrated inequalities demonstrating the fundamental properties of Gaussian measures, i.e. isoperimetric inequality and exponential tail property. For the detailed discussion, we refer to Section 3.1 in [67].

Lemma 3.11. *Let the law of X be a Gaussian measure γ on a separable Banach space E .*

(i) *Borell's inequality.* With $\sigma = m_2(X)$ given in (3.2), we have

$$\mathbb{P}\left(\left|\|X\| - \mathbb{E}\|X\|\right| > t\right) \leq 2 \exp\left\{-\frac{t^2}{2\sigma^2}\right\}.$$

(ii) *Exponential tail (Fernique).* There exists a positive constant α such that

$$\int_E e^{\alpha\|x\|^2} \gamma(dx) < \infty.$$

3.2.1 Covariance Structure

We are going to look closely on the covariance structure of a Gaussian measure $\mathcal{N}_E(a_\gamma, \mathcal{C}_\gamma)$. The following two subspaces of E^* and E respectively play an important role in shaping the Gaussian covariance structure.

From the definition, $\|f\|_{L^2(\gamma)}$ is finite for all $f \in E^*$. Hence we can introduce the following definitions.

Definition 3.12 (Reproducing kernel Hilbert space & Cameron-Martin space). For a Gaussian measure γ on a separable Banach space E , define an embedding

$$\mathcal{J} : E^* \rightarrow L^2(\gamma), \quad f \mapsto f - a_\gamma f, \quad (3.3)$$

which is continuous because of the finite moments.

(i) The *reproducing kernel Hilbert space* (RKHS) of the measure γ is defined as the closure of $\mathcal{J}(E^*)$ in $L^2(\gamma)$, denoted by E_γ^* , which self is a Hilbert space equipped with $L^2(\gamma)$ inner product.

(ii) The *Cameron-Martin space* is a subspace in E defined as

$$\mathbb{H}_\gamma := \left\{ h \in E : |h|_{\mathbb{H}_\gamma} < \infty \right\},$$

where

$$|h|_{\mathbb{H}_\gamma} := \sup \left\{ f(h) : f \in E^*, \|\mathcal{J}(f)\|_{L^2(\gamma)} \leq 1 \right\}. \quad (3.4)$$

Remark 3.13 (Gaussian covariance operator). Consider the following extension of the covariance operator \mathcal{C}_γ to E_γ^* ,

$$\begin{aligned} \mathcal{K}_\gamma : E_\gamma^* &\rightarrow E \subset (E^*)', \\ (\mathcal{K}_\gamma f)(g) &= \int_E f(x)[g(x) - a_\gamma g] \gamma(dx), \quad f \in E_\gamma^*, g \in E^*. \end{aligned} \quad (3.5)$$

Notice that, if γ is centred, \mathcal{K}_γ is just the continuous extension to E_γ^* by bounded linear transform theorem (Theorem I.7, [81]). In the case of general Gaussian measures, for any $f \in E^*$, $\mathcal{K}_\gamma f$ coincides with $\mathcal{C}_\gamma(f - a_\gamma f)$ (the original one defined in (3.1)), where $f - a_\gamma f \in E_\gamma^*$, as elements in $(E^*)'$, although f is not necessarily in E_γ^* if $a_\gamma f \neq 0$. From now on, for a Gaussian measure, γ the operator \mathcal{K}_γ is understood as in the extension of \mathcal{C}_γ as in (3.5).

Notice that the elements in E_γ^* are centred by the mapping \mathcal{J} . The following lemma shows that in fact they are also Gaussian.

Lemma 3.14. *Let the law of X be a Gaussian measure γ . For arbitrary $f \in E_\gamma^*$, $f(X)$ is a centred Gaussian random variable with variance*

$$\mathbb{E}(f(X))^2 = \int_E f^2 d\gamma = \|f\|_{L^2(\gamma)}^2.$$

Proof. It follows from the fact that the limit of a Gaussian sequence converging in probability is Gaussian. For the detail, see Lemma 2.2.8, [8]. \square

The usage of the term RKHS is not consistent in literature: it may refer to both E_γ^* and \mathbb{H}_γ . But this does not really cause any troubles in practice, as in fact E_γ^* and \mathbb{H}_γ are almost mutually replaceable, as shown in the Lemma 3.15 below.

Lemma 3.15. *Let γ be a Gaussian measure on a separable Banach space E . An element $h \in E$ belongs to \mathbb{H}_γ if and only if $h = \mathcal{K}_\gamma \widehat{h}$ for some functional $\widehat{h} \in E_\gamma^*$. In addition,*

$$|h|_{\mathbb{H}_\gamma} = \|\widehat{h}\|_{L^2(\gamma)}. \quad (3.6)$$

Consequently, $\mathcal{K}_\gamma : E_\gamma^* \rightarrow \mathbb{H}_\gamma$ is an isometric isomorphism and $\mathbb{H}_\gamma = \mathcal{K}_\gamma(E_\gamma^*)$ is a Hilbert space equipped with the inner product

$$\langle h_1, h_2 \rangle_{\mathbb{H}_\gamma} := \langle \widehat{h}_1, \widehat{h}_2 \rangle_{L^2(\gamma)}, \quad (3.7)$$

where $\mathcal{K}_\gamma \widehat{h}_i = h_i$, for $i = 1, 2$.

Proof. suppose $|h|_{\mathbb{H}_\gamma} < \infty$. Define a linear mapping

$$\Phi_h : E^* \rightarrow \mathbb{R}, \quad f \mapsto [\mathcal{J}f](h).$$

It is well-defined, since $|\Phi_h f| \leq \|\mathcal{J}f\|_{L^2(\gamma)} |h|_{\mathbb{H}_\gamma}$ from the definition (3.4), which also implies that Φ_h is continuous with respect to $L^2(\gamma)$ norm. Then by BLT theorem (Lemma A.4), there is a continuous extension of Φ_h to E_γ^* , denoted with the same symbol. Consequently, by Riesz representation theorem, there exists an element \hat{h} in E_γ^* such that $\Phi_h f = \langle f, \hat{h} \rangle_{L^2(\gamma)}$, for all $f \in E_\gamma^*$. In particular, for $f \in E^*$,

$$f(h) = \Phi_h f = \int_E (\mathcal{J}f) \hat{h} d\gamma = f(\mathcal{K}_\gamma \hat{h}).$$

Using the double dual norm (Lemma A.5), we conclude $h = \mathcal{K}_\gamma \hat{h}$, and

$$|h|_{\mathbb{H}_\gamma} = \sup\{f(h) : f \in E^*, \|\mathcal{J}f\|_{L^2(\gamma)} \leq 1\} = \|\hat{h}\|_{L^2(\gamma)}. \quad (3.8)$$

Conversely, assume $h = \mathcal{K}_\gamma \hat{h}$ for some $\hat{h} \in E_\gamma^*$. Then, for all $f \in E^*$,

$$|f(h)| = |\langle \mathcal{J}f, \hat{h} \rangle_{L^2(\gamma)}| \leq \|\mathcal{J}(f)\|_{L^2(\gamma)} \|\hat{h}\|_{L^2(\gamma)},$$

which implies $|h|_{\mathbb{H}_\gamma} < \infty$.

The isometry (3.6) follows from (3.8). Since isometry implies injectivity, the claim that \mathbb{H}_γ equipped with (3.7) is a Hilbert space follows from Lemma A.11. \square

As its name implies, the Cameron-Martin space is indeed a Hilbert space. Due to Lemma 3.15, we use the term RKHS for both \mathbb{H}_γ and E_γ^* . In most situations, it does not lead to confusions. Otherwise, we will explicitly specify whether RKHS is a subspace of E or E_γ^* .

Known from Lemma 3.8, the covariance operator can be factorized with an operator mapping E^* to a Hilbert space. Again due to Lemma 3.15, a natural choice E_γ^* surfaces for the Gaussian measure γ .

Corollary 3.16 (Canonical factorization of Gaussian covariance operators). *For a Gaussian measure γ on E , its covariance operator \mathcal{C}_γ admits the following factorization $\mathcal{C}_\gamma = \mathcal{J}^* \mathcal{J}$, where \mathcal{J} is the embedding of E^* to E_γ^* defined in (3.3), and $\mathcal{J}^* = \mathcal{K}_\gamma$, i.e. the extension of \mathcal{C}_γ to E_γ^* given in (3.5).*

The RKHS contains large information on the Gaussian measure. We present several results with different flavours below. The proofs can be found in Section 2.4 and 3.2, [8].

Seemingly contradictory, the size of RKHS \mathbb{H}_γ can be both small and big at the same time.

Proposition 3.17. *Let γ be a Gaussian measure on E . If $\dim E = \infty$, then \mathbb{H}_γ only has zero measure, i.e. $\gamma(\mathbb{H}_\gamma) = 0$. On the other hand, the topological support of γ is the closure $\overline{\mathbb{H}_\gamma}$ of \mathbb{H}_γ in E . In particular, if γ is a nondegenerate Gaussian measure, i.e. the topological support of γ is the whole space E , \mathbb{H}_γ is dense in E . Furthermore, the unit ball $U(\mathbb{H}_\gamma) = \{h \in \mathbb{H}_\gamma : |h|_{\mathbb{H}_\gamma} \leq 1\}$ is compact in E , which implies the embedding \mathbb{H}_γ*

Gaussian measures are not so relevant to the underlying space E , but rather determined by its RKHS, shown as follows.

Proposition 3.18. *Let γ be a Gaussian measure on E . If E is continuously embedded another Banach space \tilde{E} , If γ is considered to be a Gaussian measure on \tilde{E} , the Cameron-Martin space is still \mathbb{H}_γ . Let μ be another Gaussian measure on E . If $\mathbb{H}_\gamma = \mathbb{H}_\mu$ and $|h|_{\mathbb{H}_\gamma} = |h|_{\mathbb{H}_\mu}$ for all $h \in \mathbb{H}_\gamma$, then $\gamma = \mu$.*

The absolute continuity of measures is especially important for statistical study, since it implies a density, which is often desirable for the construction of statistics such as likelihood ratios, etc. For Gaussian measures, this is closely related to RKHS. Consider a Gaussian random element X with distribution γ . The distribution of the vector $X + h$ with $h \in E$ is defined by

$$\gamma_h(B) = \gamma(B - h), \quad B \in \mathcal{B}(E).$$

We call γ_h the *shift* of γ in the direction h . A shift is admissible if $\gamma_h \ll \gamma$. The proposition below shows that \mathbb{H}_γ is exactly the set of admissible shifts.

Proposition 3.19 (Cameron-Martin). *Let γ be a Gaussian measure on E and \mathbb{H}_γ be its RKHS. Consider the shifted measure $\gamma_h = \gamma(\cdot - h)$.*

If $h \in \mathbb{H}_\gamma$, then $\gamma \sim \gamma_h$, and the Radon-Nikodym density is given by

$$\frac{d\gamma_h}{d\gamma}(x) = \exp\left(\widehat{h}(x) - \frac{1}{2}|h|_{\mathbb{H}_\gamma}^2\right), \quad (3.9)$$

where $\mathcal{K}_\gamma \widehat{h} = h$; otherwise $\gamma \perp \gamma_h$.

The density (3.9) is also known as *Cameron-Martin formula*. Combining the above Proposition 3.19 with Lemma 3.10, we conclude that if $h \in \mathbb{H}_\gamma$, then

$$a_{\gamma_h} = a_\gamma + h.$$

Kullback-Leibler(KL) divergence is an important quantity used to quantify the discrepancy between probability measures, which is defined as

$$D_{KL}(\mathbb{P}; \mathbb{Q}) := \int_E \log \frac{d\mathbb{P}}{d\mathbb{Q}} d\mathbb{P}, \quad (3.10)$$

where \mathbb{P} and \mathbb{Q} are two probability measures on E . KL divergence is of great importance to Bayesian nonparametric inference, as it is used to characterise a set around the truth, whose probability mass has a direct implication on the contraction rate of the posterior distribution (see Section 4.3). For Gaussian measures, it can be calculated using RKHS norm, as shown in the lemma below.

Lemma 3.20 (Kullback-Leibler divergence of Gaussian measures). *Let γ be a centred Gaussian measure on E , and h, g be two elements from \mathbb{H}_γ . Consider the shifted measures $\gamma_h = \gamma(\cdot - h)$ and $\gamma_g = \gamma(\cdot - g)$. The KL divergence between γ_h and γ_g is given by*

$$D_{KL}(\gamma_h; \gamma_g) = \frac{1}{2}|h - g|_{\mathbb{H}_\gamma}^2.$$

Proof. From Lemma 3.15, $\mathcal{K}_\gamma : E_\gamma^* \rightarrow \mathbb{H}_\gamma$ is an isometric isomorphism. Because γ is centred, i.e. $a_\gamma = 0$, in particular we have

$$|f|_{\mathbb{H}_\gamma}^2 = \|\widehat{f}\|_{L^2(\gamma)}^2 = \langle \widehat{f}, \mathcal{K}_\gamma \widehat{f} \rangle_{E_\gamma^* \times \mathbb{H}_\gamma}.$$

Hence we have $h = \mathcal{K}_\gamma \widehat{h}$ and $g = \mathcal{K}_\gamma \widehat{g}$. Then it follows from Proposition 3.19,

$$\frac{d\gamma_h}{d\gamma_g}(x) = \frac{d\gamma_h/d\gamma}{d\gamma_g/d\gamma}(x) = \exp\left(\widehat{h}(x) - \frac{1}{2}|h|_{\mathbb{H}_\gamma}^2 - \widehat{g}(x) + \frac{1}{2}|g|_{\mathbb{H}_\gamma}^2\right).$$

Consequently,

$$\begin{aligned} D_{KL}(\gamma_h; \gamma_g) &= \int_E \log \frac{d\gamma_h}{d\gamma_g} d\gamma_h = \int_E \widehat{h}(x) - \widehat{g}(x) d\gamma_h - \frac{1}{2}|h|_{\mathbb{H}_\gamma}^2 + \frac{1}{2}|g|_{\mathbb{H}_\gamma}^2 \\ &= \mathbb{E} \widehat{h} - \mathbb{E} \widehat{g} - \frac{1}{2} \langle \widehat{h}, \widehat{h} \rangle_{L^2(\gamma)} + \frac{1}{2} \langle \widehat{g}, \widehat{g} \rangle_{L^2(\gamma)} \\ &= \frac{1}{2} \left[\langle \widehat{h}, \mathcal{K}_\gamma \widehat{h} \rangle_{E_\gamma^* \times \mathbb{H}_\gamma} - \langle \widehat{g}, \mathcal{K}_\gamma \widehat{g} \rangle_{E_\gamma^* \times \mathbb{H}_\gamma} + \langle \widehat{g}, \mathcal{K}_\gamma \widehat{g} \rangle_{L^2(\gamma)} \right] \\ &= \frac{1}{2} \langle \widehat{h} - \widehat{g}, \mathcal{K}_\gamma (\widehat{h} - \widehat{g}) \rangle_{E_\gamma^* \times \mathbb{H}_\gamma} = \frac{1}{2} |h - g|_{\mathbb{H}_\gamma}^2, \end{aligned}$$

where we use $\int_E \widehat{h}(x) d\gamma_h = \mathbb{E} \widehat{h} = \widehat{h}(h) = \langle \widehat{h}, \mathcal{K}_\gamma \widehat{h} \rangle_{E_\gamma^* \times \mathbb{H}_\gamma}$ directly from the definition. \square

Remark 3.21. All the results stated in Section 3.2 can be generalised to the *Radon* Gaussian measures on a locally convex topological space E . See Chapter 3 in [8] for the details. In particular, we mention that Proposition 3.19 and Lemma 3.20 remain valid with E being a locally convex topological space, e.g $\mathbb{R}^{\mathbb{N}}$.

3.2.2 Examples

We now present several examples of Gaussian elements.

Example 3.22 (Standard Gaussian measure on $\mathbb{R}^{\mathbb{N}}$). Let γ_n be the standard Gaussian measures on \mathbb{R} . Then the measure

$$\gamma = \bigotimes_{n \in \mathbb{N}} \gamma_n$$

is centred Gaussian on $E = \mathbb{R}^{\mathbb{N}}$. Furthermore, $E_\gamma^* \simeq \ell^2$ and $\mathbb{H}_\gamma = \ell^2$. For the proof, see Example 2.3.5 in [8].

Although at the first glance, Example 3.22 seems to be quite specific, it is the ‘only’ centred Radon Gaussian measure on locally convex spaces (in particular Banach spaces) in the sense that there always exists an isomorphism between it and another centred Radon Gaussian measure on a Banach space (see Theorem 3.4.4, [8]).

Example 3.23 (Gaussian measures on Hilbert spaces). Let $(H, \langle \cdot, \cdot \rangle)$ be a separable Hilbert space and we identify H^* with H . For simplicity, we consider the zero mean case. Let \mathcal{C} be a covariance operator. Then by spectral theorem of self-adjoint compact operators, there exists an orthonormal basis φ_k diagonalising \mathcal{C} , i.e. $\mathcal{C}\varphi_k = \kappa_k\varphi_k$. In particular, $\sum_k \kappa_k < \infty$, since \mathcal{C} is nuclear. Due to the diagonalisation, a Gaussian element X in H with distribution $\mathcal{N}_H(0, \mathcal{C})$, admits the representation

$$X = \sum_k \sqrt{\kappa_k} \varphi_k \xi_k,$$

where $\{\xi_k\}$ is a sequence of i.i.d. $\mathcal{N}_{\mathbb{R}}(0, 1)$ random variables. This representation is known as *Karhunen-Loève expansion*. Conversely, the representation also defines a Gaussian measure on H .

Now we identify the RKHS \mathbb{H} of γ using the canonical factorisation of \mathcal{C} (see Corollary 3.16). The closure of $\mathcal{J}(H^*)$ follows from explicit calculations,

$$\mathcal{J}(H^*) = H_\gamma^* = \left\{ f : f(h) = \langle f, h \rangle_H, \|\mathcal{J}(f)\|_{L^2(\gamma)}^2 := \sum_k (\sqrt{\kappa_k} f_k)^2 < \infty \right\}.$$

Then, using Lemma 3.15, we obtain

$$\mathbb{H}_\gamma = \left\{ h \in H : |h|_{\mathbb{H}}^2 = \sum_k \left(\frac{h_k}{\sqrt{\kappa_k}} \right)^2 < \infty \right\},$$

by applying change of variables.

Remark 3.24. A Gaussian measure on Hilbert space can be connected to a Hilbert scale (see Section 2.2). Using the notations from the previous example, define $\Lambda = \mathcal{C}^{-1/2}$ using spectral theorem on the domain $\text{Dom } \Lambda = \mathbb{H}_\gamma$, which is dense in X . Let $\{H_s\}_{s \in \mathbb{R}}$ be the Hilbert scale generated by Λ with $H_0 = H$. Then, we have $H_1 = \mathbb{H}_\gamma$ and $H_{-1} = H_\gamma^*$. Furthermore, $\Lambda^{-2} : H_{-1} \rightarrow H_1$ is an isometric isomorphism. In fact, Λ^{-2} is the extension \mathcal{K}_γ of covariance operator \mathcal{C} , c.f. Lemma 3.15.

Next we consider the Gaussian measure induced by Gaussian processes. Recall that a random process is a family $\{X_t\}_{t \in \mathfrak{T}}$ of random elements indexed a parameter set \mathfrak{T} , defined on a common probability space (Ω, \mathbb{P}) . A process is *Gaussian* if $(X_{t_1}, \dots, X_{t_n})$ is a \mathbb{R}^n -valued Gaussian random element, for any $\{t_1, \dots, t_n\} \subset \mathfrak{T}$. Similar to Gaussian measures, the property of Gaussian processes are completely determined by the expectation $\mathbb{E} X_t, t \in \mathfrak{T}$ and covariance $\text{Cov}(X_s, X_t), s, t \in \mathfrak{T}$.

Example 3.25 (Classical Wiener measure). Wiener process W_t on $[0, 1]$ is characterised by

$$\mathbb{E} W_t = 0, \quad \text{Cov}(W_s, W_t) = s \wedge t,$$

for $s, t \in [0, 1]$. It is well know that W_t is almost surely in $E = C[0, 1]$, the Banach space of continuous functions on $[0, 1]$, equipped with supremum norm

$\|f\|_E = \sup_{0 \leq x \leq 1} |f(x)|$. The dual space E^* is the space of sign measures of finite variation on $[0, 1]$, with duality given by

$$\langle \mu, f \rangle_{E^* \times E} := \int_{[0,1]} f d\mu.$$

Consider $W = W_{(\cdot)}$ as a random element in E . It is obvious $\mathbb{E}W = 0$. In addition,

$$\begin{aligned} \mathbb{Cov}(\langle \mu, W \rangle, \langle \nu, W \rangle) &= \mathbb{E} \int_{[0,1]^2} W(s)W(t) \mu(ds)\nu(dt) \\ &= \int_{[0,1]^2} \mathbb{E} W(s)W(t) \mu(ds)\nu(dt) = \int_{[0,1]^2} s \wedge t \mu(ds)\nu(dt) = \langle \mu, \mathcal{K}\nu \rangle, \end{aligned}$$

where

$$(\mathcal{K}\nu)(s) = \int_{[0,1]} s \wedge t \nu(dt).$$

It can be shown that the operator \mathcal{K} is indeed a covariance operator for the centred Gaussian measure on $E = C[0, 1]$ induced by Wiener process on $[0, 1]$.

The result in Example 3.25 can be extended to general Gaussian processes with continuous paths as follows. Let $E = C(\mathfrak{X}; \mathbb{R})$, where \mathfrak{X} is a compact metric space, equipped with the supremum norm. Then the dual space E^* is again given by

$$\langle \varphi_\mu, f \rangle_{E^* \times E} := \int_{\mathfrak{X}} f d\mu,$$

where μ is a signed measure with finite variation $|\mu|(E) = \|\varphi_\mu\|$ (see Chapter 6, [83]). For a Gaussian process X_t with $\mathbb{E}X_t = a(t)$ and $\mathbb{Cov}(X_s, X_t) = k(s, t)$, the distribution of $X = X_{(\cdot)}$ is a Gaussian measure on E with expectation a and covariance operator

$$(\mathcal{K}\nu)(s) = \int_{\mathfrak{X}} k(s, t) \nu(dt).$$

Conversely, every Gaussian measure on E induces a continuous Gaussian process on its dual space E^* .

Example 3.26 (Gaussian process induced by Gaussian measures). A Gaussian measure γ on a separable Banach space E always induces a continuous Gaussian process on the dual space E^* . Since E is separable, the unit ball $U(E^*)$ of its dual is metrizable for the weak* topology on $\sigma(E^*, E)$ (where E is considered as a subset of the double dual E^{**} of E). Denote by (z_k) a weak dense sequence in $U(E^*)$. For x in E , we have $\|x\| = \sup_k |z_k(x)|$. Consider each f in $U(E^*)$ as a Gaussian random variable on (E, γ) . This defines a weak* continuous Gaussian process, since for each x in E , the map $f \mapsto f(x)$ is weak* continuous.

As usual, the situation becomes more transparent, when the underlying space is a Hilbert space.

Example 3.27 (Gaussian process indexed by Hilbert space). Let X be a centred Gaussian random element on H with covariance operator \mathcal{C} . Denote the functional $\langle h, \cdot \rangle_H$ by ℓ_h . The covariance

$$\text{Cov } \ell_h(X) \ell_g(X) = \langle h, \mathcal{C}g \rangle_H = K(h, g).$$

Hence, $\{W(h) := \langle h, X \rangle_H, h \in H\}$ is a Gaussian process indexed by H with covariance function $K(h, g) = \langle h, \mathcal{C}g \rangle_H$. From the previous example, the process $\{W(h)\}_{h \in H}$ is a (weak-)continuous Gaussian process on H .

3.3 Radonification of Cylindrical Measures

In this section we continue the study of the relations between Gaussian measures and processes. We are interested in the following Gaussian process, which is of great importance in statistics.

Definition 3.28 (Isonormal process). An *isonormal* process (also known *white noise* process) $\{W(h) : h \in H\}$ indexed by a Hilbert space H is a stochastic process such that:

- (i) $W(h_1), \dots, W(h_n)$ are jointly centred Gaussian for all $h_1, \dots, h_n \in H$;
- (ii) $\text{Cov}(W(h), W(g)) = \langle h, g \rangle_H$.

Analogous to Example 3.23, the isonormal process on H can be constructed using series representation. Let $\{\varphi_k\}$ be an orthonormal basis on H and $\{\xi_k\}$ be a sequence of i.i.d. $\mathcal{N}_{\mathbb{R}}(0, 1)$ random variables. Then the process

$$W(h) := \sum_k \xi_k \langle \varphi_k, h \rangle \tag{3.11}$$

is isonormal. However, it does not induce a genuine Gaussian measure as the covariance operator $\mathcal{C} = \text{id}$ is not nuclear if $\dim H = \infty$.

We are going to demonstrate that the ‘measure’ induced by an isonormal process can be extended to a larger space such that it becomes an authentic measure, which is known as *Radonification*. To show this, we first introduce the following concepts.

The *cylindrical* sets in a Banach space E with its dual E^* has the following form

$$C = \{x \in E : (f_1(x), \dots, f_n(x)) \in C_0\}, \quad f_i \in E^*,$$

where $C_0 \in \mathcal{B}(\mathbb{R}^n)$ is a base of C . A nonnegative additive function ν is a *cylindrical measure* on E if for every continuous linear operator $\mathcal{P} : E \rightarrow \mathbb{R}^n$, the set function

$$\nu \circ \mathcal{P}^{-1} : B \mapsto \nu(\mathcal{P}^{-1}(B))$$

on the algebra $\mathcal{A}(E)$ of cylindrical sets is countably additive. Fourier transform of a cylindrical measure ν is defined similarly to the one of measures,

$$\widehat{\nu}(f) = \int_{\mathbb{R}} e^{it} \nu \circ f^{-1}(dt), \quad f \in E^*.$$

Now we see that an isonormal process from Definition 3.28 induces a cylindrical Gaussian measure.

Example 3.29 (Cylindrical Gaussian measure on Hilbert space). Let H be an infinite dimensional separable Hilbert space. For all cylindrical sets C in the form $C = P^{-1}(C_0)$, where P is an orthogonal projection onto a n -dimensional subspace E_n of E , define

$$\nu(C) = \gamma_n(C_0),$$

where γ_n is the standard Gaussian measure on E_n (identified with \mathbb{R}^n). The set function ν is called the *canonical cylindrical Gaussian measure* on E , whose Fourier transform is

$$\widehat{\nu}(f) = e^{-\frac{1}{2}\|f\|_H^2}.$$

It is straightforward from the definition that an isonormal process induces ν on E .

We recall a general result on radonification of Gaussian measures, known as Sazonov's theorem (see Theorem 3.4, [13]).

Theorem 3.30 (Sazonov). *Let H, G be two Hilbert space and ν be the canonical cylindrical Gaussian measure on H . A bounded operator $\mathcal{T} : H \rightarrow G$ is γ -radonifying if the push forward $\nu \circ \mathcal{T}^{-1}$ on G is a genuine Gaussian measure. We have, $\mathcal{T} : H \rightarrow G$ is γ -radonifying if and only if \mathcal{T} is a Hilbert-Schmidt operator.*

We outline one application of Theorem 3.30 to diagonalisable operators. The main motivation for this is to cover the white noise. A bounded operator $\mathcal{Q} : H \rightarrow H$ is called *diagonalisable*, if with an orthonormal basis $\{\varphi_k\}_{k \in \mathbb{N}}$ for H , for all $f \in H$,

$$\mathcal{Q}f = \sum_{k \in \mathbb{N}} q_k f_k \varphi_k \quad \text{with } f_k = \langle f, \varphi_k \rangle_H.$$

Let $\mathcal{Q} : H \rightarrow H$ be a self-adjoint, positive definite, and diagonalisable operator on H with a basis $\{\varphi_k\}_{k \in \mathbb{N}}$. If \mathcal{Q} is of trace class, i.e.

$$\sum_{k \in \mathbb{N}} \langle \mathcal{Q}\varphi_k, \varphi_k \rangle_H = \sum_{k \in \mathbb{N}} \|\mathcal{Q}^{1/2}\varphi_k\|_H^2 < \infty,$$

it is the covariance of a centred Gaussian measure. On the other hand, if \mathcal{Q} is not of trace class, i.e. $\sum_{k \in \mathbb{N}} q_k = \infty$, it cannot be the covariance of a Gaussian distribution on H . However, as in Example 3.27, it defines a centred Gaussian process $\{W(h) : h \in H\}$ with covariance function

$$\text{Cov}(W(h), W(g)) = \langle h, \mathcal{Q}g \rangle_H.$$

Denote

$$\mathbb{H} = \mathcal{Q}^{1/2}(H) = \left\{ f = \sum_k f_k \varphi_k : |f|_{\mathbb{H}} := \|\mathcal{Q}^{-1/2}f\|_H < \infty \right\},$$

which is a Hilbert space with the inner product $\langle \cdot, \cdot \rangle_{\mathbb{H}} = \langle \mathcal{Q}^{-1/2} \cdot, \mathcal{Q}^{-1/2} \cdot \rangle_H$ (by Lemma A.11). By restricting the previous process to \mathbb{H} , the new process $\{W(h) :$

$h \in \mathbb{H}$ is isonormal on \mathbb{H} , whose distribution is equivalent to the cylindrical measure on \mathbb{H} .

Notice that $\{e_k = \mathcal{Q}^{1/2}\varphi_k : k \in \mathbb{N}\}$ forms a complete orthonormal basis of \mathbb{H} . Introduce an inner product $\langle \cdot, \cdot \rangle_X$ on \mathbb{H} such that,

$$\langle e_j, e_k \rangle_X := w_j w_k \delta_{jk},$$

where the weights $\{w_k\}_{k \in \mathbb{N}}$ satisfy $\sum_{k \in \mathbb{N}} w_k^2 < \infty$ and δ_{jk} is the Kronecker delta. For example, any sequence of weights that converges to zero faster than a square root rate, i.e. $w_k \ll k^{-1/2}$ as k goes to infinity, satisfies the previous condition of summability. Let X be the closure of \mathbb{H} under the norm introduced by the inner product above. Then the embedding $\iota : \mathbb{H} \hookrightarrow X$ is Hilbert-Schmidt. Furthermore, because \mathcal{Q} is diagonalised with $\{\varphi_k\}$, it is easy to examine that $H \subset X$. Consequently, by Theorem 3.30, there is a centred Gaussian measure with the covariance operator $\iota \mathcal{Q} \iota^*$ on X , with the Cameron-Martin space \mathbb{H} . In particular, the technique discussed above is applicable to the white noise. Moreover, since the identity operator is diagonalisable with any orthonormal basis, we have the freedom to choose the most convenient one to work with.

Example 3.31 (White noise process). Let $H = L^2(\mathfrak{D})$ with a bounded domain $\mathfrak{D} \subset \mathbb{R}^d$. Let $\xi : \mathfrak{D} \rightarrow \mathbb{R}$ be a Gaussian process such that $\mathbb{E} \xi(s) \xi(t) = \delta(s - t)$, where δ is the Dirac delta function. Let $\langle \xi, u \rangle := \int_{\mathfrak{D}} \xi(s) u(s) ds$, where $\xi(s) ds$ is understood in the distributional sense. We have

$$\mathbb{E} \langle \xi, u \rangle \langle \xi, v \rangle = \int_{\mathfrak{D}} \int_{\mathfrak{D}} u(s) v(t) \mathbb{E} [\xi(s) \xi(t)] ds dt = \int_{\mathfrak{D}} u(t) v(t) dt = \langle u, v \rangle_{L^2},$$

for $u, v \in L^2(\mathfrak{D})$. Therefore, the Gaussian process ξ is isonormal on $L^2(\mathfrak{D})$. By the results in this section, while ξ is not a L^2 -valued Gaussian process, it can be realised in a space G such that the embedding of $L^2(\mathfrak{D})$ to G is Hilbert-Schmidt, and the stochastic integral $\int_{\mathfrak{D}} \xi(t) f(t) dt$ is well-defined. In particular,

$$\int_{\mathfrak{D}} \xi(s) \mathbb{1}_{\{s_i \leq t_i : 1 \leq i \leq d\}} ds_1 \cdots ds_d = \prod_{i \leq d} t_i$$

and

$$\mathbb{E} \langle \xi, \mathbb{1}_{\{r \leq s\}} \rangle \langle \xi, \mathbb{1}_{\{r \leq t\}} \rangle = \prod_{i \leq d} s_i \wedge t_i.$$

Consequently, $\int_0^t \xi(s) ds$ is the Brownian sheet, which is often stated as white noise is the (distributional) derivative of Brownian motion. See also Section 4.1.5 in [23].

3.4 Notes

For a systematic (and more general) treatment on Gaussian measures, we refer to [8]. Radonification of a cylindrical Gaussian measure can also be stated in the framework of nuclear spaces, in which Minlos' theorem is used as the major tool, see [47]. The fields of white noise analysis, generalised stochastic processes were originated from this approach. A closely related object is the abstract Wiener space, see [74] for details.

Chapter 4

Bayesian Nonparametrics for Gaussian Linear Models

In this chapter, we survey some basic facts of Bayesian nonparametrics, retrieved from [35]. In addition, we present a general contraction results for the Gaussian linear model with a transformed drift.

4.1 Bayesian Nonparametrics

4.1.1 Bayes' Rule

As mentioned in Section 1.3, Bayes' rule on infinite-dimensional spaces requires special care on measurability concerns. In this section, we sketch Bayes' rule in infinite-dimensional spaces. For more details, see Section 1.3 in [35].

The Bayesian procedure can be described in the following steps. First, the statistician seeks for a *prior distribution* Π to apply on the parameter space (Θ, \mathcal{S}) , where \mathcal{S} is a σ -algebra. Then, given θ , each element \mathbb{P}_θ in the statistical model $\{\mathbb{P}_\theta : \theta \in \Theta\}$ becomes a regular conditional distribution on the sample space $(\mathbb{X}, \mathcal{X})$. That means, the mapping $\Theta \times \mathcal{X} \ni (\theta, A) \mapsto \mathbb{P}_\theta(A) \in [0, 1]$ is a (*probability kernel*) from Θ to \mathbb{X} , i.e.

- (i) with fixed θ , $A \mapsto \mathbb{P}_\theta(A)$ is a probability measure on \mathbb{X} ,
- (ii) and with fixed A , $\theta \mapsto \mathbb{P}_\theta(A)$ is \mathcal{S} -measurable.

Therefore, the pair (X, θ) has a well-defined joint distribution

$$\mathbb{P}(X \in A, \theta \in B) = \int_B \mathbb{P}_\theta(A) d\Pi(\theta)$$

on the product space $(\mathbb{X} \times \Theta, \mathcal{X} \times \mathcal{S})$. Consequently, the (*Bayesian*) *marginal distribution* of X is given by

$$\mathbb{P}(A) = \int \mathbb{P}_\theta(A) d\Pi(\theta), \quad A \in \mathcal{X},$$

and the *posterior distribution*, the conditional distribution of θ given X , is

$$\Pi(B | X) = \mathbb{P}(\theta \in B | X), \quad B \in \mathcal{S}.$$

Remark 4.1. The posterior distribution is always well-defined, by Kolmogorov's definition. That means, $\Pi(\theta | X)$ is a measurable function of X such that, for a fixed $B \in \mathcal{S}$ and any $A \in \mathcal{X}$, $\mathbb{E}[\Pi(\theta | X) \mathbf{1}\{X \in A\}] = \mathbb{P}(X \in A, \theta \in B)$. However, to obtain a *regular* version of the conditional distribution, the size of the space (Θ, \mathcal{S}) cannot be too big. One sufficient condition is that Θ is a Polish space, i.e. a complete *separable* metric space, and \mathcal{S} is the Borel σ -algebra.

It is noteworthy that the posterior distribution $\Pi(\cdot | X)$ is unique up to null sets under the Bayesian marginal distribution. For a faithful Bayesian, this does not cause any problems, since it is believed that the Bayesian marginal distribution generates the data. However, in the case that the Bayesian framework is treated as an inference method, a 'true' distribution \mathbb{P}_0 generating the data X is not necessary identical to the Bayesian marginal distribution, and hence they do not have the same null sets. This phenomenon leads to indefiniteness of the posterior distribution. The pathological circumstance above is related to *misspecification* and can be excluded by the following condition,

$$\mathbb{P}_0 \ll \int \mathbb{P}_\theta \, d\Pi(\theta).$$

When the statistical model $\mathcal{E} = \{\mathbb{P}_\theta : \theta \in \Theta\}$ is dominated by a σ -finite measure μ , the Radon-Nikodym derivative provides densities p_θ related to the dominating measure μ such that $(x, \theta) \mapsto p_\theta(x)$ is measurable. Then, the density version of *Bayes' formula* is given by

$$\Pi(B | X) = \frac{\int_B p_\theta(X) \, d\Pi(\theta)}{\int p_\theta(X) \, d\Pi(\theta)}. \quad (4.1)$$

4.1.2 Bayesian Asymptotics

For a given Bayesian inferential procedure, i.e. estimating a parameter using the posterior distribution generated by Bayes' rule and a prior on the parameter space, its statistical performance can be evaluated in the asymptotic framework. Specifically, the following frequentist concepts are used to characterise posteriors.

For a sequence of experiments $\mathcal{E} = \{\mathbb{X}^{(n)}, \mathcal{X}^{(n)}, \mathbb{P}_\theta^{(n)} : \theta \in \Theta\}$, consider a prior Π on the Borel-algebra $\mathcal{B}(\Theta)$, and fix a version $\Pi_n(\cdot | X^{(n)})$ of its posterior distribution, i.e. any given choice of a regular condition distribution of θ given $X^{(n)}$ (see Remark 4.1), which is referred to as *the posterior*.

The asymptotic consistency for Bayesian procedures is defined as follows.

Definition 4.2 (Consistency). The posterior $\Pi_n(\cdot | X^{(n)})$ is (*weakly consistent*) at $\theta_0 \in \Theta$ if, for every neighbourhood U of θ_0 , in $\mathbb{P}_{\theta_0}^{(n)}$ probability,

$$\Pi_n(U^c | X^{(n)}) \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

and it is *strongly consistent* if the convergence above is $\mathbb{P}_{\theta_0}^{(n)}$ -almost surely.

The concept of contraction is useful for quantifying the speed that the posterior concentrates around the truth.

Definition 4.3 (Contraction rate). Suppose that the parameter space Θ is equipped with a pseudo-metric $d(\cdot, \cdot)$. A sequence ε_n is a *posterior contraction rate* at the parameter θ_0 if

$$\Pi_n\left(\theta : d(\theta, \theta_0) \geq M_n \varepsilon_n \mid X^{(n)}\right) \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

in $\mathbb{P}_{\theta_0}^{(n)}$, for arbitrarily slow $M_n \rightarrow \infty$.

Definition 4.3 requires some clarifications. First, the sequence ε_n is a rate rather than *the* rate, since any sequence decreasing slower also satisfies the definition. Certainly the fastest decaying sequence ε_n is of interest, but it may not be tractable in general. Therefore, ε_n is actually an *upper bound* for a targeted rate, but we will refer to it as *the* rate, with an abuse of terminology. In addition, due to the inaccessibility of minimal rates, we are satisfied with the result that the contraction rates are comparable (matching or close to) some benchmarks, for which a frequently used example is the optimal rate in minimax sense.

For statistical models in infinite-dimensional spaces, it is often the case that the unbounded sequence M_n can be replaced by a sufficiently large constant M , which gives a slightly stronger result. On the other hand, for parametric models, the unboundedness is required in order to obtain the regular rate $n^{-1/2}$. In this thesis, in all chapters we consider the strong version of Definition 4.3, i.e. M being a fixed constant, except in Chapter 7.

Some point estimators naturally emerge from the posterior distribution. An example is the centre of the smallest ball containing at least posterior mass $1/2$. Another one is the posterior mean. For further details, see Section 8.1 in [35].

In the Bayesian approach, credible sets are used to quantify uncertainty, whose counterpart in frequentist framework is confidence regions.

Definition 4.4 (Credible set). Fix a constant $1 - \gamma \in (0, 1)$. A subset $\mathcal{C}_n(X^{(n)})$ of Θ is a *credible set* for θ of credibility level $1 - \gamma$ if

$$\Pi_n\left(\theta \in \mathcal{C}_n(X^{(n)}) \mid X^{(n)}\right) \geq 1 - \gamma,$$

and the *frequentist coverage* of credible sets is given by

$$\mathbb{P}_\theta\left(\theta \in \mathcal{C}_n(X^{(n)})\right).$$

It is logical to consider the sets capturing most of the posterior mass while retaining a reasonable size at the same time. For example, if the posterior is unimodal, a ball centred at its mode can be considered as a credible set. Credibility is not the main focus of this thesis, but we will touch it in Chapter 7.

4.2 Gaussian Linear Models

In Section 1.2, we have already mentioned the following model,

$$X^{(n)} = \mathcal{A}^{(n)}\theta + \xi^{(n)}.$$

In this section, we formally define the model, using the elements introduced in Chapters 2 and 3.

There are two components in the model: the transformed signal $\mathcal{A}^{(n)}\theta$ and the noise $\xi^{(n)}$, and our interest is the original signal θ . Different combinations of $\mathcal{A}^{(n)}$ and $\xi^{(n)}$ leads to separate models. The choice of \mathcal{A} highly depends on the problem under investigation. In Part II, \mathcal{A} is a forward mapping with smoothing property, related to inverse problems, and in Part III, two items from the solution mapping of evolution equations are represented as \mathcal{A} . In the present chapter, we only assume that $\mathcal{A} : \Theta \rightarrow G$ is a bounded linear operator, from the parameter space Θ to another Hilbert space G , both of which are separable, and our focus is on the interpretation of noise. We are going to discuss two types of noise related to different observation schemes.

4.2.1 Continuous Observation

We consider the following Gaussian linear model with continuous observations. For $n \in \mathbb{N}$, let

$$\mathcal{A}^{(n)} = \mathcal{A} \quad \text{and} \quad \xi^{(n)} = \frac{1}{\sqrt{n}}\xi,$$

i.e.

$$X^{(n)} = \mathcal{A}\theta + \frac{1}{\sqrt{n}}\xi. \tag{4.2}$$

The forward operator is independent of the experiment sequence, while the noise $\xi^{(n)}$ is a fixed Gaussian element ξ scaled by $1/\sqrt{n}$, representing more information collected as $n \rightarrow \infty$. The observation is called *continuous* if we record the entire signal $\mathcal{A}\theta$ in G with noise ξ . Using the duality structure of Hilbert space G , the complete signal trajectory $\mathcal{A}\theta$ is equivalent to $(\langle \mathcal{A}\theta, g \rangle_G : g \in G)$. Concerning the noise, the conventional framework in statistics is Gaussian processes indexed by space G , which will be illustrated below. We conclude the discussion with a formula of the Kullback-Leibler distance between two observations with different signals $\mathcal{A}\theta$.

In our measurement model the observation $X^{(n)}$ will be a stochastic process $(X^{(n)}(w) : w \in G)$ such that

$$X^{(n)}(w) = \langle \mathcal{A}\theta, w \rangle_G + \frac{1}{\sqrt{n}}\xi(w), \quad w \in G, \tag{4.3}$$

where $\xi = (\xi(w) : w \in G)$ is a Gaussian process defined as follows. Let $\mathcal{Q} : G \rightarrow G$ be a bounded self-adjoint positive-definite operator. With the inner product $\langle \cdot, \cdot \rangle_G$ and norm $\|\cdot\|_G$ on G , we consider the noise ξ to be a Gaussian process indexed

by the Hilbert space G such that, for every $h, g \in G$, $\xi(h)$ is Gaussian with zero mean, i.e. $\mathbb{E} \xi(h) = 0$, and

$$\mathbb{E} \xi(h)\xi(g) = \langle h, \mathcal{Q}g \rangle_G. \quad (4.4)$$

By abuse of notation, we also call \mathcal{Q} the covariance operator of the Gaussian process ξ .

The processes $X^{(n)}$ and ξ are viewed as measurable maps in the *sample space* \mathbb{R}^G , with its product σ -field. Statistical sufficiency considerations show that the observation can also be reduced to the vector $(X^{(n)}(w_1), X^{(n)}(w_2), \dots)$, which takes values in the sample space $\mathbb{R}^{\mathbb{N}}$, for any orthonormal basis $(w_i)_{i \in \mathbb{N}}$ of G . The coordinates $X^{(n)}(w_i)$ of this vector are random variables with normal distributions such that

$$\mathbb{E} X^{(n)}(w_i) = \langle \mathcal{A}\theta, w_i \rangle_G, \quad \mathbb{E} \left(X^{(n)}(w_i) \right)^2 = \frac{1}{n} \|\mathcal{Q}^{1/2} w_i\|_G^2,$$

with covariance $\text{Cov}(X^{(n)}(w_i), X^{(n)}(w_j)) = \langle w_i, \mathcal{Q}w_j \rangle_G$ given by (4.4). If the basis $(w_i)_{i \in \mathbb{N}}$ is also orthogonal with respect to the inner product induced by \mathcal{Q} , i.e. $\langle w_i, \mathcal{Q}w_j \rangle_G = 0$ if $i \neq j$, then the variables $\xi(w_1), \xi(w_2), \dots$ are stochastically independent normal variables. In this case, $(X^{(n)}(w_1), X^{(n)}(w_2), \dots)$ is known as the *Gaussian sequence model* in statistics, albeit presently the ‘drift function’ $\mathcal{A}\theta$ involves the operator \mathcal{A} . See [11, 50] and references therein. In the rest of this thesis, we will consider the following cases.

- (i) In Part II, $\mathcal{Q} = \text{id}$ on G . ξ is an isonormal process on G and cannot be realised as a proper element of G . However, (4.3) makes perfect sense and (4.4) corresponds to $\mathbb{E} \xi(h)\xi(g) = \langle h, g \rangle_G$, the defining equation of isonormal process. We take $\mathbb{G} = \text{id}(G)$.
- (ii) In Part III, $\mathcal{Q} \neq \text{id}$ on G . The covariance structure of ξ is characterised by \mathcal{Q} . If \mathcal{Q} is of trace class, then ξ induces a centred Gaussian measure on G , whose RKHS is $\mathbb{G} = \mathcal{Q}^{1/2}(G)$ equipped with norm $\|h\|_{\mathbb{G}} = \|\mathcal{Q}^{-1/2}h\|_G$. Consequently, (4.2) defines a Borel mapping $X^{(n)}$ into G , and the interpretation $X^{(n)}(g) = \langle X^{(n)}, g \rangle_G$ leads to (4.3). Otherwise, if \mathcal{Q} is not of trace class, the same argument from the white noise case applies.

In both cases, the law of $X^{(n)}$ is dominated by ξ/\sqrt{n} is equivalent to $\mathcal{A}\theta \in \mathcal{Q}^{1/2}(G) = \mathbb{G}$. So we assume $\text{Ran } \mathcal{A} \subset \text{Ran}(\mathcal{Q}^{1/2})$. Moreover, the Kullback-Leibler distance is given in the lemmas below. For the proper Gaussian case, it is a direct consequence of Lemma 3.20.

Lemma 4.5. *Let ξ be a proper centred Gaussian on a Hilbert space G with \mathbb{G} being the RKHS. Let $\gamma_{g_1}, \gamma_{g_2}$ be the law of*

$$X_1 = g_1 + \xi \quad \text{and} \quad X_2 = g_2 + \xi$$

respectively. Then, we have

$$\mathbb{E}_{g_1} \left[\log \frac{d\gamma_{g_1}}{d\gamma_{g_2}} \right] = \frac{1}{2} \|g_1 - g_2\|_{\mathbb{G}}^2 \quad \text{and} \quad \text{Var}_{g_1} \left[\log \frac{d\gamma_{g_1}}{d\gamma_{g_2}} \right] = \|g_1 - g_2\|_{\mathbb{G}}^2,$$

where the subscript g_1 denotes that the integrals are calculated with respect to measure γ_{g_1} .

For the white noise case, it requires more elaboration. We consider the observations in the sequence model.

Lemma 4.6. *For $\theta = (\theta_1, \theta_2, \dots)$ let P_θ be the distribution of the random element $(Z_1 + \theta_1, Z_2 + \theta_2, \dots)$ in \mathbb{R}^∞ for Z_1, Z_2, \dots i.i.d. mean-zero normal variables with variance σ^2 . If $\theta \in \ell^2$, then P_θ is absolutely continuous relative to P_0 with log likelihood*

$$\log \frac{dP_\theta}{dP_0}(X_1, X_2, \dots) = \frac{1}{\sigma^2} \sum_{i=1}^{\infty} \theta_i X_i - \frac{1}{2\sigma^2} \sum_{i=1}^{\infty} \theta_i^2,$$

where the first series converges almost surely and in second mean. The expectation and variance of this variable are

$$\mathbb{E}_\theta \left[\log \frac{dP_\theta}{dP_0} \right] = \frac{1}{2\sigma^2} \sum_{i=1}^{\infty} \theta_i^2 \quad \text{and} \quad \text{Var}_\theta \left[\log \frac{dP_\theta}{dP_0} \right] = \frac{1}{\sigma^2} \sum_{i=1}^{\infty} \theta_i^2.$$

Proof. That the series converges in L^2 is clear from the fact that $\theta \in \ell^2$; the almost sure convergence next follows from the Itô-Nisio theorem. The expectation and variance of the right side are easy to compute as limits.

Write Σ_∞ for the right side of the display, and Σ_n for the expression obtained by replacing the infinite sums by the sums from 1 to n . Thus $\Sigma_n \rightarrow \Sigma_\infty$ almost surely. Since $\mathbb{E}_0 e^{2\Sigma_n} = e^{\sum_{i=1}^n \theta_i^2/\sigma^2}$ is uniformly bounded in n , it follows that e^{Σ_n} is uniformly integrable and hence converges in mean to e^{Σ_∞} . In particular, the mean of the latter variable is 1, the mean of the former variables.

It follows that the Borel measure on \mathbb{R}^∞ defined by $B \mapsto \mathbb{E}_0 1_B(X) e^{\Sigma_\infty}$ is a probability measure. For every Borel set B it is the limit of $\mathbb{E}_0 1_B(X) e^{\Sigma_n}$, which is $P_\theta(B)$ if B depends only on the first n coordinates, as e^{Σ_n} is the density of the distribution of $(Z_1 + \theta_1, \dots, Z_n + \theta_n)$ with respect to its distribution at $\theta = 0$. Since the Borel σ -field on \mathbb{R}^∞ is generated by the algebra of all cylinder sets, it follows that P_θ and the measure $B \mapsto \mathbb{E}_0 1_B(X) e^{\Sigma_\infty}$ agree. \square

Recall that a bounded operator $\mathcal{Q} : H \rightarrow H$ is called *diagonalisable*, if with an orthonormal basis $\{\varphi_k\}_{k \in \mathbb{N}}$ for H , for all $f \in H$,

$$\mathcal{Q}f = \sum_{k \in \mathbb{N}} q_k f_k \varphi_k \quad \text{with} \quad f_k = \langle f, \varphi_k \rangle_H.$$

The last result can be easily extended to the case of diagonalisable operators.

Corollary 4.7. *Let Z_k in Lemma 4.6 be independent mean-zero normal variables with variance σ_k^2 such that $\sup_k \sigma_k < \infty$. Then, the same result holds with*

$$\log \frac{dP_\theta}{dP_0}(X_1, X_2, \dots) = \sum_{i=1}^{\infty} \frac{\theta_i}{\sigma_i^2} X_i - \frac{1}{2} \sum_{i=1}^{\infty} \frac{\theta_i^2}{\sigma_i^2},$$

and

$$\mathbb{E}_\theta \left[\log \frac{dP_\theta}{dP_0} \right] = \frac{1}{2} \sum_{i=1}^{\infty} \frac{\theta_i^2}{\sigma_i^2} \quad \text{and} \quad \text{Var}_\theta \left[\log \frac{dP_\theta}{dP_0} \right] = \sum_{i=1}^{\infty} \frac{\theta_i^2}{\sigma_i^2}.$$

4.2.2 Discrete Observation

In practice, it is rare that complete information can be acquired from measurements, and only *partial* observations can be gathered. It is often designated by a natural number n , representing the amount of observations having been made. We consider the following type of partial observation. Let G be the space $L^2(\mathfrak{D})$ of square integrable functions on a bounded domain $\mathfrak{D} \subset \mathbb{R}^d$ with $d \in \mathbb{N}$. Instead of observing the entire trajectory of $\mathcal{A}\theta$ in G , we observe the process at the *design points*, a set $\mathfrak{D}_n = \{x_i\}_{i \leq n} \subset \mathfrak{D}$. To formalise the description, let

$$\mathcal{A}^{(n)} = \mathcal{E}^{(n)}\mathcal{A}, \quad \xi^{(n)} = (z_i)_{i \leq n},$$

where $\mathcal{E}^{(n)} : G \rightarrow \mathbb{R}^n$ is the evaluation operator at design points, and $(z_i)_{i \leq n}$ is the standard Gaussian in \mathbb{R}^n . The observation $X^{(n)}$ is a random vector $(X_i)_{i \leq n}$ in \mathbb{R}^n , with the coordinates given by

$$X_i = (\mathcal{A}\theta)(x_i) + z_i, \quad i = 1, \dots, n,$$

where z_i are i.i.d standard Gaussian variables.

In literature (e.g. [42]), the majority of discrete observations refers to the aforementioned type. While, less frequently, it may address the continuous observations being truncated, e.g. only the first n entries of the infinite vector $(X^{(n)}(w_k))_{k \in \mathbb{N}}$. We do not follow this convention, but it is often equivalent to the discrete observations at design points, with regularity assumptions on the signal $\mathcal{A}\theta$. For example, in Section 8.1, we will also present a concrete construction that the observation at design points implies the truncated continuous observation.

Connection to white noise model

We are going to give some heuristics to link the regression to the white noise model. Since only the observation is relevant, the shorthand notation $g = \mathcal{A}\theta$ is used. We assume that there exists a partition $\{\mathfrak{D}_i : 1 \leq i \leq n\}$ of \mathfrak{D} such that the subdomains are mutually disjoint and of equal volume $\text{vol } \mathfrak{D}_i \simeq 1/n$, and furthermore each subdomain \mathfrak{D}_i contains only one design point x_i .

The observations in white noise can be reduced to the projections on a orthonormal basis $\{v_k\}_{k \in \mathbb{N}}$ of $L^2(\mathfrak{D})$. If the basis functions are continuous as well as the drift function g , the projection can be approximated by the Riemann summation,

$$\int_{\mathfrak{D}} g(x)v(x) dx \approx \sum_{i=1}^n g(x_i)v(x_i) \text{vol } \mathfrak{D}_i.$$

Recall $\text{vol } \mathfrak{D}_i \simeq 1/n$. By plugging in the noisy observations from regression, the noise is

$$\sum_{i=1}^n z_i \frac{v(x_i)}{n} \sim \mathcal{N}_{\mathbb{R}} \left(0, \frac{1}{n} \sum_{i=1}^n v(x_i)^2 \text{vol } \mathfrak{D}_i \right).$$

On the other hand, from Example 3.31,

$$\frac{1}{\sqrt{n}}\xi(v) = \frac{1}{\sqrt{n}} \int_{\mathfrak{D}} \xi(t)v(t) dt \sim \mathcal{N}_{\mathbb{R}} \left(0, \frac{1}{n} \int_{\mathfrak{D}} v(x)^2 dt \right).$$

Both of the noises are centred and their variances are asymptotically equal.

Remark 4.8 (Continuity). It is noteworthy that the assumption on the continuity of both drift g and the test function v is crucial. As we will see in the subsequent Chapter 7 and Chapter 8, the drift $\mathcal{A}f$ needs to possess certain smoothness, which corresponds to the minimal requirement on securing continuity.

The heuristics above shows that the regression model converges to the continuous model in a certain sense, while more data are gathered. However, conversely, the point evaluations of white noise model does not lead to a regression model. This is because point evaluation is only defined for continuous functions, but the white noise ξ is only distribution-valued (see Example 3.31). One cannot talk about the values of ξ evaluated at points.

4.3 Bayesian Contraction for Gaussian Linear Models

In this section we present a general theorem on the posterior contraction for linear problems in the following form. First we introduce the smoothness class. Let $(\Theta, \langle \cdot, \cdot \rangle)$ be a separable Hilbert space. Given a set \mathfrak{S} in a finite dimensional ordered space (e.g. \mathbb{R}^m , $m \in \mathbb{N}$), assume that there exists a smoothness scale $\{\Theta_s, \langle \cdot, \cdot \rangle_s\}_{s \in \mathfrak{S}}$ with $\Theta_0 = \Theta$, such that for $s < t$, the induced norm $\|\cdot\|_t$ is strictly stronger than $\|\cdot\|_s$ and $\Theta_t \subset \Theta_s$. We consider the observation scheme as introduced in Section 4.2.1, a transform $\mathcal{A}\theta$ of θ in G with Gaussian process ξ indexed by G ,

$$X^{(n)} = \mathcal{A}\theta + \frac{1}{\sqrt{n}}\xi. \quad (4.5)$$

Let ξ be a centred Gaussian process indexed by a separable Hilbert space G with covariance operator \mathcal{Q} . Assume that $\mathcal{A} : \Theta \rightarrow G$ is a bounded linear operator satisfying

$$\text{Ran } \mathcal{A} \subset \mathbb{H} := \mathcal{Q}^{1/2}(G),$$

where \mathbb{H} is a Hilbert space equipped with the induced inner product (see Lemma A.11),

$$\langle \cdot, \cdot \rangle_{\mathbb{H}} := \langle \mathcal{Q}^{-1/2} \cdot, \mathcal{Q}^{-1/2} \cdot \rangle_G.$$

We will only consider the *non-degenerate* case, i.e. \mathbb{H} is dense in G . To link the spaces G and Θ , we introduce the following assumption.

Assumption 4.9. The family $\{\mathcal{R}_n\}_{n \in \mathbb{N}}$ of linear reconstruction operators $\mathcal{R}_n : \mathbb{H} \rightarrow \Theta$ satisfies the following properties.

- (i) With $j_n \rightarrow \infty$, let $\{W_{j_n}\}_{n \in \mathbb{N}}$ be a sequence of subspaces of G such that, for all $n \in \mathbb{N}$,

$$\dim W_{j_n} = j_n, \quad W_{j_n} \subset W_{j_{n+1}}, \quad W_{j_n} \in \mathbb{H},$$

and the kernel (i.e. the null space) of \mathcal{R}_n is $W_{j_n}^\perp$.

- (ii) Let ρ_n be a monotonically nondecreasing sequence which may go to infinity as $n \rightarrow \infty$. \mathcal{R}_n satisfies

$$\|\mathcal{R}_n\|_{G; \Theta} \simeq \rho_n. \quad (4.6)$$

(iii) For $s \in \mathfrak{S}$, let $\{\delta(j_n, s)\}_{s \in \mathfrak{S}}$ be a family of monotonically decreasing functions such that $\delta(j_n, s) \downarrow 0$ as $n \rightarrow \infty$, for all $s \in \mathfrak{S}$. We assume

$$\|\mathcal{R}_n \mathcal{A} \theta - \theta\|_{\Theta} \leq \delta(j_n, s) \|\theta\|_s \quad (4.7)$$

holds, for $\theta \in \Theta_s$, $s \in \mathfrak{S}$.

We form the posterior distribution $\Pi_n(\cdot | X^{(n)})$ as in (4.1), given a prior Π_n on the space $\Theta = \Theta_0$ and an observation $X^{(n)}$, whose conditional distribution given θ is determined by the model (4.5). We study this random distribution under the assumption that $X^{(n)}$ follows the model (4.5) for a given ‘true’ function $\theta = \theta_0$, which we assume to be an element of Θ_β in a given smoothness class $(\Theta_s)_{s \in \mathfrak{S}}$.

The result is based on an extension of the testing approach of [35] to the parameter inference of (4.5). The recovery problem is handled with the help of the reconstruction family from Assumption 4.9. We note that the reconstruction family only appears as a tool to state and derive a posterior contraction rate. In our context it does not enter into the solution of the linear problem, which is achieved through the Bayesian method.

Clearly the reconstructed signal $\theta^{(n)} = \mathcal{R}_n \mathcal{A} \theta$ is an approximation to θ , which will be better for increasing n , but increasingly complex, when ρ_n is unbounded. The following theorem uses ρ_n that balances approximation to complexity, where the complexity is implicitly determined by a testing criterion.

Theorem 4.10. *Assume $\theta_0 \in \Theta_\beta$ with $\beta \in \mathbb{R}_+^m$ such that $\mathcal{A} \Theta_\beta \subset \mathbb{H}$.*

For $\varepsilon_n \downarrow 0$ such that $n\varepsilon_n^2 \rightarrow \infty$, suppose Assumption 4.9 holds with

$$j_n \leq cn\varepsilon_n^2, \quad (4.8)$$

where c is a positive constant. In addition, for $\eta_n \geq \varepsilon_n$, assume

$$\eta_n \geq \rho_n \varepsilon_n, \quad (4.9)$$

$$\eta_n \geq \delta(j_n, \beta). \quad (4.10)$$

Let $\theta^{(n)}$ denote a reconstruction estimator $\mathcal{R}_n \mathcal{A} \theta$ to θ . Consider prior probability distributions Π on Θ satisfying

$$\Pi(\theta : \mathcal{A} \theta \in \mathbb{H}) = 1, \quad (4.11)$$

$$\Pi(\theta : \|\mathcal{A} \theta - \mathcal{A} \theta_0\|_{\mathbb{H}} < \varepsilon_n) \geq e^{-n\varepsilon_n^2}, \quad (4.12)$$

$$\Pi(\theta : \|\theta^{(n)} - \theta\|_0 > \eta_n) \leq e^{-4n\varepsilon_n^2}. \quad (4.13)$$

Then the posterior distribution in the model (4.5) contracts at the rate η_n at θ_0 , i.e. for a sufficiently large constant M we have $\Pi_n(\theta : \|\theta - \theta_0\|_0 > M\eta_n | X^{(n)}) \rightarrow 0$, in probability under the law of $X^{(n)}$ given by (4.5) with $\theta = \theta_0$.

The conditions in Theorem 4.10 deserve some explanations.

The desirable structure of priors is characterised by eqs. (4.11) to (4.13). Since we only consider dominated models in Gaussian noise, (4.11) assures the validity of the density version (4.1) of Bayes’ formula. (4.12) is the *prior mass* condition,

that is, a set around the truth should have sufficient probability mass at most exponentially decay with $n\varepsilon_n^2$. The last condition (4.13) states a requirement on the prior related to the testing approach via the constructed signals. It says that the prior mass on a set, whose elements, even without noise, cannot be accurately reconstructed quantified by η_n , should have prior mass upper bounded by a negative exponential of $n\varepsilon_n^2$.

The final contraction rate is determined with the constraints eqs. (4.8) to (4.10). The first (4.8) can be interpreted as a consideration that the reconstruction should be capped by the available information, indicated by $n\varepsilon_n^2$. The second constraint (4.9) shows the cost to recover the original signal θ from the transform \mathcal{A} , and the last one (4.10) simply states the fact that the rate cannot be faster than the recovery rate for noiseless signals.

We conclude this section with proving the main theorem of this section.

Proof of Theorem 4.10. Using Lemma 4.5, Lemma 4.6, and Corollary 4.7 from Section 4.2.1, the Kullback-Leibler divergence and variation between the distributions of $X^{(n)}$ under two functions θ and θ_0 are

$$n\|\mathcal{A}\theta - \mathcal{A}\theta_0\|_{\mathbb{H}}^2/2 \quad \text{and} \quad n\|\mathcal{A}\theta - \mathcal{A}\theta_0\|_{\mathbb{H}}^2.$$

Therefore the neighbourhoods $B_{n,2}(\theta_0, \varepsilon)$ in (8.19) of [35] contain the ball $\{\theta \in \Theta : \|\mathcal{A}\theta - \mathcal{A}\theta_0\|_{\mathbb{H}} \leq \varepsilon\}$. By assumption (4.12) this has prior mass at least $e^{-n\varepsilon_n^2}$.

Because the quotient of the left sides of (4.12) and (4.13) is $o(e^{-2n\varepsilon_n^2})$, the posterior probability of the set $\{f : \|\theta^{(n)} - \theta\|_0 > \eta_n\}$ tends to zero, by Theorem 8.20 in [35].

By a variation of Theorem 8.22 in [35] it is now sufficient to show the existence of tests τ_n such that, for some $M > 0$,

$$\mathbb{P}_{\theta_0}^{(n)} \tau_n \rightarrow 0, \quad \sup_{\substack{\theta: \|\theta - \theta_0\|_0 > M\eta_n, \\ \|\theta^{(n)} - \theta\|_0 \leq \eta_n}} \mathbb{P}_{\theta}^{(n)}(1 - \tau_n) \leq e^{-4n\varepsilon_n^2}.$$

Indeed, in the case that the prior mass condition (8.20) in Theorem 8.22 of [35] can be strengthened to (8.22), as is the case in our setup in view of (4.12), it suffices to verify (8.24) only for a single value of j . Furthermore, we can apply Theorem 8.22 with the metrics $d_n(x, y) = \|x - y\|_0 \varepsilon_n / \eta_n$ in order to reduce the restriction $d_n(\theta, \theta_{n,0}) > M\varepsilon_n$ to $\|\theta - \theta_0\|_0 > M\eta_n$.

Let $\{e_k\}_{k \leq j_n}$ be an G -orthonormal basis of W_{j_n} , and denote the G -orthogonal projection onto W_{j_n} by

$$\mathcal{P}_{j_n} : G \rightarrow W_{j_n} \subset \mathbb{H}.$$

By slight abuse of notation, define the projection \mathcal{P}_{j_n} of process ξ onto W_{j_n} by

$$\mathcal{P}_{j_n} \xi := \sum_{i \leq j_n} \xi(e_k) e_k,$$

where $\xi(e_k)$ is a zero-mean Gaussian random variable with covariance $\langle e_k, \mathcal{Q}e_k \rangle_G$. $\mathcal{P}_{j_n} \xi$ is a proper Gaussian element in the finite-dimensional space W_{j_n} , because

each $\xi(e_k)$ is Gaussian. Furthermore, notice that for any $u, v \in \mathbb{H}$,

$$\begin{aligned} \mathbb{E}\langle \mathcal{P}_{j_n} \xi, u \rangle \langle \mathcal{P}_{j_n} \xi, v \rangle &= \mathbb{E} \left[\left(\sum_{k \leq j_n} \xi(e_k) \langle u, e_k \rangle \right) \left(\sum_{l \leq j_n} \xi(e_l) \langle v, e_l \rangle \right) \right] \\ &= \sum_{k, l \leq j_n} \langle u, e_k \rangle \langle v, e_l \rangle \langle e_k, \mathcal{Q} e_l \rangle_G \\ &= \langle \mathcal{P}_n u, \mathcal{Q} \mathcal{P}_n v \rangle_G. \end{aligned}$$

Since $\mathcal{P}_n = \mathcal{P}_n^*$, $\mathcal{P}_{j_n} \xi$ is a centred Gaussian with covariance $\mathcal{P}_{j_n} \mathcal{Q} \mathcal{P}_{j_n}^*$. Thus, the following expression makes sense,

$$\mathcal{P}_{j_n} X^{(n)} = \mathcal{P}_{j_n} \mathcal{A} \theta + \frac{1}{\sqrt{n}} \mathcal{P}_{j_n} \xi.$$

Besides, because of the first property in Assumption 4.9, we have

$$f^{(n)} = \mathcal{R}_n \mathcal{A} f = \mathcal{R}_n \circ \mathcal{P}_{j_n} \mathcal{A} f, \quad \forall f \in \Theta.$$

Now we claim that

$$\mathcal{R}_n \circ \mathcal{P}_{j_n} X^{(n)} = \theta^{(n)} + \frac{1}{\sqrt{n}} \mathcal{R}_n \circ \mathcal{P}_{j_n} \xi \quad (4.14)$$

is a well-defined Gaussian random element in Θ . It suffices to show the noise $\mathcal{R}_n \circ \mathcal{P}_{j_n} \xi$ is a proper random element in Θ . Denote $\mathcal{R}_n \circ \mathcal{P}_{j_n}$ by $\widehat{\mathcal{R}}_n$. Since

$$\begin{aligned} \mathbb{E} \|\widehat{\mathcal{R}}_n \xi\|_{\Theta}^2 &= \text{Trace}(\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q} \mathcal{P}_{j_n} \mathcal{R}_n^*) = \text{Trace} \left[(\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q}^{1/2}) (\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q}^{1/2})^* \right] \\ &= \|\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q}^{1/2}\|_{HS}^2 = \sum_{k \leq j_n} \|\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q}^{1/2} e_k\|^2 \leq \|\mathcal{R}_n\|^2 \sum_{k \leq j_n} \|\mathcal{P}_{j_n} (\mathcal{Q}^{1/2} e_k)\|^2 \\ &\lesssim j_n \rho_n^2, \end{aligned}$$

where we use the properties of Schatten norms (see Proposition A.19 and Proposition A.20), and the fact that $\|\mathcal{P}_{j_n} (\mathcal{Q}^{1/2} e_k)\| \leq \|\mathcal{Q}^{1/2} e_k\|_G \leq \|\mathcal{Q}^{1/2}\|$. The preceding inequality shows that $\widehat{\mathcal{R}}_n X^{(n)}$ is indeed a proper random element in Θ . In addition, the weak second moment of the variable $\widehat{\mathcal{R}}_n \xi$ is

$$\begin{aligned} \sup_{\|\theta\|_0 \leq 1} \mathbb{E} \langle \theta, \widehat{\mathcal{R}}_n \xi \rangle_{\Theta}^2 &= \sup_{\|\theta\|_0 \leq 1} \langle \theta, \mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q} \mathcal{P}_{j_n} \mathcal{R}_n^* \theta \rangle_{\Theta} \\ &= \sup_{\|\theta\|_0 \leq 1} \|\mathcal{Q}^{1/2} \mathcal{P}_{j_n} \mathcal{R}_n^* \theta\|_G^2 = \|\mathcal{Q}^{1/2} \mathcal{P}_{j_n} \mathcal{R}_n^*\|_{\Theta; G}^2 \\ &= \|\mathcal{R}_n \mathcal{P}_{j_n} \mathcal{Q}^{1/2}\|_{G; \Theta}^2 = \|\mathcal{R}_n\|_{G; \Theta}^2 \|\mathcal{P}_{j_n} \mathcal{Q}^{1/2}\|^2 \simeq \rho_n^2, \end{aligned}$$

where the last inequality follows from the boundedness of \mathcal{P}_{j_n} and $\mathcal{Q}^{1/2}$, and Assumption 4.9.

The first inequality shows that the first moment $\mathbb{E} \|\widehat{\mathcal{R}}_n \xi\|_0$ of the variable $\|\widehat{\mathcal{R}}_n \xi\|_0$ is bounded above by $\sqrt{j_n} \rho_n$. By Borell's inequality (see Lemma 3.11),

applied to the Gaussian random variable $\widehat{\mathcal{R}}_n \xi$ in Θ_0 , we see that there exist positive constants a and b such that, for every $t > 0$,

$$\mathbb{P}\left(\|\widehat{\mathcal{R}}_n \xi\|_0 > t + a\sqrt{j_n \rho_n}\right) \leq e^{-bt^2/\rho_n^2}.$$

With $t = 2\sqrt{n}\eta_n/\sqrt{b}$, for η_n and ε_n satisfying (4.9) and (4.10) this yields, for some $a_1 > 0$,

$$\mathbb{P}\left(\|\widehat{\mathcal{R}}_n \xi\|_0 > a_1\sqrt{n}\eta_n\right) \leq e^{-4n\varepsilon_n^2}. \quad (4.15)$$

We apply this to bound the error probabilities of the tests

$$\tau_n = 1\{\|\widehat{\mathcal{R}}_n X^{(n)} - \theta_0\|_0 \geq M_0\eta_n\}, \quad (4.16)$$

where M_0 is a given constant, to be determined.

Under θ_0 , the decomposition (4.14) is valid with $\theta = \theta_0$, and hence $\widehat{\mathcal{R}}_n X^{(n)} - \theta_0 = n^{-1/2}\widehat{\mathcal{R}}_n \xi + \theta_0^{(n)} - \theta_0$. By the triangle inequality it follows that $\tau_n = 1$ implies that $n^{-1/2}\|\widehat{\mathcal{R}}_n \xi\|_0 \geq M_0\eta_n - \|\theta_0^{(n)} - \theta_0\|_0$. By (4.7), the assumption that $f_0 \in H_\beta$ implies that $\|\theta_0^{(n)} - \theta_0\|_0 \leq M_1\delta(j_n, \beta)$, for some M_1 , which is further bounded by $M_1\eta_n$, by assumption (4.10). Hence the probability of an error of the first kind satisfies

$$\mathbb{P}_{\theta_0}^{(n)} \tau_n \leq \mathbb{P}\left(\frac{1}{\sqrt{n}}\|\widehat{\mathcal{R}}_n \xi\|_0 \geq (M_0 - M_1)\eta_n\right).$$

For $M_0 - M_1 > a_1$, the right side is bounded by $e^{-4n\varepsilon_n^2}$, by (4.15).

Under θ the decomposition (4.14) gives that $\widehat{\mathcal{R}}_n X^{(n)} - \theta_0 = n^{-1/2}\widehat{\mathcal{R}}_n \xi + \theta^{(n)} - \theta_0$. By the triangle inequality $\tau_n = 0$ implies that $n^{-1/2}\|\widehat{\mathcal{R}}_n \xi\|_0 \geq \|\theta^{(n)} - \theta_0\|_0 - M_0\eta_n$. For θ such that $\|\theta - \theta_0\|_0 > M\eta_n$ and $\|\theta - \theta^{(n)}\|_0 \leq \eta_n$, we have $\|\theta^{(n)} - \theta_0\|_0 \geq (M - 1)\eta_n$. Hence the probability of an error of the second kind satisfies

$$\mathbb{P}_\theta^{(n)}(1 - \tau_n) \leq \mathbb{P}\left(\frac{1}{\sqrt{n}}\|\widehat{\mathcal{R}}_n \xi\|_0 \geq (M - 1 - M_0)\eta_n\right),$$

For $M - 1 - M_0 > a_1$, this is bounded by $e^{-4n\varepsilon_n^2}$, by (4.15).

We can first choose M_0 large enough so that $M_0 - M_1 > a_1$, and next M large enough so that $M - 1 - M_0 > a_1$, to finish the proof. \square

Part II

Inverse Problems

In Part II, under the Bayesian framework, we study the statistical linear inverse problem, i.e. the Gaussian linear model of the following form,

$$Y^{(n)} = \mathcal{A}^{(n)}f + \xi^{(n)},$$

where the operator $\mathcal{A}^{(n)}$ has no bounded inverse. The investigation is conducted focusing on the following aspects. First, we study the problem in a general framework, which in particular covers the singular value decomposition framework¹. Second, we are interested in a unified evaluation procedure for the posterior contraction of priors, both conjugate and non-conjugate. Lastly, we consider the inference of inverse problems with discrete observations.

This part is organised as follows. In Chapter 5, we formulate the inverse problem in the smoothness scales from Chapter 2. In particular, we introduce the *Galerkin* projection method, which serves as an important tool in demonstrating posterior contractions. After the problem has been properly stated, in Chapter 6 we study the continuous model contaminated by the white noise, in which the posterior contractions are obtained using a variant of the general testing approach (Theorem 4.10), without invoking any conjugacy. The rest of the chapters in this part tackle the inverse problem with discrete observations from two angles. Chapter 7 utilises the Gaussian conjugacy in linear models to study the posterior performance, and Chapter 8 extends the results obtained in Chapter 6 to the regression model.

¹See the discussion in Section 5.4 and the definition in Section 7.1.

Chapter 5

Linear Inverse Problems

5.1 Introduction

In a statistical inverse problem one observes a noisy version of a transformed signal $\mathcal{A}f$ and wishes to recover the unknown parameter f . In Part II, we consider linear inverse problems with different observation schemes. The continuously observed model is the white noise model in the given form,

$$Y^{(n)} = \mathcal{A}f + \frac{1}{\sqrt{n}}\xi, \quad (5.1)$$

where $\mathcal{A} : H \rightarrow G$ is a known bounded linear operator between separable Hilbert spaces H and G , and ξ is a stochastic ‘noise’ process, which is multiplied by the scalar ‘noise level’ $n^{-1/2}$. The white noise model represents a limiting case (in an appropriate sense) of the inverse regression model

$$Y_i = (\mathcal{A}f)(x_i) + z_i, \quad i = 1, \dots, n, \quad (5.2)$$

where z_i are independent standard normal random variables. Insights gained in inverse problems in the white noise model shed light on the behaviour of statistical procedures in the inverse regression model, which is the one encountered in actual practice, as the signal f can be typically observed only on a discrete grid of points. Both models belong to the form of Gaussian linear models introduced in Section 4.2,

$$Y^{(n)} = \mathcal{A}^{(n)}f + \xi^{(n)},$$

where $\mathcal{A}^{(n)} = \mathcal{A}$ in the white noise case and $\mathcal{A}^{(n)} = \mathcal{E}^{(n)}\mathcal{A}$ with evaluation operator $\mathcal{E}^{(n)}$ in the regression case (see Section 4.2.1 and Section 4.2.2). In this chapter, we explore the structure of operator \mathcal{A} , and in the subsequent chapters in Part II we study the relevant statistical models.

The problem is to infer f from the observation $Y^{(n)}$. To this purpose we assume that the *forward operator* \mathcal{A} is injective, but we shall be interested in the case that the inverse \mathcal{A}^{-1} , defined on the range of \mathcal{A} is not continuous (or equivalently the range of \mathcal{A} is not closed in G). The problem of recovering f from $Y^{(n)}$ is then *ill-posed*, and *regularization* methods are necessary in order to

‘invert’ the operator \mathcal{A} . These consist of constructing an approximation to \mathcal{A}^{-1} , with natural properties such as boundedness and whose domain includes the data $Y^{(n)}$, and applying this to $Y^{(n)}$. By the discontinuity of the inverse \mathcal{A}^{-1} , the noise present in the observation is necessarily multiplied, and regularization is focused on balancing the error in the approximation to \mathcal{A}^{-1} to the size of the magnified noise, in order to obtain a solution that is as close as possible to the true signal f . In this article we study this through the convergence rates of the regularized solutions to a true parameter f , as $n \rightarrow \infty$, i.e. as the noise level tends to zero. In particular, we consider contraction rates of posterior distributions resulting from a Bayesian approach to the problem.

It is also possible to consider the model (5.1) with a noise variable ξ that takes its values inside the Hilbert space G . In Section 6.7 we briefly note some results on this ‘coloured noise’ model, but our main focus is the white noise case.

The chapter is organized as follows. In Section 5.2 we introduce in greater detail our setup along with the assumptions that will be used in this part. We also present some examples for illustration. Then, an estimator, particularly suitable for the framework of inverse problems introduced in the former section, is constructed in Section 5.3. We conclude this short chapter with notes and comments in Section 5.4.

5.2 Inverse Nature

The forward operator \mathcal{A} in the model (5.1) and (5.2) is a bounded linear operator $\mathcal{A} : H \rightarrow G$ between the separable Hilbert spaces H and G , and is assumed to be smoothing. The following assumption makes this precise. This assumption is satisfied in many examples and is common in the literature (for instance [19, 38, 71]).

Recall the concept of smoothness scales from Chapter 2. In Definition 2.1 the space H is embedded as $H = H_0$ in the smoothness scale $(H_s)_{s \in \mathbb{R}}$ and hence has norm $\|\cdot\|_0$.

Assumption 5.1 (Smoothing property of \mathcal{A}). For some $\gamma > 0$ the operator $\mathcal{A} : H_{-\gamma} \rightarrow G$ is injective and bounded and, for every $f \in H_0$,

$$\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}. \quad (5.3)$$

Example 5.2 (SVD). If the operator $\mathcal{A} : H \rightarrow G$ is compact, then the positive self-adjoint operator $\mathcal{A}^*\mathcal{A} : H \rightarrow H$ possesses a countable orthonormal basis of eigenfunctions ϕ_i , which can be arranged so that the corresponding sequence of eigenvalues λ_i decreases to zero. If \mathcal{A} is injective, then all eigenvalues, whose roots are known as the *singular values* of \mathcal{A} , are strictly positive. Suppose that there exists $\gamma > 0$ such that

$$\lambda_i \simeq i^{-2\gamma}. \quad (5.4)$$

If we construct the smoothness classes $(H_s)_{s \in \mathbb{R}}$ from the basis $(\phi_i)_{i \in \mathbb{N}}$ and the numbers $b_i = i$ as in Example 2.6, then (5.3) is satisfied.

Indeed, we can write \mathcal{A} in polar decomposition as $\mathcal{A}f = U(\mathcal{A}^*\mathcal{A})^{1/2}f$, for a partial isometry $U : \text{Ran}(\mathcal{A}) \rightarrow G$, and then have $\mathcal{A}f = U \sum_i f_i \sqrt{\lambda_i} \phi_i$, so that $\|\mathcal{A}f\| = \|\sum_i f_i i^{-\gamma} \phi_i\|_0 \simeq \|f\|_{-\gamma}$.

Thus constructions using the singular value decomposition of \mathcal{A} can always be accommodated in the more general setup described in the preceding.

For more interesting illustrations of the preceding setup, consider linear differential equations of the form

$$Du(x) = f(x), \quad x \in \mathcal{D} \subset \mathbb{R}^d,$$

where D is a differential operator. Under appropriate boundary conditions, the solution u can often be expressed in terms of the Green's function associated with D , through a kernel operator

$$u(x) = \int_{\mathcal{D}} k(x, t) f(t) dt =: \mathcal{A}f(x). \quad (5.5)$$

The operator \mathcal{A} typically lifts a function $f \in L^2$ to a Sobolev space of functions, as in Example 2.5. The ill-posedness surfaces when one observes the state u with noise (which deteriorates the smoothness), and tries to recover the source function f . For illustration we include two concrete examples from the literature.

Example 5.3 (Poisson equation). The following example can be found in Sections 10.4 and 11.2 in [44]. Let $(H_s)_{s \in \mathbb{R}}$ be the periodic Sobolev spaces of (generalized) functions satisfying the boundary condition $f(0) = f(1) = 0$. Consider the following boundary problem,

$$u'' := \frac{d^2 u}{dx^2} = -f, \quad f \in H_0,$$

with the Dirichlet boundary condition: $u(0) = u(1) = 0$. The unique solution $u \in H_2$ is given by

$$u(x) = \mathcal{A}f(x) = \int_0^1 k(x, t) f(t) dt,$$

where

$$k(x, t) = \begin{cases} (1-x)t, & \text{if } x \geq t, \\ (1-t)x, & \text{otherwise.} \end{cases}$$

The operator $\mathcal{A} : H_0 \rightarrow H_0$ is Hilbert-Schmidt and hence compact, and therefore has no bounded inverse. On the other hand the inverse exists as bounded operator $\mathcal{A}^{-1} : H_2 \rightarrow H_0$, and is given by $\mathcal{A}^{-1}f = -f''$.

When $\mathcal{D} \subset \mathbb{R}$, the Sobolev norm is equivalent to $\|f\|_2 = \|f\|_0 + \|f''\|_0$ (Page 217 in [10]). Since $(\mathcal{A}f)'' = f$, we have $\|f\|_0 \leq \|\mathcal{A}f\|_2 = \|\mathcal{A}f\|_0 + \|(\mathcal{A}f)''\|_0 \leq (\|\mathcal{A}\| + 1)\|f\|_0$, i.e. $\|\mathcal{A}f\|_2 \simeq_{\mathcal{A}} \|f\|_0$.

Since the kernel is symmetric, \mathcal{A} is self-adjoint. Besides, \mathcal{A} is an isomorphism between H_2 and H_0 as shown above. Hence

$$\|\mathcal{A}f\|_0 = \sup_{\|h\|_0 \leq 1} |\langle h, \mathcal{A}f \rangle_0| = \sup_{\|h\|_0 \leq 1} |\langle \mathcal{A}h, f \rangle_0| \simeq_{\mathcal{A}} \sup_{\|h\|_2 \leq 1} |\langle h, f \rangle_0| = \|f\|_{-2},$$

by norm duality argument, for all $f \in H_0$. This shows that (5.3) holds with $\gamma = 2$.

Example 5.4 (Symm's equation [57]). Consider the Laplace equation $\Delta u = 0$ in a bounded set $\Omega \subset \mathbb{R}^2$ with boundary condition $u = g$ on the boundary $\partial\Omega$. The singular layer potential, a boundary integral

$$u(x) = -\frac{1}{\pi} \int_{\partial\Omega} h(y) \ln|x-y| ds(y), \quad x \in \Omega,$$

solves the boundary value problem if and only if the density h , belonging to the space $C(\partial\Omega)$ of continuous functions on $\partial\Omega$, solves *Symm's equation*

$$-\frac{1}{\pi} \int_{\partial\Omega} h(y) \ln|x-y| ds(y) = g(x), \quad x \in \partial\Omega. \quad (5.6)$$

Assume the boundary $\partial\Omega$ has a parametrization of the form $\{\rho(s), s \in [0, 2\pi]\}$, for some 2π -periodic analytic function $\rho : [0, 2\pi] \rightarrow \mathbb{R}^2$ such that $|\dot{\rho}(s)| > 0$ for all s . Then Symm's equation takes the following form,

$$\mathcal{A}f(z) := -\frac{1}{\pi} \int_0^{2\pi} \log|\rho(z) - \rho(s)| f(s) ds = g(\rho(z)), \quad z \in [0, 2\pi],$$

where $f(s) = h(\rho(s))|\dot{\rho}(s)|$. As shown in Theorem 3.18 from [57], the operator \mathcal{A} satisfies (5.3), with $\gamma = 1$ and $(H_s)_{s \in \mathbb{R}}$ being periodic Sobolev spaces on $[0, 2\pi]$.

The following example is an inverse problem in Hilbert scales (see Section 2.2).

Example 5.5 (Abel operator). For a given kernel function $K : (0, 1) \times (0, 1) \rightarrow \mathbb{R}$ and $\alpha \in (0, 1]$, consider the operator $A : L^2(0, 1) \rightarrow L^2(0, 1)$ given by

$$Af(x) = \frac{1}{\Gamma(\alpha)} \int_0^x (x-s)^{\alpha-1} K(x,s) f(s) ds.$$

For $K = 1$ this gives the classical *Abel operator*. Under mild smoothness conditions on K , it is shown in [39], Theorem 1, that A is smoothing (i.e. (5.3) holds) of order $\gamma = 1$ for the Sobolev scale generated by the root negative Laplacian under the Cauchy boundary condition, described in Example 2.10.

In the Bayesian setup we model a function through a prior. When a true function is known to satisfy certain boundary conditions, as in many problems involving differential forward operators, we can incorporate these in the by choosing an appropriate generating operator. For an operator \mathcal{A} defined in terms of the Laplacian and the same boundary conditions the smoothing condition (5.3) will be satisfied. The following is a another example of a pair of Λ and \mathcal{A} .

Example 5.6 (Volterra operator). Consider the operator $A : L^2((0, 1)^2) \rightarrow L^2((0, 1)^2)$ on functions $f : (0, 1)^2 \rightarrow \mathbb{R}$ on the unit square satisfying the differential equation

$$D_{x,y} \mathcal{A}f = f, \quad D_{x,y} = \frac{\partial^2}{\partial x \partial y}.$$

We can render the solution of the equation unique by imposing boundary conditions. Two solutions are given by

$$\mathcal{A}f(x, y) = \int_0^x \int_0^y f(s, t) ds dt,$$

$$\mathcal{A}_0f(x, y) = \mathcal{A}f(x, y) - \int_0^1 \mathcal{A}f(x, t) dt - \int_0^1 \mathcal{A}f(s, y) ds + \int_0^1 \int_0^1 \mathcal{A}f(s, t) ds dt.$$

The first satisfies the boundary conditions $\mathcal{A}f(x, 0) = \mathcal{A}f(0, y) = 0$, while the second is obtained from the first by subtracting its projection on the set of all functions of the form $(x, y) \mapsto g_1(x) + g_2(y)$, which forms the kernel of the differential operator. Other boundary conditions will still give different versions of the operator.

We claim that \mathcal{A}_0 is smoothing of order $\gamma = 1$ for the Hilbert scale generated by the root Λ of $D_{x,y}^2$ with Dirichlet boundary condition, while \mathcal{A} is smoothing relative to the scale of L combined with Cauchy boundary condition.

The scale under the Dirichlet boundary condition is generated by the orthogonal system of eigenfunctions $e_{k,l} : (x, y) \mapsto \sin(k\pi x) \sin(l\pi y)$, for $(k, l) \in \mathbb{N}^2$, the tensor product of the basis of the one-dimensional Dirichlet-Laplacian as in Example 2.10, with corresponding eigenvalues are $k^2l^2\pi^4$. By explicit calculation

$$\begin{aligned} \mathcal{A}e_{k,l}(x, y) &= \frac{1}{kl\pi^2} [\cos(k\pi x) \cos(l\pi y) - \cos(k\pi x) - \cos(l\pi y) + 1], \\ \mathcal{A}_0e_{k,l}(x, y) &= \frac{1}{kl\pi^2} \cos(k\pi x) \cos(l\pi y). \end{aligned}$$

The functions $(x, y) \mapsto \cos(k\pi x) \cos(l\pi y)$, for $(k, l) \in (\mathbb{N} \cup \{0\})^2$ form an orthogonal basis of $L^2((0, 1)^2)$. We conclude that for $f = \sum_{k,l} f_{k,l} e_{k,l}$,

$$\begin{aligned} \|\mathcal{A}f\|^2 &\simeq \sum_{k,l} \frac{f_{k,l}^2}{k^2l^2} + \sum_k \left(\sum_l \frac{f_{k,l}}{kl} \right)^2 + \sum_l \left(\sum_k \frac{f_{k,l}}{kl} \right)^2 + \left(\sum_{k,l} \frac{f_{k,l}}{kl} \right)^2, \\ \|\mathcal{A}_0f\|^2 &\simeq \sum_{k,l} \frac{f_{k,l}^2}{k^2l^2} \simeq \|f\|_{-1}^2, \end{aligned}$$

where $\|\cdot\|_{-1}$ refers to the scale of Λ with Dirichlet boundary condition. The first equation shows that the operator \mathcal{A} is not smoothing in this scale, but in general satisfies $\|\mathcal{A}f\| \gtrsim \|f\|_{-1}$.

On the other hand, the Cauchy boundary condition generates the system of eigenfunctions $(x, y) \mapsto \cos((k-1/2)\pi x) \cos((l-1/2)\pi y)$, for $(k, l) \in \mathbb{N}^2$. These can be seen to be also the eigenfunctions of $\mathcal{A}^*\mathcal{A}$, and hence the smoothing property of \mathcal{A} fits the SVD framework, as in Example 5.2.

The two versions \mathcal{A} and \mathcal{A}_0 possess the same inverse operator, namely the differential operator $D_{x,y}$ used for their definitions. This suggests that from the point of view of reconstructing f in the inverse problem it should not matter whether one is provided with a noisy version of either $\mathcal{A}f$ or \mathcal{A}_0f as input data, seemingly contradicting the fact that the operators are smoothing in different scales. This paradox may be resolved by considering \mathcal{A} or \mathcal{A}_0 as maps into the quotient

space $L^2((0, 1)^2)/N(D_{x,y})$, where N denotes the kernel of the operator. The map $f \mapsto [\mathcal{A}f] = [\mathcal{A}_0f]$ into the class of $\mathcal{A}f$ in this quotient space is injective and can be shown to be appropriately smoothing (see (5.10)-(5.11)), and consequently both scales can be used with both operators (cf. Remark 5.7).

Remark 5.7. For all our purposes the smoothing condition (5.3) can be relaxed to (5.10)-(5.11). This relaxation covers the situation where there exists an operator \mathcal{A}_0 that satisfies (5.3) and is a ‘version’ of \mathcal{A} in that the two operators possess a common inverse, such as when \mathcal{A} and \mathcal{A}_0 are defined to solve a differential equation with different boundary conditions. Lemma 5.9 shows that the relaxed version of the smoothing condition is then satisfied by the map $f \mapsto [\mathcal{A}f]$ of f in the class of $\mathcal{A}f$ in the quotient space $G/R(\mathcal{A} - \mathcal{A}_0)$.

5.3 Galerkin Projection

In this section we collect some (well known) results on the Galerkin method.

Consider a scale of smoothness classes $(H_s)_{s \in \mathbb{R}}$ as in Definition 2.1. Let $\mathcal{A} : H \rightarrow G$ be an injective bounded operator between separable Hilbert spaces, and let V_j be a finite-dimensional subspace of H . The Galerkin solution $f^{(j)} \in V_j$ to the image $\mathcal{A}f$ of an element f is defined as the element in V_j such that $\mathcal{A}f^{(j)}$ is equal to the orthogonal projection of $\mathcal{A}f$ onto the image space $W_j = \mathcal{A}V_j$. Thus, if $Q_j : G \rightarrow W_j$ denotes the orthogonal projection onto W_j , then the Galerkin solution can be written as

$$f^{(j)} = R_j \mathcal{A}f, \quad \text{for} \quad R_j = \mathcal{A}^{-1}Q_j,$$

where the inverse \mathcal{A}^{-1} is well defined on the linear subspace W_j .

If the operators $R_j \mathcal{A}$ are uniformly bounded with respect to j , then the convergence rate $\|f^{(j)} - f\|_0$ of the Galerkin solution to f is known to be of the same order as the distance $\|P_j f - f\|_0$ of f to its projection on V_j . (See Section 3.2 and Theorem 3.7 in [57], or the proof below.) In particular, if $f \in H_s$ and V_j satisfies (2.2), then the convergence rate is given by $\delta(j, s)$.

In order to control the stochastic noise term ξ in the observation schemes (5.1) and (5.2), it is necessary also to control the norms of the operators R_j . The following lemma summarizes the properties of the Galerkin projection needed in the proof of our main result.

Lemma 5.8. *If V_j is a finite-dimensional space as in Assumption 2.3 such that (2.2) and (2.3) hold, and $\mathcal{A} : H_0 \rightarrow G$ is a bounded linear operator satisfying $\|\mathcal{A}f\|_0 \simeq \|f\|_{-\gamma}$ for every $f \in H_0$, then the norms of the operators $R_j : G \rightarrow H_0$ and $R_j \mathcal{A} : H_0 \rightarrow H_0$ satisfy*

$$\|R_j\| \lesssim_{\mathcal{A}} \frac{1}{\delta(j, \gamma)}, \quad (5.7)$$

$$\|R_j \mathcal{A}\| \lesssim_{\mathcal{A}} 1. \quad (5.8)$$

Furthermore, for $f \in H_s$ the Galerkin solution $f^{(j)} \in V_j$ to $\mathcal{A}f$ satisfies

$$\|f^{(j)} - f\|_0 \lesssim_{\mathcal{A}} \delta(j, s) \|f\|_s. \quad (5.9)$$

Proof. For $g \in G$ we have $R_j g \in V_j$ and hence by (2.5),

$$\|R_j g\|_0 \lesssim \frac{1}{\delta(j, \gamma)} \|R_j g\|_{-\gamma} \simeq \frac{1}{\delta(j, \gamma)} \|\mathcal{A}R_j g\|_0 = \frac{1}{\delta(j, \gamma)} \|Q_j g\|_0,$$

since $\mathcal{A}R_j g = Q_j g$. Because $\|Q_j g\|_0 \leq \|g\|_0$, we conclude that $\|R_j\| \lesssim 1/\delta(j, \gamma)$.

By definition $f^{(j)} = R_j \mathcal{A}f$, and $R_j \mathcal{A}$ acts as the identity on V_j . Therefore $f^{(j)} - P_j f = R_j \mathcal{A}(f - P_j f)$, and hence

$$\|f^{(j)} - P_j f\|_0 \leq \|R_j\| \|\mathcal{A}(f - P_j f)\|_0 \simeq \|R_j\| \|f - P_j f\|_{-\gamma} \leq \|R_j\| \delta(j, \gamma) \|f\|_0,$$

by (2.4). By the preceding paragraph $\|R_j\| \delta(j, \gamma) \lesssim 1$, so that the right side is bounded above by $\|f\|_0$. By the triangle inequality

$$\|R_j \mathcal{A}f\|_0 = \|f^{(j)}\|_0 \leq \|f^{(j)} - P_j f\|_0 + \|P_j f - f\|_0 \lesssim \|f\|_0,$$

in view of the preceding display and the fact that $\|P_j f - f\|_0 \leq \|f\|_0$. This shows that $\|R_j \mathcal{A}\| \lesssim 1$.

Finally, since $f^{(j)} - f = (R_j \mathcal{A} - I)(f - P_j f)$, we have that

$$\|f^{(j)} - f\|_0 = \|(R_j \mathcal{A} - I)(f - P_j f)\|_0 \leq (\|R_j \mathcal{A}\| + 1) \|f - P_j f\|_0.$$

Inequality (5.9) follows by the boundedness of $\|R_j \mathcal{A}\|$ and (2.2). \square

As is clear from the proof, the smoothing assumption $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$ can be relaxed to the pair of inequalities

$$\|\mathcal{A}f\| \lesssim \|f\|_{-\gamma}, \quad f \perp V_j, \quad (5.10)$$

$$\|\mathcal{A}f\| \gtrsim \|f\|_{-\gamma}, \quad f \in R(R_j). \quad (5.11)$$

This helps to cover cases in which the smoothing condition is satisfied for a modification of the operator \mathcal{A} , but not \mathcal{A} itself, for example a modification taking different boundary conditions of a differential operator into account.

We introduce a *modified Galerkin solution* to $\mathcal{A}f$ to cover such a case. Let $\mathcal{A}_0, \mathcal{A} : H \rightarrow G$ be injective bounded operators between separable Hilbert spaces that possess a common inverse in the sense of existence of a linear map $B : D(B) \subset G \rightarrow H$ with domain $D(B)$ containing the linear span of the ranges of \mathcal{A}_0 and \mathcal{A} such that $B\mathcal{A}_0 = I = B\mathcal{A}$. For simplicity of notation, write $B = \mathcal{A}^- = \mathcal{A}_0^-$. Intuitively, for the inverse problem, taking $\mathcal{A}_0 f$ or $\mathcal{A}f$ as input data should be equivalent. However, it may be that \mathcal{A}_0 is smoothing in a given scale $(H_s)_{s \in \mathbb{R}}$, whereas \mathcal{A} is not. In that case we reconstruct as follows. Assume that $\Phi = \mathcal{A} - \mathcal{A}_0$ has closed range, and let $P_\Phi : G \rightarrow G$ be the orthogonal projection onto this range. Now let $Q_j : G \rightarrow G$ be the orthogonal projection onto the finite-dimensional space $(I - P_\Phi)AV_j$, and set

$$f^{(j)} = R_j \mathcal{A}f, \quad \text{for} \quad R_j = \mathcal{A}^- Q_j (I - P_\Phi). \quad (5.12)$$

Thus after removing the ‘‘irrelevant part’’ of $\mathcal{A}f$ that does not influence the inversion, we project onto the finite-dimensional space $(I - P_\Phi)AV_j$ of similarly cleaned functions $\mathcal{A}f$ with $f \in V_j$, and finally invert.

Lemma 5.9. *If V_j is a finite-dimensional space as in Assumption 2.3 such that (2.2) and (2.3) hold, and $\mathcal{A}_0, \mathcal{A} : H_0 \rightarrow G$ are bounded linear operators with common inverse satisfying $\|\mathcal{A}_0 f\| \simeq \|f\|_{-\gamma}$ for every $f \in H_0$, then the operators $R_j : G \rightarrow H_0$ and $R_j \mathcal{A} : H_0 \rightarrow H_0$ and $f^{(j)} = R_j \mathcal{A} f$ as in (5.12) satisfy (5.7), (5.8), and (5.9).*

Proof. The operator $[\mathcal{A}] : H \rightarrow G/\Phi(H)$ mapping $f \in H$ into the class of $\mathcal{A}f$ in the quotient space $G/\Phi(H)$ is one-to-one, since $[\mathcal{A}f] = 0$ implies $\mathcal{A}f \in R(\Phi)$ and hence $f = B\mathcal{A}f = 0$, since $B\Phi = 0$. Identifying $[g] \in \tilde{G} := G/\Phi(H)$ with the function $(I - P_\Phi)g$ with norm $\|[g]\|_{\tilde{G}} = \|(I - P_\Phi)g\|_G$, we see that $R_j \mathcal{A} f$ as in (5.12) is actually the Galerkin solution to $[\mathcal{A}]f$. It suffices to show that $[\mathcal{A}] : H \rightarrow \tilde{G}$ is smoothing in the sense of (5.10). Now $\|[\mathcal{A}f]\|_{\tilde{G}} = \|(I - P_\Phi)\mathcal{A}_0 f\|_G \leq \|\mathcal{A}_0 f\|_G \simeq \|f\|_{-\gamma}$, for every $f \in H$. Furthermore, for every f such that $\mathcal{A}_0 f \perp R(\Phi)$, the inequality is an equality. This is true for $f = R_j g$, since $\mathcal{A}_0 R_j g = Q_j(I - P_\Phi)g \in (I - P_\Phi)\mathcal{A}V_j$. \square

5.4 Notes

Statistical Inverse Problems

The study of statistical (nonparametric) linear inverse problems was initiated by Wahba in 1970s in [101]. The 1990s paper [24] used wavelet shrinkage methods, while around 2000, the authors of [17] investigated (5.1) in the linear partial differential equations setting, while a systematic study of Gaussian sequence models was presented in [16]. A review of work until 2008 is given in [15]. The connection of regularization methods to the Bayesian approach was recognized early on. However, the study of the recovery properties of posterior distributions was started only in [59, 60]. A review of the Bayesian approach to inverse problems, with many examples, is given in [88].

Much of the existing work on statistical inverse problems is based on the singular value decomposition (SVD) of the operator \mathcal{A} ; see, e.g., [15]. When \mathcal{A} is compact, the operator $\mathcal{A}^* \mathcal{A}$, where \mathcal{A}^* is the adjoint of \mathcal{A} , can be diagonalized with respect to an orthonormal *eigenbasis*, with eigenvalues tending to zero. The observation $Y^{(n)}$ can then be reduced to noisy observations on the Fourier coefficients of $\mathcal{A}f$ in the eigenbasis, which are multiples of the Fourier coefficients of f , and the problem is to recover the latter. In the frequentist setup thresholding or other regularization methods can be applied to reduce the weight of estimates on coefficients corresponding to smaller eigenvalues, in which the noise will overpower the signal. In the Bayesian setup one may design a prior by letting the Fourier coefficients be (independent) random variables, with smaller variances for smaller eigenvalues. These singular value methods have several disadvantages, as pointed out in [19, 24]. First, the eigenbasis functions might not be easy to compute. Second, and more importantly, these functions are directly linked to the operator \mathcal{A} , and need not be related to the function space (smoothness class) that is thought to contain the true signal f . Consequently, the parameter of interest f may not have a simple, parsimonious representation in the eigenbasis expansion, see [24]. Furthermore, it is logical to consider the series expansion of the signal f in other bases than the eigenbasis, for instance, in the situation that one can only measure

noisy coefficients of the signal f in a given basis expansion, due to a particular experimental setup. See [37, 70] for further discussion.

Deterministic Inverse Problems

There is a rich literature on inverse problems. The case that the noise ξ is a bounded *deterministic* perturbation, has been particularly well studied, and various general procedures and methods to estimate the convergence rates of regularized solutions have been proposed. See the monographs [29, 57]. The case of stochastic noise is less studied, but is receiving increasing attention.

Inverse Problems in Scales and Bayesian Approach

A canonical example are Sobolev spaces, with the operator \mathcal{A} being an integral operator. This Sobolev space setup with wavelet basis was investigated in [19, 24]. In deterministic inverse problems, a more general setup, considering \mathcal{A} that acts along nested Hilbert spaces, *Hilbert scales*, was initiated by Natterer in [71] and further developed in, amongst others, [46, 69, 70]. In the Bayesian context Hilbert scales were used in [30], under the assumption that the noise ξ is a proper Gaussian element in G , and in [1], but under rather intricate assumptions. The second question, to allow priors that are not conjugate, can also be answered under the condition of \mathcal{A} . In the linear inverse problem Gaussian priors are easy, as they lead to Gaussian posterior distributions, which can be studied by direct means. Most of the results on Bayesian inverse problems fall in this framework [1, 30, 59, 60], exceptions being [80] and [58].

Generalized Random Elements

An alternative method to give a rigorous interpretation to white noise ξ , is to embed G into a bigger space in which ξ can be realized as a Borel measurable map, or to think of ξ as a cylindrical process. See e.g., [87]. For G a set of functions on an interval, one can also realize ξ as a stochastic integral relative to Brownian motion, which takes its values in the ‘abstract Wiener space’ attached to G .

Chapter 6

Inverse Problems with Continuous Observations in Smoothness Scales

6.1 Introduction

In this chapter, we study the following linear inverse problem in white noise,

$$Y^{(n)} = \mathcal{A}f + \frac{1}{\sqrt{n}}\xi, \quad (6.1)$$

where \mathcal{A} is characterised in Section 5.2, and ξ is an isonormal Gaussian process (see Definition 3.28 and Section 4.2.1).

This chapter is organised as follows. We present a general posterior contraction theorem in Section 6.2, which is based on a testing approach using the estimator from Section 5.3. The result does not depend on any conjugacy of priors. Then, the general theorem is applied to two examples of priors: series priors in Section 6.3 and Gaussian priors in Section 6.4. Since the simple Gaussian prior is not fully adaptive, we introduce Gaussian mixture priors to obtain adaptation in Section 6.5. It is noteworthy that the priors are defined in terms of the scale, rather than the operator. In other words, the operator and the prior are assumed related, but only indirectly, through the scale. In this arrangement, priors can be chosen directly related to common bases (e.g., splines or wavelets bases) and function spaces, rather than to the operator through its singular value decomposition. Section 6.6 contains the proofs. We conclude this chapter with the discussion of several extensions of the present work in Section 6.7.

6.2 General Result

In this section we present a general theorem on posterior contraction. We form the posterior distribution $\Pi_n(\cdot | Y^{(n)})$ as in (4.1), given a prior Π on the space $H = H_0$ and an observation $Y^{(n)}$, whose conditional distribution given f is determined by the model (6.1). We study this random distribution under the assumption that $Y^{(n)}$ follows the model (6.1) for a given ‘true’ function $f = f_0$, which we assume to be an element of H_β in a given smoothness scale $(H_s)_{s \in \mathbb{R}}$, as in Definition 2.1.

The result is based on an extension of the testing approach of [35] to the inverse problem (6.1). The inverse problem is handled with the help of the Galerkin method, which is a well known strategy in numerical analysis to solve the operator equation $y = \mathcal{A}f$ for f , in particular for differential and integral operators. The Galerkin method has several variants, which are useful depending on the properties of the operator involved. Here we use the least squares method, which is of general application; for other variants and background, see e.g., [57]. In Section 5.3 we have given a self-contained derivation of the necessary inequalities, exactly in our framework. We note that the Galerkin method only appears as a tool to state and derive a posterior contraction rate. In our context it does not enter into the solution of the inverse problem, which is achieved through the Bayesian method.

Let $W_j = \mathcal{A}V_j \subset G$ be the image under \mathcal{A} of a finite-dimensional approximation space V_j linked to the smoothness scale $(H_s)_{s \in \mathbb{R}}$ as in Assumption 2.3, and let $Q_j : G \rightarrow W_j$ be the orthogonal projection onto W_j . If $\mathcal{A} : H \rightarrow G$ is injective, then \mathcal{A} is a bijection between the finite-dimensional vector spaces V_j and W_j , and hence for every $f \in H$ there exists $f^{(j)} \in V_j$ such that $\mathcal{A}f^{(j)} = Q_j \mathcal{A}f$. The element $f^{(j)}$ is called the *Galerkin solution* to $\mathcal{A}f$ in V_j . By the projection theorem in Hilbert spaces it is characterized by the property that $f^{(j)} \in V_j$ together with the orthogonality relations

$$\langle \mathcal{A}f^{(j)}, w \rangle_0 = \langle \mathcal{A}f, w \rangle_0, \quad w \in W_j. \quad (6.2)$$

The idea of the Galerkin inversion is to project the (complex) object $\mathcal{A}f$ onto the finite-dimensional space W_j , and next find the inverse image $f^{(j)}$ of the projection, in the finite-dimensional space V_j , as in the diagram:

$$\begin{array}{ccc} H_0 \ni f & \xrightarrow{\mathcal{A}} & \mathcal{A}f \in G \\ & & \downarrow Q_j \\ V_j \ni f^{(j)} & \xleftarrow{\mathcal{A}^{-1}} & Q_j \mathcal{A}f \in W_j \end{array}$$

Clearly the Galerkin solution to an element $f \in V_j$ is f itself, but in general $f^{(j)}$ is an approximation to f , which will be better for increasing j , but increasingly complex. The following theorem uses a dimension $j = j_n$ that balances approximation to complexity, where the complexity is implicitly determined by a testing criterion.

Theorem 6.1. *For smoothness classes $(H_s)_{s \in \mathbb{R}}$ as in Definition 2.1, assume that $\|\mathcal{A}f\|_0 \simeq \|f\|_{-\gamma}$ for some $\gamma > 0$, and let $f^{(j)}$ denote the Galerkin solution to $\mathcal{A}f$ relative to linear subspaces V_j associated to $(H_s)_{s \in \mathbb{R}}$ as in Assumption 2.3. Let $f_0 \in H_\beta$ for some $\beta \in (0, S)$, and for $\eta_n \geq \varepsilon_n \downarrow 0$ such that $n\varepsilon_n^2 \rightarrow \infty$, and $j_n \in \mathbb{N}$*

such that $j_n \rightarrow \infty$, and some $c > 0$, assume

$$j_n \leq cn\varepsilon_n^2, \quad (6.3)$$

$$\eta_n \geq \frac{\varepsilon_n}{\delta(j_n, \gamma)}, \quad (6.4)$$

$$\eta_n \geq \delta(j_n, \beta). \quad (6.5)$$

Consider prior probability distributions Π on H_0 satisfying

$$\Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon_n) \geq e^{-n\varepsilon_n^2}, \quad (6.6)$$

$$\Pi(f : \|f^{(j_n)} - f\|_0 > \eta_n) \leq e^{-4n\varepsilon_n^2}. \quad (6.7)$$

Then the posterior distribution in the model (6.1) contracts at the rate η_n at f_0 , i.e. for a sufficiently large constant M we have $\Pi_n(f : \|f - f_0\|_0 > M\eta_n \mid Y^{(n)}) \rightarrow 0$, in probability under the law of $Y^{(n)}$ given by (6.1) with $f = f_0$.

Proof. The Kullback-Leibler divergence and variation between the distributions of $Y^{(n)}$ under two functions f and f_0 are given by $n\|\mathcal{A}f - \mathcal{A}f_0\|^2/2$ and twice this quantity, respectively. (E.g., Lemma 8.30 in [35].) Therefore the neighbourhoods $B_{n,2}(f_0, \varepsilon)$ in (8.19) of [35] contain the ball $\{f \in H_0 : \|\mathcal{A}f - \mathcal{A}f_0\| \leq \varepsilon\}$. By assumption (6.6) this has prior mass at least $e^{-n\varepsilon_n^2}$.

Because the quotient of the left sides of (6.6) and (6.7) is $o(e^{-2n\varepsilon_n^2})$, the posterior probability of the set $\{f : \|f^{(j_n)} - f\|_0 > \eta_n\}$ tends to zero, by Theorem 8.20 in [35].

By a variation of Theorem 8.22 in [35] it is now sufficient to show the existence of tests τ_n such that, for some $M > 0$,

$$P_{f_0}^{(n)} \tau_n \rightarrow 0, \quad \sup_{\substack{f: \|f - f_0\|_0 > M\eta_n, \\ \|f^{(j_n)} - f\|_0 \leq \eta_n}} P_f^{(n)}(1 - \tau_n) \leq e^{-4n\varepsilon_n^2}.$$

Indeed, in the case that the prior mass condition (8.20) in Theorem 8.22 of [35] can be strengthened to (8.22), as is the case in our setup in view of (6.6), it suffices to verify (8.24) only for a single value of j . Furthermore, we can apply Theorem 8.22 with the metrics $d_n(f, g) = \|f - g\|_0 \varepsilon_n / \eta_n$ in order to reduce the restriction $d_n(\theta, \theta_{n,0}) > M\varepsilon_n$ to $\|f - f_0\|_0 > M\eta_n$.

Fix any orthonormal basis $(\bar{\psi}_i)_{i < j}$ of $W_j = \mathcal{A}V_j$ and define

$$\begin{aligned} \bar{Y}_j &= \sum_{i < j} Y_{\bar{\psi}_i}^{(n)} \bar{\psi}_i = \sum_{i < j} \langle \mathcal{A}f, \bar{\psi}_i \rangle \bar{\psi}_i + \frac{1}{\sqrt{n}} \sum_{i < j} \xi_{\bar{\psi}_i} \bar{\psi}_i \\ &= Q_j \mathcal{A}f + \frac{1}{\sqrt{n}} \bar{\xi}_j, \end{aligned}$$

where $\bar{\xi}_j := \sum_{i < j} \xi_{\bar{\psi}_i} \bar{\psi}_i$. The latter is a “standard normal vector in the finite-dimensional space W_j ”: because $(\xi_{\bar{\psi}_i})_{i < j}$ are i.i.d. standard normal variables, the variable $\langle \bar{\xi}_j, w \rangle = \sum_{i < j} \xi_{\bar{\psi}_i} \langle \bar{\psi}_i, w \rangle$ is $N(0, \|Q_j w\|^2)$ -distributed, for every $w \in G$.

Let the operator $R_j : G \mapsto V_j$ be defined as $R_j = \mathcal{A}^{-1}Q_j$, where \mathcal{A}^{-1} is the inverse of \mathcal{A} , which is well defined on the range $W_j = \mathcal{A}V_j$ of Q_j . Then by

definition $R_j \mathcal{A}f$ is equal to the Galerkin solution $f^{(j)}$ to $\mathcal{A}f$. By the preceding display $R_j \bar{Y}_j$ is a well-defined Gaussian random element in V_j , satisfying

$$R_j \bar{Y}_j = f^{(j)} + \frac{1}{\sqrt{n}} R_j \bar{\xi}_j. \quad (6.8)$$

The variable $R_j \bar{\xi}_j$ is a Gaussian random element in V_j with strong and weak second moments

$$\begin{aligned} \mathbb{E} \|R_j \bar{\xi}_j\|_0^2 &\leq \|R_j\|^2 \mathbb{E} \|\bar{\xi}_j\|^2 = \|R_j\|^2 \mathbb{E} \sum_{i < j} \xi_{\psi_i}^2 = \|R_j\|^2 (j-1) \lesssim \frac{j}{\delta(j, \gamma)^2}, \\ \sup_{\|f\|_0 \leq 1} \mathbb{E} \langle R_j \bar{\xi}_j, f \rangle_0^2 &= \sup_{\|f\|_0 \leq 1} \mathbb{E} \langle \bar{\xi}_j, R_j^* f \rangle^2 = \sup_{\|f\|_0 \leq 1} \|Q_j R_j^* f\|^2 \leq \|R_j^*\|^2 \lesssim \frac{1}{\delta(j, \gamma)^2}. \end{aligned}$$

In both cases the inequality on $\|R_j\| = \|R_j^*\|$ at the far right side follows from (5.7).

The first inequality shows that the first moment $\mathbb{E} \|R_j \bar{\xi}_j\|_0$ of the variable $\|R_j \bar{\xi}_j\|_0$ is bounded above by $\sqrt{j}/\delta(j, \gamma)$. By Borell's inequality (e.g. Lemma 3.1 in [67] and subsequent discussion), applied to the Gaussian random variable $R_j \bar{\xi}_j$ in H_0 , we see that there exist positive constants a and b such that, for every $t > 0$,

$$\Pr \left(\|R_j \bar{\xi}_j\|_0 > t + a \frac{\sqrt{j}}{\delta(j, \gamma)} \right) \leq e^{-bt^2 \delta(j, \gamma)^2}.$$

For $t = 2\sqrt{n}\eta_n/\sqrt{b}$ and η_n , ε_n and j_n satisfying (6.3), (6.4) and (6.5) this yields, for some $a_1 > 0$,

$$\Pr \left(\|R_{j_n} \bar{\xi}_{j_n}\|_0 > a_1 \sqrt{n}\eta_n \right) \leq e^{-4n\varepsilon_n^2}. \quad (6.9)$$

We apply this to bound the error probabilities of the tests

$$\tau_n = 1 \{ \|R_{j_n} \bar{Y}_{j_n} - f_0\|_0 \geq M_0 \eta_n \}, \quad (6.10)$$

where M_0 is a given constant, to be determined.

Under f_0 , the decomposition (6.8) is valid with $f = f_0$, and hence $R_j \bar{Y}_j - f_0 = n^{-1/2} R_j \bar{\xi}_j + f_0^{(j)} - f_0$. By the triangle inequality it follows that $\tau_n = 1$ implies that $n^{-1/2} \|R_{j_n} \bar{\xi}_{j_n}\|_0 \geq M_0 \eta_n - \|f_0^{(j_n)} - f_0\|_0$. By (5.9) the assumption that $f_0 \in H_\beta$ implies that $\|f_0^{(j_n)} - f_0\|_0 \leq M_1 \delta(j_n, \beta)$, for some M_1 , which at $j = j_n$ is further bounded by $M_1 \eta_n$, by assumption (6.5). Hence the probability of an error of the first kind satisfies

$$P_{f_0}^{(n)} \tau_n \leq \Pr \left(\frac{1}{\sqrt{n}} \|R_{j_n} \bar{\xi}_{j_n}\|_0 \geq (M_0 - M_1) \eta_n \right).$$

For $M_0 - M_1 > a_1$, the right side is bounded by $e^{-4n\varepsilon_n^2}$, by (6.9).

Under f the decomposition (6.8) gives that $R_j \bar{Y}_j - f_0 = n^{-1/2} R_j \bar{\xi}_j + f^{(j)} - f_0$. By the triangle inequality $\tau_n = 0$ implies that $n^{-1/2} \|R_{j_n} \bar{\xi}_{j_n}\|_0 \geq \|f^{(j_n)} - f_0\|_0 - M_0 \eta_n$. For f such that $\|f - f_0\|_0 > M \eta_n$ and $\|f - f^{(j_n)}\|_0 \leq \eta_n$, we have $\|f^{(j_n)} - f_0\|_0 \geq (M - 1) \eta_n$. Hence the probability of an error of the second kind satisfies

$$P_f^{(n)} (1 - \tau_n) \leq \Pr \left(\frac{1}{\sqrt{n}} \|R_{j_n} \bar{\xi}_{j_n}\|_0 \geq (M - 1 - M_0) \eta_n \right),$$

For $M - 1 - M_0 > a_1$, this is bounded by $e^{-4n\varepsilon_n^2}$, by (6.9).

We can first choose M_0 large enough so that $M_0 - M_1 > a_1$, and next M large enough so that $M - 1 - M_0 > a_1$, to finish the proof. \square

Inequality (6.6) is the usual *prior mass condition* for the ‘direct problem’ of estimating $\mathcal{A}f$ (see [33]). It determines the rate of contraction ε_n of the posterior distribution of $\mathcal{A}f$ to $\mathcal{A}f_0$. The rate of contraction η_n of the posterior distribution of f is slower due to the necessity of (implicitly) inverting the operator \mathcal{A} . The theorem shows that the rate η_n depends on the combination of the prior, through (6.7), and the inverse problem, through the various approximation rates.

Remark 6.2. It would be possible to obtain the theorem as a corollary of Theorem 2.1 in [58]. We would take the sets \mathcal{S}_n in the latter high-level result equal to the sets $\{f : \|f^{(j_n)} - f\|_0 > \eta_n\}$ appearing in (6.7). To verify the conditions of [58] for this choice, most of the preceding proof would be needed. Since the next theorem appears not to be a consequence of this approach, and its proof uses the preceding proof, we have given a direct proof instead.

The theorem applies to a true function f_0 that is ‘smooth’ of order β (i.e., $f_0 \in H_\beta$). For a prior that is constructed to give an optimal contraction rate for multiple values of β simultaneously, the theorem may not give the best result. The following theorem refines Theorem 6.1 by considering a mixture prior of the form

$$\Pi = \int \Pi_\tau dQ(\tau), \quad (6.11)$$

where Π_τ is a prior on H , for every given ‘hyperparameter’ τ running through some measurable space, and Q is a prior on this hyperparameter. The idea is to *adapt* the prior to multiple smoothness levels through the hyperparameter τ .

Theorem 6.3. *Consider the setup and assumptions of Theorem 6.1 with a prior of the form (6.11). Assume that (6.3), (6.4), (6.5) and (6.6) hold, but replace (6.7) by the pair of conditions, for numbers $\eta_{n,\tau}$ and $C > 0$ and every τ ,*

$$\Pi_\tau(f : \|f - f_0\|_0 < 2\eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2}, \quad \forall \tau \text{ with } \eta_{n,\tau} \geq C\eta_n, \quad (6.12)$$

$$\Pi_\tau(f : \|f^{(j_n)} - f\|_0 > \eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2}. \quad (6.13)$$

Then the posterior distribution in the model (6.1) contracts at the rate η_n at f_0 , i.e. for a sufficiently large constant M we have $\Pi_n(f : \|f - f_0\|_0 > M\eta_n \mid Y^{(n)}) \rightarrow 0$, in probability under the law of $Y^{(n)}$ given by (6.1) with $f = f_0$.

Proof. We take the parameter of the model as the pair (f, τ) , which receives the joint prior given by $f \mid \tau \sim \Pi_\tau$ and $\tau \sim Q$. With abuse of notation, we denote this prior also by Π . The likelihood still depends on f only, but the joint prior gives rise to a posterior distribution on the pair (f, τ) , which we also denote by $\Pi_n(\cdot \mid Y^{(n)})$, by a similar abuse of notation.

By (6.11) and eqs. (6.12) and (6.13),

$$\begin{aligned} \Pi((f, \tau) : \eta_{n,\tau} \geq C\eta_n, \|f - f_0\|_0 < 2\eta_{n,\tau}) &\leq e^{-4n\varepsilon_n^2}, \\ \Pi((f, \tau) : \|f^{(j_n)} - f\|_0 > \eta_{n,\tau}) &\leq e^{-4n\varepsilon_n^2}. \end{aligned}$$

In view of (6.6) and Theorem 8.20 in [35], the posterior probabilities of the two sets in the left sides tend to zero. As in the proof of Theorem 6.1, we can apply a variation of Theorem 8.22 in [35] to see that it is now sufficient to show the existence of tests τ_n such that, for some $M \geq 2C$,

$$P_{f_0}^{(n)} \tau_n \rightarrow 0, \quad \sup_{\substack{(f, \tau): \|f - f_0\|_0 > M\eta_n \vee 2\eta_{n, \tau}, \\ \|f^{(j_n)} - f\|_0 \leq \eta_{n, \tau}}} P_f^{(n)} (1 - \tau_n) \leq e^{-4n\varepsilon_n^2}.$$

(Note that $M\eta_n \vee 2\eta_{n, \tau} = M\eta_n$ if $\eta_{n, \tau} < C\eta_n$ and $M \geq 2C$.) We use the tests defined in (6.10), as in the proof of Theorem 6.1. The latter proof shows that the tests are consistent. We adapt the bound on the power, as follows.

By the triangle inequality $\tau_n = 0$ implies that, for (f, τ) with $\|f - f_0\|_0 > M\eta_n \vee 2\eta_{n, \tau}$ and $\|f^{(j_n)} - f\|_0 \leq \eta_{n, \tau}$,

$$\begin{aligned} n^{-1/2} \|R_{j_n} \bar{\xi}_{j_n}\|_0 &\geq \|f^{(j_n)} - f_0\|_0 - M_0 \eta_n \geq \|f - f_0\|_0 - \|f^{(j_n)} - f\|_0 - M_0 \eta_n \\ &\geq M\eta_n \vee 2\eta_{n, \tau} - \eta_{n, \tau} - M_0 \eta_n \geq (M/2 - M_0) \eta_n. \end{aligned}$$

Hence by (6.9) the probability of an error of the second kind is bounded by $e^{-4n\varepsilon_n^2}$, for M sufficiently large that $M/2 - M_0 > a_1$. \square

In a typical application of the preceding theorem the priors Π_τ for τ such that $\eta_{n, \tau} \geq C\eta_n$ will be the priors on ‘rough’ functions, with ‘intrinsic’ contraction rate $\eta_{n, \tau}$ slower than η_n . These ‘bad’ priors do not destroy the overall contraction rate, because they put little mass near the true function f_0 , by condition (6.12). It is necessary to address these priors explicitly in the conditions, because they will typically fail the approximation condition (6.7), which must be relaxed to (6.13). A further generalization might be to allow the truncation levels j_n to depend on τ , but this will not be needed for our examples.

Inspection of the proof shows that the posterior probability of the sets $\{\tau : \eta_{n, \tau} \gtrsim C\eta_n\}$ tends to zero. This means that the posterior correctly disposes of the models that are ‘too rough’, for the given true function f_0 . In general there is no similar protection against models that are too smooth, but this does not affect the contraction rate.

6.3 Random Series Priors

Suppose that $\{\phi_i\}_{i \in \mathbb{N}}$ is an orthonormal basis of $H = H_0$ that gives optimal approximation relative to the scale of smoothness classes $(H_s)_{s \in \mathbb{R}}$ in the sense that the linear spaces $V_j = \text{Span}\{\phi_i\}_{i < j}$ satisfy Assumption 2.3. Consider a prior defined as the law of the random series

$$f = \sum_{i=1}^M f_i \phi_i, \quad (6.14)$$

where M is a random variable in \mathbb{N} independent from the independent random variables f_1, f_2, \dots in \mathbb{R} .

Condition 6.4 (Random series prior). (i) The probability density function p_M of M satisfies, for some positive constants b_1, b_2 ,

$$e^{-b_1 k} \lesssim p_M(k) \lesssim e^{-b_2 k}, \quad \forall k \in \mathbb{N}.$$

(ii) The variable f_i has density $p(\cdot/\kappa_i)/\kappa_i$, for a given probability density p on \mathbb{R} and a constant $\kappa_i > 0$ such that, for some $C > 0$ and $w > 0$, $\alpha, \beta_0 > 0$,

$$p(x) \gtrsim e^{-C|x|^w}, \quad (6.15)$$

$$i^{-\beta_0/d} (\log i)^{-1/w} \lesssim \kappa_i \lesssim i^\alpha. \quad (6.16)$$

Priors of this type were studied in [4, 80], and applied to inverse problems in the SVD framework in [80] (see Section 3.1 of the latter paper for discussion). For Gaussian variables f_j and degenerate M the series (6.14) is a Gaussian process, and has been more widely studied, but we focus here on the non-Gaussian case. Since the basis $(\phi_i)_{i \in \mathbb{N}}$ used in the prior is linked to the smoothness class $(H_s)_{s \in \mathbb{R}}$, rather than to the operator \mathcal{A} , the prior is not restricted to the SVD framework. Of course, in the theorem below we do require the operator to be smoothing in the same smoothness scale, thus maintaining a link between prior and operator.

The assumption on the density p_M is mild and is satisfied, for instance, by the Poisson distribution. The assumption on the density p is mild as well, and is satisfied by many distributions with full support in \mathbb{R} , including the Gaussian and Laplace distributions. The parameter β_0 in (6.16) must be a lower bound on the smoothness of the true parameter f_0 . Apart from this, condition (6.16) is also very mild, and allows the scale parameters κ_i to tend both to zero or to infinity.

The preceding random series prior is not conjugate to the inverse problem (6.1). In general the resulting posterior distribution will not have a closed form expression, but must be computed using simulation, such as Markov chain Monte Carlo, or approximated using an optimisation method, such as variational approximation. However, the contraction rate of the posterior distribution can be established without the help of an explicit expression for the posterior distribution, as shown in the following theorem.

Theorem 6.5 (Random Series Prior). *Let $(\phi_i)_{i \in \mathbb{N}}$ be an orthonormal basis of H_0 such that the spaces $V_j = \text{Span}\{\phi_i\}_{i < j}$ satisfy Assumption 2.3 with $\delta(j, s) = j^{-s/d}$ relative to smoothness classes $(H_s)_{s \in \mathbb{R}}$ as in Definition 2.1. Assume that $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$ for some $\gamma > 0$, and let $f_0 \in H_\beta$ for some $\beta \in (0, S)$. Then, for the random series prior defined in (6.14) and satisfying Condition 6.4 with $\beta_0 \leq \beta$, and sufficiently large $\underline{M} > 0$, for $\tau = (\beta + \gamma)(1 + 2\gamma/d)/(2\beta + 2\gamma + d)$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > \underline{M} n^{-\beta/(2\beta+2\gamma+d)} (\log n)^\tau \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

The rate $n^{-\beta/(2\beta+2\gamma+d)}$ is known to be the minimax rate of estimation of a β -regular function on a d -dimensional domain, in an inverse problem with inverse parameter γ (see, e.g., [19]). The assumption that $\delta(j, s) = j^{-s/d}$ places the setup of the theorem in this setting, and hence the rate of contraction obtained in the preceding theorem is the minimax rate up to a logarithmic factor. The rate is

adaptive to the regularity of β of the true parameter, which is not used in the construction of the prior, apart from the assumption that $\beta \geq \beta_0$. (See [34] and Chapter 10 in [35] for general discussion of adaptation in the Bayesian sense.)

The proof of the theorem is deferred to Section 6.6; it will be based on Theorem 6.1.

Example 6.6 (Wavelet basis). Let p be a standard normal density, p_M a standard Poisson probability mass function, and set the scaling parameters κ_i equal to 1 (no scaling).

Consider an S -regular orthonormal wavelet basis $\{\phi_{j,k}\}$ for the space of square-integrable functions on the d -dimensional torus $(0, 2\pi]^d$. We can renumber the index (j, k) into \mathbb{N} by ordering the basis functions by their multiresolution levels, $2^{jd} + k$, and next construct the random series prior (6.14).

An S -regular orthonormal wavelet basis is known to correspond to the scale of Sobolev spaces up to smoothness level S . Therefore, by Theorem 6.5, the contraction rate of the posterior distribution is $n^{-\beta/(2\beta+2\gamma+d)}$ times a logarithmic factor whenever the operator is smoothing relative to the Sobolev scale and the true function f_0 belongs to the Sobolev space of order β , for $\beta_0 \leq \beta < S$. Thus the posterior distributions are adaptive up to a logarithmic factor to the scale of Sobolev spaces of orders between β_0 and S .

For increasing $\beta \geq S$ the rate given by the theorem still improves. However, the ‘regularity’ β defined by the scale $(H_s)_{s \in \mathbb{R}}$ may then not coincide with the Sobolev scale.

6.4 Gaussian Priors

If the function f in (6.1) is equipped with a Gaussian prior, then the corresponding posterior distribution will be Gaussian as well. Furthermore, the posterior mean will then be equal to the solution found by the method of Tikhonov-type regularization (see e.g. [30, 59, 88]). Although this allows to study the posterior mean and the full posterior distribution by direct methods, in this section we derive the rate of posterior contraction from the general result Theorem 6.1. An advantage of this approach is that the proof can be extended to mixtures of Gaussian priors. Taking mixtures is important to obtain optimal recovery rates for true functions of different smoothness levels. See Section 6.5.

A Gaussian prior on the Hilbert space $H = H_0$ is determined by a mean, which we shall take equal to zero, and a covariance operator. To connect the prior to a smoothness scale $(H_s)_{s \in \mathbb{R}}$ as defined in Definition 2.1, it is natural to assume that the latter forms a *Hilbert scale*, which may be viewed a smoothness scale with additional structure. For reference we include a short summary on Hilbert scales. Extended discussions of Hilbert scales in the context of regularization theory can be found e.g. in Chapter 8 of [29], and a general treatment of the subject in [62].

Centred Gaussian distributions on a separable Hilbert space correspond bijectively to covariance operators. By definition a random variable F with values in H_0 is Gaussian if $\langle F, g \rangle_0$ is normally distributed, for every $g \in H_0$, and it has zero mean if these variables have zero means. The variances of these variables can then

be written as

$$\mathbb{E}\langle F, g \rangle_0^2 = \langle Cg, g \rangle_0,$$

for a linear operator $C : H_0 \rightarrow H_0$, called the *covariance operator*. A covariance operator C is necessarily self-adjoint, nonnegative, and of *trace class*, i.e., $\sum_{i \in \mathbb{N}} \langle C\phi_i, \phi_i \rangle < \infty$, for some (and then every) orthonormal basis $(\phi_i)_{i \in \mathbb{N}}$ of H_0 ; and every operator with these properties generates a Gaussian distribution.

In the setting of a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by the operator L it is natural to choose a Gaussian prior with covariance operator of the form $L^{-2\alpha}$, for some $\alpha > 0$. If L^{-1} has eigenvalues λ_j , then this operator is of trace class if $\sum_{j \in \mathbb{N}} \lambda_j^{-2\alpha} < \infty$. Thus α must be chosen big enough for the Gaussian prior to exist as a ‘proper’ prior on H_0 . For instance, if $\lambda_j \simeq j^{-1/d}$, then every choice $\alpha > d/2$ yields a proper prior.

This leads to the following theorem on posterior contraction rates for Gaussian priors, the proof of which is given in Section 6.6.

Theorem 6.7 (Gaussian Prior). *Consider a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by an operator L as in the preceding such that $L^{-1} : H_0 \rightarrow H_0$ is compact with eigenvalues λ_j satisfying $\lambda_j \simeq j^{-1/d}$. Suppose the operator $\mathcal{A} : H_0 \rightarrow G$ satisfies $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$, assume that $f_0 \in H_\beta$, for some $\beta > 0$, and let the prior be zero-mean Gaussian with covariance operator $L^{-2\alpha}$, for some $\alpha > d/2$. Then the posterior distribution satisfies, for sufficiently large $M > 0$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-((\alpha-d/2) \wedge \beta)/(2\alpha+2\gamma)} \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

If F is distributed according to the prior in the preceding theorem, then $L^s F$ is also zero-mean Gaussian distributed, with covariance operator $L^{2s-2\alpha}$, which has eigenvalues $j^{-(2\alpha-2s)/d}$. For $s < \alpha - d/2$, this operator is of trace class and hence $L^s F$ is a proper random variable in H_0 . In other words, the distribution of F gives probability 1 to $L^{-s} H_0 = H_s$, for every $s < \alpha - d/2$. The prior in the preceding theorem can therefore be interpreted as being ‘almost’ of regularity $\alpha - d/2$. The rate $n^{-((\alpha-d/2) \wedge \beta)/(2\alpha+2\gamma)}$ is therefore comparable to the rate obtained in Theorem 3.5 in [80] and Theorem 4.1 in [59] (without scaling parameter), except that the parameter α in the latter references is denoted presently by $\alpha - d/2$.

An improvement of the present theorem is that the covariance operator of the Gaussian prior is not directly linked to the operator \mathcal{A} , but only weakly so by (5.3). For example, we may construct a prior by a random series (see Theorem I.23 in Appendix I.6, [35]), in any basis corresponding to the smoothness scale. We illustrate this below by using the wavelet basis for an inverse problem given by a differential operator, after first noting that the singular value setup is covered as well.

Example 6.8 (SVD). The scale of smoothness classes constructed in Example 2.6 and Example 5.2 is the Hilbert scale attached to the operator L given by $Lf = \sum_{i \in \mathbb{N}} b_i f_i \phi_i$ defined on the domain of functions $f = \sum_{i \in \mathbb{N}} f_i \phi_i$, with $\sum_{i \in \mathbb{N}} b_i^2 f_i^2 < \infty$. Under assumption (5.4) this operator can also be expressed as $L = (\mathcal{A}^* \mathcal{A})^{-1/(2\gamma)}$, and depends on the operator \mathcal{A} through its eigenfunctions. A Gaussian prior with

covariance operator $L^{-2\alpha}$ corresponds to modelling the coefficients f_i relative to the basis ϕ_i as independent zero-mean normal variables F_i with variances $b_i^{-2\alpha}$. This follows, because in that case $\mathbb{E}\langle F, g \rangle_0^2 = \sum_{i \in \mathbb{N}} b_i^{-2\alpha} g_i^2 = \langle L^{-2\alpha} g, g \rangle_0^2$, for every $g \in H_0$.

Thus in this case the prior coincides with the ones in the literature studied under the SVD framework, e.g. [59, 60]. In the present more general setting L need not be directly linked to \mathcal{A} , except that the operator must possess the smoothing property Assumption 5.1.

Example 6.9 (Sobolev scales, wavelet prior). Let $\{\phi_{j,k}\}_{(j,k) \in \Lambda}$, be an S -regular orthonormal wavelet basis in $L^2(\mathbb{T})$, on $\mathbb{T} := (0, 2\pi]$. Let $f_{j,k} = \int_{\mathbb{T}} f(x) \phi_{j,k}(x) dx$ be the wavelet coefficients of a function f . By Parseval's identity, the map $U : f \mapsto \{f_{j,k}\}$ is a unitary operator $U : L^2(\mathbb{T}) \rightarrow \ell^2(\Lambda)$. The multiplication operator $m : \{f_{j,k}\} \mapsto \{2^j f_{j,k}\}$ on $\ell^2(\Lambda)$ has s -th power given by $m^s : \{f_{j,k}\} \mapsto \{2^{js} f_{j,k}\}$. Then $L := U^* m U$ has s -th power $L^s := U^* m^s U$ and generates a Hilbert scale $(H_s)_{s \in \mathbb{R}}$. For $f \in H_s$, we have

$$\|f\|_{H_s(\mathbb{T})}^2 = \sum_{j=0}^{\infty} 2^{2js} \sum_{k=0}^{2^j-1} f_{j,k}^2.$$

This norm can be shown to be equivalent to the standard Sobolev norm, for $0 \leq s < S$.

The Gaussian prior with covariance operator $L^{-2\alpha}$ can be represented by a random series of the form

$$F = \sum_{(j,k) \in \Lambda} F_{j,k} \phi_{j,k},$$

where $F_{j,k} \sim \mathcal{N}(0, 2^{-2j\alpha})$ are independent random variables. This prior corresponds to the Hilbert scale, but does not refer to an operator \mathcal{A} . For instance, the eigenbasis of the operator in Example 5.4 is the Fourier basis (see [57]), and not the wavelet basis. Thus we have constructed a Gaussian prior that is not related to the eigenbasis, but attains the same contraction rate.

It may be noted that the scale $(H_s)_{s \in \mathbb{R}}$ is well defined for every $s \in \mathbb{R}$, and with the preceding prior Theorem 6.7 is applicable to the full scale, and gives a contraction rate relative to the scale, which is optimal when $\beta = \alpha - d/2$. However, the scale agrees with the Sobolev scale only for $\beta < S$, and hence the optimality is in the Sobolev sense only if $\beta < S$. This restriction is typical when working with an approximation scheme such as wavelets or splines. One can of course choose a suitably large value of S , or may mix over multiple wavelet bases, as in the next section.

As mentioned in Section 6.1, there are many works on Bayesian inverse problems with Gaussian priors. The setup of the preceding theorem is similar to [1, 30], arguably closer to [1]. While we mainly treat the white noise case, our results can be extended to cover the noise structure in [1], and hence also cover the model in [30]. On the other hand, we differ from [1] in the following sense. First, unlike Assumption 3.1 in [1], our characterization of the smoothing property of the operator \mathcal{A} , i.e. Assumption 5.1, is simple, and in principle, our setup can also

be extended to severely ill-posed problems, see Section 6.7. Second, our proof strategy is different, as we do not use Gaussian conjugacy, which is the main tool in [1]. This also allows us to obtain posterior contraction rates for non-conjugate priors in Section 6.3, and for Gaussian mixtures in Section 6.5.

6.5 Gaussian Mixtures

The posterior contraction rate resulting from a zero-mean Gaussian prior with covariance operator $L^{-2\alpha}$, as considered in Section 6.4, is equal to the minimax rate $n^{-\beta/(2\beta+2\gamma+d)}$ (see [19]) only when $\alpha - d/2 = \beta$, i.e., when the prior smoothness $\alpha - d/2$ matches the true smoothness β . By mixing over Gaussian priors of varying smoothness the minimax rate can often be obtained simultaneously for a range of values β (cf. [61], [98], [89]). In this section we consider mixtures of the mean-zero Gaussian priors with covariance operators $\tau^2 L^{-2\alpha}$ over the ‘hyperparameter’ τ . Thus the prior Π is the distribution of τF , where F is a zero-mean Gaussian variable in H_0 with covariance operator $L^{-2\alpha}$, as in Section 6.4, and τ is an independent scale parameter. The variable $1/\tau^a$ may be taken to possess a Gamma distribution for some given $0 < a \leq 2$, or, more generally, should satisfy the following mild condition.

Condition 6.10. The distribution Q of τ has support $[0, \infty)$ and satisfies

$$\begin{cases} -\log Q((t, 2t)) \lesssim t^{-2}, & \text{as } t \downarrow 0, \\ -\log Q((t, 2t)) \lesssim t^{d/(\alpha-d/2)}, & \text{as } t \rightarrow \infty. \end{cases}$$

Theorem 6.11 (Gaussian mixture prior). *Consider a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by an operator L as in the preceding such that $L^{-1} : H_0 \rightarrow H_0$ is compact with eigenvalues λ_j satisfying $\lambda_j \simeq j^{-1/d}$. Suppose the operator $\mathcal{A} : H_0 \rightarrow G$ satisfies $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$, assume that $f_0 \in H_\beta$, for some $\beta \in (0, \alpha]$, and let the prior be a mixture of the zero-mean Gaussian distributions with covariance operators $\tau^2 L^{-2\alpha}$ over the parameters τ equipped with a prior satisfying Condition 6.10, for some $\alpha > d/2$. Then the posterior distribution satisfies, for sufficiently large $M > 0$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-\beta/(2\beta+2\gamma+d)} \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

The proof is given in Section 6.6.

6.6 Proofs

6.6.1 Proof of Theorem 6.5

The theorem is a corollary to Theorem 6.1 and uses arguments as in the proof of Proposition 3.2 in [80].

First we determine ε_n to satisfy the prior mass condition (6.6) of the direct problem. Let P_j be the projection onto the linear span of the first $j - 1$ basis

elements ϕ_i . By the assumption on \mathcal{A} and the triangle inequality, for any $i_n \in \mathbb{N}$,

$$\begin{aligned} \|\mathcal{A}f - \mathcal{A}f_0\| &\lesssim \|f - f_0\|_{-\gamma} \lesssim \|f - P_{i_n}f_0\|_{-\gamma} + \|P_{i_n}f_0 - f_0\|_{-\gamma} \\ &\lesssim \|f - P_{i_n}f_0\|_{-\gamma} + \delta(i_n, \gamma)\delta(i_n, \beta)\|f_0\|_{\beta}, \end{aligned} \quad (6.17)$$

by (2.4), if $0 \leq \beta, \gamma < S$. Here $\delta(i_n, \gamma)\delta(i_n, \beta) = i_n^{-(\gamma+\beta)/d} \simeq \varepsilon_n$ if $i_n \simeq \varepsilon_n^{-d/(\gamma+\beta)}$.

By the orthogonality of the basis (ϕ_i) , the function ϕ_j is orthogonal to the space V_j spanned by $(\phi_i)_{i < j}$. Hence $P_j\phi_j = 0$, so that $\|\phi_j\|_{-\gamma} \leq \delta(j, \gamma)\|\phi_j\|_0 \lesssim j^{-\gamma/d}$, for every j , by (2.4). Consequently, for $f = \sum_{i=1}^{i_n-1} f_i\phi_i \in V_{i_n}$ and $f_0 = \sum_i f_{0,i}\phi_i$, by the triangle inequality,

$$\|f - P_{i_n}f_0\|_{-\gamma} \lesssim \sum_{i=1}^{i_n-1} |f_i - f_{0,i}| i^{-\gamma/d}.$$

It follows that there exists a constant $a > 0$ such that

$$\begin{aligned} \Pi(f : \|f - P_{i_n}f_0\|_{-\gamma} < a\varepsilon) &\geq \Pi\left(\left((f_i), M\right) : \sum_{i=1}^{i_n-1} |f_i - f_{0,i}| i^{-\gamma/d} < \varepsilon, M = i_n - 1\right) \\ &\geq \prod_{i=1}^{i_n} \Pi\left(f_i : |f_i - f_{0,i}| < \frac{\varepsilon i^{\gamma/d}}{i_n}\right) \Pi(M = i_n - 1) \\ &\geq \prod_{i=1}^{i_n} \int_0^{\varepsilon i^{\gamma/d}/(\kappa_i i_n)} p\left(x + \frac{f_{0,i}}{\kappa_i}\right) dx e^{-b_1 i_n}, \end{aligned}$$

in view of Condition 6.4. By (6.15) of the latter assumption, the integral $\int_0^r p(x + \mu) dx$ is bounded below by a constant times $re^{-C(r+|\mu|)^w}$. It follows that for ε such that $\varepsilon i^{\gamma/d}/(\kappa_i i_n) \leq 1$, for $i \leq i_n$, the preceding display is lower bounded by a multiple of

$$\varepsilon^{i_n} \left[\prod_{i=1}^{i_n} \frac{i^{\gamma/d}}{\kappa_i i_n} \right] \exp\left[-C \sum_{i=1}^{i_n} \left(1 + \frac{|f_{0,i}|}{\kappa_i}\right)^w\right] e^{-b_1 i_n}.$$

By (6.16), we have $i^{\gamma/d}/\kappa_i \gtrsim (1/i)^{\gamma/d-\alpha}$, which is bounded below by 1 if $\gamma/d - \alpha \geq 0$ and by $(1/i_n)^{\alpha-\gamma/d}$ otherwise, and hence always by $(1/i_n)^\alpha$. This shows that the first term in square brackets is bounded below by $(a_2/i_n^{\alpha+1})^{i_n}$, for some $a_2 > 0$. Since $f_0 \in H_\beta$, by assumption, the norm duality (2.1) gives that $|f_{0,i}| = |\langle f_0, \phi_i \rangle_0| \leq \|f_0\|_\beta \|\phi_i\|_{-\beta} \lesssim i^{-\beta/d}$. Together with (6.16) this gives that $|f_{0,i}|/\kappa_i \lesssim i^{(\beta_0-\beta)/d} (\log i)^{1/w} \leq (\log i)^{1/w}$, whence minus the exponent in the second term in square brackets is bounded by a multiple of $i_n(1 + (\log i_n)^{1/w})^w$. We conclude that there exists a constant $a_3 > 0$ such that

$$\Pi(f : \|f - P_{i_n}f_0\|_{-\gamma} < a\varepsilon) \geq \varepsilon^{i_n} e^{-a_3 i_n \log i_n} e^{-b_1 i_n},$$

for every $\varepsilon > 0$ such that $\varepsilon i^{\gamma/d}/(\kappa_i i_n) \leq 1$, for every $i \leq i_n$. Since $i^{\gamma/d}/\kappa_i \lesssim i^{(\gamma+\beta_0)/d} (\log i)^{1/w}$, again by (6.16), a sufficient condition for the latter is that $\varepsilon i_n^{(\gamma+\beta_0)/d} (\log i_n)/i_n \leq 1$.

Combining this with (6.17), we see that (6.6) is satisfied for ε_n such that there exists i_n with

$$i_n^{-(\gamma+\beta)/d} \lesssim \varepsilon_n, \quad i_n \log i_n \lesssim n\varepsilon_n^2, \quad \varepsilon_n i_n^{(\gamma+\beta_0)/d} (\log i_n) \leq i_n.$$

This leads to the rates

$$\varepsilon_n \simeq (\log n/n)^{(\beta+\gamma)/(2\beta+2\gamma+d)}, \quad i_n \simeq (n/\log n)^{d/(2\beta+2\gamma+d)}.$$

(The third requirement is easily satisfied and remains inactive.) We can choose a sufficiently large proportionality constant in \simeq when defining ε_n , so that (6.6) is satisfied for ε_n , since the left and right sides of (6.6) are increasing and decreasing in ε_n , respectively.

Since the Galerkin projection $f^{(j)}$ is equal to f itself if $f \in V_j$, we have that $\|f^{(j_n)} - f\|_0 = 0$ for the random series $f = \sum_{i=1}^M f_i \phi_i$ if $M < j_n$. By (ii) of Condition 6.4 it follows that, for some $b'_2 > 0$ and every $\eta_n > 0$,

$$\Pi(f : \|f^{(j_n)} - f\|_0 > \eta_n) \leq \Pi(M \geq j_n) \leq e^{-b'_2 j_n}.$$

Hence (6.7) is satisfied for $j_n = n\varepsilon_n^2/(4b'_2)$. Thus we choose

$$j_n \simeq n^{d/(2\beta+2\gamma+d)} (\log n)^{(2\beta+2\gamma)/(2\beta+2\gamma+d)},$$

with a sufficiently large constant in \simeq . Then (6.3) is satisfied and it remains to solve η_n from (6.4) and (6.5). This leads to the inequalities

$$\begin{aligned} \eta_n &\geq \varepsilon_n j_n^{\gamma/d} \simeq n^{-\beta/(2\beta+2\gamma+d)} (\log n)^{(1+2\gamma/d)(\beta+\gamma)/(2\beta+2\gamma+d)}, \\ \eta_n &\geq j_n^{-\beta/d} \simeq n^{-\beta/(2\beta+2\gamma+d)} (\log n)^{-\beta(2\beta+2\gamma)/((2\beta+2\gamma+d)d)}. \end{aligned}$$

The rate is the maximum of the rates at the right hand sides, which coincides with the first rate. This concludes the proof.

6.6.2 Proof of Theorem 6.7

The theorem is a corollary to Theorem 6.1. The main tasks are to determine ε_n satisfying the prior mass condition (6.6) of the direct problem, and next to identify η_n from the prior mass condition (6.7) and the other conditions.

The first task is achieved in the following lemma.

Lemma 6.12. *Under the assumptions of Theorem 6.7, for $f_0 \in H_\beta$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon) \lesssim \begin{cases} \varepsilon^{-d/(\alpha+\gamma-d/2)}, & \text{if } d/2 < \alpha \leq \beta + d/2, \\ \varepsilon^{-(2\alpha-2\beta)/(\beta+\gamma)}, & \text{if } \alpha > \beta + d/2. \end{cases} \quad (6.18)$$

Proof. Since by assumption $\|\mathcal{A}f - \mathcal{A}f_0\| \simeq \|f - f_0\|_{-\gamma}$, the probability in the left side is the decentered small ball probability $\Pi(f : \|f - f_0\|_{-\gamma} < a\varepsilon)$ of the Gaussian random variable F distributed according to the prior and viewed as map into $H_{-\gamma} \supset H_0$, for some $a > 0$. Because F has covariance operator $L^{-2\alpha}$ as

a map in H_0 , its reproducing kernel Hilbert space (or Cameron-Martin space) \mathbb{H} (which does not depend on its range space) is equal to the range of $L^{-\alpha}$ under the norm $\|L^{-\alpha}h\|_{\mathbb{H}} = \|h\|_0$ (see e.g., Example I.14 of [35]). Since $L^{-\alpha} : H_0 \rightarrow H_\alpha$ is a norm isometry, by (iii) of Proposition 2.8, this is the Hilbert space H_α with its natural norm $\|\cdot\|_\alpha$. The left side of (6.18) is therefore up to constants equivalent to

$$\inf_{h \in H_\alpha: \|h - f_0\|_{-\gamma} < \varepsilon} \|h\|_\alpha^2 - \log \Pi(\|f\|_{-\gamma} < \varepsilon). \quad (6.19)$$

See [64, 65, 99], or Section 11.2, in particular, Proposition 11.19 in [35].

By (2.4) $\|P_j f_0 - f_0\|_{-\gamma} \lesssim \delta(j, \gamma) \delta(j, \beta) \|f_0\|_\beta$, which is bounded above by ε for $j \simeq \varepsilon^{-d/(\beta+\gamma)}$. Thus for this value of j the first term in (6.19) is bounded above by

$$\|P_j f_0\|_\alpha \lesssim \begin{cases} \|P_j f_0\|_\beta, & \text{if } \alpha \leq \beta, \\ 1/\delta(j, \alpha - \beta) \|P_j f_0\|_\beta, & \text{if } \alpha > \beta \end{cases}$$

by (2.11). Here $\|P_j f_0\|_\beta \leq \|P_j f_0 - f_0\|_\beta + \|f_0\|_\beta \leq (\delta(j, 0) + 1) \|f_0\|_\beta$, by (2.10). It follows that the contribution of the decentering in (6.19) is of order 1 if $\alpha \leq \beta$ and is bounded above by a term of order $\varepsilon^{-2(\alpha-\beta)/(\beta+\gamma)}$ if $\alpha > \beta$.

By Lemma 2.22, the metric entropy $\log N(\varepsilon, \{f \in H_\alpha : \|f\|_\alpha \leq 1\}, \|\cdot\|_{-\gamma})$ is of the order $\varepsilon^{-d/(\alpha+\gamma)}$. Hence, by [64] (see Lemma 6.2 in [100]),

$$-\log \Pi(\|f\|_{-\gamma} < \varepsilon) \simeq \varepsilon^{-d/(\alpha+\gamma-d/2)}.$$

Finally, the assertion of the lemma follows from discussion by cases. \square

It follows that (6.6) is satisfied for

$$\varepsilon_n \geq n^{-(\beta \wedge (\alpha-d/2) + \gamma)/(2\alpha+2\gamma)}. \quad (6.20)$$

The next step of the proof is to bound the prior probability in (6.7).

Lemma 6.13. *Under the assumptions of Theorem 6.7, there exist $a, b > 0$, such that for every $j \in \mathbb{N}$ and $t > 0$,*

$$\Pi(f : \|f^{(j)} - f\|_0 > t + aj^{1/2-\alpha/d}) \leq e^{-bt^2 j^{2\alpha/d}}.$$

Proof. We have $f^{(j)} - f = (R_j \mathcal{A} - I)f$, for $R_j = \mathcal{A}^{-1}Q_j$. Therefore, the probability on the left concerns the random variable $(R_j \mathcal{A} - I)F$, if F is a variable distributed according to the prior Π . Since F is zero-mean normal with covariance operator $L^{-2\alpha}$, this variable is zero-mean Gaussian with covariance operator $(R_j \mathcal{A} - I)L^{-2\alpha}(R_j \mathcal{A} - I)^*$. We shall compute the weak and strong second moments of the variable $(R_j \mathcal{A} - I)F$, and next apply Borell's inequality for the norm of a Gaussian variable to obtain the exponential bound.

Because $\langle (R_j \mathcal{A} - I)F, g \rangle_0 = \langle F, (R_j \mathcal{A} - I)^* g \rangle_0$ is zero-mean Gaussian with variance $\|L^{-\alpha}(R_j \mathcal{A} - I)^* g\|_0^2 = \|(R_j \mathcal{A} - I)^* g\|_{-\alpha}^2$, the weak second moment of $(R_j \mathcal{A} - I)F$ is given by

$$\sup_{\|g\|_0 \leq 1} \mathbb{E} \langle (R_j \mathcal{A} - I)F, g \rangle_0^2 = \sup_{\|g\|_0 \leq 1} \|(R_j \mathcal{A} - I)^* g\|_{-\alpha}^2.$$

By the norm duality (2.1), the right side is equal to

$$\sup_{\|g\|_0 \leq 1} \sup_{\|f\|_\alpha \leq 1} \langle f, (R_j \mathcal{A} - I)^* g \rangle_0^2 \leq \sup_{\|f\|_\alpha \leq 1} \|(R_j \mathcal{A} - I)f\|_0^2 \lesssim \delta(j, \alpha)^2.$$

in view of (5.9).

The strong second moment of the Gaussian variable $(R_j \mathcal{A} - I)F$ is equal to the trace of its covariance operator. As $\text{Trace}(S^*S) = \sum_i \|S\phi_i\|^2 = \sum_i \sum_j \langle S\phi_i, \phi_j \rangle^2 = \sum_i \|S^*\phi_i\|^2$, for any orthonormal basis (ϕ_i) and operator S , we have

$$\mathbb{E}\|(R_j \mathcal{A} - I)F\|_0^2 = \sum_{i \in \mathbb{N}} \|(R_j \mathcal{A} - I)L^{-\alpha}\phi_i\|_0^2.$$

For the orthonormal basis of eigenfunctions of L^{-1} and V_j the span of the first $j-1$ of these eigenfunctions, as in Proposition 2.9, $L^{-\alpha}V_j \subset V_j$, and hence $(R_j \mathcal{A} - I)L^{-\alpha}\phi_i$ vanishes for $i < j$. For $i \geq j$ the latter element is the difference $g^{(j)} - g$ of the Galerkin solution $g^{(j)}$ to $g = L^{-\alpha}\phi_i$. Therefore, by (5.9) the preceding display is bounded above by a multiple of

$$\sum_{i \geq j} \delta(i, \alpha)^2 \|L^{-\alpha}\phi_i\|_\alpha^2 = \sum_{i \geq j} \delta(i, \alpha)^2 \|\phi_i\|_0^2 \lesssim j^{1-2\alpha/d},$$

where we used the estimate $\sum_{i > j} i^{-b} \leq j^{1-b}/(b-1)$, for $b > 1$.

Since the first moment of $\|(R_j \mathcal{A} - I)F\|_0$ is bounded by the root of its second moment, the lemma follows by Borell's inequality (see e.g. Lemma 3.1 and subsequent discussion in [67]). \square

For $t^2 = 4n\varepsilon_n^2/(bj_n^{2\alpha/d})$ and $j = j_n$ the bound in the preceding lemma becomes $e^{-4n\varepsilon_n^2}$. Hence (6.7) is satisfied for

$$\eta_n \gtrsim \sqrt{n\varepsilon_n} j_n^{-\alpha/d} + j_n^{1/2-\alpha/d}.$$

Here we choose ε_n the minimal solution that satisfies the direct prior mass condition (6.6), given in (6.20). Next we solve for η_n under the constraints (6.4) and (6.5). The first of these constraints, $j_n \leq n\varepsilon_n^2$, shows that the first term on the right side of the preceding display always dominates the second term. Therefore, we obtain the requirements $j_n \leq n\varepsilon_n^2$ and

$$\begin{aligned} \eta_n &\geq \sqrt{n} n^{-(\beta \wedge (\alpha-d/2) + \gamma)/(2\alpha+2\gamma)} j_n^{-\alpha/d}, \\ \eta_n &\geq n^{-(\beta \wedge (\alpha-d/2) + \gamma)/(2\alpha+2\gamma)} j_n^{\gamma/d}, \\ \eta_n &\geq j_n^{-\beta/d}. \end{aligned}$$

Depending on the relation between α and $\beta + d/2$, two situations need to be discussed separately.

- (i) $\alpha \leq \beta + d/2$. We choose $j_n \simeq n^{d/(2\alpha+2\gamma)} = n\varepsilon_n^2$ and then see that the first two requirements in the preceding display both reduce to $\eta_n \geq n^{-(\alpha-d/2)/(2\alpha+2\gamma)}$, while the third becomes $\eta_n \geq n^{-\beta/(2\alpha+2\gamma)}$ and becomes inactive.
- (ii) $\alpha > \beta + d/2$. We choose $j_n \simeq n^{d/(2\alpha+2\gamma)} \leq n\varepsilon_n^2$, and then see that all three requirements reduce to $\eta_n \geq n^{-\beta/(2\alpha+2\gamma)}$.

Finally, we apply Theorem 6.1 to complete the proof.

6.6.3 Proof of Theorem 6.11

Let Π_τ denote the zero-mean Gaussian distribution on H with covariance operator $\tau^2 L^{-2\alpha}$ (where $\alpha > d/2$).

Lemma 6.14. *Under the assumptions of Theorem 6.11, for $f_0 \in H_\beta$ and $\beta \leq \alpha$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi_\tau(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon) \lesssim \frac{1}{\tau^2} \left(\frac{1}{\varepsilon}\right)^{(2\alpha-2\beta)/(\beta+\gamma)} + \left(\frac{\tau}{\varepsilon}\right)^{d/(\alpha+\gamma-d/2)}.$$

Lemma 6.15. *Under the assumptions of Theorem 6.11, for $f_0 \in H_\beta$ and $\beta \leq \alpha$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi_\tau(f : \|f\|_0 < \varepsilon) \gtrsim \left(\frac{\tau}{\varepsilon}\right)^{d/(\alpha-d/2)}.$$

Lemma 6.16. *Under the assumptions of Theorem 6.11, there exist $a, b > 0$ such that, for every $j \in \mathbb{N}$ and $x, \tau > 0$,*

$$\Pi_\tau(f : \|f^{(j)} - f\|_0 > \tau x + \tau a j^{1/2-\alpha/d}) \leq e^{-bx^2 j^{2\alpha/d}}$$

Proofs. The proof of the first lemma follows the same lines as the proof of Lemma 6.12, except that now the Cameron-Martin space of the measure Π_τ on $H_{-\gamma}$ is H_α equipped with the norm $\|\cdot\|_{\mathbb{H}} = \frac{1}{\tau} \|\cdot\|_\alpha$ rather than its natural norm. The second lemma follows similarly, but considers the centered probability only. The third lemma is immediate from Lemma 6.13 as Π_τ is the law of τF , for F the Gaussian variable with the law Π as in the latter lemma, and the map $f \mapsto f^{(j)} - f$ is linear. \square

As preparation for the proof of Theorem 6.11, we first show that the minimax rate can be obtained by a Gaussian prior with the deterministic scaling, dependent on β , given by

$$\tau_n = n^{(\alpha-d/2-\beta)/(2\beta+2\gamma+d)}. \quad (6.21)$$

Theorem 6.17. *Assume the conditions on the Hilbert scale, the forward operator A and the true parameter f_0 in Theorem 6.7 hold. Suppose that the priors Π are zero-mean Gaussian with covariance operators $\tau_n^2 L^{-2\alpha}$ with τ_n as given in (6.21) and $\alpha > d/2$. Then for $\beta \leq \alpha$, the posterior distribution satisfies, for sufficiently large $M > 0$,*

$$\Pi_n(f : \|f - f_0\|_0 > Mn^{-\beta/(2\beta+2\gamma+d)} \mid Y^{(n)}) \stackrel{P_{f_0}^{(n)}}{\xrightarrow{}} 0.$$

Proof. The theorem is a corollary to Theorem 6.1. The proof follows the same lines as the proof of Theorem 6.7. By Lemma 6.14, inequality (6.6) is satisfied for

$$\varepsilon_n \gtrsim n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}.$$

By Lemma 6.16, inequality (6.7) is satisfied for

$$\eta_n \gtrsim \tau_n (\sqrt{n} \varepsilon_n j_n^{-\alpha/d} + j_n^{1/2-\alpha/d}).$$

We choose $j_n \simeq n\varepsilon_n^2$, and the minimal solution $\varepsilon_n = n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}$ to the second last display. It is then straightforward to verify that (6.4), (6.5) and (6.7) are satisfied for $\eta_n \simeq n^{-\beta/(2\beta+2\gamma+d)}$. \square

Theorem 6.11 is a corollary of Theorem 6.3, with the choices

$$\begin{aligned} \eta_n &\simeq n^{-\beta/(2\beta+2\gamma+d)}, & \varepsilon_n &\simeq n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}, \\ j_n &\simeq n\varepsilon_n^2 = n^{d/(2\beta+2\gamma+d)}. \end{aligned}$$

Conditions (6.3), (6.4), and (6.5) are satisfied for these choices. It remains to verify (6.6), and eqs. (6.12) and (6.13).

For ease of notation, for the moment, define η_n and ε_n as in the preceding display, with exact equality (i.e., with the constant set equal 1). Let τ_n be the ‘optimal’ scaling rate defined in (6.21).

Verification of (6.6). For $\tau \simeq \tau_n$ and $\varepsilon \simeq \varepsilon_n$ as given and $\beta \leq \alpha$, both terms in the right side of Lemma 6.14 are of the order $n\varepsilon_n^2$. The lemma yields, for $\tau_n \leq \tau \leq 2\tau_n$ and some constant $a_1 > 0$,

$$-\log \Pi_\tau(f : \|Af - Af_0\| < \varepsilon_n) \leq a_1 n\varepsilon_n^2.$$

This shows that

$$\begin{aligned} \Pi(f : \|Af - Af_0\| < \varepsilon_n) &= \int_0^\infty \Pi_\tau(f : \|Af - Af_0\| < \varepsilon_n) dQ(\tau) \\ &\geq e^{-a_1 n\varepsilon_n^2} Q(\tau_n, 2\tau_n). \end{aligned}$$

If $\alpha - d/2 < \beta$, then $\tau_n \rightarrow 0$, and Condition 6.10 on Q gives that

$$-\log Q(\tau_n, 2\tau_n) \lesssim \tau_n^{-2} = n^{(2\beta-2\alpha+d)/(2\beta+2\gamma+d)} \leq n^{d/(2\beta+2\gamma+d)} = n\varepsilon_n^2,$$

if $\beta \leq \alpha$. If $0 < \beta < \alpha - d/2$, then $\tau_n \rightarrow \infty$, and Condition 6.10 on Q gives that

$$\begin{aligned} -\log Q(\tau_n, 2\tau_n) &\lesssim \tau_n^{d/(\alpha-d/2)} = n^{(d(\alpha-d/2-\beta)/(\alpha-d/2)(2\beta+2\gamma+d))} \\ &\leq n^{d/(2\beta+2\gamma+d)} = n\varepsilon_n^2. \end{aligned}$$

Finally if $\alpha - d/2 = \beta$, then $\tau_n = 1$ and $Q(\tau_n, 2\tau_n) \gtrsim 1$. Thus in all three cases $Q(\tau_n, 2\tau_n)$ is bounded below by a power of $e^{-n\varepsilon_n^2}$. Combining this with the preceding, we see that $\Pi(f : \|Af - Af_0\| \leq \varepsilon_n) \geq e^{-a_2 n\varepsilon_n^2}$, for some positive constant a_2 , which we can take bigger than 1. Then (6.6) is satisfied for ε_n equal to $\sqrt{a_2}$ times the current ε_n .

Verification of (6.12). Lemma 6.15 gives that

$$\Pi_\tau(f : \|f - f_0\|_0 < 2\eta_{n,\tau}) \leq \Pi_\tau(f : \|f\|_0 < 2\eta_{n,\tau}) \leq e^{-a_3(\tau/\eta_{n,\tau})^{d/(\alpha-d/2)}},$$

for some constant a_3 . This is bounded above by $e^{-4a_2 n\varepsilon_n^2}$ if

$$\eta_{n,\tau} = 2a_4 \tau n^{(d/2-\alpha)/(2\beta+2\gamma+d)} = 2a_4 \tau \eta_n / \tau_n,$$

for a sufficiently small constant $a_4 > 0$.

Verification of (6.13). Choosing $x = a_4\eta_n/\tau_n = \eta_{n,\tau}/(2\tau)$ in Lemma 6.16, we see that the left side of (6.13) is bounded above by $e^{-4a_2n\varepsilon_n^2}$ if j_n satisfies

$$a_4j_n^{1/2-\alpha/d} \leq a_4\eta_n/\tau_n, \quad \text{and} \quad ba_4^2(\eta_n/\tau_n)^2j_n^{2\alpha/d} \geq 4a_2n\varepsilon_n^2.$$

Both inequalities become equalities for j_n of the order $j_n \simeq n^{d/(2\beta+2\gamma+d)}$, as indicated at the beginning of the proof. Since $1/2-\alpha/d < 0$ and $2\alpha/d > 0$, the left side of the first inequality is decreasing in j_n and the left side of second inequality is increasing. Thus both inequalities are satisfied for $j_n = a_5n^{d/(2\beta+2\gamma+d)}$ and a sufficiently large constant a_5 .

Finally we choose ε_n and j_n in Theorem 6.3 equal to $\sqrt{a_2}$ and a_5 times the orders indicated at the beginning of the proof. Then (6.3) is satisfied, and (6.4) and (6.5) are satisfied if η_n is chosen of the indicated order times a sufficiently large constant.

6.7 Discussion and Comments

In this section we comment on the present setup and discuss directions in which the results in this chapter can be extended.

Coloured Noise

We have examined the case that the noise ξ in model (6.1) is white noise. Statistical estimation in the case that the noise is a proper centred Gaussian random element in G , as studied in [30], is easier in terms of minimax rates (if in both cases the noise is scaled to the same unit), as this would imply that the noise is less variable. By inspection of our proofs one sees that the concentration inequalities that drive the testing criterion remain valid if the covariance operator of the noise is bounded above by the identity, as is assumed in [1, 7]. As a consequence, the proof of Theorem 6.1 goes through and the theorem remains valid, as do the corollaries in the later sections. However, for truly coloured noise the result may be suboptimal, as one may expect a faster posterior contraction rate, which will incorporate the decrease of the noise variance in certain directions. The methods of the present chapter can be adapted to this case as long as the covariance operator fits the scale of smoothness classes, as in [30]. A sharp result in full generality may be difficult to attain, as it will be the outcome of the interaction of the directions of decrease in the noise, the true parameter and the prior.

Approximation Numbers of Embeddings

In the corollaries to the main result we have assumed that the approximation numbers $\delta(j, s)$ of the canonical embedding $\iota : H_s \rightarrow H_0$ are of polynomial order $j^{-s/d}$. This order matches the approximation numbers of Sobolev spaces on d -dimensional, bounded domains, and seems common. Other decay rates do arise, e.g., an exponential rate in severely ill-posed problems (as in the heat equation considered in [60]), or a logarithmic rate (as in [14]). The general Theorem 6.1 remains valid, but its corollaries must be adapted. For Gaussian priors in logarithmic or exponential scales, this is relatively straightforward using the general

theory of approximation numbers, which relates these to singular values and metric entropy. Some results can be found in Section 10.4.

Chapter 7

Inverse Problems with Discrete Observations: Gaussian Conjugacy

7.1 Introduction

Linear inverse problems have been studied since long in the statistical and numerical analysis literature; see, e.g., [3, 7, 15, 16, 19, 24, 56, 57, 101], and references therein. Emphasis in these works has been on the signal-in-white noise model,

$$Y = Af + \varepsilon W, \quad (7.1)$$

where the parameter of interest f lies in some infinite-dimensional function space, A is a linear operator with values in a possibly different space, W is white noise, and ε is the noise level. Applications of linear inverse problems include, e.g., computerized tomography, see [72], partial differential equations, see [54], and scattering theory, see [20].

Arguably, in practice one does not have access to a full record of observations on the unknown function f as in the idealised model (7.1), but rather one indirectly observes it at a finite number of points. This statistical setting can be conveniently formalised as follows: let the signal of interest f be an element in a Hilbert space H_1 of functions defined on a compact interval $[0, 1]$. The forward operator A maps f to another Hilbert space H_2 . We assume that H_1, H_2 are subspaces of $L^2([0, 1])$, typically collections of functions of certain smoothness as specified in the later sections, and that the design points are chosen deterministically,

$$\left\{ x_i = \frac{i}{n} \right\}_{i=1, \dots, n}. \quad (7.2)$$

Assuming continuity of Af and defining

$$Y_i = Af(x_i) + \xi_i, \quad i = 1, \dots, n, \quad (7.3)$$

with ξ_i i.i.d. standard Gaussian random variables, our observations are the pairs $(x_i, Y_i)_{i \leq n}$, and we are interested in estimating f . A prototype example we think of is the case when A is the solution operator in the Dirichlet problem for the heat equation acting on the initial condition f ; see Example 7.8 below for details.

Model (7.3) is related to the inverse regression model studied e.g. in [5] and [6]. Although the setting we consider is somewhat special, our contribution is arguably the first one to study from a theoretical point of view a nonparametric Bayesian approach to estimation of f in the inverse problem setting with partial observations (see [35] for a monographic treatment of modern Bayesian nonparametrics). In the context of the signal-in-white noise model (7.1), a nonparametric Bayesian approach has been studied thoroughly in [59] and [60], and techniques from these works will turn out to be useful in our context as well. Our results will deal with derivation of posterior contraction rates and study of asymptotic frequentist coverage of Bayesian credible sets. A posterior contraction rate can be thought of as a Bayesian analogue of a convergence rate of a frequentist estimator, cf. [33] and [35]. Specifically, we will show that as the sample size $n \rightarrow \infty$, the posterior distribution concentrates around the ‘true’ parameter value, under which data have been generated, and hence our Bayesian approach is consistent and asymptotically recovers the unknown ‘true’ f . The rate at which this occurs will depend on the smoothness of the true parameter and the prior and the ill-posedness degree of the problem. Correct combinations of these values lead to optimal posterior contraction rates (up to logarithmic factors). Furthermore, a Bayesian approach automatically provides uncertainty quantification in parameter estimation through the spread of the posterior distribution, specifically by means of posterior credible sets. We will give an asymptotic frequentist interpretation of these sets in our context. In particular, we will see that the frequentist coverage will depend on a combination of smoothness of the true parameter and the prior, and the ill-posedness of the problem. Oversmoothing priors lead to zero coverage, while undersmoothing priors produce highly conservative results.

The chapter is organized as follows: in Section 7.2, we give a detailed description of the problem, introduce the singular value decomposition and convert the model (7.3) into an equivalent truncated sequence model that is better amenable to our theoretical analysis. We show how a Gaussian prior in this sequence model leads to a Gaussian posterior and give an explicit characterisation of the latter. Our main results on posterior contraction rates and Bayesian credible sets are given in Section 7.3, followed by simulation examples in Section 7.4 that illustrate our theoretical results. Section 7.5 contains the proofs of the main theorems, while the technical lemmas used in the proofs are collected in Section 7.6.

7.1.1 Notation

The notational conventions we use in this work are the following: definitions are marked by the $:=$ symbol; $|\cdot|$ denotes the absolute value and $\|\cdot\|_H$ indicates the norm related to the space H ; $\langle \cdot, \cdot \rangle_H$ is understood as the canonical inner product in the inner product space H ; subscripts are omitted when there is no danger of confusion; $\mathcal{N}(\mu, \Sigma)$ denotes the Gaussian distribution with mean μ and covariance operator Σ ; subscripts \mathcal{N}_n and \mathcal{N}_H may be used to emphasize the fact that the distribution is defined on the space \mathbb{R}^n or on the abstract space H ; $\text{Cov}(\cdot, \cdot)$ denotes the covariance or the covariance operator, depending on the context

7.2 Sequence model

7.2.1 Singular value decomposition

We impose a common assumption on the forward operator A from the literature on inverse problems, see, e.g., [3], [7] and [15].

Assumption 7.1. Operator A is injective and compact.

It follows that A^*A is also compact and in addition self-adjoint. Hence, by the spectral theorem for self-adjoint compact operators, see [21], we have a representation $A^*Af = \sum_{k \in \mathbb{N}} a_k^2 f_k \varphi_k$, where $\{\varphi_k\}$ and $\{a_k\}$ are the eigenbasis on H_1 and eigenvalues, respectively, (corresponding to the operator A^*A), and $f_k = \langle f, \varphi_k \rangle$ are the Fourier coefficients of f . This decomposition of A^*A is known as the singular value decomposition (SVD), and $\{a_k\}$ are also called singular values.

It is easy to show that the conjugate basis $\psi_k := A\varphi_k/a_k$ of the orthonormal basis $\{\varphi_k\}_k$ is again an orthonormal system in H_2 and gives a convenient basis for $\text{Ran } A$, the range of A in H_2 . Furthermore, the following relations hold (see [3]),

$$A\varphi_k = a_k\psi_k, \quad A^*\psi_k = a_k\varphi_k. \quad (7.4)$$

Recall a standard result (see, e.g., [44]): a Hilbert space H is isometric to ℓ^2 , and Parseval's identity $\|f\|_{\ell^2}^2 := \sum_k |f_k|^2 = \|f\|_H^2$ holds; here f_k are the Fourier coefficients with respect to some known and fixed orthonormal basis.

We will employ the eigenbasis $\{\varphi_k\}$ of A^*A to define the Sobolev space of functions. This will define the space in which the unknown function f resides.

Definition 7.2. We say f is in the Sobolev space S^β with smoothness parameter $\beta \geq 0$, if it can be written as $f = \sum_{k=1}^\infty f_k \varphi_k$ with $f_k = \langle f, \varphi_k \rangle$, and if its norm $\|f\|_\beta := (\sum_{k=1}^\infty f_k^2 k^{2\beta})^{1/2}$ is finite.

Remark 7.3. The above definition agrees with the classical definition of the Sobolev space if the eigenbasis is the trigonometric basis, see, e.g., [95]. With a fixed basis, which is always the case in this article, one can identify the function f and its Fourier coefficients $\{f_k\}$. Thus, we use S^β to denote both the function space and the sequence space. For example, it is easy to verify that $S^0 = \ell^2$ (correspondingly $S^0 = L^2$), $S^\beta \subset \ell^2$ for any nonnegative β , and $S^\beta \subset \ell^1$ when $\beta > 1/2$.

Recall that $Af = \sum a_i f_i \psi_i$. Then we have $Af \in S^{\beta+p}$ if $a_k \asymp k^{-p}$, and $Af \in S^\infty := \cap_{k \in \mathbb{N}} S^k$, if a_k decays exponentially fast. Such a lifting property is beneficial in the forward problem, since it helps to obtain a smooth solution. However, in the context of inverse problems it leads to a difficulty in recovery of the original signal f , since information on it is washed out by smoothing. Hence, in the case of inverse problems one does not talk of the lifting property, but of ill-posedness, see [15].

Definition 7.4. An inverse problem is called mildly ill-posed, if $a_k \asymp k^{-p}$ as $k \rightarrow \infty$, and extremely ill-posed, if $a_k \asymp e^{-k^s p}$ with $s \geq 1$ as $k \rightarrow \infty$, where p is strictly positive in both cases.

In the rest of the article, we will confine ourselves to the following setting.

Assumption 7.5. The unknown true signal f in (7.3) satisfies $f \in S^\beta \subset H_1$ for $\beta > 0$. Furthermore, the ill-posedness is of one of the two types in Definition Definition 7.4.

Remark 7.6. As an immediate consequence of the lifting property, we have $H_2 \subset H_1$.

We conclude this section with two canonical examples of the operator A .

Example 7.7 (mildly ill-posed case: Volterra operator [59]). The classical Volterra operator $A : L^2[0, 1] \rightarrow L^2[0, 1]$ and its adjoint A^* are

$$Af(x) = \int_0^x f(s) ds, \quad A^*f(x) = \int_x^1 f(s) ds.$$

The eigenvalues, eigenfunctions of A^*A and the conjugate basis are given by

$$\begin{aligned} a_i^2 &= \frac{1}{(i - 1/2)^2 \pi^2}, \\ \varphi_i(x) &= \sqrt{2} \cos((i - 1/2)\pi x), \\ \psi_i(x) &= \sqrt{2} \sin((i - 1/2)\pi x), \end{aligned}$$

for $i \geq 1$.

Example 7.8 (extremely ill-posed case: heat equation [60]). Consider the Dirichlet problem for the heat equation:

$$\begin{aligned} \frac{\partial}{\partial t} u(x, t) &= \frac{\partial^2}{\partial x^2} u(x, t), \quad u(x, 0) = f(x), \\ u(0, t) &= u(1, t) = 0, \quad t \in [0, T], \end{aligned} \tag{7.5}$$

where $u(x, t)$ is defined on $[0, 1] \times [0, T]$ and $f(x) \in L^2[0, 1]$ satisfies $f(0) = f(1) = 0$. The solution of (7.5) is given by

$$u(x, t) = \sqrt{2} \sum_{k=1}^{\infty} f_k e^{-k^2 \pi^2 t} \sin(k\pi x) =: Af(x),$$

where $\{f_k\}$ are the coordinates of f in the basis $\{\sqrt{2} \sin(k\pi x)\}_{k \geq 1}$.

For the solution map A , the eigenvalues of A^*A are $e^{-k^2 \pi^2 t}$, the eigenbasis and conjugate basis coincide and $\varphi_k(x) = \psi_k(x) = \sqrt{2} \sin(k\pi x)$.

7.2.2 Equivalent formulation

In this subsection we develop a sequence formulation of the model (7.3), which is very suitable for asymptotic Bayesian analysis. First, we briefly discuss the relevant results that provide motivation for our reformulation of the problem.

In Examples 7.7 and 7.8, the sine and cosine bases form the eigenbasis. In fact, the Fourier basis (trigonometric polynomials) frequently arises as an eigenbasis for various operators, e.g. in the case of differentiation, see [27], or circular deconvolution, see [16]. For simplicity, we will use Fourier basis as a primary example in the rest of the article. Possible generalization to other bases is discussed in Remark 7.10.

Restriction of our attention to the Fourier basis is motivated by its special property: discrete orthogonality. The next lemma illustrates this property for the sine basis (Example 7.8).

Lemma 7.9 (discrete orthogonality). *Let $\{\psi_k\}_{k \in \mathbb{N}}$ be the sine basis, i.e.*

$$\psi_k(x) = \sqrt{2} \sin(k\pi x), \quad k = 1, 2, 3, \dots$$

Then:

(i.) *Discrete orthogonality holds:*

$$\langle \psi_j, \psi_k \rangle_d := \frac{1}{n} \sum_{i=1}^n \psi_j(i/n) \psi_k(i/n) = \delta_{jk}, \quad j, k = 1, \dots, n-1. \quad (7.6)$$

Here δ_{jk} is the Kronecker delta.

(ii.) *Fix $l \in \mathbb{N}$. For any fixed $1 \leq k \leq n-1$ and all $j \in \{ln, ln+1, \dots, (l+1)n-1\}$, there exists only one $\bar{k} \in \{1, 2, \dots, n-1\}$ depending only on the parity of l , such that for $\tilde{j} = ln + \bar{k}$, the equality*

$$|\langle \psi_{\tilde{j}}, \psi_k \rangle_d| = 1 \quad (7.7)$$

holds, while $\langle \psi_{\tilde{j}}, \psi_k \rangle_d = 0$ for all $\tilde{j} = ln + \tilde{k}$ such that $\tilde{k} \neq \bar{k}$, $\tilde{k} \in \{1, 2, \dots, n-1\}$.

Remark 7.10. For other trigonometric bases, discrete orthogonality can also be attained. Thus, the conjugate eigenbasis in Example 7.7 is discretely orthogonal with design points $\{(i-1/2)/n\}_{i=1, \dots, n}$. We refer to [2] and references therein for details. With some changes in the arguments, our asymptotic statistical results still remain valid with such modifications of design points compared to (7.2). We would like to stress the fact that restricting attention to bases with discrete orthogonality property does constitute a loss of generality. However, there exist classical bases other than trigonometric bases that are discretely orthogonal (possibly after a suitable modification of design points). See, for instance, [78] for an example of Lagrange polynomials.

Motivated by the observations above, we introduce our central assumption on the basis functions.

Assumption 7.11. Given the design points $\{x_i\}_{i=1, \dots, n}$ in (7.2), we assume the conjugate basis $\{\psi_k\}_{k \in \mathbb{N}}$ of the operator A in (7.3) possesses the following properties:

(i.) for $1 \leq j, k \leq n-1$,

$$\langle \psi_j(x), \psi_k(x) \rangle_d := \frac{1}{n} \sum_{i=1}^n \psi_j(x_i) \psi_k(x_i) = \delta_{jk}$$

(ii.) For $1 \leq k \leq n-1$ and $j \in \{ln, \dots, (l+1)n-1\}$ with fixed $l \in \mathbb{N}$, there exists only one $\tilde{j} = ln + \bar{k}$, such that $0 < |\langle \psi_{\tilde{j}}, \psi_k \rangle_d| < M$, where M is a fixed constant, and \bar{k} depends on the parity of l only. For other $j \neq \tilde{j}$, $|\langle \psi_j, \psi_k \rangle_d| = 0$.

Using the shorthand notation

$$f = \sum_j f_j \varphi_j = \sum_{j=1}^{n-1} f_j \varphi_j + \sum_{j \geq n} f_j \varphi_j =: f^n + f^r,$$

we obtain for $k = 1, \dots, n-1$ that

$$\begin{aligned} U_k &= \frac{1}{n} \sum_{i=1}^n Y_i \psi_k(x_i) = \langle Af^n, \psi_k \rangle_d + \langle Af^r, \psi_k \rangle_d + \frac{1}{n} \sum_{i=1}^n \xi_i \psi_k(x_i) \\ &= a_k f_k + R_k + \frac{1}{\sqrt{n}} \zeta_k, \end{aligned} \tag{7.8}$$

where

$$R_k := R_k(f) = \langle Af^r, \psi_k \rangle_d, \quad \zeta_k := \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i \psi_k(x_i).$$

By Assumption 7.11, we have

$$|R_k| = |\langle Af^r, \psi_k \rangle_d| \leq \sum_{j \geq n} a_j |f_j| |\langle \psi_j, \psi_k \rangle_d| = \sum_{l=1}^{\infty} a_{ln+\bar{k}} |f_{ln+\bar{k}}|, \tag{7.9}$$

which leads to (via Cauchy-Schwarz)

$$R_k^2(f) \leq \left(\sum_{l=1}^{\infty} a_{ln+\bar{k}}^2 (ln + \bar{k})^{-2\beta} \right) \|f\|_{\beta}^2.$$

Hence, for a mildly ill-posed problem, i.e. $a_k \asymp k^{-p}$, the following bound holds, uniformly in the ellipsoid $\{f : \|f\|_{\beta} \leq K\}$,

$$\begin{aligned} \sup_{f: \|f\|_{\beta} \leq K} R_k^2(f) &\lesssim \sum_{l=1}^{\infty} (ln)^{-2\beta-2p} = n^{-2(\beta+p)} \sum_{l=1}^{\infty} l^{-2(\beta+p)} \\ &\asymp n^{-2(\beta+p)} = o(1/n), \end{aligned} \tag{7.10}$$

for any $1 \leq k \leq n-1$ when $\beta + p > 1/2$.

If the problem is extremely ill-posed, i.e. $a_k \asymp e^{-k^s p}$, we use the inequality

$$R_k^2(f) \leq \left(\sum_{j \geq n} a_j |f_j| \right)^2 \leq \left(\sum_{j \geq n} a_j^2 \right) \|f^r\|^2.$$

Since $a_j \asymp \exp(-pj^s) \leq \exp(-pj)$, it follows that $\sum_{j \geq n} a_j^2$ is up to a constant bounded from above by $\exp(-2pn)$. Hence

$$\sup_{f: \|f\|_{\beta} \leq K} R_k^2(f) \lesssim \exp(-2pn) \ll 1/n. \quad (7.11)$$

In [59], [60], the Gaussian prior $\Pi = \otimes_{i \in \mathbb{N}} \mathcal{N}(0, \lambda_i)$ is employed on the coordinates of the eigenbasis expansion of f . If $\lambda_i = \rho_n^2 i^{-1-2\alpha}$, the sum $\sum_{i \in \mathbb{N}} \lambda_i = \rho_n^2 \sum_{i \in \mathbb{N}} i^{-1-2\alpha}$ is convergent, and hence this prior is the law of a Gaussian element in H_1 .

In our case, we consider the same type of the prior with an additional constraint that only the first $n-1$ components of the prior are non-degenerate, i.e. $\Pi = (\otimes_{i < n} \mathcal{N}(0, \lambda_i)) \times (\otimes_{i \geq n} \mathcal{N}(0, 0))$, where λ_i is as above. In addition, we assume the prior on f is independent of the noise ζ_k , $k = 1, \dots, n-1$, in (7.8). With these assumptions in force, we see $\Pi(R_k = 0) = 1$, for $k = 1, \dots, n-1$. Furthermore, the posterior can be obtained from the product structure of the model and the prior via the normal conjugacy,

$$\Pi(f|U^n) = \otimes_{k \in \mathbb{N}} \mathcal{N}(\hat{f}_k, \sigma_k^2), \quad (7.12)$$

$$\text{with } \hat{f}_k = \frac{na_k \lambda_k \mathbb{1}_{\{k < n\}}}{na_k^2 \lambda_k + 1} U_k, \quad \sigma_k^2 = \frac{\lambda_k \mathbb{1}_{\{k < n\}}}{na_k^2 \lambda_k + 1}.$$

We also introduce

$$\hat{f} = \mathbb{E}(f|U^n) = (\mathbb{E}(f_k|U_k)) = (\hat{f}_k)_{k \in \mathbb{N}} = (b_k U_k)_{k \in \mathbb{N}}, \quad (7.13)$$

where $b_k = \frac{na_k \lambda_k \mathbb{1}_{\{k < n\}}}{na_k^2 \lambda_k + 1}$. We conclude this section with a useful fact that will be applied in later sections:

$$\hat{f}_k = b_k U_k = b_k \left(a_k f_k + R_k + \frac{\zeta_k}{\sqrt{n}} \right) = \mathbb{E} \hat{f}_k + \tau_k \zeta_k, \quad (7.14)$$

where $\mathbb{E} \hat{f}_k = a_k b_k f_k + b_k R_k$ and $\tau_k = b_k / \sqrt{n}$.

7.3 Main results

7.3.1 Contraction rates

In this section, we determine the rate at which the posterior distribution concentrates on shrinking neighbourhoods of the ‘true’ parameter f_0 as the sample size n grows to infinity.

Assume the observations in (7.3) have been collected under the parameter value $f_0 = \sum_{k \in \mathbb{N}} f_{0,k} \varphi_k$. Thus our observations $(U_k)_{k < n}$ given in (7.8) have the law $\otimes_{k < n} \mathcal{N}(a_k f_{0,k} + R_k, 1/n)$. We will use the notation $\Pi_n(\cdot|U)$ to denote the posterior distribution given in (7.12).

Theorem 7.12 (Posterior contraction: mildly ill-posed problem). *If the problem is mildly ill-posed as $a_k \asymp k^{-p}$ with $p > 0$, the true parameter $f_0 \in S^\beta$ with $\beta > 0$, and furthermore $\beta + p > 1/2$, by letting $\lambda_k = \rho_n^2 k^{-1-2\alpha}$ with $\alpha > 0$ and any positive ρ_n satisfying $\rho_n^2 n \rightarrow \infty$, we have, for any $K > 0$ and $M_n \rightarrow \infty$,*

$$\sup_{\|f_0\|_\beta \leq K} \mathbb{E}_{f_0} \Pi_n (f : \|f - f_0\|_{H_1} \geq M_n \varepsilon_n | U^n) \rightarrow 0,$$

where

$$\varepsilon_n = \varepsilon_{n,1} \vee \varepsilon_{n,2} = (\rho_n^2 n)^{-\beta/(2\alpha+2p+1) \wedge 1} \vee \rho_n (\rho_n^2 n)^{-\alpha/(2\alpha+2p+1)}. \quad (7.15)$$

In particular,

- (i.) if $\rho_n = 1$, then $\varepsilon_n = n^{-(\alpha \wedge \beta)/(2\alpha+2p+1)}$;
- (ii.) if $\beta \leq 2\alpha + 2p + 1$ and $\rho_n \asymp n^{(\alpha-\beta)/(2\beta+2p+1)}$, then $\varepsilon_n = n^{-\beta/(2\beta+2p+1)}$;
- (iii.) if $\beta > 2\alpha + 2p + 1$, then for every scaling ρ_n , $\varepsilon_n \gg n^{-\beta/(2\beta+2p+1)}$.

Thus we recover the same posterior contraction rates as obtained in [59], at the cost of an extra constraint $\beta + p > 1/2$. The frequentist minimax convergence rate for mildly ill-posed problems in the white noise setting with $\varepsilon = n^{-1/2}$ is $n^{-\beta/(2\beta+2p+1)}$, see [15]. We will compare our result to this rate. Our theorem states that in case (i.) the posterior contraction rate reaches the frequentist optimal rate if the regularity of the prior matches the truth ($\beta = \alpha$) and the scaling factor ρ_n is fixed. Alternatively, as in case (ii.), the optimal rate can also be attained by proper scaling, provided a sufficiently regular prior is used. In all other cases the contraction rate is slower than the minimax rate. Our results are similar to those in [59] in the white noise setting. The extra constraint $\beta + p > 1/2$ that we have in comparison to that work demands an explanation. As (7.10) shows, the size of negligible terms $R_k(f_0)$ in (7.8) decreases as the smoothness $\beta + p$ of the transformed signal Af_0 increases. In order to control R_k , a minimal smoothness of Af_0 is required. The latter is guaranteed if $p + \beta \geq 1/2$, for it is known that in that case Af_0 will be at least continuous, while it may fail to be so if $p + \beta < 1/2$, see [95].

Remark 7.13. The control on $R_k(f_0)$ from (7.9) depends on the fact that the eigenbasis possesses the properties in Assumption 7.11. If instead of Assumption 7.11 (ii.) one only assumes $|\langle \psi_j, \psi_k \rangle| \leq 1$ for any $k \leq n - 1$ and $j \geq n$, the constraint on the smoothness of Af_0 has to be strengthened to $\beta + p \geq 1$ in order to obtain the same results as in Theorem 7.12, because the condition $\beta + p \geq 1$ guarantees that the control on $R_k(f_0)$ in (7.10) remains valid.

Now we consider the extremely ill-posed problem. The following result holds.

Theorem 7.14 (Posterior contraction: extremely ill-posed problem). *Let the problem be extremely ill-posed as $a_k \asymp e^{-pk^s}$ with $s \geq 1$, and let the true parameter $f_0 \in S^\beta$ with $\beta > 0$. Let $\lambda_k = \rho_n^2 k^{-1-2\alpha}$ with $\alpha > 0$ and any positive ρ_n satisfying $\rho_n^2 n \rightarrow \infty$. Then*

$$\sup_{\|f_0\|_\beta \leq K} \mathbb{E}_{f_0} \Pi_n (f : \|f - f_0\|_{H_1} \geq M_n \varepsilon_n | U^n) \rightarrow 0,$$

for any $K > 0$ and $M_n \rightarrow \infty$, where

$$\varepsilon_n = \varepsilon_{n,1} \vee \varepsilon_{n,2} = (\log(\rho_n^2 n))^{-\beta/s} \vee \rho_n (\log(\rho_n^2 n))^{-\alpha/s}. \quad (7.16)$$

In particular,

(i.) if $\rho_n = 1$, then $\varepsilon_n = (\log n)^{-(\alpha \wedge \beta)/s}$,

(ii.) if $n^{-1/2+\delta} \lesssim \rho_n \lesssim (\log n)^{(\alpha-\beta)/s}$ for some $\delta > 0$, then $\varepsilon_n = (\log n)^{-\beta/s}$.

Furthermore, if $\lambda_k = \exp(-\alpha k^s)$ with $\alpha > 0$, the following contraction rate is obtained: $\varepsilon_n = (\log n)^{-\beta/s}$.

Since the frequentist minimax estimation rate in extremely ill-posed problems in the white noise setting is $(\log n)^{-\beta/s}$ (see [15]), Theorem 7.14 shows that the optimal contraction rates can be reached by suitable choice of the regularity of the prior, or by using an appropriate scaling. In contrast to the mildly ill-posed case, we have no extra requirement on the smoothness of Af_0 . The reason is obvious: because the signal is lifted to S^∞ by the forward operator A , the term (7.11) converges to zero exponentially fast, implying that $R_k(f_0)$ in (7.8) is always negligible.

7.3.2 Credible sets

In the Bayesian paradigm, the spread of the posterior distribution is a common measure of uncertainty in parameter estimates. In this section we study the frequentist coverage of Bayesian credible sets in our problem.

When the posterior is Gaussian, it is customary to consider credible sets centered at the posterior mean, which is what we will also do. In addition, because in our case the covariance operator of the posterior distribution does not depend on the data, the radius of the credible ball is determined by the credibility level $1 - \gamma$ and the sample size n . A credible ball centred at the posterior mean \hat{f} from (7.13) is given by

$$\hat{f} + B(r_{n,\gamma}) := \{f \in H_1 : \|f - \hat{f}\|_{H_1} \leq r_{n,\gamma}\}, \quad (7.17)$$

where the radius $r_{n,\gamma}$ is determined by the requirement that

$$\Pi_n(\hat{f} + B(r_{n,\gamma}) | U^n) = 1 - \gamma. \quad (7.18)$$

By definition, the frequentist coverage or confidence of the set (7.17) is

$$\mathbb{P}_{f_0}(f_0 \in \hat{f} + B(r_{n,\gamma})), \quad (7.19)$$

where the probability measure is the one induced by the law of U^n given in (7.8) with $f = f_0$. We are interested in the asymptotic behaviour of the coverage (7.19) as $n \rightarrow \infty$ for a fixed f_0 uniformly in Sobolev balls, and also along a sequence f_0^n changing with n .

The following two theorems hold.

Theorem 7.15 (Credible sets: mildly ill-posed problem). *Assume the same assumptions as in Theorem 7.12 hold, and let $\tilde{\beta} = \beta \wedge (2\alpha + 2p + 1)$. The asymptotic coverage of the credible set (7.17) is*

- (i.) 1, uniformly in $\{f_0 : \|f_0\|_{\beta} \leq 1\}$, if $\rho_n \gg n^{(\alpha - \tilde{\beta})/(2\tilde{\beta} + 2p + 1)}$;
- (ii.) 1, for every fixed $f_0 \in S^\beta$, if $\beta < 2\alpha + 2p + 1$ and $\rho_n \asymp n^{(\alpha - \tilde{\beta})/(2\tilde{\beta} + 2p + 1)}$; c , along some f_0^n with $\sup_n \|f_0^n\|_{\beta} < \infty$, if $\rho_n \asymp n^{(\alpha - \tilde{\beta})/(2\tilde{\beta} + 2p + 1)}$ (any $c \in [0, 1)$).
- (iii.) 0, along some f_0^n with $\sup_n \|f_0^n\|_{\beta} < \infty$, if $\rho_n \ll n^{(\alpha - \tilde{\beta})/(2\tilde{\beta} + 2p + 1)}$.

Theorem 7.16 (Credible sets: extremely ill-posed problem). *Assume the setup of Theorem 7.14. Then if $\lambda_k = \rho_n^2 k^{-1 - 2\alpha}$ with $\alpha > 0$ and any positive ρ_n satisfying $\rho_n^2 n \rightarrow \infty$, the asymptotic coverage of the credible set (7.17) is*

- (i.) 1, uniformly in $\{f_0 : \|f_0\|_{S^\beta} \leq 1\}$, if $\rho_n \gg (\log n)^{(\alpha - \beta)/2}$;
- (ii.) 1, uniformly in f_0 with $\|f_0\|_{\beta} \leq r$ with r small enough;
1, for any fixed $f_0 \in S^\beta$,
provided the condition $\rho_n \asymp (\log n)^{(\alpha - \beta)/s}$ holds;
- (iii.) 0, along some f_0^n with $\sup_n \|f_0^n\|_{\beta} < \infty$, if $\rho_n \lesssim (\log n)^{(\alpha - \beta)/s}$.

Moreover, if $\lambda_k = e^{-\alpha^s}$ with $\alpha > 0$ and any positive ρ_n satisfying $\rho_n^2 n \rightarrow \infty$, the asymptotic coverage of the credible set (7.17) is

- (iv.) 0, for every f_0 such that $|f_{0,i}| \gtrsim e^{-ci^s/2}$ for some $c < \alpha$.

For the two theorems in this section, the most intuitive explanation is offered by the case $\rho_n \equiv 1$. The situations (i.), (ii.) and (iii.) correspond to $\alpha < \beta$, $\alpha = \beta$ and $\alpha > \beta$, respectively. The message is that the oversmoothing prior ((iii.) in Theorem 7.15 and (iii.), (iv.) in Theorem 7.16) leads to disastrous frequentist coverage of credible sets, while the undersmoothing prior ((i.) in both theorems) delivers very conservative frequentist results (coverage 1). With the right regularity of the prior (case (ii.)), the outcome depends on the norm of the true parameter f_0 . Our results are thus similar to those obtained in the white noise setting in [59] and [60].

7.4 Simulation examples

In this section we carry out a small-scale simulation study illustrating our theoretical results. Examples we use to that end are those given in Subsection 7.2.1. These were also used in simulations in [59] and [60].

In the setting of Example 7.7, we use the following true signal,

$$f_0(x) = \sum_{i=1}^{\infty} f_{0,i} \varphi_i(x) \text{ with } f_{0,k} = k^{-3/2} \sin(k). \quad (7.20)$$

It is easy to check that $f_0 \in S^1$.

In the setup of Example 7.8, the initial condition is assumed to be

$$f_0(x) = 4x(x-1)(8x-5). \quad (7.21)$$

One can verify that in this case

$$f_{0,k} = \frac{8\sqrt{2}(13 + 11(-1)^k)}{\pi^3 k^3},$$

and $f_0 \in S^\beta$ for any $\beta < 5/2$.

First, we generate noisy observations $\{Y_i\}_{i=1,\dots,n}$ from our observation scheme (7.3) at design points $x_i = \frac{i-1/2}{n}$ in the case of Volterra operator, and $x_i = i/n$ in the case of the heat equation. Next, we apply the transform described in (7.8) and obtain transformed observations $\{U_i\}_{i=1,\dots,n-1}$. Then, by (7.12), the posterior of the coefficients with the eigenbasis φ_i is given by

$$f_k|U^n \sim \mathcal{N}\left(\frac{na_k\lambda_k\mathbb{1}_{\{k < n\}}}{na_k^2\lambda_k + 1}U_k, \frac{\lambda_k\mathbb{1}_{\{k < n\}}}{na_k^2\lambda_k + 1}\right).$$

Figures 7.1 and 7.2 display plots of (estimated) 95% L_2 -credible bands for different sample sizes and different priors. For all priors we assume $\rho_n \equiv 1$, and use different smoothness degrees α , as shown in the titles of the subplots. In addition, the columns from left to right corresponds to 10^3 , 10^4 and 10^5 observations. The (estimated) credible bands are obtained by generating 1000 realizations from the posterior and retaining 95% of them that are closest in the L^2 -distance to the posterior mean.

Two simulations reflect several similar facts. First, because of the difficulty due to the inverse nature of the problem, the recovery of the true signal is relatively slow, as the posteriors for the sample size 10^3 are still rather diffuse around the true parameter value. Second, it is evident that undersmoothing priors (the top rows in the figures) deliver conservative credible bands, but still capture the truth. On the other hand, oversmoothing priors lead to overconfident, narrow bands, failing to actually express the truth (bottom rows in the figures). As already anticipated due to a greater degree of ill-posedness, recovery of the initial condition in the heat equation case is more difficult than recovery of the true function in the case of the Volterra operator. Finally, we remark that qualitative behaviour of the posterior in our examples is similar to the one observed in [59] and [60]; for larger samples sizes n , discreteness of the observation scheme does not appear to have a noticeably adversary effect compared to the fully observed case in [59] and [60].

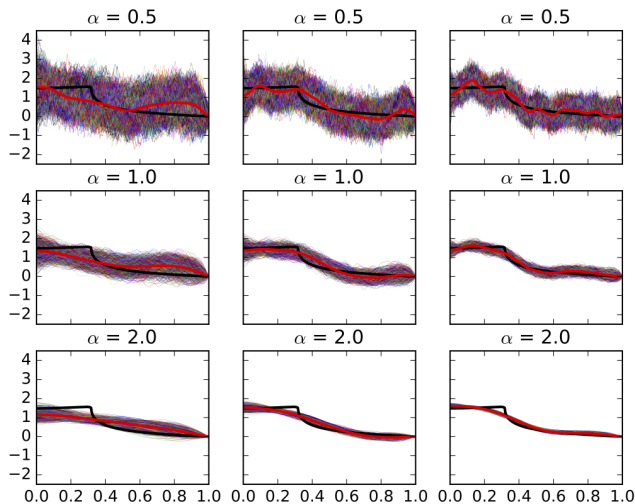


Figure 7.1: Realizations of the posterior mean (red) and 950 of 1000 draws from the posterior (colored thin lines) with smallest L^2 distance to the posterior mean. From left to right columns, the posterior is computed based on sample size 10^3 , 10^4 and 10^5 respectively. The true parameter (black) is of smoothness $\beta = 1$ and given by coefficients $f_{0,k} = k^{-3/2} \sin(k)$.

7.5 Proofs

7.5.1 Proof of Lemma 7.9

This proof is a modification of the one of Lemma 1.7 in [95]. With the following temporary definitions $a := e^{i\pi \frac{i}{n}}$ and $b := e^{i\pi \frac{k}{n}}$, using Euler's formula, we have

$$\begin{aligned}
 \langle \psi_j, \psi_k \rangle_d &= -\frac{1}{2n} \sum_{s=1}^n (a^s - a^{-s})(b^s - b^{-s}) \\
 &= -\frac{1}{2n} \sum_{s=1}^n [(ab)^s - (a/b)^s - (a/b)^{-s} + (ab)^{-s}], \\
 &= -\frac{1}{2n} \left[\underbrace{\sum_{s=1}^n (ab)^s}_A - \underbrace{\sum_{s=1}^n (a/b)^s}_B - \underbrace{\sum_{s=1}^n (a/b)^{-s}}_C + \underbrace{\sum_{s=1}^n (ab)^{-s}}_D \right].
 \end{aligned} \tag{7.22}$$

Furthermore,

$$ab = e^{i\pi \frac{j+k}{n}}, \quad \frac{a}{b} = e^{i\pi \frac{j-k}{n}}.$$

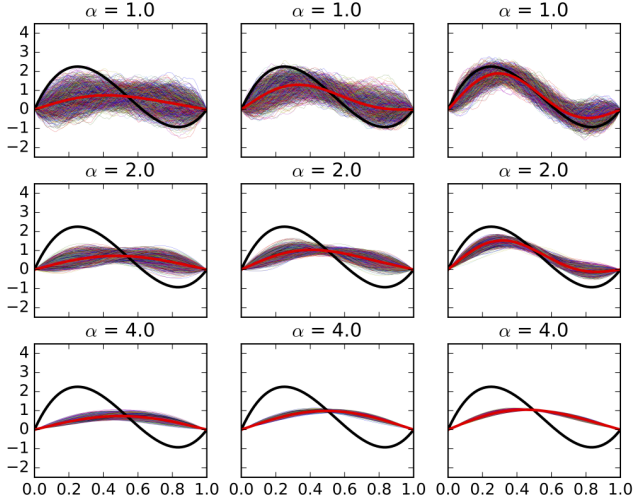


Figure 7.2: Realizations of the posterior mean (red) and 950 of 1000 draws from the posterior (colored thin lines) with smallest L^2 distance to the posterior mean. From left to right columns, the posterior is computed based on sample size 10^3 , 10^4 and 10^5 respectively. The true parameter (black) is of smoothness β for any $\beta < 5/2$ and given by (7.21).

Observe that when $ab \neq 1$, we have

$$A = \frac{ab(1 - (ab)^n)}{1 - ab}, \quad D = \frac{1 - (ab)^{-n}}{ab - 1}, \quad A + D = \frac{ab(1 - (ab)^n) - (1 - (ab)^{-n})}{1 - ab}.$$

Similarly, if $a/b \neq 1$,

$$B + C = \frac{(a/b)(1 - (a/b)^n) - (1 - (a/b)^{-n})}{1 - (a/b)}.$$

We fix $1 \leq k \leq n - 1$ and discuss different situations depending on j .

(I.) $1 \leq j \leq n - 1$ and $j + k \neq n$.

Since $n \neq j + k < 2n$, we always have $ab = e^{i\pi \frac{j+k}{n}} \neq 1$, and the terms A and D can be calculated as above. Similarly, since $-n < j - k < n$, $a/b = 1$ only when $j = k$. Moreover, $j + k$ and $j - k$ have the same parity, and so $j = k$ is only possible if $j + k$ is even.

(i.) $j + k$ is even.

In this case, $(ab)^n = 1$. This leads to $A = D = 0$.

Further, if $j = k$, we have $a/b = b/a = 1$ and $B = C = n$. Otherwise, if $j \neq k$, we have $a/b \neq 1$ and $(a/b)^n = 1 = (b/a)^n$ (since $j - k$ is even),

and so

$$B = \frac{a/b(1 - (a/b)^n)}{1 - a/b} = 0, \quad C = 0,$$

which implies (7.22) equals 1.

(ii.) $j + k$ is odd. We have $(ab)^n = (a/b)^n = -1$, which results in $A + D = B + C = -2$, and so (7.22) equals 0.

(II.) $1 \leq j < n$ and $j + k = n$. We have $ab = -1$. Arguing as above, if n is odd, $A + D = -2$ and $B + C = -2$. If n is even, $A = D = 0$ and $B = C = n\delta_{jk}$.

The remaining cases follow the same arguments, and hence we omit the (lengthy and elementary) calculations.

(III.) $j = ln$ with $l \in \mathbb{N}$.

It can be shown that $A + D = B + C$ always holds.

(IV.) $j \in \{ln + 1, \dots, (l + 1)n - 1\}$.

When l is even, one obtains $\langle \psi_j, \psi_k \rangle_d = \delta_{\tilde{j}k}$, where $\tilde{j} = j - ln$. Otherwise, for odd l , $\langle \psi_j, \psi_k \rangle_d = -\delta_{\tilde{j}k}$ where $\tilde{j} = (l + 1)n - j$.

7.5.2 Proof of Theorem 7.12

In this proof we use the notation $\|\cdot\| = \|\cdot\|_{H_1} = \|\cdot\|_{\ell^2}$. To show

$$\sup_{\|f_0\|_{\beta} \leq K} \mathbb{E}_{f_0} \Pi_n (f : \|f - f_0\| \geq M_n \varepsilon_n | U^n) \rightarrow 0,$$

we first apply Markov's inequality,

$$M_n^2 \varepsilon_n^2 \Pi_n (f : \|f - f_0\|^2 \geq M_n^2 \varepsilon_n^2 | U^n) \leq \int \|f - f_0\|^2 d\Pi_n(f | U^n).$$

From (7.12) and the bias-variance decomposition,

$$\int \|f - f_0\|^2 d\Pi_n(f | U^n) = \|\hat{f} - f_0\|^2 + \|\sigma\|^2,$$

where $\sigma = (\sigma_k)_k$ is given in (7.12). Because σ is deterministic,

$$\mathbb{E}_{f_0} [\Pi_n (f : \|f - f_0\| \geq M_n \varepsilon_n | U^n)] \leq \frac{1}{M_n^2 \varepsilon_n^2} \left(\mathbb{E}_{f_0} \|\hat{f} - f_0\|^2 + \|\sigma\|^2 \right).$$

Since $M_n \rightarrow \infty$ is assumed, it suffices to show that the terms in brackets are bounded by a constant multiple of ε_n^2 uniformly in f_0 in the Sobolev ellipsoid.

Using (7.14), we obtain

$$\mathbb{E}_{f_0} \|\hat{f} - f_0\|^2 = \|\mathbb{E}_{f_0} \hat{f} - f_0\|^2 + \|\tau\|^2 = \|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 + \|f_0^n\|^2 + \|\tau\|^2,$$

where $\tau = (\tau_k)_k$ given in (7.14) and

$$\begin{aligned} f_0^n &= (f_{0,1}, \dots, f_{0,n-1}, 0, \dots), \\ f_0^r &= (0, \dots, 0, f_{0,n}, f_{0,n+2}, \dots). \end{aligned}$$

We need to obtain a uniform upper bound over the ellipsoid $\{f_0 : \|f_0\|_\beta \leq K\}$ for

$$\|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 + \|f_0^r\|^2 + \|\tau\|^2 + \|\sigma\|^2. \quad (7.23)$$

We have

$$\begin{aligned} \|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 &= \sum_{k=1}^{n-1} \left(\frac{na_k^2 \lambda_k}{na_k^2 \lambda_k + 1} f_{0,k} + \frac{na_k \lambda_k}{na_k^2 \lambda_k + 1} R_k - f_{0,k} \right)^2 \\ &\lesssim \underbrace{\sum_{k=1}^{n-1} \frac{1}{(na_k^2 \lambda_k + 1)^2} f_{0,k}^2}_{A_1} + n \sup_{k < n} R_k^2 \underbrace{\sum_{k=1}^{n-1} \frac{na_k^2 \lambda_k^2}{(na_k^2 \lambda_k + 1)^2}}_{A_2}, \end{aligned} \quad (7.24)$$

and

$$\|f_0^r\|^2 = \sum_{k \geq n} f_{0,k}^2, \quad \|\tau\|^2 = \sum_{k=1}^{n-1} \frac{na_k^2 \lambda_k^2}{(na_k^2 \lambda_k + 1)^2} = A_2, \quad \|\sigma\|^2 = \sum_{k=1}^{n-1} \frac{\lambda_k}{na_k^2 \lambda_k + 1}.$$

Recall that we write (7.15) as $\varepsilon_n = \varepsilon_{n,1} \vee \varepsilon_{n,2}$. The statements (i.)–(iii.) follow by elementary calculations. Specifically, in (ii.) the given ρ_n is the best scaling, as it gives the fastest rate. From [59] (see the argument below (7.3) on page 21), A_1 is bounded by a fixed multiple of $(\varepsilon_{n,1})^2$, and $\|\tau\|^2, \|\sigma\|^2$ are bounded by multiples of $(\varepsilon_{n,2})^2$. Hence, to show that the rate is indeed (7.15), it suffices to show that $n \sup_{k < n} R_k^2 A_2$ and $\|f_0^r\|^2$ can be bounded by a multiple of $(\varepsilon_n)^2$ uniformly in the ellipsoid $\{f_0 : \|f_0\|_\beta \leq K\}$. Since $A_2 = \|\tau\|^2$, to that end it is sufficient to show that $\sup_{k < n} n R_k^2 = O(1)$, and that $\|f_0^r\|^2 = O(\varepsilon_n)^2$.

Since $f_0 \in S^\beta$, we have the following straightforward bound,

$$\|f_0^r\|^2 \leq n^{-2\beta} \sum_{k \geq n} f_{0,k}^2 k^{2\beta} \leq n^{-2\beta} \|f_0\|_\beta^2 \lesssim n^{-2\beta},$$

which is uniform in $\{f_0 : \|f_0\|_\beta \leq K\}$. By comparing to the rates in the statements (ii.)–(iii.), it is easy to see that $n^{-2\beta}$ is always negligible with respect to ε_n^2 .

Proving $\sup_{k < n} n R_k^2 = O(1)$ is equivalent to showing $\sup_{k < n} R_k^2 = O(1/n)$; but the latter has been already proved in (7.10). Notice that we actually obtained a sharper bound $\sup_{k < n} n R_k^2 = o(1)$ than the one necessary for our purposes in this proof. However, this sharper bound will be used in the proof of Theorem 7.15. By taking supremum over f_0 , we thus have

$$\sup_{\|f_0\|_{S^\beta} \leq K} \left(\|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 + \|f_0^r\|^2 \right) \lesssim \varepsilon_n^2 + n^{-2\beta} \lesssim \varepsilon_n^2, \quad (7.25)$$

with which we conclude that up to a multiplicative constant, (7.23) is bounded by ε_n^2 uniformly over the ellipsoid $\sup_{\|f_0\|_\beta \leq K}$. This completes the proof.

7.5.3 Proof of Theorem 7.14

We start by generalizing Theorem 3.1 in [60]. Following the same lines as in the proof of that theorem and using Lemma 7.17, 7.18, 7.19, 7.20 in Section 7.6 of the present paper instead of analogous technical results in [60], the statement of Theorem 3.1 in [60] can be extended from $s = 2$ to a general $s \geq 1$, for which the posterior rate is given by (7.16), or $\varepsilon_n = \varepsilon_{n,1} \vee \varepsilon_{n,2}$ in short.

In our model, we again obtain (7.23) and also that a fixed multiple of $(\varepsilon_{n,1})^2$ is an upper bound of A_1 , and that $\|\tau\|^2, \|\sigma\|^2$ can be bounded from above by fixed multiples of $(\varepsilon_{n,1})^2$.

Now as in the proof of Theorem 7.12 in Section 7.5.2, we will show that $\sup_{\|f_0\|_\beta \leq K} (\|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 + \|f_0^r\|^2)$ can be bounded by a fixed multiple of $(\varepsilon_n)^2$ by proving that $\sup_{k \leq n} nR_k^2 = O(1)$. By (7.11), $n(R_k)^2 \lesssim \exp(-2pn)n$, and the righthand side converges to zero. Therefore,

$$\sup_{\|f_0\|_\beta \leq K} \left(\|\mathbb{E}_{f_0} \hat{f} - f_0^n\|^2 + \|f_0^r\|^2 \right) \lesssim \varepsilon_n.$$

Parts (i.) and (ii.) of the statement of the theorem are obtained by direct substitutions, using the fact that $\log n \ll n$. Notice that if $\rho_n \gtrsim (\log n)^{(\alpha-\beta)/s}$, the rate ε_n deteriorates and is dominated by the second term in (7.16).

For the case $\lambda_k = \exp(-\alpha k^s)$, the argument follows the same lines as in Section 5.1 in [60], and our arguments above.

7.5.4 Proof of Theorem 7.15

The proof runs along the same lines as the proof of Theorem 4.2 in [59]. We will only show the main steps here.

In Section 7.2.2, we have shown that the posterior distribution is $\otimes_{k \in \mathbb{N}} \mathcal{N}(\hat{f}_k, \sigma_k^2)$, the radius $r_{n,\gamma}$ in (7.17) satisfies $\mathbb{P}_{X_n}(X_n < r_{n,\gamma}^2) = 1 - \gamma$, where X_n is a random variable distributed as the square norm of an $\otimes_{k \in \mathbb{N}} \mathcal{N}(\hat{f}_k, \sigma_k^2)$ variable. Let $T = (\tau_k^2)_{k \in \mathbb{N}}$. Under (7.8), the variable \hat{f} is distributed as $\mathcal{N}_{H_1}(\mathbb{E}_{f_0} \hat{f}, T) := \otimes_{k \in \mathbb{N}} \mathcal{N}(\mathbb{E}_{f_0} \hat{f}_k, \tau_k^2)$. Hence the coverage (7.19) can be rewritten as

$$\mathbb{P}_{W_n}(\|W_n + \mathbb{E}_{f_0} \hat{f} - f_0\|_{H_1} \leq r_{n,\gamma}), \quad (7.26)$$

where $W_n \sim \mathcal{N}_{H_1}(0, T)$. Denote $V_n = \|W_n\|_{H_1}^2$ and observe that one has in distribution

$$X_n = \sum_{1 \leq i < n} \sigma_i^2 Z_i^2, \quad V_n = \sum_{1 \leq i < n} \tau_i^2 Z_i^2$$

for $\{Z_i\}$ independent standard Gaussian random variables with

$$\sigma_i^2 = \frac{\lambda_i}{na_i^2 \lambda_i + 1}, \quad \tau_i^2 = \frac{na_i^2 \lambda_i^2}{(na_i^2 \lambda_i + 1)^2}.$$

By the same argument as in [59], one can show that the standard deviations of X_n and V_n are negligible with respect to their means,

$$\mathbb{E}X_n \asymp \rho_n^2 (\rho_n^2 n)^{-2\alpha/(2\alpha+2p+1)}, \quad \mathbb{E}V_n \asymp \rho_n^2 (\rho_n^2 n)^{-2\alpha/(2\alpha+2p+1)}, \quad (7.27)$$

and the difference of their means,

$$\mathbb{E}(X_n - V_n) \asymp \rho_n^2 (\rho_n^2 n)^{-2\alpha/(2\alpha+2p+1)}.$$

Since $X_n \geq V_n$, the distributions of X_n and V_n are asymptotically separated, i.e. $\mathbb{P}(V_n \leq v_n \leq X_n) \rightarrow 1$ for some v_n , e.g. $v_n = \mathbb{E}(V_n + X_n)/2$. Since $r_{n,\gamma}^2$ are $1 - \gamma$ quantiles of X_n , we also have $\mathbb{P}(V_n \leq r_{n,\gamma}^2(1 + o(1))) \rightarrow 1$. In addition, by (7.27),

$$r_{n,\gamma}^2 \asymp \rho_n^2 (\rho_n^2 n)^{-2\alpha/(2\alpha+2p+1)}.$$

Introduce

$$B_n := \sup_{\|f_0\|_\beta \lesssim 1} \|\mathbb{E}_{f_0} \hat{f} - f_0\|_{H_1} = \sup_{\|f_0\|_\beta \lesssim 1} \left(\|\mathbb{E}_{f_0} \hat{f} - f_0^n\|_{H_1} + \|f_0^n\|_{H_1} \right). \quad (7.28)$$

It follows from the arguments for (7.10) in the proof of Theorem 7.12 that

$$B_n \lesssim \varepsilon_{n,1} \vee (\sqrt{n} R \varepsilon_{n,2}),$$

where $R = \sup_{k < n} R_k \lesssim n^{-(p+\beta)}$. Now apply the argument on the lower bound from Lemma 8.1 in [59] (with $q = \beta, t = 0, u = 2\alpha + 2p + 1, v = 2, N = \rho_n^2 n$) to obtain that $B_n \gtrsim \varepsilon_{n,1}$. Thus we have

$$\varepsilon_{n,1} \lesssim B_n \lesssim \varepsilon_{n,1} \vee (\sqrt{n} R \varepsilon_{n,2}).$$

We consider separate cases. In case (i.), substituting the corresponding ρ_n into the expression of $\varepsilon_{n,1}$ and $\varepsilon_{n,2}$, we have $\varepsilon_{n,1} \ll \varepsilon_{n,2}$. By (7.10), $B_n \lesssim \varepsilon_{n,1} \vee (\sqrt{n} R \varepsilon_{n,2}) \ll \varepsilon_{n,2} \asymp r_{n,\gamma}$. This leads to

$$\begin{aligned} \mathbb{P}(\|W_n + \mathbb{E}_{f_0} \hat{f} - f_0\|_{H_1} \leq r_{n,\gamma}) &\geq \mathbb{P}(\|W_n\|_{H_1} \leq r_{n,\gamma} - B_n) \\ &= \mathbb{P}(V_n \leq r_{n,\gamma}^2(1 + o(1))) \rightarrow 1 \end{aligned} \quad (7.29)$$

uniformly in the set $\{f_0 : \|f_0\|_\beta \lesssim 1\}$.

In case (iii.), the given ρ_n leads to $\varepsilon_{n,1} \gg \varepsilon_{n,2}$ and consequently $B_n \gg r_{n,\gamma}$. Hence,

$$\mathbb{P}(\|W_n + \mathbb{E}_{f_0} \hat{f}^n - f_0^n\|_{H_1} \leq r_{n,\gamma}) \leq \mathbb{P}(\|W_n\|_{H_1} \geq B_n - r_{n,\gamma}) \rightarrow 0,$$

for any f_0^n (nearly) attaining the supremum.

In case (ii.), we have $B_n \asymp r_{n,\gamma}$. If $\beta < 2\alpha + 2p + 1$, by Lemma 8.1 in [59] the bias $\mathbb{E}_{f_0} \hat{f} - f_0$ at a fixed f_0 is of strictly smaller order than B_n . Following the argument of case (i.), the asymptotic coverage can be shown to converge to 1.

For existence of a sequence along which the coverage is $c \in [0, 1)$, we only give a sketch of the proof here; the details can be filled in as in [59].

The coverage (7.26) with f_0 replaced by f_0^n tends to c , if for $b_n = \mathbb{E}_{f_0} \hat{f}^n - f_0^n$ and z_c a standard normal quantile,

$$\frac{\|W_n + b_n\|_{H_1}^2 - \mathbb{E}\|W_n + b_n\|_{H_1}^2}{\text{sd}\|W_n + b_n\|_{H_1}^2} \rightsquigarrow \mathcal{N}(0, 1), \quad (7.30)$$

$$\frac{r_{n,\gamma}^2 - \mathbb{E}\|W_n + b_n\|_{H_1}^2}{\text{sd}\|W_n + b_n\|_{H_1}^2} \rightarrow z_c, \quad (7.31)$$

Since W_n is centred Gaussian $\mathcal{N}_{H_1}(0, T)$, (7.31) can be expressed as

$$\frac{r_{n,\gamma}^2 - \mathbb{E}V_n - \sum_{i=1}^{n-1} b_{n,i}^2}{\sqrt{\text{var } V_n + 4 \sum_{i=1}^{n-1} \tau_{i,n}^2 b_{n,i}^2}} \rightarrow z_c. \quad (7.32)$$

Here $\{b_{n,i}\}$ has exactly one nonzero entry depending on the smoothness cases $\beta \leq 2\alpha + 2p + 1$ and $\beta > 2\alpha + 2p + 1$. The nonzero entry, which we call b_{n,i_n} , has the following representation, with d_n to be yet determined,

$$b_{n,i_n}^2 = r_{n,\gamma}^2 - \mathbb{E}V_n - d_n \text{sd } V_n.$$

Since $r_{n,\gamma}^2, \mathbb{E}V_n$ and $r_{n,\gamma}^2 - \mathbb{E}V_n$ have the same order and $\text{sd } V_n$ is of strictly smaller order, one can show that the lefthand side of (7.32) is equivalent to

$$\frac{d_n \text{sd } V_n}{\sqrt{\text{var } V_n + 4\tau_{i_n,n}^2 (r_{n,\gamma}^2 - \mathbb{E}V_n)(1 + o(1))}},$$

for bounded or slowly diverging d_n . Then (7.32) can be obtained by discussing different smoothness cases separately, by a suitable choice of i_n, d_n .

To prove the asymptotic normality in (7.30), the numerator can be written as

$$\|W_n + b_n\|_{H_1}^2 - \mathbb{E}\|W_n + b_n\|_{H_1}^2 = \sum_i \tau_{i,n}^2 (Z_i^2 - 1) + 2b_{n,i_n} \tau_{i_n,n} Z_{i_n}.$$

Next one applies the arguments as in [59].

7.5.5 Proof of Theorem 7.16

This proof is almost identical to the proof of Theorem 2.2 in [60]. We supply the main steps.

Following the same arguments as in the proof of Theorem 7.15, we obtain

$$\begin{aligned} \mathbb{E}X_n &\asymp \rho_n^2 (\log(\rho_n^2 n))^{-2\alpha/s} \gg \text{sd } X_n \asymp \rho_n^2 (\log(\rho_n^2 n))^{-1/(2s)-2\alpha/s}, \\ \mathbb{E}V_n &\asymp \rho_n^2 (\log(\rho_n^2 n))^{-1/s-2\alpha/s} \asymp \text{sd } V_n, \end{aligned}$$

as in the proof of Theorem 2.2 in [60]. This leads to

$$r_{n,\gamma}^2 \asymp \rho_n^2 (\log(\rho_n^2 n))^{-2\alpha/s},$$

and furthermore,

$$\mathbb{P}(V_n \leq \delta r_{n,\gamma}^2) = \mathbb{P}\left(\frac{V_n - \mathbb{E}V_n}{\text{sd } V_n} \leq \frac{\delta r_{n,\gamma}^2 - \mathbb{E}V_n}{\text{sd } V_n}\right) \rightarrow 1,$$

for every $\delta > 0$.

Similar to Theorem 7.15, the bounds on the square norm B_n (defined in (7.28)) of the bias are known: upper bound from the proof of Theorem 7.14, and lower bound from Lemma 7.17,

$$\varepsilon_{n,1} \lesssim B_n \lesssim \varepsilon_{n,1} \vee (\sqrt{n}R\varepsilon_{n,2}),$$

where $\varepsilon_{n,1}, \varepsilon_{n,2}$ are given in (7.16), and $\sqrt{n}R$ satisfies the bound (7.11).

In case (i.), $B_n \ll r_{n,\gamma}$, and hence (7.29) applies. The rest of the results can be obtained in a similar manner.

7.6 Auxiliary lemmas

The following lemmas are direct generalisations of the case $s = 2$ in the Appendix of [60] to a general s . They can be easily proved by simple adjustments of the original proofs in [60], and we only state the results.

Lemma 7.17 (Lemma 6.1 in [60]). *For $q \in \mathbb{R}$, $u \geq 0$, $v > 0$, $t + 2q \geq 0$, $p > 0$, $0 \leq r < pv$ and $s \geq 1$,*

$$\sup_{\|f\|_{S^q} \leq 1} \sum_{i=1}^{\infty} \frac{f_i^2 i^{-t} e^{-ri^s}}{(1 + Ni^{-u} e^{-pi^s})^v} \asymp N^{-r/p} (\log N)^{-t/s - 2q/s + ru/ps},$$

as $N \rightarrow \infty$.

In addition, for any fixed $f \in S^q$,

$$N^{r/p} (\log N)^{t/s + 2q/s - ru/ps} \sum_{i=1}^{\infty} \frac{f_i^2 i^{-t} e^{-ri^s}}{(1 + Ni^{-u} e^{-pi^s})^v} \rightarrow 0,$$

as $N \rightarrow \infty$.

Lemma 7.18 (Lemma 6.2 in [60]). *For $t, u \geq 0$, $v > 0$, $p > 0$, $0 < r < vp$ and $s \geq 1$, as $N \rightarrow \infty$,*

$$\sum_{i=1}^{\infty} \frac{i^{-t} e^{-ri^s}}{(1 + Ni^{-u} e^{-pi^s})^v} \asymp N^{-r/p} (\log N)^{-t/s + ru/ps}.$$

If $r = 0$ and $t > 1$, while other assumptions remain unchanged,

$$\sum_{i=1}^{\infty} \frac{i^{-t} e^{-ri^s}}{(1 + Ni^{-u} e^{-pi^s})^v} \asymp (\log N)^{-(t+1)/s}.$$

Lemma 7.19 (Lemma 6.4 in [60]). *Assume $s \geq 1$. Let I_N be the solution in i to $Ni^{-u} e^{-pi^s} = 1$, for $u \geq 0$ and $p > 0$. Then*

$$I_N \sim \left(\frac{1}{p} \log N \right)^{1/s}$$

Lemma 7.20 (Lemma 6.5 in [60]). *Let $s \geq 1$. As $K \rightarrow \infty$, we have*

(i.) for $a > 0$ and $b \in \mathbb{R}$,

$$\int_1^K e^{ax^s} x^b dx \sim \frac{1}{as} e^{aK^s} K^{b-s+1},$$

(ii.) for $a, b, K > 0$,

$$\int_K^\infty e^{-ax^s} x^{-b} dx \leq \frac{1}{as} e^{-aK^s} K^{-b-s+1}.$$

Chapter 8

Inverse Problems with Discrete Observations in Smoothness Scales

In this chapter, we continue the study in Chapter 6. In practice one usually does not have access to a ‘continuous’ observation Y , but only records noisy samples of the unknown function $\mathcal{A}f$ at a finite number of locations in its domain. This is the situation that we consider in the present chapter. To cope with discrete observations, we consider the operator \mathcal{A} introduced in Section 5.2 with a specified domain. Assume that $\mathcal{A} : H \rightarrow G$ is a linear mapping from a Hilbert space H into a pre-Hilbert space $G = \mathcal{L}^2(\mathfrak{D})$ of square-integrable functions $g : \mathfrak{D} \rightarrow \mathbb{R}^d$ on a bounded domain $\mathfrak{D} \subset \mathbb{R}^d$. Then, the observation is formed as follows. For a given set of design points

$$\mathfrak{D}_n := \{x_1, \dots, x_n\} \subset \mathfrak{D}. \quad (8.1)$$

we observe the vector $Y^n = (Y_1, \dots, Y_n)$ defined by

$$Y_i = \mathcal{A}f(x_i) + Z_i, \quad i = 1, \dots, n, \quad (8.2)$$

with Z_i i.i.d. standard normal random variables. We wish to estimate f from the observations $(x_i, Y_i)_{1 \leq i \leq n}$.

This chapter extends results for the white noise model obtained in Chapter 6 to the case of discrete observations. Although we repeat some necessary definitions, we refer to the mentioned chapter for further examples and discussion. This chapter is organized as follows. Section 8.1 demonstrates a procedure to reconstruct continuous signals from discrete observations. With the help of the just mentioned reconstruction procedure and the projection method from Section 5.3, we present a general contraction theorem for the regression model in Section 8.2. Then, the same priors considered in Chapter 6 are studied: series priors and Gaussian priors in Section 8.3 and Section 8.4, respectively. In addition, Gaussian mixture priors are introduced to obtain adaptation in Section 8.5. In the end, Section 8.6 contains the proofs.

8.1 Signal Reconstruction

The sampling scheme (8.2) collects discrete data, but the operator \mathcal{A} acts on a continuous function, which we wish to estimate on its full domain. In this section we describe an interpolation technique that maps discrete signals to the continuous domain. A similar technique has been used in the context of proving asymptotic equivalence of the white noise model and nonparametric regression; see [82] and the references therein. The assumptions are also inspired by the theory from the field of numerical analysis, see [12].

The range space G of the operator $\mathcal{A} : H \rightarrow G$ is a collection of functions $g : \mathcal{D} \rightarrow \mathbb{R}$, equipped with a pre-inner product $\langle \cdot, \cdot \rangle$. The design points (8.1) give rise to a “discrete” semi-inner product and semi-norm on G given by

$$\langle g, h \rangle_n := \frac{1}{n} \sum_{i=1}^n g(x_i)h(x_i) \quad \|g\|_n := \sqrt{\frac{1}{n} \sum_{i=1}^n g^2(x_i)}.$$

(The notation $\|\cdot\|_n$ clashes with the notation $\|\cdot\|_s$ for the norms of the smoothness scales, but this should not lead to confusion as n will never appear as smoothness level.)

For every $n \in \mathbb{N}$ we fix an n -dimensional subspace $\widetilde{W}_n \subset G$ with two properties.

Assumption 8.1 (Interpolation).

- (i) There exist constants $0 < C_1 < C_2 < \infty$, independent of n , such that

$$C_1 \|w\| \leq \|w\|_n \leq C_2 \|w\|, \quad w \in \widetilde{W}_n. \quad (8.3)$$

- (ii) For every $g \in G$ the unique element $\mathcal{I}_n g$ of \widetilde{W}_n that interpolates g at the design points, i.e. $\mathcal{I}_n g(x_i) = g(x_i)$, for every $i = 1, \dots, n$, satisfies, for every s in some interval (s_d, S_d) ,

$$\|\mathcal{I}_n g - g\| \lesssim \delta_d(n, s) \|g\|_s. \quad (8.4)$$

Condition (8.3) requires that the discrete and continuous norms be equivalent on the subspace \widetilde{W}_n , whereas (8.4) ensures that the subspaces \widetilde{W}_n have good approximation properties under discretization for smooth functions. In this condition $\|\cdot\|_s$ are the norms of a smoothness scale $(G_s)_{s \in \mathbb{R}}$ as in Definition 2.1, in which the space G is embedded as $G = G_0$, and the approximation numbers $\delta_d(n, s)$ will often be the same as the approximation rates $\delta(n, s)$ in Assumption 2.3. However, the approximation (8.4) is typically not true for every smoothness level $s > 0$, but only for s in a range (s_d, S_d) . For instance, for Sobolev scales the lower bound s_d is typically equal to $d/2$ for d the dimension of the domain \mathcal{D} of the functions in G , and the upper bound S_d is the regularity of the basis elements used to define the scale.

Condition (8.3) implies that the set \widetilde{W}_n is also n -dimensional over the design points, so that the interpolation $\mathcal{I}_n g$ indeed exists and is unique. In Lemma 8.3 it

will be seen to be also the orthogonal projection of $g \in G$ onto $\widetilde{W}_n \subset G$ relative to the *discrete* inner product $\langle \cdot, \cdot \rangle_n$.

Fix an arbitrary orthonormal basis $e_{1,n}, \dots, e_{n,n}$ of \widetilde{W}_n relative to the discrete inner product $\langle \cdot, \cdot \rangle_n$, and given the discrete data (Y_1, \dots, Y_n) as in (8.2), define

$$Y^{(n)} = \sum_{i=1}^n \frac{1}{n} \sum_{j=1}^n Y_j e_{i,n}(x_j) e_{i,n}. \quad (8.5)$$

This embeds the discrete data as a ‘continuous signal’ into the space $\widetilde{W}_n \subset G$. If the observations satisfy (8.2), then the continuous observation $Y^{(n)}$ can be decomposed as

$$Y^{(n)} = \sum_{i=1}^n \langle Af, e_{i,n} \rangle_n e_{i,n} + \frac{1}{n} \sum_{j=1}^n Z_j \sum_{i=1}^n e_{i,n}(x_j) e_{i,n} = \mathcal{I}_n Af + \frac{1}{\sqrt{n}} \xi^{(n)},$$

where $\xi^{(n)}$ is a Gaussian random variable with values in the space in \widetilde{W}_n . Loosely speaking, as $n \rightarrow \infty$ the operators \mathcal{I}_n should tend to the identity operator, and the mean of the signal $Y^{(n)}$ should become more representative of the full signal Af . As shown in the following lemma the noise $\xi^{(n)}$ remains bounded as $n \rightarrow \infty$.

Lemma 8.2. *The variable $\xi^{(n)}$ defined in the preceding display is a Gaussian random element in $\widetilde{W}_n \subset G$ with mean zero. Under (8.3) its covariance operator is up to multiplicative constants that do not depend on n bounded below and above by the orthogonal projection $\tilde{Q}_n : G \rightarrow \widetilde{W}_n$ relative to the continuous inner product $\langle \cdot, \cdot \rangle$.*

Proof. For $g \in G$ we can write $\langle \xi^{(n)}, g \rangle = n^{-1/2} \sum_{j=1}^n Z_j \sum_{i=1}^n e_{i,n}(x_j) \langle e_{i,n}, g \rangle$. Clearly the expectation of this variable vanishes, while the variance is given by

$$\begin{aligned} \text{Var}(\langle \xi^{(n)}, g \rangle) &= \frac{1}{n} \sum_{j=1}^n \left(\sum_{i=1}^n e_{i,n}(x_j) \langle e_{i,n}, g \rangle \right)^2 \\ &= \frac{1}{n} \sum_{k=1}^n \sum_{l=1}^n \sum_{j=1}^n e_{k,n}(x_j) \langle e_{k,n}, g \rangle e_{l,n}(x_j) \langle e_{l,n}, g \rangle \\ &= \sum_{k=1}^n \sum_{l=1}^n \langle e_{k,n}, e_{l,n} \rangle_n \langle e_{k,n}, g \rangle \langle e_{l,n}, g \rangle = \left\| \sum_{k=1}^n \langle e_{k,n}, g \rangle e_{k,n} \right\|_n^2. \end{aligned}$$

The right side is the square of the norm $\|g_n\|$ of the vector $g_n = (g_{1,n}, \dots, g_{n,n})$ of continuous coefficients $g_{j,n} = \langle g, e_{j,n} \rangle$ of g relative to the discrete basis. The orthogonal projection $\tilde{Q}_n g$ of g onto \widetilde{W}_n relative to the continuous inner product can be written in terms of the discrete basis $e_{1,n}, \dots, e_{n,n}$ as $\tilde{Q}_n g = \sum_{i=1}^n \alpha_i e_{i,n}$, for $\alpha = (\alpha_1, \dots, \alpha_n)^T = \Sigma_n^{-1} g$ and Σ_n the Gram matrix $(\langle e_{i,n}, e_{j,n} \rangle)$. Hence $\|\tilde{Q}_n g\|^2 = \alpha^T \Sigma_n \alpha = g_n^T \Sigma_n^{-1} g_n$. Because Σ_n is bounded above and below by a multiple of the identity, by Lemma 8.3 below, it follows that $\|\tilde{Q}_n g\| \simeq \|g_n\|$. \square

Lemma 8.3. *Suppose that (8.3) holds. Then*

- (i) For every $g \in G$ the function $\mathcal{I}_n g$ is the orthogonal projection of g onto $\widetilde{W}_n \subset G$ relative to the discrete inner product $\langle \cdot, \cdot \rangle_n$.
- (ii) $C_1 \|\mathcal{I}_n g\| \leq \|g\|_n \leq C_2 \|\mathcal{I}_n g\|$, for every $g \in G$.
- (iii) The Gram matrix $(\langle e_{i,n}, e_{j,n} \rangle)_{i,j=1..n}$ of any basis $e_{1,n}, \dots, e_{n,n}$ of \widetilde{W}_n that is orthonormal relative to the discrete inner product $\langle \cdot, \cdot \rangle_n$ is bounded below and above by the identity, up to multiplicative constants that do not depend on n .

Proof. For an arbitrary orthonormal basis $e_{1,n}, \dots, e_{n,n}$ of \widetilde{W}_n relative to the discrete inner product $\langle \cdot, \cdot \rangle_n$, the matrix $(e_{j,n}(x_i))_{i,j=1..n}$ has orthogonal columns of Euclidean length \sqrt{n} and hence can be represented as $\sqrt{n}O$ for an orthogonal matrix O . The interpolation in \widetilde{W}_n of a function g at the design points is the function $\sum_{j=1}^n \alpha_j e_{j,n}$ with the coefficients $\alpha = (\alpha_1, \dots, \alpha_n)^T$ satisfying $\sqrt{n}O\alpha = (g(x_i))_{i=1..n}$. The unique solution to the latter equation is $\alpha = (1/\sqrt{n})O^T(g(x_i))_{i=1..n}$, which can be seen to be equal to $(\langle g, e_{i,n} \rangle_n)_{i=1..n}$. Thus the interpolation is indeed the orthogonal projection $\mathcal{I}_n g = \sum_{i=1}^n \langle g, e_{i,n} \rangle_n e_{i,n}$.

As the functions g and $\mathcal{I}_n g$ coincide at the design points, clearly $\|g\|_n = \|\mathcal{I}_n g\|_n$, for any $g \in G$, and this is equivalent to the continuous norm $\|\mathcal{I}_n g\|_n$, by (8.3).

To prove the third statement note that the square discrete and continuous norms of $\sum_{i=1}^n \alpha_i e_{i,n}$ are given by $\alpha^T \alpha$ and $\alpha^T \Sigma \alpha$, respectively, for Σ the Gram matrix of the basis functions $e_{i,n}$ relative to the continuous inner product. By (8.3) these norms are proportional and hence the eigenvalues of Σ are bounded from below and above. \square

The following two examples exhibit suitable discretization spaces, both with equidistant design points.

Example 8.4 (Trigonometric Polynomials). This example is adapted from Section 2.3 in [82]. Let $\mathcal{D} = \mathfrak{I}^d = (0, 1]^d$, for $d \in \mathbb{N}$, and consider the set of $n = m^d$ design points $\mathcal{D}_n = \{k/m\}_{k \in \{1, \dots, m\}^d}$, for a given odd natural number m . In this case, the Fourier system with $i = \sqrt{-1}$,

$$e_k(x) = e^{i2\pi \langle k, x \rangle_{\mathbb{R}^d}}, \quad k = (k_1, \dots, k_d) \in \mathbb{Z}^d,$$

is not only orthonormal in the continuous space $L^2(\mathfrak{I}^d)$, but also with respect to the discrete inner product $\langle \cdot, \cdot \rangle_n$, i.e.

$$\langle e_j, e_k \rangle_n = \begin{cases} 1, & \text{if } j_l \equiv k_l \pmod{m}, \forall l \in \{1, \dots, d\}, \\ 0, & \text{otherwise.} \end{cases} \quad (8.6)$$

The scale of isotropic Sobolev spaces $H_s(\mathfrak{I}^d)$ is defined in terms of the Fourier coefficients $f_k = \int_{\mathfrak{I}^d} f(x) e_k(x) dx$ of functions $f \in L^2(\mathfrak{I}^d)$, as (for $|k|$ any norm on \mathbb{R}^d)

$$H_s(\mathfrak{I}^d) := \left\{ f \in L^2(\mathfrak{I}^d) : \|f\|_{H_s}^2 := \sum_{k \in \mathbb{Z}^d} (1 + |k|)^{2s} |f_k|^2 < \infty \right\}.$$

For smoothness levels $s \in \mathbb{N}$, this norm is equivalent to the canonical Sobolev norm $\sum_{|l|_1 \leq s} \|D^l f\|_{L^2(\mathfrak{T}^d)}$.

The spaces V_j obtained as the linear span of the basis elements, ordered suitably, satisfy Assumption 2.3. Due to (8.6), the space $\widetilde{W}_n = \text{Span}\{e_k : |k|_\infty \leq (m-1)/2\}$ satisfies (8.3) with $C_1 = C_2 = 1$.

As noted in [82], the following estimates hold for $f \in H_s(\mathfrak{T}^D)$,

$$\begin{aligned} \|f - Q_n f\|_{L^2} &\lesssim n^{-s/d} \|f\|_{H_s}, \\ \|Q_n f - \mathcal{I}_n f\| &\lesssim_d n^{-s/d} \|f\|_{H_s}. \quad \text{if } s > d/2. \end{aligned}$$

Here Q_n is the orthogonal projection on \widetilde{W}_n . Consequently (8.4) is fulfilled for $s > s_d = d/2$.

Example 8.5 (Wavelets). This example is adapted from Section 3.3 in [82]. Let $\mathcal{D} = \mathfrak{T}^d = (0, 1]^d$, for $d \in \mathbb{N}$, and consider the design points $\mathcal{D}_n = \{k2^{-j}\}_{k \in \{1, \dots, 2^j\}^d}$, where $n = 2^{jd}$ for some $j \in \mathbb{N}$. We consider a multiresolution analysis $\{V_j\}_{j \geq 0}$ on $L^2(\mathfrak{T}^d)$ obtained by periodization and tensor products. Let $\tilde{\phi}$ be a standard orthonormal scaling function of an S -regular multiresolution analysis for $L^2(\mathbb{R})$, with compact support in $[S-1, S]$. In particular, the polynomial exactness condition is satisfied: $\sum_{k \in \mathbb{Z}} k^q \tilde{\phi}(x-k) - x^q$ is a polynomial of maximal degree $q-1$ for $q \in [0, S-1]$. As shown in [82], the functions

$$e_{j,k}(x_1, \dots, x_d) = \sum_{m \in \mathbb{Z}^d} 2^{jd/2} \prod_{i=1}^d \tilde{\phi}(2^j x_i - k_i + 2^j m_i),$$

are well defined and form an orthonormal basis in $L^2(\mathfrak{T}^d)$. Furthermore, for $\widetilde{W}_n := V_j = \text{Span}\{e_{j,k} \mid k \in \{1, \dots, 2^j\}^d\}$ with $n = 2^{jd} \geq 2S-1$, conditions (8.3) is satisfied with constants C_1, C_2 that depend only on $\tilde{\phi}$. Moreover, for the functions $e_{j,k}$ belong to the Besov space $B_{2,2}^s(\mathbb{T})$, for $s < S$, and, for every f in this Besov space and $d/2 < s < S$,

$$\|f - \mathcal{I}_n f\|_{L^2} \lesssim n^{-s/d} \|f\|_{B_{2,2}^s}.$$

Thus (8.4) is satisfied, with the smoothness scale $(H_s)_{s \in \mathbb{R}}$ taken equal to the canonical Sobolev spaces (i.e. Besov spaces $B_{2,2}^s$) on \mathfrak{T}^d .

Other examples of suitable discretization spaces are provided by orthogonal polynomials, for instance the systems of Legendre, Chebyshev, or Jacobi polynomials, etc., for suitably chosen design points. First, $(H_s(\mathfrak{T}))_{s \in \mathbb{R}}$ being canonical Sobolev spaces on $\mathfrak{T} = (0, 1]$ satisfies Assumption 2.3 This is due to the standard Sturm-Liouville theory (see 5.2 in [12]): the polynomials form infinitely differentiable orthogonal bases in $L^2(\mathfrak{T})$. Second Assumption 8.1 is satisfied with Gaussian quadrature points as design points (Section 5.3 in [12]). These results can be extended to the multivariate domains by using tensor products. See Chapter 5 in [12] for more information.

8.2 General Contraction Rates

In this section we present a general theorem on posterior contraction. We form the posterior distribution $\Pi_n(\cdot | Y^n)$ as in (4.1), given a prior Π on the space $H = H_0$ and an observation $Y^n = (Y_1, \dots, Y_n)$, whose conditional distribution given f is determined by the model (8.2). We study this random distribution under the assumption that Y^n follows the model (8.2) for a given ‘true’ function $f = f_0$, which we assume to be an element of H_β in a given smoothness scale $(H_s)_{s \in \mathbb{R}}$, as in Definition 2.1.

The theorem is stated in terms of the Galerkin solution to the continuous inverse problem, which is defined as follows. (See e.g., [57] for a general introduction to the Galerkin method and Section 5.3 for a self-contained derivation of the necessary inequalities, exactly in our framework.) Let $W_j = AV_j \subset G$ be the image under the operator \mathcal{A} of a finite-dimensional approximation space V_j linked to the smoothness scale $(H_s)_{s \in \mathbb{R}}$ as in Assumption 2.3, and let $Q_j : G \rightarrow W_j$ be the orthogonal projection onto W_j . If $A : H \rightarrow G$ is injective, then A is a bijection between the finite-dimensional vector spaces V_j and W_j , and hence for every $f \in H$ there exists $f^{(j)} \in V_j$ such that $Af^{(j)} = Q_j Af$. The element $f^{(j)}$ is called the *Galerkin solution* to Af in V_j , and is an approximation to f that is more accurate, but also more complex, for larger j .

In our current setting we have no access to the continuous function Af , but must reconstruct f from the discrete approximation to $\mathcal{A}_n f$, for $\mathcal{A}_n = \mathcal{I}_n A$, and \mathcal{I}_n the interpolation operator defined in Section 8.1. Thus we shall use the Galerkin solution $f^{(j,n)} = A^{-1}Q_j \mathcal{A}_n f$ to the interpolation $\mathcal{A}_n f$ of the discrete signal. This *discrete Galerkin solution* is illustrated in the following diagram

$$\begin{array}{ccc}
 H \ni f & \xrightarrow{\mathcal{A}_n = \mathcal{I}_n A} & \mathcal{A}_n f \in \widetilde{W}_n \subset G \\
 & & \downarrow Q_j \\
 H \supset V_j \ni f^{(j,n)} & \xleftarrow{A^{-1}} & Q_j \mathcal{A}_n f \in W_j \subset G
 \end{array}$$

In this scheme the space \widetilde{W}_n , used to construct the continuous interpolation, may or may not be equal to $W_n = AV_n$. Setting it equal to W_n simplifies the scheme, but then the interpolation properties in Assumption 8.1 must be verified for AV_n .

Theorem 8.6. *For smoothness classes $(H_s)_{s \in \mathbb{R}}$ as in Definition 2.1, assume that the operator $\mathcal{A} : H_0 \rightarrow G$ satisfies $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$, for some $\gamma > 0$. Let $f^{(j,n)}$ denote the discrete Galerkin solution to $\mathcal{A}_n f = \mathcal{I}_n \mathcal{A}f$ relative to linear subspaces V_j associated to $(H_s)_{s \in \mathbb{R}}$ as in Assumption 2.3 and interpolation spaces \widetilde{W}_n satisfying (8.3)-(8.4) from Assumption 8.1. Let $f_0 \in H_\beta$ and $Af_0 \in G_{\beta+\gamma}$ for some $\beta \in (s_d - \gamma, S_d - \gamma)$, and for $\eta_n \geq \varepsilon_n \downarrow 0$ such that $n\varepsilon_n^2 \rightarrow \infty$, and $j_n \in \mathbb{N}$ such that $j_n \rightarrow \infty$, and some*

$c > 0$, assume

$$j_n \leq cn\varepsilon_n^2, \quad (8.7)$$

$$\eta_n \geq \frac{\varepsilon_n}{\delta(j_n, \gamma)}, \quad (8.8)$$

$$\eta_n \geq \delta(j_n, \beta) \vee \frac{\delta_d(n, \beta + \gamma)}{\delta(j_n, \gamma)}. \quad (8.9)$$

Consider prior probability distributions Π on H_0 satisfying

$$\Pi(f \in H : \|\mathcal{A}f - \mathcal{A}f_0\|_n < \varepsilon_n) \gtrsim e^{-n\varepsilon_n^2}, \quad (8.10)$$

$$\Pi(f \in H : \|f^{(j_n, n)} - f\| > \eta_n) \lesssim e^{-4n\varepsilon_n^2}. \quad (8.11)$$

Then the posterior distribution in the model (8.2) contracts at the rate η_n at f_0 , i.e. for a sufficiently large constant M we have $\Pi_n(f : \|f - f_0\| > M\eta_n \mid Y_1, \dots, Y_n) \rightarrow 0$, in probability if Y_1, \dots, Y_n follow (8.2) with $f = f_0$.

Proof. The Kullback-Leibler divergence and variation between the (multivariate-normal) distributions of (Y_1, \dots, Y_n) under two functions f and f_0 are given by $n\|\mathcal{A}f - \mathcal{A}f_0\|_n^2/2$ and twice this quantity, respectively. Therefore the neighbourhoods $B_{n,2}(f_0, \varepsilon_n)$ in (8.19) of [35] contain the balls $\{f \in H : \|\mathcal{A}f - \mathcal{A}f_0\|_n \leq \varepsilon_n\}$. By assumption (8.10) this has prior mass at least $\exp(-n\varepsilon_n^2)$.

Because the quotient of the left sides of (8.11) and (8.10) is $o(\exp(-2n\varepsilon_n^2))$, the posterior probability of the set $\{f : \|f^{(j_n, n)} - f\| > \eta_n\}$ tends to zero, by Theorem 8.20 in [35].

By a variation of Theorem 8.22 in [35] it is now sufficient to show the existence of tests τ_n such that, for some $M > 0$,

$$P_{f_0}^{(n)}\tau_n \rightarrow 0, \quad \sup_{\substack{\|f - f_0\| > M\eta_n, \\ \|f^{(j_n, n)} - f\| \leq \eta_n}} P_f^{(n)}(1 - \tau_n) \leq e^{-4n\varepsilon_n^2}.$$

Define the operator $\mathcal{R}_j : G \mapsto V_j$ by $\mathcal{R}_j = \mathcal{A}^{-1}Q_j$, where $Q_j : G \rightarrow W_j$ is the orthogonal projection onto $W_j = \mathcal{A}V_j$ and \mathcal{A}^{-1} is the inverse of \mathcal{A} (restricted to W_j). Then we shall employ the tests

$$\tau_n = 1\{\|\mathcal{R}_{j_n}Y^{(n)} - f_0\| \geq M_0\eta_n\}, \quad (8.12)$$

where M_0 is a given constant, to be determined.

By definition $f^{(j, n)} = \mathcal{R}_j\mathcal{A}_nf$ is equal to the discrete Galerkin solution to \mathcal{A}_nf . For $Y^{(n)}$ defined in (8.5) we have

$$\mathcal{R}_jY^{(n)} = \mathcal{R}_j\mathcal{A}_nf + \frac{1}{\sqrt{n}}\mathcal{R}_j\xi^{(n)} = f^{(j, n)} + \frac{1}{\sqrt{n}}\mathcal{R}_j\xi^{(n)}. \quad (8.13)$$

The variable $\mathcal{R}_j\xi^{(n)} = \mathcal{R}_jQ_j\xi^{(n)}$ is a centered Gaussian random element in V_j with strong and weak second moments

$$\mathbb{E}\|\mathcal{R}_j\xi^{(n)}\|^2 \leq \|\mathcal{R}_j\|^2\mathbb{E}\|Q_j\xi^{(n)}\|^2 \lesssim \|\mathcal{R}_j\|^2j \lesssim \frac{j}{\delta(j, \gamma)^2},$$

$$\sup_{\|f\| \leq 1} \mathbb{E}\langle \mathcal{R}_j\xi^{(n)}, f \rangle^2 = \sup_{\|f\| \leq 1} \mathbb{E}\langle \xi^{(n)}, \mathcal{R}_j^*f \rangle_G^2 \lesssim \sup_{\|f\| \leq 1} \|\mathcal{R}_j^*f\|_G^2 \leq \|\mathcal{R}_j^*\|^2 \lesssim \frac{1}{\delta(j, \gamma)^2}.$$

In both cases the inequality on $\|\mathcal{R}_j\| = \|\mathcal{R}_j^*\|$ at the far right side follows from (5.7), and we also use that, by Lemma 8.2, the covariance operator of $\xi^{(n)}$ is bounded above by a multiple of the projection onto \widetilde{W}_n , and hence the identity, so that the covariance operator of $Q_j \xi^{(n)}$ is bounded above by a multiple of Q_j .

The first inequality shows that the first moment $\mathbb{E}\|\mathcal{R}_j \xi^{(n)}\|$ of the variable $\|\mathcal{R}_j \xi^{(n)}\|$ is bounded above by $\sqrt{j}/\delta(j, \gamma)$. By Borell's inequality (e.g. Lemma 3.1 in [67] and subsequent discussion), applied to the Gaussian random variable $\mathcal{R}_j \xi^{(n)}$ in H_0 , we see that, there exist positive constants a and b such that, for every $t > 0$,

$$\Pr\left(\|\mathcal{R}_j \xi^{(n)}\| > t + a \frac{\sqrt{j}}{\delta(j, \gamma)}\right) \leq e^{-bt^2 \delta(j, \gamma)^2}.$$

For $t = 2\sqrt{n}\eta_n/\sqrt{b}$ and η_n , ε_n and j_n satisfying (8.7), (8.8) and (8.9) this yields, for some $a_1 > 0$,

$$\Pr\left(\|\mathcal{R}_{j_n} \xi^{(n)}\| > a_1 \sqrt{n}\eta_n\right) \leq e^{-4n\varepsilon_n^2}. \quad (8.14)$$

We apply this to bound the two error probabilities of the tests τ_n .

Under f_0 the decomposition (8.13) is valid with $f = f_0$, and hence $\mathcal{R}_j Y^{(n)} - f_0 = n^{-1/2} \mathcal{R}_j \xi^{(n)} + f_0^{(j, n)} - f_0$. By the triangle inequality it follows that $\tau_n = 1$ implies that $n^{-1/2} \|\mathcal{R}_{j_n} \xi^{(n)}\| \geq M_0 \eta_n - \|f_0^{(j, n)} - f_0\|$. By the triangle inequality followed by (5.7) and (8.4), and (5.9),

$$\begin{aligned} \|f_0^{(j, n)} - f_0\| &\leq \|\mathcal{R}_j\| \|\mathcal{I}_n \mathcal{A} f_0 - \mathcal{A} f_0\| + \|\mathcal{R}_j \mathcal{A} f_0 - f_0\| \\ &\lesssim \frac{\delta_d(n, \beta + \gamma)}{\delta(j_n, \gamma)} \|\mathcal{A} f_0\|_{\beta + \gamma} + \delta(j_n, \beta) \|f_0\|_{\beta} \leq M_1 \eta_n, \end{aligned}$$

by assumption (8.9). Hence the probability of an error of the first kind satisfies

$$P_{f_0}^{(n)} \tau_n \leq \Pr\left(\frac{1}{\sqrt{n}} \|\mathcal{R}_{j_n} \xi^{(n)}\| \geq (M_0 - M_1) \eta_n\right),$$

For $M_0 - M_1 > a_1$, the right side is bounded by $e^{-4n\varepsilon_n^2}$, by (8.14).

Under f the decomposition (8.13) gives that $\mathcal{R}_j Y^{(n)} - f_0 = n^{-1/2} \mathcal{R}_j \xi^{(n)} + f^{(j, n)} - f_0$. By the triangle inequality $\tau_n = 0$ implies that $n^{-1/2} \|\mathcal{R}_{j_n} \xi^{(n)}\| \geq \|f^{(j_n, n)} - f_0\| - M_0 \eta_n$. For f such that $\|f - f_0\| > M \eta_n$ and $\|f - f^{(j_n, n)}\| \leq \eta_n$, we have $\|f^{(j_n, n)} - f_0\| \geq (M - 1) \eta_n$. Hence the probability of an error of the second kind satisfies

$$P_f^{(n)} (1 - \tau_n) \leq \Pr\left(\frac{1}{\sqrt{n}} \|\mathcal{R}_{j_n} \xi^{(n)}\| \geq (M - 1 - M_0) \eta_n\right),$$

For $M - 1 - M_0 > a_1$, this is bounded by $e^{-4n\varepsilon_n^2}$, by (8.14).

We can first choose M_0 large enough so that $M_0 - M_1 > a_1$, and next M large enough so that $M - 1 - M_0 > a_1$, to finish the proof. \square

The theorem has a similar form as Theorem 6.1 obtained in Chapter 6 in the case of observation of a continuous signal (in white noise). Some interpretations of the theorem from the previous chapter are also applicable to the current one.

For completeness, we repeat them here. Inequality (8.10) is the usual *prior mass condition* for the ‘direct problem’ of estimating $\mathcal{A}f$ at the design points (see [33]). It determines the rate of contraction ε_n of the posterior distribution of $\mathcal{A}f$ to $\mathcal{A}f_0$ relative to the discrete seminorm $\|\cdot\|_n$. The rate of contraction η_n of the posterior distribution of f is slower due to the necessity of (implicitly) inverting the operator \mathcal{A} . The theorem shows that the rate η_n depends on the combination of the prior, through (8.11), and the inverse problem, through the various approximation rates. The factor $\delta_d(n, \beta + \gamma) / \delta(j_n, \gamma)$ in (8.9), which arises from having discrete observations only, will typically be negligible relative to $\delta(j_n, \beta)$. The Galerkin projection $f^{(j,n)}$ in (8.11) now incorporates the errors of both inversion and discretisation.

Same adaptation by mixture priors (c.f. Theorem 6.11) can also be achieved. The following theorem refines Theorem 8.6 by considering a mixture prior of the form

$$\Pi = \int \Pi_\tau dQ(\tau), \quad (8.15)$$

where Π_τ is a prior on H , for every given ‘hyperparameter’ τ running through some measurable space, and Q is a prior on this hyperparameter. The idea is to *adapt* the prior to multiple smoothness levels through the hyperparameter τ .

Theorem 8.7. *Consider the setup and assumptions of Theorem 8.6 with a prior of the form (8.15). Assume that (8.7), (8.8), (8.9) and (8.10) hold, but replace (8.11) by the pair of conditions, for numbers $\eta_{n,\tau}$ and $C > 0$ and every τ ,*

$$\Pi_\tau(f : \|f - f_0\| < 2\eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2}, \quad \forall \tau \text{ with } \eta_{n,\tau} \geq C\eta_n, \quad (8.16)$$

$$\Pi_\tau(f : \|f^{(j_n,n)} - f\| > \eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2}. \quad (8.17)$$

Then the posterior distribution in the model (8.2) contracts at the rate η_n at f_0 , i.e. for a sufficiently large constant M we have $\Pi_n(f : \|f - f_0\| > M\eta_n \mid Y_1, \dots, Y_n) \rightarrow 0$, in probability if Y_1, \dots, Y_n follow (8.2) with $f = f_0$.

Proof. We take the parameter of the model as the pair (f, τ) , which receives the joint prior given by $f \mid \tau \sim \Pi_\tau$ and $\tau \sim Q$. With abuse of notation, we denote this prior also by Π . The likelihood still depends on f only, but the joint prior gives rise to a posterior distribution on the pair (f, τ) , which we also denote by $\Pi_n(\cdot \mid Y^n)$, by a similar abuse of notation.

By (8.15) and eqs. (8.16) and (8.17),

$$\Pi((f, \tau) : \eta_{n,\tau} \geq C\eta_n, \|f - f_0\| < 2\eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2},$$

$$\Pi((f, \tau) : \|f^{(j_n,n)} - f\| > \eta_{n,\tau}) \leq e^{-4n\varepsilon_n^2}.$$

In view of (8.10) and Theorem 8.20 in [35], the posterior probabilities of the two sets in the left sides tend to zero. As in the proof of Theorem 8.6, we can apply a variation of Theorem 8.22 in [35] to see that it is now sufficient to show the existence of tests τ_n such that, for some $M \geq 2C$,

$$P_{f_0}^{(n)} \tau_n \rightarrow 0, \quad \sup_{\substack{(f,\tau) : \|f - f_0\| > M\eta_n \vee 2\eta_{n,\tau}, \\ \|f^{(j_n,n)} - f\| \leq \eta_{n,\tau}}} P_f^{(n)} (1 - \tau_n) \leq e^{-4n\varepsilon_n^2}.$$

(Note that $M\eta_n \vee 2\eta_{n,\tau} = M\eta_n$ if $\eta_{n,\tau} < C\eta_n$ and $M \geq 2C$.) We use the tests defined in (8.12), as in the proof of Theorem 8.6. The latter proof shows that the tests are consistent. The bound on the power can be adapted same as in the proof of Theorem 6.11, and hence the detail is omitted here. \square

In a typical application of the preceding theorem the priors Π_τ for τ such that $\eta_{n,\tau} \geq C\eta_n$ will be the priors on rough functions, with ‘intrinsic’ contraction rate $\eta_{n,\tau}$ slower than η_n . These ‘bad’ priors do not destroy the overall contraction rate, because they put little mass near the true function f_0 , by condition (8.16). It is necessary to address these priors explicitly in the conditions, because they will typically fail the approximation condition (8.11), which must be relaxed to (8.17).

In Theorem 8.6 and Theorem 8.7, it is sufficient for the operator \mathcal{A} to satisfy the following two properties: lifting property $\|\mathcal{A}f\| \simeq \|f\|_{-\gamma}$ with some $\gamma > 0$ and $\mathcal{A}f_0 \in G_{\beta+\gamma}$ if $f_0 \in H_\beta$. However, when analysing particular priors, the aforementioned conditions on \mathcal{A} may not be strong enough for verifying the conditions (8.10) and (8.11), while they can be verified with the following condition, which is stronger but often satisfied in practice, e.g. Example 5.4. More examples can be found in [57].

Assumption 8.8 (Smoothing property of \mathcal{A}). For some $\gamma > 0$ the operator $\mathcal{A} : H_{s-\gamma} \rightarrow G_s$ is injective and bounded for every $s \geq 0$, and, for every $f \in H_0 \cap H_{s-\gamma}$,

$$\|\mathcal{A}f\|_s \simeq \|f\|_{s-\gamma}. \quad (8.18)$$

8.3 Random Series Priors

In this section we study the performance of the random series prior from Section 6.3 to the inverse regression model (8.2). We briefly recall the set-up of the prior and the further details are referred to Section 6.3.

Suppose that $\{\phi_i\}_{i \in \mathbb{N}}$ is an orthonormal basis of $H = H_0$ that gives optimal approximation relative to the scale of smoothness classes $(H_s)_{s \in \mathbb{R}}$ in the sense that the linear spaces $V_j = \text{Span}\{\phi_i\}_{i < j}$ satisfy Assumption 2.3. Consider a prior defined as the law of the random series

$$f = \sum_{i=1}^M f_i \phi_i, \quad (8.19)$$

where M is a random variable in \mathbb{N} independent from the independent random variables f_1, f_2, \dots in \mathbb{R} .

Condition 8.9 (Random series prior). (i) The probability density function p_M of M satisfies, for some positive constants b_1, b_2 ,

$$e^{-b_1 k} \lesssim p_M(k) \lesssim e^{-b_2 k}, \quad \forall k \in \mathbb{N}.$$

(ii) The variable f_i has density $p(\cdot/\kappa_i)/\kappa_i$, for a given probability density p on \mathbb{R} and a positive constant κ_i such that, for some $C > 0$ and $0 < v < w < \infty$,

$\beta_0 > 0$ and $\alpha > 0$,

$$e^{-C|x|^w} \lesssim p(x) \lesssim e^{-C|x|^v}, \quad (8.20)$$

$$i^{-\beta_0/d} \lesssim \kappa_i \lesssim i^\alpha. \quad (8.21)$$

The same contraction rate is obtained for the random series prior in the inverse regression model. The discussion on the result is referred to the discussion of Theorem 6.5.

Theorem 8.10 (Random Series Prior). *Let $(\phi_i)_{i \in \mathbb{N}}$ be an orthonormal basis of H_0 such that the spaces $V_j = \text{Span}\{\phi_i\}_{i < j}$ satisfy Assumption 2.3 with $\delta(j, t) = j^{-t/d}$ relative to smoothness classes $(H_s)_{s \in \mathbb{R}}$ as in Definition 2.1 and sufficiently large t to be specified. Suppose that Assumption 8.1 holds with $\delta_d(n, s) = n^{-s/d}$, the operator \mathcal{A} satisfies Assumption 8.8, and let $f_0 \in H_\beta$ for some $\beta \in (0, S]$ and $\beta + \gamma > s_d$. Then, for the random series prior defined in (8.19) and satisfying Condition 8.9 with $\beta_0 \leq \beta$, and sufficiently large $M > 0$, for $\tau = (\beta + \gamma)(1 + 2\gamma/d)/(2\beta + 2\gamma + d)$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-\beta/(2\beta+2\gamma+d)} (\log n)^\tau \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

8.4 Gaussian Priors

In this section we study the posterior contractions of Gaussian priors in the inverse regression model (8.2). Recall the definition of a Gaussian prior from Section 6.4. Centred Gaussian distributions on a separable Hilbert space correspond bijectively to covariance operators. By definition a random variable F with values in H_0 is Gaussian if $\langle F, g \rangle_0$ is normally distributed, for every $g \in H_0$, and it has zero mean if these variables have zero means. The variances of these variables can then be written as

$$\mathbb{E} \langle F, g \rangle_0^2 = \langle Cg, g \rangle_0,$$

for a linear operator $C : H_0 \rightarrow H_0$, called the *covariance operator*. A covariance operator C is necessarily self-adjoint, nonnegative, and of *trace class*, i.e., $\sum_{i \in \mathbb{N}} \langle C\phi_i, \phi_i \rangle < \infty$, for some (and then every) orthonormal basis $(\phi_i)_{i \in \mathbb{N}}$ of H_0 ; and every operator with these properties generates a Gaussian distribution.

In the setting of a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by the operator L it is natural to choose a Gaussian prior with covariance operator of the form $L^{-2\alpha}$, for some $\alpha > 0$. If L^{-1} has eigenvalues λ_j , then this operator is of trace class if $\sum_{j \in \mathbb{N}} \lambda_j^{-2\alpha} < \infty$. Thus α must be chosen big enough for the Gaussian prior to exist as a ‘proper’ prior on H_0 . For instance, if $\lambda_j \simeq j^{-1/d}$, then every choice $\alpha > d/2$ yields a proper prior.

This leads to the following theorem on posterior contraction rates for Gaussian priors, the proof of which is given in Section 8.6.

Theorem 8.11 (Gaussian Prior). *Consider a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by an operator L as in the preceding such that $L^{-1} : H_0 \rightarrow H_0$ is compact with eigenvalues λ_j satisfying $\lambda_j \simeq j^{-1/d}$. Suppose that Assumption 8.1 holds with*

$\delta_d(n, s) = n^{-s/d}$, the operator $\mathcal{A} : H_0 \rightarrow G$ satisfies Assumption 8.8, $f_0 \in H_\beta$, for some $\beta > 0$ such that $\beta + \gamma > s_d$, and let the prior be zero-mean Gaussian with covariance operator $L^{-2\alpha}$, for some $\alpha > (d - \gamma) \vee d/2$. Then the posterior distribution satisfies, for sufficiently large $M > 0$,

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-((\alpha-d/2) \wedge \beta)/(2\alpha+2\gamma)} \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

The results are comparable to Theorem 6.7, and hence we refer to Section 6.4 for the discussion.

8.5 Gaussian Mixtures

The posterior contraction rate resulting from a zero-mean Gaussian prior with covariance operator $L^{-2\alpha}$, as considered in Section 8.4, is equal to the minimax rate $n^{-\beta/(2\beta+2\gamma+d)}$ (see [19]) only when $\alpha - d/2 = \beta$, i.e., when the prior smoothness $\alpha - d/2$ matches the true smoothness β . As shown in Section 6.5, by mixing over Gaussian priors of varying smoothness the minimax rate can often be obtained simultaneously for a range of values β . In this section we consider the same mixture of the mean-zero Gaussian priors with covariance operators $\tau^2 L^{-2\alpha}$ over the ‘hyperparameter’ τ , from Section 6.5. Thus the prior Π is the distribution of τF , where F is a zero-mean Gaussian variable in H_0 with covariance operator $L^{-2\alpha}$, as in Section 8.4, and τ is an independent scale parameter satisfying Condition 6.10. We repeat the condition below for the reader’s convenience.

Condition 8.12. The distribution Q of τ has support $[0, \infty)$ and satisfies

$$\begin{cases} -\log Q((t, 2t)) \lesssim t^{-2}, & \text{as } t \downarrow 0, \\ -\log Q((t, 2t)) \lesssim t^{d/(\alpha-d/2)}, & \text{as } t \rightarrow \infty. \end{cases}$$

Theorem 8.13 (Gaussian mixture prior). *Consider a Hilbert scale $(H_s)_{s \in \mathbb{R}}$ generated by an operator L as in the preceding such that $L^{-1} : H_0 \rightarrow H_0$ is compact with eigenvalues λ_j satisfying $\lambda_j \simeq j^{-1/d}$. Suppose the operator $\mathcal{A} : H_s \rightarrow G_{s+\gamma}$ satisfies Assumption 8.8, assume that $f_0 \in H_\beta$, for some $(d/2 - \gamma) \vee 0 < \beta \leq \alpha$, and let the prior be a mixture of the zero-mean Gaussian distributions with covariance operators $\tau^2 L^{-2\alpha}$ over the parameters τ equipped with a prior satisfying Condition 8.12, for some $\alpha > (d - \gamma) \vee d/2$. Then the posterior distribution satisfies, for sufficiently large $M > 0$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-\beta/(2\beta+2\gamma+d)} \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

The proof is given in Section 8.6.

8.6 Proofs

8.6.1 Proof of Theorem 8.10

The theorem is a corollary to Theorem 8.6. We shall verify the conditions with

$$\varepsilon_n \simeq (\log n/n)^{(\beta+\gamma)/(2\beta+2\gamma+d)}, \quad j_n \simeq n^{d/(2\beta+2\gamma+d)} (\log n)^{(2\beta+2\gamma)/(2\beta+2\gamma+d)}.$$

Let P_j be the orthogonal projection of H on the linear span of the first $j-1$ basis elements ϕ_j , and define an additional sequence of integers by

$$i_n \simeq (n/\log n)^{d/(2\beta+2\gamma+d)}.$$

By the orthogonality of the basis (ϕ_i) , the function ϕ_j is orthogonal to the space V_j spanned by $(\phi_i)_{i < j}$. Hence $P_j \phi_j = 0$, so that $\|\phi_j\|_{-s} \leq \delta(j, s) \|\phi_j\|_0 \lesssim j^{-s/d}$, for every j and $s \geq 0$, by (2.4). The same estimate is also true for $0 < -s < S$, directly by assumption (2.3). Therefore, by the triangle inequality, we have

$$\|f\|_s \lesssim \sum_j |f_j| j^{s/d}, \quad \text{if } f = \sum_j f_j \phi_j.$$

Furthermore, since $f_0 \in H_\beta$ by assumption, the norm duality (2.1) gives that $|f_{0,i}| = |\langle f_0, \phi_i \rangle_0| \leq \|f_0\|_\beta \|\phi_i\|_{-\beta} \lesssim i^{-\beta/d}$.

First we verify the prior condition (8.10) of the direct problem. By Lemma 8.3, $\|\mathcal{A}f - \mathcal{A}f_0\|_n \simeq \|\mathcal{I}_n(\mathcal{A}f - \mathcal{A}f_0)\|$. By several applications of the triangle inequality, since $\beta + \gamma \in (s_d, S_d)$ and $\|\mathcal{A}f\|_{\beta+\gamma} \simeq \|f\|_\beta$,

$$\begin{aligned} \|\mathcal{I}_n(\mathcal{A}f - \mathcal{A}f_0)\| &\leq \|\mathcal{I}_n \mathcal{A}f - \mathcal{A}f\| + \|\mathcal{I}_n \mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\| \\ &\lesssim \delta_d(n, \beta + \gamma) (\|f\|_\beta + \|f_0\|_\beta) + \|f - P_{i_n} f_0\|_{-\gamma} + \|f_0 - P_{i_n} f_0\|_{-\gamma} \\ &\lesssim \delta_d(n, \beta + \gamma) \|f - P_{i_n} f_0\|_\beta + \|f - P_{i_n} f_0\|_{-\gamma} \\ &\quad + \delta_d(n, \beta + \gamma) (\|P_{i_n} f_0\|_\beta + \|f_0\|_\beta) + \delta(i_n, \gamma) \delta(i_n, \beta) \|f_0\|_\beta, \end{aligned}$$

by (2.4). The last term is of the order $\delta(i_n, \gamma) \delta(i_n, \beta) = i_n^{-(\gamma+\beta)/d} \simeq \varepsilon_n$, while the second last term is bounded above by $\delta_d(n, \beta + \gamma) (\|f_0\|/\delta(i_n, \beta) + \|f_0\|_\beta) \ll \varepsilon_n$, if $\beta + \gamma > d/2$. For $f = \sum_{i=1}^{i_n-1} f_i \phi_i \in V_{i_n}$ the sum of the first two terms is bounded above by

$$\delta_d(n, \beta + \gamma) \sum_{i=1}^{i_n-1} |f_i - f_{0,i}| i^{\beta/d} + \sum_{i=1}^{i_n-1} |f_i - f_{0,i}| i^{-\gamma/d} \lesssim \sum_{i=1}^{i_n-1} |f_i - f_{0,i}| i^{-\gamma/d}.$$

Same as shown in Section 6.6.1, the right side of this equation is bounded above by ε_n with prior probability at least $e^{-n\varepsilon_n^2/2}$. Since also $\Pi(M = i_n - 1) \geq e^{-b_1 i_n} \geq e^{-n\varepsilon_n^2/2}$, it follows that (8.10) is satisfied.

Next we verify (8.11). Since $\Pi(M \geq j_n) \leq e^{-b_2 j_n} \leq e^{-4n\varepsilon_n^2}$, by Condition 8.9, we may intersect the event in (8.11) with the event $M < j_n$. For $M < j$ the random series $f = \sum_{i=1}^M f_i \phi_i$ is contained in V_j and the Galerkin approximation

$\mathcal{R}_j \mathcal{A}f$ of f is exact. Since $f^{(j,n)} = \mathcal{R}_j \mathcal{I}_n \mathcal{A}f$, the triangle inequality followed by (8.4) give, for $s + \gamma \in (s_d, S_d)$,

$$\begin{aligned} \|f^{(j,n)} - f\| &\leq \|\mathcal{R}_j \mathcal{I}_n \mathcal{A}f - \mathcal{R}_j \mathcal{A}f\| + \|\mathcal{R}_j \mathcal{A}f - f\| \\ &\leq \frac{\delta_d(n, s + \gamma)}{\delta(j, \gamma)} \|f\|_s + \|\mathcal{R}_j \mathcal{A}f - f\| \\ &\leq \frac{\delta_d(n, s + \gamma)}{\delta(j, \gamma) \delta(j, s)} \|f\| + 0, \end{aligned}$$

if $f \in V_j$, by (2.3), for s such that $f \in H_s$. We conclude that it suffices to prove that

$$\Pi\left(\sum_{i=1}^{j_n-1} f_i^2 > \eta_n^2 (j_n/n)^{2\gamma+2s}\right) \leq e^{-4n\varepsilon_n^2}.$$

With the given choices of η_n and j_n , for some $a \in \mathbb{R}$,

$$\eta_n^2 (j_n/n)^{2\gamma+2s} = n^{(4(s+\gamma)(\beta+\gamma)-2d\beta)/(2\beta+2\gamma+d)} (\log n)^a.$$

For sufficiently large s this is an arbitrary high power of n . We have that

$$\mathbb{E} \sum_{i=1}^{j_n-1} f_i^2 \simeq \sum_{i=1}^{j_n-1} \kappa_i^2 \lesssim j_n^{2\alpha+1}.$$

$$\mathbb{E} \left| \sum_{i=1}^{j_n-1} (f_i^2 - \mathbb{E} f_i^2) \right| \lesssim \sqrt{\sum_{i=1}^{j_n-1} \kappa_i^4} \lesssim j_n^{2\alpha+1/2}.$$

By the tail bound on the density p we further have that the $\psi_{v/2}$ Orlicz norm of f_i^2 is bounded above by κ_i^2 . Therefore,

$$\left\| \sum_{i=1}^{j_n-1} (f_i^2 - \mathbb{E} f_i^2) \right\|_{\psi_{v/2}} \lesssim \begin{cases} j_n^{2\alpha+1/2} + (\log j_n)^{2/v} \max_{i < j_n} \kappa_i^2, & \text{if } v \leq 2, \\ j_n^{2\alpha+1/2} + (\sum_{i < j_n} \kappa_i^{2q})^{1/q}, & \text{if } 2 < v \leq 4, \end{cases}$$

where q is conjugate to $v/2$. In both cases the first term $j_n^{2\alpha+1/2}$ dominates. So provided

$$n^{(4(s+\gamma)(\beta+\gamma)-2d\beta)/(2\beta+2\gamma+d)} (\log n)^a \gtrsim j_n^{2\alpha+1}$$

we can first center $\sum_{i=1}^{j_n-1} f_i^2$ at mean zero, and next bound the tail of the centered variable with the help of the Orlicz norm. This will give a bound of the type

$$1/\psi_{v/2} \left(\frac{n^{(4(s+\gamma)(\beta+\gamma)-2d\beta)/(2\beta+2\gamma+d)} (\log n)^a}{j_n^{2\alpha+1/2}} \right).$$

Under the preceding display the quotient is a positive power n^t of n times a logarithmic factor and we obtain a bound of the form

$$e^{-(n^t)^{v/2}}.$$

Here t can be arbitrarily large by choosing s large. So condition (8.11) holds provided s can be chosen sufficiently large in the interpolation inequality.

8.6.2 Proof of Theorem 8.11

The theorem is a corollary to Theorem 8.6. The main tasks are to determine ε_n satisfying the prior mass condition (8.10) of the direct problem, and next to identify η_n from the prior mass condition (8.11) and the other conditions.

Similar to the proof of Theorem 8.10, the prior mass condition (8.10) is decomposed into different components, which will be studied separately. Recall that by Lemma 8.3, $\|\mathcal{A}f - \mathcal{A}f_0\|_n \simeq \|\mathcal{I}_n(\mathcal{A}f - \mathcal{A}f_0)\|$. By triangle inequality, we have

$$\|\mathcal{I}_n(\mathcal{A}f - \mathcal{A}f_0)\| \leq \|\mathcal{I}_n\mathcal{A}f - \mathcal{A}f\| + \|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\|.$$

In the preceding display, the last term is same as the prior mass condition for white noise model and has been studied in Section 6.6.2. We recall the result in the following lemma.

Lemma 8.14 (Lemma 6.12). *Under the assumptions of Theorem 8.11, for $f_0 \in H_\beta$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon) \lesssim \begin{cases} \varepsilon^{-d/(\alpha+\gamma-d/2)}, & \text{if } d/2 < \alpha \leq \beta + d/2, \\ \varepsilon^{-(2\alpha-2\beta)/(\beta+\gamma)}, & \text{if } \alpha > \beta + d/2. \end{cases}$$

Hence, It follows that $\Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\|_0 < \varepsilon_n/2) \geq e^{-n\varepsilon_n^2}$ with

$$\varepsilon_n \leq \begin{cases} 2^{d/(2\alpha+2\gamma)} n^{-(\alpha+\gamma-d/2)/(2\alpha+2\gamma)}, & \text{if } d/2 < \alpha \leq \beta + d/2, \\ 2^{(2\alpha-2\beta)/(2\alpha+2\gamma)} n^{-(\beta+\gamma)/(2\alpha+2\gamma)}, & \text{if } \alpha > \beta + d/2. \end{cases}$$

The conditions of ε_n given above can be further simplified into

$$\varepsilon_n \geq C 2^{d/(2\gamma+d)} n^{-(\beta \wedge (\alpha-d/2) + \gamma)/(2\alpha+2\gamma)}, \quad (8.22)$$

where $2^{d/(2\gamma+d)} < 2$ is independent of α and β and C will be determined below.

Since $\beta + \gamma > d/2 = s_d$, we have $\|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| \leq \delta_d(n, \beta + \gamma) \|f_0\|_\beta \simeq n^{-(\beta+\gamma)/d}$. Therefore, by selecting a sufficiently large constant C in the preceding display such that the right hand side of (8.22) is upper bounded by ε_n , we obtain that $\Pi(f : \|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\|_0 < \varepsilon_n) \geq e^{-n\varepsilon_n^2}$.

Using a basic probability property, $P(A \cap B) = P(A) + P(B) - P(A \cup B) \geq P(B) - P(A^c)$,

$$\begin{aligned} & \Pi(f : \|\mathcal{I}_n\mathcal{A}f - \mathcal{A}f\| + \|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\|_0 < \varepsilon_n) \\ & \geq \Pi(f : \|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\|_0 < \varepsilon_n/2, \quad \|\mathcal{I}_n\mathcal{A}f - \mathcal{A}f\| < \varepsilon_n/2) \\ & \geq \Pi(f : \|\mathcal{I}_n\mathcal{A}f_0 - \mathcal{A}f_0\| + \|\mathcal{A}f - \mathcal{A}f_0\|_0 < \varepsilon_n/2) - \Pi(f : \|\mathcal{I}_n\mathcal{A}f - \mathcal{A}f\| \geq \varepsilon_n/2) \\ & \geq e^{-n\varepsilon_n^2} - \Pi(f : \|\mathcal{I}_n\mathcal{A}f - \mathcal{A}f\| \geq \varepsilon_n/2). \end{aligned}$$

Since the prior of f is a centred Gaussian distribution with covariance operator $L^{-2\alpha}$, f is in H_s almost surely for $s < \alpha - d/2$. Consequently, $\mathcal{A}f \in G_{s+\gamma}$ almost surely. Since α is chosen in the range $\alpha > (d - \gamma) \vee d/2$, which implies $\alpha + \gamma > d$,

and consequently there exists an s satisfying $d/2 = s_d < s + \gamma < S_d$. Hence, (8.4) holds and we have

$$\|\mathcal{I}_n \mathcal{A}f - \mathcal{A}f\| \leq \delta_d(n, s + \gamma) \|f\|_s = \delta_d(n, s + \gamma) \|L^s f\|_0.$$

It leads to

$$\Pi(f : \|\mathcal{I}_n \mathcal{A}f - \mathcal{A}f\| \geq \varepsilon_n/2) \leq \Pi\left(f : \|L^s f\| \geq \frac{\varepsilon_n}{2\delta_d(n, s + \gamma)}\right).$$

Therefore, it suffices to show that, for the transformed centred Gaussian random variable $L^s F$ with covariance operator $L^{-2(\alpha-s)}$ (where F is centred Gaussian with covariance $L^{-2\alpha}$),

$$\Pi(\|F\| > r_n) < e^{-n\varepsilon_n^2},$$

where $r_n = \frac{\varepsilon_n}{2\delta_d(n, s + \gamma)}$ and ε_n is as given in (8.22).

For F , Since $\mathbb{E}\|F\|^2 = \sum_{i \in \mathbb{N}} i^{-\frac{\alpha-s}{d/2}} < \infty$, the first moment of $\mathbb{E}\|F\|$ is also bounded by Jensen's inequality. In particular, the upper bound of $\mathbb{E}\|F\|$ is independent of n . The weak second moment is given by

$$\sigma = \sup_{\|h\|_0 \leq 1} \mathbb{E}\langle F, h \rangle_0^2 = \sup_{\|h\|_0 \leq 1} \|h\|_{-(\alpha-s)}^2.$$

By the norm duality (2.1), the right side is equal to

$$\sup_{\|h\|_0 \leq 1} \sup_{\|f\|_{\alpha-s} \leq 1} \langle f, h \rangle_0^2 \leq \sup_{\|f\|_{\alpha-s} \leq 1} \|f\|_0^2 \leq 1.$$

Then by Borell's inequality, when $r_n > \mathbb{E}\|F\|$,

$$\begin{aligned} \Pi(\|F\| > r_n) &\leq \Pi(\|F\| - \mathbb{E}\|F\| > r_n - \mathbb{E}\|F\|) \\ &\leq e^{-\frac{(r_n - \mathbb{E}\|F\|)^2}{2\sigma^2}} \leq e^{-r_n(r_n - \mathbb{E}\|F\|)} \ll e^{-n\varepsilon_n^2}, \end{aligned}$$

where we use the fact that $n\varepsilon_n^2 \ll r_n \rightarrow \infty$ since $\delta_d(n, s) \simeq n^{-s/d}$ and $s + \gamma > d/2$.

Combining above results, we have shown that the (8.10)

$$\Pi(f \in H : \|\mathcal{A}f - \mathcal{A}f_0\|_n < \varepsilon_n) \geq \frac{1}{2} e^{-n\varepsilon_n^2}$$

is satisfied with ε_n given in (8.22).

The next step of the proof is to bound the prior probability in (8.11). The following lemma is a modification of Lemma 8.2 in [43].

Lemma 8.15. *Under the assumptions of Theorem 8.11, there exist $a, b > 0$, such that for every $j \in \mathbb{N}$ and $t > 0$,*

$$\Pi(f : \|f^{(j)} - f\|_0 > t + aj^{1/2-\alpha/d}) \leq e^{-bt^2 j^{2\alpha/d}}.$$

Proof. We have $f^{(j)} - f = (\mathcal{R}_j \mathcal{A}_n - I)f$, for $\mathcal{R}_j = \mathcal{A}^{-1}Q_j$ and $\mathcal{A}_n = \mathcal{I}_n \mathcal{A}$. Therefore, the probability on the left concerns the random variable $(\mathcal{R}_j \mathcal{A}_n - I)F$, if F is a variable distributed according to the prior Π . Since F is zero-mean normal with covariance operator $L^{-2\alpha}$, this variable is zero-mean Gaussian with covariance operator $(\mathcal{R}_j \mathcal{A}_n - I)L^{-2\alpha}(\mathcal{R}_j \mathcal{A}_n - I)^*$. We shall apply Borell's inequality to obtain the exponential bound, after computing the weak and strong second moments of the variable $(\mathcal{R}_j \mathcal{A}_n - I)F$.

Because $\langle (\mathcal{R}_j \mathcal{A}_n - I)F, g \rangle_0 = \langle F, (\mathcal{R}_j \mathcal{A}_n - I)^* g \rangle_0$ is zero-mean Gaussian with variance $\|L^{-\alpha}(\mathcal{R}_j \mathcal{A}_n - I)^* g\|_0^2 = \|(\mathcal{R}_j \mathcal{A}_n - I)^* g\|_{-\alpha}^2$, the weak second moment of $(\mathcal{R}_j \mathcal{A}_n - I)F$ is given by

$$\sup_{\|g\|_0 \leq 1} \mathbb{E} \langle (\mathcal{R}_j \mathcal{A}_n - I)F, g \rangle_0^2 = \sup_{\|g\|_0 \leq 1} \|(\mathcal{R}_j \mathcal{A}_n - I)^* g\|_{-\alpha}^2.$$

By the norm duality (2.1), the right side is equal to

$$\begin{aligned} \sup_{\|g\|_0 \leq 1} \sup_{\|f\|_\alpha \leq 1} \langle f, (\mathcal{R}_j \mathcal{A}_n - I)^* g \rangle_0^2 &\leq \sup_{\|f\|_\alpha \leq 1} \|(\mathcal{R}_j \mathcal{A}_n - I)f\|_0^2 \\ &\lesssim \left(\frac{\delta(n, \alpha + \gamma)}{\delta(j, \gamma)} + \delta(j, \alpha) \right)^2 \lesssim \delta(j, \alpha)^2, \end{aligned}$$

when $\delta(j, s) = j^{-s/d}$, $n \gg j$ and (5.9).

The strong second moment of the Gaussian variable $(\mathcal{R}_j \mathcal{A}_n - I)F$ is equal to the trace of its covariance operator. As

$$\text{Trace}(S^* S) = \sum_i \|S\phi_i\|^2 = \sum_i \sum_j \langle S\phi_i, \phi_j \rangle^2 = \sum_i \|S^* \phi_i\|^2 = \|S\|_{HS}^2,$$

we have

$$\begin{aligned} \mathbb{E} \|(\mathcal{R}_j \mathcal{I}_n \mathcal{A} - I)F\|^2 &= \|(I - \mathcal{R}_{j_n} \mathcal{I}_n \mathcal{A})L^{-\alpha}\|_{HS}^2 \\ &\leq \|\mathcal{R}_{j_n}\|^2 \|\mathcal{I}_n - I : G_{\alpha+\gamma} \rightarrow G_0\|_{HS}^2 \|\mathcal{A}\|^2 \\ &\quad + \|(I - \mathcal{R}_{j_n} \mathcal{A}) : H_\alpha \rightarrow H_0\|_{HS}^2 \|L^{-\alpha} : H_0 \rightarrow H_\alpha\|^2 \\ &\lesssim_{(\mathcal{A}, L)} \|\mathcal{R}_{j_n}\|^2 \|\mathcal{I}_n - I : G_{\alpha+\gamma} \rightarrow G_0\|_{HS}^2 + \|(I - \mathcal{R}_{j_n} \mathcal{A}) : H_\alpha \rightarrow H_0\|_{HS}^2, \end{aligned} \quad (8.23)$$

where the first inequality is because of an elementary application of triangle inequality and a norm estimation of the operators in Hilbert spaces (see Proposition A.20).

In the following argument we will use some results from approximation number (and more generally s-number) in Hilbert spaces. The necessary material is collected below, and for the detail we refer to Section 2.4. Recall that the l th *approximation number* of a bounded linear operator $T : X \rightarrow Y$ between normed spaces is defined as

$$a_l(T : G \rightarrow H) = \inf_{U : \text{Rank } U < l} \|T - U\|_{X \rightarrow Y}$$

where the infimum is taken over all linear operators $U : X \rightarrow Y$ of rank (i.e., dimension of the range space) strictly less than l , and the norm on the right is the

operator norm $\|T - U\|_{X \rightarrow Y} = \sup_{f: \|f\|_X \leq 1} \|(T - U)f\|_Y$. Approximation number is in fact an example of a more general concept called *s-numbers*, which also include singular values. We will use the following fact: on Hilbert spaces there is only one s-number, i.e. singular value and approximation number are identical to each other.

Introduce a temporary notation $S := (I - \mathcal{R}_{j_n} \mathcal{A}) : H_\alpha \rightarrow H_0$, whose operator norm is $\|S : H_\alpha \rightarrow H_0\| \lesssim j_n^{-\alpha/d}$ by (5.9). Since $I : H_\alpha \rightarrow H_0$ is compact and $\mathcal{R}_{j_n} \mathcal{A}$ is of finite rank, S is compact as well. Hence it has singular values $s_l = a_l((I - \mathcal{R}_{j_n} \mathcal{A}) : H_\alpha \rightarrow H_0)$, and in particular, the first singular value is bounded by its operator norm, i.e. $s_1 \lesssim j_n^{-\alpha/d}$. Since $a_l(I : H_\alpha \rightarrow H_0) \simeq \delta(l, \alpha)$ (see the discussion following (2.6)), we have

$$a_l(S : H_\alpha \rightarrow H_0) \gtrsim \delta(j_n + l, \alpha),$$

which is because that the infimum on the right hand side is taken with all operators with rank less than $j_n + l$, while on the left hand side there is a fixed part $\mathcal{R}_{j_n} \mathcal{A}$ and the infimum is only taken over the operators with rank less than l . On the other hand, with the operator $U = \mathcal{R}_{j_n+l} \mathcal{A} - \mathcal{R}_{j_n} \mathcal{A}$ of rank l , we have

$$a_l(S : H_\alpha \rightarrow H_0) \leq \|I - \mathcal{R}_{j_n} \mathcal{A} - U\| \lesssim \delta(j_n + l, \alpha).$$

Combining the previous inequalities, we obtain

$$s_l = a_l(S : H_\alpha \rightarrow H_0) \simeq \delta(j_n + l, \alpha),$$

which leads to

$$\|(I - \mathcal{R}_{j_n} \mathcal{A}) : H_\alpha \rightarrow H_0\|_{HS}^2 = \sum_{l=1}^{\infty} (j_n + l)^{-2\alpha/d} = \sum_{l=j_n+1}^{\infty} l^{-2\alpha/d} \leq j_n^{1-2\alpha/d},$$

where we used the estimate $\sum_{i>j_n} i^{-b} \leq j_n^{1-b}/(b-1)$ for $b > 1$.

With the same argument above,

$$\|\mathcal{I}_n - I : G_{\alpha+\gamma} \rightarrow G_0\|_{HS}^2 \leq n^{1-2(\alpha+\gamma)/d}.$$

Since $\alpha + \gamma > d/2$, $j_n \leq n\varepsilon_n^2 < n$ and $\|\mathcal{R}_{j_n}\| \lesssim 1/\delta(j_n, \gamma) = j_n^{-\gamma/d}$, the first term in (8.23) is of order strictly smaller than the second term.

Since the first moment of $\|(\mathcal{R}_j \mathcal{A}_n - I)F\|_0$ is bounded by the root of its second moment, the lemma follows by Borell's inequality (see e.g. Lemma 3.1 and subsequent discussion in [67]). \square

For $t^2 = n\varepsilon_n^2/(4bj_n^{2\alpha/d})$ the bound in the preceding lemma becomes $e^{-4n\varepsilon_n^2}$. Hence (8.11) is satisfied for

$$\eta_n \gtrsim \sqrt{n\varepsilon_n}/j_n^{\alpha/d} + j_n^{1/2-\alpha/d}.$$

Here we choose ε_n the minimal solution that satisfies the direct prior mass condition (8.10), given in (8.22). Next we solve for η_n under the constraints, (8.8) and

(8.9). The first of these constraints, $j_n \leq n\varepsilon_n^2$, shows that the first term on the right side of the preceding display always dominates the second term. Therefore, we obtain the requirements $j_n \leq n\varepsilon_n^2$ and

$$\begin{aligned}\eta_n &\geq \sqrt{n} n^{-(\beta \wedge (\alpha - d/2) + \gamma)/(2\alpha + 2\gamma)}, \\ \eta_n &\geq n^{-(\beta \wedge (\alpha - d/2) + \gamma)/(2\alpha + 2\gamma)} j_n^{\gamma/d}, \\ \eta_n &\geq j_n^{-\beta/d}.\end{aligned}$$

Depending on the relation between α and $\beta + d/2$, two situations need to be discussed separately.

(i) $\alpha < \beta + d/2$. We choose $j_n \simeq n^{d/(2\alpha + 2\gamma)} = n\varepsilon_n^2$ and then see that the first two requirements in the preceding display both reduce to $\eta_n \geq n^{-(\alpha - d/2)/(2\alpha + 2\gamma)}$, while the third becomes $\eta_n \geq n^{-\beta/(2\alpha + 2\gamma)}$ and becomes inactive.

(ii) $\alpha > \beta + d/2$. We choose $j_n \simeq n^{d/(2\alpha + 2\gamma)} \leq n\varepsilon_n^2$, and then see that all three requirements reduce to $\eta_n \geq n^{-\beta/(2\alpha + 2\gamma)}$.

Finally apply Theorem 8.6 to complete the proof.

8.6.3 Proof of Theorem 8.13

Let Π_τ denote the zero-mean Gaussian distribution on H with covariance operator $\tau^2 L^{-2\alpha}$ (where $\alpha > d/2$). The following lemmas are the counterparts of Lemmas 6.14 to 6.16.

Lemma 8.16. *Under the assumptions of Theorem 8.13, for $f_0 \in H_\beta$ and $(d/2 - \gamma) \vee 0 < \beta \leq \alpha$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi_\tau(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon) \lesssim \frac{1}{\tau^2} \left(\frac{1}{\varepsilon}\right)^{(2\alpha - 2\beta)/(\beta + \gamma)} + \left(\frac{\tau}{\varepsilon}\right)^{d/(\alpha + \gamma - d/2)}.$$

Lemma 8.17. *Under the assumptions of Theorem 8.13, for $f_0 \in H_\beta$ and $(d/2 - \gamma) \vee 0 < \beta \leq \alpha$, as $\varepsilon \downarrow 0$,*

$$-\log \Pi_\tau(f : \|f\|_0 < \varepsilon) \gtrsim \left(\frac{\tau}{\varepsilon}\right)^{d/(\alpha - d/2)}.$$

Lemma 8.18. *Under the assumptions of Theorem 8.13, there exist $a, b > 0$ such that, for every $j \in \mathbb{N}$ and $x, \tau > 0$,*

$$\Pi_\tau(f : \|f^{(j)} - f\|_0 > \tau x + \tau a j^{1/2 - \alpha/d}) \leq e^{-bx^2 j^{2\alpha/d}}$$

Proofs. The proofs are identical to the one in Section 6.6.3, and hence they are omitted. \square

As preparation for the proof of Theorem 8.13, we first show that the minimax rate can be obtained by a Gaussian prior with the deterministic scaling, dependent on β , given by

$$\tau_n = n^{(\alpha - d/2 - \beta)/(2\beta + 2\gamma + d)}. \quad (8.24)$$

Theorem 8.19. *Assume the conditions on the Hilbert scale, the forward operator \mathcal{A} and the true parameter f_0 in Theorem 8.11 hold. Suppose that the priors Π are zero-mean Gaussian with covariance operators $\tau_n^2 L^{-2\alpha}$ with τ_n as given in (8.24) and $\alpha > (d - \gamma) \vee d/2$. Then for $(d/2 - \gamma) \vee 0 < \beta \leq \alpha$, the posterior distribution satisfies, for sufficiently large $M > 0$,*

$$\Pi_n \left(f : \|f - f_0\|_0 > Mn^{-\beta/(2\beta+2\gamma+d)} \mid Y^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0.$$

Proof. The theorem is a corollary to Theorem 8.6. The proof follows the same lines as the proof of Theorem 8.11. By Lemma 8.16, inequality (8.10) is satisfied for

$$\varepsilon_n \gtrsim n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}.$$

By Lemma 8.18, inequality (8.11) is satisfied for

$$\eta_n \gtrsim \tau_n (\sqrt{n} \varepsilon_n j_n^{-\alpha/d} + j_n^{1/2-\alpha/d}).$$

We choose $j_n \simeq n\varepsilon_n^2$, and the minimal solution $\varepsilon_n = n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}$ to the second last display. It is then straightforward to verify that (8.8), (8.9) and (8.11) are satisfied for $\eta_n \simeq n^{-\beta/(2\beta+2\gamma+d)}$. \square

Theorem 8.13 is a corollary of Theorem 8.7, with the choices

$$\begin{aligned} \eta_n &\simeq n^{-\beta/(2\beta+2\gamma+d)}, & \varepsilon_n &\simeq n^{-(\beta+\gamma)/(2\beta+2\gamma+d)}, \\ j_n &\simeq n\varepsilon_n^2 = n^{d/(2\beta+2\gamma+d)}. \end{aligned}$$

Conditions (8.7), (8.8), and (8.9) are satisfied for these choices. It remains to verify (8.10), and (8.17)–(8.16).

For ease of notation, for the moment, define η_n and ε_n as in the preceding display, with exact equality (i.e., with the constant set equal 1). Let τ_n be the ‘optimal’ scaling rate defined in (8.24).

Verification of (8.10). For $\tau \simeq \tau_n$ and $\varepsilon \simeq \varepsilon_n$ as given and $\beta \leq \alpha$, both terms in the right side of Lemma 8.16 are of the order $n\varepsilon_n^2$. The lemma yields, for $\tau_n \leq \tau \leq 2\tau_n$ and some constant $a_1 > 0$,

$$-\log \Pi_\tau(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon_n) \leq a_1 n\varepsilon_n^2.$$

This shows that

$$\begin{aligned} \Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon_n) &= \int_0^\infty \Pi_\tau(f : \|\mathcal{A}f - \mathcal{A}f_0\| < \varepsilon_n) dQ(\tau) \\ &\geq e^{-a_1 n\varepsilon_n^2} Q(\tau_n, 2\tau_n). \end{aligned}$$

If $\alpha - d/2 < \beta$, then $\tau_n \rightarrow 0$, and Condition 8.12 on Q gives that

$$-\log Q(\tau_n, 2\tau_n) \lesssim \tau_n^{-2} = n^{(2\beta-2\alpha+d)/(2\beta+2\gamma+d)} \leq n^{d/(2\beta+2\gamma+d)} = n\varepsilon_n^2,$$

if $\beta \leq \alpha$. If $0 < \beta < \alpha - d/2$, then $\tau_n \rightarrow \infty$, and Condition 8.12 on Q gives that

$$\begin{aligned} -\log Q(\tau_n, 2\tau_n) &\lesssim \tau_n^{d/(\alpha-d/2)} = n^{(d(\alpha-d/2-\beta)/(\alpha-d/2)(2\beta+2\gamma+d))} \\ &\leq n^{d/(2\beta+2\gamma+d)} = n\varepsilon_n^2. \end{aligned}$$

Finally if $\alpha - d/2 = \beta$, then $\tau_n = 1$ and $Q(\tau_n, 2\tau_n) \gtrsim 1$. Thus in all three cases $Q(\tau_n, 2\tau_n)$ is bounded below by a power of $e^{-n\varepsilon_n^2}$. Combining this with the preceding, we see that $\Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\|_n \leq \varepsilon_n) \geq e^{-a_2 n \varepsilon_n^2}$, for some positive constant a_2 , which we can take bigger than 1. Then (8.10) is satisfied for ε_n equal to $\sqrt{a_2}$ times the current ε_n .

Verification of (8.16). Lemma 8.17 gives that

$$\Pi_\tau(f : \|f - f_0\|_0 < 2\eta_{n,\tau}) \leq \Pi_\tau(f : \|f\|_0 < 2\eta_{n,\tau}) \leq e^{-a_3(\tau/\eta_{n,\tau})^{d/(\alpha-d/2)}},$$

for some constant a_3 . This is bounded above by $e^{-4a_2 n \varepsilon_n^2}$ if

$$\eta_{n,\tau} = 2a_4 \tau n^{(d/2-\alpha)/(2\beta+2\gamma+d)} = 2a_4 \tau \eta_n / \tau_n,$$

for a sufficiently small constant $a_4 > 0$.

Verification of (8.17). Choosing $x = a_4 \eta_n / \tau_n = \eta_{n,\tau} / (2\tau)$ in Lemma 8.18, we see that the left side of (8.17) is bounded above by $e^{-4a_2 n \varepsilon_n^2}$ if j_n satisfies

$$a j_n^{1/2-\alpha/d} \leq a_4 \eta_n / \tau_n, \quad \text{and} \quad b a_4^2 (\eta_n / \tau_n)^2 j_n^{2\alpha/d} \geq 4a_2 n \varepsilon_n^2.$$

Both inequalities become equalities for j_n of the order $j_n \simeq n^{d/(2\beta+2\gamma+d)}$, as indicated at the beginning of the proof. Since $1/2 - \alpha/d < 0$ and $2\alpha/d > 0$, the left side of the first inequality is decreasing in j_n and the left side of second inequality is increasing. Thus both inequalities are satisfied for $j_n = a_5 n^{d/(2\beta+2\gamma+d)}$ and a sufficiently large constant a_5 .

Finally we choose ε_n and j_n in Theorem 8.7 equal to $\sqrt{a_2}$ and a_5 times the orders indicated at the beginning of the proof. Then (8.7) is satisfied, and (8.8) and (8.9) are satisfied if η_n is chosen of the indicated order times a sufficiently large constant.

Part III

Evolution Equations

Heuristically speaking, a *stochastic evolution equation* in infinite dimensions describes a stochastic dynamical¹ system (a random process) whose trajectory (path) is in an infinite dimensional space. It naturally arises when the states of the dynamical system are infinite dimensional. Infinite dynamical systems pervasively exist in many quantitative fields. Examples of infinite dimensional evolution equations include population dynamics from biology, time dependent field equations emerged from physics, evolution of financial instruments such as the term structure of interest rates, whose details can be found in the introduction chapters in [23, 85]. Stochastic evolution equations are often used interchangeably² with *stochastic partial differential equations* (SPDEs), and we will also adopt this convention.

In this part, our goal is to use the Bayesian nonparametric approach to recover the parameters in the model, which is formally defined in Section 9.5, under the small noise asymptotic regime. That is, for $t \in [0, T]$, as n goes to infinity, we continuously observe the solution $X(t)$ of the SPDE

$$\begin{cases} dX^{(n)}(t) + \mathcal{L}X^{(n)}(t) dt = f(t) dt + \frac{1}{\sqrt{n}}B dW(t) \\ X^{(n)}(0) = u \in H \end{cases}, \quad (\text{III.1})$$

where u is the initial condition. We will address that the parameter estimations of SPDEs are usually inverse problems, to which the general mechanism developed in Section 4.3 can be applied with some modifications. Consequently, the contraction rates are obtained for the recovery of the initial condition u and the drift f . For the purpose of streamlining the arguments, we only work with Gaussian priors, however it is noteworthy that the proof does not rely on any properties of the prior on conjugacy or Gaussianity, and hence, in principle, the method is applicable to other types of priors as well, such as random series priors.

The part is organized as follows. In Chapter 9, we introduce the necessary mathematical components to formalize the model (III.1). After that, we study the Bayesian statistical inference for SPDE in Chapter 10. The recovery of the initial condition u is examined Section 10.1, and the recovery of drift f is studied in Section 10.2.

¹Differential dynamical system.

²Although arguably stochastic evolution equation is a broader term than SPDE.

Chapter 9

Linear Evolution Equations

Stochastic partial differential equations (SPDEs) provide a powerful toolkit to study dynamical systems with stochastic nature. In the last few decades, it has gained much interest due to the wide range of applications in physics and finance. One way to study SPDEs is using the theory of linear evolution equations in infinite-dimensional spaces. This analytical approach treats dynamical systems as vector-valued ordinary differential functions, whose state spaces are function spaces. In this chapter we collect the basic results from this approach. For a detailed treatment on the subject, we refer to the standard reference [23] and the more recent monographs [68, 85].

First we recall some standard definitions from stochastic processes. Fix a probability space $(\Omega, \mathcal{F}, \mathbb{P})$, and a positive number $T < \infty$.

Definition 9.1. With an index set $\mathcal{J} \subset \mathbb{R}_0^+$, a *filtration* $\{\mathcal{F}_t\}_{\mathcal{J}}$ is a family of sub σ -algebras of \mathcal{F} such that $\mathcal{F}_s \subset \mathcal{F}_t$ if $s < t$. The following conditions are known as the *usual conditions*.

- (i) Completeness: $A \in \mathcal{F}_0$, for all $A \in \mathcal{F}$ such that $\mathbb{P}(A) = 0$.
- (ii) Right continuity: $\mathcal{F}_t = \mathcal{F}_{t+} := \bigcap_{s \in \mathcal{J}: s > t} \mathcal{F}_s$, for all $t \in \mathcal{J}$.

A filtration that satisfies the usual conditions is also called *normal*.

Definition 9.2. Let $(E, \|\cdot\|)$ be a separable Banach space. An E -valued process $\{M(t) : t \geq 0\}$ on $(\Omega, \mathcal{F}, \mathbb{P})$ with filtration $\{\mathcal{F}_t\}_{0 \leq t \leq T}$ is a \mathcal{F}_t -*martingale* if

- (i) $\mathbb{E}\|M(t)\| < \infty$, for all $t \geq 0$,
- (ii) $M(t)$ is \mathcal{F}_t -measurable for all $t \geq 0$,
- (iii) and $\mathbb{E}(M(t)|\mathcal{F}_s) = M(s)$ almost surely in \mathbb{P} , for all $0 \leq s \leq t < \infty$.

In addition, denote the space of all E -valued continuous square integrable martingales by $\mathcal{M}_T^2(E)$, which is a Banach space equipped with the norm

$$\|M\|_{\mathcal{M}_T^2} := \sup_{t \in [0, T]} (\mathbb{E}\|M(t)\|^2)^{1/2} = (\mathbb{E}\|M(T)\|^2)^{1/2}.$$

For the proof, see, e.g. Proposition 3.10, [23].

9.1 \mathcal{Q} -Wiener Processes

The characteristic differing stochastic evolution systems from the deterministic ones is the appearance of stochastic noise. As the state of the system is infinite dimensional, the noise is also expected to be infinite dimensional. \mathcal{Q} -Wiener processes, a generalization of classical Wiener processes to vector-valued processes, will be used to model the infinite dimensional noise.

Let U be a *separable* infinite-dimensional Hilbert space, and $\mathcal{Q} \in L(U)$ be a bounded linear operator on U satisfying, for any $u, v \in U$,

- (i) nonnegative: $\langle \mathcal{Q}u, u \rangle \geq 0$,
- (ii) symmetric: $\langle \mathcal{Q}u, v \rangle = \langle u, \mathcal{Q}v \rangle$,
- (iii) nuclear: $\text{Trace } \mathcal{Q} < \infty$.

Definition 9.3. On the probability space $(\Omega, \mathcal{F}, \mathbb{P})$, a (standard) \mathcal{Q} -Wiener process is a U -valued process $\{W(t) : t \in [0, T]\}$ satisfying

- (i) $W(0) = 0$,
- (ii) W is continuous almost surely in \mathbb{P} ,
- (iii) W has independent increments, i.e.

$$W(t_1), W(t_2) - W(t_1), \dots, W(t_n) - W(t_{n-1}),$$

are independent, for all $0 \leq t_1 < \dots < t_n \leq T$,

- (iv) for all $0 \leq s \leq t \leq T$, $W(t) - W(s)$ is distributed as a centred Gaussian on U with covariance operator $(t - s)\mathcal{Q}$.

\mathcal{Q} is called the covariance operator of \mathcal{Q} -Wiener process.

Remark 9.4. For a \mathcal{Q} -Wiener process $\{W(t) : t \in [0, T]\}$, there always exists a normal filtration $\{\mathcal{F}_t\}_{0 \leq t \leq T}$ such that:

1. $W(t)$ is adapted to $\{\mathcal{F}_t\}_{0 \leq t \leq T}$, i.e. $W(t)$ is \mathcal{F}_t -measurable for all $0 \leq t \leq T$,
2. and $W(t) - W(s)$ is independent of \mathcal{F}_s for all $0 \leq s \leq t \leq T$.

See Proposition 2.1.13 in [68].

Similar to the Karhunen-Loève expansion of Gaussian elements, the \mathcal{Q} -Wiener process has a concrete representation.

Proposition 9.5 (Presentation of \mathcal{Q} -Wiener Process). *Let $\{e_i\}_{i \in \mathbb{N}}$ be the eigenfunctions of \mathcal{Q} with the corresponding eigenvalues $\{q_i\}_{i \in \mathbb{N}}$. A U -valued process $W(t)$ is \mathcal{Q} -Wiener if and only if there exists a sequence of independent ordinary Wiener processes $\{W_i : i \in \mathcal{J}\}$ with $\mathcal{J} = \{i : q_i > 0\}$ on $(\Omega, \mathcal{F}, \mathbb{P})$ such that*

$$W(t) = \sum_{i \in \mathcal{J}} \sqrt{q_i} e_i W_i(t), \quad t \in [0, T], \quad (9.1)$$

which converges in $L^2(\mathbb{P})$, and always has a \mathbb{P} -a.s. continuous version, i.e.

$$\mathbb{P}[W(t) \in C([0, T], U)] = 1.$$

Proof. See Proposition 2.1.10 in [68]. □

9.2 Stochastic Integrals in Hilbert Spaces

A H -valued stochastic integral with respect to $W^{\mathcal{Q}}(t)$ can be developed in a similar manner as for the scalar valued ordinary stochastic integral. We outline the construction of the integral and summarize a few useful results. A detailed treatment can be found in Chapter 2 in [68] and Section 4.2 & 4.3 in [23].

Similar to ordinary stochastic integrals, the construction relies on a space of martingales and a class of elementary processes. Recall that the space $\mathcal{M}_T^2(H)$ of all H -valued continuous square integrable martingales $\{M(t) : t \in [0, T]\}$ is a Banach space equipped with the norm

$$\|M\|_{\mathcal{M}_T^2} := \sup_{t \in [0, T]} (\mathbb{E}\|M(t)\|^2)^{1/2} = (\mathbb{E}\|M(T)\|^2)^{1/2}.$$

The standard machinery to construct stochastic integrals can be summarised as follows.

- (I) For a class \mathcal{E} of elementary processes, define a linear mapping

$$\text{Int} : \mathcal{E} \ni \Phi \rightarrow \int_0^t \Phi(s) dW(s) =: \Phi \cdot W(t),$$

which is an element in $\mathcal{M}_T^2(H)$.

- (II) Find a norm on \mathcal{E} such that $\text{Int} : \mathcal{E} \rightarrow \mathcal{M}_T^2(H)$ is an isometry. Since $\mathcal{M}_T^2(H)$ is complete, Int is extended to the abstract completion $\bar{\mathcal{E}}$ of \mathcal{E} . Furthermore, an explicit representation is found for $\bar{\mathcal{E}}$.

- (III) Further extend the integral to local martingale by localization.

For completeness of the exposition, we also discuss some detail on the procedure. The readers familiar with ordinary stochastic integrals will immediately notice the similarity. Fix a probability space $(\Omega, \mathcal{F}, \mathbb{P})$ with a normal filtration, and denote $dt \otimes \mathbb{P}$ by \mathbb{P}_T , where dt is the Lebesgue measure. Let $W^{\mathcal{Q}}$ be the \mathcal{Q} -Wiener process introduced in Section 9.1. In this section, $L(U, H)$ denotes the space of the bounded operators from U to H .

- (I) An $L(U, H)$ -valued process $\Phi(t), t \in [0, T]$ with normal filtration is *elementary* if

$$\Phi(t) = \sum_{m=0}^{k-1} \Phi_m \mathbb{1}_{(t_m, t_{m+1}]}(t), \quad t \in [0, T],$$

for some $0 < t_0 < \dots < t_k = T$, $\Phi_m : \Omega \rightarrow L(U, H)$ is \mathcal{F}_{t_m} -measurable with respect to the strong Borel σ -algebra on $L(U, H)$, $0 \leq m \leq k - 1$, and each

Φ_m takes only a finite number of values in $L(U, H)$. The stochastic integral of elementary processes is defined by

$$\text{Int}(\Phi)(t) := \int_0^t \Phi(s) dW(s) := \sum_{m=0}^{k-1} \Phi_m(W(t_{m+1} \wedge t) - W(t_m \wedge t)), \quad t \in [0, T],$$

and it belongs to $\mathcal{M}_T^2(H)$.

(II) Recall that for $\Phi \in L(U, H)$, $\Phi \circ \mathcal{Q}^{1/2}$ is a Hilbert-Schmidt operator in $S_2(U, H)$ (see Section A.3). The previous stochastic integral satisfies the following *Itô isometry*,

$$\left\| \int_0^\cdot \Phi(s) dW(s) \right\|_{\mathcal{M}_T^2}^2 = \mathbb{E} \left(\int_0^T \|\Phi(s) \circ \mathcal{Q}^{1/2}\|_{HS}^2 ds \right) =: \|\Phi\|_T^2. \quad (9.2)$$

Notice that $\|\cdot\|_T$ is only a semi-norm on \mathcal{E} , and two elementary processes belonging to one equivalent class are not necessarily \mathbb{P}_T -a.e. equal, because the equivalence only occurs on $\mathcal{Q}^{1/2}(U)$ \mathbb{P}_T -a.e.

Because of the Itô isometry,

$$\text{Int} : (\mathcal{E}, \|\cdot\|_T) \rightarrow (\mathcal{M}_T^2, \|\cdot\|_{\mathcal{M}_T^2})$$

is an isomorphism, and consequently there uniquely exists an isometric extension of mapping Int to the abstract completion $\bar{\mathcal{E}}$ of \mathcal{E} with respect to $\|\cdot\|_T$.

Now we prepare for the explicit presentation of $\bar{\mathcal{E}}$. Let $U_0 := \mathcal{Q}^{1/2}(U)$, which is a Hilbert space with the induced inner product $\langle \cdot, \cdot \rangle_0 := \langle \mathcal{Q}^{-1/2} \cdot, \mathcal{Q}^{-1/2} \cdot \rangle$ (see Lemma A.11). We denote the space of all Hilbert-Schmidt operators¹ from U_0 to H by L_2^0 , which is a separable Hilbert space equipped with the norm

$$\|\Phi\|_{L_2^0}^2 = \|\Phi \mathcal{Q}^{1/2}\|_{HS}^2 = \text{Trace}[(\Phi \mathcal{Q}^{1/2})(\Phi \mathcal{Q}^{1/2})^*]. \quad (9.3)$$

In particular, the space $L(U, H)$ can be embedded into L_2^0 by restricting the domain of operators to U_0 . Let $L(U, H)_0 = \{T|_{U_0} : T \in L(U, H)\}$ denote the space of the restricted operators. Then, $L(U, H)_0 \subset L_2^0$ and

$$\|\Phi\|_T = \left[\mathbb{E} \int_0^T \|\Phi(s)\|_{L_2^0}^2 ds \right]^{1/2},$$

for $\Phi \in \mathcal{E}$.

The final claim is that $\bar{\mathcal{E}}$ is given by

$$\begin{aligned} \mathcal{N}_W^2(0, T; H) &:= \{\Phi : [0, T] \times \Omega \rightarrow L_2^0 \mid \Phi \text{ is predictable and } \|\Phi\|_T < \infty\} \\ &= L^2([0, T] \times \Omega, dt \otimes \mathbb{P}; L_2^0), \end{aligned}$$

which we also write $\mathcal{N}_W^2(H)$ and \mathcal{N}_W^2 for brevity. It is not difficult to see $\mathcal{E} \subset \mathcal{N}_W^2$. In addition, because \mathcal{N}_W^2 is in fact a Bochner space (see Section 1.2.b in [49]), it is complete. The preceding claim is proved by showing that \mathcal{E} is a dense subset of \mathcal{N}_W^2 (Proposition 2.3.8 in [68]).

¹The consistent notation is S_2^U with $U = \mathcal{Q}^{1/2}(U)$. Here we adopt the conventional notations from the SPDE literature.

(III) The integrals can be further extended to local martingales with the following class of integrands,

$$\mathcal{N}_W(0, T; H) := \left\{ \Phi : [0, T] \times \Omega \rightarrow L_2^0 \mid \Phi \text{ is predictable and } \mathbb{P} \left(\int_0^T \|\Phi(s)\|_{L_2^0}^2 ds < \infty \right) = 1 \right\}.$$

The detail is omitted here, since we will only consider the martingale case. We only remark that the obvious relation $\mathcal{N}_W^2(0, T; H) \subset \mathcal{N}_W(0, T; H)$, and some results below will be stated with the more general space.

Following the procedure described above, the isometric extension of Int to \mathcal{N}_W^2 ,

$$\text{Int} : \Phi \in \mathcal{N}_W^2 \mapsto \int_0^t \Phi(s) dW^{\mathcal{Q}}(s), \quad t \in [0, T],$$

defines the stochastic integral.

Remark 9.6. Notice that the stochastic integral is defined with a class of predictable processes. However, in this thesis we will only consider the case of deterministic integrands.

Not surprisingly, the H -valued stochastic integral shares many similar properties as the ordinary stochastic integral. For example, the integration is interchangeable with other linear operations, as stated in Lemma 9.7 below.

Lemma 9.7 (Lemma 2.4.1 in [68]). *Let $\Phi \in \mathcal{N}_W(H)$ and $\mathcal{T} \in L(H, \tilde{H})$, where \tilde{H} is another separable Hilbert space. Then the process $\mathcal{T}\Phi(t), t \in [0, T]$, is an element of $\mathcal{N}_W(\tilde{H})$, and*

$$\mathcal{T} \left(\int_0^T \Phi(s) dW(s) \right) = \int_0^T \mathcal{T}(\Phi(s)) dW(s)$$

\mathbb{P} -almost surely.

In addition, the stochastic integral also admits a series representation, as in the following lemma.

Lemma 9.8 (Proposition 2.4.5 in [68]). *If $\Phi \in \mathcal{N}_W^2(H)$, then*

$$\int_0^t \Phi(s) dW(s) = \sum_{i \in \mathbb{N}} \sqrt{q_i} \int_0^t \Phi(s)(e_i) dW_i(s), \quad t \in [0, T],$$

\mathbb{P} -almost surely, where q_i, e_i, W_i are as in Proposition 9.5 and the sum on the right-hand side converges in $L^2(\mathbb{P})$ and realises in $C([0, T], U)$ almost surely.

For the detailed properties of stochastic integrals, we refer to Section 2.4 in [68] and Section 4.3 in [23].

9.3 Extension of Stochastic Integrals

In the previous section, we have developed the integration theory under the condition that the covariance \mathcal{Q} is a trace class operator. In this section we are going to relax the condition, and consequently, the class of integrators of the stochastic integral is extended.

9.3.1 Cylindrical Wiener Process

We first extend the H -valued Wiener process.

Recall that the Wiener process in Section 9.1 admits the representation

$$W(t) = \sum_{k \in \mathbb{N}} \sqrt{q_k} W_k(t) e_k, \quad t \in [0, T],$$

where $\{e_k\}_{k \in \mathbb{N}}$ is an orthonormal basis for $U_0 = \mathcal{Q}^{1/2}(U)$ and $W_k(t), k \in \mathbb{N}$, are independent real-valued Wiener processes. The convergence of the series in $L^2(\mathbb{P})$ is due to the fact that the inclusion $U_0 \subset U$ is Hilbert-Schmidt (see Proposition 2.1.10, [68]).

Now let \mathcal{Q} be an operator satisfying the properties in Section 9.1 but not necessarily nuclear, that is, $\text{Trace } \mathcal{Q} = \sum_k q_k < \infty$ is not required. Let $U_0 = \mathcal{Q}^{1/2}(U)$ with the induced inner product and $\{e_k\}_{k \in \mathbb{N}}$ be an orthonormal basis for U_0 . Let U_1 be another Hilbert space such that there exists a Hilbert-Schmidt embedding $\mathcal{J} : U_0 \rightarrow U_1$.

Remark 9.9. The space U_1 always exists. For example, let $U_1 = U$ and define

$$\mathcal{J} : U_0 \rightarrow U, \quad f \mapsto \sum_{k \in \mathbb{N}} \rho_k \langle f, e_k \rangle_U e_k,$$

with a sequence $\{\rho_k\}_{k \in \mathbb{N}}$ such that $\sum_k \rho_k^2 < \infty$. Then, \mathcal{J} is injective and Hilbert-Schmidt, c.f. Section 3.3.

Then, there exists a Wiener process on the larger Hilbert space U_1 , associated with the previously introduced \mathcal{Q} .

Proposition 9.10. *Under the previously introduced condition, define $\mathcal{Q}_1 := \mathcal{J}\mathcal{J}^*$. Then, \mathcal{Q}_1 is nonnegative, symmetric and nuclear. The series*

$$W(t) = \sum_{k \in \mathbb{N}} W_k(t) \mathcal{J} e_k, \quad t \in [0, T],$$

converges in $\mathcal{M}_T^2(U_1)$ and defines \mathcal{Q}_1 -Wiener process on U_1 . Furthermore,

$$\mathcal{Q}_1^{1/2}(U_1) = \mathcal{J}(U_0),$$

and for all $u_0 \in U_0$,

$$\|u_0\|_0 = \|\mathcal{Q}_1^{-1/2} \mathcal{J} u_0\|_1 =: \|\mathcal{J} u_0\|_{\mathcal{Q}_1^{1/2}(U_1)},$$

i.e. $\mathcal{J} : U_0 \rightarrow \mathcal{Q}_1^{1/2}(U_1)$ is an isometry.

The process $\{W(t) : t \in [0, T]\}$ introduced in Proposition 9.10 is called a *cylindrical Wiener process* with *covariance* \mathcal{Q} in U (while it does not realise in U).

9.3.2 Stochastic Integral with Cylindrical Wiener Process

Let $W(t)$ be a cylindrical Wiener process with covariance \mathcal{Q} as previously introduced. Known from Section 9.2, a process $\Phi(t)$ is integrable with respect to $W(t)$ if it is predictable, $\Phi(t) \in L_2(\mathcal{Q}_1^{1/2}(U_1), H)$ for all $t \in [0, T]$, and

$$\mathbb{P} \left(\int_0^T \|\Phi(s)\|_{L_2(\mathcal{Q}_1^{1/2}(U_1), H)}^2 ds < \infty \right) = 1.$$

We adopt all the notations from Proposition 9.10 and recall $\mathcal{Q}_1^{1/2}(U_1) = \mathcal{J}(U_0)$. Then by the polarisation identity, for all $u, v \in U_0$,

$$\langle \mathcal{J}u, \mathcal{J}v \rangle_{\mathcal{Q}_1^{1/2}(U)} = \langle u, v \rangle_0,$$

which implies that $\mathcal{J}e_k$ is an orthonormal basis for $\mathcal{Q}_1^{1/2}(U_1)$. Consequently, we have

$$\Phi \in L_2^0 = L_2(\mathcal{Q}^{1/2}(U), H) \iff \Phi \circ \mathcal{J}^{-1} \in L_2(\mathcal{Q}_1^{1/2}(U_1), H),$$

because

$$\|\Phi\|_{L_2^0}^2 = \sum_k \langle \Phi e_k, \Phi e_k \rangle = \sum_k \langle \Phi \circ \mathcal{J}^{-1} \mathcal{J}e_k, \Phi \circ \mathcal{J}^{-1} \mathcal{J}e_k \rangle = \|\Phi \circ \mathcal{J}^{-1}\|_{L_2(\mathcal{Q}_1(U_0), H)}^2.$$

Then, for the cylindrical process $W(t)$, the integral can be defined as

$$\int_0^t \Phi(s) dW(s) := \int_0^t \Phi(s) \circ \mathcal{J}^{-1} dW(s), \quad t \in [0, T].$$

Apparently, the previously defined stochastic integral holds with the same class of integrands from Section 9.2. Therefore, the stochastic integral has been extended to cylindrical Wiener processes.

We end this section with several remarks. First, the integral is actually independent from the space U_1 . This is because the Itô isometry (9.2) together with (9.3) is independent from the bigger space U_1 . Second, when $\text{Trace } \mathcal{Q} < \infty$, we can simply take \mathcal{J} to be the identity map $\text{id} : U_0 \rightarrow U$, and then the definition coincides with the integral developed Section 9.2. In both statements above, the important fact is that the stochastic integral is uniquely defined with the space $\mathcal{Q}^{1/2}(U) (= \mathcal{Q}_1^{1/2}(U_1))$, which is the reproducing kernel Hilbert space of the Gaussian measure $\mathcal{N}_U(0, \mathcal{Q})$ (the radonified Gaussian measure $\mathcal{N}_{U_1}(0, \mathcal{Q}_1)$, see Section 3.3).

9.4 Deterministic Evolution Equations

We start with introducing the dynamical systems without noise. Let H be a separable Hilbert space. An evolution equation in H is given by

$$\begin{cases} u'(t) + \mathcal{L}u(t) = f(t), \\ u(0) = u_0 \in H, \end{cases} \quad (9.4)$$

where $u : [0, T] \rightarrow H$ is a H -valued function, $\mathcal{L} : D(\mathcal{L}) \subset H \rightarrow H$ is a linear operator, unbounded in general, with a dense domain $D(\mathcal{L})$ in H , and $u'(t)$ is the strong derivative of $u(t)$, i.e.

$$u'(t) = \lim_{h \rightarrow 0} \frac{u(t+h) - u(t)}{h},$$

where the limit is taken in the topology of H .

The existence and uniqueness of (9.4) with the initial condition U_0 is known as the *deterministic* abstract (nonhomogeneous) Cauchy problem. It can be answered in the language of C_0 -semigroup theory. When \mathcal{L} is a infinitesimal generator of a C_0 -semigroup $S(\cdot)$ in H and $U_0 \in L^p([0, T]; H)$, $p \in [1, \infty]$, the *strict* solution u of problem (9.4), i.e. a function u that belongs to $W^{1,p}([0, T]; H) \cap L^p([0, T]; D(\mathcal{L}))$ and satisfies (9.4), is given by the following *variation of constant formula*,

$$u(t) = S(t)U_0 + \int_0^t S(t-s)f(s) ds. \tag{9.5}$$

See Chapter 4 in [75] for the details.

9.5 Solutions of SPDEs

We now formally define the model (III.1). Assume that all Hilbert spaces in this section are separable. Recall the model

$$dX(t) + \mathcal{L}X(t) dt = f(t) dt + \frac{1}{\sqrt{n}} B dW^{\mathcal{Q}}(t), \tag{9.6}$$

with the initial condition $X(0) = u \in H$.

The items in (9.6) are defined as follows. $f : [0, T] \rightarrow H$ is a H -valued measurable function, $\mathcal{L} : D(\mathcal{L}) \subset H \rightarrow H$ is a densely defined unbounded linear operator, $B : U \rightarrow H$ is a linear operator, $W^{\mathcal{Q}}$ is a \mathcal{Q} -Wiener process, and u is a deterministic vector in H .

One key component in the development of SPDE is the stochastic integral with respect to \mathcal{Q} -Wiener process in Hilbert spaces, covered in Section 9.2. In addition, some language of functional analysis is also standard in this field. To ease the reading effort, we collect the minimal material on the syntax of functional analysis in Appendix A.

Concrete examples covered by the above model are the parabolic evolution equations in a bounded domain $\mathcal{D} \subset \mathbb{R}^d$ with a smooth boundary $\partial\mathcal{D}$, with time interval $[0, T]$. Time-space domain and boundary is denoted by $\mathcal{D}_T := [0, T] \times \mathcal{D}$ and $\Gamma_T := [0, T] \times \partial\mathcal{D}$ respectively. Let \mathcal{L} be a strongly elliptic differential operator of order $2m$. That is, given

$$\mathcal{L}(x, D) = \sum_{|k|_1 \leq 2m} a_k(x) D^k,$$

where $k = (k_1, \dots, k_d)$ are multi-indices, and $D^k := D_{x_1}^{k_1} \dots D_{x_d}^{k_d}$, its principle part $\mathcal{L}(x, D)' = \sum_{|k|_1=2m} a_k(x) D^k$ satisfies

$$(-1)^m \mathcal{L}(x, z)' \geq c|z|^{2m}.$$

We are going to introduce the assumptions which guarantee the unique existence of a solution for the model (9.6). We start with the items characterizing the deterministic dynamics.

Suppose $f \in L^2([0, T]; H)$, where H is a Hilbert space of functions defined on a compact domain $\mathcal{D} \subset \mathbb{R}^d$. Assume that there exists a densely defined positive self-adjoint operator, $\Lambda : D(\Lambda) \subset H \rightarrow H$ that has an eigensystem, i.e. eigenfunctions $\{\varphi_k\}_{k \in \mathbb{N}^d}$ with corresponding eigenvalues $\{\lambda_k\}_{k \in \mathbb{N}^d}$.

Remark 9.11 (Index of eigensystem). By Proposition 5.12 in [84], the existence of eigensystem is equivalent to the existence of purely discrete spectrum. Hence the index set can be chosen \mathbb{N} . Instead, we use \mathbb{N}^d to make the formula more naturally fit later use, and occasionally we switch back to \mathbb{N} , e.g. when using spectral integral representation as in the next paragraph.

If \mathcal{L} is a positive function $g(\Lambda)$ of the operator Λ , \mathcal{L} admits a spectral integral representation in terms of the spectral measure E_λ of Λ , i.e. for any $h \in D(\mathcal{L})$,

$$\mathcal{L}h = g(\Lambda)h = \int f(\lambda) dE_\lambda h = \int_0^\infty f(\lambda) dE_\lambda h,$$

where the last equality is due to the positivity of \mathcal{L} . Because of Remark 9.11, upon reindexing, the spectral integral can be written into the series form,

$$\int_0^\infty g(\lambda) dE_\lambda h = \sum_{k \in \mathbb{N}^d} g(\lambda_k) \langle \varphi_k, h \rangle \varphi_k.$$

We will mainly use the spectral integral form in the latter text as it slightly shortens the notation. By Theorem 6.14 in [84], $-\mathcal{L}$ generates a strongly continuous contraction semigroup $S(t)$, $0 \leq t < \infty$, which also admits a spectral integral representation, i.e.

$$e^{-t\mathcal{L}} := S(t) = \int_0^\infty e^{-tg(\lambda)} E_\lambda = \sum_{k \in \mathbb{N}^d} e^{-tg(\lambda_k)} \langle \varphi_k, \cdot \rangle \varphi_k. \quad (9.7)$$

The discussion above can be summarized as follows.

Assumption 9.12. We impose the following conditions on the model (9.6). The function f is in $L^2([0, T]; H)$, where H is a Hilbert space of functions defined on a compact domain $\mathcal{D} \subset \mathbb{R}^d$. There exists a densely defined positive self-adjoint operator, $\Lambda : D(\Lambda) \subset H \rightarrow H$ that has an eigensystem, i.e. eigenfunctions $\{\varphi_k\}_{k \in \mathbb{N}^d}$ with corresponding eigenvalues $\{\lambda_k = |k|_1\}_{k \in \mathbb{N}^d}$. Furthermore, the operator $\mathcal{L} = \Lambda^{(\nu)}$ is a function of Λ defined via the eigensystem, i.e. for all φ_k ,

$$\mathcal{L}\varphi_k = \left(\sum_{i=1}^d k_i^\nu \right) \varphi_k =: \ell_k \varphi_k,$$

where the constant $\nu \in \mathbb{N}$.

Remark 9.13. Recall $|\cdot|_p$ is the p -norm on \mathbb{R}^d . From Lemma 2.19, for any $q \in [0, \infty]$,

$$\ell_k = |\ell^\nu|_{1 \simeq d} |k^\nu|_{2 \simeq d} |k|_q^\nu,$$

where the constant is universally between 1 and d . We will frequently use this equivalence in the subsequent sections.

Now we move to the stochastic part. To facilitate the statistical investigation, we impose the following relationship between the operator \mathcal{L} characterizing the deterministic dynamics and the operators B and \mathcal{Q} expressing the stochastic propagation.

Assumption 9.14. Let $\{\lambda_k, \varphi_k\}_{k \in \mathbb{N}^d}$ and $\{q_k, e_k\}_{k \in \mathbb{N}^d}$ be the eigensystems of Λ and \mathcal{Q} , respectively. In addition to Assumption 9.12, we assume that the operator B in (9.6) is linear² from U to H and is diagonalized by $\{e_k; \varphi_k\}$, i.e. for all $k \in \mathbb{N}^d$, $Be_k = b_k \varphi_k$. Furthermore, we assume that for the eigenvalues b_k, q_k and λ_k , the following conditions hold:

- (i) Initial condition. With a positive constant $d/2 < \mu < \infty$,

$$\frac{b_k^2 q_k}{\ell_k} \simeq_d |k|^{-2\mu}. \quad (9.8)$$

- (ii) Drift. The operator B is bounded and with $p = \mu - \nu/2 \geq 0$, where μ is given in (i),

$$b_k^2 q_k \simeq_d |k|^{-2p}, \quad (9.9)$$

where the involved constant is independent of k .

Remark 9.15. The case (ii) in Assumption 9.14 is a consequence of (i) and Assumption 9.12. Situation (i) can be relaxed to $b_k^2 q_k / \ell_k$ being bounded from above and below by a polynomial decay of $|k|$. This will not affect the results for the recovery of initial condition, because the severe ill-posedness induces exponential decay, see Section 10.3.1. However, for the recovery of the drift, the rate will then depend on both of the upper and lower bounds. To simplify the exposition, we confine ourselves to the special case (9.9).

Under Assumption 9.14, with $\|T\|_{L^2}^2 = \text{Trace}[T\mathcal{Q}^{1/2}(T\mathcal{Q}^{1/2})^*]$, we have

$$\begin{aligned} & \int_0^T \|S(t)B\|_{L^2}^2 dt \\ &= \int_0^T \text{Trace}[S(t)B\mathcal{Q}B^*S^*(t)] dr = \sum_{i \in \mathbb{N}^d} \frac{b_i^2 q_i}{2\ell_i} (1 - e^{-2\ell_i T}) \\ &\leq \sum_{i \in \mathbb{N}^d} \frac{b_i^2 q_i}{\ell_i} \simeq_d \sum_{k \in \mathbb{N}^d} |k|^{-2\mu} \simeq \sum_{j \in \mathbb{N}} \sum_{|k|_1=j} j^{-2\mu} \\ &\simeq \sum_{j \in \mathbb{N}} \left(\binom{j+d}{d} - \binom{j+d-1}{d} \right) j^{-2\mu} \simeq \sum_{j \in \mathbb{N}} j^{d-1-2\mu} < \infty. \end{aligned} \quad (9.10)$$

²Not necessarily bounded.

We now introduce the concept of solutions to SPDEs. According to Theorem 5.4 in [23], under Assumption 9.12 and the assumption that the semigroup $S(t)$ satisfies (9.10), the unique analytical³ *weak* solution of (9.6) is given by the *stochastic* variation of constant formula, for $t \in [0, T]$,

$$X(t) = S(t)u + \int_0^t S(t-s)f(s) ds + \frac{1}{\sqrt{n}} \int_0^t S(t-s)B dW^{\mathcal{Q}}(s). \quad (9.11)$$

Remark 9.16. (9.11) is defined as the *mild solution* of (9.6), see page 161 in [23]. However, under Assumption 9.12 and (9.10), the weak solution and mild solution coincide (Theorem 6.7, [23]). For SPDEs, there are several concepts of solutions, which are not always equivalent. For details, we refer to the relevant sections in [23] and Appendix G in [68].

Using the covariance formula given in Theorem 5.2 in [23], the covariance of the stochastic integral is given by

$$\begin{aligned} \text{Trace Cov} \left(\int_0^t S(t-s)B dW^{\mathcal{Q}}(s) \right) &= \text{Trace} \left(\int_0^t S(t-s)BQB^*S^*(t-s)ds \right) \\ &= \sum_{k \in \mathbb{N}^d} \frac{b_k^2 q_k}{2\ell_k} \frac{1 - e^{-2\ell_k T}}{e^{2\ell_k T}} \leq \sum_{k \in \mathbb{N}^d} \frac{b_k^2 q_k}{\ell_k} < \infty, \end{aligned}$$

under Assumption 9.14. As a consequence, the stochastic component in the weak solution is a proper Gaussian process in H . This is in contrast to the well-known white noise model, which is almost surely not in H .

Now we discuss the series representation of (9.11). Let $\{\varphi_i\}_{i \in \mathbb{N}^d}$ be the orthonormal basis of H from assumption 9.12. Under Assumptions 9.12 and 9.14, the three items in (9.11) can be represented as follows.

Since $u \in H$, we have

$$u = \sum_{k \in \mathbb{N}^d} u_k \varphi_k, \quad \text{and} \quad S(t)u = \sum_{k \in \mathbb{N}^d} e^{-t\ell_k} u_k \varphi_k.$$

Since $L^2([0, T]; H) \cong H \otimes L^2([0, T]; \mathbb{R})$ (see Section 2.3.1), the drift term can be written in the form of $f(t) = \sum_{k \in \mathbb{N}^d} \varphi_k f_i(t)$. Consequently,

$$\int_0^t S(t-s)f(s) ds = \sum_{k \in \mathbb{N}^d} \varphi_k \int_0^t e^{-\ell_k(t-s)} f_k(s) ds.$$

Regarding the stochastic integral in (9.11), by Lemma 9.8, it admits the following representation

$$\sum_{k \in \mathbb{N}^d} \frac{b_k \sqrt{q_k}}{\sqrt{n}} \varphi_k \int_0^t e^{-\ell_k(t-s)} dW_k(s), \quad (9.12)$$

where $W_k, k \in \mathbb{N}^d$, are independent standard real-valued Wiener processes.

³It means that the notion ‘weak’ is in the PDE sense.

In summary, (9.11) admits the series representation

$$X^{(n)}(t) = \sum_{k \in \mathbb{N}^d} X_k^{(n)}(t) \varphi_k \quad (9.13)$$

with

$$X_k^{(n)}(t) = e^{-t\ell_k} u_k + \int_0^t e^{-(t-s)\ell_k} f_k(s) ds + \frac{b_k \sqrt{q_k}}{\sqrt{n}} \int_0^t e^{-(t-s)\ell_k} dW_k(s),$$

where $u_k \in \mathbb{R}$ and $f_k \in L^2[0, T]$ are from

$$u = \sum_{k \in \mathbb{N}^d} u_k \varphi_k \quad \text{and} \quad f(t) = \sum_{k \in \mathbb{N}^d} \varphi_k f_k(t)$$

respectively, and W_k are independent Wiener process on $[0, T]$.

9.6 Notes

Assumption 9.14 imposes constraints on the structure of the stochastic integral, which is necessary to obtain solutions as H -valued processes. Consider the following example. Recall that the (negative) Laplacian $-\Delta$ in \mathbb{R}^d has eigenvalues λ_k of the order $|k|^2$. As a consequence, a stochastic heat equation with space-time white noise, i.e. $U = H$ and $B = Q = I$, does not have a H -valued weak solution when the underlying domain is not one-dimensional. In general, in order to obtain a regular (H -valued) solution, the regularity property of the composition of $S(t)BQ^{1/2}$ needs to compensate the deterioration with growth of dimensions.

On the other hand, the requirement on B and Q to obtain well defined statistics is not as restrictive as to obtain unique solutions. While the solution to (9.6) only exists in a larger space than H , meaningful statistics for the terms $S(t)u$ and $\int_0^t S(t-s)f(s) ds$ of interest still can be obtained, see [51] for the detail.

In the end, the Q -Wiener process on U might be considered superfluous, as the stochastic integral can always be treated as $\int_0^t S(t-s) d\widetilde{W}(s)$, where \widetilde{W} is a $\widetilde{Q} = BQB^*$ -Wiener process on H .

Chapter 10

Bayesian Inference for Linear Evolution Equations

To prepare the main focus of this chapter, the inference for linear evolution equations, we first recall some key results from Chapter 9. In this chapter, we always assume the following condition, in which we choose a concrete but widely applicable example for H .

Condition 10.1. With $H = L^2(\mathfrak{D})$, Assumption 9.12 and Assumption 9.14 are satisfied.

A stochastic linear evolution equation is given by

$$\begin{cases} dX(t) + \mathcal{L}X(t) dt = f(t) dt + \frac{1}{\sqrt{n}} B dW^{\mathcal{Q}}(t) \\ X(0) = u \in H \end{cases}, \quad (10.1)$$

where u is the initial condition and $f : [0, T] \rightarrow H$ is the drift. Under Condition 10.1, the *mild* solution of (10.1) exists and is given by

$$X(t) = S(t)u + \int_0^t S(t-s)f(s) ds + \frac{1}{\sqrt{n}} \int_0^t S(t-s)B dW^{\mathcal{Q}}(s). \quad (10.2)$$

Furthermore, since the functions u and f in (10.1) admit the following representations,

$$u(x) = \sum_{k \in \mathbb{N}^d} u_k \varphi_k(x), \quad \text{and} \quad f(x, t) = \sum_{k \in \mathbb{N}^d} f_k(t) \varphi_k(x),$$

where $\{\varphi_k\}_{k \in \mathbb{N}^d}$ is the eigenbasis of \mathcal{L} . The solution (10.2) admits a series representation

$$X^{(n)}(t) = \sum_{k \in \mathbb{N}^d} X_k^{(n)}(t) \varphi_k, \quad (10.3)$$

whose coefficients are real-valued processes

$$X_k^{(n)}(t) = e^{-t\ell_k} u_k + \int_0^t e^{-(t-s)\ell_k} f_k(s) ds + \frac{b_k \sqrt{q_k}}{\sqrt{n}} \int_0^t e^{-(t-s)\ell_k} dW_k(s), \quad (10.4)$$

where $W_k(t)$ are independent standard Wiener processes, and the other constants are from the aforementioned assumptions.

In this chapter, we investigate the Bayesian approach to the recovery of the parameters u and f in (10.1). In Section 10.1, we study the recovery of initial condition u and in Section 10.2 we investigate the inference for drift f . In each section, we start with introducing a Gaussian prior that is tailored to the problem. Then, the contraction rates are proved using the general framework developed in Chapter 4. It is worthwhile to mention that our proofs do not rely on the conjugacy nor other Gaussian properties of the prior. The contraction rates for other priors can also be obtained using the same argument, namely verifying the conditions in Theorem 4.10.

10.1 Recovery of the Initial condition

In this section, suppose that all other parameters except the initial condition u are known. With no loss of generality, we can assume that $f(t) = 0$, for all $t \in [0, T]$.

Condition 10.2 (Observation of Final Value). Fix $T > 0$. For $n \in \mathbb{N}$, we observe the solution $X^{(n)}(T)$ of (10.1) at time T , i.e.

$$X^{(n)}(T) = S(T)u + \frac{1}{\sqrt{n}} \int_0^T S(T-s)B dW^{\mathcal{Q}}(s). \quad (10.5)$$

10.1.1 Spatial Gaussian Priors

Since the operator \mathcal{L} governs the spatial status of the evolution system, it is natural to consider a smoothness class that adapts to the structure of $\mathcal{L} = \Lambda^{(\nu)}$.

Centred Gaussian distributions on a separable Hilbert space correspond bijectively to covariance operators. By definition a random variable F with values in H is Gaussian if $\langle F, g \rangle$ is normally distributed, for every $g \in H$, and it has zero mean if these variables have zero means. The variances of these variables can then be written as

$$\mathbb{E}\langle F, g \rangle^2 = \langle Cg, g \rangle,$$

for a linear operator $C : H \rightarrow H$, called the *covariance operator*. A covariance operator C is necessarily self-adjoint, nonnegative, and of *trace class*, i.e., $\sum_{k \in \mathbb{N}^d} \langle C\varphi_k, \varphi_k \rangle < \infty$, for some (and then every) orthonormal basis $(\varphi_k)_{k \in \mathbb{N}^d}$ of H ; and every operator with these properties generates a Gaussian distribution.

Since the spatial regularity is characterized by isotropic Sobolev spaces H_s , we introduce the following Gaussian priors, which are fully adapted to the spatial smoothness.

Given a multi-index $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{R}_+^d$, a spatial Gaussian prior is the law of

$$u = \sum_{k \in \mathbb{N}^d} u_k \varphi_k, \quad \text{with } u_k \stackrel{\text{independent}}{\sim} \mathcal{N}_{\mathbb{R}}(0, |k^\alpha|^{-2}). \quad (10.6)$$

From Lemma 2.18, we conclude that $\alpha > \frac{d}{2}$ guarantees a ‘proper’ Gaussian prior on H .

10.1.2 Contraction Rate for Initial Condition

The observation (10.5) given in Condition 10.2 can be rewritten as a general linear problem,

$$X^{(n)} = \mathcal{A}u + \frac{\xi}{\sqrt{n}},$$

where

$$\mathcal{A} = S(T), \quad \xi = \int_0^T S(T-s)B dW^{\mathcal{Q}}(s),$$

and ξ is a proper Gaussian random element in H with the covariance operator

$$\Sigma = \int_0^T S(T-r)BQB^*S(T-r) dr = \sum_{k \in \mathbb{N}^d} \sigma_k^2 \langle \varphi_k, \cdot \rangle \varphi_k,$$

with

$$\sigma_k^2 := \frac{b_k^2 q_k}{2\ell_k} (1 - e^{-2\ell_k T}) \simeq \frac{b_k^2 q_k}{2\ell_k}. \quad (10.7)$$

Because of Assumption 9.14, we have

$$\|u\|_{\mathbb{H}_\xi}^2 = \sum_{k \in \mathbb{N}^d} \frac{u_k^2}{\sigma_k^2} \simeq \sum_{k \in \mathbb{N}^d} |k|^{2\mu} u_k^2, \quad (10.8)$$

which implies $\mathbb{H}_\xi = H_\mu$.

Notice that due to (9.7) and Assumption 9.12, the operator $\mathcal{A} = S(T)$ possesses a smoothing property and the recovery of u from (10.5) is in fact an inverse problem. By applying a general contraction result (Theorem 4.10) for inverse problems, modified from Theorem 3.1 in [43], the contraction rate of the spatial Gaussian prior from Section 10.1.1 for the recovery of initial condition u from (10.5) is obtained in the following theorem.

Theorem 10.3 (Gaussian Prior for the Initial Condition Recovery).

With $\mathbf{s} = (s, \dots, s) \in \mathbb{R}^d$, $s > 0$, let $\{H_{\mathbf{s}}\}_{\mathbf{s}}$ be the isotropic smoothness class introduced in Section 2.3.1 with the orthonormal basis $\{\varphi_k\}_{k \in \mathbb{N}^d}$ from Assumption 9.12 and with $\lambda_{k_i} = k_i$, $i = 1, \dots, d$. Consider the prior given in Section 10.1.1 with $\alpha > (\nu + d)/2$, and $X^{(n)}$ be given in (10.5). For any $u_0 \in H_\beta$ with $\beta > 0$, the posterior distribution satisfies, for sufficiently large $M > 0$,

$$\Pi_n \left(u : \|u - u_0\|_{L^2} > M(\log n)^{-s} \mid X^{(n)} \right) \rightarrow 0$$

in $\mathbb{P}_{u_0}^{(n)}$, with

$$s = \frac{1}{\nu} \left[\left(\alpha - \frac{\nu + d}{2} \right) \wedge \beta \right].$$

The contraction rate is of logarithmic order, because of the exponential smoothing property of the semigroup $S(T)$. Similar phenomenon has also been discovered in the recovery of the initial condition in white noise, c.f. [60] and Chapter 7. A noteworthy observation is that the rate obtained in Theorem 10.3 is identical to the

rate for the white noise case, see e.g. Theorem 7.14. In other words, the ‘smoother’ noise in (10.5), which in contrary to the white noise realises as a proper Gaussian element, does not lead to a faster rate. This is because the extreme ill-posedness from ‘inverting’ $S(T)$ predominantly resolves the logarithmic rate, and any noise with RKHS H_μ with $\mu \in \mathbb{R}^+$ will not improve the order of the rate.

10.2 Recovery of the Drift

In this section, suppose that all other parameters except the drift

$$f : [0, T] \rightarrow H$$

are known. With no loss of generality, we can assume that $u = 0$.

Condition 10.4 (Indirect Observation of Drift Term). Fix $T > 0$. For $n \in \mathbb{N}$, we observe continuously the solution $X^{(n)}(t)$ of (10.1) for $0 \leq t \leq T$, i.e.

$$X^{(n)}(t) = \int_0^t S(t-s)f(s) ds + \frac{1}{\sqrt{n}} \int_0^t S(t-s)B dW^\mathcal{Q}(s). \quad (10.9)$$

10.2.1 Spatial-Temporal Gaussian Priors

For the recovery of drift terms, the priors necessarily need to sit in the function space $L^2([0, T]; H)$, where $H = L^2(\mathcal{D})$. Since $L^2([0, T]; H) \cong H \otimes L^2([0, T]; \mathbb{R}) \cong L^2(\mathcal{D}_T)$, one may introduce a Gaussian prior on the space $L^2(\mathcal{D}_T)$ following the same procedure in the previous paragraphs. However, as mentioned in Section 2.3.1, it may be of interest to distinguish the smoothness in each spatial and temporal directions.

We introduce zero Gaussian priors on $L^2([0, T]; H)$ using series expansion. In order to do that, we fix an orthonormal basis $\{\psi_k\}_k$ of $L^2([0, T])$.

From Section 2.3.1, recall that given the orthonormal basis $\{\psi_l\}_{l \in \mathbb{N}}$ of $L^2([0, T]; \mathbb{R})$,

$$\{\tilde{\varphi}_{k,l}\}_{(i,j) \in \mathbb{N}^d \times \mathbb{N}} = \{\varphi_k \otimes \psi_l\}_{(i,j) \in \mathbb{N}^d \times \mathbb{N}}$$

is an orthonormal basis of $L^2([0, T]; H)$, of which any function $f(x, t)$ admits the representation

$$\sum_{(k,l) \in \mathbb{N}^d \times \mathbb{N}} f_{k,l} \varphi_k(x) \psi_l(t) \quad \text{with} \quad \|\{f_k\}\|_{\ell^2(\mathbb{N}^{d+1})} = \sum_{k \in \mathbb{N}^{d+1}} f_k^2 < \infty.$$

Given a multi-index $\alpha = (\alpha_1, \dots, \alpha_{d+1}) \in \mathbb{R}_+^{d+1}$ and $\beta_* = (p, \dots, p, 0) \in \mathbb{R}_+^{d+1}$, the spatial-temporal Gaussian prior is the law of

$$f(x, t) = \sum_{(k,l) \in \mathbb{N}^d \times \mathbb{N}} f_{k,l} \tilde{\varphi}_{k,l}(x, t) = \sum_{(k,l) \in \mathbb{N}^d \times \mathbb{N}} f_{k,l} \varphi_k(x) \psi_l(t), \quad (10.10)$$

with

$$f_k \stackrel{\text{independent}}{\sim} \mathcal{N}_{\mathbb{R}}(0, |k^{\beta_*}|^{-2} |k^\alpha|^{-2}),$$

i.e. f_k are independent zero mean Gaussian random variables with variances

$$|k^{\beta_*}|^{-2}|k^\alpha|^{-2} = \left(1 + \sum_{i \leq d} k_i^{2p}\right)^{-2} \left(\sum_{i \leq d+1} k_i^{2\alpha_i}\right)^{-2}.$$

By Lemma 2.18, when $\mathcal{H}(\alpha) > (d+1)/2$, the prior has sample paths that are H_{β_*} -valued almost surely.

We conclude this section with the following remark. Using the basis constructed above, the function f_k of (10.3) can also be expressed as

$$f_k(t) = \sum_{l \in \mathbb{N}} f_{k,l} \psi_l(t).$$

10.2.2 Contraction Rate for Drift Recovery

In this section we study the performance of Bayesian methods in the recovery of the drift term $f \in L^2([0, T]; H)$. As shown in Condition 10.4, i.e. (10.9)

$$X^{(n)}(t) = \int_0^t S(t-s)f(s) ds + \frac{1}{\sqrt{n}} \int_0^t S(t-s)B dW^{\mathcal{Q}}(s),$$

the noise is a vector-valued Gaussian process, whose RKHS is determined by the operator-valued Kernel $S(t-s)$. This imposes a challenging technicality, as some analytical tools such as the operator version of Mercer theorem is required in order to obtain a workable structure of RKHS. However, a unique characteristic of the observation under discussion is that the same (operator-valued) integral kernel $k(t, s) = S(t-s)$ is applied to both the drift and the noise. This property offers us a workaround to avoid the aforementioned difficulty: whitening the process. To be specific, by a proper transform of the signal, we will show that the observation (10.9) along its spatial basis is statistically equivalent to a sequence version of the white noise model (Section 10.3.2.1), the latter of which can be further related to a Gaussian $(d+1)$ -dimensional sequence model (Section 10.3.2.2). As a consequence, the problem is reduced to standard nonparametric estimation without inverse nature, which is a multi-dimensional problem because the underlying space-time domain is a compact set in \mathbb{R}^{d+1} .

Now we show the contraction rates of the Gaussian prior Section 10.2.1 in the recovery of a Drift term.

Theorem 10.5 (Gaussian Prior for the Drift Recovery).

Let $\beta_* = (p, \dots, p, 0) \in \mathbb{R}_+^{d+1}$. For any $f_0 \in H_{\beta}$ with $\beta > \beta_*$, where H_{β} is an anisotropic smoothness class defined in Section 2.3.1 with the orthonormal basis $\{\varphi_k \otimes \psi_l\}_{(i,j) \in \mathbb{N}^d \times \mathbb{N}}$ such that $\{\varphi_k\}_{k \in \mathbb{N}^d}$ from Assumption 9.12, an orthonormal basis $\{\psi_l\}_{l \in \mathbb{N}}$ of $L^2([0, T])$, and $\lambda_{k_i} = k_i$, $i = 1, \dots, d+1$. Let the prior be zero-mean spatial-temporal Gaussian proposed in Section 10.2.1 with $\mathcal{H}(\alpha) > (d+1)/2$, and $X^{(n)}$ be the observations in the form of (10.9). The posterior distribution satisfies, for sufficiently large $M > 0$,

$$\Pi_n \left(f : \|f - f_0\|_{H_{\beta_*}} > Mn^{-s} \mid X^{(n)} \right) \xrightarrow{P_{f_0}^{(n)}} 0,$$

where

$$s = \left(\frac{\mathcal{H}(\alpha) - (d+1)/2}{2\mathcal{H}(\alpha)} \right) \wedge \left(\frac{1}{2 + 2 \sup_{i \leq d+1} \frac{(\alpha_i - \beta_i) \vee 0}{\beta_i}} \right).$$

In particular, when $\alpha_i = \frac{2\mathcal{H}(\beta) + d + 1}{2\mathcal{H}(\beta)} \beta_i$ for each $1 \leq i \leq d+1$, the two items in the expression of s above are balanced and

$$s = \frac{\mathcal{H}(\beta)}{2\mathcal{H}(\beta) + d + 1}.$$

Remark 10.6. The contraction rate is given in the norm of smoothness class H_{β_*} . With a proper choice of the basis functions, such as Fourier basis, the space H_{β_*} can be connected to certain type of multidimensional Sobolev spaces.

10.3 Proofs

10.3.1 Proofs in Section 10.1

The theorem is a corollary to Theorem 4.10. The main tasks are to determine ε_n satisfying the prior mass condition (4.12) of the direct problem, and next to identify η_n from the prior mass condition (4.13) and the other conditions.

The first task is achieved in the following lemma.

Lemma 10.7. For $f_0 \in H_s$, the prior Π from Section 10.1.1, as $\varepsilon \downarrow 0$,

$$-\log \Pi(f : \|\mathcal{A}f - \mathcal{A}f_0\|_{\mathbb{H}_\varepsilon} < \varepsilon) \lesssim \left(\log \frac{1}{\varepsilon} \right)^r, \quad (10.11)$$

where

$$r = \left(2 \frac{\alpha - \beta}{\nu} \right) \vee \left(\frac{\nu + d}{\nu} \right).$$

Remark 10.8. Since $\nu \in \mathbb{N}$ from Assumption 9.12 and $d \in \mathbb{N}$, $r > 1$.

Proof. The probability in the left side is the decentred small ball probability $\Pi(g : \|g - g_0\|_{\mathbb{H}_\varepsilon} < a\varepsilon)$ of the Gaussian random variable $G = \mathcal{A}F$ distributed according to the prior under the linear transform \mathcal{A} . Symbolically we denote the covariance operator (which is diagonal) of F by Λ_F . Due to the property of Gaussian measure, the random element G is also a centred Gaussian random element in H with the covariance operator

$$\mathcal{A}\Lambda_F\mathcal{A}^* : h = \sum_{k \in \mathbb{N}^d} h_k \varphi_k \mapsto \sum_{k \in \mathbb{N}^d} e^{-2T|k^\nu|} |k^\alpha|^{-2} h_k \varphi_k. \quad (10.12)$$

The RKHS \mathbb{H}_G of G is given by

$$\left\{ g = \sum_{k \in \mathbb{N}} g_k \varphi_k \in H : \|g\| = \sum_{k \in \mathbb{N}^d} e^{2T|k^\nu|} |k^\alpha|^2 g_k^2 < \infty \right\}. \quad (10.13)$$

It is convenient to work with the following norm of \mathbb{H}_G ,

$$\|g\|_{\mathbb{H}_G}^2 := \sum_{k \in \mathbb{N}^d} e^{2T|k|^\nu} |k|^{2\alpha} g_k^2,$$

which is equivalent to the norm in (10.13).

Recall (10.8), we have $\mathbb{H}_\xi = H_\mu$ (as sets) and $\|h\|_{\mathbb{H}_\xi} \simeq \|h\|_{H_\mu}$. Due to the exponential smoothing property of \mathcal{A} , for any $f \in H$ and $\mathbf{s} \in \mathbb{R}_+^d$, we have $\mathcal{A}f \in H_{\mathbf{s}}$. Hence $\Pr(G \in \mathbb{H}_\xi) = 1$. Hence the distribution of G can be considered as a Gaussian measure on \mathbb{H}_ξ with RKHS \mathbb{H}_G .

The left side of (10.11) is therefore up to constants equivalent to

$$\inf_{g \in \mathbb{H}_G: \|g - g_0\|_{\mathbb{H}_\xi} < \varepsilon} \|g\|_{\mathbb{H}_G}^2 - \log \Pi(\|g\|_{\mathbb{H}_\xi} < \varepsilon). \quad (10.14)$$

See [64, 65, 99], or Section 11.2, in particular, Proposition 11.19 in [35].

Let P_j be the H_0 -orthonormal projection to the j basis $\{\varphi_k\}_{|k|_\infty \leq j^{1/d}}$. Since \mathcal{A} and P_j commute, using Remark 9.13, we have

$$\begin{aligned} & \|P_j \mathcal{A}f_0 - \mathcal{A}f_0\|_{\mathbb{H}_\xi}^2 \\ &= \sum_{|k|_\infty \geq j^{1/d}} e^{-2T\ell_k} |k|^{2\mu} f_{0,k}^2 \simeq_d \sum_{|k|_\infty \geq j^{1/d}} e^{-2T|k|^\nu} |k|^{2(\mu-\beta)} (|k|^{2\beta} f_{0,k}^2) \\ &\lesssim_d \exp\left(-2Tj^{\nu/d}\right) j^{-2(\beta-\mu)/d} \|f_0\|_\beta^2, \end{aligned} \quad (10.15)$$

and hence $\|P_j \mathcal{A}f_0 - \mathcal{A}f_0\|_{\mathbb{H}_\xi}$ is bounded above by ε for $j \simeq_T (-\log \varepsilon)^{d/\nu}$. By substituting this value of j into

$$\begin{aligned} \|P_j \mathcal{A}f_0\|_{\mathbb{H}_G}^2 &= \sum_{|k|_\infty \leq j^{1/d}} |k|^{2\alpha} f_{0,k}^2 \lesssim_d \sum_{|k|_\infty \leq j^{1/d}} |k|^{2\alpha-2\beta} |k|^{2\beta} f_{0,k}^2 \\ &\leq j^{2\frac{(\alpha-\beta)\vee 0}{d}} \|f_0\|_\beta^2, \end{aligned}$$

we conclude that the first term in (10.14) is bounded above by $(-\log \varepsilon)^{2\frac{\alpha-\beta}{\nu}\vee 0}$.

For the second term in (10.14), by Corollary 10.15, the metric entropy

$$\log N(\varepsilon, \{g \in \mathbb{H}_G : \|g\|_{\mathbb{H}_G} \leq 1\}, \|\cdot\|_{\mathbb{H}_\xi}) \simeq (-\log \varepsilon)^{(\nu+d)/\nu}.$$

Hence, by [64] (see Lemma 6.2 in [100]),

$$-\log \Pi(\|\mathcal{A}f\|_{\mathbb{H}_\xi} < \varepsilon) \simeq \left(\log \frac{1}{\varepsilon}\right)^{(\nu+d)/\nu}.$$

Finally, the assertion of the lemma follows from combining the above results. \square

It follows that (4.12) is satisfied for any ε_n such that

$$e^{-(\log \frac{1}{\varepsilon_n})^r} \geq e^{-n\varepsilon_n^2},$$

where $r > 1$ is given in the lemma above. Let $x = -\log \varepsilon_n$. The preceding display can be rewritten into

$$\frac{2}{r} x e^{\frac{2}{r}x} = x^* e^{x^*} \leq \frac{2}{r} n^{1/r}.$$

The following lemma is useful for the proof.

Lemma 10.9. *Let $W(x)$ be the Lambert W function, i.e. the inversion of the mapping $[e^{-1}, \infty) \ni x \mapsto xe^x$. We have, when $x \rightarrow \infty$,*

$$W(x) \sim \log x - \log \log x.$$

Proof. From the identity $x = W(x)e^{W(x)}$, $W(x)$ is an increasing function with respect to x and $W(e) = 1$. In addition, we also have the following identities,

$$W(x) = \log\left(\frac{x}{W(x)}\right) \quad \text{and} \quad W(x) + \log W(x) = \log x.$$

From now on only consider the case $x > e$. The second relation in the last display implies that $W(x) < \log x$. Hence,

$$W(x) = \log\left(\frac{x}{W(x)}\right) > \log\left(\frac{x}{\log x}\right) = \log x - \log \log x.$$

On the other hand,

$$W(x) = \log\left(\frac{x}{\log\left(\frac{x}{W(x)}\right)}\right) < \log\left(\frac{x}{\log\left(\frac{x}{\log x}\right)}\right) = \log x - \log \log x + \log \log \log x.$$

Since $\log \log x \gg \log \log \log x$ as $x \rightarrow \infty$, the proof is complete. \square

The asymptotic expansion of the Lambert W function implies that x such that

$$x = \frac{r}{2}W\left(\frac{2}{r}n^{1/r}\right) \sim \log \frac{\sqrt{n}}{(\log n)^{r/2}}$$

satisfies the last inequality above. Consequently, (4.12) is satisfied with

$$\varepsilon_n \simeq \frac{(\log n)^{r/2}}{\sqrt{n}}, \quad (10.16)$$

Now we construct the reconstruction operator $\mathcal{R}_n : H \rightarrow H$. For $g \in H$, we consider the following truncation regularizer,

$$\mathcal{R}_n g := \sum_{|k|_\infty < j_n^{1/d}} e^{T\ell_k} g_k \varphi_k, \quad (10.17)$$

where $j_n \rightarrow \infty$ as $n \rightarrow \infty$. Consequently, from Assumption 9.12, with some positive constant c , we have

$$\|\mathcal{R}_n\| = \sup_{|k|_\infty < j_n^{1/d}} \exp(|k|^\nu |T|) = e^{c j_n^{\nu/d} T}, \quad (10.18)$$

and (4.6) is satisfied with $\rho_n = e^{c j_n^{\nu/d} T}$.

Besides, for $u_0 = \sum_{k \in \mathbb{N}^d} u_{0,k} \varphi_k \in H_\beta$, by Lemma 2.20, we have

$$\|\mathcal{R}_n \mathcal{A} u_0 - u_0\|^2 = \sum_{|k|_\infty \geq j_n^{1/d}} u_{0,k}^2 \lesssim_d j_n^{-2\beta/d} \|u_0\|_{H_\beta}^2. \quad (10.19)$$

The next step of the proof is to bound the prior probability in (4.13).

Lemma 10.10. *Let \mathcal{R}_n be given as (10.17) and the corresponding $j_n = j$. There exist $a, b > 0$, such that for every $j \in \mathbb{N}$ and $t > 0$,*

$$\Pi(f : \|\mathcal{R}_n \mathcal{A} f - f\|_0 > t + a j^{1/2 - \alpha/d}) \leq e^{-bt^2 j^{2\alpha/d}}.$$

Proof. Let $f^{(n)} = \mathcal{R}_n \mathcal{A} f$. Therefore, the probability on the left concerns the random variable $(\mathcal{R}_n \mathcal{A} - I)F$, if F is a variable distributed according to the prior Π . Since F is zero-mean normal with a covariance operator symbolically denoted by Λ_F , this variable is zero-mean Gaussian with covariance operator $(\mathcal{R}_n \mathcal{A} - I)\Lambda_F(\mathcal{R}_n \mathcal{A} - I)^*$. We shall compute the weak and strong second moments of the variable $(\mathcal{R}_n \mathcal{A} - I)F$, and next apply Borell's inequality for the norm of a Gaussian variable to obtain the exponential bound.

Since

$$(\mathcal{R}_n \mathcal{A} - I)F = - \sum_{|k|_\infty \geq j^{1/d}} F_k \varphi_k$$

with $F_k \stackrel{i.i.d.}{\sim} \mathcal{N}_{\mathbb{R}}(0, |k^\alpha|^{-2})$, we have

$$\langle (\mathcal{R}_n \mathcal{A} - I)F, g \rangle = - \sum_{|k|_\infty \geq j^{1/d}} F_k g_k,$$

for arbitrary $g = \sum_{k \in \mathbb{N}^d} g_k \varphi_k$. Then, the weak second moment of $(\mathcal{R}_n \mathcal{A} - I)F$ is given by

$$\sup_{\|f\|_0 \leq 1} \mathbb{E} \langle (\mathcal{R}_n \mathcal{A} - I)F, f \rangle^2 \leq \sup_{\sum_k f_k^2 \leq 1} \sum_{|k|_\infty > j^{1/d}} |k^\alpha|^{-2} f_k^2 \simeq j^{-2\alpha/d}.$$

The strong second moment of the Gaussian variable $(\mathcal{R}_n \mathcal{A} - I)F$ is

$$\mathbb{E} \|(\mathcal{R}_n \mathcal{A} - I)F\|^2 = \sum_{|k|_\infty \geq j^{1/d}} |k|^{-2\alpha}.$$

In the proof of Lemma 2.18, we have shown that for the hypercubes

$$C_n = \{k \in \mathbb{N}^d : k_i \lesssim n^{1/d}, i = 1, \dots, d\},$$

the increment $C_n \setminus C_{n-1}$ covers index points of the order n^{d-1} . Hence, the strong second moment can be bounded by

$$\sum_{|k|_\infty \geq j^{1/d}} |k|^{-2\alpha} \leq_d \sum_{l \geq j^{1/d}} \sum_{k \in [C_l \setminus C_{l-1}]} \prod_{i \leq d} k_i^{-2\alpha/d} \simeq \sum_{l \geq j^{1/d}} l^{d-1-2\alpha} \leq j^{(d-2\alpha)/d},$$

where we used the estimate $\sum_{i > j} i^{-b} \leq j^{1-b}/(b-1)$, for $b > 1$.

Since the first moment of $\|(\mathcal{R}_n \mathcal{A} - I)F\|_0$ is bounded by the root of its second moment, the lemma follows by Borell's inequality (see e.g. Lemma 3.1 and subsequent discussion in [67]). \square

Let $t^2 = 4n\varepsilon_n^2/(bj_n^{2\alpha/d})$. Then, $t \gtrsim j_n^{1/2-\alpha/d}$, due to the constraint (4.8), i.e. $j_n \lesssim n\varepsilon_n^2$. Substituting t into Lemma 10.10,

$$\Pi(f : \|\mathcal{R}_n \mathcal{A}f - f\|_0^2 > 4n\varepsilon_n^2/(bj_n^{2\alpha/d})) \leq e^{-4n\varepsilon_n^2},$$

which together with (10.16) implies

$$\eta_n \gtrsim j_n^{-\alpha/d}(\log n)^{r/2}.$$

In addition, the constraints eqs. (4.9) and (4.10) impose

$$\begin{aligned} \eta_n &\gtrsim e^{Tcj_n^{\nu/d}} \frac{(\log n)^{r/2}}{\sqrt{n}}, \\ \eta_n &\gtrsim j_n^{-\beta/d}, \end{aligned}$$

where the right-hand side of the inequalities are given by (10.18) and (10.19).

We need to determine j_n in order to solve for η_n . Since $d/\nu < r$, $(\log n)^{d/\nu} \ll n\varepsilon_n^2$. Hence, we can choose $j_n = \tilde{c}^{d/\nu}(\log n)^{d/\nu} \lesssim n\varepsilon_n^2$, with \tilde{c} such that $\tilde{c}Tc < 1/2$. Substituting j_n into the preceding constraints leads to

$$\begin{aligned} \eta_n &\gtrsim (\log n)^{-\frac{2\alpha-\nu-d}{2\nu}}, \\ \eta_n &\gtrsim n^{-(\frac{1}{2}-\tilde{c}Tc)}(\log n)^{r/2}, \\ \eta_n &\gtrsim (\log n)^{-\beta/\nu}. \end{aligned}$$

The second inequality above is negligible compared to the other two. The theorem follows from Theorem 4.10.

10.3.2 Proofs in Section 10.2

The major step of the proof can be summarized as follows.

- (1) Section 10.3.2.1. Consider the sequence of scalar processes $\{X_k^{(n)}\}_{k \in \mathbb{N}^d}$ from (9.13). Using a sequence of transforms $\{\mathcal{T}_k\}_{\mathbb{N}^d}$, the sequence of processes $\{X_k^{(n)}\}_{k \in \mathbb{N}^d}$ is whitened in time.
- (2) Section 10.3.2.2. The signal f to recover is also isometrically transformed into \tilde{f} . The transformed observation can be expressed with a Gaussian sequence,

$$\tilde{X}_{k,l}^{(n)} = f_{k,l} + \frac{1}{\sqrt{n}} \tilde{\xi}_{k,l} \in \mathbb{R}, \quad (k, l) \in \mathbb{N}^d \times \mathbb{N},$$

where the covariance structure of $\tilde{\xi}_{k,l}$ is determined by the operators B and \mathcal{Q} .

- (3) Section 10.3.2.3. As the final preparation for the proof of Theorem 10.5, we establish a posterior contraction rate for the multi-dimensional white noise model.

- (4) Section 10.3.2.4 Using the result obtained in Section 10.3.2.3, and an isometric property possessed by the prior and the noise $\tilde{\xi}$, we finally conclude the proof of Theorem 10.5.

Recall that the eigenbasis $\{\varphi_k\}_{k \in \mathbb{N}^d}$ of \mathcal{L} is an orthonormal basis in space and $\{\psi_k\}_{k \in \mathbb{N}}$ is an orthonormal basis in time, and denote their tensor product by $\{\tilde{\varphi}_{k,l} = \varphi_k \otimes \psi_l\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$.

10.3.2.1 Whitening Ornstein-Uhlenbeck processes

Due to (9.13), the observation (10.9) is equivalent to the following (functional) sequence model, for $k \in \mathbb{N}^d$, we observe

$$X_k^{(n)}(t) = \int_0^t e^{-\ell_k(t-s)} f_k(s) ds + \frac{b_k \sqrt{q_k}}{\sqrt{n}} \int_0^t e^{-\ell_k(t-s)} dW_k(s), \quad (10.20)$$

in the product space $(L^2([0, T]; \mathbb{R}))^{\mathbb{N}^d}$ with the product measure $\bigotimes_{k \in \mathbb{N}^d} (\mu_{W_k})$, where μ_{W_k} are the probability measures induced by the processes W_k given in (9.12), which are mutually independent Wiener processes.

Notice that the real-valued processes $X_k^{(n)}(t)$ is Ornstein-Uhlenbeck processes. We are going to convert them into the standard white noise model, and start with introducing a useful function together with its inverse. For $\lambda > 0$ and $t \in [0, T]$, define

$$\vartheta(t) = \frac{\log(2\lambda t + 1)}{2\lambda} \quad \text{and} \quad \vartheta^{-1}(t) = \frac{e^{\lambda t} - 1}{2\lambda}, \quad (10.21)$$

where ϑ^{-1} is well-defined since $\vartheta : [0, T) \rightarrow \mathbb{R}^+$ is bijective. With function ϑ , we can define the following transform

$$(\mathcal{T}g)(t) := \sqrt{2\lambda t + 1} (g \circ \vartheta)(t), \quad t \in [0, \vartheta(T)], \quad (10.22)$$

for any continuous function g on $[0, T]$.

Using the newly defined transform \mathcal{T} , the noise can be whitened as follows.

Lemma 10.11. *Let $W(t)$ be a Brownian motion, i.e. a standard real-valued Wiener process, and $\vartheta(t)$ be given in (10.21). If*

$$\xi(t) = \int_0^t e^{-\lambda(t-s)} dW(s),$$

then $(\mathcal{T}\xi)(t) = \sqrt{2\lambda t + 1} (\xi \circ \vartheta)(t)$ is a Brownian motion.

Proof. The process $M(t) = e^{\lambda t} \xi(t)$ is a continuous martingale whose quadratic variation is

$$[M]_t = \int_0^t e^{2\lambda s} ds = \frac{e^{2\lambda t} - 1}{2\lambda} = \vartheta^{-1}(t).$$

Consequently, $M \circ \vartheta(t)$ is a continuous martingale with

$$[M \circ \vartheta]_t = [M]_{\vartheta(t)} = t.$$

Therefore,

$$M \circ \vartheta(t) = e^{\lambda\vartheta(t)} \xi \circ \vartheta(t) = \sqrt{2\lambda t + 1} (\xi \circ \vartheta)(t)$$

is a Brownian motion. \square

Similarly, the transform (10.22) can be applied to the deterministic integral.

Lemma 10.12. *Assume $f \in L^2[0, T]$. Let*

$$F(t) = \int_0^t e^{-\lambda(t-s)} f(s) ds$$

and $\tilde{f}(u)$ be the transform of function f such that

$$\tilde{f}(u) = \frac{f \circ \vartheta(u)}{\sqrt{2\lambda u + 1}}. \quad (10.23)$$

Then, the following statements hold.

(i) $(\mathcal{T}F)(t) = \sqrt{2\lambda t + 1} (F \circ \vartheta)(t) = \int_0^t \tilde{f}(u) du.$

(ii) If $\{\psi_k\}_k$ is an orthonormal basis for $L^2[0, T]$, then,

$$\int_0^{\vartheta(T)} \tilde{\psi}_k \tilde{\psi}_l du = \delta_{kl} \quad \text{and} \quad \int_0^{\vartheta(T)} \tilde{f} \tilde{\psi}_k du = \int_0^T f \psi_k du.$$

Proof. Since $f \in L^2$, F is continuous. The first statement follows from

$$\begin{aligned} (\mathcal{T}F)(t) &= \sqrt{2\lambda t + 1} \int_0^{\vartheta(t)} e^{-\lambda(\vartheta t - s)} f(s) ds \\ &\stackrel{s=\vartheta(u)}{=} \sqrt{2\lambda t + 1} e^{-\lambda\vartheta(t)} \int_0^t e^{\lambda\vartheta(u)} f \circ \vartheta(u) \vartheta'(u) du \\ &= \int_0^t f \circ \vartheta(u) \frac{1}{\lambda} (e^{\lambda\vartheta(u)})' du = \int_0^t \frac{f \circ \vartheta(u)}{\sqrt{2\lambda u + 1}} du, \end{aligned}$$

where $(\cdot)'$ denotes the ordinary derivative.

The next statement is obtained by changing variables. The first equation follows from

$$\int_0^{\vartheta(T)} \tilde{\psi}_k \tilde{\psi}_l du = \int_0^{\vartheta(T)} \frac{\psi_k(\vartheta(u)) \psi_l(\vartheta(u))}{2\lambda u + 1} du \stackrel{s=\vartheta(u)}{=} \int_0^T \psi_k(s) \psi_l(s) ds = \delta_{kl}.$$

The same argument also applies to the second one. \square

Applying the transform \mathcal{T} defined in (10.22) to

$$X(t) = \int_0^t e^{-\lambda(t-s)} f(s) ds + c \int_0^t e^{-\lambda(t-s)} dW(s), \quad t \in [0, T],$$

we obtain

$$\tilde{X}(t) := \mathcal{T}X(t) = \int_0^t \tilde{f}(s) ds + c\tilde{W}(t), \quad t \in [0, \vartheta(T)],$$

where \tilde{f} is given in (10.23) and $\tilde{W}(t)$ is a Brownian motion.

Now given an orthonormal basis $\{\psi_l\}_{l \in \mathbb{N}}$ of $L^2[0, T]$, we can form, for $l \in \mathbb{N}$,

$$\int_0^{\vartheta(T)} \tilde{\psi}_l d\tilde{X}(s) = \int_0^T f(s)\psi_l(s) ds + c \int_0^{\vartheta(T)} \tilde{\psi}_l d\tilde{W} = f_l + cz_l, \quad (10.24)$$

where z_l are i.i.d. standard Gaussian.

10.3.2.2 Complete Sequence Model

Now consider the independent signals $X_k^{(n)}$ as given in (10.20). Define \mathcal{T}_k as the transform (10.22) from the previous section, with $\lambda = \ell_k$ and $c_k = (\sqrt{q_k}b_k)/\sqrt{n}$. Then, we can transform the signals into

$$\tilde{X}_k^{(n)}(t) := \mathcal{T}_k X_k^{(n)}(t) = \int_0^t \tilde{f}_k(s) ds + c_k \tilde{W}_k(t), \quad (10.25)$$

where $\tilde{W}_k(t)$ are independent Brownian motions. Because of (10.24), we can form observations

$$\tilde{X}_{k,l}^{(n)} = f_{k,l} + \frac{b_k \sqrt{q_k}}{\sqrt{n}} z_{k,l} \in \mathbb{R}, \quad (k, l) \in \mathbb{N}^d \times \mathbb{N},$$

where $z_{k,l} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1)$. In fact, these are the observations of the coordinates of the following multidimensional Gaussian sequence model,

$$\tilde{X}^{(n)} = \tilde{f} + \frac{1}{\sqrt{n}} \tilde{\xi}, \quad (10.26)$$

where $\tilde{f} = \{f_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}} \in \ell^2(\mathbb{N}^d \times \mathbb{N})$ are the coefficients of f in the series representation with basis $\{\varphi_k \otimes \psi_l\}_{k \in \mathbb{N}^d, l \in \mathbb{N}}$ and $\tilde{\xi} = \{\tilde{\xi}_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$ is a random vector in $\mathbb{R}^{\mathbb{N}^{d+1}}$ whose entries are independent zero mean Gaussian random variables with variance $\left(\sum_{i \leq d} k_i^p\right)^{-2}$. This variance is determined by the decay of $b_k \sqrt{q_k}$ under Assumption 9.14.

10.3.2.3 Gaussian Posterior Contraction for Multi-dimensional White Noise Model

In this subsection, as the final preparation for proving Theorem 10.5, we prove the posterior contraction rate of the equivalent prior of Section 10.2.1 on sequence spaces, for a simple Gaussian sequence model. To be precise, we consider the following situation.

Let $\alpha \in \mathbb{R}_+^m$ be a multi-index. Consider a prior as the law of $F = \{F_k\}_{k \in \mathbb{N}^m}$, where F_k are independent centred real Gaussian random variables with variance $|k^\alpha|^{-2}$. We impose $\mathcal{H}(\alpha) > m/2$ so that the prior are almost surely realised in the space $\ell^2(\mathbb{N}^m)$.

Lemma 10.13. *Consider the observation is given by the multi-dimensional white noise model,*

$$X_k^{(n)} = f_k + \frac{1}{\sqrt{n}} z_k, \quad \text{for } k \in \mathbb{N}^m, \quad (10.27)$$

where z_k are independent standard Gaussian random variables. If true parameter $f_0 = \{f_{0,k}\}_{k \in \mathbb{N}^m} \in h_\beta$, where h_β is a Sobolev ellipsoid defined in Section 2.3.1 equipped with norm (2.17), then the Gaussian prior above behaves, as $\varepsilon \downarrow 0$,

$$-\log \Pi(f : \|f - f_0\|_{\ell^2} < \varepsilon) \lesssim \varepsilon^{-r_1} \vee \varepsilon^{-r_2} \quad (10.28)$$

with

$$r_1 = 2 \sup_{i \leq m} \frac{(\alpha_i - \beta_i) \vee 0}{\beta_i} \quad \text{and} \quad r_2 = \frac{m}{\mathcal{H}(\alpha) - m/2}.$$

Proof. Notice that the RKHS \mathbb{H}_F of the prior F is $h_\alpha \subset \ell^2$. The left side of (10.28) is up to constants equivalent to

$$\inf_{f \in h_\alpha : \|f - f_0\|_{\ell^2} < \varepsilon} \|f\|_{h_\alpha}^2 - \log \Pi(\|f\|_{\ell^2} < \varepsilon). \quad (10.29)$$

See [64, 65, 99], or Section 11.2, in particular, Proposition 11.19 in [35].

Let P_N be the truncation of a sequence to $\{k < N : k_i < N_i, i \leq m\}$ with $N = (N_1, \dots, N_m)$. Applying Lemma 2.20, we obtain $\|P_N f_0 - f_0\|_{\ell^2} \leq \varepsilon$, if $N_i \gtrsim \varepsilon^{-1/\beta_i}$ for all $i \leq m$.

Taking $N_i \simeq \varepsilon^{-1/\beta_i}$, $i \leq m$, an upper bound on the first term in (10.28) is obtained as,

$$\begin{aligned} \|P_N f_0\|_{\mathbb{H}_F}^2 &= \sum_{k < N} |k^\alpha|^2 f_{0,k}^2 \leq \sum_{k < N} \left[\sum_{i \leq m} k_i^{2\beta_i} k_i^{2(\alpha_i - \beta_i) \vee 0} \right] f_{0,k}^2 \\ &\lesssim d \sum_{k < N} |k^\beta|^2 \left[\sum_{i \leq m} k_i^{2(\alpha_i - \beta_i) \vee 0} \right] f_{0,k}^2 \leq \left[\sum_{i \leq m} N_i^{2(\alpha_i - \beta_i) \vee 0} \right] \|f_0\|_{h_\beta}^2 \\ &\simeq \sum_{i \leq m} \varepsilon^{-\frac{2(\alpha_i - \beta_i) \vee 0}{\beta_i}} \|f_0\|_{h_\beta}^2 \leq d \varepsilon^{-\sup_{i \leq m} \frac{2(\alpha_i - \beta_i) \vee 0}{\beta_i}} \|f_0\|_{h_\beta}^2 \end{aligned}$$

For the second term (small ball probability) in (10.28), by Corollary 2.24, the metric entropy $\log N(\varepsilon, \{f \in h_\beta : \|f\|_{h_\beta} \leq 1\}, \|\cdot\|_{\ell^2})$ is of the order $\varepsilon^{-m/\mathcal{H}(\alpha)}$. Hence, under the condition $\mathcal{H}(\alpha) > d/2$, by [64] (see Lemma 6.2 in [100]),

$$-\log \Pi(\|f\|_{\ell^2} < \varepsilon) \simeq \varepsilon^{-\frac{m}{\mathcal{H}(\alpha) - m/2}}.$$

□

According to equations (1.2) and (1.3) in [99], the minimal ε_n satisfying

$$-\log \Pi(f : \|f - f_0\|_{\ell^2} < \varepsilon_n) \lesssim n\varepsilon_n^2$$

is the posterior contraction rate of the Gaussian prior considered in this section. By direct calculation, it can be shown that when for all $1 \leq i \leq m$,

$$\alpha_i = \left(\frac{2\mathcal{H}(\beta) + m}{2\mathcal{H}(\beta)} \right) \beta_i,$$

the rates r_1 and r_2 in (10.28) are balanced to $r_1 = r_2 = m/\mathcal{H}(\beta)$ and the posterior contraction rate reaches the minimax rate (see [53]), i.e.

$$\varepsilon_n \simeq n^{-\frac{\mathcal{H}(\beta)}{2\mathcal{H}(\beta)+m}}.$$

10.3.2.4 Proof of Theorem 10.5

After the long preparation, the proof of Theorem 10.5 simply follows from assembling all the results obtained up to now.

Recall $\{\tilde{\varphi}_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}} = \{\varphi_k \otimes \psi_l\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$ is a fixed basis of $L^2(\mathfrak{D}_T)$, satisfying the assumptions in Section 9.5. For a function f in $L^2(\mathfrak{D}_T)$, its coefficients in the basis $\{\tilde{\varphi}_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$ are denoted by $\tilde{f} = \{f_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$. Recall that for the norms $\|\cdot\|_{L^2}$, $\|\cdot\|_{H_\beta}$, $\|\cdot\|_{\ell^2}$, $\|\cdot\|_{h_\beta}$ from Section 2.3, we have the isometries

$$\|f\|_{L^2} = \|\tilde{f}\|_{\ell^2}, \quad \|f\|_{H_\beta} = \|\tilde{f}\|_{h_\beta},$$

implying that it is sufficient to show the convergence of the coefficients in the sequence space.

Consider the change of variables $\hat{f} = \{|k^{\beta_*}| f_{k,l}\}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}}$. The model (10.26) can be rewritten into,

$$\hat{X}_{k,l}^{(n)} = \hat{f}_{k,l} + \frac{1}{\sqrt{n}} z_{k,l}, \quad \text{for } k \in k, l \}_{(k,l) \in \mathbb{N}^d \times \mathbb{N}},$$

where $z_{k,l}$ are independent standard Gaussian random variables. Notice that the prior of f induces a prior of \hat{f} in $\ell^2(\mathbb{N}^d \times \mathbb{N})$. Therefore, using Lemma 10.13, we obtain the posterior contraction rate of the induced prior in $\|\cdot\|_{\ell^2(d+1)}$.

For the preceding change of variables, an isometry $\|\hat{f}\|_{\ell^2} = \|\tilde{f}\|_{h_{\beta_*}}$ holds, given \tilde{f} is in h_{β_*} . The isometry implies that rates of \hat{f} relative to $\|\cdot\|_{\ell^2}$ can be translated to the rates of \tilde{f} relative to $\|\cdot\|_{h_{\beta_*}}$. Consequently, the result can be translated to the rate in $\|\cdot\|_{H_{\beta_*}}$ and the proof is complete.

10.4 Entropy Number with Non-Polynomial Rates

In Section 2.4, we have shown the estimate of metric numbers of the embedding $\iota : H_{s+\underline{t}} \rightarrow H_t$ with the scale. The same argument can be applied to $\iota : \tilde{H} \rightarrow H_t$ where \tilde{H} is another Hilbert space contained in the smoothness class $\{H_t\}_t$, see the lemma below.

Lemma 10.14. *Given \mathbb{H}_G from (10.13) and the isotropic Sobolev spaces $\{H_s\}_{s \in \mathbb{R}}$ defined in Section 2.3.1. Then for the canonical embedding $\iota : \mathbb{H}_G \rightarrow H_s$, when j is large enough, the entropy number is of the order*

$$e^{-c_1 j^{\nu/(\nu+d)}} \leq e_j(\iota : \mathbb{H}_G \rightarrow H_s) \leq e^{-c_2 j^{\nu/(\nu+d)}}, \quad (10.30)$$

where c_1, c_2 are universal positive constants.

Proof. The proof follows the same argument from Section 2.4 (also see the Appendix B in [43]). The singular values of $\iota : \mathbb{H}_G \rightarrow H_s$ are of order $e^{-Tj^{\nu/d}} j^{-(\alpha-s)/d}$, for which the upper bound is obtained using the same argument in (10.15) and the lower bound is obtained by taking the unit vector φ_k such that $k_i \simeq j^{1/d}$. Consequently, the approximation numbers $a_j(\iota : \mathbb{H}_G \rightarrow H_s)$ have the same order. In particular, $a_j(\iota : \mathbb{H}_G \rightarrow H_s) = O(e^{-Tj^{\nu/d}})$. By the second example in Section 3 of [103], we obtain the final statement of the lemma. \square

Because $\varepsilon \mapsto H(\varepsilon, \iota)$ is the inverse mapping of $j \mapsto e_j(\iota)$, we obtain the corollary below.

Corollary 10.15 (Metric entropy). *The metric entropy of the unit ball of \mathbb{H}_G , given in (10.13), in isotropic Sobolev spaces $\{H_s\}_{s \in \mathbb{R}}$ is given by*

$$H(\varepsilon, \iota) := \log N\left(\varepsilon, \{g \in \mathbb{H}_G : \|g\|_{\mathbb{H}_G} \leq 1\}, \|\cdot\|_{H_s}\right) \sim \left(\log \frac{1}{\varepsilon}\right)^{\frac{\nu+d}{\nu}},$$

as $\varepsilon \downarrow 0$.

Appendix

Appendix A

Mathematical Tools

In this appendix we collect the mathematical elements, mainly from operator theory, that serve as the underlying language and building blocks for this thesis. They are from well established fields and can be found in textbooks and monographs. Hence, results will be present, and proofs are referred to the literature.

Operators are ubiquitous in this thesis, as one main component of the Gaussian linear model, the transform \mathcal{A} , is an operator. In particular, compact operators is of great importance, which is demonstrated by the following examples. First, the ill-posedness in a large class of linear inverse problems is characterised as the compactness of transform operators. Second, a Gaussian measure is a proper probability measure (instead of a generalised stochastic process) only when its covariance operator is of trace class, which is necessarily compact. Third, an element in a compact space can be well approximated by a finite-dimensional subspace, and the error estimate is closely related to the compactness. Besides the aforementioned cases, there are other places where the compactness is leveraged.

In this section, we collect the necessary information on operator theory, with special attention to compact operators. All the materials are standard and can be found in many textbooks, e.g. [102].

First let us summarize the common notations for operators. Let X, Y be normed spaces over the field \mathbb{R} . A linear operator \mathcal{T} from X to Y is a linear mapping from the *domain* of \mathcal{T} , i.e. a subspace of X denoted by $\text{Dom } \mathcal{T}$, into Y . The image of \mathcal{T} is called *range*, i.e. $\text{Ran } \mathcal{T} = \mathcal{T}(\text{Dom } \mathcal{T}) = \{\mathcal{T}f : f \in \text{Dom } \mathcal{T}\}$. A linear operator from X to \mathbb{R} is a linear *functional*. The notation $\mathcal{T} : X \rightarrow Y$ is understood as $\text{Dom } \mathcal{T} = X$ and $\text{Ran } \mathcal{T} \subseteq Y$, unless the domain is given explicitly.

An operator is *injective* precisely when $\mathcal{T}f = 0$ implies $f = 0$. For an injective operator, the *inverse* \mathcal{T}^{-1} of \mathcal{T} is given by

$$\text{Dom } \mathcal{T}^{-1} = \text{Ran } \mathcal{T}, \quad \mathcal{T}^{-1}g = f, \quad \text{for } g = \mathcal{T}f \in \text{Ran } \mathcal{T}.$$

The space of bounded linear operators from X to Y is denoted as $B(X, Y)$, i.e.

$$B(X, Y) := \left\{ \mathcal{T} : X \rightarrow Y \mid \text{linear and } \|\mathcal{T}\|_{X \rightarrow Y} := \sup_{h \in X: \|h\| \leq 1} \|\mathcal{T}h\|_Y < \infty \right\},$$

where $\|\cdot\|_{X \rightarrow Y}$ is the operator norm and $\|\cdot\|$ may be used if no danger. If $X = Y$, we write $L(X)$.

Definition A.1 (Adjoint). Let X, Y be Banach spaces. The *adjoint* \mathcal{T}^* of a densely defined (not necessarily bounded) linear operator $\mathcal{T} : X \rightarrow Y$ is the operator uniquely determined by

$$\begin{aligned} \mathcal{T}^*y^* &= x^*, \\ \langle y^*, \mathcal{T}x \rangle &= \langle x^*, x \rangle, \quad \forall x \in \text{Dom } \mathcal{T}. \end{aligned}$$

An densely defined operator $\mathcal{S} : X \rightarrow X$ is *self-adjoint* if $\text{Dom } \mathcal{S} = \text{Dom } \mathcal{S}^*$ and $\langle \mathcal{S}h, g \rangle = \langle h, \mathcal{S}^*g \rangle$, for all $h, g \in \text{Dom } \mathcal{S}$.

Remark A.2. If X and Y are Hilbert spaces, the dual spaces X^* and Y^* can be identified with the original space by Riesz representation theorem. If the operator $\mathcal{T} : X \rightarrow Y$ is bounded, then the definition above is equivalent to the standard definition of adjoints on Hilbert spaces, that there exists a unique operator $\mathcal{T}^* : Y \rightarrow X$ such that

$$\langle \mathcal{T}x, y \rangle_Y = \langle x, \mathcal{T}^*y \rangle_X,$$

for all $x \in X$ and $y \in Y$.

Definition A.3 (Positivity). An operator \mathcal{T} on a Hilbert space is called *positive*, denoted by $\mathcal{T} \geq 0$, if $\langle \mathcal{T}h, h \rangle \geq 0$, for all $h \in \text{Dom } \mathcal{T}$. For two positive operators \mathcal{S}, \mathcal{T} , we write $\mathcal{S} \geq \mathcal{T}$ if $\text{Dom } \mathcal{S} \subset \text{Dom } \mathcal{T}$ and $\mathcal{S} - \mathcal{T} \geq 0$ on $\text{Dom } \mathcal{S}$. We also write $\mathcal{S} = \mathcal{T}$ if $\mathcal{S} \geq \mathcal{T}$ and $\mathcal{T} \geq \mathcal{S}$.

If the above properties hold up to independent constants, then we use the notations $\mathcal{S} \lesssim \mathcal{T}$ and $\mathcal{S} \simeq \mathcal{T}$.

A.1 Miscellaneous Lemmas

In this section, we collect a few useful lemmas.

The following lemma is known the bounded linear transform (BLT) theorem, (see Theorem I.7, [81]).

Lemma A.4 (BLT theorem). *Let \mathcal{T} be a bounded linear operator from $(X, \|\cdot\|_X)$ to a complete normed space Y . Then there exists a unique bounded extension $\tilde{\mathcal{T}}$ of \mathcal{T} from the completion of X under $\|\cdot\|_X$ to Y .*

The following lemma is a direct consequence of Hahn-Banach theorem.

Lemma A.5. *Given a normed space $(E, \|\cdot\|)$ with its topological dual E^* , the following holds*

$$\|x\| = \sup_{f \in U(E^*)} |\langle f, x \rangle|.$$

Using positivity, we have another characterisation of operator norms on Hilbert spaces.

Lemma A.6. *Let \mathcal{T} be an positive element in $B(H)$. Then,*

$$\|\mathcal{T}\| = \sup_{\|x\| \leq 1} \langle \mathcal{T}x, x \rangle.$$

The following result provides the soundness to Gelfand triples.

Lemma A.7. *Let G and H be two Banach spaces such that G is a dense subset of H , and the embedding $\iota : G \rightarrow H, g \mapsto g$ is continuous. Then, the following hold.*

(i) *The inclusion mapping $\tilde{\iota} : H^* \rightarrow G^*$, $\ell \mapsto \ell|_G$, where $\ell|_G$ is the restriction of ℓ to set G , is continuous. In particular,*

$$\langle \ell, g \rangle_{H^* \times H} = \langle \ell|_G, g \rangle_{G^* \times G}, \quad \forall \ell \in H^*, \forall g \in G. \quad (\text{A.1})$$

(ii) *H^* is dense in G^* , if G is reflexive.*

In particular, if H is a Hilbert space and G is reflexive, we have

$$G \subset H = H^* \subset G^*.$$

Proof. First we show the continuity of $\tilde{\iota}$. Notice that for all $g \in G$, $\|g\|_H \lesssim \|g\|_G$, because of the continuity of ι . For any $\ell \in H^*$, we have

$$|\ell(g)| \lesssim \|\ell\|_H \|g\|_G.$$

Let $\tilde{\ell}$ be the restriction of ℓ to the subset $G \subset H$. Then, $\tilde{\ell} \in G^*$ such that

$$\tilde{\ell}(g) = \ell(g), \quad \forall g \in G, \quad (\text{A.2})$$

and

$$\|\tilde{\ell}\|_{G^*} \leq \|\ell\|_{H^*}, \quad \forall \ell \in H^*.$$

In addition, $\tilde{\ell} = 0$ implies $\ell = 0$. This is because of (A.2) and the density of G in H . Hence the inclusion mapping $\tilde{\iota} : \ell \rightarrow \tilde{\ell}$ is injective and continuous, and (A.1) holds.

Now we are going to show that H^* is dense in G^* by contradiction. If the statement is not true, then the closure of H^* in G^* is a proper closed subspace of G^* . By Hahn-Banach theorem, there exists a non-zero functional $\varphi_g \in (G^*)^*$ such that $\varphi_g(\tilde{\ell}) = 0$ for all $\tilde{\ell} \in \tilde{\iota}(H^*) \subset G^*$. Because of reflexivity, the functional can be identified with an element $g \in G$, such that $\varphi_g(\tilde{\ell}) = \tilde{\ell}(g) = 0$, for all $\tilde{\ell} \in \tilde{\iota}(H^*) \subset G^*$. Due to (A.1), $\ell(g) = 0$, for all $\ell \in H^*$. Since $g \in G \subset H$, it implies that $g = 0$, which contradicts to $\varphi_g \neq 0$. □

An embedding of Hilbert spaces naturally gives rise to an isometric isomorphism, which is useful in several occasions in this thesis. Meanwhile, it also shares some similar flavour of Lemma A.7.

Lemma A.8. *Assume that Hilbert space H is a dense subspace of Hilbert space X such that $\|h\|_H \geq \|h\|_X$, for all $h \in H$, and let the canonical embedding be*

$$\iota : H \rightarrow X, \quad h \mapsto h.$$

Then,

$$\mathcal{U} = (\iota^*)^{-1/2} : \text{Dom } \mathcal{U} \subset X \rightarrow X,$$

where $\text{Dom } \mathcal{U} = H$, is an isometric isomorphism, i.e. $\|\mathcal{U}h\|_X = \|h\|_H$.

Proof. This proof is adopted from Theorem IV.1.12, [63].

Since ι is compact, so is $\mathcal{S} = \iota^*$. Furthermore, $\mathcal{S} : X \rightarrow X$ is self-adjoint and positive, and $\text{Ran } \mathcal{S} \subset H$. We can define a self-adjoint operator $\mathcal{T} = \mathcal{S}^{-1}$ on domain $\text{Dom } \mathcal{T} = \text{Ran } \mathcal{S} \subset H$, such that

$$\langle h, g \rangle_H = \langle \mathcal{T}h, g \rangle_X, \tag{A.3}$$

for all $h \in \text{Dom } \mathcal{T}$ and $g \in H$.

Using spectral theorem, define an operator $\mathcal{U} = (\mathcal{T})^{1/2}$, whose domain $\text{Dom } \mathcal{U}$ is the closure of $\text{Dom } \mathcal{T}$ with respect to the norm

$$\|\mathcal{U}h\|_X = \sqrt{\langle \mathcal{T}h, h \rangle_X} = \|h\|_H.$$

The domain $\text{Dom } \mathcal{U}$ a closed set in H . We are going to show in fact $\text{Dom } \mathcal{U} = H$ by contradiction. Assume $\text{Dom } \mathcal{U} \subsetneq H$. Then by Hahn-Banach theorem, there exists an element $h_0 \in H$ such that $\langle g, h_0 \rangle_H = 0$ for all $g \in \text{Dom } \mathcal{U}$, and in particular, all $g \in \text{Dom } \mathcal{T}$. Due to (A.3), we have $\langle \mathcal{U}g, h_0 \rangle_X = 0$, for all $g \in \text{Dom } \mathcal{U}$. Since $\text{Ran } \mathcal{U} = X$, we conclude that $h_0 = 0$, which leads to a contradiction. \square

A.2 Pseudo-Inverse

Let H be a Hilbert space and G be a normed space.

Definition A.9. Let $\mathcal{T} \in B(H, G)$, and $\text{Ker } \mathcal{T} = \{h \in H \mid \mathcal{T}h = 0\}$. The *pseudo-inverse* is defined as:

$$\mathcal{T}^{-1} := (\mathcal{T}|_{(\text{Ker } \mathcal{T})^\perp})^{-1} : \mathcal{T}((\text{Ker } \mathcal{T})^\perp) = \mathcal{T}(H) \rightarrow (\text{Ker } \mathcal{T})^\perp,$$

which is bijective by construction.

Remark A.10. For $g \in \mathcal{T}(H)$, one can let $\mathcal{T}^{-1}g \in H$ be the solution of operator equation $\mathcal{T}h = g$ with the minimal norm. This gives an equivalent definition of pseudo-inverse.

When \mathcal{T} has a genuine inverse, it induces an inner product on its image, as shown in the following lemma.

Lemma A.11. *Suppose that $\mathcal{T} : H \rightarrow G$ is an injective linear operator. Then, the range $\mathcal{T}(H) : \text{Ran } \mathcal{T} \subset G$ of \mathcal{T} is a Hilbert space equipped with the inner product*

$$\langle x, y \rangle_{\mathcal{T}(H)} := \langle \mathcal{T}^{-1}x, \mathcal{T}^{-1}y \rangle_H, \quad x, y \in \mathcal{T}(H). \tag{A.4}$$

In addition, $\mathcal{T} : H \rightarrow T$ is bounded and its adjoint is $\mathcal{T}^* = \mathcal{T}^{-1}$.

Proof. Since \mathcal{T} is injective, the inner product given above is well-defined. Let $\{x_n\}$ be a Cauchy sequence in T . Then $x_n = \mathcal{T}h_n$ with $h_n \in H$ and $\{h_n\}$ is Cauchy as well because $\|h_m - h_n\|_H = \|x_m - x_n\|_T$. Since H is Hilbert, there exists a h such that $h_n \rightarrow h$. Therefore, there exists a vector $x = \mathcal{T}h \in T$ and

$$\|h_n - h\|_H = \|x_n - x\|_T \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

The boundedness and the adjoint follow directly from the construction. □

The previous lemma has the following implications on the pseudo-inverse.

Corollary A.12. *Let $\mathcal{T} \in B(H)$ and \mathcal{T}^{-1} be the pseudo-inverse of \mathcal{T} . Then, the following statements hold.*

- (i) $\mathcal{T}(H)$ is an Hilbert space equipped with the inner product induced by \mathcal{T} given in (A.4).
- (ii) If $\{e_k\}_k$ is an orthonormal basis of $(\text{Ker } \mathcal{T})^\perp$, then $\{\mathcal{T}e_k\}_k$ is an orthonormal basis for the Hilbert space $\mathcal{T}(H)$ in (i).

The proof of the following lemma is given in Proposition C.0.5, [68].

Lemma A.13. *Let H_1, H_2 and G be Hilbert spaces, and let $\mathcal{T}_1 \in B(H_1, G)$ and $\mathcal{T}_2 \in B(H_2, G)$. If $\|\mathcal{T}_1^*g\|_1 = \|\mathcal{T}_2^*g\|_2$ for all $g \in G$, then for all*

$$g \in \text{Ran } \mathcal{T}_1 = \text{Ran } \mathcal{T}_2,$$

we have $\|\mathcal{T}_1^{-1}g\|_1 = \|\mathcal{T}_2^{-1}g\|_2$.

Proposition A.14. *Let $\mathcal{T} \in B(H, G)$ and $\mathcal{Q} = \mathcal{T}\mathcal{T}^* \in B(G)$. Then, for all*

$$g \in \text{Ran } \mathcal{Q}^{1/2} = \text{Ran } \mathcal{T},$$

we have $\|\mathcal{Q}^{-1/2}g\| = \|\mathcal{T}^{-1}g\|_H$.

Proof. Since \mathcal{Q} is self-adjoint, the square root is defined via the spectral theorem. Furthermore, for all $g \in G$, we have

$$\left\| (\mathcal{Q}^{1/2})^*g \right\|^2 = \left\| \mathcal{Q}^{1/2}g \right\|^2 = \langle g, \mathcal{Q}g \rangle = \|\mathcal{T}^*g\|_H.$$

The rest follows from the previous lemma. □

A.3 Compact Operators

Let H, G be two Hilbert spaces. We denote the inner product of H by $\langle \cdot, \cdot \rangle_H$, or simply $\langle \cdot, \cdot \rangle$ when there is no confusion.

Definition A.15 (Compact operators). An operator $\mathcal{T} : H \rightarrow G$ is compact if for any bounded sequence $\{h_n\}$ in H , the sequence $\{\mathcal{T}h_n\}$ in G contains a convergent subsequence. The space of compact operators in $L(H, G)$, equipped with the operator norm, is denoted by $S_\infty(H, G)$.

If $\mathcal{T} : H \rightarrow G$ is an compact operator, then $\mathcal{T}^*\mathcal{T}$ is compact, self-adjoint, and *non-negative*, i.e. $\langle \mathcal{T}^*\mathcal{T}h, h \rangle \geq 0$ for all $h \in H$. The *absolute value* of \mathcal{T} is defined with the equality $|\mathcal{T}| = (\mathcal{T}^*\mathcal{T})^{1/2}$, where $(\mathcal{T}^*\mathcal{T})^{1/2}$ is the unique non-negative square root of $\mathcal{T}^*\mathcal{T}$ (see Theorem 7.4 in [102]). The positive eigenvalues of $|\mathcal{T}|$ are called the *singular values* of \mathcal{T} . In fact, singular values $\{s_j(\mathcal{T})\}$ encode great information about the operator \mathcal{T} .

Theorem A.16 (Singular Value Decomposition (Theorem 7.6, [102])). *Let $\mathcal{T} : H \rightarrow G$ be a compact operator and $\{s_j\}$ denote the (possibly finite) non-decreasing sequence of the singular values of \mathcal{T} . There exists orthonormal sequence (h_j) from H and (g_j) from G such that for all $h \in H$ and $g \in G$,*

$$\begin{aligned} \mathcal{T}h &= \sum_j s_j \langle h, h_j \rangle g_j, & \mathcal{T}^*g &= \sum_j s_j \langle g, g_j \rangle h_j, \\ |\mathcal{T}|h &= \sum_j s_j \langle h, h_j \rangle h_j, & |\mathcal{T}^*|g &= \sum_j s_j \langle g, g_j \rangle g_j. \end{aligned}$$

The (h_j) and (g_j) are the eigenvectors of $|\mathcal{T}|$ and $|\mathcal{T}^|$, respectively. In particular, $\mathcal{T}, |\mathcal{T}|, \mathcal{T}^*$ and $|\mathcal{T}^*|$ have the same singular values.*

With the help of singular values, we can define the following spaces of operators.

Definition A.17 (Schatten Class). For a compact operator $\mathcal{T} : H \rightarrow G$, the *p-Schatten norm*, $p \in [1, \infty)$, is defined with its singular values $\{s_j(\mathcal{T})\}$,

$$\|\mathcal{T}\|_p := \left(\sum_j |s_j(\mathcal{T})|^p \right)^{1/p}.$$

We denote by $S_p(H, G)$ the set of compact operators with finite *p-Schatten norm*.

Remark A.18. It is not difficult to show that $\|\mathcal{T}\|_\infty = s_1(\mathcal{T})$ from the definition.

The *p-Schatten norms* are authentic norms satisfying the triangle inequality. S_p spaces are similar to L^p spaces. For example, the spaces S_p are Banach spaces and in particular, S_2 is a Hilbert space. A version of the Hölder inequality also holds in S_p spaces. These properties are summarised in the following propositions.

Proposition A.19 (Lemma 10 and 14, XI.9, [25]). *Let $1 \leq p \leq p' \leq \infty$ and let $1/p + 1/q = 1$.*

- (i) *For $\mathcal{T}, \mathcal{U} \in S_p$, we have $\|\mathcal{T} + \mathcal{U}\|_p \leq \|\mathcal{T}\|_p + \|\mathcal{U}\|_p$.*
- (ii) *S_p is complete under the norm $\|\cdot\|_p$. S_2 is an inner product space.*
- (iii) *$S_p \subset S_{p'}$ and $\|h\|_p \geq \|h\|_{p'}$, if $h \in S_p$.*
- (iv) *For $\mathcal{T} \in S_p$ and $\mathcal{U} \in S_q$, then $\|\mathcal{T}\mathcal{U}\|_1 \leq \|\mathcal{T}\|_p \|\mathcal{U}\|_q$.*

Proposition A.20 (Theorem 7.8, [102]).

(i) If $\mathcal{T} \in S_p$ and $\mathcal{U} \in S_q$ ($p, q \in [1, \infty)$) and $1/r = 1/p + 1/q$, then

$$\|\mathcal{T}\mathcal{U}\|_r \leq 2^{1/r} \|\mathcal{T}\|_p \|\mathcal{U}\|_q.$$

(ii) If $\mathcal{T} \in L(H_1, H_2)$ and $\mathcal{U} \in S_p(H_0, H_1)$, then

$$\|\mathcal{T}\mathcal{U}\|_p \leq \|\mathcal{T}\| \|\mathcal{U}\|_p.$$

For $\mathcal{T} \in S_p(H_1, H_2)$ and $\mathcal{U} \in L(H_0, H_1)$, the corresponding assertion also holds.

Schatten class contains two important sets of compact operators.

- $S_2(H, G)$ is identical to the set of *Hilbert-Schmidt* operators. Namely, for any Hilbert-Schmidt operator $\mathcal{T} : H \rightarrow G$,

$$\|\mathcal{T}\|_2 = \|\mathcal{T}\|_{HS} := \sum_{i \in \mathcal{I}} \|\mathcal{T}\varphi_i\|_G^2 < \infty,$$

where $\{\varphi_i\}_{i \in \mathcal{I}}$ is an orthonormal basis in H .

- $S_1(H, G)$ is identical to the set of *trace class* operators. Namely, for any operator $\mathcal{T} : H \rightarrow G$ of trace class, we have

$$\|\mathcal{T}\|_1 = \text{Trace } \mathcal{T} := \sum_{i \in \mathcal{I}} \langle \sqrt{\mathcal{T}^* \mathcal{T}} \varphi_i, \varphi_i \rangle_H < \infty,$$

where $\{\varphi_i\}_{i \in \mathcal{I}}$ is an orthonormal basis in H .

The two equalities above are obvious. Since the Hilbert-Schmidt norm and trace are both independent of the basis $\{\varphi_i\}_{i \in \mathcal{I}}$, the equalities are obtained by taking the basis to be the eigenbasis of $|\mathcal{T}|$. Furthermore, the following can be derived directly from (iii) in the previous proposition,

$$\|\mathcal{T}\| \leq \|\mathcal{T}\|_{HS} = \|\mathcal{T}^*\|_{HS} \leq \text{Trace } \mathcal{T},$$

where the equality of Hilbert-Schmidt norms can be found in Theorem 6.9 in [102].

References

- [1] S. AGAPIOU, S. LARSSON, AND A. M. STUART, *Posterior contraction rates for the bayesian approach to linear ill-posed inverse problems*, Stochastic Processes and their Applications, 123 (2013), pp. 3828 – 3860.
- [2] A. AKANSU AND H. AGIRMAN-TOSUN, *Generalized discrete fourier transform with nonlinear phase*, Signal Processing, IEEE Transactions on, 58 (2010), pp. 4547–4556.
- [3] P. ALQUIER, E. GAUTIER, AND G. STOLTZ, *Inverse Problems and High-Dimensional Estimation: Stats in the Château Summer School, August 31 - September 4, 2009*, Lecture Notes in Statistics, Springer, 2011.
- [4] J. ARBEL, G. GAYRAUD, AND J. ROUSSEAU, *Bayesian optimal adaptive estimation using a sieve prior*, Scandinavian Journal of Statistics, 40 (2013), pp. 549–570.
- [5] M. BIRKE, N. BISSANTZ, AND H. HOLZMANN, *Confidence bands for inverse regression models*, Inverse Problems, 26 (2010), p. 115020.
- [6] N. BISSANTZ, H. DETTE, AND K. PROKSCH, *Model checks in inverse regression models with convolution-type operators*, Scandinavian Journal of Statistics, 39 (2012), pp. 305–322.
- [7] N. BISSANTZ, T. HOHAGE, A. MUNK, AND F. RUYMGAART, *Convergence rates of general regularization methods for statistical inverse problems and applications*, SIAM Journal on Numerical Analysis, 45 (2007), pp. 2610–2636.
- [8] V. BOGACHEV, *Gaussian Measures*, Mathematical surveys and monographs, American Mathematical Society, 1998.
- [9] V. I. BOGACHEV, *Measure Theory*, vol. Volume 2, Springer, 1 ed., 2007.
- [10] H. BREZIS, *Functional Analysis, Sobolev Spaces and Partial Differential Equations*, Universitext, Springer New York, 2010.
- [11] L. D. BROWN AND M. G. LOW, *Asymptotic equivalence of nonparametric regression and white noise*, Ann. Statist., 24 (1996), pp. 2384–2398.

- [12] C. CANUTO, M. HUSSAINI, A. QUARTERONI, AND T. ZANG, *Spectral Methods: Fundamentals in Single Domains*, Scientific Computation, Springer Berlin Heidelberg, 2010.
- [13] R. CARMONA AND M. TEHRANCHI, *Interest Rate Models: an Infinite Dimensional Stochastic Analysis Perspective*, Springer Finance, Springer Berlin Heidelberg, 2007.
- [14] I. CASTILLO AND R. NICKL, *Nonparametric bernstein-von mises theorems in gaussian white noise*, Ann. Statist., 41 (2013), pp. 1999–2028.
- [15] L. CAVALIER, *Nonparametric statistical inverse problems*, Inverse Problems, 24 (2008), p. 034004.
- [16] L. CAVALIER AND A. TSYBAKOV, *Sharp adaptation for inverse problems with random noise*, Probability Theory and Related Fields, 123 (2002), pp. 323–354.
- [17] P.-L. CHOW, I. A. IBRAGIMOV, AND R. Z. KHASHMINSKII, *Statistical approach to some ill-posed problems for linear partial differential equations*, Probability Theory and Related Fields, 113 (1999), pp. 421–441.
- [18] A. COHEN, *Numerical Analysis of Wavelet Methods*, Studies in mathematics and its applications, Elsevier, 2003.
- [19] A. COHEN, M. HOFFMANN, AND M. REISS, *Adaptive wavelet galerkin methods for linear inverse problems*, SIAM Journal on Numerical Analysis, 42 (2004), pp. 1479–1501.
- [20] D. COLTON AND R. KRESS, *Inverse Acoustic and Electromagnetic Scattering Theory*, Applied Mathematical Sciences, Springer New York, 2012.
- [21] J. CONWAY, *A Course in Functional Analysis*, Graduate Texts in Mathematics, Springer, 1990.
- [22] H. CRAMER, *Mathematical methods of statistics*, Princeton paperbacks Princeton landmarks in mathematics and physics, Princeton University Press, 1999.
- [23] G. DA PRATO AND J. ZABCZYK, *Stochastic Equations in Infinite Dimensions*, Encyclopedia of Mathematics and its Applications, Cambridge University Press, 2014.
- [24] D. L. DONOHO, *Nonlinear solution of linear inverse problems by wavelet-vaguelette decomposition*, Applied and Computational Harmonic Analysis, 2 (1995), pp. 101–126.
- [25] N. DUNFORD AND J. SCHWARTZ, *Linear Operators, Part II: Spectral theory, self-adjoint operator in Hilbert space*, Pure and Applied Mathematics, Interscience Publishers, 1963.

-
- [26] R. EDMUNDS, *Inequalities between entropy and approximation numbers of compact maps*, *Zeitschrift für Analysis und ihre Anwendungen*, 7 (1988), pp. 223–227.
- [27] S. EFROMOVICH, *Simultaneous sharp estimation of functions and their derivatives*, *Ann. Statist.*, 26 (1998), pp. 273–278.
- [28] K. ENGEL, S. BRENDLE, R. NAGEL, M. CAMPITI, T. HAHN, G. METAFUNE, G. NICKEL, D. PALLARA, C. PERAZZOLI, A. RHANDI, ET AL., *One-Parameter Semigroups for Linear Evolution Equations*, Graduate Texts in Mathematics, Springer New York, 2006.
- [29] H. ENGL, M. HANKE, AND A. NEUBAUER, *Regularization of Inverse Problems*, Mathematics and Its Applications, Springer Netherlands, 2000.
- [30] J.-P. FLORENS AND A. SIMONI, *Regularizing priors for linear inverse problems*, *Econometric Theory*, 32 (2016), pp. 71–121.
- [31] M. GADELLA AND F. GÓMEZ, *Dirac formulation of quantum mechanics: Recent and new results*, *Reports on Mathematical Physics*, 59 (2007), pp. 127 – 143.
- [32] I. GEL'FAND AND N. VILENKIN, *Generalized Functions, Volume 4*, AMS Chelsea Publishing, American Mathematical Society, 2016.
- [33] S. GHOSAL, J. K. GHOSH, AND A. W. VAN DER VAART, *Convergence rates of posterior distributions*, *The Annals of Statistics*, 28 (2000), pp. 500–531.
- [34] S. GHOSAL, J. LEMBER, AND A. VAN DER VAART, *Nonparametric Bayesian model selection and averaging*, *Electron. J. Stat.*, 2 (2008), pp. 63–89.
- [35] S. GHOSAL AND A. VAN DER VAART, *Fundamentals of Nonparametric Bayesian Inference*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 2017.
- [36] E. GINÉ AND R. NICKL, *Mathematical Foundations of Infinite-Dimensional Statistical Models*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 2016.
- [37] A. GOLDENSHLUGER AND S. V. PEREVERZEV, *Adaptive estimation of linear functionals in hilbert scales from indirect white noise observations*, *Probability Theory and Related Fields*, 118 (2000), pp. 169–186.
- [38] ———, *On adaptive inverse estimation of linear functionals in hilbert scales*, *Bernoulli*, 9 (2003), pp. 783–807.
- [39] R. GORENFLO AND M. YAMAMOTO, *Operator-theoretic treatment of linear Abel integral equations of first kind*, *Japan J. Indust. Appl. Math.*, 16 (1999), pp. 137–161.

- [40] G. G. GOULD, *The spectral representation of normal operators on a rigged hilbert space*, Journal of the London Mathematical Society, s1-43 (1968), pp. 745–754.
- [41] G. GRUBB, *Distributions and Operators*, Graduate Texts in Mathematics, Springer New York, 2010.
- [42] S. GUGUSHVILI, A. VAN DER VAART, AND D. YAN, *Bayesian inverse problems with partial observations*, Transactions of A. Razmadze Mathematical Institute, 172 (2018), pp. 388 – 403.
- [43] S. GUGUSHVILI, A. VAN DER VAART, AND D. YAN, *Bayesian linear inverse problems in regularity scales*, Ann. Inst. H. Poincaré Probab. Statist., (2019 (accepted)).
- [44] M. HAASE, *Functional Analysis: An Elementary Introduction*, Graduate Studies in Mathematics, Amer Mathematical Society, 2014.
- [45] D. HAROSKE AND H. TRIEBEL, *Distributions, Sobolev Spaces, Elliptic Equations*, EMS Monographs in mathematics, European Mathematical Society, 2008.
- [46] M. HEGLAND, *Variable hilbert scales and their interpolation inequalities with applications to tikhonov regularization*, Applicable Analysis, 59 (1995), pp. 207–223.
- [47] T. HIDA, *Brownian motion*, Applications of mathematics, Springer-Verlag, 1980.
- [48] H. HOLDEN, B. OKSENDAL, J. UBOE, AND T. ZHANG, *Stochastic Partial Differential Equations: A Modeling, White Noise Functional Approach*, Probability and Its Applications, Birkhäuser Boston, 2013.
- [49] T. HYTÖNEN, J. VAN NEERVEN, M. VERAAR, AND L. WEIS, *Analysis in Banach Spaces: Volume I: Martingales and Littlewood-Paley Theory*, Ergebnisse der Mathematik und ihrer Grenzgebiete. 3. Folge / A Series of Modern Surveys in Mathematics, Springer International Publishing, 2016.
- [50] I. IBRAGIMOV AND R. HAS’MINSKII, *Statistical Estimation: Asymptotic Theory*, Stochastic Modelling and Applied Probability, Springer New York, 2013.
- [51] I. A. IBRAGIMOV AND R. V. KHAS’MINSKII, *Estimation problems for coefficients of stochastic partial differential equations. part ii*, Theory of Probability & Its Applications, 44 (2000), pp. 469–494.
- [52] I. A. IBRAGIMOV AND R. Z. KHAS’MINSKII, *Estimation problems for coefficients of stochastic partial differential equations. part i*, Theory of Probability & Its Applications, 43 (1999), pp. 370–387.
- [53] Y. INGSTER AND N. STEPANOVA, *Estimation and detection of functions from anisotropic sobolev classes*, Electron. J. Statist., 5 (2011), pp. 484–506.

-
- [54] V. ISAKOV, *Inverse Problems for Partial Differential Equations*, Applied Mathematical Sciences, Springer New York, 2013.
- [55] E. M. J. L. LIONS, *Non-Homogeneous Boundary Value Problems and Applications: Vol. 1*, Die Grundlehren der mathematischen Wissenschaften 181, Springer-Verlag Berlin Heidelberg, 1 ed., 1972.
- [56] J. KAIPIO AND E. SOMERSALO, *Statistical and Computational Inverse Problems*, Applied Mathematical Sciences, Springer New York, 2006.
- [57] A. KIRSCH, *An Introduction to the Mathematical Theory of Inverse Problems*, Applied Mathematical Sciences, Springer, 2011.
- [58] B. KNAPIK AND J.-B. SALOMOND, *A general approach to posterior contraction in nonparametric inverse problems*, Bernoulli, 24 (2018), pp. 2091–2121.
- [59] B. KNAPIK, A. VAN DER VAART, AND J. VAN ZANTEN, *Bayesian inverse problems with gaussian priors*, The Annals of Statistics, 39 (2011), pp. 2626–2657.
- [60] ———, *Bayesian recovery of the initial condition for the heat equation*, Communications in Statistics - Theory and Methods (2013), 42 (2013), pp. 1294–1313.
- [61] B. T. KNAPIK, B. T. SZABÓ, A. W. VAN DER VAART, AND J. H. VAN ZANTEN, *Bayes procedures for adaptive inference in inverse problems for the white noise model*, Probab. Theory Related Fields, 164 (2016), pp. 771–813.
- [62] S. KREIN AND Y. PETUNIN, *Scales of banach spaces*, Russian Mathematical Surveys, 21 (1966), p. 85.
- [63] S. KREIN AND E. SEMENOV, *Interpolation of Linear Operators*, Translations of Mathematical Monographs, American Mathematical Society, 2002.
- [64] J. KUELBS AND W. LI, *Metric entropy and the small ball problem for Gaussian measures*, J. Funct. Anal., 116 (1993), pp. 133–157.
- [65] J. KUELBS, W. LI, AND W. LINDE, *The Gaussian measure of shifted balls*, Probab. Theory Related Fields, 98 (1994), pp. 143–162.
- [66] L. LE CAM, *Asymptotic Methods in Statistical Decision Theory*, Springer Series in Statistics, Springer-Verlag New York, 1 ed., 1986.
- [67] M. LEDOUX AND M. TALAGRAND, *Probability in Banach spaces*, vol. 23, Springer-Verlag, Berlin, 1991.
- [68] W. LIU AND M. RÖCKNER, *Stochastic Partial Differential Equations: An Introduction*, Universitext, Springer International Publishing, 2015.
- [69] B. A. MAIR AND F. H. RUYMGAART, *Statistical inverse estimation in hilbert scales*, SIAM Journal on Applied Mathematics, 56 (1996), pp. 1424–1444.

- [70] P. MATHÉ AND S. V. PEREVERZEV, *Optimal discretization of inverse problems in hilbert scales. regularization and self-regularization of projection methods*, SIAM Journal on Numerical Analysis, 38 (2001), pp. 1999–2021.
- [71] F. NATTERER, *Error bounds for tikhonov regularization in hilbert scales*, Applicable Analysis, 18 (1984), pp. 29–37.
- [72] F. NATTERER, *The Mathematics of Computerized Tomography*, Classics in Applied Mathematics, Society for Industrial and Applied Mathematics, 2001.
- [73] A. NEUBAUER, *When do sobolev spaces form a hilbert scale?*, Proceedings of the American Mathematical Society, 103 (1988), pp. 557–562.
- [74] D. NUALART, *The Malliavin Calculus and Related Topics*, Probability and Its Applications, Springer Berlin Heidelberg, 2006.
- [75] A. PAZY, *Semigroups of linear operators and applications to PDEs*, Applied Mathematical Sciences, Springer, springer ed., 1992.
- [76] A. PIETSCH, *s-numbers of operators in banach spaces*, Studia Mathematica, 51 (1974), pp. 201–223.
- [77] A. PIETSCH, *Eigenvalues and S-Numbers*, Mathematik und ihre Anwendungen in Physik und Technik, Akademische Verlagsgesellschaft Geest & Portig, 1987.
- [78] A. QUARTERONI, R. SACCO, AND F. SALERI, *Numerical Mathematics*, Texts in Applied Mathematics, Springer, 2010.
- [79] A. QUARTERONI AND A. VALLI, *Numerical Approximation of Partial Differential Equations*, Springer Series in Computational Mathematics, Springer Berlin Heidelberg, 2009.
- [80] K. RAY, *Bayesian inverse problems with non-conjugate priors*, Electron. J. Statist., 7 (2013), pp. 2516–2549.
- [81] M. REED AND B. SIMON, *Methods of Modern Mathematical Physics: Functional analysis*, no. vol. I in Methods of Modern Mathematical Physics, Academic Press, 1980.
- [82] M. REISS, *Asymptotic equivalence for nonparametric regression with multivariate and random design*, Ann. Statist., 36 (2008), pp. 1957–1982.
- [83] W. RUDIN, *Real and Complex Analysis*, Mathematics series, McGraw-Hill, 1987.
- [84] K. SCHMÜDGEN, *Unbounded Self-adjoint Operators on Hilbert Space*, Graduate Texts in Mathematics, Springer Netherlands, 2012.
- [85] B. L. R. SERGEY V. LOTOTSKY, *Stochastic Partial Differential Equations*, Universitext, Springer International Publishing, 2017.

-
- [86] S. E. SHREVE, *Stochastic calculus for finance II: Continuous-time models*, Springer Finance, Springer, 1st ed. 2004. corr. 2nd printing ed., 2004.
- [87] A. V. SKOROHOD, *Integration in Hilbert Space*, Ergebnisse der Mathematik und ihrer Grenzgebiete 79, Springer-Verlag Berlin Heidelberg, 1 ed., 1974.
- [88] A. M. STUART, *Inverse problems: A bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–559.
- [89] B. SZABÓ, A. VAN DER VAART, AND H. VAN ZANTEN, *Empirical bayes scaling of gaussian priors in the white noise model*, Electron. J. Statist., 7 (2013), pp. 991–1018.
- [90] H. TRIEBEL, *Interpolation Theory, Function Spaces, Differential Operators*, Carnegie-Rochester Conference Series on Public Policy, North-Holland Publishing Company, 1978.
- [91] ———, *Theory of Function Spaces III*, Monographs in Mathematics, Birkhäuser Basel, 2006.
- [92] ———, *Function Spaces and Wavelets on Domains*, EMS tracts in mathematics, European Mathematical Society, 2008.
- [93] ———, *Theory of Function Spaces*, Modern Birkhäuser Classics, Springer Basel, 2010.
- [94] ———, *Theory of Function Spaces II*, Modern Birkhäuser Classics, Springer Basel, 2010.
- [95] A. TSYBAKOV, *Introduction to Nonparametric Estimation*, Springer Series in Statistics, Springer, 2008.
- [96] N. VAKHANIA, V. TARIELADZE, AND S. CHOBANYAN, *Probability Distributions on Banach Spaces*, Mathematics and its Applications, Springer Netherlands, 1987.
- [97] A. VAN DER VAART, *Asymptotic Statistics*, Cambridge Series in Statistical and Probabilistic Mathematics, Cambridge University Press, 2000.
- [98] A. VAN DER VAART, *Bayesian regularization*, in Proceedings of the International Congress of Mathematicians. Volume IV, Hindustan Book Agency, New Delhi, 2010, pp. 2370–2385.
- [99] A. W. VAN DER VAART AND J. H. VAN ZANTEN, *Rates of contraction of posterior distributions based on gaussian process priors*, The Annals of Statistics, 36 (2008), pp. 1435–1463.
- [100] ———, *Reproducing kernel Hilbert spaces of Gaussian priors*, vol. Volume 3 of Collections, Institute of Mathematical Statistics, Beachwood, Ohio, USA, 2008, pp. 200–222.

- [101] G. WAHBA, *Practical approximate solutions to linear operator equations when the data are noisy*, SIAM Journal on Numerical Analysis, 14 (1977), pp. 651–667.
- [102] J. WEIDMANN, *Linear Operators in Hilbert Spaces*, Graduate Texts in Mathematics, Springer New York, 1980.
- [103] Z. ZHENG, *Some relations between entropy and approximation numbers*, Science in China Series A: Mathematics, 42 (1999), pp. 478–487.

Index

- Abel operator, 68
- approximation number, 31
- approximation property, 18

- Bayes' formula, 50
- Bayes' rule, 8, 49

- Cameron-Martin
 - formula, 42
 - space, 40
- compact operator, 175
- consistency, 3, 50
- consistent, 50
- contraction rate, 51
- covariance, 146
- covariance operator, 37
 - Gaussian, 40
- credible set, 51
- cylindrical
 - sets, 46
- cylindrical
 - σ -algebra, 36
 - Gaussian measure, 47

- discrete orthogonality, 99
- distribution
 - Bayesian marginal, 49
 - posterior, 50
 - prior, 49

- entropy number, 31
- evolution equations, *see* SPDEs

- factorization, 38
 - Gaussian, 41
- filtration, 141
- forward mappings, 5

- forward operator, 65
- Fourier transform, 36
 - Gaussian, 39

- Galerkin projection, 70
- Gaussian linear model, 4, 52
- Gaussian measure, 38
- Gaussian process, 44
- Gaussian sequence model, 53
- Gelfand triple, 23

- harmonic mean, 29
- heat equation, 98
- Hilbert scale, 20
- Hilbert-Schmidt, 177

- ill-posed, 7
 - extremely, 97
 - mildly, 97
- interpolation, 116
- inverse problems, 7

- Karhunen-Loève expansion, 44
- kernel
 - probability, 49
- Kullback-Leibler divergence, 42

- martingale, 141
- mean, 37
- measure
 - classical Wiener, 44
 - shift, 42
- metric entropy, 32
 - spatial-temporal spaces, 168
- moment
 - strong, 37
 - weak, 37

- norm duality, 18
- observation
 - (indirect) drift term, 156
 - continuous, 52
 - discrete, 55
 - final value, 154
- Poisson equation, 67
- priors
 - Gaussian, 82, 125
 - Gaussian mixture, 85, 126
 - random series, 80, 124
 - spatial Gaussian, 154
 - spatial-temporal Gaussian, 156
- pseudo-inverse, 174
- Radon measure, 36
- Radonification, 46
- reconstruction operator, 56
- regularization, 7, 65
- reproducing kernel Hilbert space, 39
- singular value decomposition, 97, 176
- singular values, 31, 97
- smoothing property, 66
- smoothness class, 17
 - anisotropic, 28
 - isotropic, 29
- smoothness scale, 17
 - multi-dimensional, 30
- Sobolev scale, 22
- Sobolev smoothness, 3
- SPDEs, 7
- Symm's equation, 68
- tensor product, 26
- topological support, 37
- trace class, 177
- trigonometric polynomials, 118
- Volterra operator, 68, 98
- wavelets, 119
- white noise process, 48
- Wiener process
 - cylindrical, 146
 - \mathcal{Q} , 142

Summary

This dissertation studies the Asymptotics of Bayesian nonparametric inference for Gaussian linear models. The models are in the form of

$$\text{Observations} = \text{Transformed signal} + \text{Gaussian noise,} \\ \text{e.g. a smoothed parameter}$$

and the goal is to recover the original signal using Bayesian methods.

In Part I we collect the materials that are essential for the development of this thesis. In particular, we provide theory for the transformed-signal-in-noise model. The goal is to estimate the (original) signal with as few assumptions on the parameter space as possible, and to recover the signal with fast rates. The parameter (i.e., signal) is assumed to possess a certain level of smoothness, which will be quantified in this part of the thesis. The noise structure of the statistical model is also of importance, because it largely determines the difficulty of the recovery problem. To recover the signal we will employ Bayesian methods, which are formally defined in Chapter 4.

In Chapter 2 smoothness classes are introduced to quantify the measure of regularity of the parameter space. As first example we consider smoothness scales. A smoothness scale $\{H_s\}_{s \in \mathbb{R}}$ is a collection of Hilbert spaces H_s that are nested by their norms, and this collection has an additional property which is referred to as a norm duality. A particular interesting subclass of smoothness scales are Hilbert scales. A Hilbert scale is a smoothness scale where the index s naturally defines a generating operator. The resulting generating operator establishes a link between the Hilbert scale and the covariance operator of Gaussian priors and the Gaussian noise. Another type of a smoothness class is used to describe the anisotropic smoothness in higher dimensions, which are used to study stochastic evolution equations, see Part III. We also define approximation numbers, which describe the approximation properties of the smoothness classes, and connect them to metric entropies.

In Chapter 3 we formally introduce Gaussian measures in infinite dimensional spaces. Before doing so, we first review general probability measures on Banach spaces. Once Gaussian measures are formally defined, we show that they have desirable properties. Of particular interest is the covariance structure of Gaussian measures, which we demonstrate with several examples. In addition, we also briefly discuss cylindrical measures and radonification in this chapter.

In Chapter 4 we formally introduce Bayesian inference in infinite dimensional spaces. First Bayes' rule on general function spaces is defined, and the basic

concepts of asymptotic analysis such as consistency and contraction rates are introduced. Bayes rule naturally induces an estimation procedure, which will be applied to Gaussian linear models. In particular, we consider the case with continuous, and the case with discrete observations. This chapter concludes with a general theorem on posterior contraction rates, which provides the general approach throughout this thesis.

Part II investigates linear inverse problems with Gaussian noise, which allows us to use the theory developed for Gaussian linear models. In the setting of inverse problems, one observes a smoothed signal contaminated by noise. The smoothing makes the problem difficult, and naive methods to recover the unsmoothed signal become inappropriate, because the noise dominates due to the inverse operator. One way to overcome this problem is by regularising the inverse operator, which can be accomplished by using Bayesian methods. In this part we investigate inverse problems under two observation schemes: continuously and discretely observed inverse problems, which are commonly known as the white noise model and regression.

Chapter 5 serves as an introduction to linear inverse problems. These problems are introduced both heuristically and mathematically rigorously, and several examples are given. This chapter concludes with the Galerkin method, a relatively general projection method, which is a simple but convenient tool for demonstrating contraction rates for inverse problems used in the remainder of this part.

In Chapter 6, we study the asymptotic performance of Bayesian nonparametric methods for linear inverse problems with continuous observations. A general theorem for continuously observed inverse problems is proved by modifying the proof of the general theorem on contraction rates that is derived in Chapter 4. In general adaptivity is not guaranteed. The additional conditions necessary for adaptivity are given in the second theorem in this chapter. The first theorem is then used to study the performance of two types of priors: random series priors and Gaussian priors. Random series priors lead to posteriors that contract at the minimax rate, up to a logarithmic factor, whereas Gaussian priors are in general suboptimal. To improve on the latter case, we introduce a Gaussian mixture prior, by setting a hyperprior on Gaussian priors, resulting in a procedure that is both adaptive and optimal in the minimax sense, which is shown with the second general theorem in this chapter.

Chapter 7 is our first attempt to tackle discretely observed inverse problems, by using Gaussian conjugacy in linear problems. The word ‘discretely’ refers to the fact that the observations are only recorded at certain points, in contrast to the continuous case, where the entire trajectory is observed. A sequence formulation for discretely observed inverse problems is derived based on the singular value decomposition of the smoothing operator. The main results include the contraction rates, and the credible sets for mildly and extremely ill-posed inverse problems. The results are further exemplified by simulations.

Chapter 8 is our second effort to analyse the discretely observed inverse problems. In this chapter, we use the results of Chapter 6 to examine the model of interest. For this we consider a signal reconstruction technique that builds continuous trajectories by interpolating discrete points, for which we also derive error estimates. This is followed by extending the general theorems from Chapter 6

that are valid for inverse problems with continuous observations to the discretely observed models. This leads to derivation for the contraction rates of random series priors, Gaussian priors, and Gaussian mixtures.

Part III studies the Bayesian nonparametric inference for stochastic evolution equations. The state space of the stochastic evolution equation is assumed to be infinite-dimensional and this dynamical system is assumed to be affected by a random noise process. In this part the dynamics are assumed to be driven by a deterministic component described by a linear partial differential equation, and a stochastic component described by additive space-time Gaussian noise. This model can be considered as a lifting of the white noise model to infinite-dimensional state space. The goal in this part is to derive theory for Bayesian inference concerning the drift and initial condition of the evolution equations.

Chapter 9 provides a general overview of stochastic evolution equations. We start with defining Q -Wiener processes, a generalisation of Brownian motion to infinite-dimensional space. These Q -Wiener processes allow us to formulate a class of stochastic integrals in Hilbert spaces, which is extended to a larger class of integrands. These integrals combined with the theory of deterministic evolution equations are then used to define stochastic evolution equations formally. We also demonstrate that the solutions to these equations can be represented by stochastic processes.

Chapter 10 studies the asymptotic performance of Gaussian priors in the recovery of the initial conditions and drifts of stochastic evolution equations. To obtain the results on posterior contraction, we modify the general approach established in Chapter 4. For the recovery of initial conditions, we introduce a spatial Gaussian prior and derive the rate at which the posterior distribution concentrates around the true initial condition. This task can be viewed as a generalisation of the extremely ill-posed problem from Chapter 7. The recovery of the drift component using Bayesian methods is also investigated for which we construct another type of Gaussian priors. These priors are developed to identify the anisotropic smoothness of signals in higher dimensions. A crucial insight is that we can apply a transformation to the observations so the model is converted into a multi-indexed Gaussian sequence model. This allows us to derive results for the contraction rates.

Samenvatting

In dit proefschrift wordt de theorie achter het asymptotische gedrag van Bayesiaanse inferentie voor niet-parametrische Gaussische lineaire modellen ontwikkeld. De modellen hebben de vorm van

$$\text{Observaties} = \begin{array}{l} \text{Getransformeerde signaal} \\ \text{zoals een gladgestreken parameter} \end{array} + \text{Gaussische ruis,}$$

en het doel is om het originelt signaal te reconstrueren met behulp van de Bayesiaanse methoden.

In Deel I beschrijven we de benodigde concepten voor het ontwikkelen van de theorie in dit proefschrift. In het bijzonder richten we ons op de theorie voor het getransformeerde-sigitaal-in-ruis model. Het doel is om het (originele) signaal terug te schatten met zo min mogelijke aannames op de parameter ruimte, en met een hoge convergentiesnelheid. Aangenomen wordt dat de parameter (signaal) aan bepaalde gladheidsvoorwaarden voldoet, welke gekwantificeerd worden in het eerste gedeelte van dit proefschrift. De ruisstructuur van het statistisch model is ook van belang, omdat het een grote invloed heeft op de moeilijkheid van het schattingsprobleem. Voor het schatten gebruiken we Bayesiaanse methoden, welke formeel gedefinieerd worden in Hoofdstuk 4.

In Hoofdstuk 2 worden gladheidsklassen geïntroduceerd om de mate van regulariteit van de parameterruimte te kwantificeren. Een eerste gladheidsklasse die we bespreken zijn zogeheten gladheidsschalen. Een gladheidsschaal $\{H_s\}_{s \in \mathbb{R}}$ is een collectie van Hilbertruimtes H_s die op basis van de normen is genest, en welke voldoet aan een zogeheten norm dualiteit. Een bijzonder subklasse van gladheidsschalen zijn Hilbertschalen. Een Hilberschaal is een gladheidsschaal waarbij de index s op natuurlijke wijze een genererende operator definieert. Met de resulterende genererende operator kan er een link gelegd worden tussen de Hilbertschaal en de covariantie operator van Gaussische a priori verdelingen en Gaussische ruis. Een ander type gladheidsklasse dat aan bod komt in dit hoofdstuk wordt gebruikt om de anisotropische gladheid in hogere dimensies te beschrijven. Deze laatste gladheidsklasse wordt gebruikt om stochastische evolutie vergelijkingen te bestuderen in Hoofdstuk III. We definiëren ook benaderingsgetallen die de benaderingseigenschappen van gladheidsklassen beschrijven, en we relateren deze aan metrische entropiën.

In Hoofdstuk 3 geven we een formele beschrijving van Gaussische maten in oneindig dimensionale ruimtes. Eerst geven we een overzicht van kansmaten op

Banachruimtes in het algemeen. Nadat we Gaussische maten formeel hebben gedefinieerd, laten we zien dat ze gunstige eigenschappen hebben. Van bijzonder belang is de covariantie structuur van Gaussische maten, wat we benadrukken met een aantal voorbeelden. Daarnaast geven we ook een kort overzicht van cilindrische maten en radonificatie in dit hoofdstuk.

In Hoofdstuk 4 geven we een formele beschrijving van Bayesiaanse inferentie in oneindig dimensionale ruimtes. Hiervoor introduceren we de regel van Bayes voor algemene functieruimtes, en we bespreken de basisconcepten van de asymptotiek zoals consistentie en convergentiesnelheid. De regel van Bayes leidt tot een natuurlijke schattingsprocedure, welke we toepassen op Gaussische lineaire modellen. We zullen in het bijzonder geïnteresseerd zijn in het geval met continue observaties, en het geval met discrete observaties. Het hoofdstuk wordt afgesloten met een algemene stelling over a posteriori convergentiesnelheden met een algemene bewijsmethode welke de basis vormt voor soortgelijke resultaten in de rest van dit proefschrift.

In Deel II richten we op lineaire inverse problemen met Gaussische ruis, waarvoor de theorie van Gaussische lineaire modellen toepasbaar is. In dit geval observeren we een gladgestreken signaal dat verontreinigd is door ruis. Door het gladstrijken van het signaal wordt het schatten bemoeilijkt. De naïeve benadering resulteert in slechte schattingen van het origineel signaal, doordat de ruis opgeblazen wordt door de inverse operator. Een oplossing is om de inverse operator te regulariseren, wat op natuurlijke wijze gedaan kan worden met de Bayesiaanse methode. In dit gedeelte van de proefschrift bestuderen we inverse problemen voor twee observatieregimes: Continu en discreet geobserveerde inverse problemen die ook bekend staan als het witte ruis model en regressie.

In Hoofdstuk 5 worden lineaire inverse problemen heuristisch en wiskundig rigoreus beschreven. Om het geheel tastbaar te maken geven we ook enkele voorbeelden. In dit hoofdstuk wordt ook de relatief algemene projectie van Galerkin belicht, omdat we hiermee vrij gemakkelijk convergentiesnelheden voor inverse problemen kunnen afleiden zoals beschreven is in de rest van Deel II.

In Hoofdstuk 6 beschrijven we het asymptotisch gedrag van niet-parametrische Bayesiaanse methoden toegepast op lineaire inverse problemen met continue observaties. Een algemene stelling voor continu geobserveerde inverse problemen wordt bewezen door een aanpassing in het bewijs van de algemene stelling over convergentiesnelheden dat afgeleid is in Hoofdstuk 4. In het algemeen is adaptiviteit niet gegarandeerd, want daar zijn extra condities voor nodig welke afgeleid en beschreven zijn in de tweede stelling in dit hoofdstuk. De eerste stelling wordt gebruikt om het asymptotisch gedrag van twee Bayesiaanse procedures te bestuderen: De procedure op basis van a priori verdelingen die beschreven zijn als gerandomiseerde reeksen, en de procedure op basis van Gaussische a priori verdelingen. De a priori verdelingen die beschreven worden als gerandomiseerde reeksen resulteren in a posteriori verdelingen die, op een vermenigvuldiging met een logaritmische factor na, convergeren met de minimax convergentiesnelheid. De Gaussische a priori verdelingen leiden echter tot a posteriori verdelingen die met een suboptimale convergentiesnelheid krimpen. Om het laatste geval te verbeteren introduceren wij mengvormen van Gaussische a priori verdelingen, door een hogere orde a priori verdeling te kiezen op de Gaussische a priori verdelingen, wat resulteert in

een procedure die adaptief en optimaal is in de minimax zin. Dit laatste wordt bewezen met behulp van de tweede stelling van dit hoofdstuk.

In Hoofdstuk 7 beschrijven we onze eerste aanpak van discreet geobserveerde inverse problemen. Dit doen we door gebruik te maken van de Gaussische conjugatie eigenschap in lineaire problemen. Hier wordt “discreet” gebruikt om aan te geven dat de observaties alleen op discrete punten zijn gemeten, terwijl in het continue geval het gehele pad geobserveerd is. Een reeksformulering van de discreet geobserveerde inverse problemen is afgeleid op basis van een singulierewaardenontbinding van de gladstrijkende operator. De hoofresultaten beschrijven de convergentiesnelheden en de a posteriori plausibele verzamelingen voor milde en extreem singuliere inverse problemen. De resultaten worden geïllustreerd met simulaties.

In Hoofdstuk 7 beschrijven we onze tweede aanpak van discreet geobserveerde inverse problemen. In dit hoofdstuk gebruiken we de resultaten van Hoofdstuk 6 om het model te bestuderen. Hiervoor gebruiken we een signaalreconstructietechniek dat continue paden construeert door de discrete punten te interpoleren, waarvoor de interpolatiefout ook geschat kan worden. We veralgemeniseren dan de stelling van Hoofdstuk 6, welke alleen geldt voor inverse problemen met continue observaties, zodat het ook toepasbaar is op discrete geobserveerde inverse problemen. Dit resulteert in convergentiesnelheden voor a priori verdelingen die beschreven worden als gerandomiseerde reeksen, Gaussische a priori verdelingen, en mengvormen van Gaussische a priori verdelingen.

In Deel III richten we ons op niet-parametrisch Bayesiaanse inferentie voor stochastische evolutie vergelijkingen. Er wordt aangenomen dat de toestandsruimte van de stochastische evolutie vergelijking oneindig dimensionaal is, en dat de bijbehorende dynamische systeem beïnvloed wordt door een ruisproces. In dit gedeelte wordt er van uitgegaan dat de dynamica gedreven wordt door een deterministische component die beschreven wordt door een lineaire partiële differentiaalvergelijking, en een stochastische component die beschreven wordt door een additieve Gaussische ruis over zowel de toestandsruimte als de tijdsdimensie. Dit model kan gezien worden als een bevordering van het witte ruis model naar een oneindig dimensionale toestandsruimte. Het doel in dit gedeelte van de proefschrift is om theorie te ontwikkelen voor Bayesiaanse inferentie voor de driftterm en de beginwaarden van de evolutie vergelijkingen.

In Hoofdstuk 9 geven we een algemene overzicht van stochastische evolutie vergelijkingen. Eerst definiëren we \mathcal{Q} -Wienerprocessen, een veralgemenisering van Brownse beweging naar oneindig dimensionale ruimtes. Met deze \mathcal{Q} -Wienerprocessen kunnen we een klasse van stochastische integralen over Hilbertruimtes beschrijven, welke daarna ook veralgemeniseerd wordt zodat het toepasbaar is op een grotere klasse van integranden. Deze integralen gecombineerd met de theorie van deterministische evolutie vergelijkingen worden dan gebruikt om stochastische evolutie vergelijkingen formeel te definiëren. We laten ook zien dat de oplossingen van deze vergelijkingen gerepresenteerd kunnen worden als stochastische processen.

In Hoofdstuk 10 beschrijven we het asymptotisch gedrag van de Bayesiaanse schattingsprocedure voor de beginwaarden en de driftterm van de stochastische evolutie vergelijkingen op basis van Gaussische a priori verdelingen. Om convergentiesnelheden af te leiden passen we de algemene aanpak die beschreven is

in Hoofdstuk 4 aan. Om de beginwaarden te schatten construeren we een type spatiële Gaussische a priori verdelingen, waarvoor we ook de snelheid kunnen afleiden waarmee de a posteriori verdelingen zich concentreert om de ware beginwaarden. Dit schattingsprobleem kan gezien worden als een veralgemenisering van het extreme singuliere inverse probleem dat beschreven wordt in Hoofdstuk 7. Voor de Bayesiaanse schattingsprocedure van de driftterm construeren we een ander type Gaussische a priori verdelingen. Deze a priori verdelingen zijn ontwikkeld om de anisotropische gladheid van signalen in hoger dimensies te identificeren. Een cruciaal punt is dat we de observaties dusdanig kunnen transformeren zodat het model geconverteerd wordt in een multi-index Gaussische reeks model. Met dit inzicht kunnen we de convergentiesnelheden afleiden.

Acknowledgement

I am in debt to many people for their support during my PhD study.

First, I would like to thank Prof. Aad van der Vaart and Dr Shota Gugushvili for offering the PhD position to me. I am most thankful for their guidance and supervision to my PhD research. In addition, they also gave me large freedom to choose research topics and tolerated my stubbornness, a lot. Furthermore, I am also grateful to Aad for his generous support on academic activities, and to Shota for his suggestions on interesting workshops, and exotic spirits.

Second, I also owe many thanks to my colleagues at NN Re. They showed great understanding and indulged me with the freedom to work on my thesis aside to my regular job. I am particularly grateful to them for sharing the workloads, especially the reporting burden.

Third, to my friends, I am appreciated to your company in the past years. I will always cherish the memories shared by us all: the conferences, not only academic discussions but also drinks in the evening; the gaming nights, online and offline; trips, by plane or by bike; badminton games, competitive or recreational; talks, discussions, and coffee breaks... It is truly a regrettable fact that we only shared the memories of a part of the journey. Many of you have relocated to other countries or even continents. Especially, I would like to thank Alexander Ly, for sharing the passion for mathematics and other more lively subjects, and the assistance on revising and translating the summary of this thesis. And of course, as always, I feel indebted and grateful to my beloved Xi for her companionship, understanding, and support, without whom I could not finalize this thesis.

In the end, to my family, no word can nearly express my gratitude. I am in great debt to my grandparents, who had made the right decisions at several points in my lifetime so that it became possible for me to enter the academic world, as well as my mother, who has been always supportive and protective in my entire life.

Curriculum Vitae

Dong Yan was born on January 10, 1989 in Yangxi, Guangdong, People's Republic of China.

After graduating from Shenzhen Middle School, China, in 2007 Dong began his bachelor study at Delft University of Technology, the Netherlands. In 2010, he earned a Bachelor of Science in Aerospace Engineering (cum laude). He was awarded a van Effen scholarship, which allowed him to continue his Master's study at the same university. He completed the Master's programmes in Applied Mathematics and Aerospace Engineering.

In February 2014, Dong started his PhD project, funded by European Research Council, at the Mathematical Institute, Leiden University. He conducted the research concentrated on Bayesian nonparametric inference for Gaussian linear models. The project was supervised by Aad van der Vaart and Shota Gugushvili. During his PhD, Dong also lectured or assisted multiple courses, and attended several conferences and workshops.

In November 2017, Dong joined NN Group, a Dutch insurance and asset management company. He worked as a quantitative analyst at NN Re, the reinsurance business unit.

